

Linear Regression

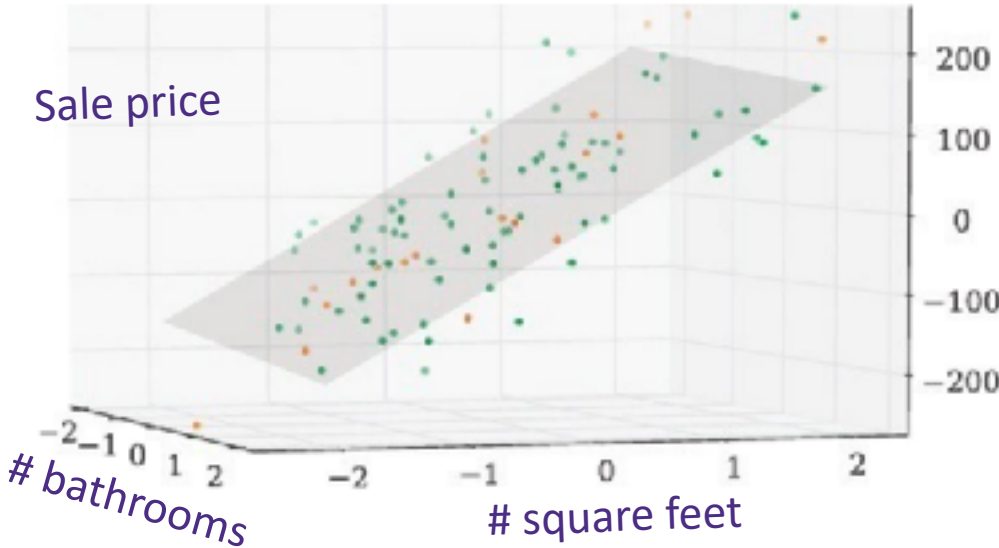


The regression problem, d-dimensions

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ House sale price *from*

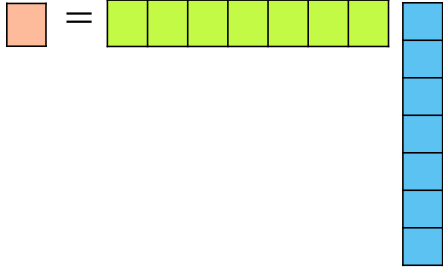
$x =$ {# sq. ft., # baths, date of sale, etc.}



Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

Hypothesis/Model: linear

$y_i = x_i^T w + \epsilon_i$ $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

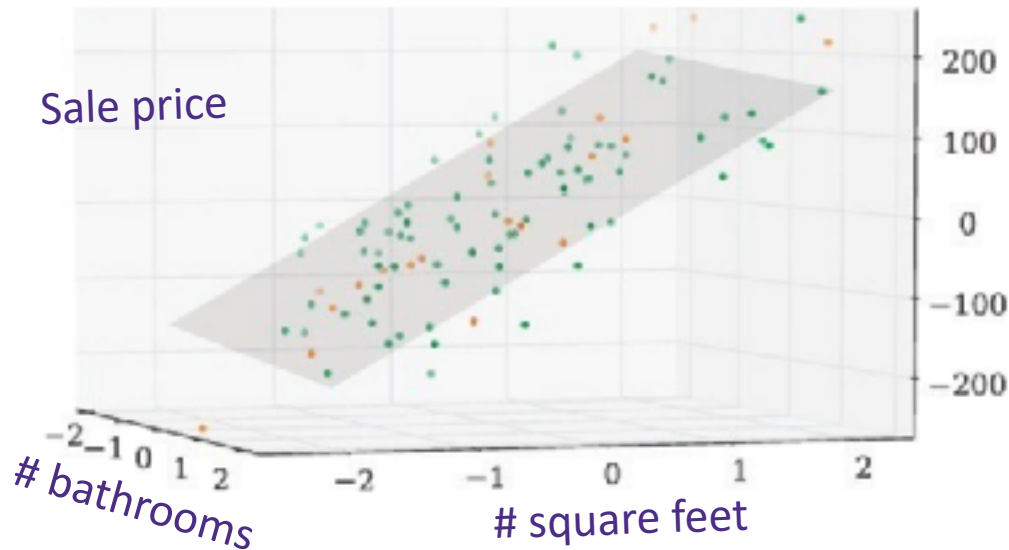


The regression problem, d-dimensions

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ House sale price from

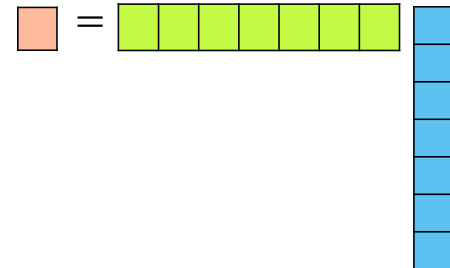
$x = \{\# \text{ sq. ft.}, \# \text{ baths}, \text{date of sale}, \text{etc.}\} \quad x \in \mathbb{R}^d$



Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

Hypothesis/Model: linear

$$y_i = x_i^T w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

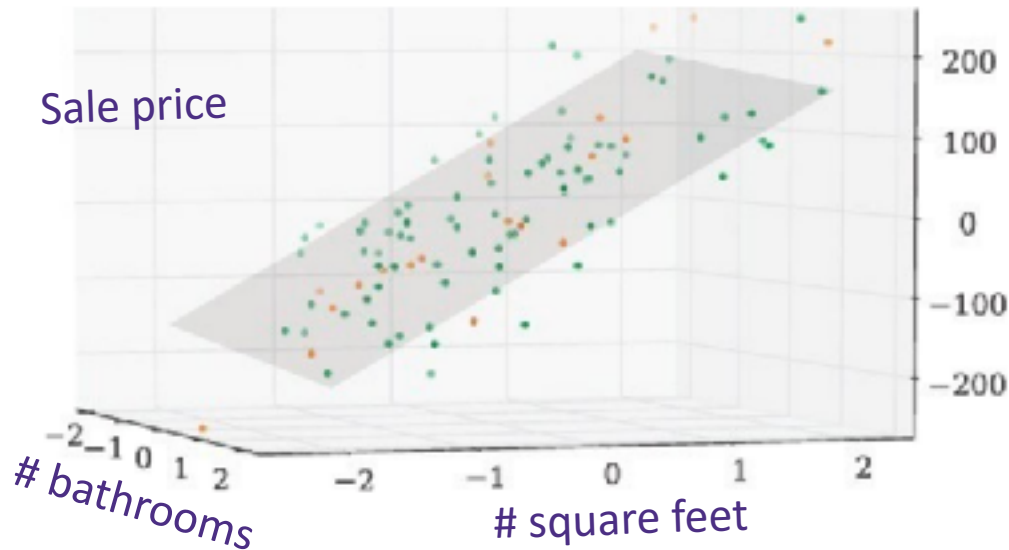


The regression problem, d-dimensions

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ House sale price from

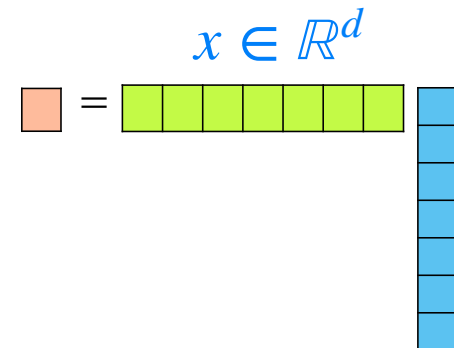
$x = \{\# \text{ sq. ft.}, \# \text{ baths}, \text{date of sale, etc.}\} \quad x \in \mathbb{R}^d$



Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

Hypothesis/Model: linear

$$y_i = x_i^T w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

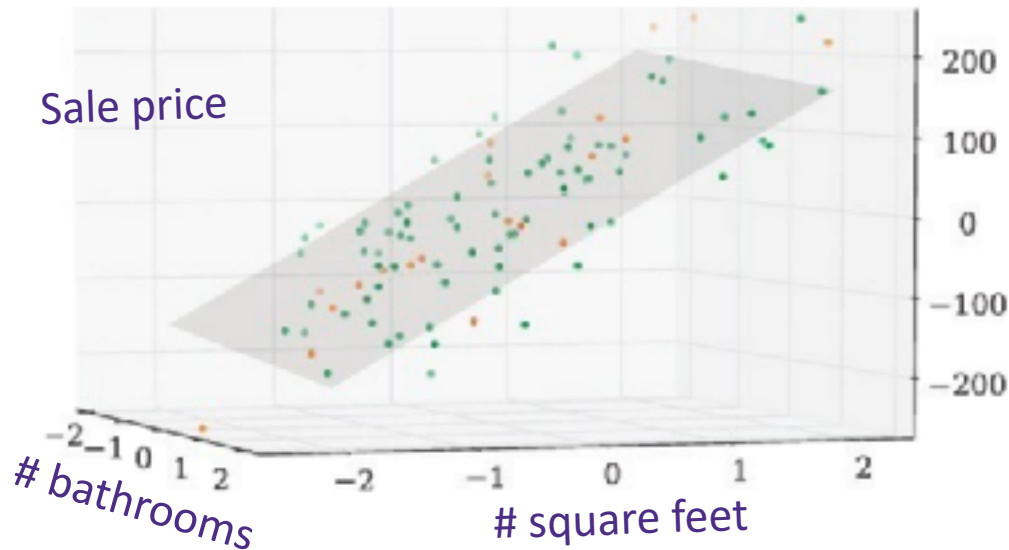


The regression problem, d-dimensions

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ House sale price from

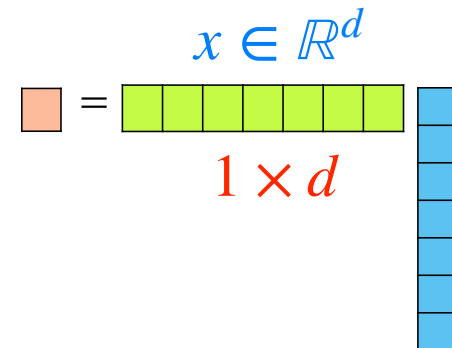
$x = \{\# \text{ sq. ft.}, \# \text{ baths}, \text{date of sale, etc.}\} \quad x \in \mathbb{R}^d$



Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

Hypothesis/Model: linear

$$y_i = x_i^T w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

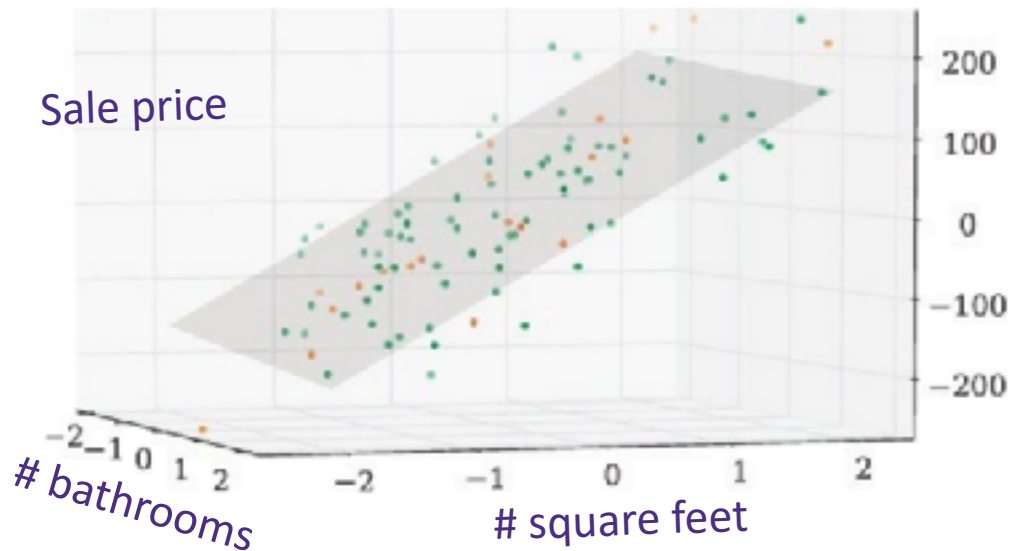


The regression problem, d-dimensions

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ House sale price from

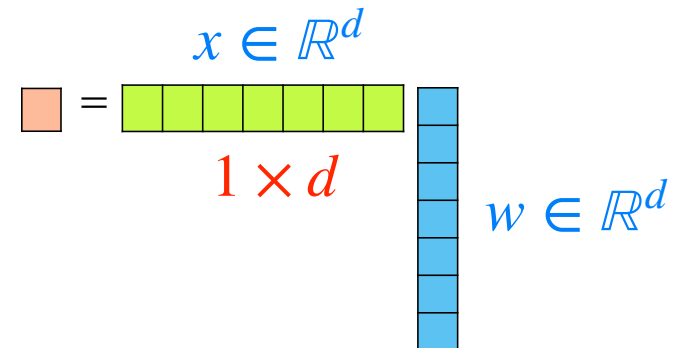
$x = \{\# \text{ sq. ft.}, \# \text{ baths}, \text{date of sale}, \text{etc.}\} \quad x \in \mathbb{R}^d$



Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

Hypothesis/Model: linear

$$y_i = x_i^T w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

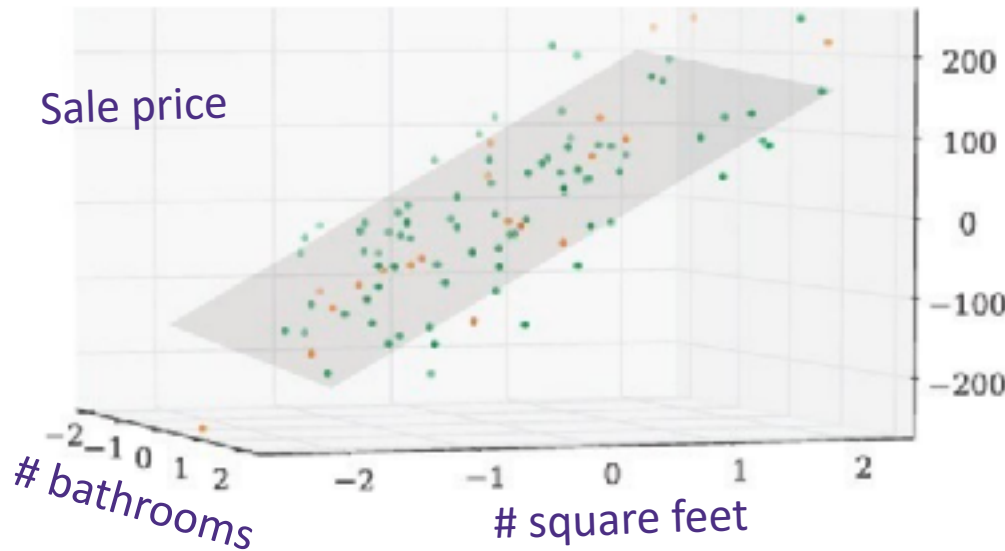


The regression problem, d-dimensions

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ House sale price from

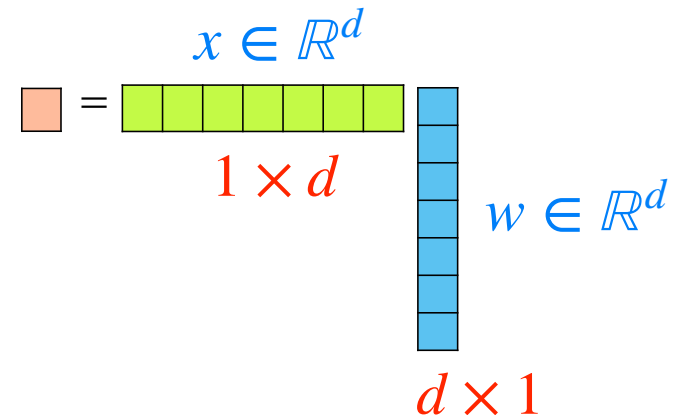
$x = \{\# \text{ sq. ft.}, \# \text{ baths}, \text{date of sale, etc.}\} \quad x \in \mathbb{R}^d$



Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

Hypothesis/Model: linear

$$y_i = x_i^T w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

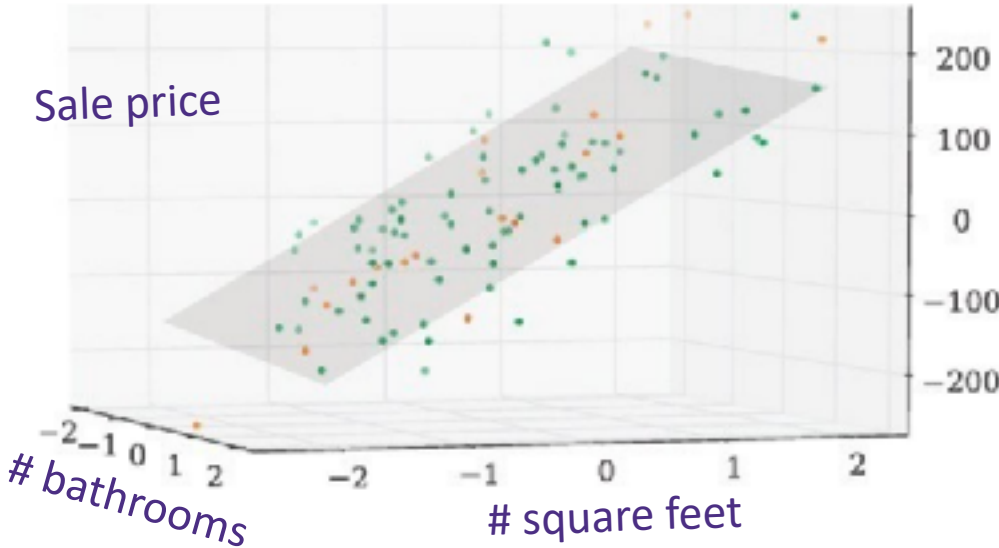


The regression problem, d-dimensions

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ House sale price from

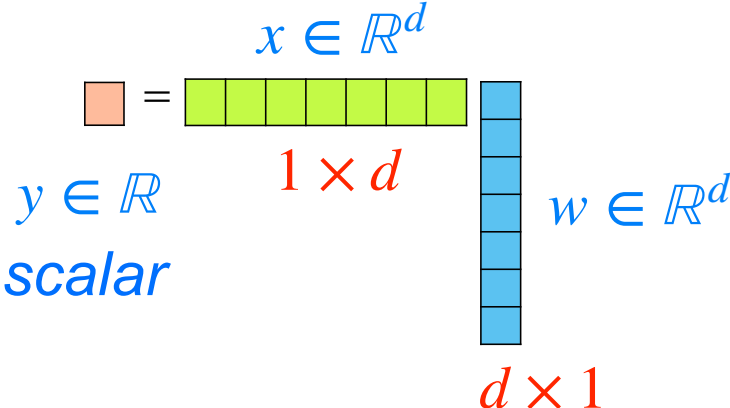
$x = \{\# \text{ sq. ft.}, \# \text{ baths, date of sale, etc.}\} \quad x \in \mathbb{R}^d$



Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

Hypothesis/Model: linear

$y_i = x_i^T w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

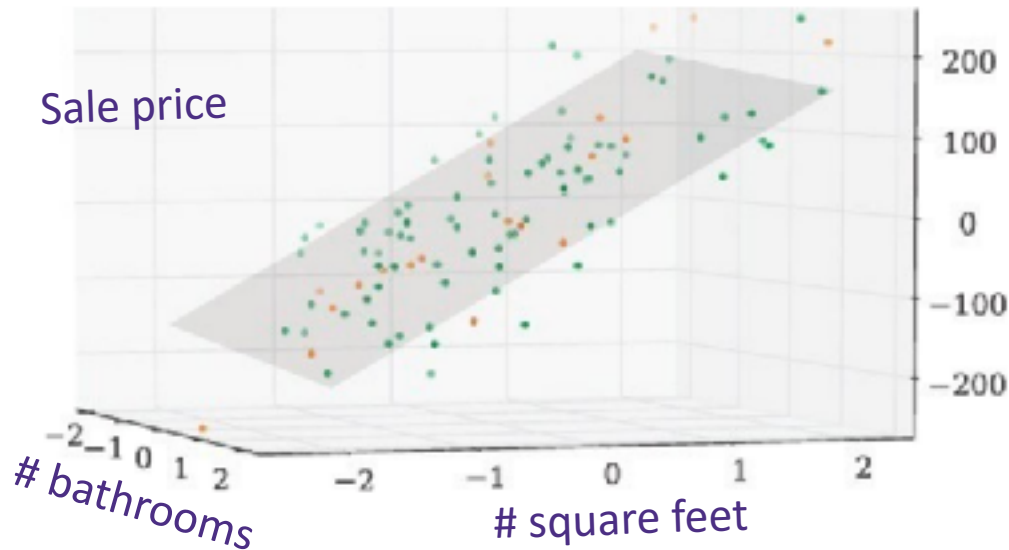


The regression problem, d-dim

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ House sale price *from*

$x =$ {# sq. ft., # baths, date of sale, etc.}



Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

Hypothesis/Model: linear

$$y_i = x_i^T w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

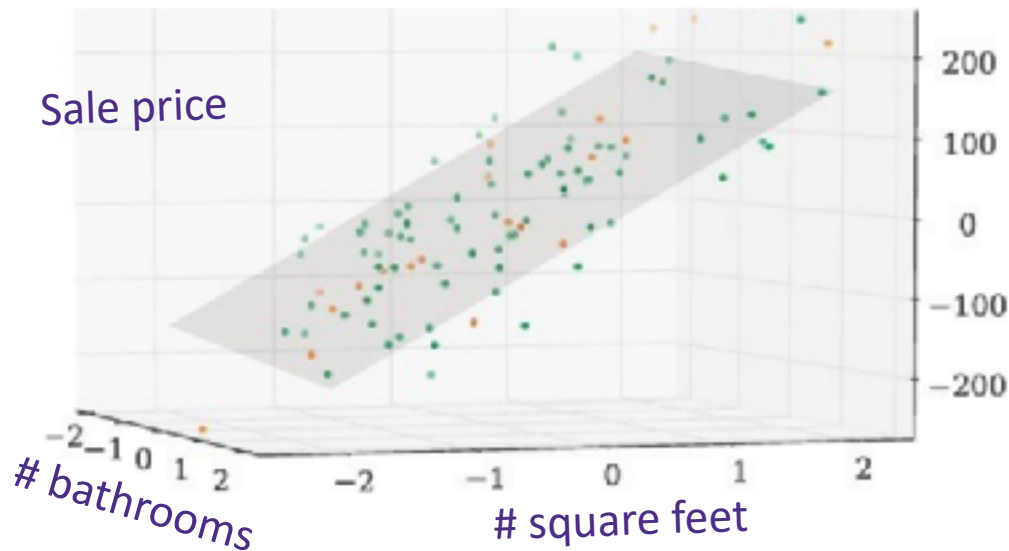
$$p(y|x, w, \sigma) =$$

The regression problem, d-dim

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ House sale price *from*

$x =$ {# sq. ft., # baths, date of sale, etc.}



Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

Hypothesis/Model: linear

$$y_i = x_i^T w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

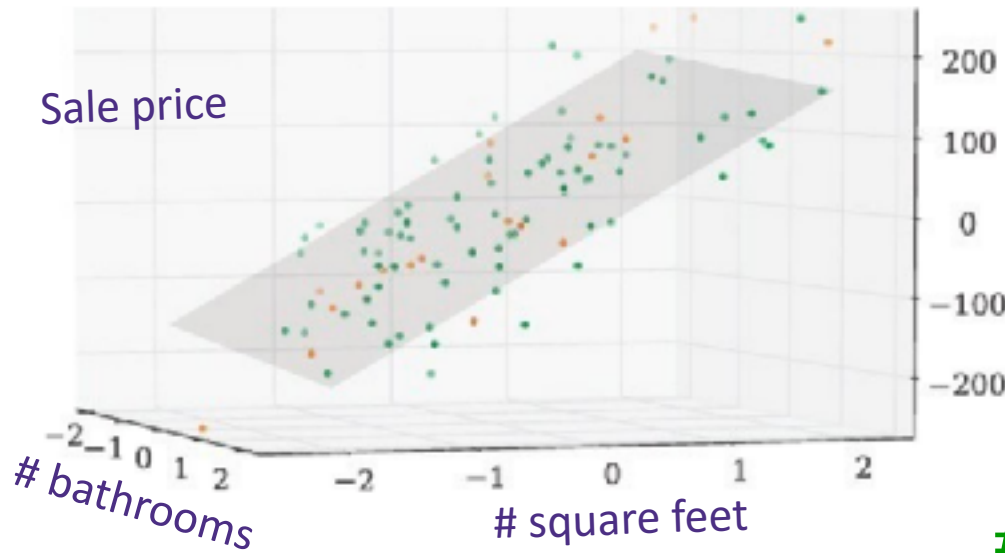
$$p(y|x, w, \sigma) = \mathcal{N}(x^T w, \sigma^2)$$

The regression problem, d-dim

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ House sale price *from*

$x =$ {# sq. ft., # baths, date of sale, etc.}



Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

Hypothesis/Model: linear

$$y_i = x_i^T w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

$$p(y|x, w, \sigma) = \mathcal{N}(x^T w, \sigma^2)$$

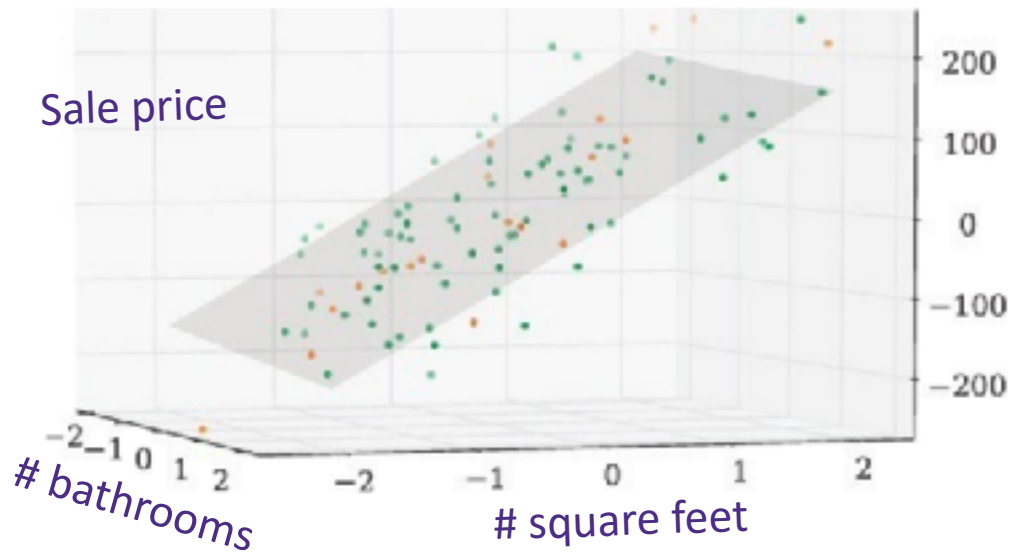
go back to formula sheet,
retrieve PDF of Gaussian...

The regression problem, d-dim

Given past sales data on [zillow.com](https://www.zillow.com), predict:

$y =$ House sale price *from*

$x =$ {# sq. ft., # baths, date of sale, etc.}



Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

Hypothesis/Model: linear

$$y_i = x_i^T w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

$$p(y|x, w, \sigma) = \mathcal{N}(x^T w, \sigma^2)$$

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^T w)^2/2\sigma^2}$$

What are we optimizing for then?

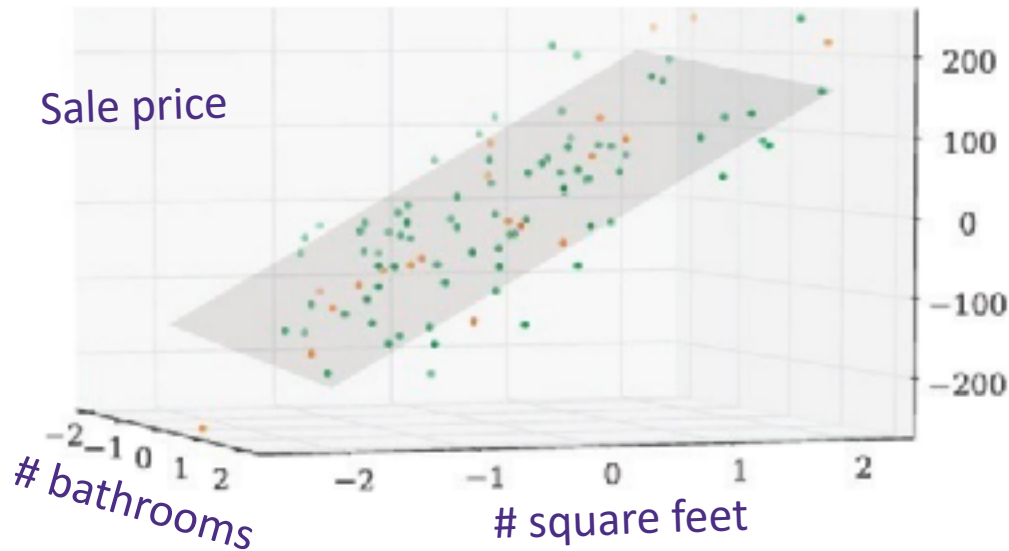
We discussed minimizing squared error/finding the best linear fit

Loss: least squares solution

$$\min_w \sum_{i=1}^n (y_i - x_i w)^2$$

But we could also do max likelihood!

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$$



Maximizing log-likelihood for linear regression

Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$$

Likelihood: $P(\mathcal{D}|w, \sigma) = \prod_{i=1}^n p(y_i|x_i, w, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2}$

Maximizing log-likelihood for linear regression

Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$$

Likelihood: $P(\mathcal{D}|w, \sigma) = \prod_{i=1}^n p(y_i|x_i, w, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2}$

Yay, we have a new likelihood! Now, how do we find MLE?

Maximizing log-likelihood, step 1

Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$$

Likelihood: $P(\mathcal{D}|w, \sigma) = \prod_{i=1}^n p(y_i|x_i, w, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2}$

$$\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2} \right) \# \text{ take log}$$

Maximizing log-likelihood, step 1

Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$$

Likelihood: $P(\mathcal{D}|w, \sigma) = \prod_{i=1}^n p(y_i|x_i, w, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2}$

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2} \right)$ # take log

Maximizing log-likelihood, step 1

Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$$

Likelihood: $P(\mathcal{D}|w, \sigma) = \prod_{i=1}^n p(y_i|x_i, w, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2}$

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

Maximizing log-likelihood, step 1

Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$$

Likelihood: $P(\mathcal{D}|w, \sigma) = \prod_{i=1}^n p(y_i|x_i, w, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2}$

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

what is log of product?

Maximizing log-likelihood, step 1

Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$$

Likelihood: $P(\mathcal{D}|w, \sigma) = \prod_{i=1}^n p(y_i|x_i, w, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2}$

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

what is log of product?

$$\log AB =$$

Maximizing log-likelihood, step 1

Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$$

Likelihood: $P(\mathcal{D}|w, \sigma) = \prod_{i=1}^n p(y_i|x_i, w, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2}$

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

what is log of product?

$$\log AB = \log A + \log B$$

Maximizing log-likelihood, step 1

Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$$

Likelihood: $P(\mathcal{D}|w, \sigma) = \prod_{i=1}^n p(y_i|x_i, w, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2}$

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^\top w)^2}{2\sigma^2} \right) \right]$$

what is log of product?
 $\log AB = \log A + \log B$

Maximizing log-likelihood, step 2

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i - x_i^T w)^2 / 2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T w)^2}{2\sigma^2} \right) \right]$$

Maximizing log-likelihood, step 2

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i - x_i^T w)^2 / 2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T w)^2}{2\sigma^2} \right) \right]$$

pull out terms that don't depend on i

Maximizing log-likelihood, step 2

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i - x_i^T w)^2 / 2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T w)^2}{2\sigma^2} \right) \right]$$

$$= n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

pull out terms that don't depend on i

Maximizing log-likelihood, step 2

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i - x_i^T w)^2 / 2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T w)^2}{2\sigma^2} \right) \right]$$

$$= n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

pull out terms that don't depend on i

we're going to max wrt w , so

Maximizing log-likelihood, step 2

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i - x_i^T w)^2 / 2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T w)^2}{2\sigma^2} \right) \right]$$

$$= n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

pull out terms that don't depend on i

we're going to max wrt w , so
get rid of terms that don't depend on w

Maximizing log-likelihood, step 2

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i - x_i^T w)^2 / 2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T w)^2}{2\sigma^2} \right) \right]$$

$$= n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

pull out terms that don't depend on i

we're going to max wrt w , so
get rid of terms that don't depend on w

Maximizing log-likelihood, step 2

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i - x_i^T w)^2 / 2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T w)^2}{2\sigma^2} \right) \right]$$

$$= n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

pull out terms that don't depend on i

we're going to max wrt w , so
get rid of terms that don't depend on w

$$- \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

Maximizing log-likelihood, step 2

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i - x_i^T w)^2 / 2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T w)^2}{2\sigma^2} \right) \right]$$

$$= n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

pull out terms that don't depend on i

we're going to max wrt w , so
get rid of terms that don't depend on w

$$- \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

We found our simplified LL!

Maximizing log-likelihood, step 2

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i - x_i^T w)^2 / 2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T w)^2}{2\sigma^2} \right) \right]$$

$$= n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

pull out terms that don't depend on i

we're going to max wrt w , so
get rid of terms that don't depend on w

$$- \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

We found our simplified LL!

Now what do?

Maximizing log-likelihood, step 2

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i - x_i^T w)^2 / 2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T w)^2}{2\sigma^2} \right) \right]$$

$$= n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

pull out terms that don't depend on i

we're going to max wrt w , so
get rid of terms that don't depend on w

$$\arg \max_w - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

We found our simplified LL!

Now what do?

Maximizing log-likelihood, step 2

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i - x_i^T w)^2 / 2\sigma^2} \right)$ # take log

first manipulate log likelihood to make it easier to work with

$$= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - x_i^T w)^2}{2\sigma^2} \right) \right]$$

$$= n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

pull out terms that don't depend on i

we're going to max wrt w , so
get rid of terms that don't depend on w

$$\arg \max_w - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2}$$

We found our simplified LL!

Now what do?

Maximizing a negative?

Maximizing log-likelihood, optimization

Training Data: $x_i \in \mathbb{R}^d$
 $\{(x_i, y_i)\}_{i=1}^n$ $y_i \in \mathbb{R}$

$$p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$$

Likelihood: $P(\mathcal{D}|w, \sigma) = \prod_{i=1}^n p(y_i|x_i, w, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2}$

Maximize (wrt w): $\log P(\mathcal{D}|w, \sigma) = \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_i-x_i^\top w)^2/2\sigma^2} \right)$

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

Maximizing log-likelihood → minimizing mean squared error (MSE)

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

Maximizing log-likelihood → minimizing mean squared error (MSE)

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

- Wow, minimizing squared error seems like a reasonable thing to do to fit a predictor!

Maximizing log-likelihood → minimizing mean squared error (MSE)

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

- Wow, minimizing squared error seems like a reasonable thing to do to fit a predictor!
 - When noise is Gaussian, maximizing likelihood of our linear model is equivalent to minimizing the sum of squared errors. We get this from cranking through the math.

Maximizing log-likelihood → minimizing mean squared error (MSE)

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

- Wow, minimizing squared error seems like a reasonable thing to do to fit a predictor!
 - When noise is Gaussian, maximizing likelihood of our linear model is equivalent to minimizing the sum of squared errors. We get this from cranking through the math.
 - However, if we had chosen a different noise distribution, we'd get a different estimator. E.g. Laplacian noise -> minimize absolute error

Maximizing log-likelihood → minimizing mean squared error (MSE)

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

- Wow, minimizing squared error seems like a reasonable thing to do to fit a predictor!
 - When noise is Gaussian, maximizing likelihood of our linear model is equivalent to minimizing the sum of squared errors. We get this from cranking through the math.
 - However, if we had chosen a different noise distribution, we'd get a different estimator. E.g. Laplacian noise -> minimize absolute error
 - Which one is better depends on assumptions about our data. Gaussian has no long tails, mean squared error incurs high penalty for outliers

Maximizing log-likelihood → minimizing mean squared error (MSE)

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

- Wow, minimizing squared error seems like a reasonable thing to do to fit a predictor!
 - When noise is Gaussian, maximizing likelihood of our linear model is equivalent to minimizing the sum of squared errors. We get this from cranking through the math.
 - However, if we had chosen a different noise distribution, we'd get a different estimator. E.g. Laplacian noise -> minimize absolute error
 - Which one is better depends on assumptions about our data. Gaussian has no long tails, mean squared error incurs high penalty for outliers
 - **TL;DR: Model of data ← math → what to optimize**

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

So how do we actually optimize this / fit our model?

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

So how do we actually optimize this / fit our model?

Set gradient=0, solve for w

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

So how do we actually optimize this / fit our model?

Set gradient=0, solve for w

Recall: $w \in \mathbb{R}^d$

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

So how do we actually optimize this / fit our model?

Set gradient=0, solve for w

Recall: $w \in \mathbb{R}^d$

$$\nabla_w f(x) = \begin{bmatrix} \frac{\partial}{\partial w_0} f(x) \\ \frac{\partial}{\partial w_1} f(x) \\ \vdots \end{bmatrix}$$

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

So how do we actually optimize this / fit our model?

Set gradient=0, solve for w

Recall: $w \in \mathbb{R}^d$

$$\nabla_w f(x) = \begin{bmatrix} \frac{\partial}{\partial w_0} f(x) \\ \frac{\partial}{\partial w_1} f(x) \\ \vdots \end{bmatrix}$$

Gradient wrt w is a vector of partial derivatives

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

So how do we actually optimize this / fit our model?

Set gradient=0, solve for w

Recall: $w \in \mathbb{R}^d$

$$\nabla_w f(x) = \begin{bmatrix} \frac{\partial}{\partial w_0} f(x) \\ \frac{\partial}{\partial w_1} f(x) \\ \vdots \end{bmatrix}$$

$$\nabla_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

Gradient wrt w is a vector of partial derivatives

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

So how do we actually optimize this / fit our model?

Set gradient=0, solve for w

Recall: $w \in \mathbb{R}^d$

$$\nabla_w f(x) = \begin{bmatrix} \frac{\partial}{\partial w_0} f(x) \\ \frac{\partial}{\partial w_1} f(x) \\ \vdots \end{bmatrix}$$

$$\nabla_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

Gradient wrt w is a vector of partial derivatives

Chain rule

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

So how do we actually optimize this / fit our model?

Set gradient=0, solve for w

$$\# \text{ Recall: } w \in \mathbb{R}^d \quad \nabla_w f(x) = \begin{bmatrix} \frac{\partial}{\partial w_0} f(x) \\ \frac{\partial}{\partial w_1} f(x) \\ \vdots \end{bmatrix}$$

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

Gradient wrt w is a vector of partial derivatives

Chain rule $\nabla_w f(x)^2 = 2f(x) \cdot \nabla_w f(x)$

Maximizing log-likelihood wrt w

So how do we actually optimize this / fit our model?

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

Set gradient=0, solve for w

Recall: $w \in \mathbb{R}^d$

$$\nabla_w f(x) = \begin{bmatrix} \frac{\partial}{\partial w_0} f(x) \\ \frac{\partial}{\partial w_1} f(x) \\ \vdots \end{bmatrix}$$

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

Gradient wrt w is a vector of partial derivatives

$$= \sum_{i=1}^n 2(y_i - x_i^T w)^1 \nabla_w (y_i - x_i^T w)$$

Chain rule $\nabla_w f(x)^2 = 2f(x) \cdot \nabla_w f(x)$

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

So how do we actually optimize this / fit our model?

Set gradient=0, solve for w

$$\# \text{ Recall: } w \in \mathbb{R}^d \quad \nabla_w f(x) = \begin{bmatrix} \frac{\partial}{\partial w_0} f(x) \\ \frac{\partial}{\partial w_1} f(x) \\ \vdots \end{bmatrix}$$

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

Gradient wrt w is a vector of partial derivatives

$$= \sum_{i=1}^n 2(y_i - x_i^T w)^1 \nabla_w (y_i - x_i^T w)$$

Chain rule $\nabla_w f(x)^2 = 2f(x) \cdot \nabla_w f(x)$

“Matrix Cookbook” identities (cheatsheet fodder)

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

So how do we actually optimize this / fit our model?

Set gradient=0, solve for w

$$\# \text{ Recall: } w \in \mathbb{R}^d \quad \nabla_w f(x) = \begin{bmatrix} \frac{\partial}{\partial w_0} f(x) \\ \frac{\partial}{\partial w_1} f(x) \\ \vdots \end{bmatrix}$$

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

Gradient wrt w is a vector of partial derivatives

$$= \sum_{i=1}^n 2(y_i - x_i^T w)^1 \nabla_w (y_i - x_i^T w)$$

Chain rule $\nabla_w f(x)^2 = 2f(x) \cdot \nabla_w f(x)$

"Matrix Cookbook" identities (cheatsheet fodder)

$$\nabla_w x^T w =$$

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

So how do we actually optimize this / fit our model?

Set gradient=0, solve for w

$$\# \text{ Recall: } w \in \mathbb{R}^d \quad \nabla_w f(x) = \begin{bmatrix} \frac{\partial}{\partial w_0} f(x) \\ \frac{\partial}{\partial w_1} f(x) \\ \vdots \end{bmatrix}$$

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

Gradient wrt w is a vector of partial derivatives

$$= \sum_{i=1}^n 2(y_i - x_i^T w)^1 \nabla_w (y_i - x_i^T w)$$

Chain rule $\nabla_w f(x)^2 = 2f(x) \cdot \nabla_w f(x)$

"Matrix Cookbook" identities (cheatsheet fodder)

$$\nabla_w x^T w = x$$

Maximizing log-likelihood wrt w

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

So how do we actually optimize this / fit our model?

Set gradient=0, solve for w

$$\# \text{ Recall: } w \in \mathbb{R}^d \quad \nabla_w f(x) = \begin{bmatrix} \frac{\partial}{\partial w_0} f(x) \\ \frac{\partial}{\partial w_1} f(x) \\ \vdots \end{bmatrix}$$

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

Gradient wrt w is a vector of partial derivatives

$$= \sum_{i=1}^n 2(y_i - x_i^T w)^1 \nabla_w (y_i - x_i^T w)$$

Chain rule $\nabla_w f(x)^2 = 2f(x) \cdot \nabla_w f(x)$

$$= \sum_{i=1}^n 2(y_i - x_i^T w)(-x_i)$$

"Matrix Cookbook" identities (cheatsheet fodder)

$$\nabla_w x^T w = x$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n 2(y_i - x_i^T w)(-x_i)$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n 2(y_i - x_i^T w)(-x_i) = 0$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n 2(y_i - x_i^T w)(-x_i) = 0 \quad \# \text{ check dimensions!}$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n 2(y_i - x_i^T w) \underbrace{(-x_i)}_{d \times 1} = 0 \quad \# \text{ check dimensions!}$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underset{1 \times 1}{2(y_i - x_i^T w)} \underset{d \times 1}{(-x_i)} = 0 \quad \# \text{ check dimensions!}$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n 2(y_i - x_i^T w) \begin{matrix} (-x_i) \\ 1 \times 1 \quad 1 \times 1 \quad d \times 1 \end{matrix} = 0 \quad \# \text{ check dimensions!}$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)}_{1 \times 1} \underbrace{(-x_i)}_{d \times 1} = 0 \quad \# \text{ check dimensions!}$$

$d \times 1$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)}_{d \times 1} \underbrace{(-x_i)}_{d \times 1} = 0 \quad \# \text{ check dimensions!}$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)}_{d \times 1} \underbrace{(-x_i)}_{d \times 1} = 0$$

check dimensions!
 $\vec{0} \in \mathbb{R}^d$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\begin{aligned} \nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 &= \sum_{i=1}^n 2(y_i - x_i^T w)(-x_i) = 0 && \# \text{ check dimensions!} \\ & \quad \underbrace{1 \times 1 \quad 1 \times 1 \quad d \times 1}_{d \times 1} \quad d \times 1 && \vec{0} \in \mathbb{R}^d \\ &= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0 \end{aligned}$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\begin{aligned}\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 &= \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)}_{1 \times 1} \underbrace{(-x_i)}_{d \times 1} = 0 \\ &= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0\end{aligned}$$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

Now for some tricks

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\begin{aligned}\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 &= \sum_{i=1}^n 2(y_i - x_i^T w)(-x_i) = 0 \\ &= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0\end{aligned}$$

$\underbrace{1 \times 1 \quad 1 \times 1 \quad d \times 1}_{d \times 1} \quad d \times 1$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

Now for some tricks

$$x^T w = a \quad \# \text{ scalar}$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\begin{aligned}\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 &= \sum_{i=1}^n 2(y_i - x_i^T w)(-x_i) = 0 \\ &= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0\end{aligned}$$

1×1 1×1 $d \times 1$ $d \times 1$
 $d \times 1$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

Now for some tricks

$$x^T w = a \quad \# \text{ scalar}$$

$$a \vec{v} = \vec{v} a$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)(-x_i)}_{d \times 1} = 0$$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

$$= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0$$

Now for some tricks

$$= 2 \sum_{i=1}^n [-x_i y_i + x_i x_i^T w] = 0$$

$$x_i^T w = a \quad \# \text{ scalar}$$

$$a \vec{v} = \vec{v} a$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)}_{d \times 1} \underbrace{(-x_i)}_{d \times 1} = 0$$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

$$= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0$$

Now for some tricks

$$= 2 \sum_{i=1}^n [-x_i y_i + x_i x_i^T w] = 0$$

$$x_i^T w = a \quad \# \text{ scalar}$$

$$a \vec{v} = \vec{v} a$$

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w$$

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)}_{d \times 1} \underbrace{(-x_i)}_{d \times 1} = 0$$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

$$= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0$$

Now for some tricks

$$= 2 \sum_{i=1}^n [-x_i y_i + x_i x_i^T w] = 0$$

$$x_i^T w = a \quad \# \text{ scalar}$$

$$a \vec{v} = \vec{v} a$$

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w$$

How can we isolate w?

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)}_{d \times 1} \underbrace{(-x_i)}_{d \times 1} = 0$$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

$$= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0$$

Now for some tricks

$$= 2 \sum_{i=1}^n [-x_i y_i + x_i x_i^T w] = 0$$

$$x_i^T w = a \quad \# \text{ scalar}$$

$$a \vec{v} = \vec{v} a$$

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w$$

How can we isolate w?

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)(-x_i)}_{d \times 1} = 0$$

1×1 1×1 $d \times 1$ $d \times 1$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

$$= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0$$

Now for some tricks

$$= 2 \sum_{i=1}^n [-x_i y_i + x_i x_i^T w] = 0$$

$$x_i^T w = a \quad \# \text{ scalar}$$

$$a \vec{v} = \vec{v} a$$

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w$$

$d \times d$

How can we isolate w?

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)(-x_i)}_{d \times 1} = 0$$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

$$= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0$$

Now for some tricks

$$= 2 \sum_{i=1}^n [-x_i y_i + x_i x_i^T w] = 0$$

$$x_i^T w = a \quad \# \text{ scalar}$$

$$a \vec{v} = \vec{v} a$$

Reminder

$x^T x$ = inner product or dot product 1x1

$x x^T$ = outer product dx d

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w$$

How can we isolate w?

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)(-x_i)}_{d \times 1} = 0$$

1×1 1×1 $d \times 1$ $d \times 1$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

$$= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0$$

Now for some tricks

$$= 2 \sum_{i=1}^n [-x_i y_i + x_i x_i^T w] = 0$$

$$x^T w = a \quad \# \text{ scalar}$$

$$a \vec{v} = \vec{v} a$$

Reminder

$x^T x$ = inner product or dot product 1×1

$x x^T$ = outer product $d \times d$

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w$$

$d \times d$

A matrix

How can we isolate w?

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)(-x_i)}_{d \times 1} = 0$$

1×1 1×1 $d \times 1$ $d \times 1$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

$$= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0$$

Now for some tricks

$$= 2 \sum_{i=1}^n [-x_i y_i + x_i x_i^T w] = 0$$

$$x_i^T w = a \quad \# \text{ scalar}$$

$$a \vec{v} = \vec{v} a$$

Reminder

$x^T x$ = inner product or dot product 1x1

$x x^T$ = outer product dx d

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w$$

$d \times d$

A matrix

How can we isolate w?

Take the inverse!

Maximizing log-likelihood, setting gradient = 0

Set gradient=0, solve for w

$$\nabla_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \sum_{i=1}^n \underbrace{2(y_i - x_i^T w)(-x_i)}_{d \times 1} = 0$$

1×1 1×1 $d \times 1$ $d \times 1$

check dimensions!

$$\vec{0} \in \mathbb{R}^d$$

$$= 2 \sum_{i=1}^n [-y_i x_i + x_i^T w x_i] = 0$$

Now for some tricks

$$= 2 \sum_{i=1}^n [-x_i y_i + x_i x_i^T w] = 0$$

$$x^T w = a \quad \# \text{ scalar}$$

$$a \vec{v} = \vec{v} a$$

Reminder

$x^T x$ = inner product or dot product 1x1

$x x^T$ = outer product dx d

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w^{-1}$$

$d \times d$

A matrix

How can we isolate w?

Take the inverse!

Maximizing log-likelihood, solving for w

Set gradient=0, solve for w

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w^{-1}$$

A matrix

How can we isolate w?

Take the inverse!

Maximizing log-likelihood, solving for w

Set gradient=0, solve for w

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w^{-1}$$

A matrix

How can we isolate w?

Take the inverse!

$$AA^{-1} = I$$

Maximizing log-likelihood, solving for w

Set gradient=0, solve for w

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w^{-1}$$

A matrix

How can we isolate w?

Take the inverse!

$$AA^{-1} = I$$

$$\left(\sum_{i=1}^n x_i x_i^T \right)^{-1} \sum_{i=1}^n x_i y_i = \hat{w}_{\text{MLE}}$$

Maximizing log-likelihood, solving for w

Set gradient=0, solve for w

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w$$

A matrix

How can we isolate w?

Take the inverse!

$$AA^{-1} = I$$

$$\left(\sum_{i=1}^n x_i x_i^T \right)^{-1} \sum_{i=1}^n x_i y_i = \hat{w}_{\text{MLE}}$$

Maximizing log-likelihood, solving for w

Set gradient=0, solve for w

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w$$

A matrix

How can we isolate w?

Take the inverse!

$$AA^{-1} = I$$

$$\left(\sum_{i=1}^n x_i x_i^T \right)^{-1} \sum_{i=1}^n x_i y_i = \hat{w}_{\text{MLE}}$$

The equation for fitting our model parameters w to our data

Maximizing log-likelihood, solving for w

Set gradient=0, solve for w

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w$$

A matrix

How can we isolate w?

Take the inverse!

$$AA^{-1} = I$$

$$\left(\sum_{i=1}^n x_i x_i^T \right)^{-1} \sum_{i=1}^n x_i y_i = \hat{w}_{\text{MLE}}$$

The equation for fitting our model parameters w to our data

- What constraints does taking the inverse place on our data?
 - Matrix must be invertible
 - $n \geq d$

Maximizing log-likelihood, solving for w

Set gradient=0, solve for w

$$\sum_{i=1}^n x_i y_i = \left(\sum_{i=1}^n x_i x_i^T \right) w^{-1}$$

A matrix

How can we isolate w?

Take the inverse!

$$AA^{-1} = I$$

$$\left(\sum_{i=1}^n x_i x_i^T \right)^{-1} \sum_{i=1}^n x_i y_i = \hat{w}_{MLE}$$

The equation for fitting our model parameters w to our data

- What constraints does taking the inverse place on our data?
 - Matrix must be invertible
 - $n \geq d$
- For this model we can analytically derive how to find the optimal weights. Will not always be true for more complex models

Maximizing log-likelihood, final form

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

Set gradient=0, solve for w

$$\hat{w}_{MLE} = \left(\sum_{i=1}^n x_i x_i^\top \right)^{-1} \sum_{i=1}^n x_i y_i$$

The regression problem in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

We ❤️ matrix notation!

d : # of features

n : # of examples/datapoints

The regression problem in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

We ❤️ matrix notation!

d : # of features

n : # of examples/datapoints

$$x \in \mathbb{R}^d$$

The regression problem in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

We ❤️ matrix notation!

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \begin{matrix} d \times 1 \\ d \times 1 \\ d \times 1 \end{matrix}$$

d : # of features

n : # of examples/datapoints

$$x \in \mathbb{R}^d$$

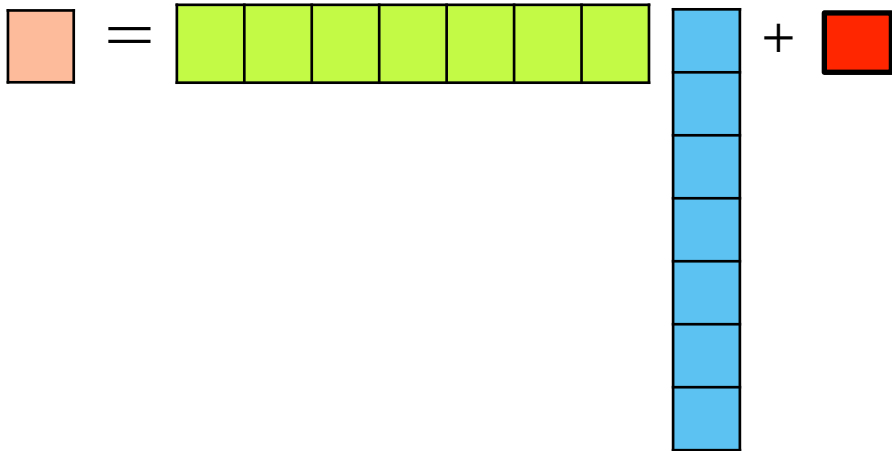
The regression problem in matrix notation, translating

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

$d \times 1$ $d \times 1$ $d \times 1$ d : # of features
 n : # of examples/datapoints
 $x \in \mathbb{R}^d$

Previously... $y_i = x_i^\top w + \epsilon_i$



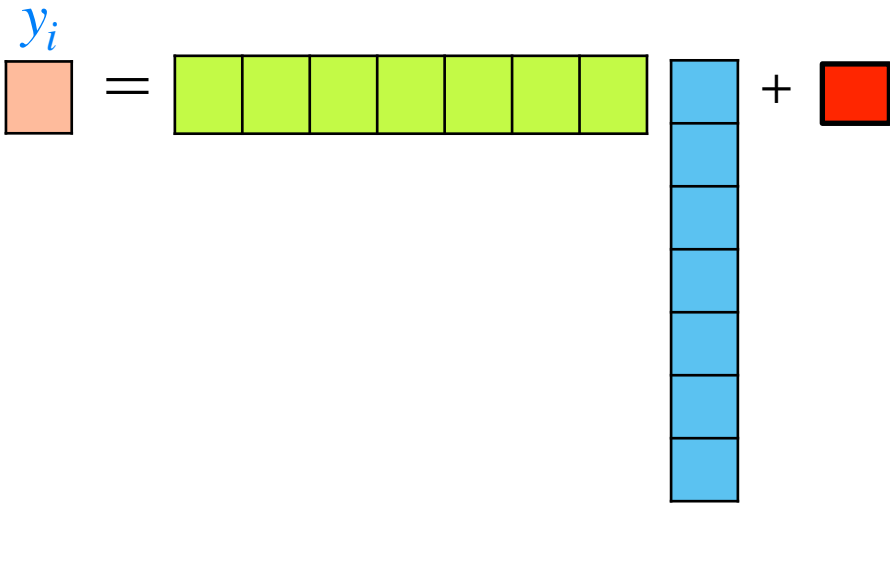
The regression problem in matrix notation, translating

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \quad \begin{array}{l} d \times 1 \\ d \times 1 \\ d \times 1 \end{array} \quad \begin{array}{l} d : \# \text{ of features} \\ n : \# \text{ of examples/datapoints} \end{array}$$

$x \in \mathbb{R}^d$

Previously... $y_i = x_i^T w + \epsilon_i$



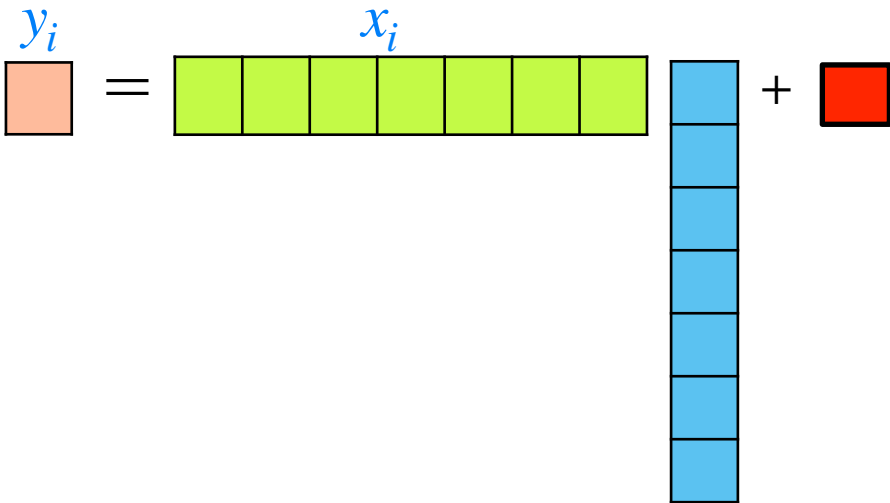
The regression problem in matrix notation, translating

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

$d \times 1$ $d \times 1$ $d \times 1$ d : # of features
 n : # of examples/datapoints
 $x \in \mathbb{R}^d$

Previously... $y_i = x_i^T w + \epsilon_i$



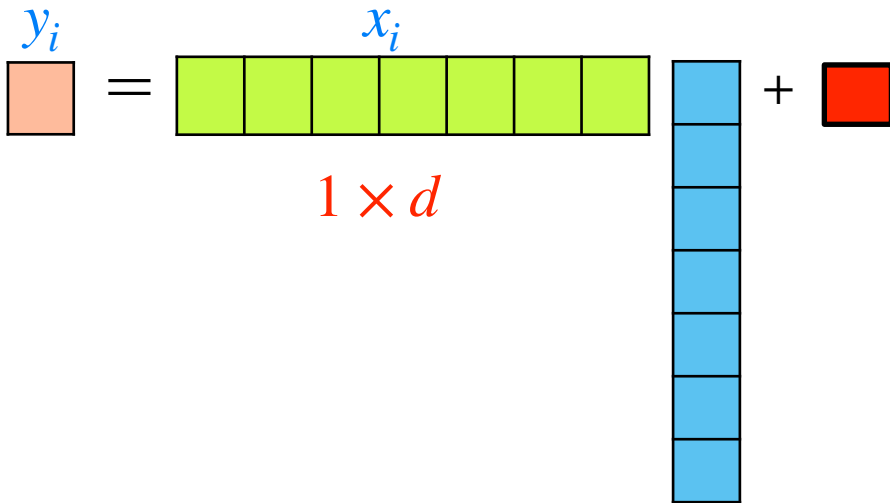
The regression problem in matrix notation, translating

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \quad \begin{array}{l} d \times 1 \\ d \times 1 \\ d \times 1 \end{array} \quad \begin{array}{l} d : \# \text{ of features} \\ n : \# \text{ of examples/datapoints} \end{array}$$

$x \in \mathbb{R}^d$

Previously... $y_i = x_i^T w + \epsilon_i$



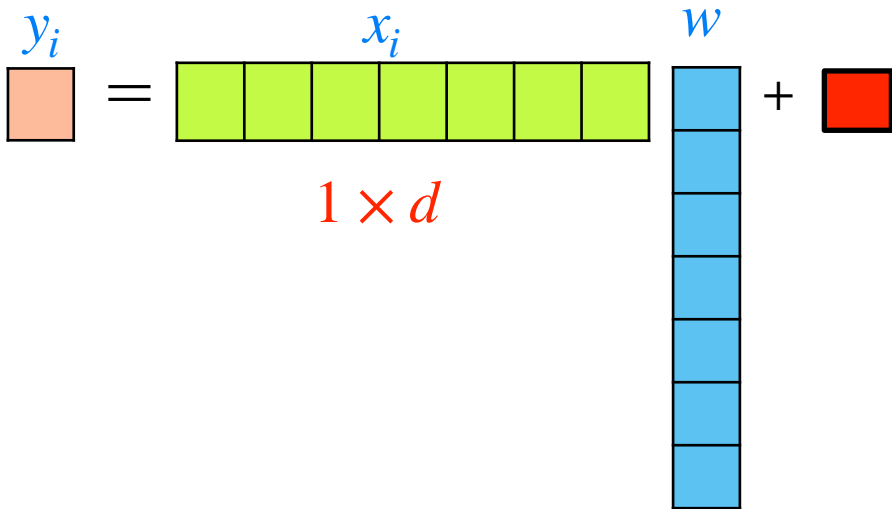
The regression problem in matrix notation, translating

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \quad \begin{array}{l} d \times 1 \\ d \times 1 \\ d \times 1 \end{array} \quad \begin{array}{l} d : \# \text{ of features} \\ n : \# \text{ of examples/datapoints} \end{array}$$

$x \in \mathbb{R}^d$

Previously... $y_i = x_i^T w + \epsilon_i$



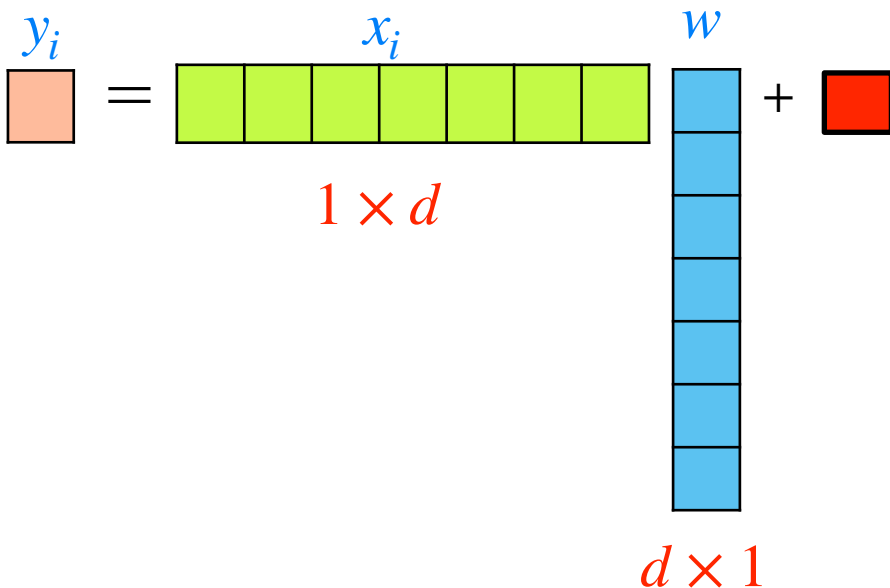
The regression problem in matrix notation, translating

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \quad \begin{array}{l} d \times 1 \\ d \times 1 \\ d \times 1 \end{array} \quad \begin{array}{l} d : \# \text{ of features} \\ n : \# \text{ of examples/datapoints} \end{array}$$

$x \in \mathbb{R}^d$

Previously... $y_i = x_i^T w + \epsilon_i$



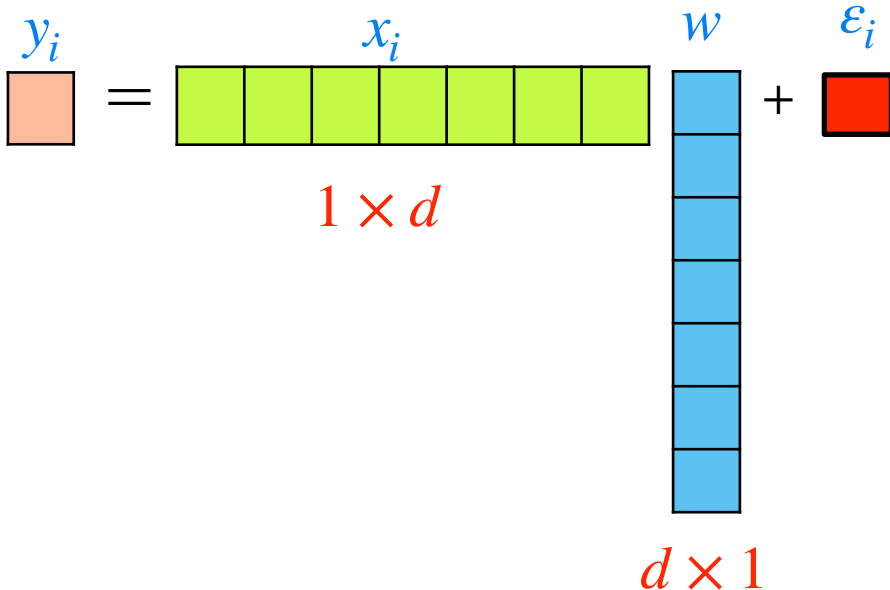
The regression problem in matrix notation, translating

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \quad \begin{array}{l} d \times 1 \\ d \times 1 \\ d \times 1 \end{array} \quad \begin{array}{l} d : \# \text{ of features} \\ n : \# \text{ of examples/datapoints} \end{array}$$

$x \in \mathbb{R}^d$

Previously... $y_i = x_i^T w + \epsilon_i$



The regression problem in matrix notation, translated

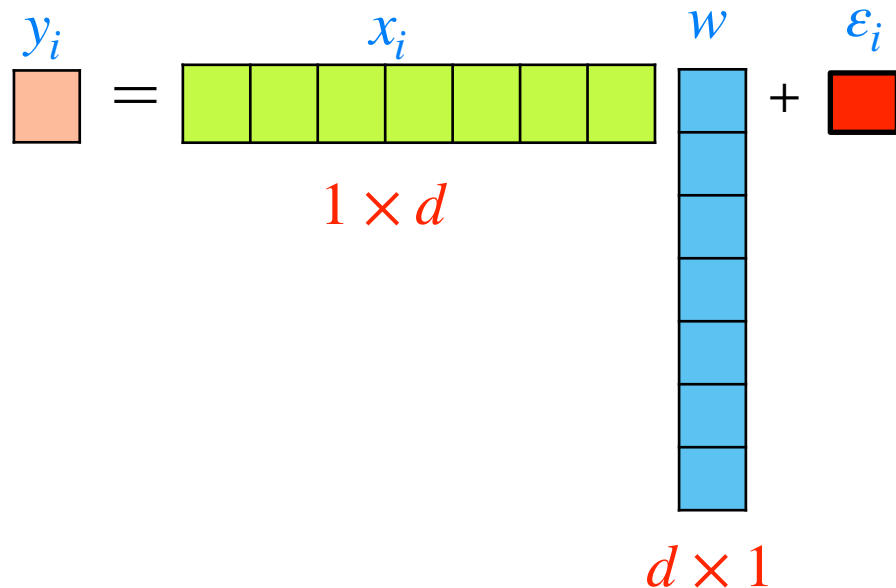
$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

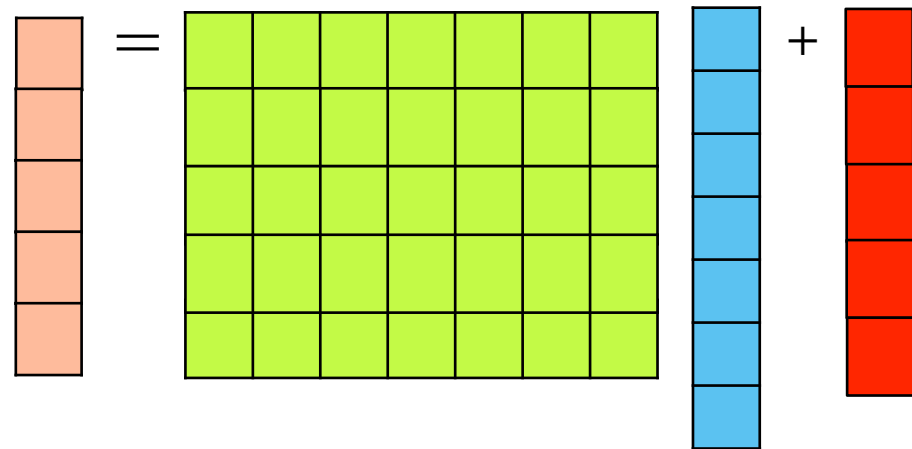
d : # of features

n : # of examples/datapoints

Previously... $y_i = x_i^\top w + \epsilon_i$



Now: $y = \mathbf{X}w + \epsilon$



The regression problem in matrix notation, translated

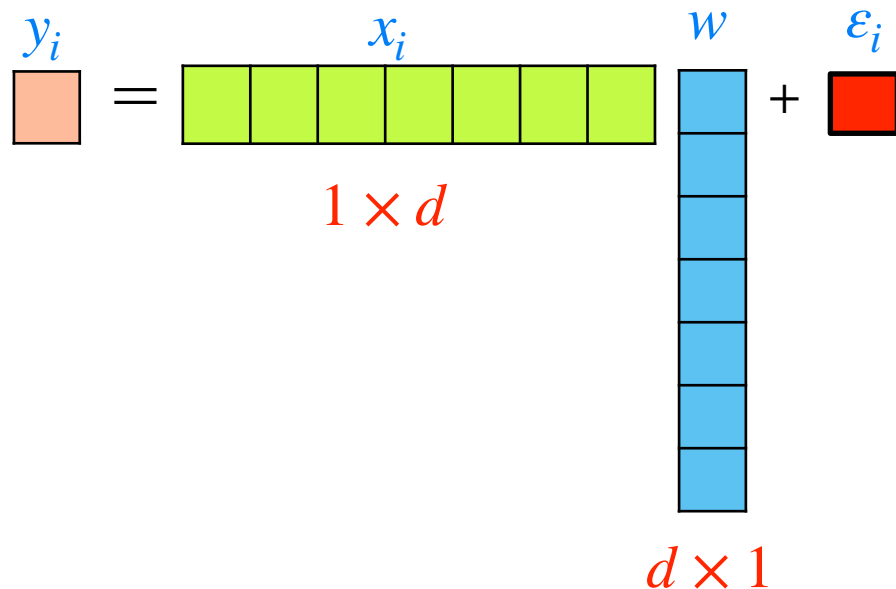
$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

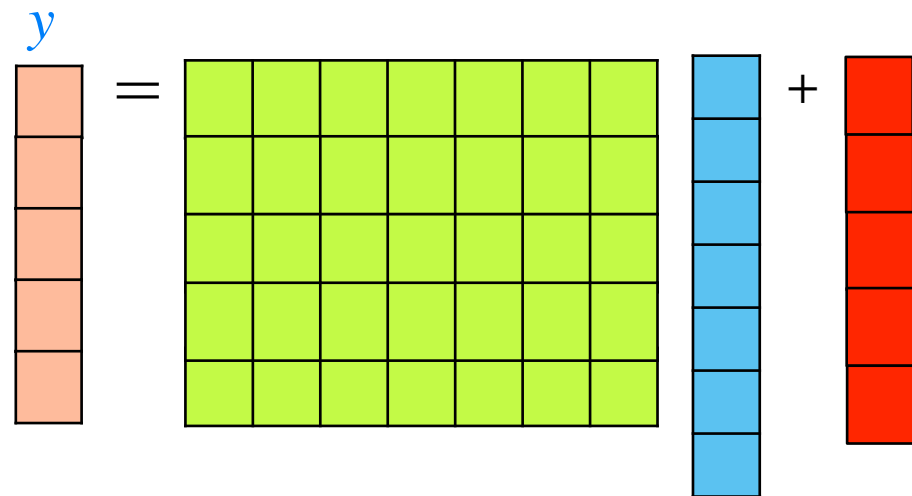
d : # of features

n : # of examples/datapoints

Previously... $y_i = x_i^\top w + \epsilon_i$



Now: $y = \mathbf{X}w + \epsilon$



The regression problem in matrix notation, translated

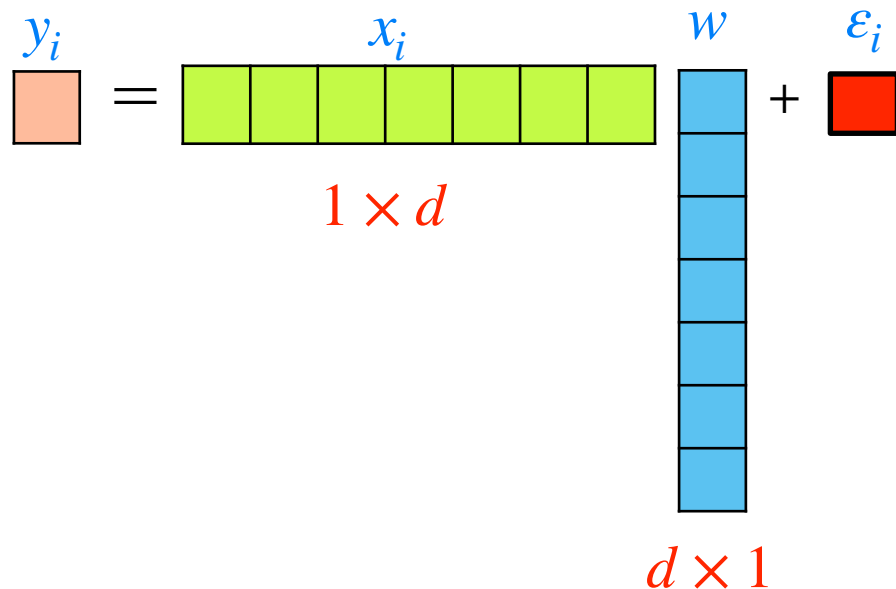
$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

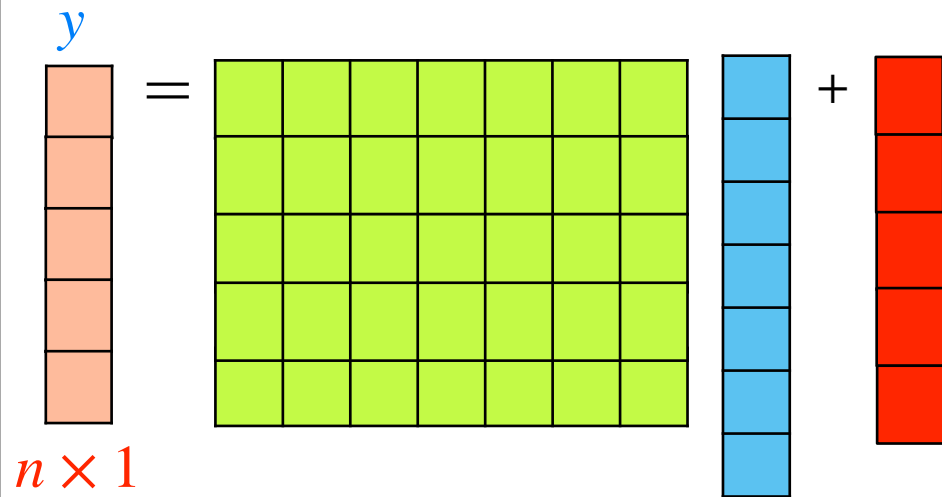
d : # of features

n : # of examples/datapoints

Previously... $y_i = x_i^\top w + \epsilon_i$



Now: $y = \mathbf{X}w + \epsilon$



The regression problem in matrix notation, translated

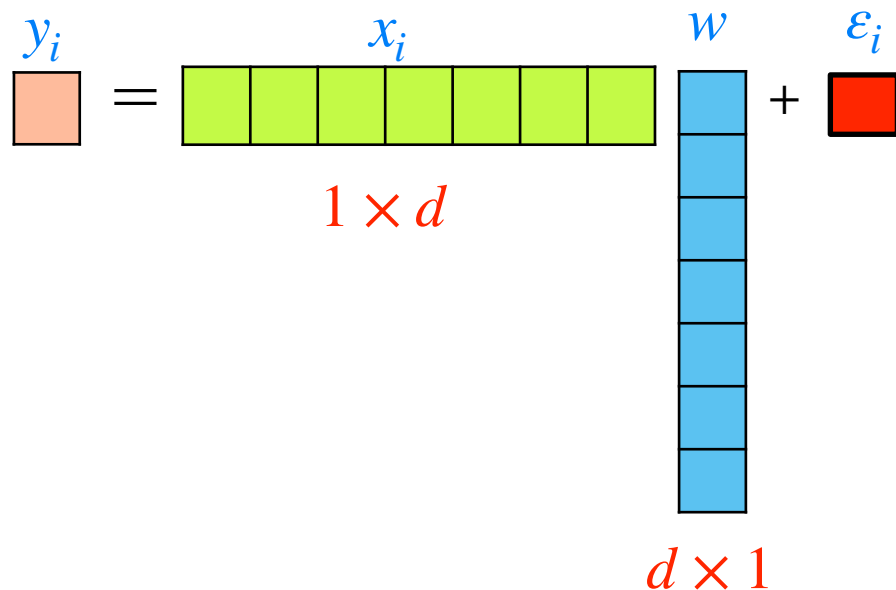
$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

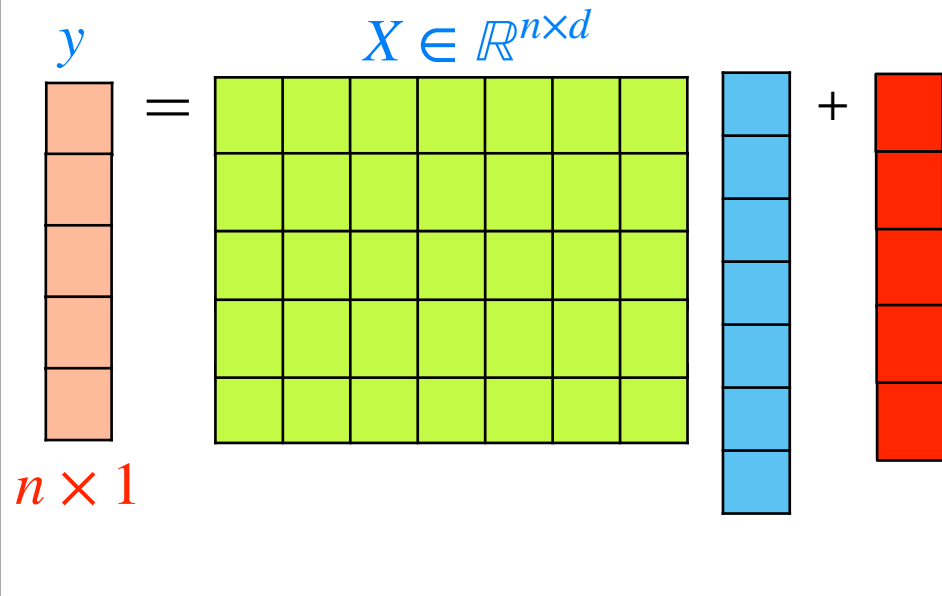
d : # of features

n : # of examples/datapoints

Previously... $y_i = x_i^\top w + \epsilon_i$



Now: $y = \mathbf{X}w + \epsilon$



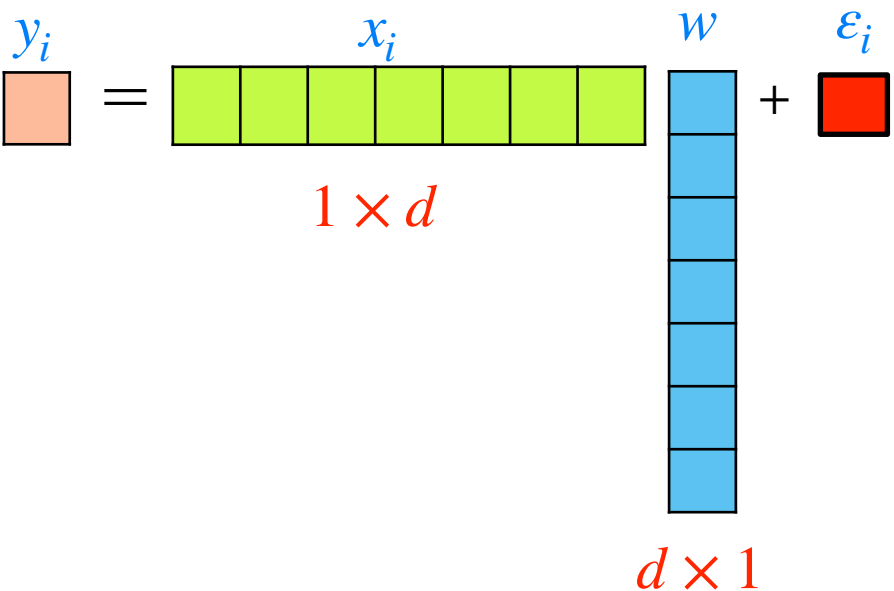
The regression problem in matrix notation, translated

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2$$

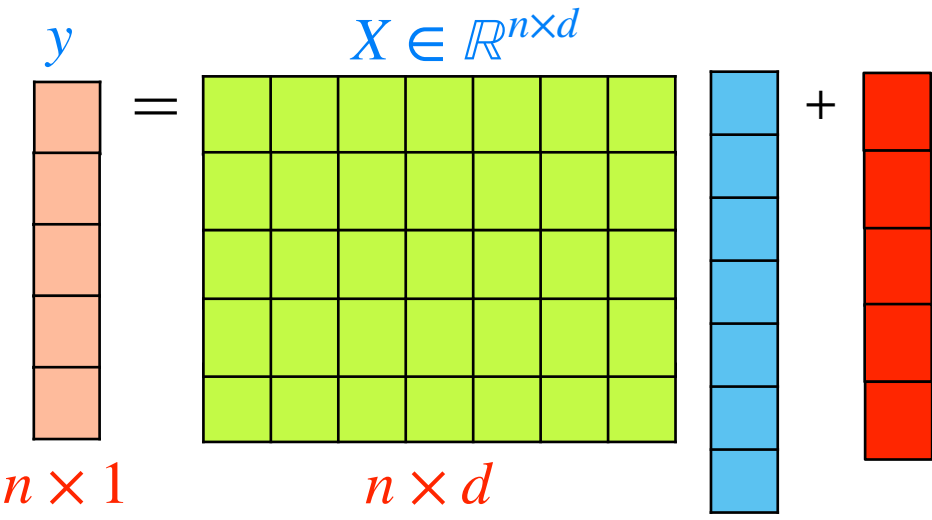
$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

d : # of features
n : # of examples/datapoints

Previously... $y_i = x_i^T w + \epsilon_i$



Now: $y = Xw + \epsilon$



The regression problem in matrix notation, translated

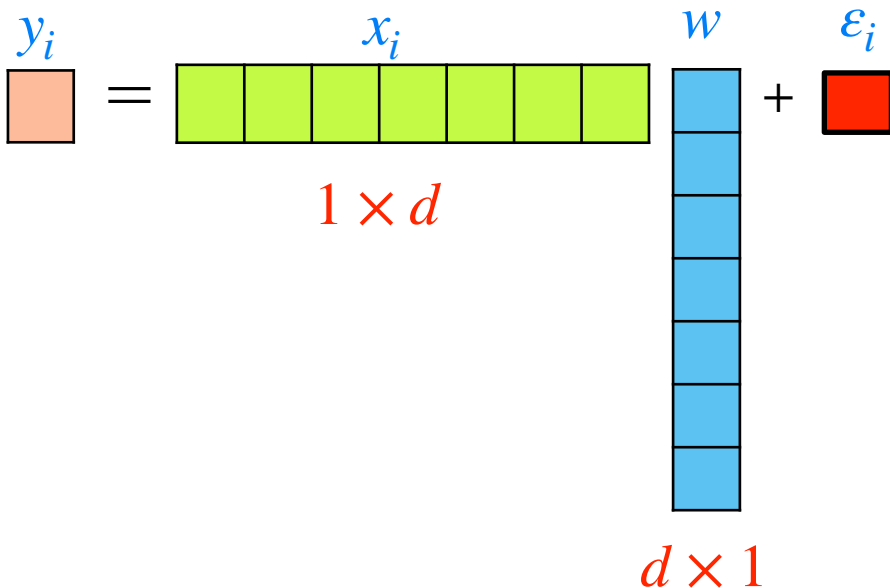
$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

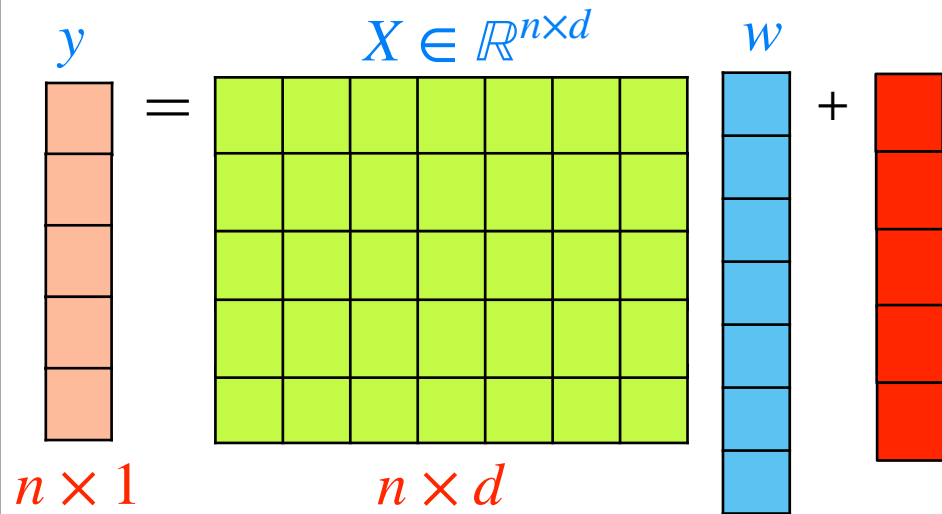
d : # of features

n : # of examples/datapoints

Previously... $y_i = x_i^\top w + \epsilon_i$



Now: $y = \mathbf{X}w + \epsilon$



The regression problem in matrix notation, translated

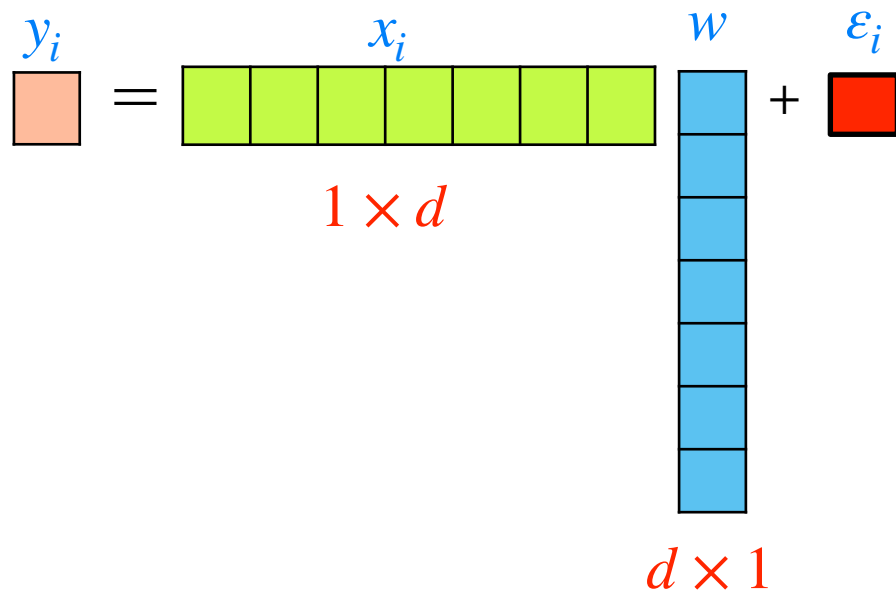
$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

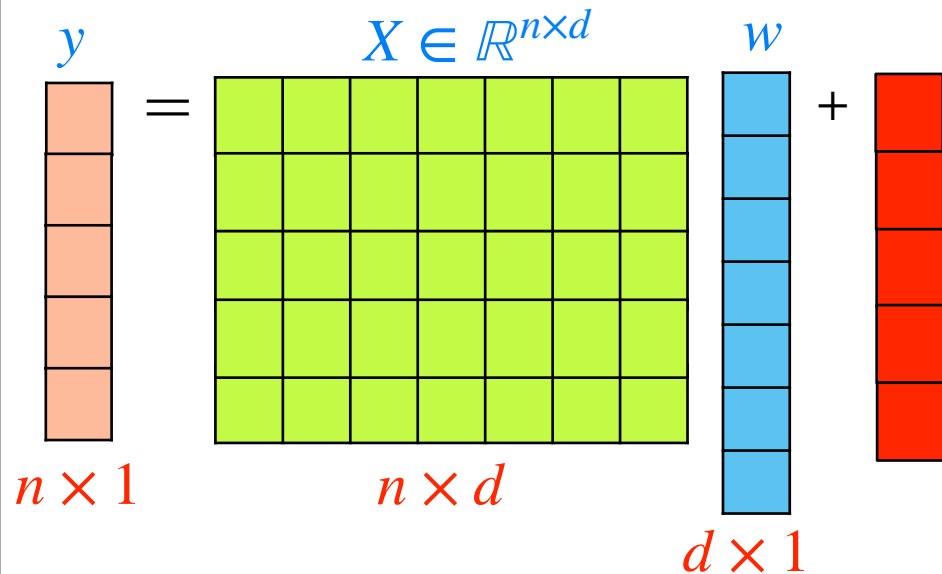
d : # of features

n : # of examples/datapoints

Previously... $y_i = x_i^\top w + \epsilon_i$



Now: $y = \mathbf{X}w + \epsilon$



The regression problem in matrix notation, translated

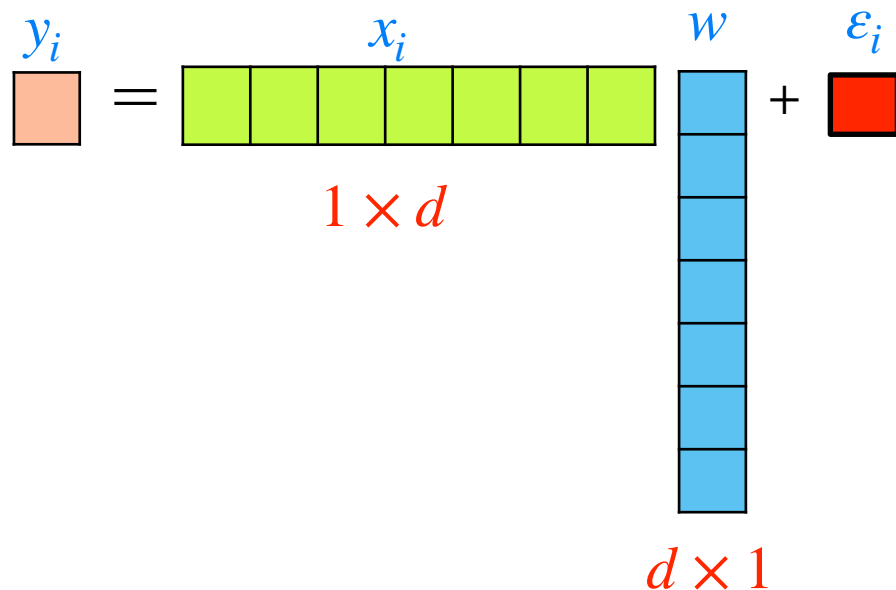
$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

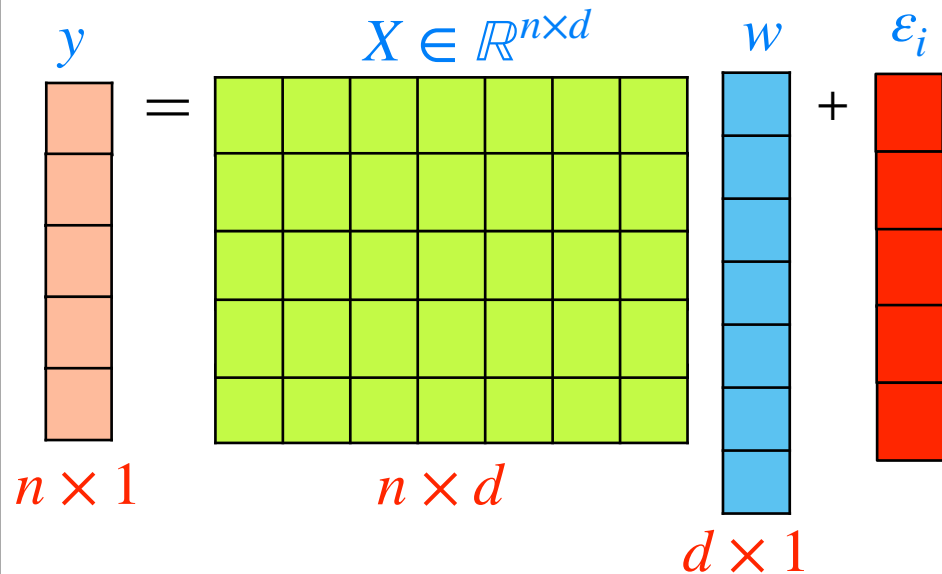
d : # of features

n : # of examples/datapoints

Previously... $y_i = x_i^\top w + \epsilon_i$



Now: $y = \mathbf{X}w + \epsilon$



The regression problem in matrix notation, translated

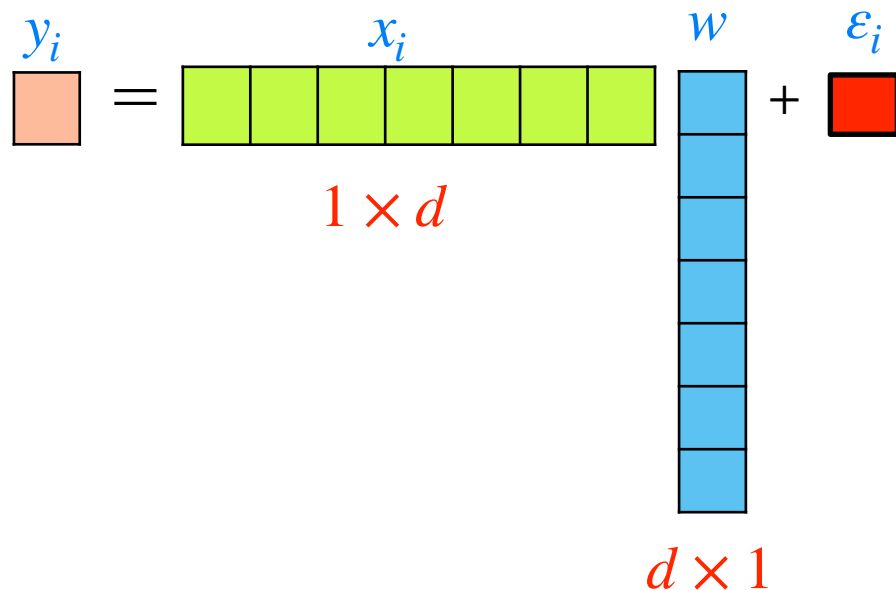
$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

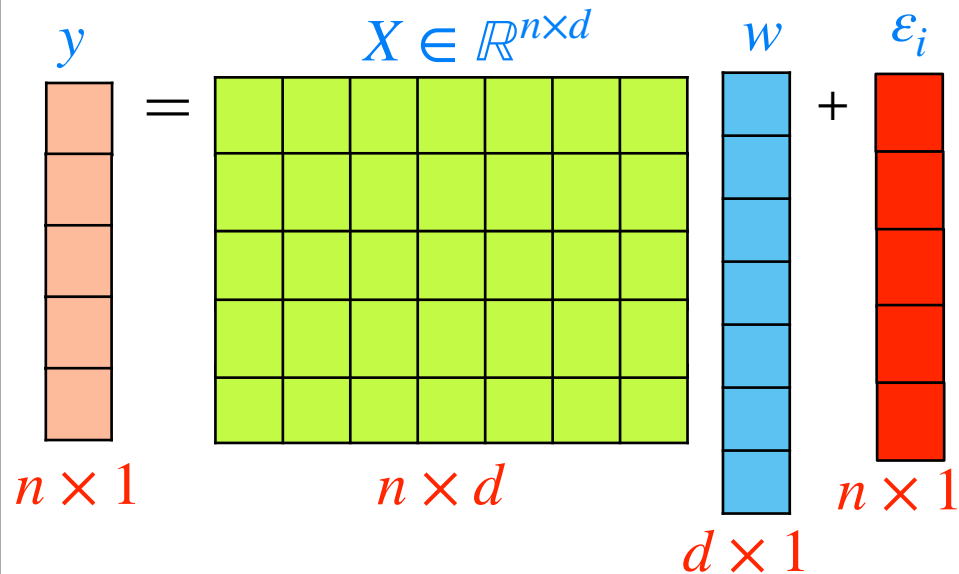
d : # of features

n : # of examples/datapoints

Previously... $y_i = x_i^\top w + \epsilon_i$



Now: $y = \mathbf{X}w + \epsilon$



Least squares in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

d : # of features

n : # of examples/datapoints

Now: $\mathbf{y} = \mathbf{X}w + \epsilon$

Least squares in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2 \quad \# \text{ How to make matrix-y?}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

d : # of features

n : # of examples/datapoints

$$\# \text{ Now: } \mathbf{y} = \mathbf{X}w + \epsilon$$

Least squares in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2 \quad \# \text{ How to make matrix-y?}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

d : # of features

n : # of examples/datapoints

$$\# \text{ Now: } \mathbf{y} = \mathbf{X}w + \epsilon$$

$$\# \text{ More identities: } \ell_2 \text{ norm: } \|z\|_2 = \sqrt{\sum_{i=1}^n z_i^2} = \sqrt{z^\top z}$$

Least squares in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2 \quad \# \text{ How to make matrix-y?}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

d : # of features

n : # of examples/datapoints

$$\# \text{ Now: } \mathbf{y} = \mathbf{X}w + \epsilon$$

$$\# \text{ More identities: } \ell_2 \text{ norm: } \|z\|_2 = \sqrt{\sum_{i=1}^n z_i^2} = \sqrt{z^\top z}$$

$$\text{Let } z_i = (y_i - x_i^\top w)$$

Least squares in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2 \quad \# \text{ How to make matrix-y?}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

d : # of features

n : # of examples/datapoints

$$\# \text{ Now: } \mathbf{y} = \mathbf{X}w + \epsilon$$

$$\# \text{ More identities: } \ell_2 \text{ norm: } \|z\|_2 = \sqrt{\sum_{i=1}^n z_i^2} = \sqrt{z^\top z}$$

$$\text{Let } z_i = (y_i - x_i^\top w) \quad \text{Then}$$

Least squares in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2 \quad \# \text{ How to make matrix-y?}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}$$

d : # of features

n : # of examples/datapoints

$$\# \text{ Now: } \mathbf{y} = \mathbf{X}w + \epsilon$$

$$\# \text{ More identities: } \ell_2 \text{ norm: } \|z\|_2 = \sqrt{\sum_{i=1}^n z_i^2} = \sqrt{z^\top z}$$

$$\text{Let } z_i = (y_i - x_i^\top w) \quad \text{Then}$$

$$\hat{w}_{LS} = \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2$$

Least squares in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2 \quad \# \text{ How to make matrix-y?}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

d : # of features

n : # of examples/datapoints

$$\# \text{ Now: } \mathbf{y} = \mathbf{X}w + \epsilon$$

$$\# \text{ More identities: } \ell_2 \text{ norm: } \|z\|_2 = \sqrt{\sum_{i=1}^n z_i^2} = \sqrt{z^\top z}$$

$$\text{Let } z_i = (y_i - x_i^\top w) \quad \text{Then}$$

$$\begin{aligned} \hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= \arg \min_w (\mathbf{y} - \mathbf{X}w)^\top (\mathbf{y} - \mathbf{X}w) \end{aligned}$$

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw) \quad \# \text{ More identities:}$$

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w y^T y - y^T Xw - w^T X^T y + w^T X^T Xw$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w y^T y - y^T Xw - w^T X^T y + w^T X^T Xw$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - y^T Xw - w^T X^T y + w^T X^T Xw$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - y^T Xw - w^T X^T y + w^T X^T Xw$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underline{y^T Xw} - w^T X^T y + w^T X^T Xw$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underline{y^T Xw}_{1 \times n} - w^T X^T y + w^T X^T Xw$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underline{y^T Xw}_{1 \times n \quad n \times d} - w^T X^T y + w^T X^T Xw$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underline{y^T Xw} - w^T X^T y + w^T X^T Xw$$

$1 \times n \quad n \times d \quad d \times 1$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T X}_{1 \times n} \underbrace{w}_{n \times d} - w^T X^T y + \underbrace{w^T X^T X}_{1 \times 1} w$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T X}_{1 \times n} \underbrace{w}_{n \times d} - \underbrace{w^T X^T}_{d \times 1} y + \underbrace{w^T X^T X w}_{1 \times 1}$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T X}_{1 \times n} \underbrace{w}_{n \times d} - \underbrace{w^T X^T}_{1 \times d} y + \underbrace{w^T X^T X w}_{1 \times 1}$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T X}_{1 \times n \ n \times d} w - \underbrace{w^T X^T}_{1 \times d \ d \times n} y + \underbrace{w^T X^T X w}_{1 \times 1}$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T X}_{1 \times n \ n \times d} w - \underbrace{w^T X^T}_{1 \times d \ d \times n} y + \underbrace{w^T X^T X w}_{1 \times 1}$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T X}_{1 \times n \quad n \times d} w - \underbrace{w^T X^T}_{1 \times d \quad d \times n} y + \underbrace{w^T X^T X}_{1 \times 1} w$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T X}_{1 \times n \quad n \times d} w - \underbrace{w^T X^T}_{1 \times d \quad d \times n} y + \underbrace{w^T X^T X}_{1 \times 1} w$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

$$s \in \mathbb{R} \rightarrow s = s^T$$

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T X}_{1 \times n \ n \times d \ d \times 1} w - \underbrace{(w^T X^T y)}_{1 \times d \ d \times n \ n \times 1} + \underbrace{w^T X^T X w}_{1 \times 1 \quad 1 \times 1}$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

$$s \in \mathbb{R} \rightarrow s = s^T$$

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T X}_{1 \times n \quad n \times d} w - \underbrace{(w^T X^T y)}_{1 \times d \quad d \times n \quad n \times 1} + \underbrace{w^T X^T X w}_{1 \times 1 \quad 1 \times 1}$$

$$= \arg \min_w - y^T Xw - y^T Xw + w^T X^T Xw$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

$$s \in \mathbb{R} \rightarrow s = s^T$$

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T Xw}_{1 \times n \quad n \times d \quad d \times 1} - \underbrace{(w^T X^T y)^T}_{1 \times d \quad d \times n \quad n \times 1} + w^T X^T Xw$$

1×1
 1×1

$$= \arg \min_w -y^T Xw - y^T Xw + w^T X^T Xw$$

$$= \arg \min_w -2y^T Xw + w^T X^T Xw$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

$$s \in \mathbb{R} \rightarrow s = s^T$$

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T Xw}_{1 \times n \quad n \times d \quad d \times 1} - \underbrace{(w^T X^T y)^T}_{1 \times d \quad d \times n \quad n \times 1} + w^T X^T Xw$$

1×1
 1×1

$$= \arg \min_w -y^T Xw - y^T Xw + w^T X^T Xw$$

$$= \arg \min_w -2y^T Xw + w^T X^T Xw$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

$$s \in \mathbb{R} \rightarrow s = s^T$$

Pretty simple!

Solving least squares in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

As before, first simplify the function to make it easier to work with

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T Xw}_{1 \times n \quad n \times d \quad d \times 1} - \underbrace{(w^T X^T y)^T}_{1 \times d \quad d \times n \quad n \times 1} + w^T X^T Xw$$

1×1
 1×1

$$= \arg \min_w -y^T Xw - y^T Xw + w^T X^T Xw$$

$$= \arg \min_w -2y^T Xw + w^T X^T Xw$$

More identities:

$$(ABC)^T = C^T B^T A^T$$

Minimizing wrt w , so?

Check dimensions

$$s \in \mathbb{R} \rightarrow s = s^T$$

Pretty simple!

So now what?

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

Take derivative, set to 0!

$$\nabla_w [-2\mathbf{y}^T \mathbf{X}w + w^T \mathbf{X}^T \mathbf{X}w] = 0$$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

Take derivative, set to 0!

$$\nabla_w [-2\mathbf{y}^T \mathbf{X}w + w^T \mathbf{X}^T \mathbf{X}w] = 0$$

Useful matrix gradient identities

(A) $\nabla_w x^T w = x$

(B) $\nabla_w x^T A w = A^T x$

(C) $\nabla_w w^T A w = (A + A^T)w$

Quadratic form

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [-2\mathbf{y}^T \mathbf{X}w + w^T \mathbf{X}^T \mathbf{X}w] = 0$$

Take derivative, set to 0!

Useful matrix gradient identities

(A) $\nabla_w x^T w = x$

(B) $\nabla_w x^T A w = A^T x$

(C) $\nabla_w w^T A w = (A + A^T)w$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [\underline{-2y^T Xw + w^T X^T Xw}] = 0$$

Take derivative, set to 0!

Useful matrix gradient identities

$$(A) \quad \nabla_w x^T w = x$$

$$(B) \quad \nabla_w x^T A w = A^T x$$

$$(C) \quad \nabla_w w^T A w = (A + A^T)w$$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [\underline{-2y^T Xw + w^T X^T Xw}] = 0$$

(B)

Take derivative, set to 0!

Useful matrix gradient identities

$$(A) \quad \nabla_w x^T w = x$$

$$(B) \quad \nabla_w x^T A w = A^T x$$

$$(C) \quad \nabla_w w^T A w = (A + A^T) w$$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [\underbrace{-2\mathbf{y}^T \mathbf{X}w}_{(B)} + \underbrace{w^T \mathbf{X}^T \mathbf{X}w}] = 0$$

Take derivative, set to 0!

Useful matrix gradient identities

$$(A) \quad \nabla_w x^T w = x$$

$$(B) \quad \nabla_w x^T A w = A^T x$$

$$(C) \quad \nabla_w w^T A w = (A + A^T)w$$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [\underbrace{-2\mathbf{y}^T \mathbf{X}w}_{(B)} + \underbrace{w^T \mathbf{X}^T \mathbf{X}w}] = 0$$

Take derivative, set to 0!

Useful matrix gradient identities

$$(A) \quad \nabla_w x^T w = x$$

$$(B) \quad \nabla_w x^T A w = A^T x$$

$$(C) \quad \nabla_w w^T A w = (A + A^T)w$$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Symmetric: $(X^T X)^T = X^T X$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [\underbrace{-2\mathbf{y}^T \mathbf{X}w}_{(B)} + \underbrace{w^T \mathbf{X}^T \mathbf{X}w}_{(C)}] = 0$$

Take derivative, set to 0!

Useful matrix gradient identities

$$(A) \quad \nabla_w x^T w = x$$

$$(B) \quad \nabla_w x^T A w = A^T x$$

$$(C) \quad \nabla_w w^T A w = (A + A^T)w$$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Symmetric: $(X^T X)^T = X^T X$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [\underbrace{-2\mathbf{y}^T \mathbf{X}w}_{(B)} + \underbrace{w^T \mathbf{X}^T \mathbf{X}w}_{(C)}] = 0$$

$$= -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}w = 0$$

Take derivative, set to 0!

Useful matrix gradient identities

$$(A) \quad \nabla_w x^T w = x$$

$$(B) \quad \nabla_w x^T A w = A^T x$$

$$(C) \quad \nabla_w w^T A w = (A + A^T)w$$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Symmetric: $(\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X}$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [\underbrace{-2\mathbf{y}^T \mathbf{X}w}_{(B)} + \underbrace{w^T \mathbf{X}^T \mathbf{X}w}_{(C)}] = 0$$

$$= -\cancel{2\mathbf{X}^T \mathbf{y}} + 2\mathbf{X}^T \mathbf{X}w = 0$$

Take derivative, set to 0!

Useful matrix gradient identities

$$(A) \quad \nabla_w x^T w = x$$

$$(B) \quad \nabla_w x^T A w = A^T x$$

$$(C) \quad \nabla_w w^T A w = (A + A^T)w$$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Symmetric: $(\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X}$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [\underbrace{-2\mathbf{y}^T \mathbf{X}w}_{(B)} + \underbrace{w^T \mathbf{X}^T \mathbf{X}w}_{(C)}] = 0$$

$$= -\cancel{2\mathbf{X}^T \mathbf{y}} + \cancel{2\mathbf{X}^T \mathbf{X}w} = 0$$

Take derivative, set to 0!

Useful matrix gradient identities

$$(A) \quad \nabla_w x^T w = x$$

$$(B) \quad \nabla_w x^T A w = A^T x$$

$$(C) \quad \nabla_w w^T A w = (A + A^T)w$$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Symmetric: $(\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X}$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [\underbrace{-2\mathbf{y}^T \mathbf{X}w}_{(B)} + \underbrace{w^T \mathbf{X}^T \mathbf{X}w}_{(C)}] = 0$$

$$= -\cancel{2\mathbf{X}^T \mathbf{y}} + \cancel{2\mathbf{X}^T \mathbf{X}w} = 0$$

$$\mathbf{X}^T \mathbf{X}w = \mathbf{X}^T \mathbf{y}$$

Take derivative, set to 0!

Useful matrix gradient identities

$$(A) \quad \nabla_w x^T w = x$$

$$(B) \quad \nabla_w x^T A w = A^T x$$

$$(C) \quad \nabla_w w^T A w = (A + A^T)w$$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Symmetric: $(\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X}$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [\underbrace{-2\mathbf{y}^T \mathbf{X}w}_{(B)} + \underbrace{w^T \mathbf{X}^T \mathbf{X}w}_{(C)}] = 0$$

$$= -\cancel{2\mathbf{X}^T \mathbf{y}} + \cancel{2\mathbf{X}^T \mathbf{X}w} = 0$$

$$\mathbf{X}^T \mathbf{X}w = \mathbf{X}^T \mathbf{y}$$

$$\hat{w}_{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Take derivative, set to 0!

Useful matrix gradient identities

$$(A) \quad \nabla_w x^T w = x$$

$$(B) \quad \nabla_w x^T A w = A^T x$$

$$(C) \quad \nabla_w w^T A w = (A + A^T)w$$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Symmetric: $(\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X}$

Least squares: take the gradient and set to 0

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$

$$\nabla_w [\underbrace{-2\mathbf{y}^T \mathbf{X}w}_{(B)} + \underbrace{w^T \mathbf{X}^T \mathbf{X}w}_{(C)}] = 0$$

$$= -\cancel{2\mathbf{X}^T \mathbf{y}} + \cancel{2\mathbf{X}^T \mathbf{X}w} = 0$$

$$\mathbf{X}^T \mathbf{X}w = \mathbf{X}^T \mathbf{y}$$

$$\hat{w}_{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Take derivative, set to 0!

Useful matrix gradient identities

$$(A) \quad \nabla_w x^T w = x$$

$$(B) \quad \nabla_w x^T A w = A^T x$$

$$(C) \quad \nabla_w w^T A w = (A + A^T)w$$

Quadratic form

If A is symmetric ($A = A^T$),

Then $\nabla_w w^T A w = 2Aw$

Symmetric: $(\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X}$

Least squares solution in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

d : # of features

n : # of examples/datapoints

More identities: ℓ_2 norm: $\|z\|_2 = \sqrt{\sum_{i=1}^n z_i^2} = \sqrt{z^\top z}$

$$\begin{aligned} \hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= \arg \min_w (\mathbf{y} - \mathbf{X}w)^\top (\mathbf{y} - \mathbf{X}w) \end{aligned}$$

$$\hat{w}_{LS} = \hat{w}_{MLE} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Least squares solution in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

d : # of features

n : # of examples/datapoints

More identities: ℓ_2 norm: $\|z\|_2 = \sqrt{\sum_{i=1}^n z_i^2} = \sqrt{z^\top z}$

$$\begin{aligned} \hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= \arg \min_w (\mathbf{y} - \mathbf{X}w)^\top (\mathbf{y} - \mathbf{X}w) \end{aligned}$$

$$\hat{w}_{LS} = \hat{w}_{MLE} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

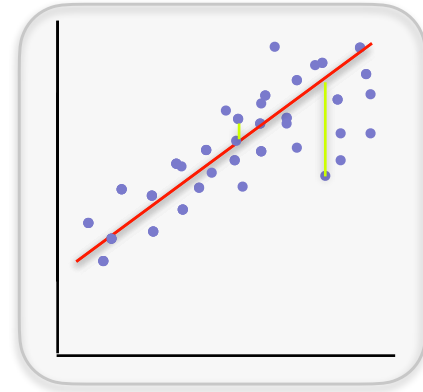
Why LS? “Ordinary Least Squares”

Regression... what about an intercept?

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

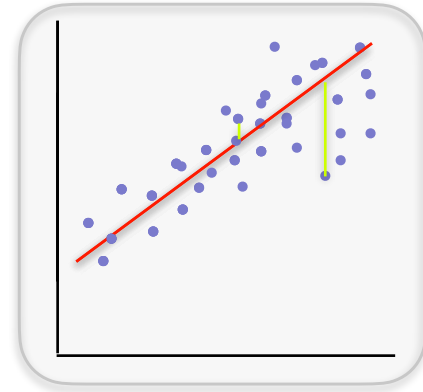
Regression... what about an intercept?

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$



Regression... what about an intercept?

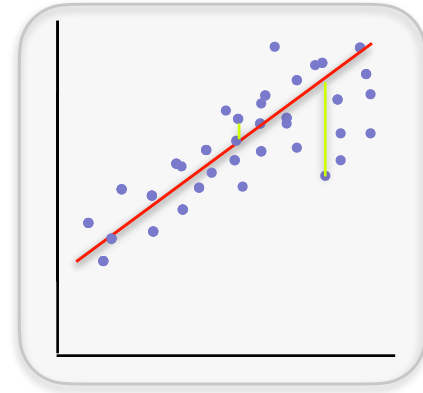
$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$



What about an offset?

Regression... what about an intercept?

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

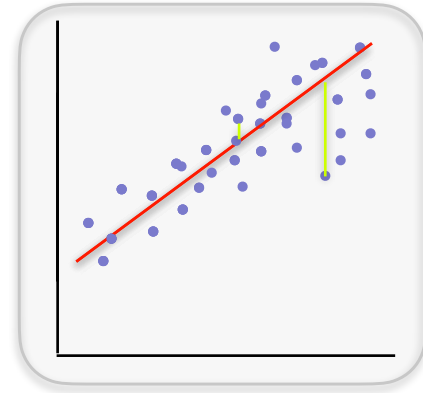


What about an offset?

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w,b} \sum_{i=1}^n (y_i - (x_i^T w + b))^2$$

Regression... what about an intercept?

$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$



What about an offset?

$$\begin{aligned}\hat{w}_{LS}, \hat{b}_{LS} &= \arg \min_{w,b} \sum_{i=1}^n (y_i - (x_i^T w + b))^2 \\ &= \arg \min_{w,b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2\end{aligned}$$

Dealing with an offset

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2$$

$$\mathbf{X}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{X}^T \mathbf{1} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{1}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{1}^T \mathbf{1} = \mathbf{1}^T \mathbf{y}$$

If $\mathbf{X}^T \mathbf{1} = 0$ (i.e., if each feature is mean-zero) then

$$\hat{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i$$

Dealing with an offset

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \|\mathbf{y} - (\mathbf{X}w + \mathbf{1}b)\|_2^2$$

$$\mathbf{X}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{X}^T \mathbf{1} = \mathbf{X}^T \mathbf{y}$$

$$\mathbf{1}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{1}^T \mathbf{1} = \mathbf{1}^T \mathbf{y}$$

If $\mathbf{X}^T \mathbf{1} = 0$,

$$\hat{w}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i$$

In general, when $\mathbf{X}^T \mathbf{1} \neq 0$,

$$\mu = \frac{1}{n} \mathbf{X}^T \mathbf{1}$$

$$\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{1}\mu^T$$

$$\hat{w}_{LS} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

$$\hat{b}_{LS} = \frac{1}{n} \mathbf{1}^T \mathbf{y} - \mu^T \hat{w}_{LS}$$

How do we make predictions from this model?

$$\hat{\mathbf{w}}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i$$

A new house is about to be listed. What should it sell for?

How do we make predictions from this model?

$$\hat{\mathbf{w}}_{LS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i$$

A new house is about to be listed. What should it sell for?

$$\hat{y}_{\text{new}} = x_{\text{new}}^T \hat{\mathbf{w}}_{LS} + \hat{b}_{LS}$$

Process

Decide on a **model** for the likelihood function $f(x; \theta)$

Find the function which fits the data best

Choose a loss function- least squares

Pick the function which minimizes loss on data

Use function to make prediction on new examples

Process

Decide on a **model**: $y_i = x_i^T w + b + \epsilon_i$

Choose a loss function - least squares, or log likelihood

Pick the function which minimizes loss on data

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \sum_{i=1}^n (y_i - (x_i^T w + b))^2$$

Use function to make prediction on new examples

$$\hat{y}_{\text{new}} = x_{\text{new}}^T \hat{w}_{LS} + \hat{b}_{LS}$$

Analysis of Error- Unbiased

$$\text{if } y_i = x_i^T w + \epsilon_i \quad \text{and} \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2) \quad \mathbf{Y} = \mathbf{X}w + \epsilon$$

$$\begin{aligned} \hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w + \epsilon) \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \end{aligned}$$

Maximum Likelihood Estimator is unbiased:

$$\begin{aligned} \mathbb{E}[\hat{w} - w | X] &= \mathbb{E}[w + (X^T X)^{-1} X^T \epsilon - w | X] && \text{Plug-in } \uparrow \\ &= \mathbb{E}[(X^T X)^{-1} X^T \epsilon | X] && w-w = 0 \\ &= (X^T X)^{-1} X^T \mathbb{E}[\epsilon | X] && \mathbb{E}[af(X) | X] = f(X)\mathbb{E}[a] \\ &= 0 && \mathbb{E}[\epsilon | X] = 0 \end{aligned}$$

Analysis of Error- Covariance

Covariance is:

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

Analysis of Error- Covariance

Covariance is:

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

$$\mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])^T (\hat{w} - \mathbb{E}[\hat{w}])]$$

Analysis of Error- Covariance

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

Covariance is:

$$\begin{aligned}\mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])^T (\hat{w} - \mathbb{E}[\hat{w}])] \\ = \mathbb{E}[(\hat{w} - w)^T (\hat{w} - w)]\end{aligned}$$

Analysis of Error- Covariance

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

Covariance is:

$$\begin{aligned}& \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])^T (\hat{w} - \mathbb{E}[\hat{w}])] \\ &= \mathbb{E}[(\hat{w} - w)^T (\hat{w} - w)] \\ &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon]^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)\end{aligned}$$

Analysis of Error- Covariance

Covariance is:

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

$$\begin{aligned}\mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])^T (\hat{w} - \mathbb{E}[\hat{w}])] \\ &= \mathbb{E}[(\hat{w} - w)^T (\hat{w} - w)] \\ &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon]^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)\end{aligned}$$

$$(\mathbf{ABC})^T = \mathbf{C}^T \mathbf{B}^T \mathbf{A}^T$$

Analysis of Error- Covariance

Covariance is:

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

$$\begin{aligned}& \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])^T (\hat{w} - \mathbb{E}[\hat{w}])] \\ &= \mathbb{E}[(\hat{w} - w)^T (\hat{w} - w)] \\ &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon]^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon) \\ &= \mathbb{E}[\epsilon^T \mathbf{X}^T ((\mathbf{X}^T \mathbf{X})^{-1})^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)]\end{aligned}$$

$$(\mathbf{ABC})^T = \mathbf{C}^T \mathbf{B}^T \mathbf{A}^T$$

Analysis of Error- Covariance

Covariance is:

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

$$\begin{aligned}& \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])^T (\hat{w} - \mathbb{E}[\hat{w}])] \\ &= \mathbb{E}[(\hat{w} - w)^T (\hat{w} - w)] \\ &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon]^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)] \\ &= \mathbb{E}[\epsilon^T \mathbf{X}^T ((\mathbf{X}^T \mathbf{X})^{-1})^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)]\end{aligned}$$

$(ABC)^T = C^T B^T A^T$
 $A^{TT} = A, (\mathbf{X}^T \mathbf{X})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1T}$

Analysis of Error- Covariance

Covariance is:

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

$$\begin{aligned}& \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])^T (\hat{w} - \mathbb{E}[\hat{w}])] \\ &= \mathbb{E}[(\hat{w} - w)^T (\hat{w} - w)] \\ &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon]^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)] \\ &= \mathbb{E}[\epsilon^T \mathbf{X}^T ((\mathbf{X}^T \mathbf{X})^{-1})^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)] \quad (ABC)^T = C^T B^T A^T \\ &= \mathbb{E}[\epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] \quad A^{TT} = A, (\mathbf{X}^T \mathbf{X})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1T}\end{aligned}$$

Analysis of Error- Covariance

Covariance is:

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

$$\mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])^T (\hat{w} - \mathbb{E}[\hat{w}])]$$

$$= \mathbb{E}[(\hat{w} - w)^T (\hat{w} - w)]$$

$$= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon]^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)]$$

$$= \mathbb{E}[\epsilon^T \mathbf{X}^T ((\mathbf{X}^T \mathbf{X})^{-1})^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)] \quad (\mathbf{ABC})^T = \mathbf{C}^T \mathbf{B}^T \mathbf{A}^T$$

$$\mathbf{A}^{TT} = \mathbf{A}, (\mathbf{X}^T \mathbf{X})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1T}$$

$$= \mathbb{E}[\epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon]$$

$$(\mathbf{X}^T \mathbf{X})^{-1} = \mathbf{X}^{-1} \mathbf{X}^{T^{-1}}$$

Analysis of Error- Covariance

Covariance is:

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

$$\begin{aligned}& \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])^T (\hat{w} - \mathbb{E}[\hat{w}])] \\ &= \mathbb{E}[(\hat{w} - w)^T (\hat{w} - w)] \\ &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon]^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)] \\ &= \mathbb{E}[\epsilon^T \mathbf{X}^T ((\mathbf{X}^T \mathbf{X})^{-1})^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)] && (ABC)^T = C^T B^T A^T \\ &= \mathbb{E}[\epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] && A^{TT} = A, (\mathbf{X}^T \mathbf{X})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1T} \\ &= \mathbb{E}[\epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] && (\mathbf{X}^T \mathbf{X})^{-1} = \mathbf{X}^{-1} \mathbf{X}^{T-1} \\ &= \mathbb{E}[\epsilon^T (\mathbf{X} \mathbf{X}^T)^{-1} \epsilon]\end{aligned}$$

Analysis of Error- Covariance

Covariance is:

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

$$\begin{aligned}& \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])^T (\hat{w} - \mathbb{E}[\hat{w}])] \\ &= \mathbb{E}[(\hat{w} - w)^T (\hat{w} - w)] \\ &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon]^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)] \\ &= \mathbb{E}[\epsilon^T \mathbf{X}^T ((\mathbf{X}^T \mathbf{X})^{-1})^T ((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon)] \quad (ABC)^T = C^T B^T A^T \\ &= \mathbb{E}[\epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] \quad A^{TT} = A, (\mathbf{X}^T \mathbf{X})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1T} \\ &= \mathbb{E}[\epsilon^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon] \quad (\mathbf{X}^T \mathbf{X})^{-1} = \mathbf{X}^{-1} \mathbf{X}^{T^{-1}} \\ &= \mathbb{E}[\epsilon^T (\mathbf{X} \mathbf{X}^T)^{-1} \epsilon] \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}\end{aligned}$$

Analysis of Error

$$\text{if } y_i = x_i^T w + \epsilon_i \quad \text{and} \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2) \quad \mathbf{Y} = \mathbf{X}w + \epsilon$$

$$\begin{aligned} \hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w + \epsilon) \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \end{aligned}$$

$$\mathbb{E}[\hat{w}_{MLE}] = w$$

$$\text{Cov}(\hat{w}_{MLE}) = \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])(\hat{w} - \mathbb{E}[\hat{w}])^T] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

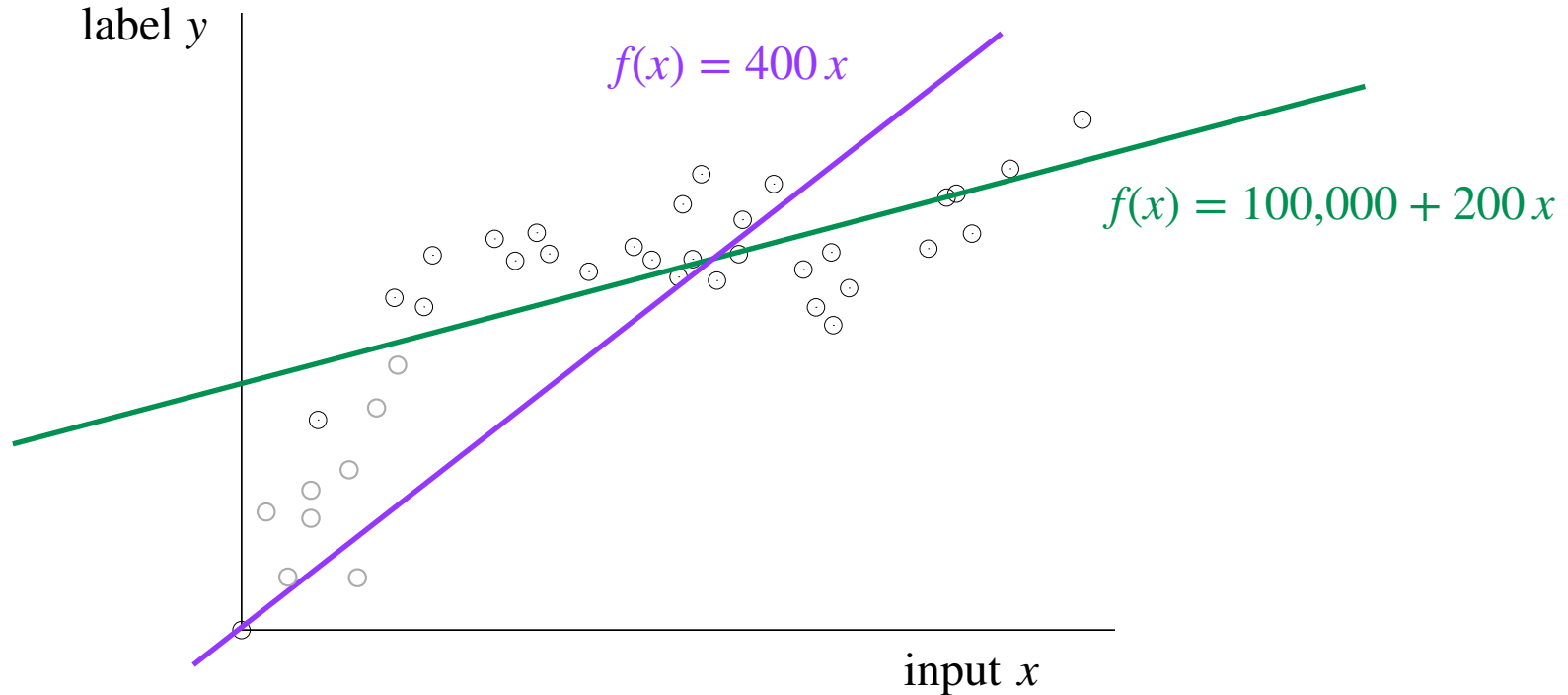
$$\hat{w}_{MLE} \sim \mathcal{N}(w, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

Linear Regression

Generalized Linear/ Polynomial



Recap: Linear Regression



- In general high-dimensions, we fit a linear model with intercept $y_i \simeq w^T x_i + b$, or equivalently $y_i = w^T x_i + b + \epsilon_i$ with model parameters $(w \in \mathbb{R}^d, b \in \mathbb{R})$ that minimizes ℓ_2 -loss

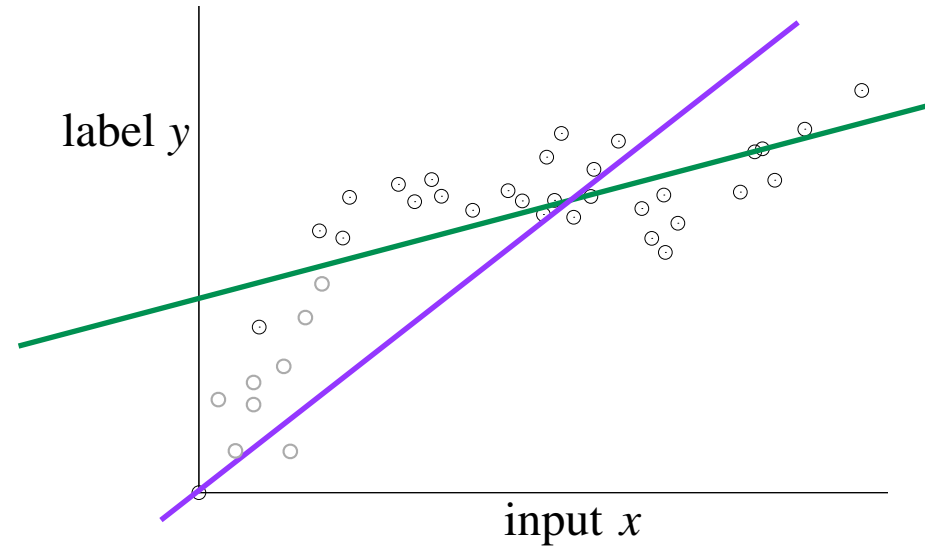
$$\mathcal{L}(w, b) = \sum_{i=1}^n \underbrace{(y_i - (w^T x_i + b))^2}_{\text{error } \epsilon_i} \quad \# \text{ MSE}$$

Quadratic regression in 1-dimension

- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter (b, w_1) :**

- $\hat{y}_i = \underline{b} + \underline{w_1 x_i}$



Quadratic regression in 1-dimension

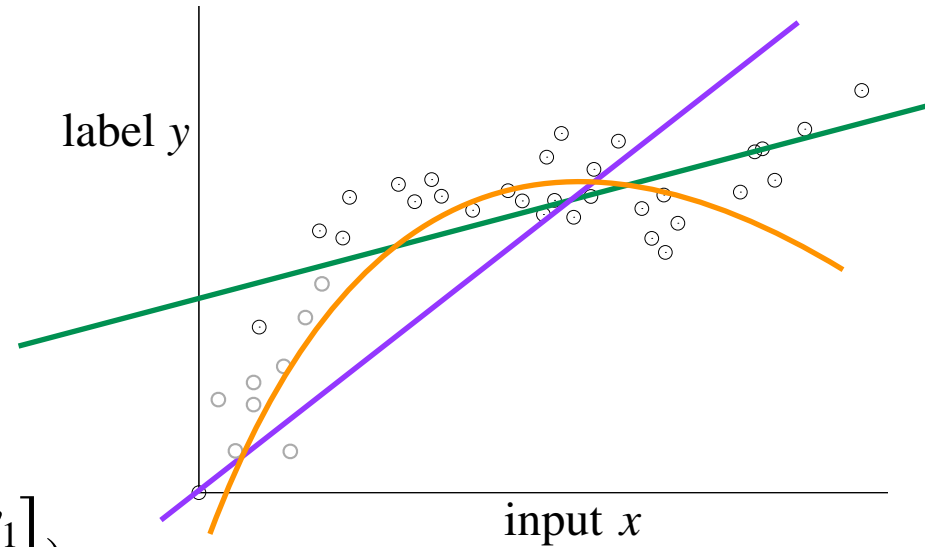
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter (b, w_1) :**

- $\hat{y}_i = \underline{b} + \underline{w_1 x_i}$

- **Quadratic model with parameter $(b, w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix})$:**

- $\hat{y}_i = \underline{b} + \underline{w_1 x_i + w_2 x_i^2}$



Quadratic regression in 1-dimension

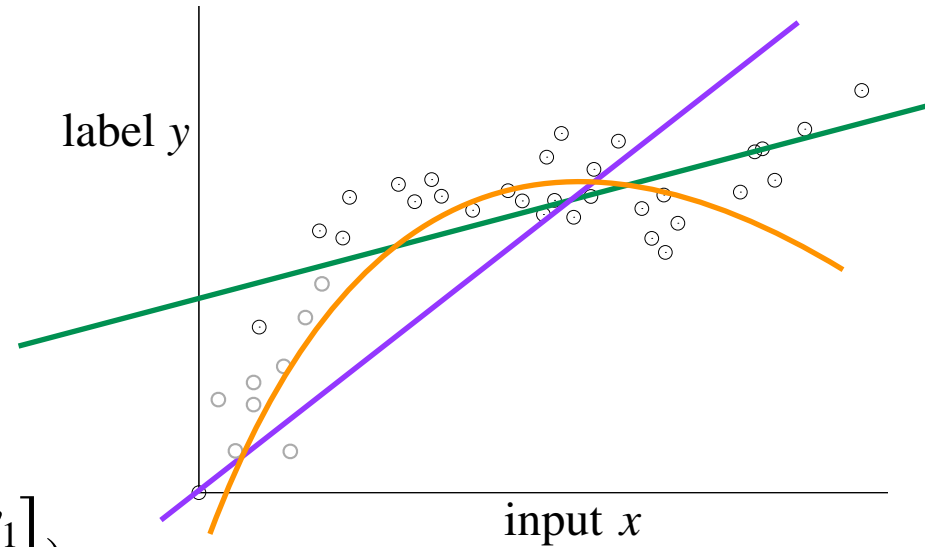
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter (b, w_1) :**

- $\hat{y}_i = \underline{b} + \underline{w_1 x_i}$

- **Quadratic model with parameter $(b, w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix})$:**

- $\hat{y}_i = \underline{b} + \underline{w_1 x_i} + \underline{w_2 x_i^2}$



$$h(x_i) = \begin{bmatrix} 1 \\ x_i \\ x_i^2 \end{bmatrix}$$

Quadratic regression in 1-dimension

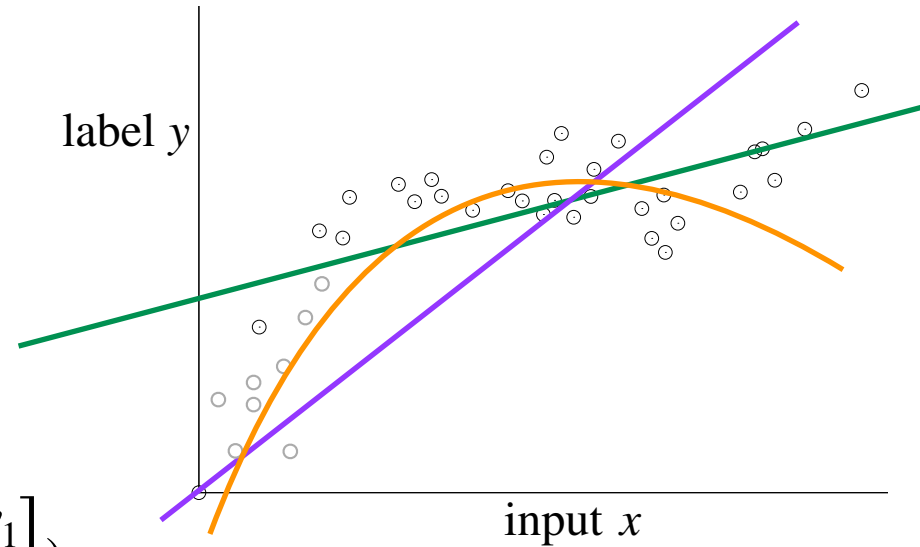
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter (b, w_1) :**

- $\hat{y}_i = b + w_1 x_i$

- **Quadratic model with parameter $(b, w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix})$:**

- $\hat{y}_i = b + w_1 x_i + w_2 x_i^2$



$$h(x_i) = \begin{bmatrix} 1 \\ x_i \\ x_i^2 \end{bmatrix} \quad \hat{y} = h(x_i)^T \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \end{bmatrix}$$

Quadratic regression in 1-dimension

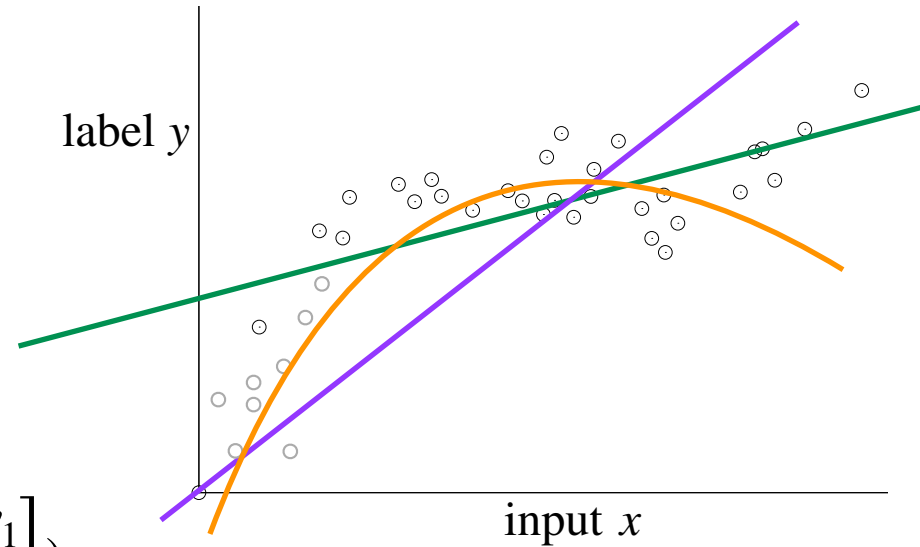
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter (b, w_1) :**

- $\hat{y}_i = b + w_1 x_i$

- **Quadratic model with parameter $(b, w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix})$:**

- $\hat{y}_i = b + w_1 x_i + w_2 x_i^2$



$$h(x_i) = \begin{bmatrix} 1 \\ x_i \\ x_i^2 \end{bmatrix} \quad \hat{y} = h(x_i)^T \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \end{bmatrix}$$

Still linear regression, but on quadratic features

Quadratic regression in 1-dimension

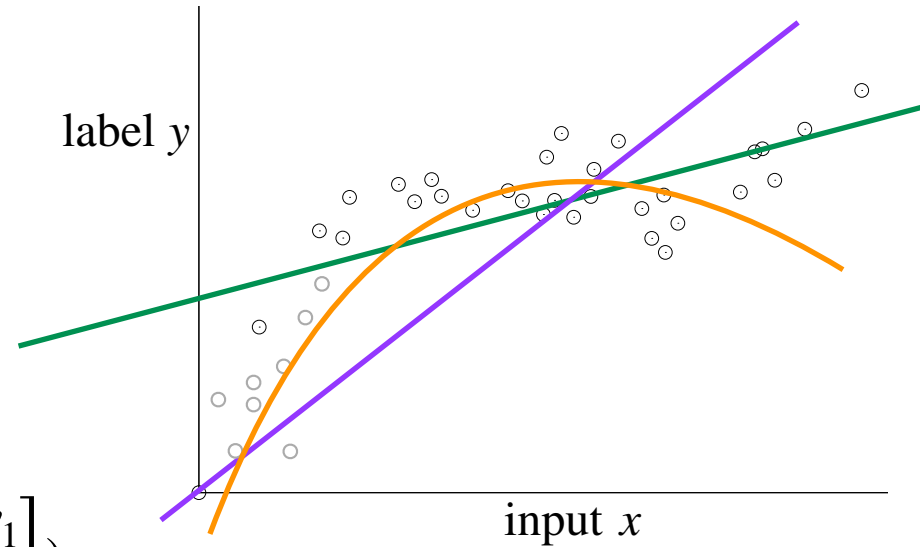
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter (b, w_1) :**

- $\hat{y}_i = b + w_1 x_i$

- **Quadratic model with parameter $(b, w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix})$:**

- $\hat{y}_i = b + w_1 x_i + w_2 x_i^2$



$$h(x_i) = \begin{bmatrix} 1 \\ x_i \\ x_i^2 \end{bmatrix} \quad \hat{y} = h(x_i)^T \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \end{bmatrix}$$

Still linear regression, but on quadratic features

Linear layer fit to non-linear features

Quadratic regression in 1-dimension

- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter (b, w_1) :**

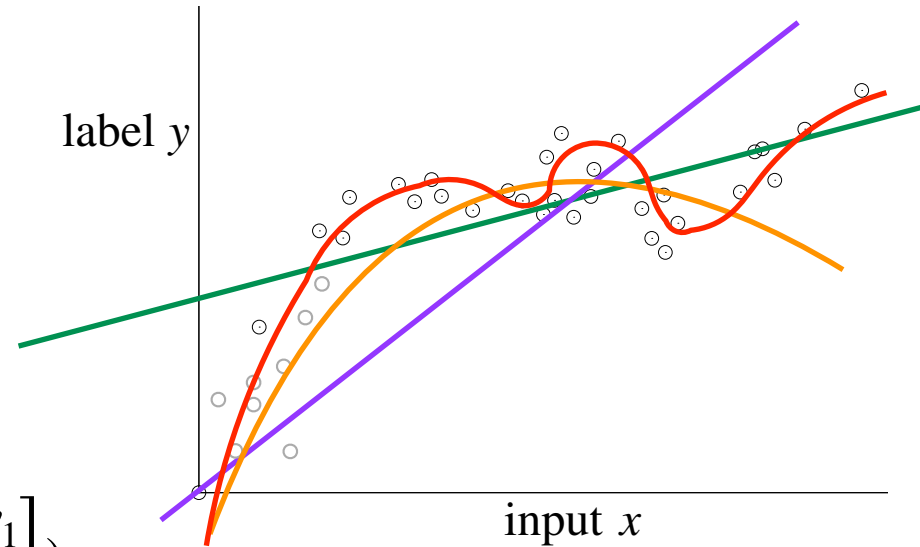
- $\hat{y}_i = \underline{b} + \underline{w_1 x_i}$

- **Quadratic model with parameter $(b, w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix})$:**

- $\hat{y}_i = \underline{b} + \underline{w_1 x_i} + \underline{w_2 x_i^2}$

- **Degree-p polynomial model with p parameters**

- $\hat{y}_i = \underline{b} + \underline{w_1 x_i} + \underline{w_2 x_i^2} + \dots + \underline{w_p x_i^p}$



$$h(x_i) = \begin{bmatrix} 1 \\ x_i \\ x_i^2 \end{bmatrix} \quad \hat{y} = h(x_i)^T \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix}$$

Still linear regression, but on quadratic features

Linear layer fit to non-linear features

Quadratic regression in 1-dimension

- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **Linear model with parameter (b, w_1) :**

- $\hat{y}_i = b + w_1 x_i$

- **Quadratic model with parameter $(b, w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix})$:**

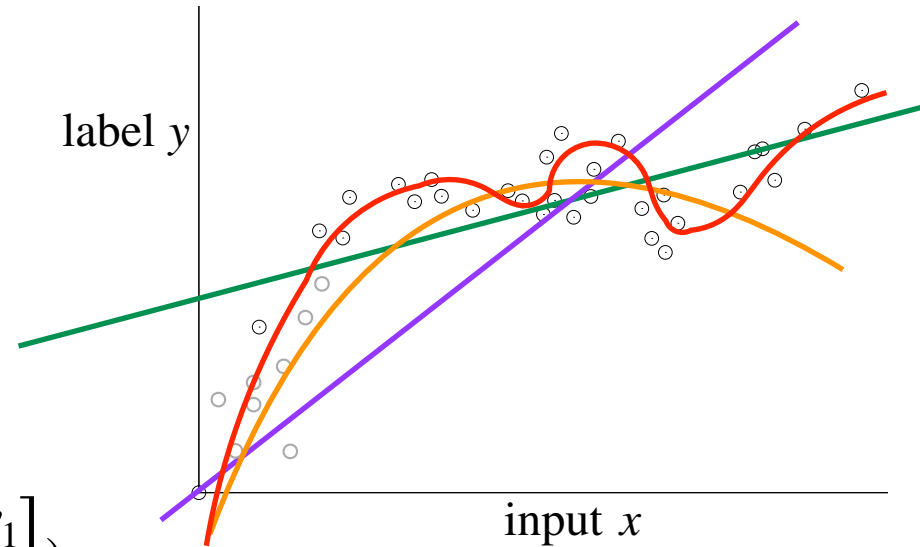
- $\hat{y}_i = b + w_1 x_i + w_2 x_i^2$

- **Degree-p polynomial model with p parameters**

- $\hat{y}_i = b + w_1 x_i + w_2 x_i^2 + \dots + w_p x_i^p$

- **General p-features with parameter $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$:**

- $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$

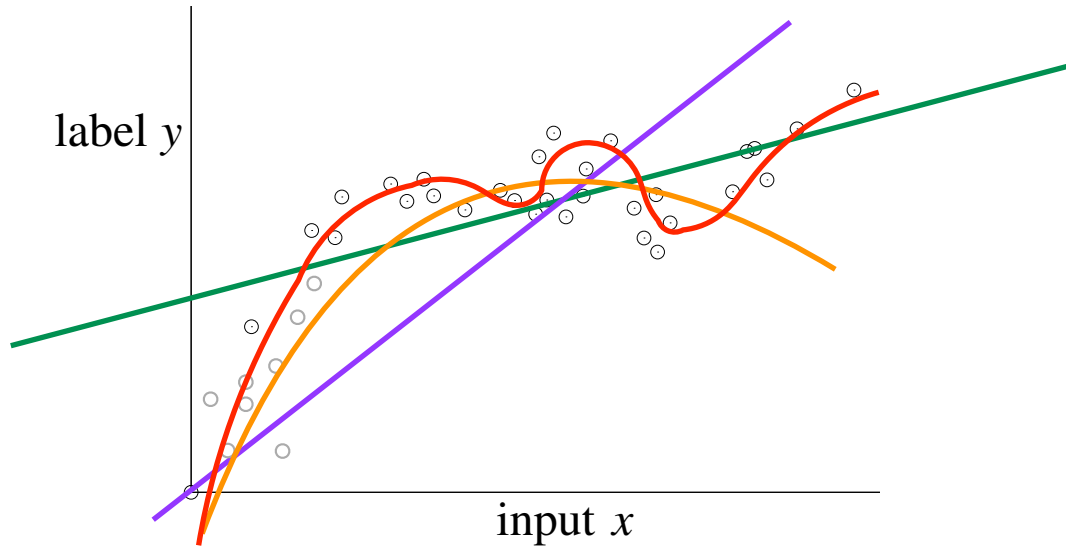


$$h(x_i) = \begin{bmatrix} 1 \\ x_i \\ x_i^2 \end{bmatrix} \quad \hat{y} = h(x_i)^T \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix}$$

Still linear regression, but on quadratic features

Linear layer fit to non-linear features

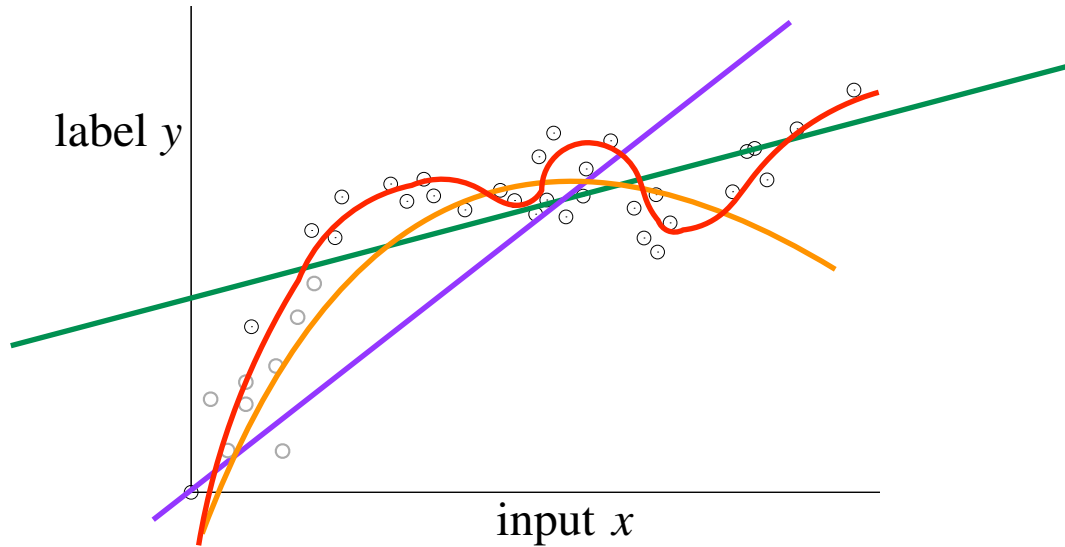
Quadratic regression in 1-dimension



Which line gives us the best fit?

- **Green** & **purple**?
- **Orange**?
- **Red**?

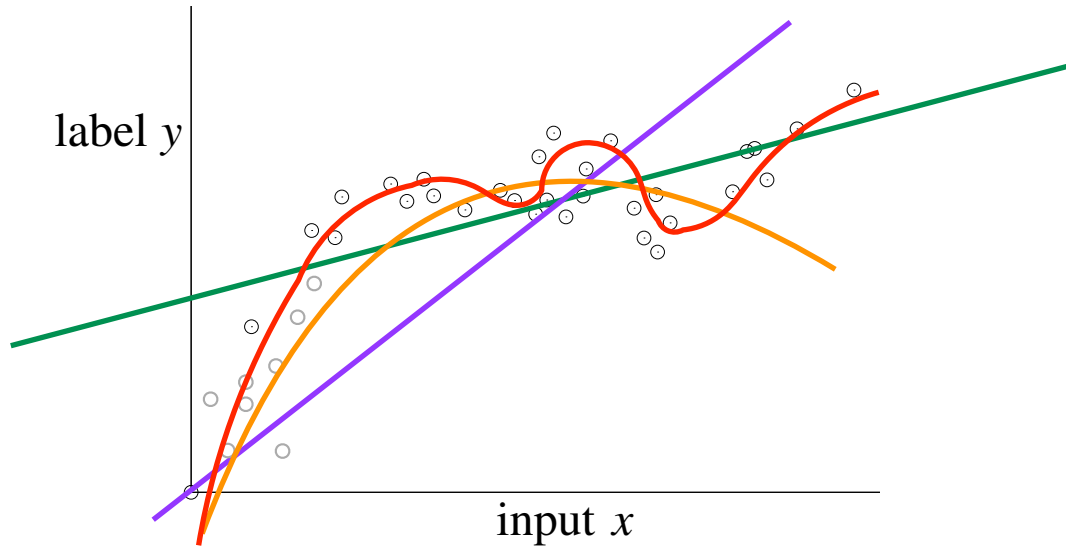
Quadratic regression in 1-dimension



Which line gives us the best fit?

- **Green & purple?** Can't capture variance. **Underfitting**
- **Orange?**
- **Red?**

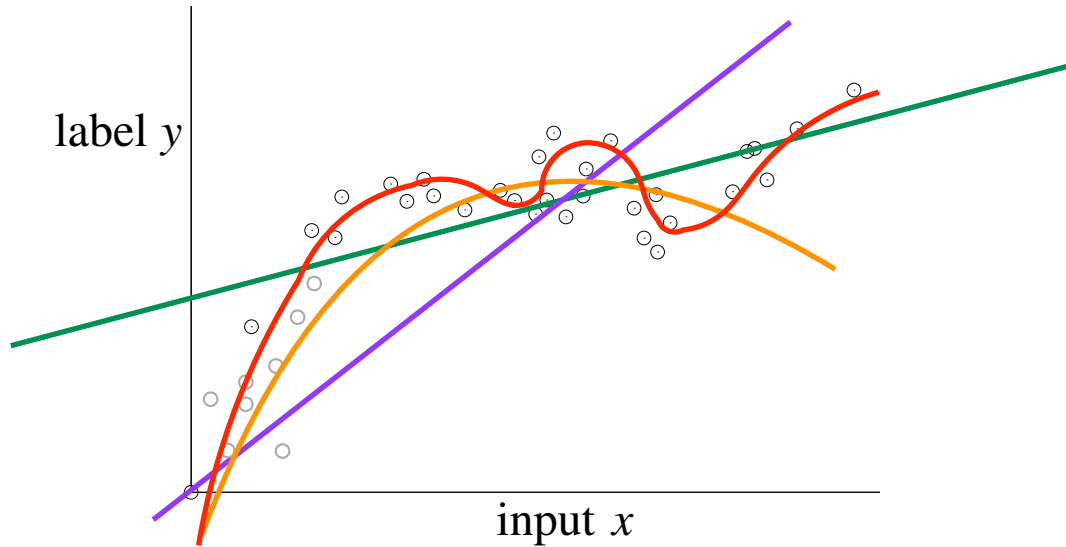
Quadratic regression in 1-dimension



Which line gives us the best fit?

- **Green & purple?** Can't capture variance. **Underfitting**
- **Orange?** Can't capture variance. **Underfitting**
- **Red?**

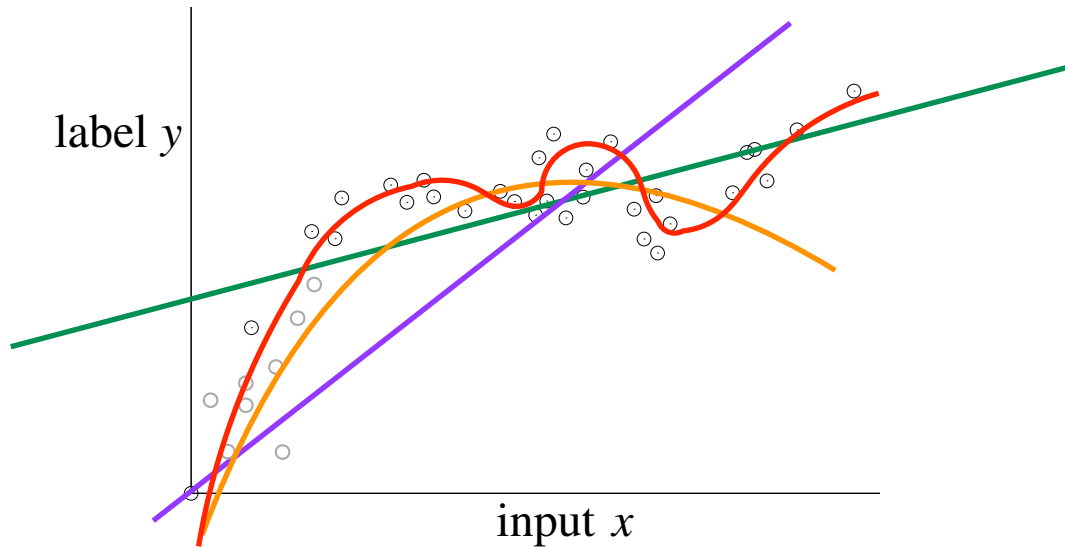
Quadratic regression in 1-dimension



Which line gives us the best fit?

- **Green** & **purple**? Can't capture variance. **Underfitting**
- **Orange**? Can't capture variance. **Underfitting**
- **Red**? Just right?

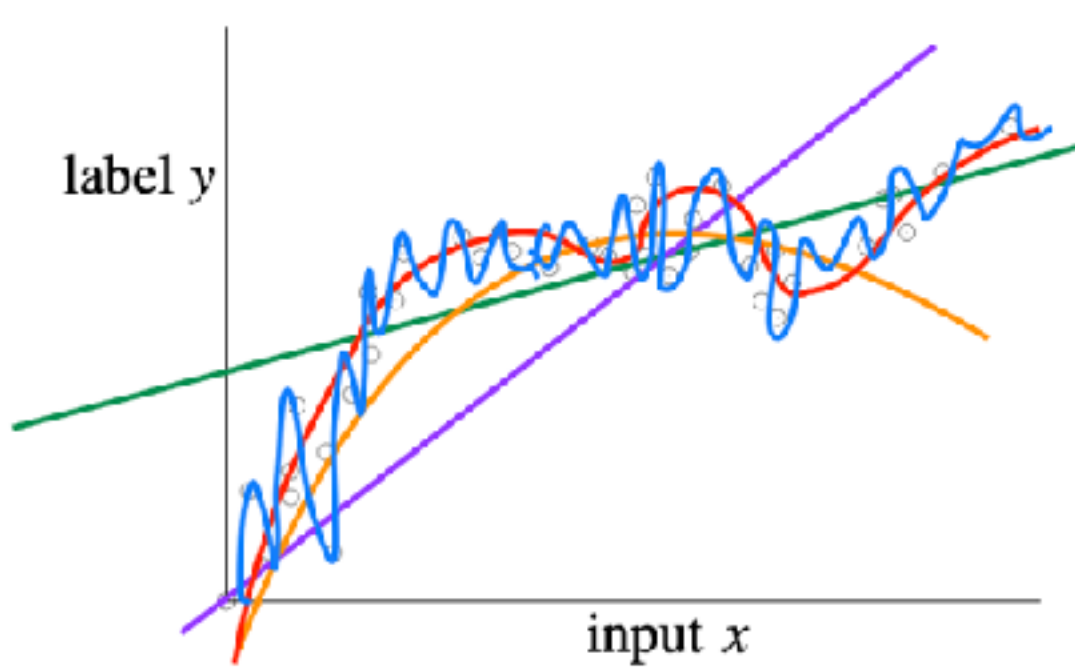
Quadratic regression in 1-dimension



Which line gives us the best fit?

- **Green** & **purple**? Can't capture variance. **Underfitting**
- **Orange**? Can't capture variance. **Underfitting**
- **Red**?
Just right?
Could we take it too far?

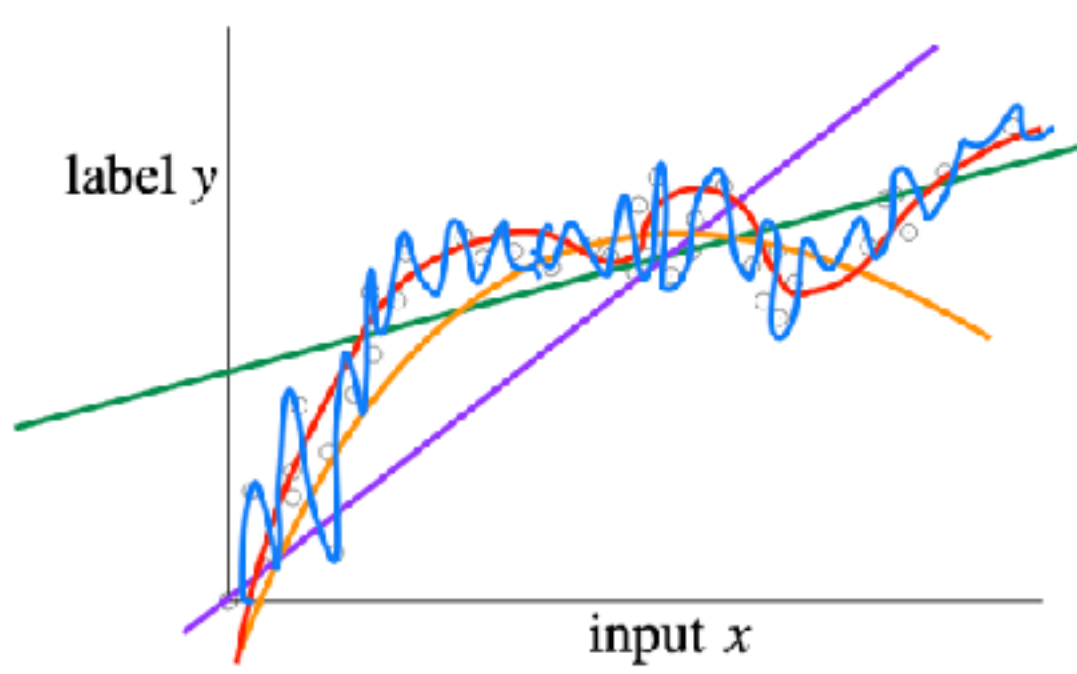
Quadratic regression in 1-dimension



Which line gives us the best fit?

- **Green** & **purple**? Can't capture variance. **Underfitting**
- **Orange**? Can't capture variance. **Underfitting**
- **Red**? Just right?
- **Blue**?

Quadratic regression in 1-dimension



Which line gives us the best fit?

- **Green** & **purple**? Can't capture variance. **Underfitting**
- **Orange**? Can't capture variance. **Underfitting**
- **Red**? Just right?
- **Blue**? Fits random fluctuations / measurement errors in training dataset. **Overfitting.**

Quadratic regression in 1-dimension

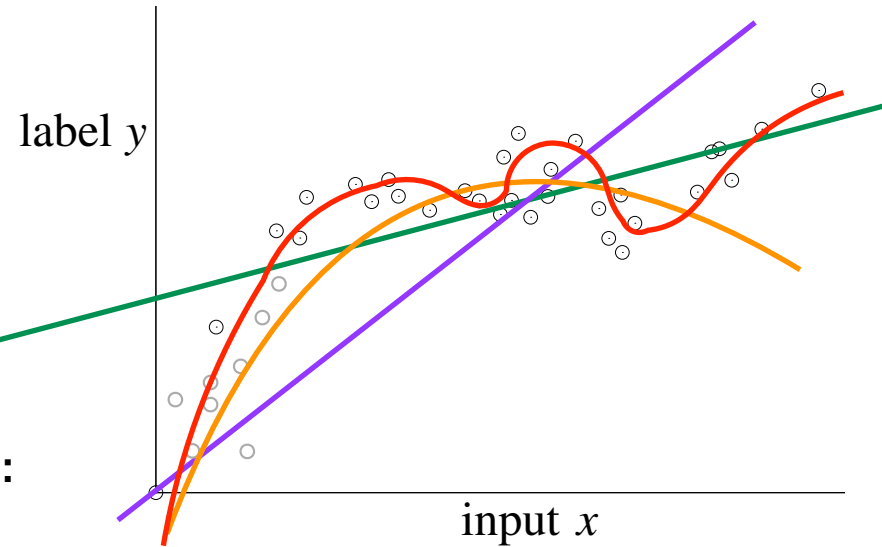
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter** $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$:

- $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$

Note: h can be arbitrary non-linear functions!

$$h(x) = \left[\log(x), x^2, \sin(x), \sqrt{x} \right]^\top$$



Breakout discussion:

Quadratic regression in 1-dimension

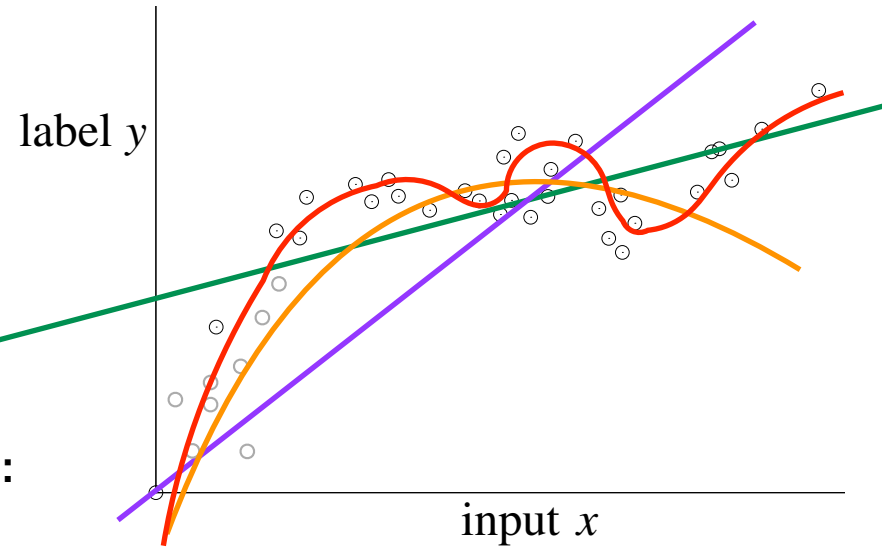
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter** $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$:
 - $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$

Note: h can be arbitrary non-linear functions!

$$h(x) = \left[\log(x), x^2, \sin(x), \sqrt{x} \right]^\top$$

Transform features however



Breakout discussion:

Quadratic regression in 1-dimension

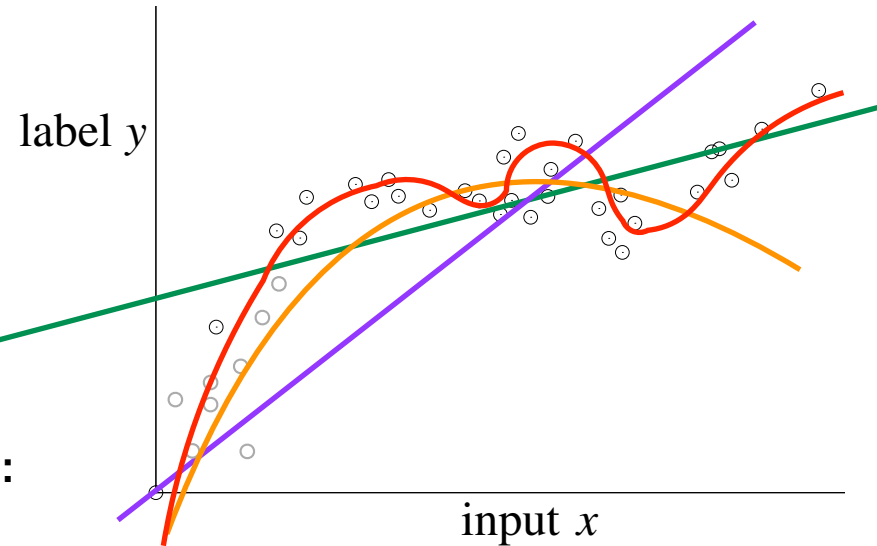
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter** $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$:
 - $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$

Note: h can be arbitrary non-linear functions!

$$h(x) = \left[\log(x), x^2, \sin(x), \sqrt{x} \right]^T$$

Transform features however



Breakout discussion:

- When would you use $\sin(x)$?
- How should I transform the house feature: zip code?

Quadratic regression in 1-dimension

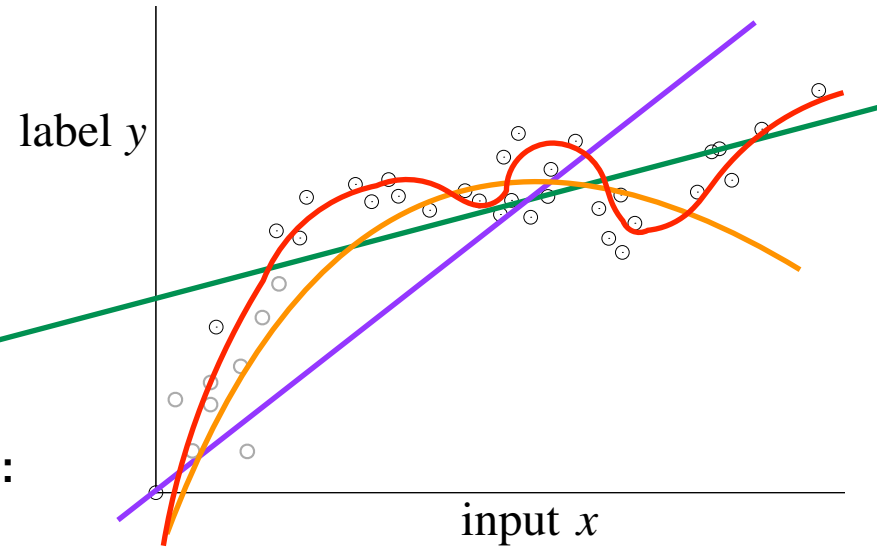
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter** $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$:
 - $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$

Note: h can be arbitrary non-linear functions!

$$h(x) = \left[\log(x), x^2, \sin(x), \sqrt{x} \right]^T$$

Transform features however



Breakout discussion:

- When would you use $\sin(x)$?
 - # Cyclic data e.g. weather
- How should I transform the house feature: zip code?

Quadratic regression in 1-dimension

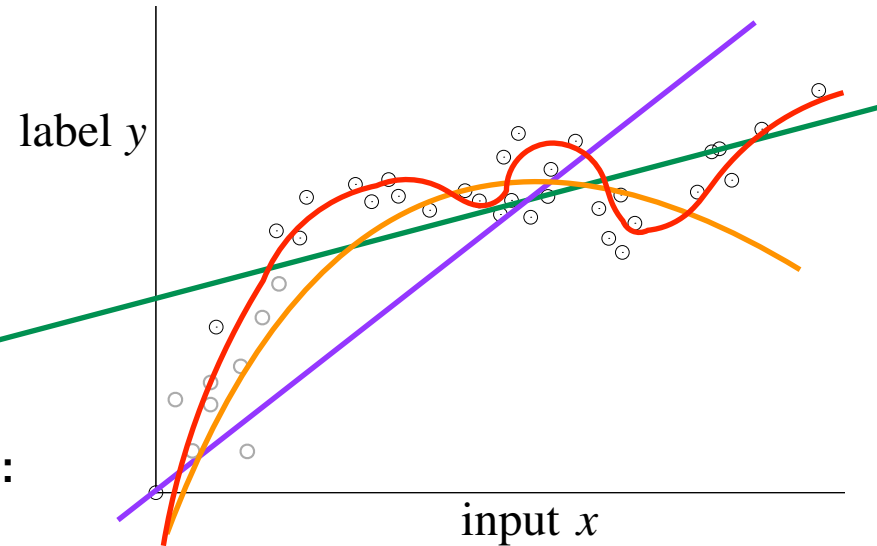
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter** $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$:
 - $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$

Note: h can be arbitrary non-linear functions!

$$h(x) = \left[\log(x), x^2, \sin(x), \sqrt{x} \right]^T$$

Transform features however



Breakout discussion:

- When would you use $\sin(x)$?
 - # Cyclic data e.g. weather
- How should I transform the house feature: zip code?
 - # Lat / long

Quadratic regression in 1-dimension

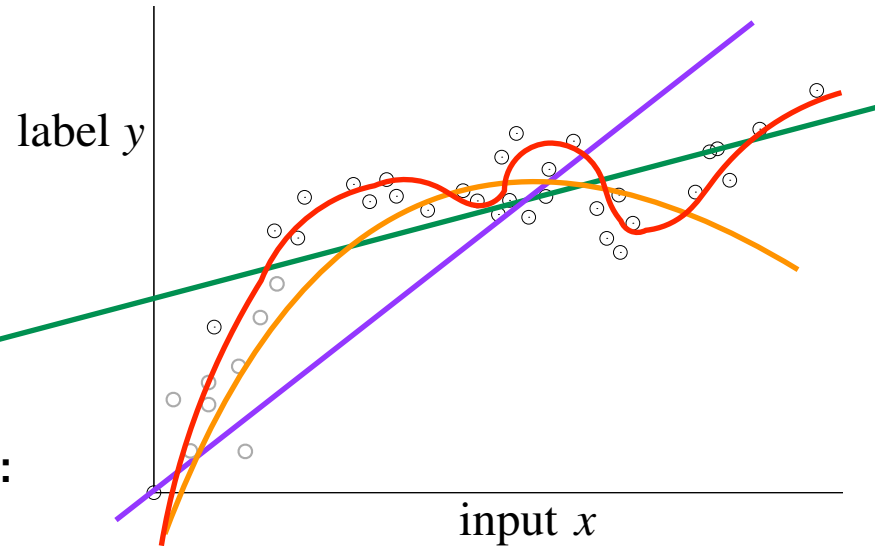
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter** $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$:
 - $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$

Note: h can be arbitrary non-linear functions!

$$h(x) = \left[\log(x), x^2, \sin(x), \sqrt{x} \right]^T$$

Transform features however



Breakout discussion:

- When would you use $\sin(x)$?
 - # Cyclic data e.g. weather
- How should I transform the house feature: zip code?
 - # Lat / long

This is the art of feature engineering

Quadratic regression in 1-dimension

- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

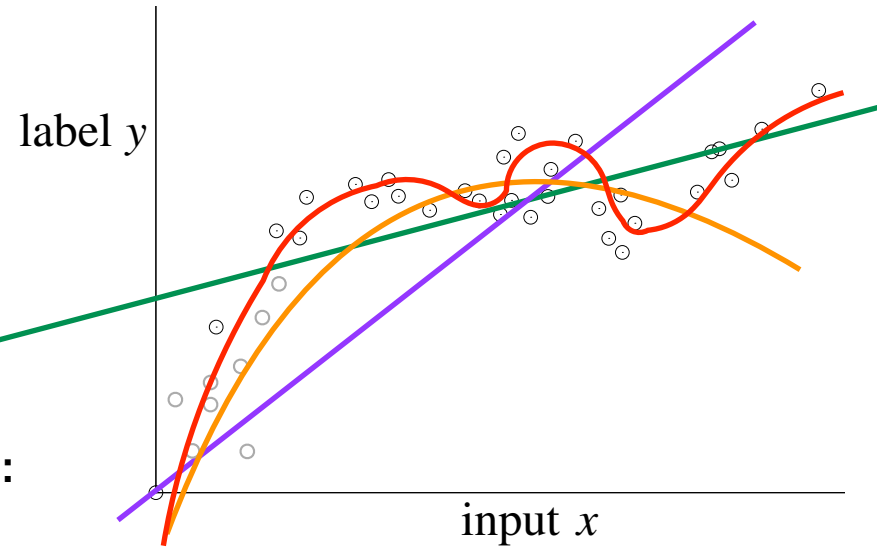
- **General p-features with parameter** $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$:
 - $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$

Note: h can be arbitrary non-linear functions!

$$h(x) = \left[\log(x), x^2, \sin(x), \sqrt{x} \right]^\top$$

Transform features however

How will I know if I have good features?



Breakout discussion:

- When would you use $\sin(x)$?

Cyclic data e.g. weather

- How should I transform the house feature: zip code?

Lat / long

This is the art of feature engineering

Quadratic regression in 1-dimension

- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter** $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$:
 - $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$

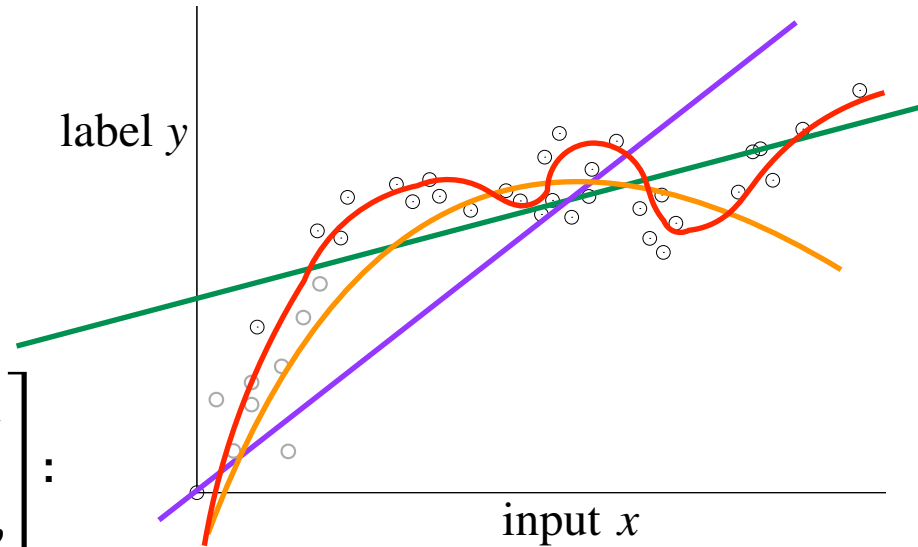
Note: h can be arbitrary non-linear functions!

$$h(x) = \left[\log(x), x^2, \sin(x), \sqrt{x} \right]^T$$

Transform features however

How will I know if I have good features?

PLOT PLOT PLOT PLOT PLOT



Breakout discussion:

- When would you use $\sin(x)$?

Cyclic data e.g. weather

- How should I transform the house feature: zip code?

Lat / long

This is the art of feature engineering

Quadratic regression in 1-dimension

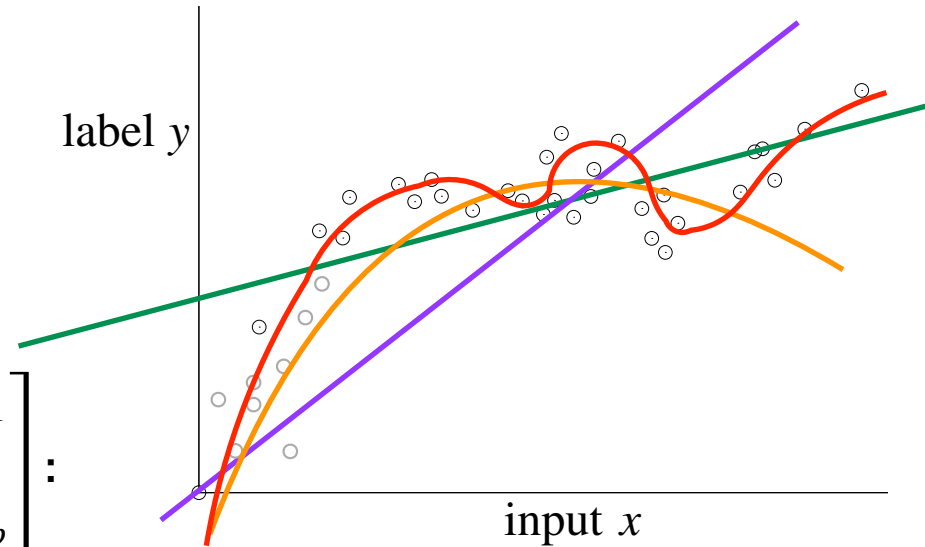
- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter $w =$**

$$\begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix} :$$

- $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$

How do we learn w ?



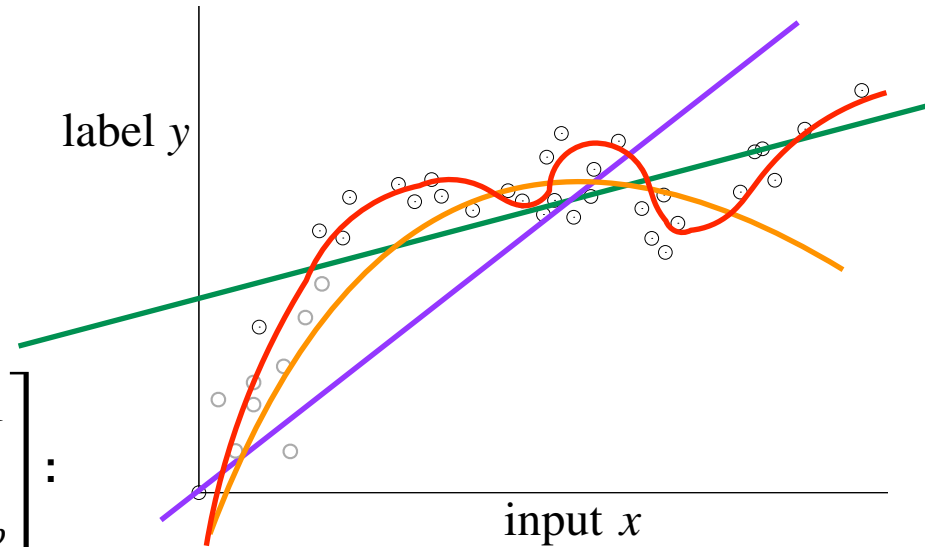
Quadratic regression in 1-dimension

- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter $w =$**

$$\begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix} :$$

- $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$



How do we learn w ?

$$\mathbf{H} = \begin{bmatrix} - & - & h(x_1)^\top & - & - \\ & & \vdots & & \\ - & - & h(x_n)^\top & - & - \end{bmatrix} \in \mathbb{R}^{n \times p}$$

$$\hat{w} = \arg \min_w \|\mathbf{H}w - \mathbf{y}\|_2^2$$

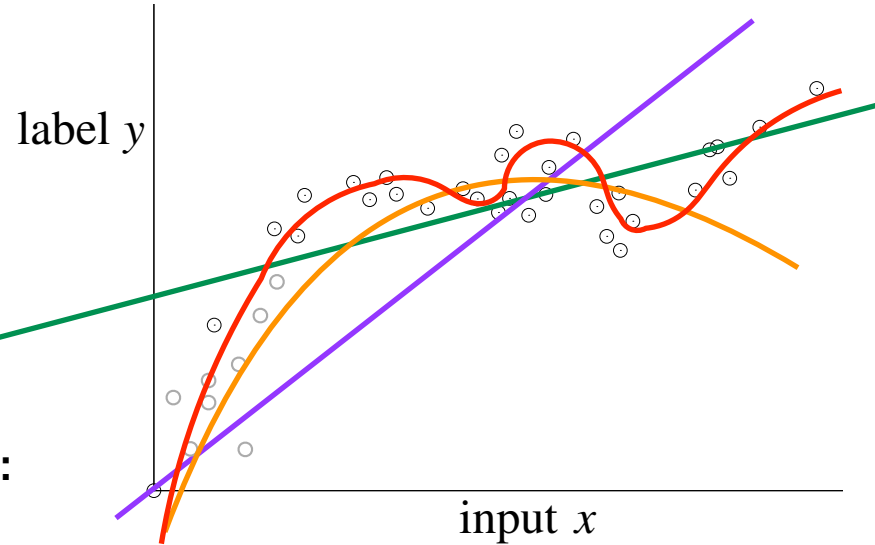
For a new test point x , predict
 $\hat{y} = \langle \hat{w}, h(x) \rangle$

Quadratic regression in 1-dimension

- **Data:** $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$

- **General p-features with parameter** $w = \begin{bmatrix} w_1 \\ \vdots \\ w_p \end{bmatrix}$:

- $\hat{y}_i = \langle w, h(x_i) \rangle$ where $h : \mathbb{R} \rightarrow \mathbb{R}^p$



How do we learn w ?

$$\mathbf{H} = \begin{bmatrix} - & - & h(x_1)^T & - & - \\ & & \vdots & & \\ - & - & h(x_n)^T & - & - \end{bmatrix} \in \mathbb{R}^{n \times p}$$

$$\hat{w} = \arg \min_w \|\mathbf{H}w - \mathbf{y}\|_2^2$$

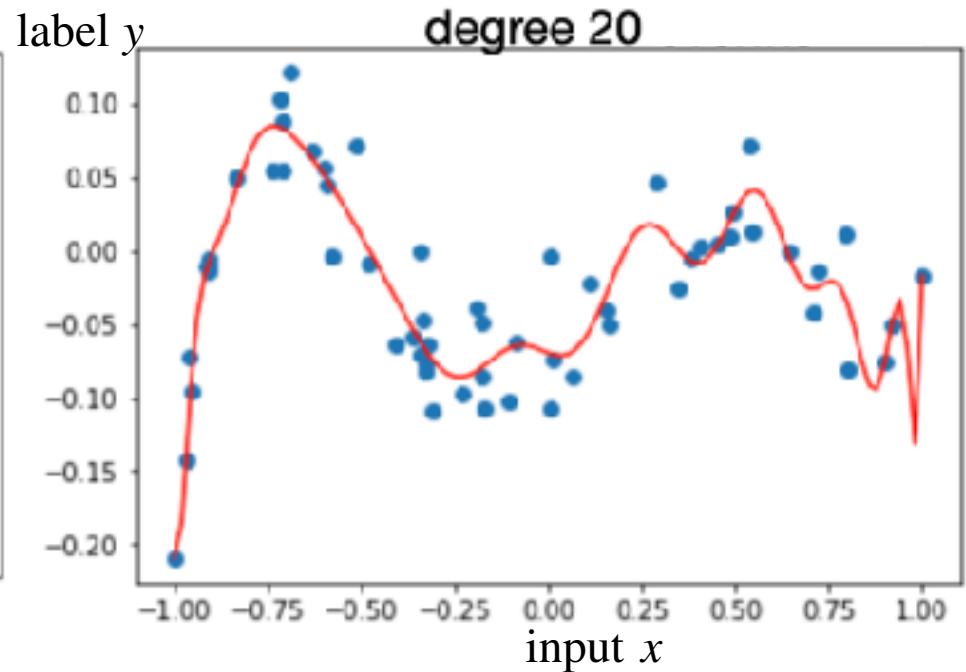
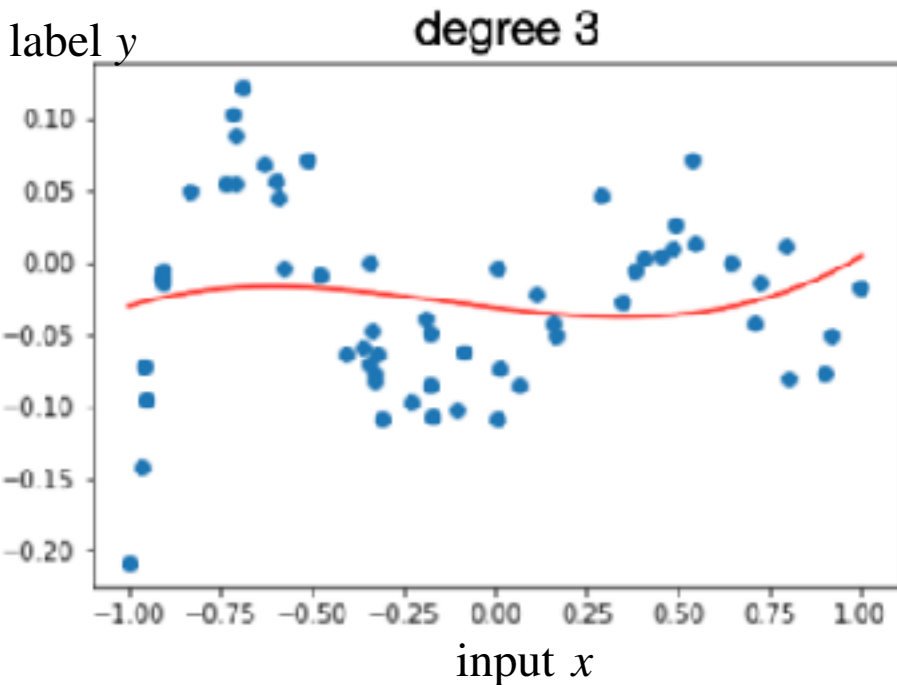
For a new test point x , predict
 $\hat{y} = \langle \hat{w}, h(x) \rangle$

$$\hat{w}_{\text{MLE}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{y}$$

- Training and test dataset split
 - Cross validation
-

Which p should we choose?

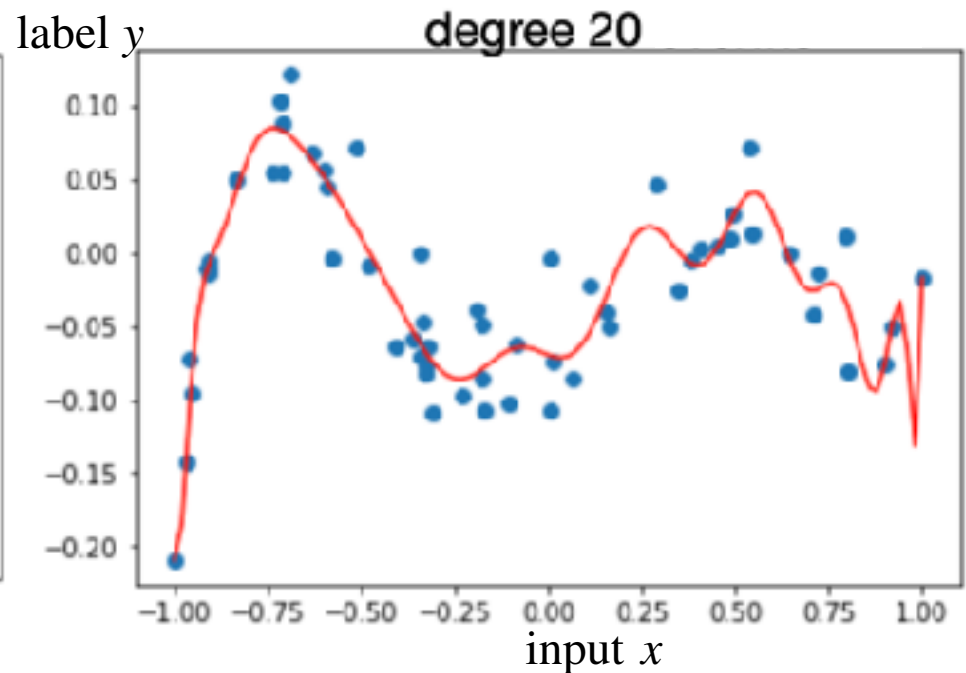
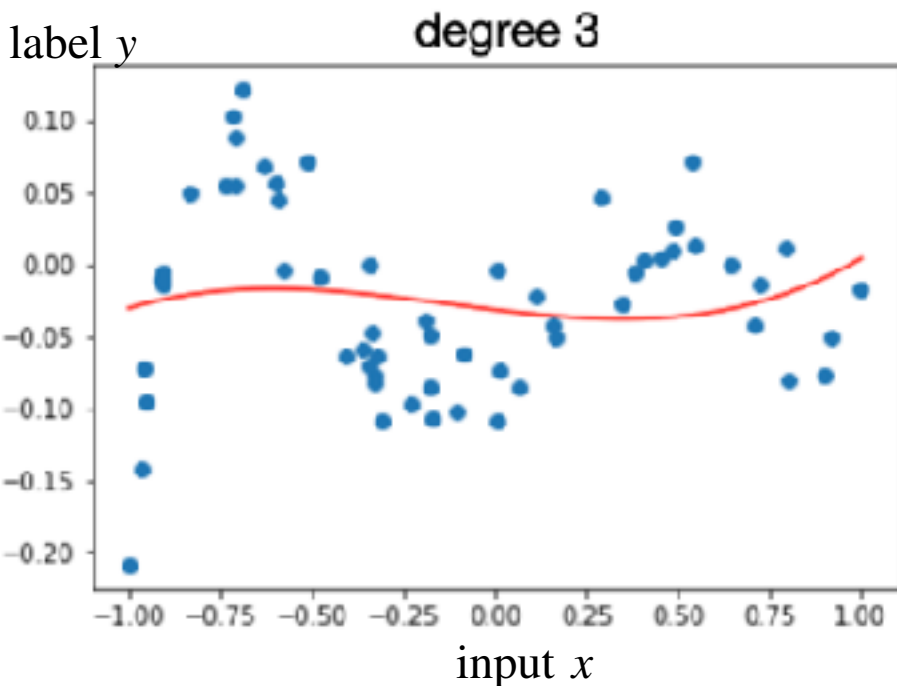
- First instance of class of models with different representation power = model complexity



- How do we determine which is better model?

Which p should we choose?

- First instance of class of models with different representation power = model complexity

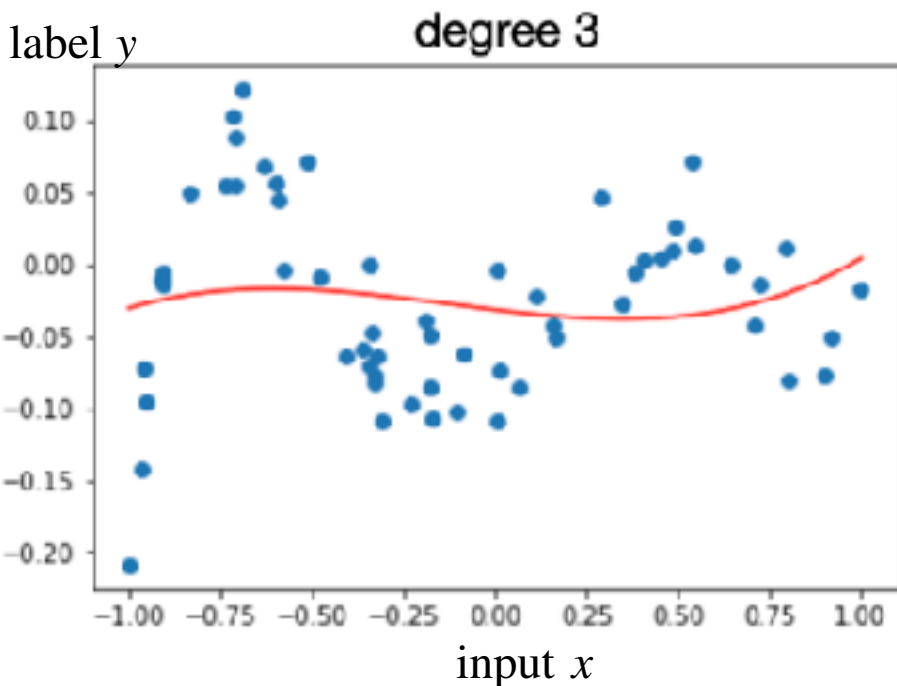


Underfitting: not fitting

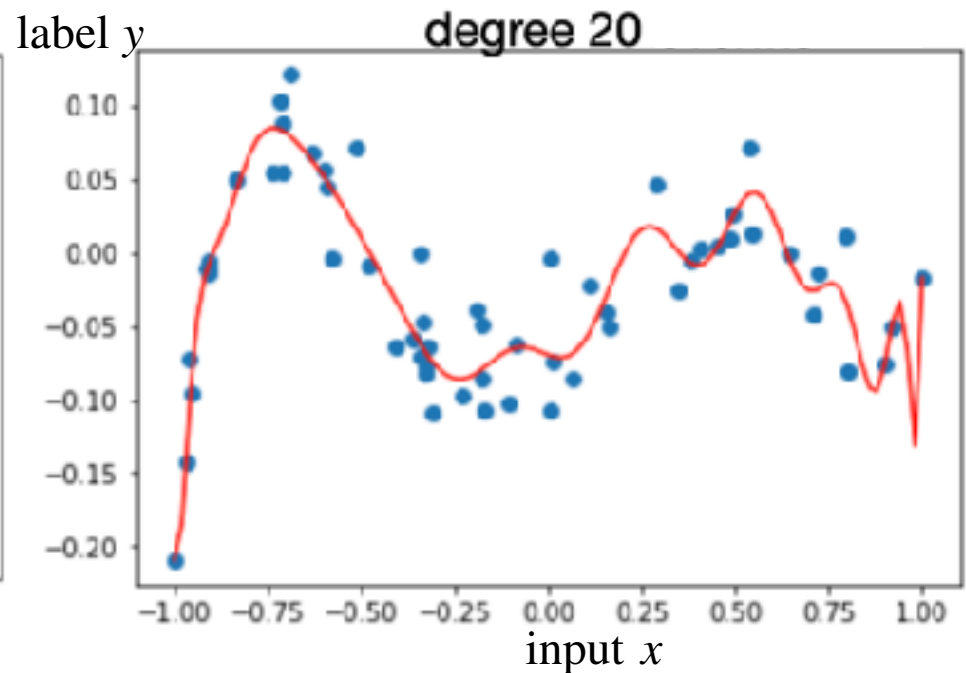
- How do we determine which is better model?

Which p should we choose?

- First instance of class of models with different representation power = model complexity



Underfitting: not fitting

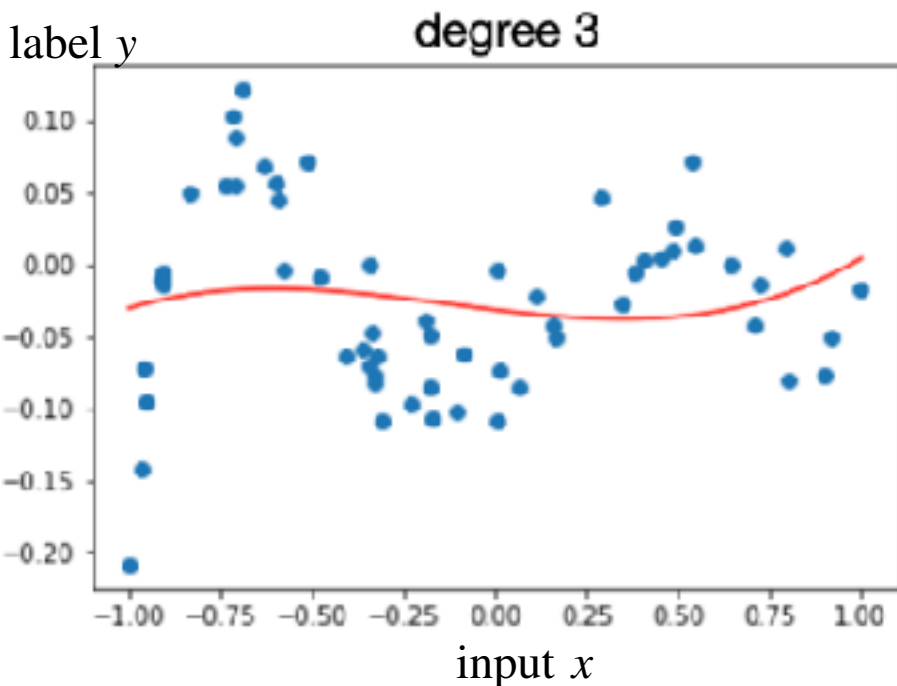


Overfitting: fitting “idiosyncrasies” of training data that won’t generalize

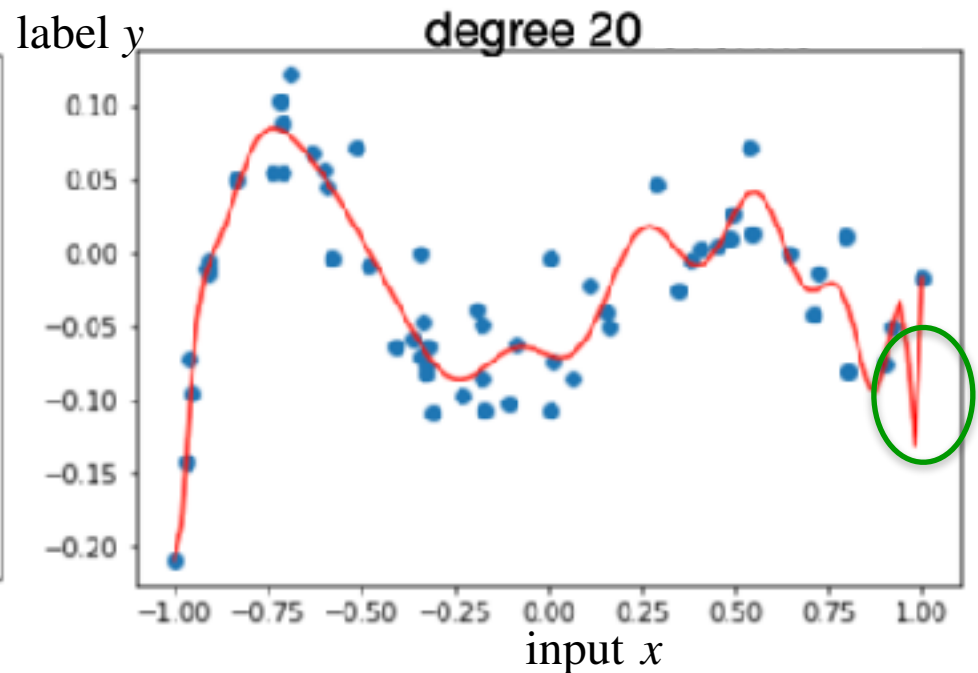
- How do we determine which is better model?

Which p should we choose?

- First instance of class of models with different representation power = model complexity



Underfitting: not fitting



Overfitting: fitting “idiosyncrasies” of training data that won’t generalize

- How do we determine which is better model?

Generalization

- we say a predictor **generalizes** if it performs as well on unseen data as on training data (we will formalize shortly)
- the data used to train a predictor is **training data** or **in-sample data**
- we want the predictor to work on **out-of-sample data**
- we say a predictor **fails to generalize** if it performs well on in-sample data but does not perform well on out-of-sample data

Generalization

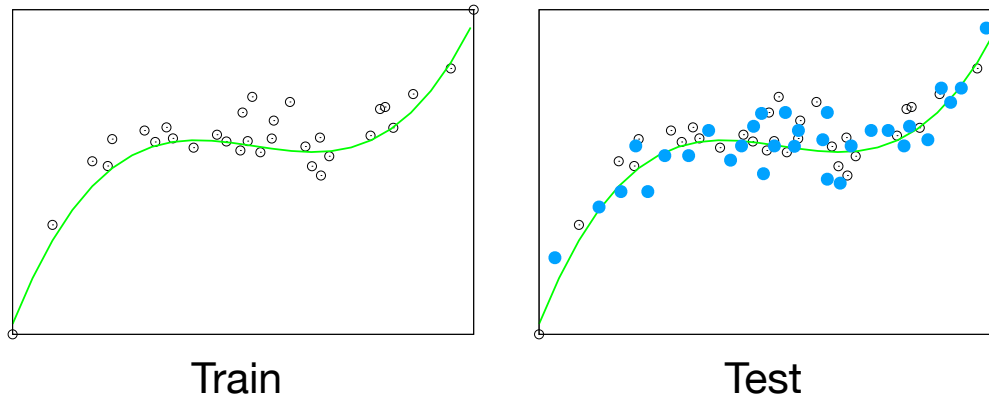
Test data

- we say a predictor **generalizes** if it performs as well on unseen data as on training data (we will formalize shortly)
- the data used to train a predictor is **training data** or **in-sample data**
- we want the predictor to work on **out-of-sample data**
- we say a predictor **fails to generalize** if it performs well on in-sample data but does not perform well on out-of-sample data

Generalization

Consider the following:

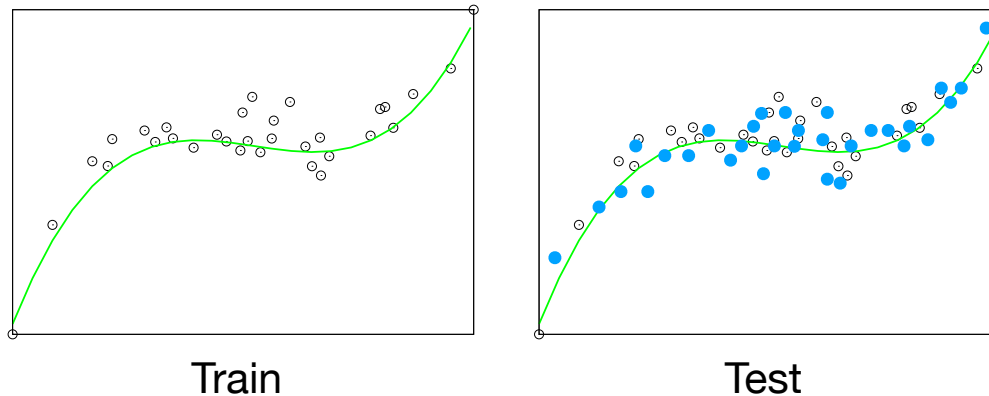
- **train** a cubic predictor on 32 (**in-sample**) white circles: Mean Squared Error (MSE) 174
- **predict** label y for 30 (**out-of-sample**) blue circles: MSE 192
- **Would you conclude this predictor generalizes effectively?**



Generalization

Consider the following:

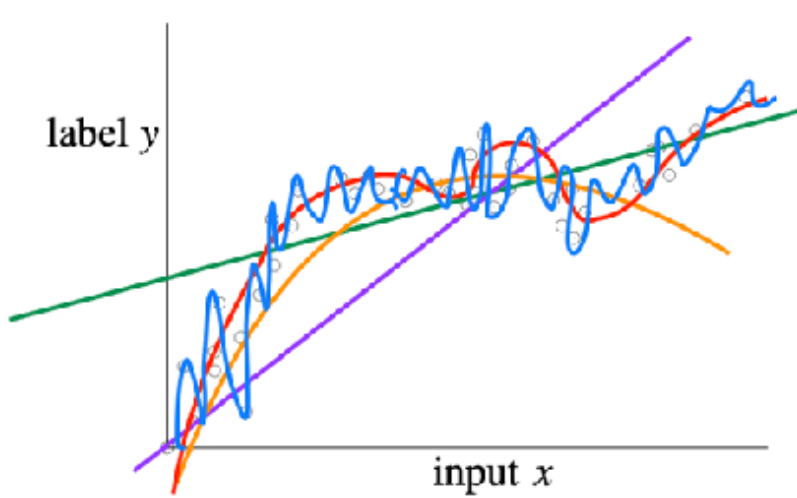
- **train** a cubic predictor on 32 (**in-sample**) white circles: Mean Squared Error (MSE) 174
- **predict** label y for 30 (**out-of-sample**) blue circles: MSE 192
- **Would you conclude this predictor generalizes effectively?**
 - **Yes**, as in-sample MSE \simeq out-of-sample MSE



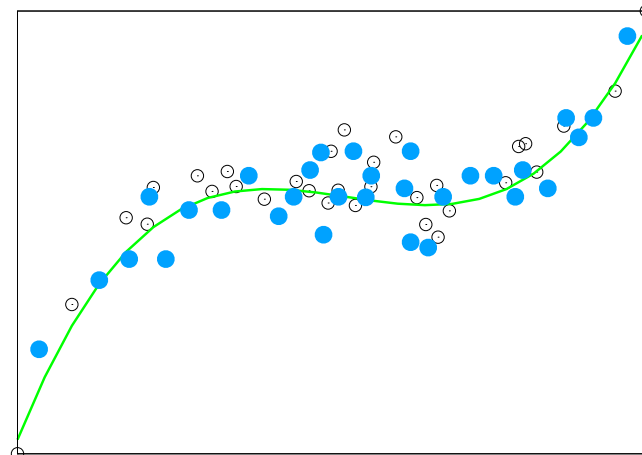
Generalization

Consider the following:

- **train** a cubic predictor on 32 (**in-sample**) white circles: Mean Squared Error (MSE) = 2
- **predict** label y for 30 (**out-of-sample**) blue circles: MSE = 356
- **Would you conclude this predictor generalizes effectively?**
 - **No**, as in-sample MSE \ll out-of-sample MSE



Train



Split the data into **training** and **testing**

- a way to mimic how the predictor performs on unseen data
- given a single dataset $S = \{(x_i, y_i)\}_{i=1}^n$
- we split the dataset into two: training set and test set (e.g., 90/10)
- **training set** used to train the model

- minimize $\mathcal{L}_{\text{train}}(w) = \frac{1}{|S_{\text{train}}|} \sum_{i \in S_{\text{train}}} (y_i - x_i^T w)^2$ # Fit the params of the model

- **test set** used to evaluate the model

- $\mathcal{L}_{\text{test}}(w) = \frac{1}{|S_{\text{test}}|} \sum_{i \in S_{\text{test}}} (y_i - x_i^T w)^2$

- this assumes that test set is similar to unseen data
- **test set should never be used in training or picking unknowns!!**

Split the data into **training** and **testing**

- a way to mimic how the predictor performs on unseen data
- given a single dataset $S = \{(x_i, y_i)\}_{i=1}^n$
- we split the dataset into two: training set and test set (e.g., 90/10)
- **training set** used to train the model

- minimize $\mathcal{L}_{\text{train}}(w) = \frac{1}{|S_{\text{train}}|} \sum_{i \in S_{\text{train}}} (y_i - x_i^T w)^2$

Fit the params
of the model

- **test set** used to evaluate the model

- $\mathcal{L}_{\text{test}}(w) = \frac{1}{|S_{\text{test}}|} \sum_{i \in S_{\text{test}}} (y_i - x_i^T w)^2$

- this assumes that test set is similar to unseen data

- **test set should never be used in training or picking unknowns!!**

Split the data into **training** and **testing**

- a way to mimic how the predictor performs on unseen data
- given a single dataset $S = \{(x_i, y_i)\}_{i=1}^n$
- we split the dataset into two: training set and test set (e.g., 90/10)
- **training set** used to train the model

- minimize $\mathcal{L}_{\text{train}}(w) = \frac{1}{|S_{\text{train}}|} \sum_{i \in S_{\text{train}}} (y_i - x_i^T w)^2$ # Fit the params of the model

- **test set** used to evaluate the model

- $\mathcal{L}_{\text{test}}(w) = \frac{1}{|S_{\text{test}}|} \sum_{i \in S_{\text{test}}} (y_i - x_i^T w)^2$

- this assumes that test set is similar to unseen data

- **test set should never be used in training or picking unknowns!!**



Split the data into **training** and **testing**

- a way to mimic how the predictor performs on unseen data
- given a single dataset $S = \{(x_i, y_i)\}_{i=1}^n$
- we split the dataset into two: training set and test set (e.g., 90/10)
- **training set** used to train the model

- minimize $\mathcal{L}_{\text{train}}(w) = \frac{1}{|S_{\text{train}}|} \sum_{i \in S_{\text{train}}} (y_i - x_i^T w)^2$ # Fit the params of the model

- **test set** used to evaluate the model

- $\mathcal{L}_{\text{test}}(w) = \frac{1}{|S_{\text{test}}|} \sum_{i \in S_{\text{test}}} (y_i - x_i^T w)^2$

- this assumes that test set is similar to unseen data

- **test set should never be used in training or picking unknowns!!**



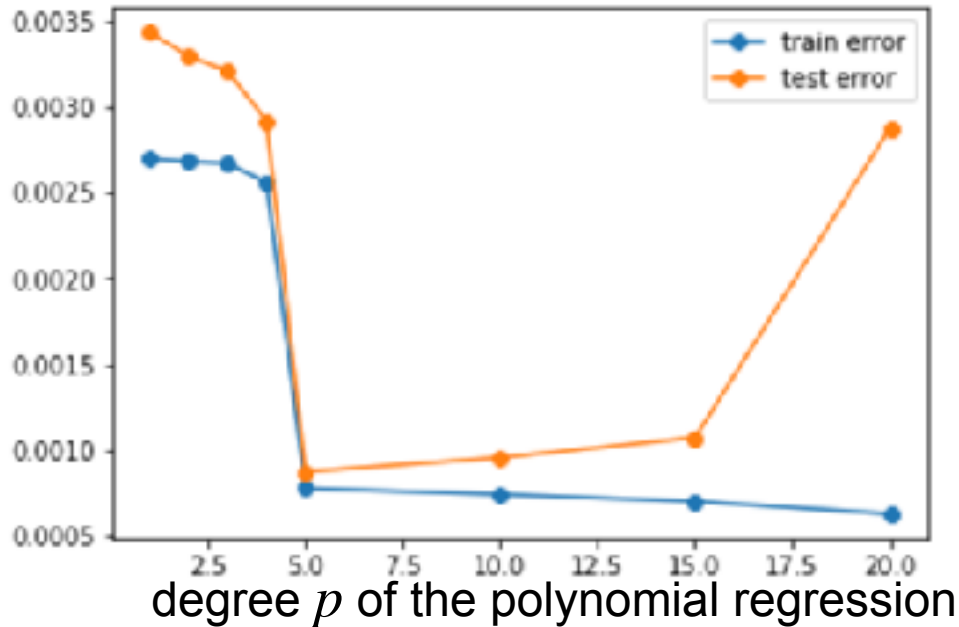
This is cheating and will lead you to the wrong conclusions. Big statistical no no

Demo: Polynomial Regression

See [colab notebook demo](#)

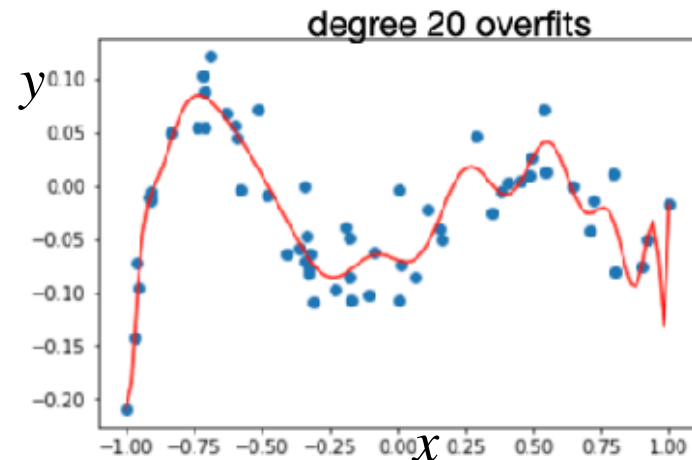
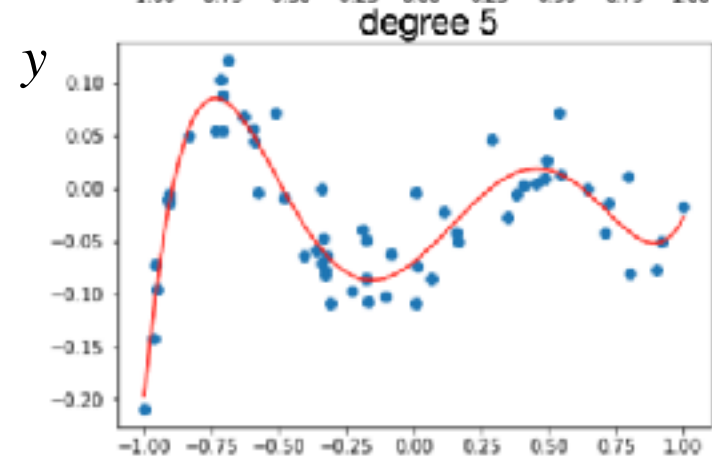
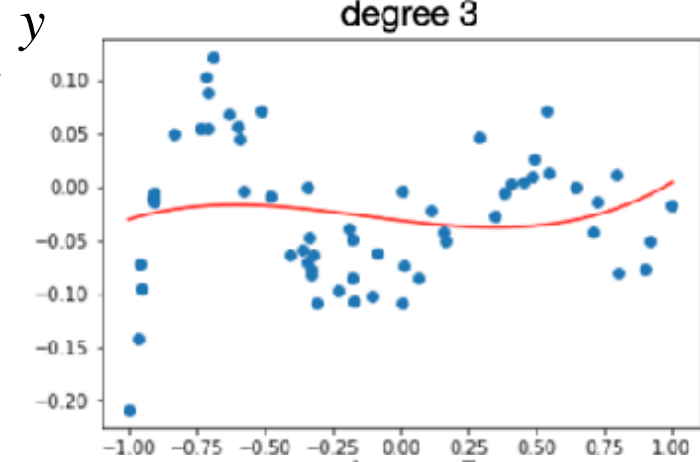
Train/test error vs. complexity

Error



- Degree $p = 5$, since it achieves **minimum test error**
- **Train error** monotonically decreases with model complexity
- **Test error** has a U shape. There is a best value

test set should never be used in training or picking degree



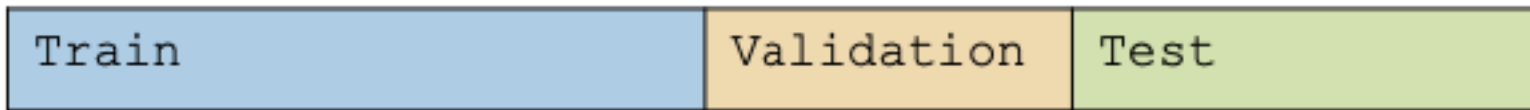
Cross-Validation

How... How... How???????

- > How do we pick the number of basis functions...
- > We could use the test data, but...

Validation Set

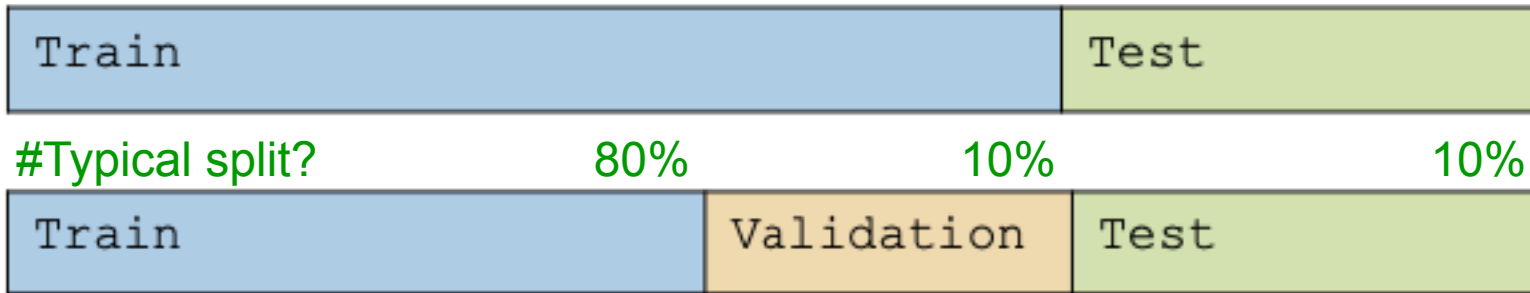
- > Simplest approach is to add another dataset partition called the **validation set**



- > Each set has its own role
 - > **Train:** Used to train a model of each model complexity
 - > **Validation:** Used as a hold-out to evaluate each model. Generally choose model complexity with lowest validation error.
 - > **Test:** Once the model is chosen, can get an estimate of future error by test. This would be what you report. **Only do this once!**

Validation Set

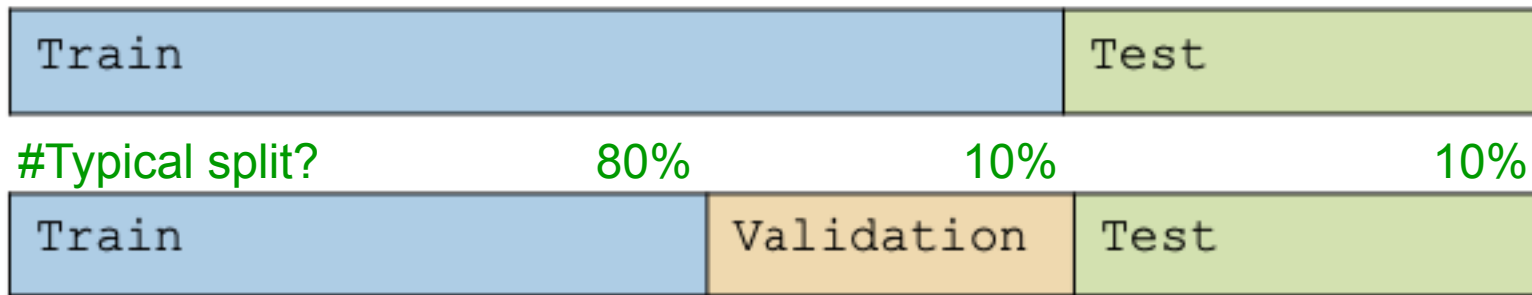
- > Simplest approach is to add another dataset partition called the **validation set**



- > Each set has its own role
 - > **Train:** Used to train a model of each model complexity
 - > **Validation:** Used as a hold-out to evaluate each model. Generally choose model complexity with lowest validation error.
 - > **Test:** Once the model is chosen, can get an estimate of future error by test. This would be what you report. **Only do this once!**

Validation Set

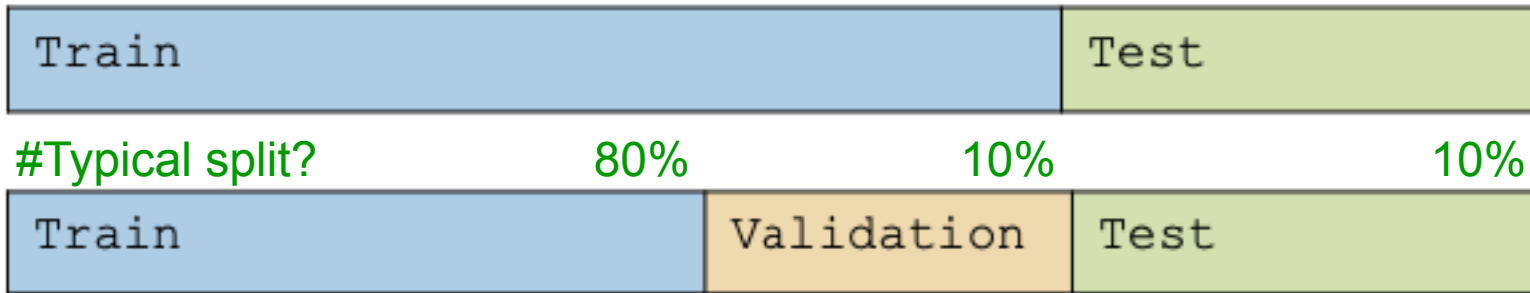
- > Simplest approach is to add another dataset partition called the **validation set**



- > Each set has its own role
 - # Fit model params like w
 - > **Train:** Used to train a model of each model complexity
 - > **Validation:** Used as a hold-out to evaluate each model. Generally choose model complexity with lowest validation error.
 - > **Test:** Once the model is chosen, can get an estimate of future error by test. This would be what you report. **Only do this once!**

Validation Set

- > Simplest approach is to add another dataset partition called the **validation set**



- > Each set has its own role
 - # Fit model params like w
 - > **Train:** Used to train a model of each model complexity
 - > **Validation:** Used as a hold-out to evaluate each model. Generally choose model complexity with lowest validation error. # Choose *hyperparameters*, like polynomial degree p
 - > **Test:** Once the model is chosen, can get an estimate of future error by test. This would be what you report. **Only do this once!**

(LOO) Leave-one-out cross validation

- > Consider a validation set with 1 example:
 - D – training data
 - $D \setminus j$ – training data with j th data point (x_j, y_j) moved to validation set
- > Learn model $f_{D \setminus j}$ with $D \setminus j$ dataset
- > Estimate true error as squared error on predicting y_j :
 - Unbiased estimate of error $\text{error}_{\text{true}}(f_{D \setminus j})!$

(LOO) Leave-one-out cross validation

- > Consider a validation set with 1 example:
 - D – training data # n samples
 - $D \setminus j$ – training data with j th data point (x_j, y_j) moved to validation set
- > Learn model $f_{D \setminus j}$ with $D \setminus j$ dataset
- > Estimate true error as squared error on predicting y_j :
 - Unbiased estimate of error $\text{error}_{\text{true}}(f_{D \setminus j})!$

(LOO) Leave-one-out cross validation

- > Consider a validation set with 1 example:
 - D – training data # n samples
 - $D \setminus j$ – training data with j th data point (x_j, y_j) moved to validation set
- > Learn model $f_{D \setminus j}$ with $D \setminus j$ dataset # $n-1$ samples
- > Estimate true error as squared error on predicting y_j :
 - Unbiased estimate of error $\text{error}_{\text{true}}(f_{D \setminus j})!$

(LOO) Leave-one-out cross validation

- > Consider a validation set with 1 example:
 - D – training data
 - $D \setminus j$ – training data with j th data point (x_j, y_j) moved to validation set
- > Learn model $f_{D \setminus j}$ with $D \setminus j$ dataset
- > Estimate true error as squared error on predicting y_j :
 - Unbiased estimate of error_{true}($f_{D \setminus j}$)!
- > LOO cross validation: Average over all data points j :
 - For each data point you leave out, learn a new model $f_{D \setminus j}$
 - Estimate error as:

$$\text{error}_{LOO} = \frac{1}{n} \sum_{j=1}^n (y_j - f_{D \setminus j}(x_j))^2$$

LOO cross validation is (almost) unbiased estimate!

- > When computing LOOCV error, we only use $N-1$ data points
 - So it's not quite an estimate of true error of learning with N data points
 - Usually pessimistic, though – learning with less data typically gives worse answer
- > LOO is almost unbiased! Why not use LOO error for model selection!
 - E.g., picking degree



LOO cross validation is (almost) unbiased estimate!

- > When computing LOOCV error, we only use $N-1$ data points
 - So it's not quite an estimate of true error of learning with N data points
 - Usually pessimistic, though – learning with less data typically gives worse answer
- > LOO is almost unbiased! Why not use LOO error for model selection!
 - E.g., picking degree



You just fit n models!!

Computational cost of LOO

- > **Suppose you have 100,000 data points**
- > **You implemented a great version of your learning algorithm**
 - **Learns in only 1 second**
- > **Computing LOO will take about 1 day!!!**
 -

Use k -fold cross validation

- > Randomly divide training data into k equal parts
 - D_1, \dots, D_k
- > For each i
 - Learn model $f_{D \setminus D_i}$ using data point not in D_i
 - Estimate error of $f_{D \setminus D_i}$ on validation set D_i :



$$\text{error}_{D_i} = \frac{1}{|D_i|} \sum_{(x_j, y_j) \in D_i} (y_j - f_{D \setminus D_i}(x_j))^2$$

Use k -fold cross validation

> Randomly divide training data into k equal parts

- D_1, \dots, D_k

> For each i

- Learn model $f_{D \setminus D_i}$ using data point not in D_i
- Estimate error of $f_{D \setminus D_i}$ on validation set D_i :

$$\text{error}_{D_i} = \frac{1}{|D_i|} \sum_{(x_j, y_j) \in D_i} (y_j - f_{D \setminus D_i}(x_j))^2$$

> k -fold cross validation error is average over data splits:

$$\text{error}_{k\text{-fold}} = \frac{1}{k} \sum_{i=1}^k \text{error}_{D_i}$$

> k -fold cross validation properties:

- Much faster to compute than LOO
- More (pessimistically) biased – using much less data, only $n(k-1)/k$
- Usually, $k = 10$

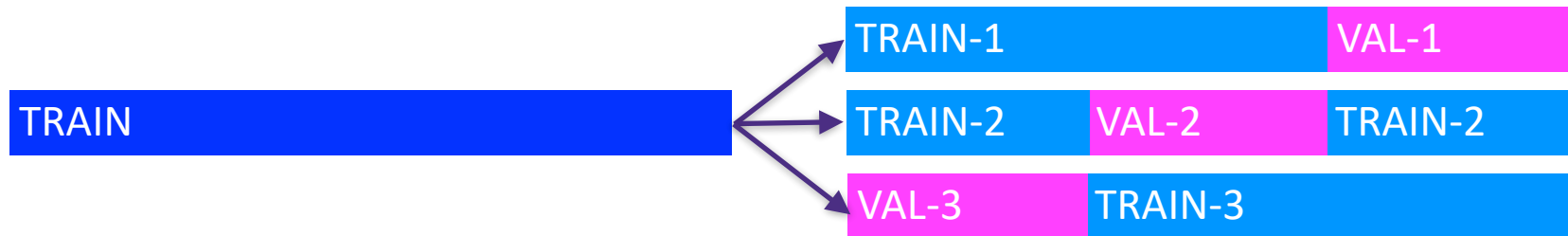
1	2	3	4	5
Train	Train	Validation	Train	Train

Recap

- > Given a dataset, begin by splitting into



- > Model selection: Use k-fold cross-validation on **TRAIN** to train predictor and choose magic parameters such as degree



- > Model assessment: Use **TEST** to assess the accuracy of the model you output
 - **Never ever ever ever ever train or choose parameters based on the test data**

Example 1

- > You wish to predict the stock price of zoom.us given historical stock price data
- > You use all daily stock price up to Jan 1, 2020 as **TRAIN** and Jan 2, 2020 - April 13, 2020 as **TEST**
- > What's "wrong" with this procedure?

Example 1

- > You wish to predict the stock price of zoom.us given historical stock price data
- > You use all daily stock price up to Jan 1, 2020 as **TRAIN** and Jan 2, 2020 - April 13, 2020 as **TEST**
- > What's "wrong" with this procedure?
 - ML assumes data is i.i.d.; is this?

Example 1

- > You wish to predict the stock price of zoom.us given historical stock price data
- > You use all daily stock price up to Jan 1, 2020 as **TRAIN** and Jan 2, 2020 - April 13, 2020 as **TEST**
- > What's "wrong" with this procedure?
 - ML assumes data is i.i.d.; is this?
 - Nope! Future data is o.o.d. (out of distribution) from the training data

Example 1

- > You wish to predict the stock price of zoom.us given historical stock price data
- > You use all daily stock price up to Jan 1, 2020 as **TRAIN** and Jan 2, 2020 - April 13, 2020 as **TEST**
- > What's “wrong” with this procedure?
 - ML assumes data is i.i.d.; is this?
 - Nope! Future data is o.o.d. (out of distribution) from the training data
 - Why is “wrong” in square quotes?

Example 1

- > You wish to predict the stock price of zoom.us given historical stock price data
- > You use all daily stock price up to Jan 1, 2020 as **TRAIN** and Jan 2, 2020 - April 13, 2020 as **TEST**
- > What's "wrong" with this procedure?
 - ML assumes data is i.i.d.; is this?
 - Nope! Future data is o.o.d. (out of distribution) from the training data
 - Why is "wrong" in square quotes?
 - You might still want to do this. It's economically valuable, if you could.

Example 2

- > **Given 10,000-dimensional data and n examples, we pick a subset of 50 dimensions that have the highest correlation with labels in the entire dataset:**

50 indices j that have largest

$$\frac{|\sum_{i=1}^n x_{i,j} y_i|}{\sqrt{\sum_{i=1}^n x_{i,j}^2}}$$

- > **After picking our 50 features, we then break data into train and test dataset.**
- > **We train linear regression on these selected features on the training set. We compute the test error and report it**
- > **What's wrong with this procedure?**

Example 2

- > **Given 10,000-dimensional data and n examples, we pick a subset of 50 dimensions that have the highest correlation with labels in the entire dataset:**

50 indices j that have largest

$$\frac{|\sum_{i=1}^n x_{i,j} y_i|}{\sqrt{\sum_{i=1}^n x_{i,j}^2}}$$

- > **After picking our 50 features, we then break data into train and test dataset.**
- > **We train linear regression on these selected features on the training set. We compute the test error and report it**
- > **What's wrong with this procedure?**
 - **Need to hold out the test data first, otherwise you're contaminating your test data.**

Example 2

- > **Given 10,000-dimensional data and n examples, we pick a subset of 50 dimensions that have the highest correlation with labels in the entire dataset:**

50 indices j that have largest

$$\frac{|\sum_{i=1}^n x_{i,j} y_i|}{\sqrt{\sum_{i=1}^n x_{i,j}^2}}$$

- > **After picking our 50 features, we then break data into train and test dataset.**
- > **We train linear regression on these selected features on the training set. We compute the test error and report it**
- > **What's wrong with this procedure?**
 - Need to hold out the test data first, otherwise you're contaminating your test data.
 - This happened in my first ML project. Thought the accuracy was 90%. When I fixed, it went down to 60%.

Example 2

- > **Given 10,000-dimensional data and n examples, we pick a subset of 50 dimensions that have the highest correlation with labels in the entire dataset:**

50 indices j that have largest

$$\frac{|\sum_{i=1}^n x_{i,j} y_i|}{\sqrt{\sum_{i=1}^n x_{i,j}^2}}$$

- > **After picking our 50 features, we then break data into train and test dataset.**
- > **We train linear regression on these selected features on the training set. We compute the test error and report it**
- > **What's wrong with this procedure?**
 - Need to hold out the test data first, otherwise you're contaminating your test data.
 - This happened in my first ML project. Thought the accuracy was 90%. When I fixed, it went down to 60%.
 - Top 50 features might be redundant with each other.

Recap

- > Learning is...
 - Collect some data
 - > E.g., housing info and sale price
 - Randomly split dataset into TRAIN, VAL, and TEST
 - > E.g., 80%, 10%, and 10%, respectively
 - Choose a hypothesis class or model
 - > E.g., linear with non-linear transformations
 - Choose a loss function
 - > E.g., least squares on TRAIN
 - Choose an optimization procedure
 - > E.g., set derivative to zero to obtain estimator, cross-validation on VAL to pick num. features
- > Justifying the accuracy of the estimate
 - > E.g., report TEST error

Bias-Variance Tradeoff

Optimal Prediction

Goal: Predict $Y \in \mathbb{R}$ **given** $X \in \mathbb{R}^d$ **if** $(X, Y) \sim P_{XY}$

Find function η that minimizes

$$\mathbb{E}_{XY} [(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X = x] \right]$$

(Hint: for any x , $\eta(x) = c_x$ where c_x minimizes $\mathbb{E}_{Y|X} [(Y - c_x)^2 | X = x]$)

Optimal Prediction

True distribution, don't actually have access to it

Goal: Predict $Y \in \mathbb{R}$ **given** $X \in \mathbb{R}^d$ **if** $(X, Y) \sim P_{XY}$

Find function η that minimizes

$$\mathbb{E}_{XY} [(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X = x] \right]$$

(Hint: for any x , $\eta(x) = c_x$ where c_x minimizes $\mathbb{E}_{Y|X} [(Y - c_x)^2 | X = x]$)

Optimal Prediction

True distribution, don't actually have access to it

Goal: Predict $Y \in \mathbb{R}$ **given** $X \in \mathbb{R}^d$ **if** $(X, Y) \sim P_{XY}$

Find function η that minimizes # Tower rule: $E[Y] = E_x[E[Y|X]]$

$$\mathbb{E}_{XY}[(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x] \right]$$

(Hint: for any x , $\eta(x) = c_x$ where c_x minimizes $\mathbb{E}_{Y|X}[(Y - c_x)^2 | X = x]$)

Optimal Prediction

True distribution, don't actually have access to it

Goal: Predict $Y \in \mathbb{R}$ **given** $X \in \mathbb{R}^d$ **if** $(X, Y) \sim P_{XY}$

Find function η that minimizes # Tower rule: $E[Y] = E_x[E[Y|X]]$

$$\mathbb{E}_{XY} [(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X = x] \right]$$

(Hint: for any x , $\eta(x) = c_x$ where c_x minimizes $\mathbb{E}_{Y|X} [(Y - c_x)^2 | X = x]$)

$$\frac{d}{dc} \mathbb{E}_{Y|X} [(Y - c)^2 | X = x] = 0$$

Optimal Prediction

True distribution, don't actually have access to it

Goal: Predict $Y \in \mathbb{R}$ **given** $X \in \mathbb{R}^d$ **if** $(X, Y) \sim P_{XY}$

Find function η that minimizes # Tower rule: $E[Y] = E_x[E[Y|X]]$

$$\mathbb{E}_{XY} [(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X = x] \right]$$

(Hint: for any x , $\eta(x) = c_x$ where c_x minimizes $\mathbb{E}_{Y|X} [(Y - c_x)^2 | X = x]$)

$$\frac{d}{dc} \mathbb{E}_{Y|X} [(Y - c)^2 | X = x] = 0$$

$$= \mathbb{E}_{Y|X} [2(Y - c)(-1) | X = x] = 0$$

Optimal Prediction

True distribution, don't actually have access to it

Goal: Predict $Y \in \mathbb{R}$ **given** $X \in \mathbb{R}^d$ **if** $(X, Y) \sim P_{XY}$

Find function η that minimizes # Tower rule: $\mathbb{E}[Y] = \mathbb{E}_x[\mathbb{E}[Y|X]]$

$$\mathbb{E}_{XY} [(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X = x] \right]$$

(Hint: for any x , $\eta(x) = c_x$ where c_x minimizes $\mathbb{E}_{Y|X} [(Y - c_x)^2 | X = x]$)

$$\frac{d}{dc} \mathbb{E}_{Y|X} [(Y - c)^2 | X = x] = 0$$

$$= \mathbb{E}_{Y|X} [2(Y - c)(-1) | X = x] = 0$$

$$= -2[\mathbb{E}_{Y|X} [Y | X = x] - c] = 0$$

Optimal Prediction

True distribution, don't actually have access to it

Goal: Predict $Y \in \mathbb{R}$ **given** $X \in \mathbb{R}^d$ **if** $(X, Y) \sim P_{XY}$

Find function η that minimizes # Tower rule: $\mathbb{E}[Y] = \mathbb{E}_x[\mathbb{E}[Y|X]]$

$$\mathbb{E}_{XY} [(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X = x] \right]$$

(Hint: for any x , $\eta(x) = c_x$ where c_x minimizes $\mathbb{E}_{Y|X} [(Y - c_x)^2 | X = x]$)

$$\frac{d}{dc} \mathbb{E}_{Y|X} [(Y - c)^2 | X = x] = 0$$

$$= \mathbb{E}_{Y|X} [2(Y - c)(-1) | X = x] = 0$$

$$= -2[\mathbb{E}_{Y|X} [Y | X = x] - c] = 0$$

$$c_x = \mathbb{E}_{Y|X} [Y | X = x] = \eta(x)$$

Optimal Prediction

True distribution, don't actually have access to it

Goal: Predict $Y \in \mathbb{R}$ **given** $X \in \mathbb{R}^d$ **if** $(X, Y) \sim P_{XY}$

Find function η that minimizes # Tower rule: $\mathbb{E}[Y] = \mathbb{E}_x[\mathbb{E}[Y|X]]$

$$\mathbb{E}_{XY} [(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X = x] \right]$$

(Hint: for any x , $\eta(x) = c_x$ where c_x minimizes $\mathbb{E}_{Y|X} [(Y - c_x)^2 | X = x]$)

$$\frac{d}{dc} \mathbb{E}_{Y|X} [(Y - c)^2 | X = x] = 0$$

$$= \mathbb{E}_{Y|X} [2(Y - c)(-1) | X = x] = 0$$

$$= -2[\mathbb{E}_{Y|X} [Y | X = x] - c] = 0$$

$$c_x = \mathbb{E}_{Y|X} [Y | X = x] = \eta(x)$$

Optimal Prediction

Goal: Predict $Y \in \mathbb{R}$ **given** $X \in \mathbb{R}^d$ **if** $(X, Y) \sim P_{XY}$

Find function η that minimizes

$$\mathbb{E}_{XY} [(Y - \eta(X))^2] = \mathbb{E}_X \left[\mathbb{E}_{Y|X} [(Y - \eta(x))^2 | X = x] \right]$$

(Hint: for any x , $\eta(x) = c_x$ where c_x minimizes $\mathbb{E}_{Y|X} [(Y - c_x)^2 | X = x]$)

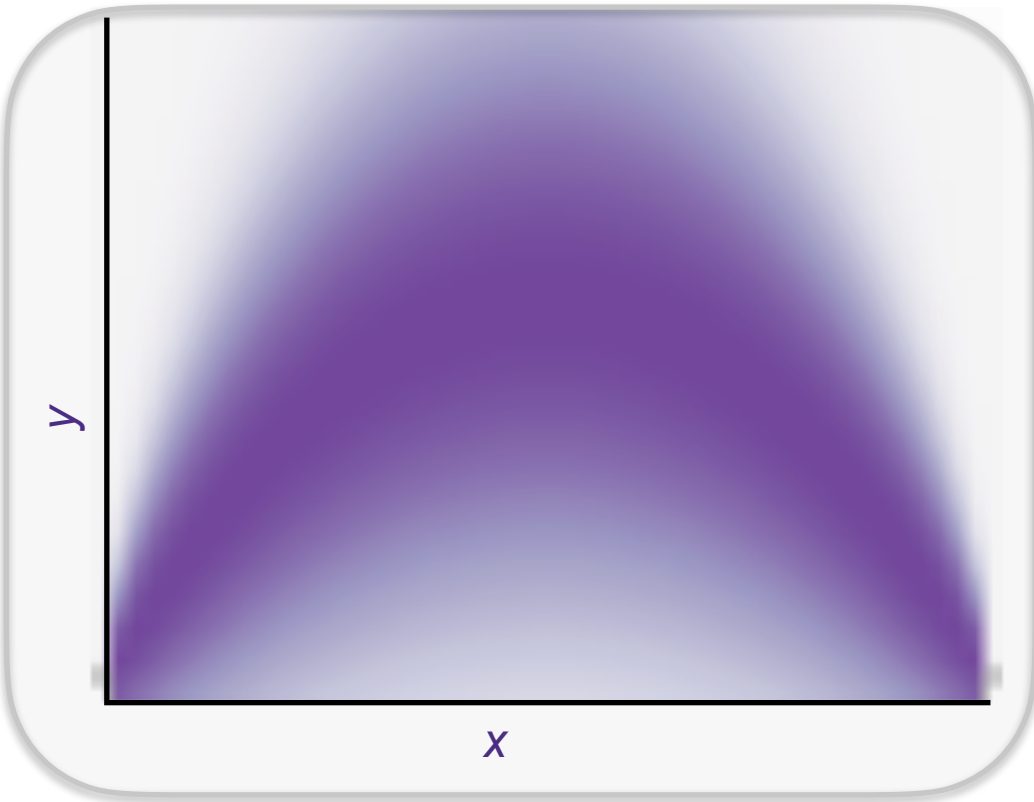
$$\begin{aligned} 0 &= \frac{d}{dc_x} \mathbb{E}_{Y|X} [(Y - c_x)^2 | X = x] \\ &= \mathbb{E}_{Y|X} \left[\frac{d}{dc_x} (Y - c_x)^2 | X = x \right] \\ &= \mathbb{E}_{Y|X} [-2(Y - c_x) | X = x] = -2\mathbb{E}_{Y|X} [Y | X = x] + 2c_x \end{aligned}$$

Squared Error Optimal Predictor: $\eta(x) = \mathbb{E}_{Y|X} [Y | X = x]$

Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

$$P_{XY}(X = x, Y = y)$$

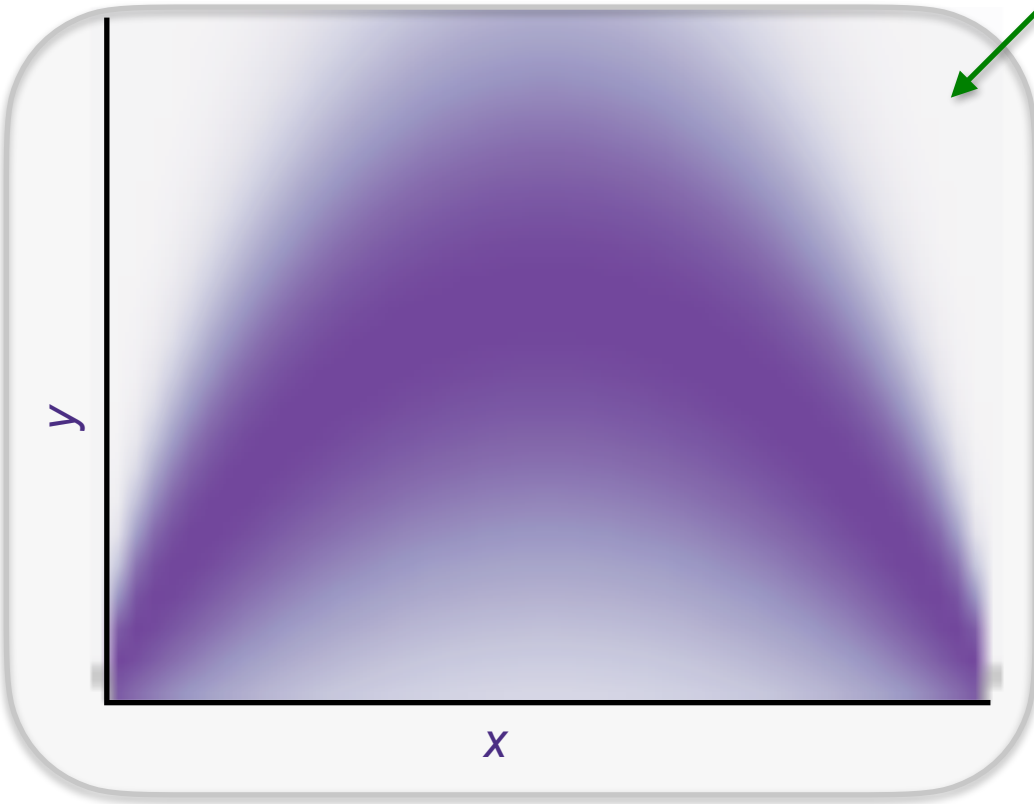


Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

$$P_{XY}(X = x, Y = y)$$

Low density

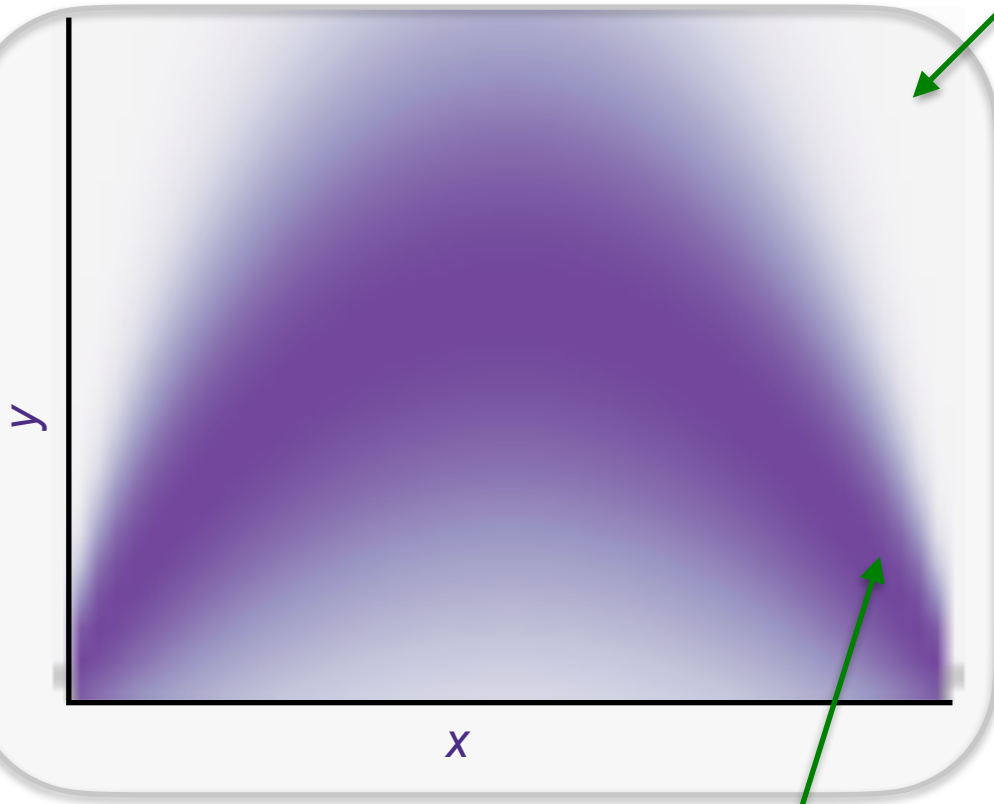


Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

$$P_{XY}(X = x, Y = y)$$

Low density

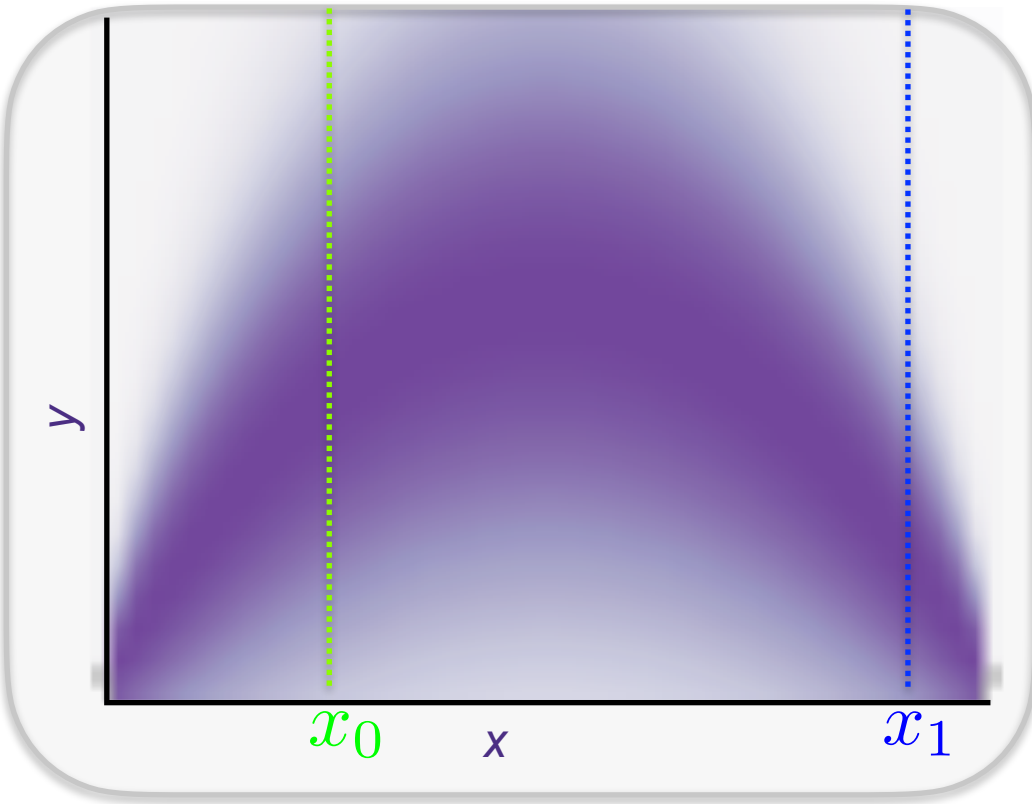


High density

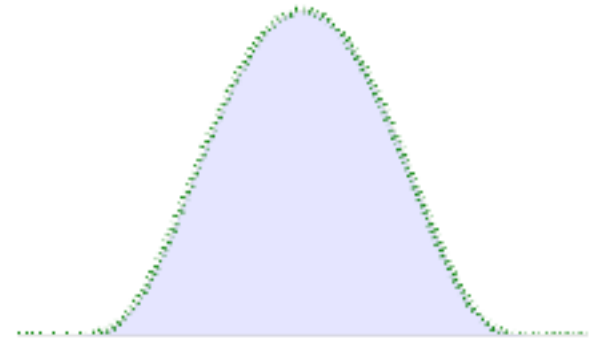
Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

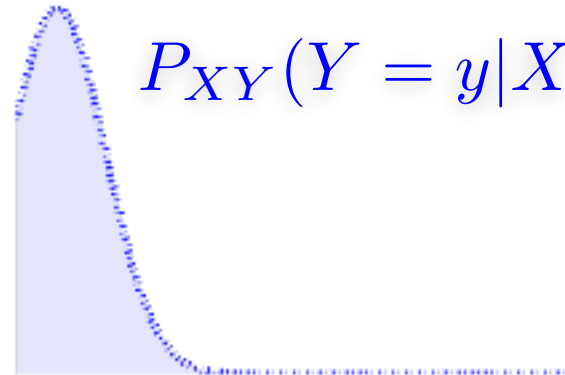
$$P_{XY}(X = x, Y = y)$$



$$P_{XY}(Y = y | X = x_0)$$



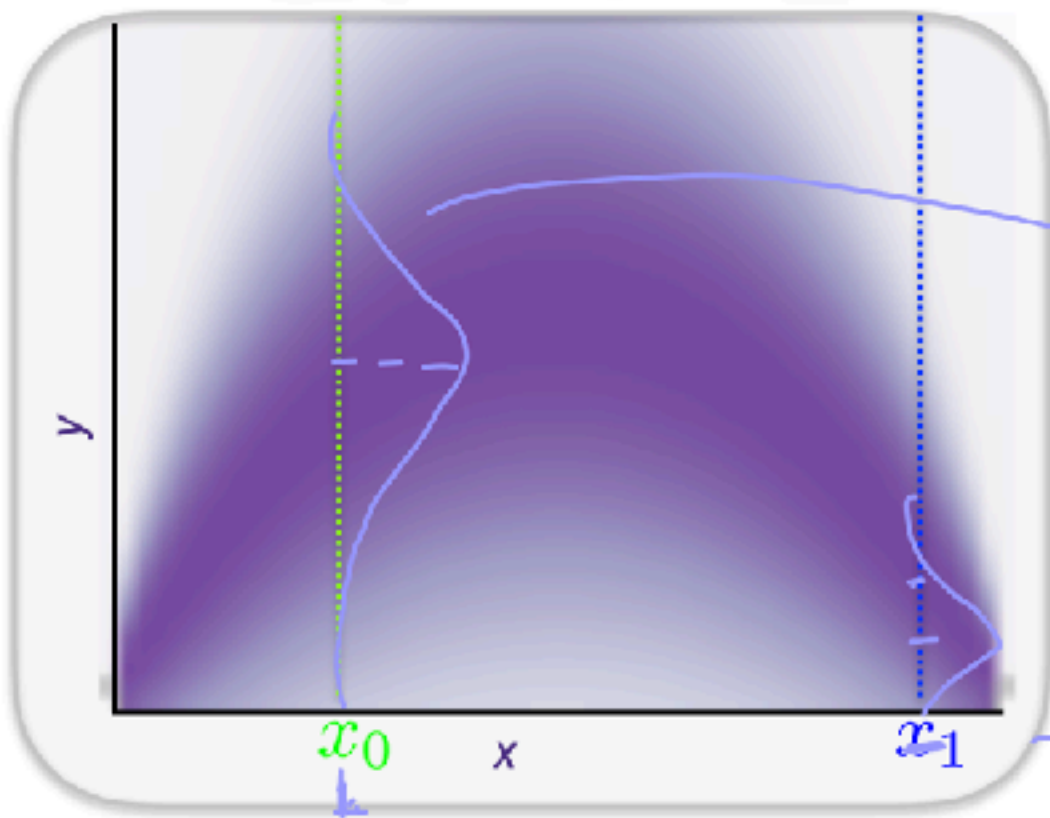
$$P_{XY}(Y = y | X = x_1)$$



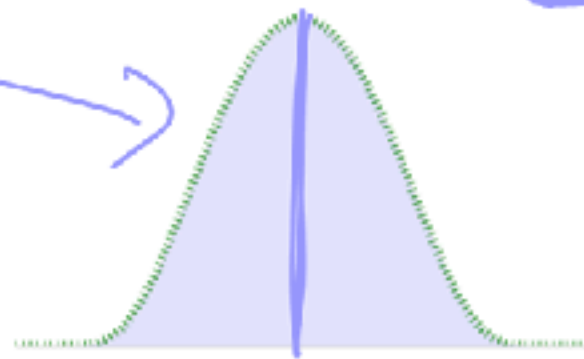
Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

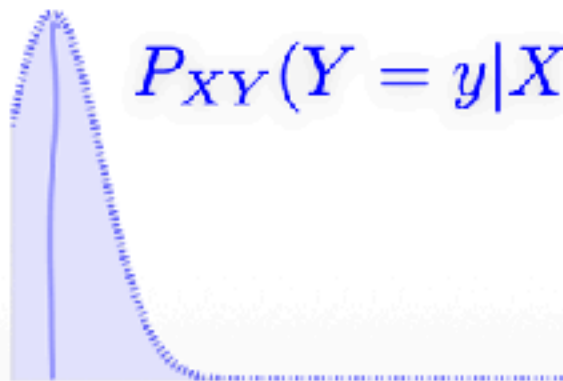
$$P_{XY}(X = x, Y = y)$$



$$P_{XY}(Y = y | X = \underline{x_0})$$



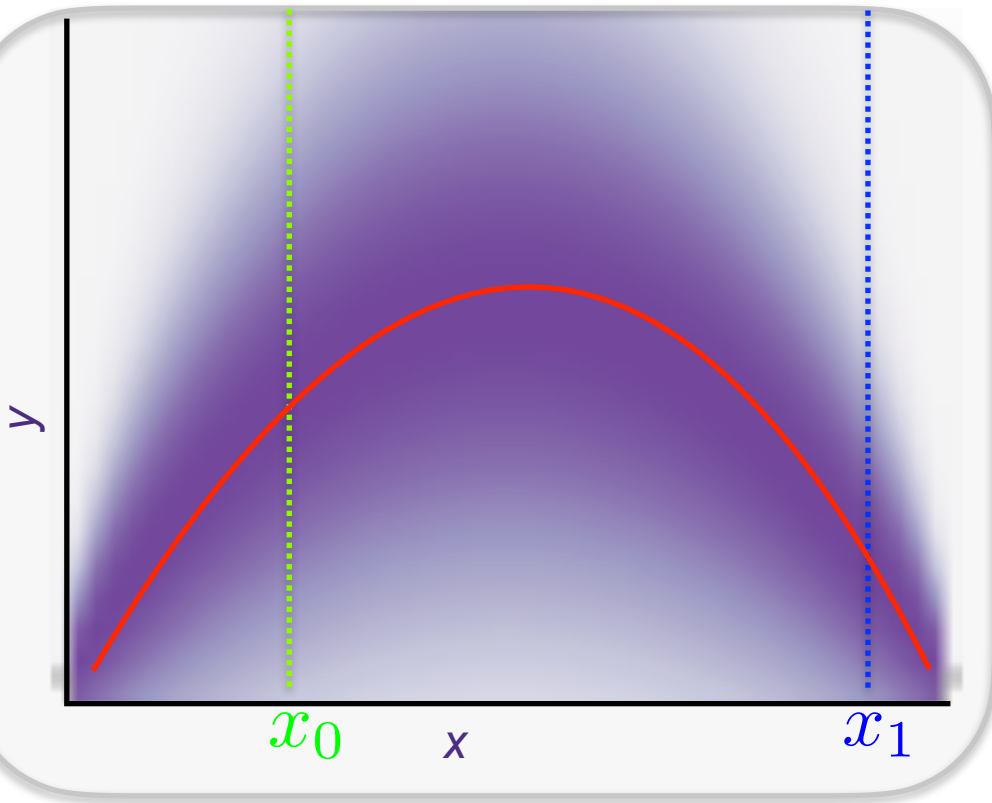
$$P_{XY}(Y = y | X = x_1)$$



Statistical Learning

$$\mathbb{E}_{XY}[(Y - \eta(X))^2]$$

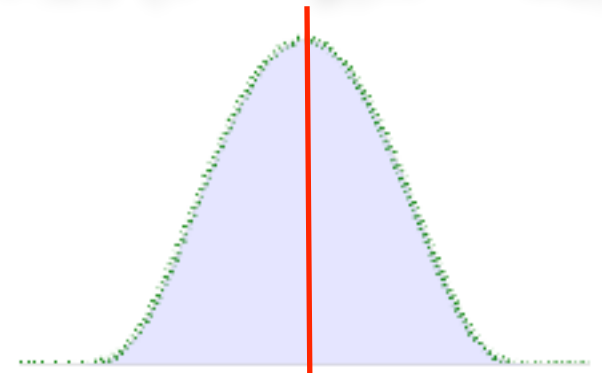
$$P_{XY}(X = x, Y = y)$$



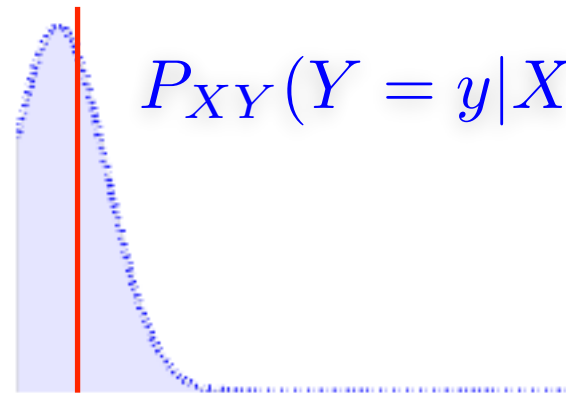
Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

$$P_{XY}(Y = y|X = x_0)$$

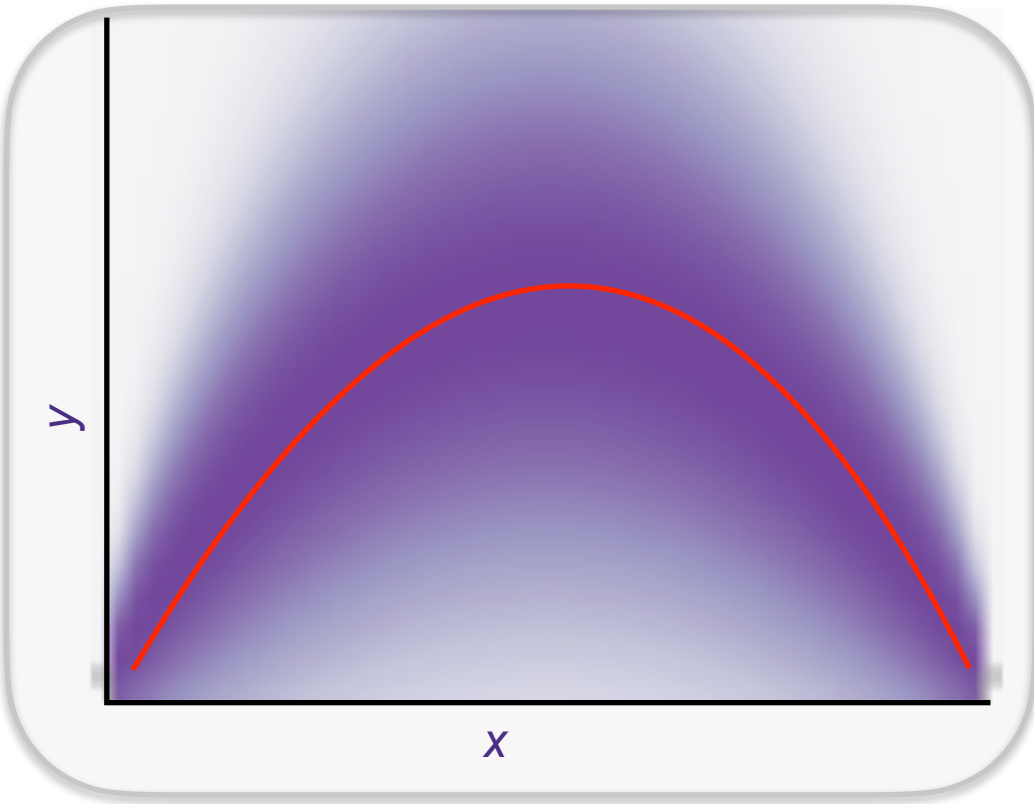


$$P_{XY}(Y = y|X = x_1)$$



Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

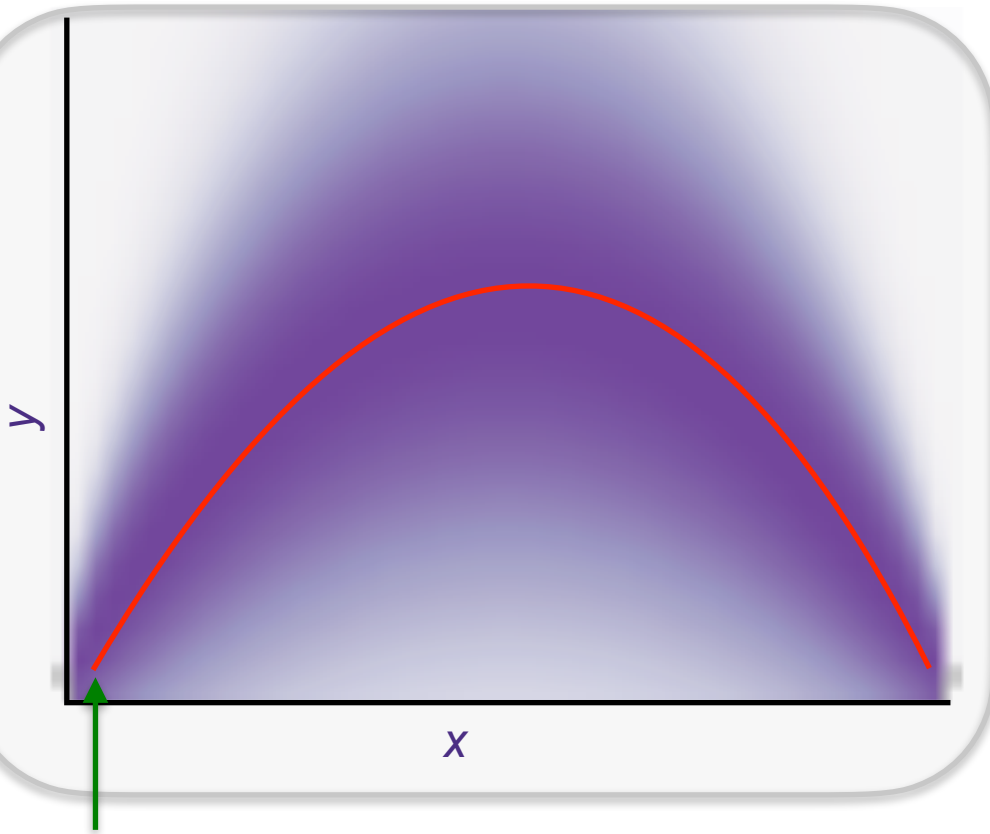
$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$

Ideally, we want to find:

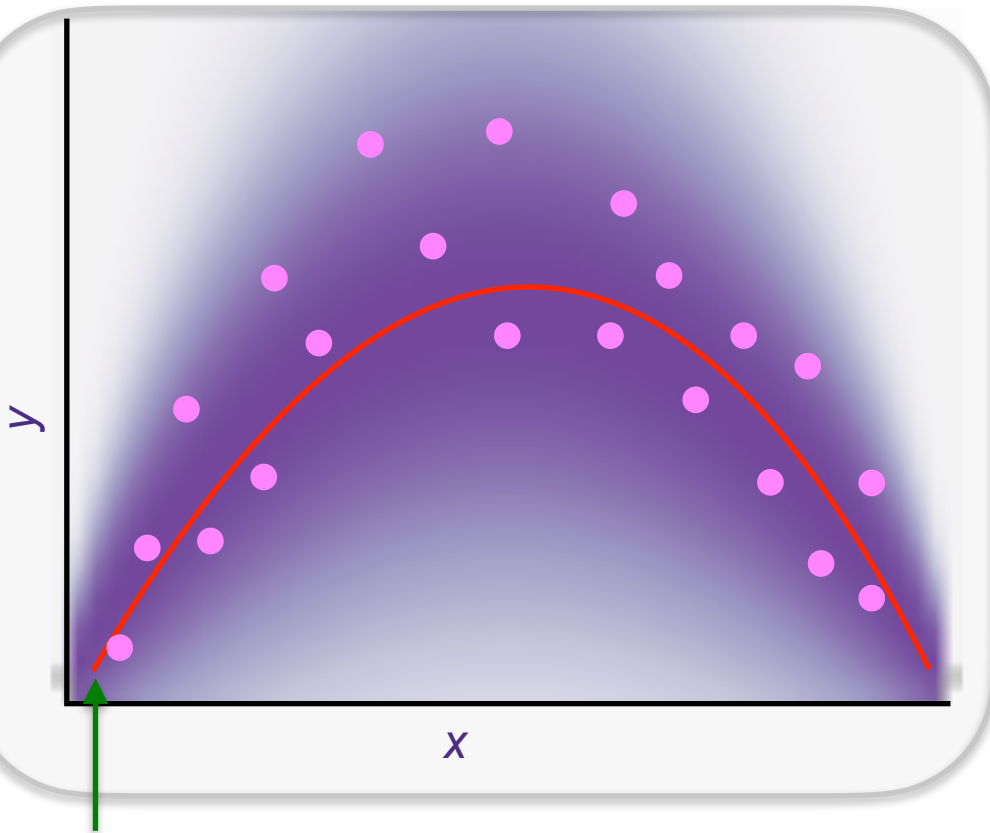
$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$



$\# \eta(x)$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



$\eta(x)$

Ideally, we want to find:

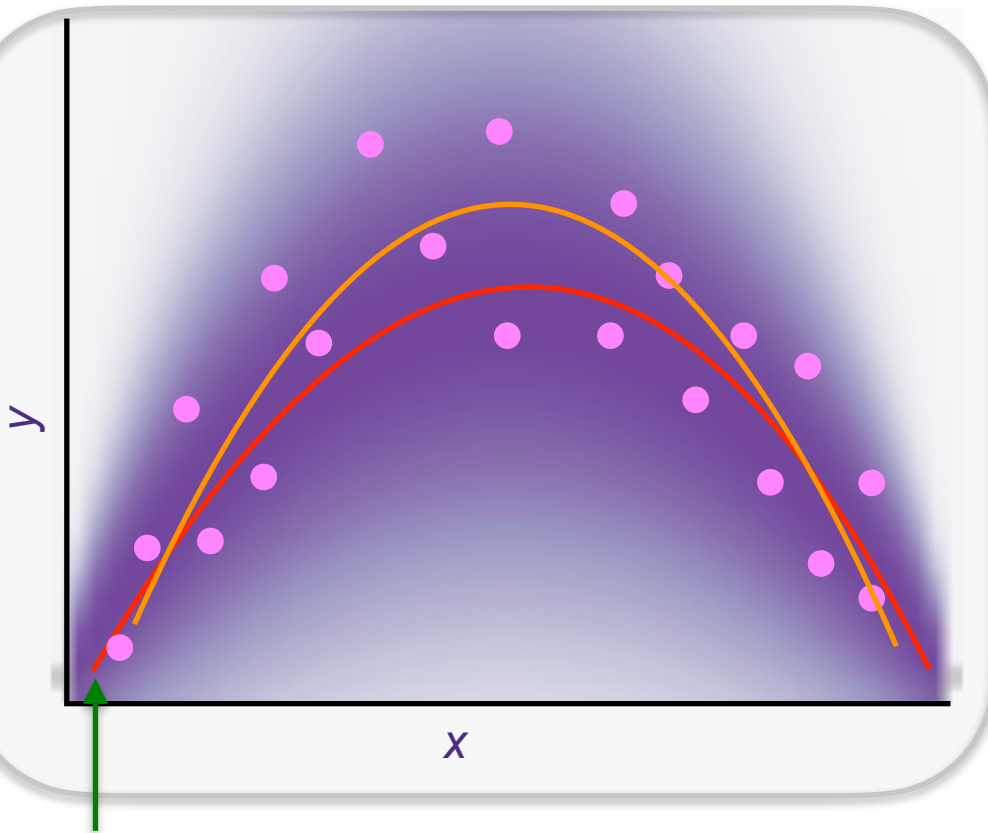
$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:

$$(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY} \quad \text{for } i = 1, \dots, n$$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



$\eta(x)$

Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

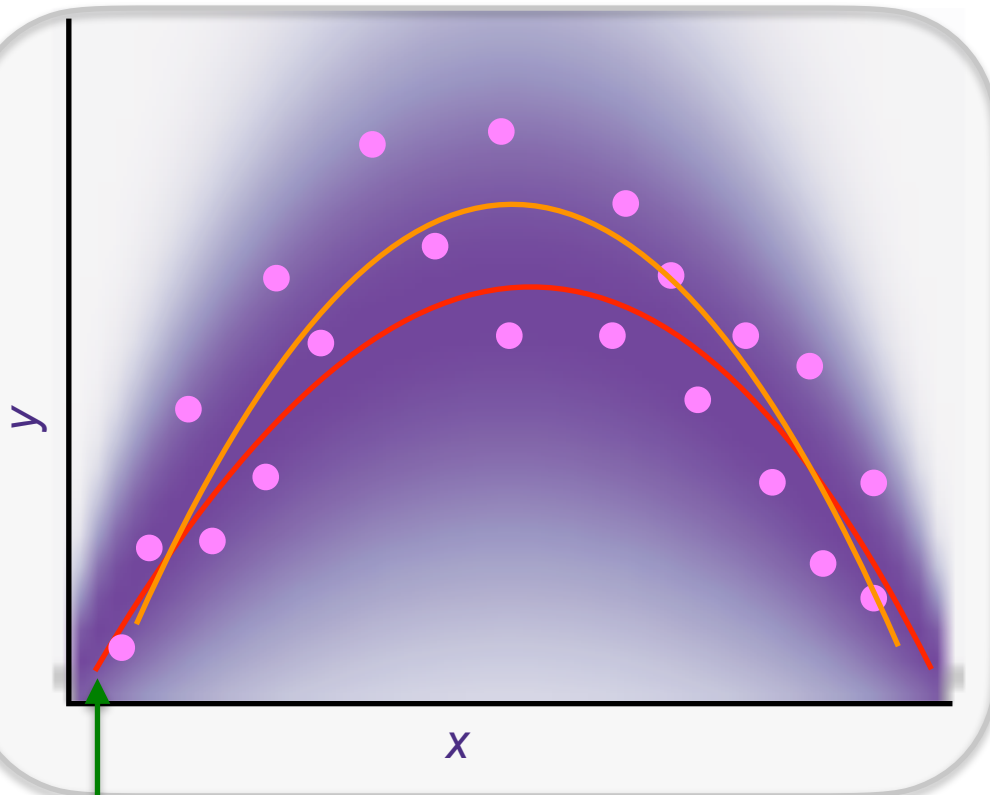
But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

and are restricted to a
function class (e.g., linear)
so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



$\eta(x)$

Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

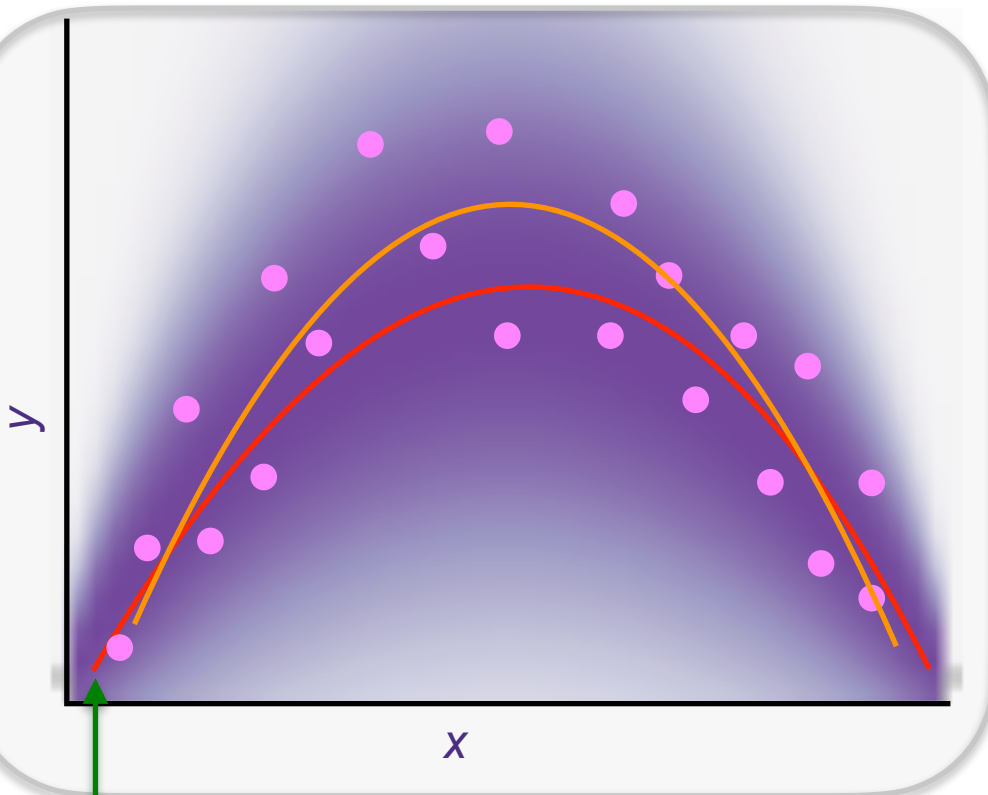
and are restricted to a
function class (e.g., linear)
so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$\hat{f} \neq \eta(x)$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



$\eta(x)$

Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

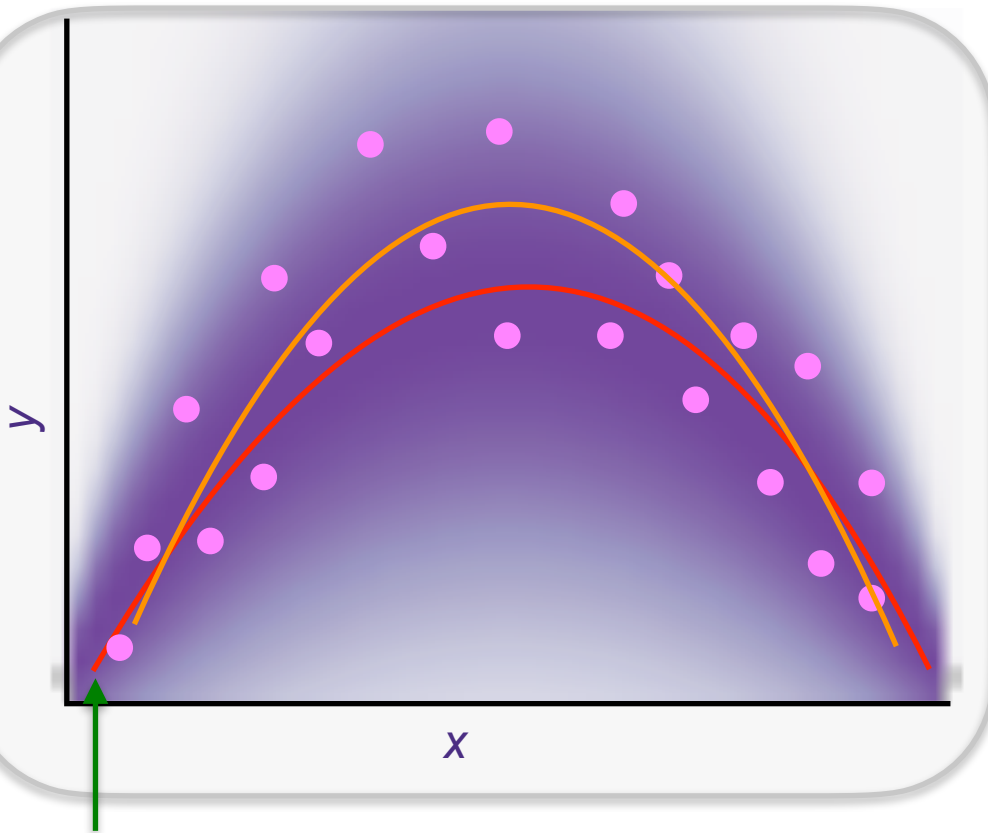
and are restricted to a
function class (e.g., linear)
so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$\hat{f} \neq \eta(x)$ Why?

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



$\eta(x)$

Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

and are restricted to a
function class (e.g., linear)
so we compute:

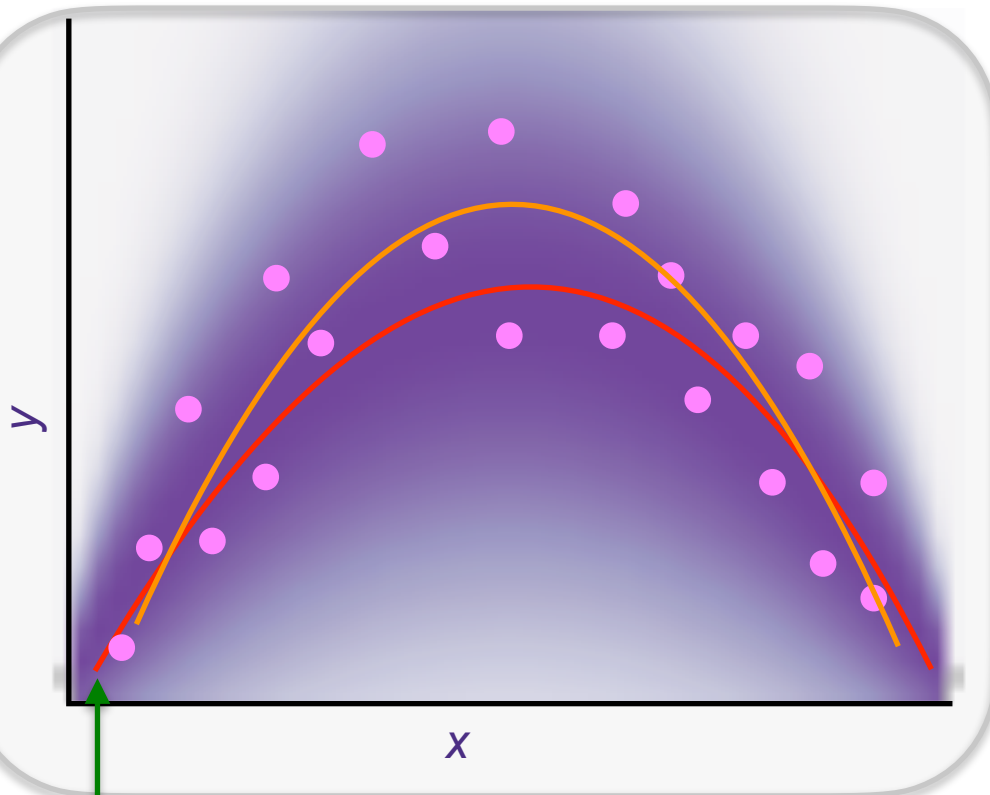
$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$\hat{f} \neq \eta(x)$ Why?

n is limited

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



$\eta(x)$

Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

and are restricted to a
function class (e.g., linear)
so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$\hat{f} \neq \eta(x)$ Why?

n is limited # could have chosen
wrong model class \mathcal{F}

Empirical Reconstruction Error (ERE)

Empirical Reconstruction Error (ERE)

What is the $\mathbb{E}[\hat{f}]$?

Empirical Reconstruction Error (ERE)

What is the $\mathbb{E}[\hat{f}]$?

$$\mathbb{E}_{XY} \left[\frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \right]$$

Empirical Reconstruction Error (ERE)

What is the $\mathbb{E}[\hat{f}]$?

$$\mathbb{E}_{XY}\left[\frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2\right]$$

By linearity of expectation

Empirical Reconstruction Error (ERE)

What is the $\mathbb{E}[\hat{f}]$?

$$\mathbb{E}_{XY}\left[\frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2\right]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y_i - f(x_i))^2] \quad \# \text{ By linearity of expectation}$$

Empirical Reconstruction Error (ERE)

What is the $\mathbb{E}[\hat{f}]$?

$$\mathbb{E}_{XY} \left[\frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y_i - f(x_i))^2]$$

By linearity of expectation

by IID: $\mathbb{E}[x_i]$ for any $i = \mathbb{E}[x] \forall x$

Empirical Reconstruction Error (ERE)

What is the $\mathbb{E}[\hat{f}]$?

$$\mathbb{E}_{XY} \left[\frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y_i - f(x_i))^2] \quad \# \text{ By linearity of expectation}$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y - f(x))^2] \quad \text{by IID: } \mathbb{E}[x_i] \text{ for any } i = \mathbb{E}[x] \forall x$$

Empirical Reconstruction Error (ERE)

What is the $\mathbb{E}[\hat{f}]$?

$$\mathbb{E}_{XY} \left[\frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y_i - f(x_i))^2]$$

By linearity of expectation

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y - f(x))^2]$$

by IID: $\mathbb{E}[x_i]$ for any $i = \mathbb{E}[x] \forall x$

Empirical Reconstruction Error (ERE)

What is the $\mathbb{E}[\hat{f}]$?

$$\mathbb{E}_{XY} \left[\frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y_i - f(x_i))^2] \quad \# \text{ By linearity of expectation}$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y - f(x))^2] \quad \text{by IID: } \mathbb{E}[x_i] \text{ for any } i = \mathbb{E}[x] \forall x$$

$$= \mathbb{E}[(y - f(x))^2]$$

Empirical Reconstruction Error (ERE)

What is the $\mathbb{E}[\hat{f}]$?

$$\mathbb{E}_{XY} \left[\frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y_i - f(x_i))^2] \quad \# \text{ By linearity of expectation}$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y - f(x))^2] \quad \text{by IID: } \mathbb{E}[x_i] \text{ for any } i = \mathbb{E}[x] \forall x$$

$$= \mathbb{E}[(y - f(x))^2] \quad \# \text{ Yet again, it works out to the mean squared error!}$$

Empirical Reconstruction Error (ERE)

What is the $\mathbb{E}[\hat{f}]$?

$$\mathbb{E}_{XY} \left[\frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y_i - f(x_i))^2] \quad \# \text{ By linearity of expectation}$$

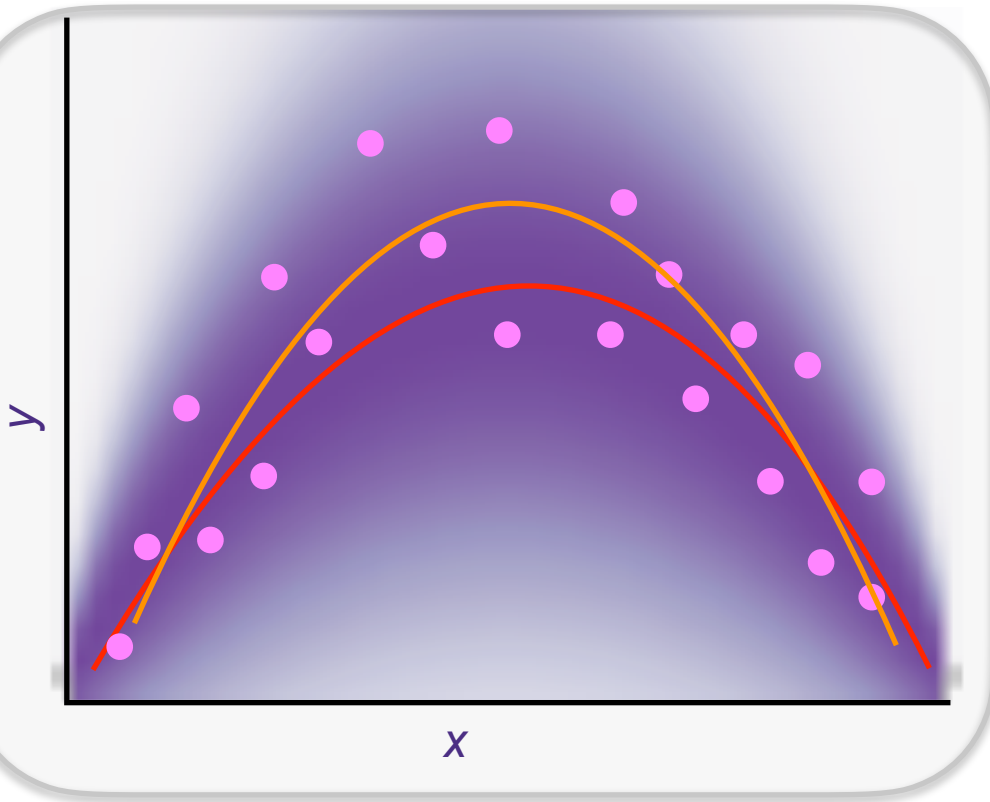
$$= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y - f(x))^2] \quad \text{by IID: } \mathbb{E}[x_i] \text{ for any } i = \mathbb{E}[x] \forall x$$

$$= \mathbb{E}[(y - f(x))^2] \quad \# \text{ Yet again, it works out to the mean squared error!}$$

Minimizing ERE \rightarrow minimizing squared error loss in general
(*in expectation....*)

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

and are restricted to a
function class (e.g., linear)
so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

We care about future predictions: $\mathbb{E}_{XY}[(Y - \hat{f}(X))^2]$

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

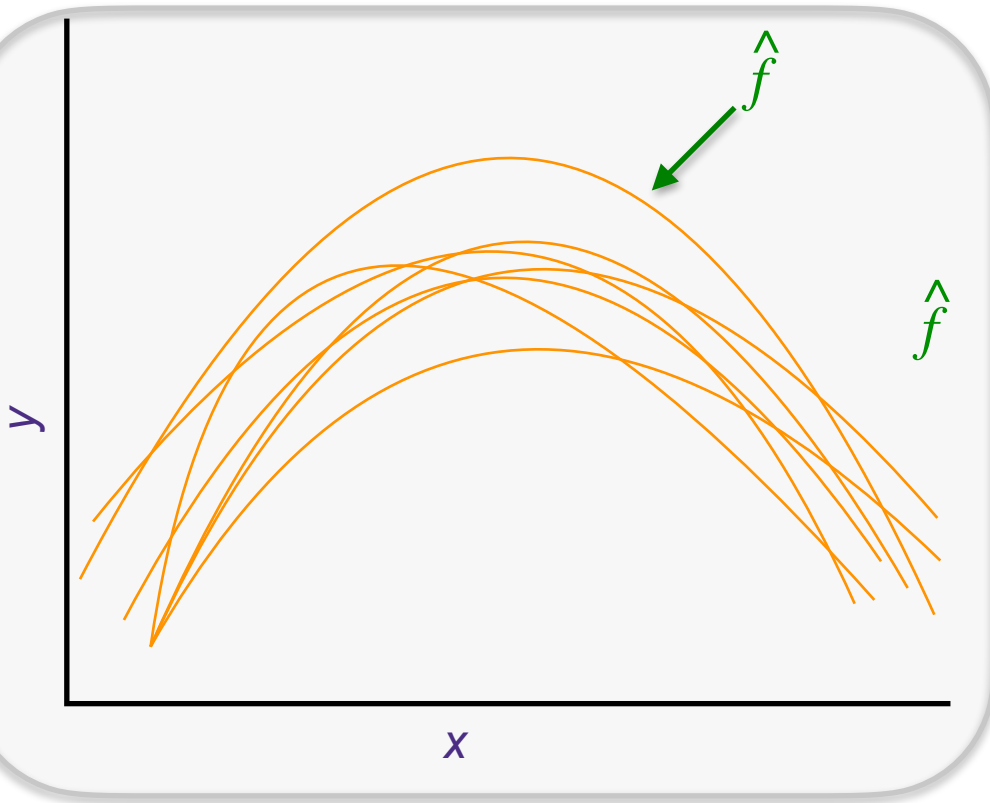
and are restricted to a
function class (e.g., linear)
so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

Each draw $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ results in different \hat{f}

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

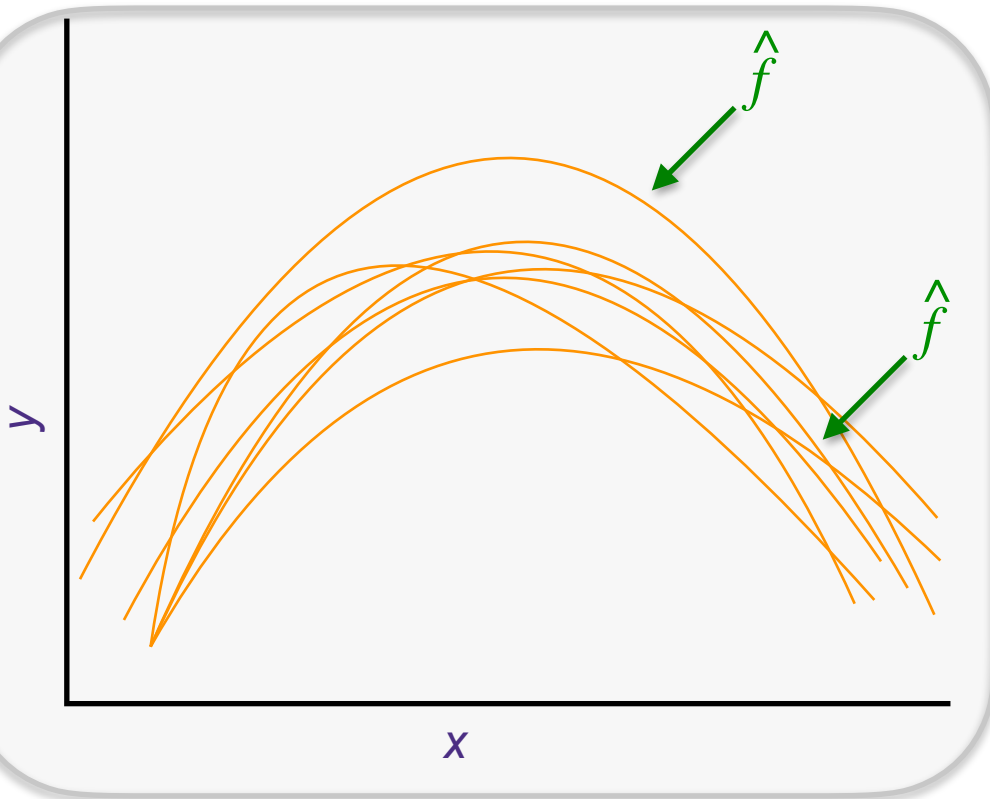
and are restricted to a
function class (e.g., linear)
so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

Each draw $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ results in different \hat{f}

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

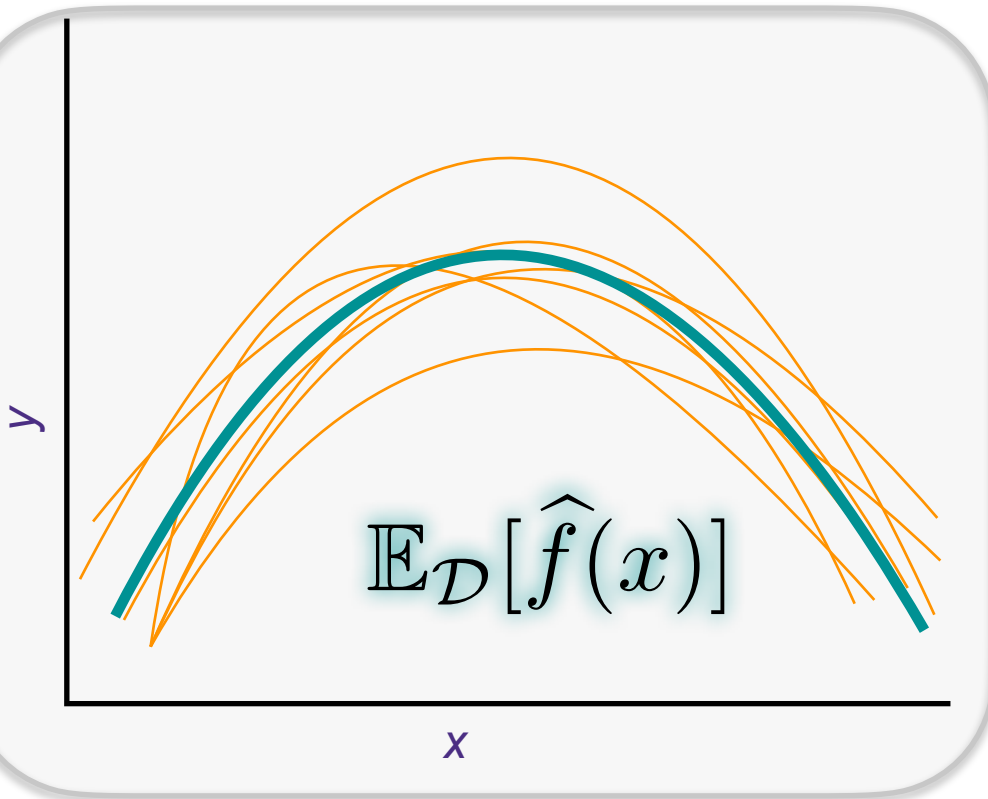
and are restricted to a function class (e.g., linear) so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

Each draw $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ results in different \hat{f}

Statistical Learning

$$P_{XY}(X = x, Y = y)$$



Ideally, we want to find:

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

But we only have samples:
 $(x_i, y_i) \stackrel{i.i.d.}{\sim} P_{XY}$ for $i = 1, \dots, n$

and are restricted to a function class (e.g., linear) so we compute:

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

Each draw $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ results in different \hat{f}

Bias-Variance Tradeoff

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x] \qquad \hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \hat{f}_{\mathcal{D}}(x))^2]|X = x] = \mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \eta(x) + \eta(x) - \hat{f}_{\mathcal{D}}(x))^2]|X = x]$$

Bias-Variance Tradeoff

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x] \quad \hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\begin{aligned} \mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \hat{f}_{\mathcal{D}}(x))^2]|X = x] &= \mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \eta(x) + \eta(x) - \hat{f}_{\mathcal{D}}(x))^2]|X = x] \\ &= \mathbb{E}_{Y|X} \left[\mathbb{E}_{\mathcal{D}}[(Y - \eta(x))^2 + 2(Y - \eta(x))(\eta(x) - \hat{f}_{\mathcal{D}}(x)) \right. \\ &\quad \left. + (\eta(x) - \hat{f}_{\mathcal{D}}(x))^2] | X = x \right] \\ &= \underbrace{\mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x]}_{\text{irreducible error}} + \underbrace{\mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{learning error}} \end{aligned}$$

irreducible error

Caused by stochastic label noise

learning error

Caused by either using too “simple” of a model or not enough data to learn the model accurately

Bias-Variance Tradeoff

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x] \qquad \hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\underline{\mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2]} = \mathbb{E}_{\mathcal{D}}[(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]$$

Bias-Variance Tradeoff

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x] \qquad \hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

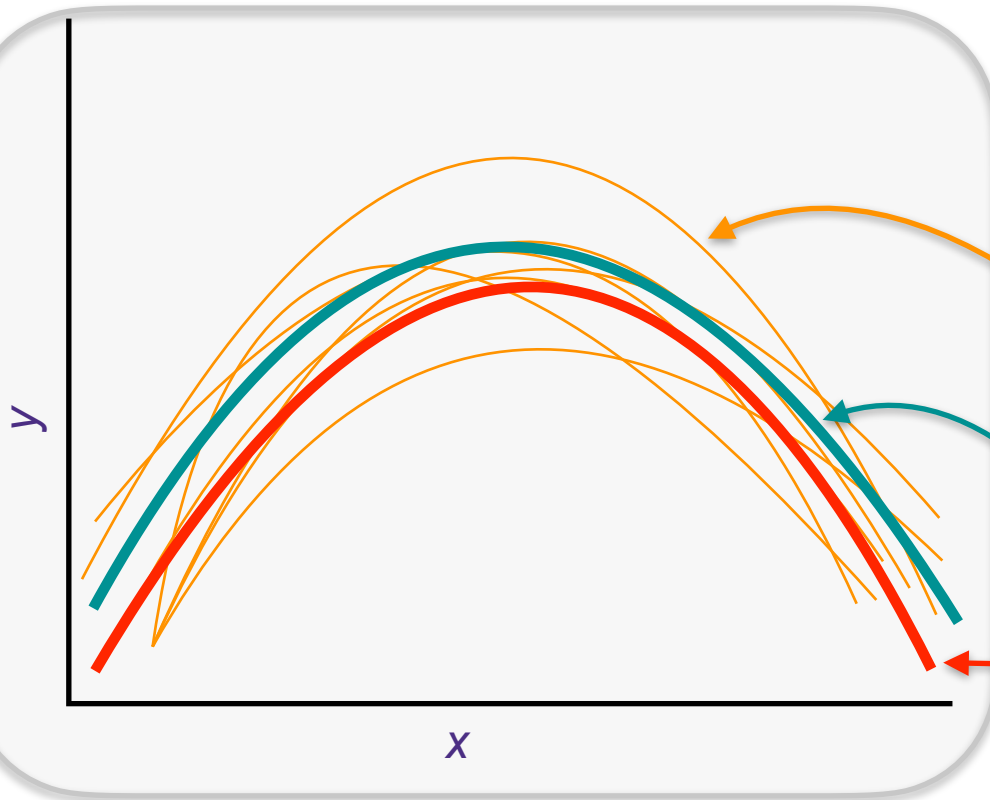
$$\begin{aligned} \underline{\mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2]} &= \mathbb{E}_{\mathcal{D}}[(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] + \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2] \\ &= \mathbb{E}_{\mathcal{D}}[(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2 + 2(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x)) \\ &\quad + (\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2] \\ &= \underline{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2} + \underline{\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]} \end{aligned}$$

biased squared

variance

Statistical Learning

$$\underbrace{\mathbb{E}_{\mathcal{D}}[(\eta(x) - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{learning error}} = \underbrace{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2}_{\text{biased squared}} + \underbrace{\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{variance}}$$



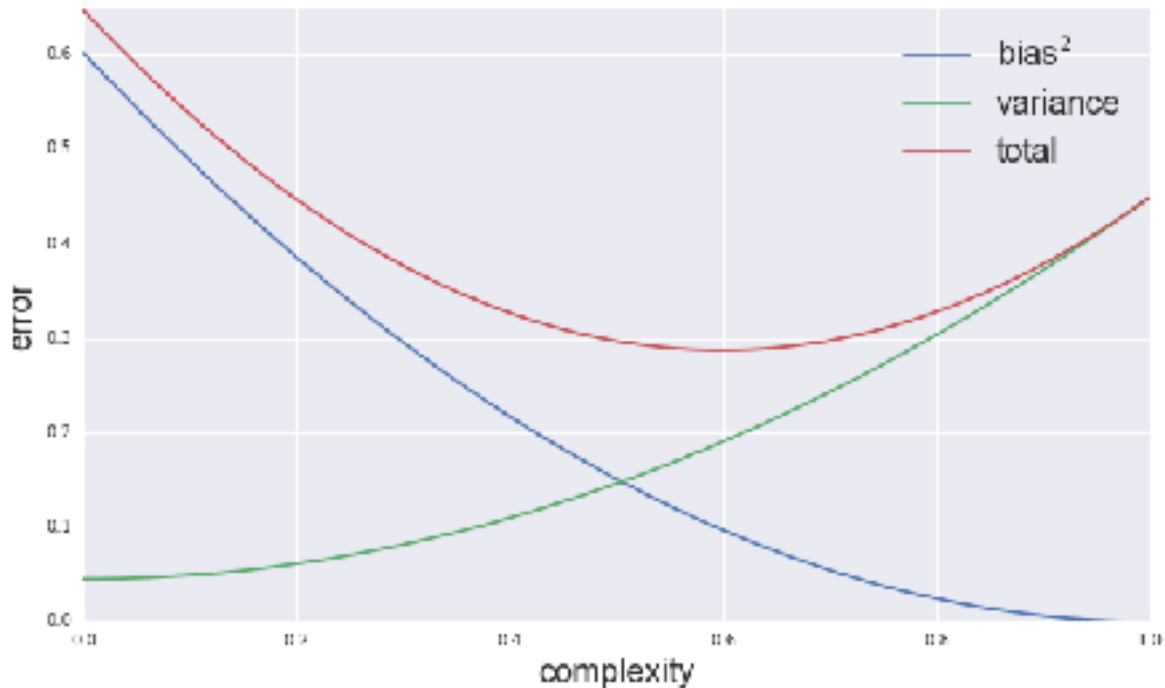
$$\hat{f} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\mathbb{E}_{\mathcal{D}}[\hat{f}(x)]$$

$$\eta(x) = \mathbb{E}_{Y|X}[Y|X = x]$$

Bias-Variance Tradeoff

$$\mathbb{E}_{Y|X}[\mathbb{E}_{\mathcal{D}}[(Y - \hat{f}_{\mathcal{D}}(x))^2] | X = x] = \underbrace{\mathbb{E}_{Y|X}[(Y - \eta(x))^2 | X = x]}_{\text{irreducible error}} + \underbrace{(\eta(x) - \mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)])^2}_{\text{biased squared}} + \underbrace{\mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[\hat{f}_{\mathcal{D}}(x)] - \hat{f}_{\mathcal{D}}(x))^2]}_{\text{variance}}$$



Bias-Variance Demo

See [colab notebook](#)

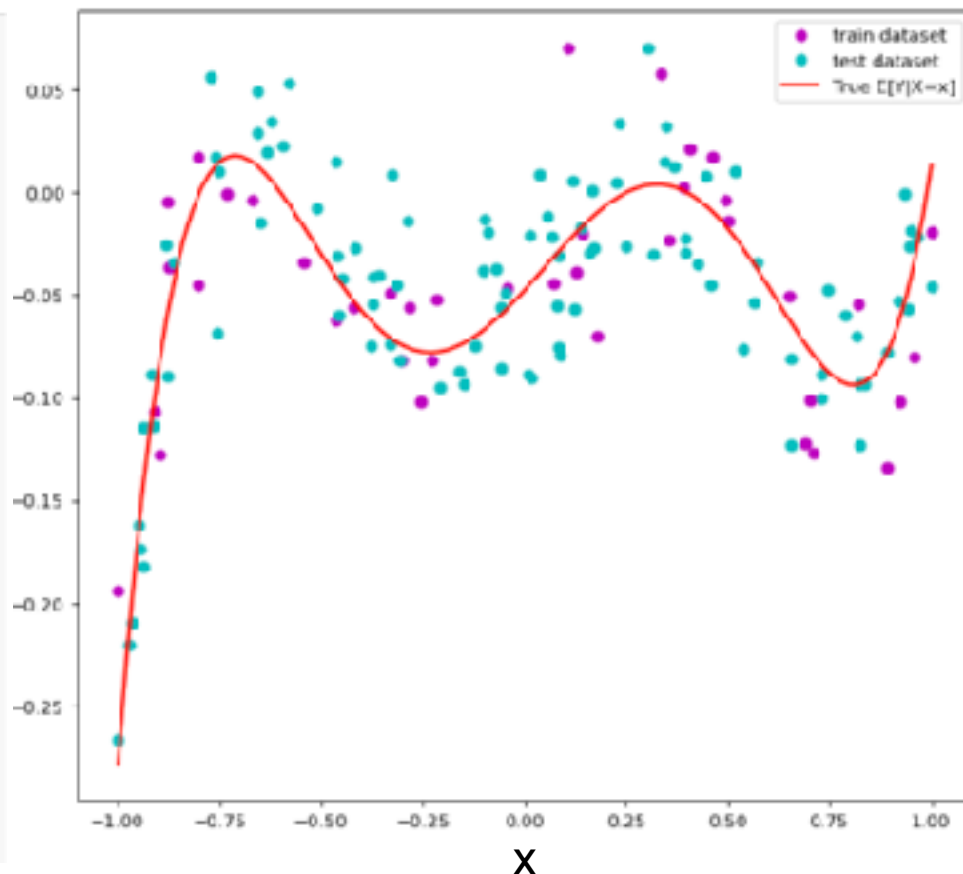
Bias-variance tradeoff

Let's define a ground truth $y(x) = E[Y|X=x]$ to be a degree-5 polynomial.

Then, define $P(Y|X=x) = \mathcal{N}(x, \sigma^2)$

```
In [1]: 1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # define eta = E[Y|X] as a 5-th order polynomial
5 # (this is the same function we used in demo_polynomial.ipynb)
6 eta = lambda x: (x-.99)*(x-.4)*(x-.25)*(x+.5)*(x+.8)
7
8 # Generate samples from P(Y,X)
9 def sample_Pxy(n, noise_sigma=0.03):
10
11     # Sample inputs x from P(x).
12     x = np.random.uniform(-1,1,n)
13     x = np.sort(x)
14     x[0]=-1
15     x[-1]=1
16
17     # Draw samples from P(Y|X)
18     y = eta(x) + noise_sigma*np.random.randn(n)
19
20     return x,y
21
22 # Generate training data.
23 n_train = 40 # sample size
24 x, y = sample_Pxy(n=n_train)
25
26 # Generate test data.
27 n_test = 100
28 X, y_ = sample_Pxy(n=n_test)
29
30 # Plot the samples and ground truth.
31 t = np.linspace(-1,1,100)
32 y0 = eta(t)
33 fig=plt.figure(figsize=(9, 8), dpi= 80, facecolor='w', edgecolor='k')
34 plt.plot(x,y,'mo',label='train dataset')
35 plt.plot(X,y_,'co',label='test dataset')
36 plt.plot(t,y0,'r-', linewidth=2, label='True E[Y|X=x]')
37 plt.legend()
38 plt.show()
39
```

y



X

In typical scenarios, we only have one set of samples, which we separate into S_{test} and S_{train}

- it is critical that those two sets do not overlap and the sets are chosen randomly to ensure that the test error is independent of the training error, and also the test samples are coming from the same distribution as the new samples that will come in the future

However, in order to understand how the test error behaves (theoretically), we consider the expected test error, and call it true error: i.e. $\mathcal{L}_{\text{true}} = E[\mathcal{L}_{\text{test}}]$

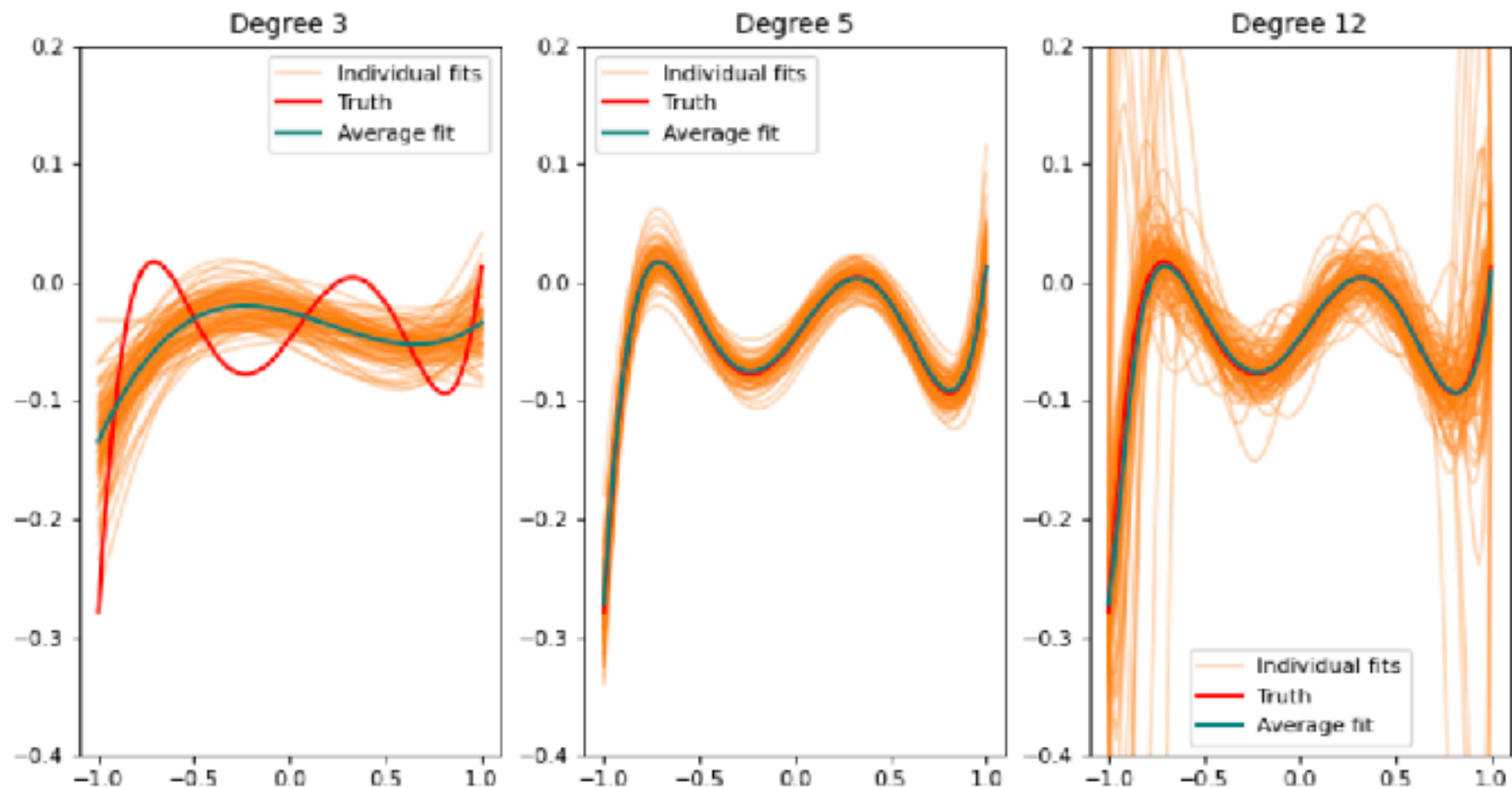
- test error is an **unbiased** estimate of the true error
- true error is unobservable (we cannot compute it, given finite samples)
- but we care the most about the true error
- so we use test error as a surrogate or an approximation

In order to compute the true error, we simulate a process where we get many fresh samples, and train new predictor each time with the fresh set of samples. It is important to understand that the resulting predictor $f_{S_{\text{train}}}(\cdot)$ is a random function, where the randomness comes from the fresh random training data set S_{train} . We will draw many such random functions, plot them, and see how test, train, true errors behave.

```

In [2]:
1 num_runs = 100
2 yhat_simple = np.zeros(num_runs, len(t))
3 yhat_justleft = np.zeros(num_runs, len(t))
4 yhat_justright = np.zeros(num_runs, len(t))
5
6 def poly_features(x, degree):
7     """
8     Generates a matrix where each column corresponds to x raised to a power from 0 to 'degree'.
9
10    Parameters:
11    x (numpy array): The input data.
12    degree (int): The maximum degree of the polynomial features.
13
14    Returns:
15    numpy array: A matrix where the columns are 1, x, x^2, ..., x^degree.
16    """
17    return np.vstack([x**i for i in range(degree + 1)])
18
19 run=0
20 while run < num_runs:
21     n = 40 # sample size
22     x,y = sample_data(n)
23
24     # degree=3 polynomial linear regression
25     p=3
26     X = poly_features(x, degree=p)
27     w = np.linalg.pinv(np.matmul(X.T,X)).dot(np.matmul(X.T,y))
28     T_ = poly_features(t, degree=p)
29     yhat_simple[int(run)] = np.matmul(T_,w)
30     y_ = np.matmul(X,w)
31     X_ = poly_features(x_, degree=p)
32     y_hat_ = np.matmul(X_,w)
33     w_ = w
34
35     # degree=5 polynomial linear regression
36     p=5
37     X = poly_features(x, degree=p)
38     w = np.linalg.pinv(np.matmul(X.T,X)).dot(np.matmul(X.T,y))
39     T_ = poly_features(t, degree=p)
40     yhat_justright[int(run)] = np.matmul(T_,w)
41     y_ = np.matmul(X,w)
42     X_ = poly_features(x_, degree=p)
43     y_hat_ = np.matmul(X_,w)
44     w_ = w
45
46     # degree=12 polynomial linear regression
47     p=12
48     X = poly_features(x, degree=p)
49     w = np.linalg.pinv(X.T.dot(X)).dot(X.T.dot(y))
50     T_ = poly_features(t, degree=p)
51     yhat_complex[int(run)] = np.matmul(T_,w)
52     y_ = np.matmul(X,w)
53     X_ = poly_features(x_, degree=p)
54     y_hat_ = np.matmul(X_,w)
55
56     run=run+1
57
58 def build_plot(yhat, titles):
59     ave_fit = np.zeros(np.size(t))
60     for i in range(1, num_runs):
61         plt.plot(t, yhat[int(i*run)], color='tab:orange', alpha=0.3)
62         ave_fit = ave_fit + yhat[int(i*run)]
63     plt.plot(t, yhat[0], 'tab:orange', alpha=0.3, label='Individual fits')
64     plt.plot(t, y, 'r', linestyle=2, label='Truth')
65     plt.plot(t, (1/num_runs)*ave_fit, color='tab:green', linestyle=2, label='Average fit')
66     plt.legend()
67     plt.title(titles)
68     axes = plt.gca()
69     axes.set_ylim([-0.4, 0.2])
70
71 fig=plt.figure(figsize=(12, 6), dpi=80, facecolor='w', edgecolor='k')
72 plt.subplot(1,3,1)
73 build_plot(yhat_simple, 'Degree 3')
74 plt.subplot(1,3,2)
75 build_plot(yhat_justright, 'Degree 5')
76 plt.subplot(1,3,3)
77 build_plot(yhat_complex, 'Degree 12')
78 plt.show()
79

```



Takeaways

- True function has degree 5, with additive noise
- Degree 3 fit has very high bias because the teal line, which is the average of the fits, is very different from true red line (because $3 < 5$). Each individual fit (orange) varies about the average fit (teal) moderately indicating moderate variance. The overall error ($\text{bias}^2 + \text{variance}$) of each individual fit is high.
- Degree 12 has very low bias because the teal line is almost equal to the red (because $12 > 5$). But each individual fit varies a lot about its average teal line indicating high variance. The overall error ($\text{bias}^2 + \text{variance}$) of each individual fit is high.
- Degree 5 has very low bias because the teal line is almost equal to the red (because $5 = 5$). And each individual fit varies only a little about its average teal line indicating small variance. The overall error ($\text{bias}^2 + \text{variance}$) of each individual fit is low.