

CSE 446

Linear Regression

Natasha Jaques

- *A ML workhorse*
- *Easy/fast to fit*
- *Interpretable*

Recap: Maximum Likelihood Estimation

Observe X_1, X_2, \dots, X_n drawn IID from $f(x; \theta)$ for some “true” $\theta = \theta_*$

Likelihood function $L_n(\theta) = \prod_{i=1}^n f(X_i; \theta)$

Log-Likelihood function $l_n(\theta) = \log(L_n(\theta)) = \sum_{i=1}^n \log(f(X_i; \theta))$

Maximum Likelihood Estimator (MLE) $\hat{\theta}_{MLE} = \arg \max_{\theta} L_n(\theta)$

Recap: MLE to learn Gaussian parameters

- MLE:

$$\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\sigma}^2_{MLE} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{MLE})^2$$

- MLE for the variance of a Gaussian is **biased**

$$\mathbb{E}[\hat{\sigma}^2_{MLE}] \neq \sigma^2$$

- Unbiased variance estimator:

$$\hat{\sigma}^2_{unbiased} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu}_{MLE})^2$$

Recap: Maximum Likelihood Estimation

ML “algorithms” we’ve learned so far:

- Maximum Likelihood Estimation (MLE)
 - Fit a Bernoulli distribution (coin flips)
 - Fit a Gaussian distribution (μ, σ)
- Now: fit a linear predictor $x \rightarrow y$, where y is a continuous variable

Maximum Likelihood Estimation

Observe X_1, X_2, \dots, X_n drawn IID from $f(x; \theta)$ for some “true” $\theta = \theta_*$

Likelihood function $L_n(\theta) = \prod_{i=1}^n f(X_i; \theta)$

prob of dataset given
model parameters

Log-Likelihood function $l_n(\theta) = \log(L_n(\theta)) = \sum_{i=1}^n \log(f(X_i; \theta))$

Maximum Likelihood Estimator (MLE) $\hat{\theta}_{MLE} = \arg \max_{\theta} L_n(\theta)$

- *We will use the same recipe as last time, to fit Linear Regression with MLE*
- *What part do we need to change?*

The MLE is a “recipe” that begins with a *model* for data $f(x; \theta)$

Regression



The regression problem, 1-dimensional

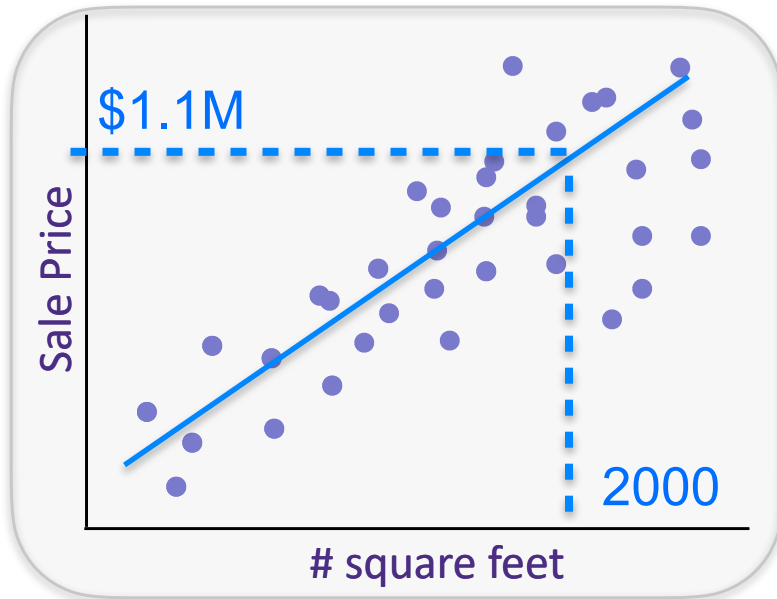
Given past sales data on [zillow.com](https://www.zillow.com), predict: # Goal: learn $p(y|x)$

y = House sale price *from*

x = {# sq. ft.}

now we have label y

previously we just had x



Training Data:

$$\{(x_i, y_i)\}_{i=1}^n$$

$$x_i \in \mathbb{R}$$

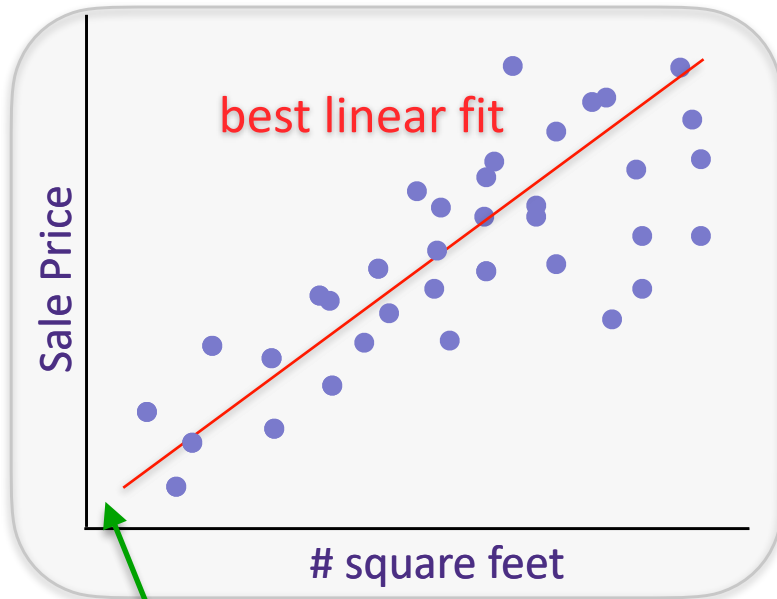
$$y_i \in \mathbb{R}$$

Fit a function to our data, 1-d

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

$y =$ House sale price *from*

$x =$ {# sq. ft.}



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

Hypothesis/Model: linear

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

$w \in \mathbb{R}$ slope of line

- why the noise?

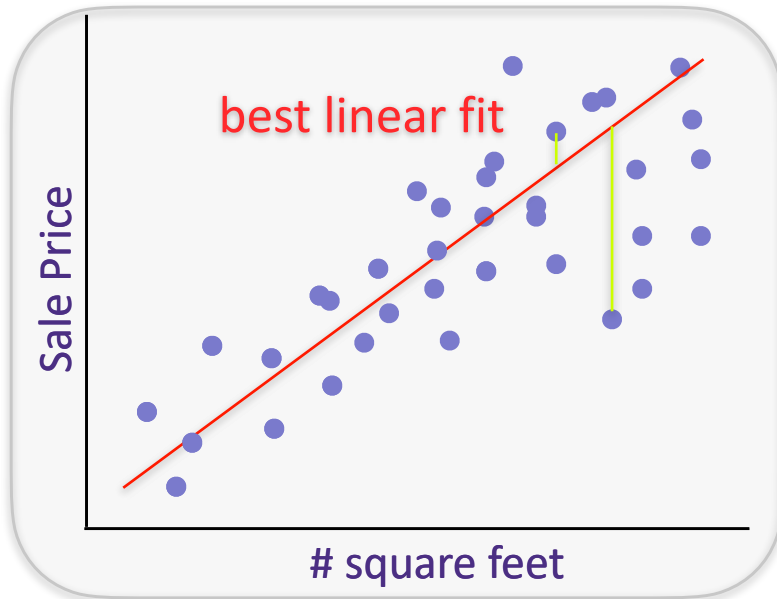
(no intercept for now to make the math easier)

Fit a function to our data, 1-d

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

y = House sale price *from*

x = {# sq. ft.}



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

Hypothesis/Model: linear

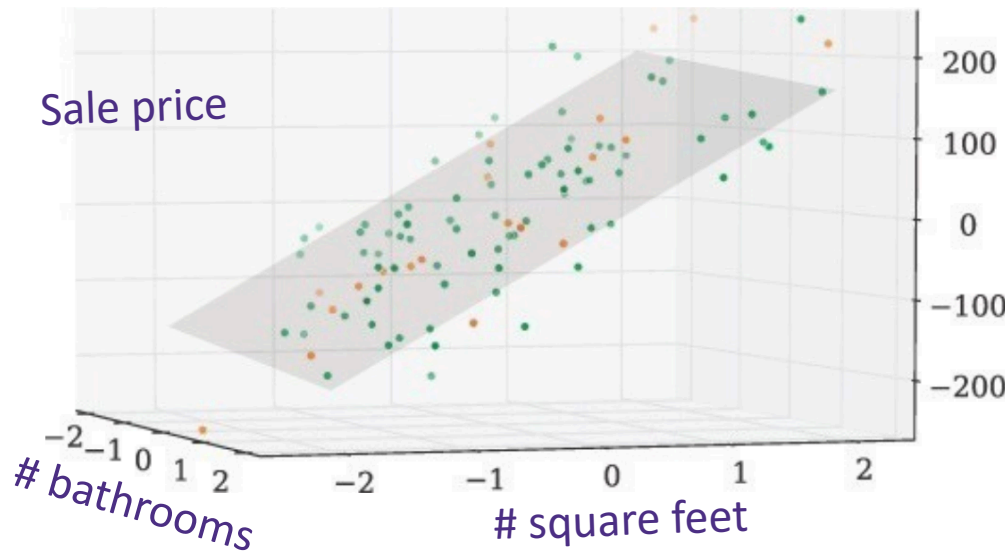
$$y_i = x_i w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

The regression problem, d-dim

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

y = House sale price from

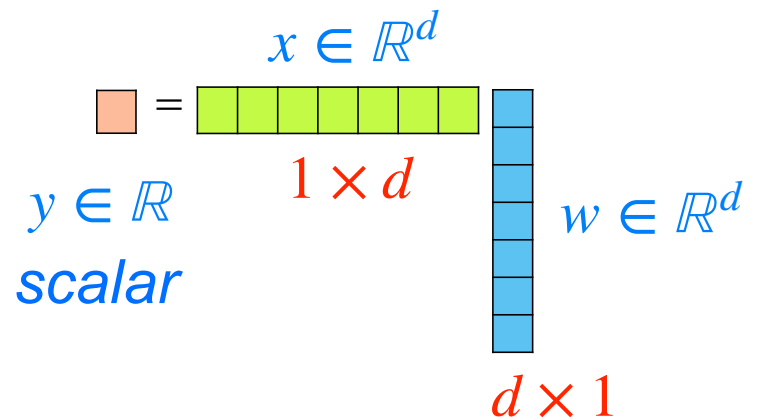
$x = \{\# \text{ sq. ft.}, \# \text{ baths}, \text{date of sale, etc.}\} \quad x \in \mathbb{R}^d$



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

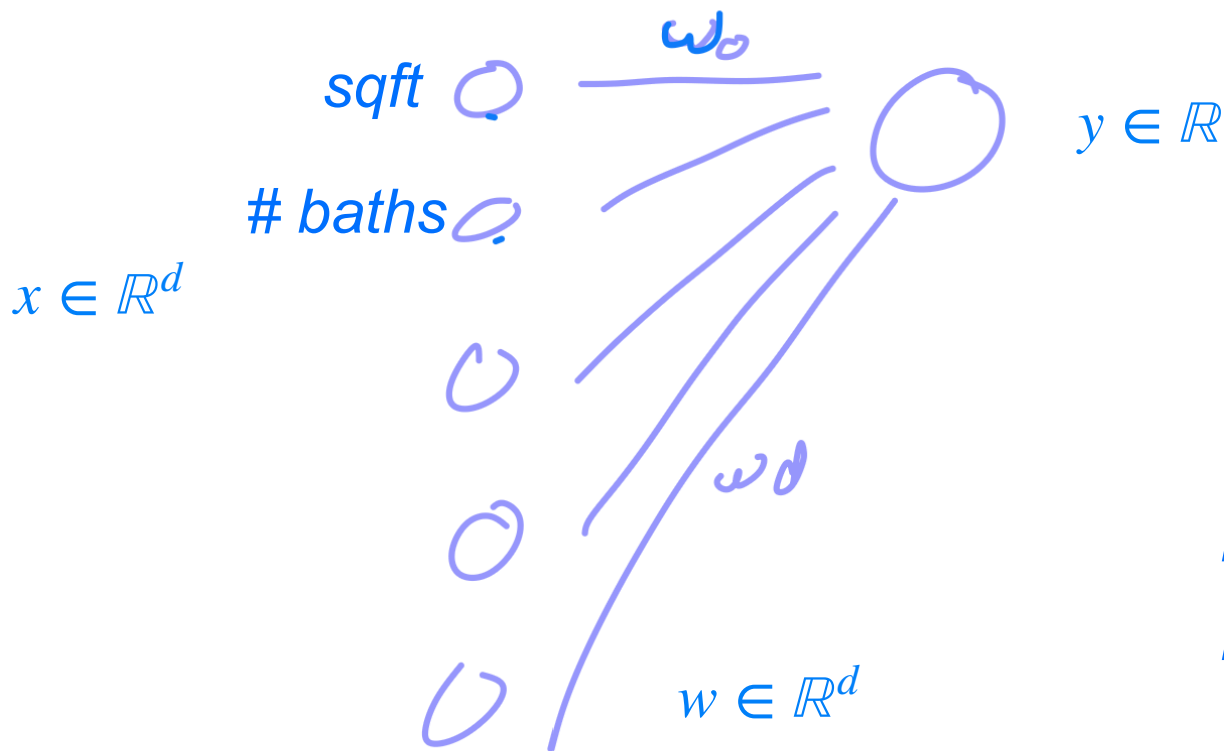
Hypothesis/Model: linear

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



Linear regression with d features

What simple model is this equivalent to?



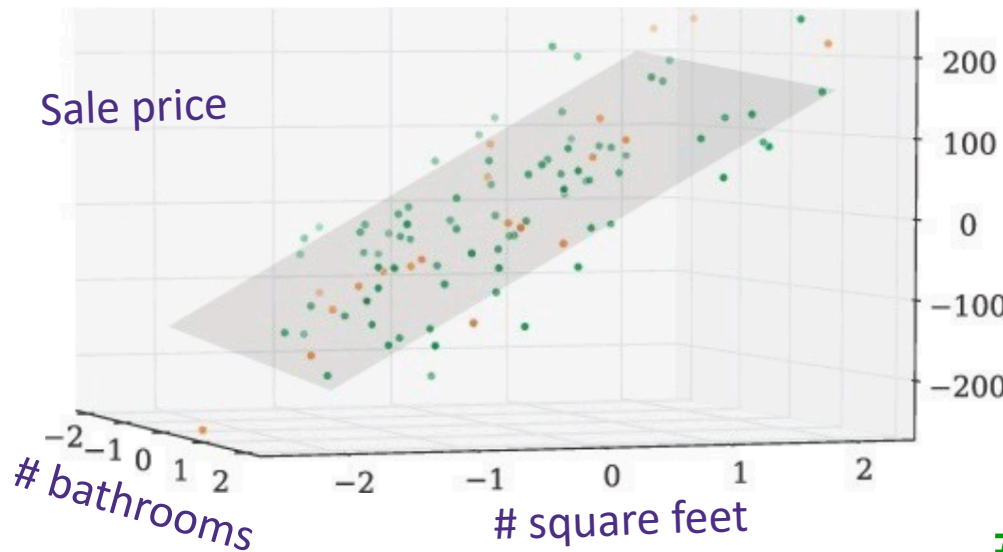
**A 1 layer
neural
network!**

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

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

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

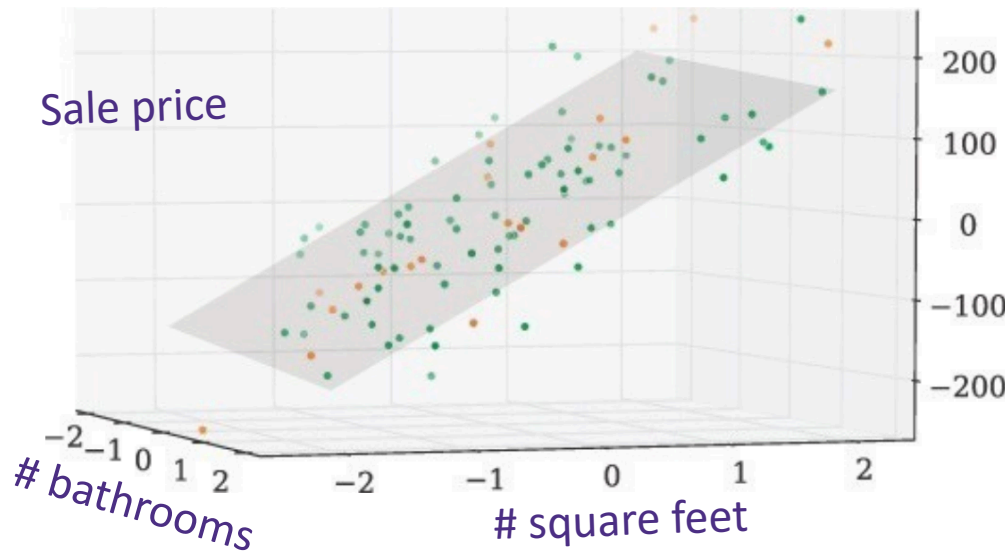
go back to formula sheet,
retrieve PDF of Gaussian...

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

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

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

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

Maximizing log-likelihood

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

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

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

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

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

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!
 - 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

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

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

"Matrix Cookbook" identities (cheatsheet fodder)

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

Maximizing log-likelihood

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

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

A matrix

How can we isolate w?

Take the inverse!

Maximizing log-likelihood

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

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

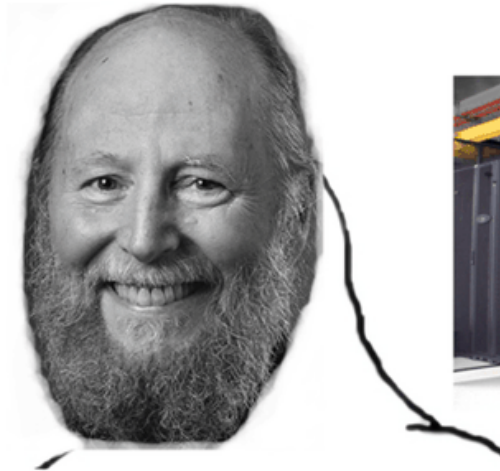
Set gradient=0, solve for w

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

The regression problem in matrix notation

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

We ❤️ matrix notation!



haha gpus go brrrrr

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

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

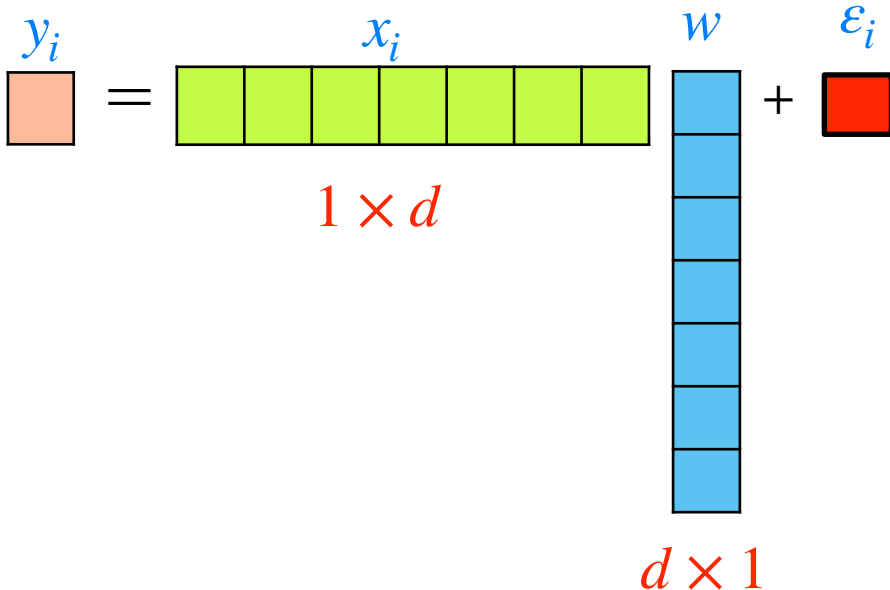
The regression problem in matrix notation

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

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

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

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



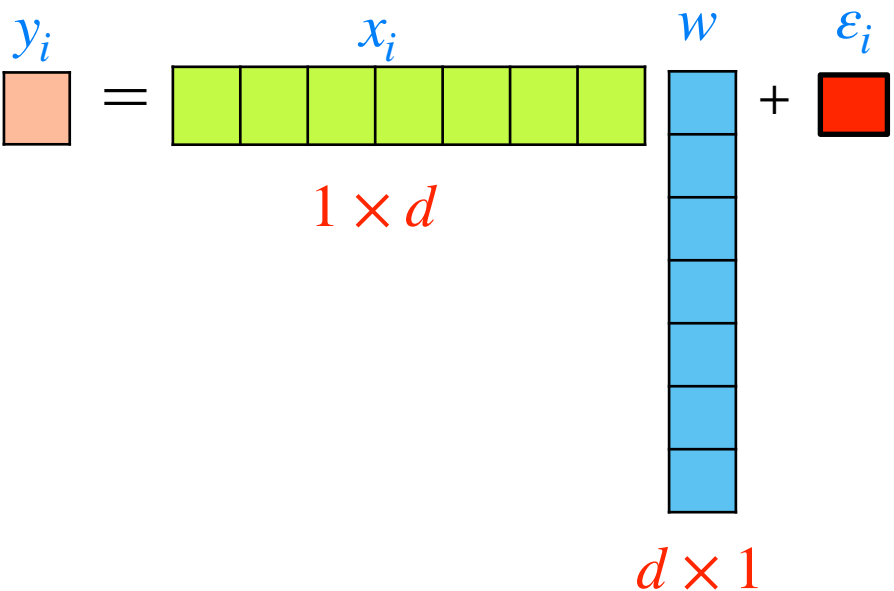
The regression problem in matrix notation

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

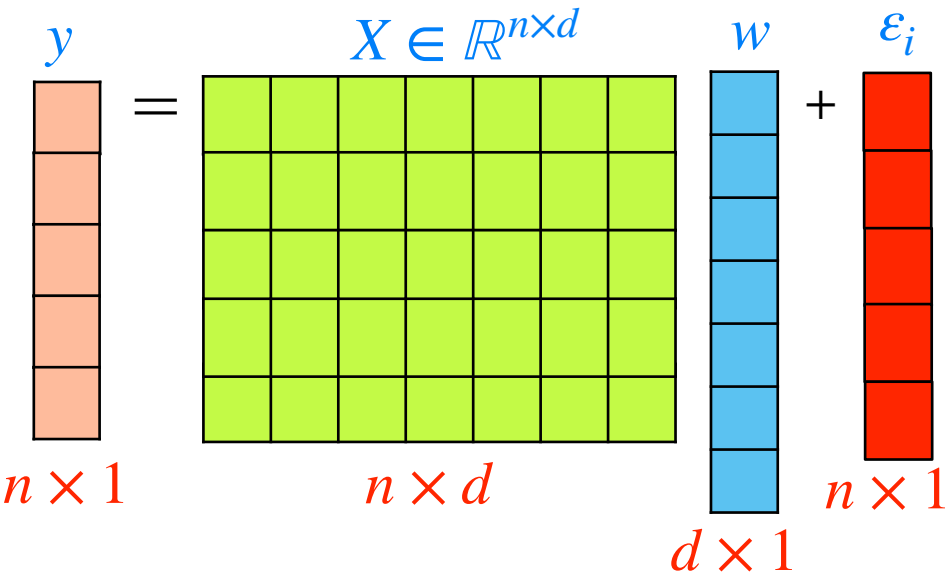
$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix}$$

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

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

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

d : # of features

n : # of examples/datapoints

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

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

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

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

The regression problem 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?

The regression problem in matrix notation

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

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

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

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

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

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

The regression problem in matrix notation

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

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

d : # of features

n : # of examples/datapoints

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

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

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

Why LS? “Ordinary Least Squares”

The regression problem in matrix notation



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

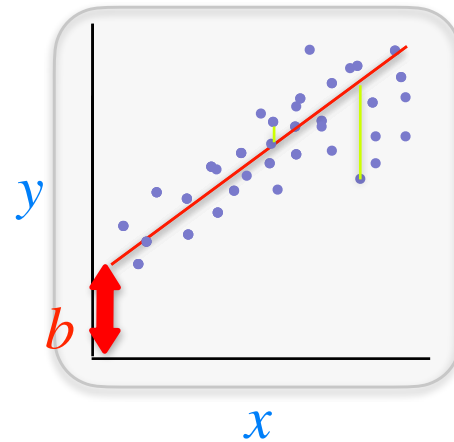


$$\hat{w}_{MLE} = \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

The regression problem in matrix notation

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

What about an offset?



$$\hat{y}_i = x_i^T w + b$$

Prediction for a single point

$$\left(\hat{y}_i - y_i\right)^2$$

Still want to minimize squared error

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

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

$n \times 1$

$b \in \mathbb{R}$

Dealing with an offset

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

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

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

Could think of this as system of equations...

$$\begin{bmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}^T \mathbf{1} \\ \mathbf{1}^T \mathbf{X} & \mathbf{1}^T \mathbf{1} \end{bmatrix} \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} \mathbf{X}^T \mathbf{y} \\ \mathbf{1}^T \mathbf{y} \end{bmatrix}$$

$$\begin{bmatrix} w \\ b \end{bmatrix} = \text{np.linalg.solve}() \quad \# \text{ There's a better way...}$$

Dealing with an offset

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

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

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

$$\tilde{x}_i = \begin{bmatrix} x_i \\ 1 \end{bmatrix} \rightarrow \text{Scalar} \in \mathbb{R}$$

$$\tilde{\omega} = \begin{bmatrix} w \\ b \end{bmatrix}$$

$$\hat{y}_i = \tilde{x}_i^T \tilde{\omega}$$

$$\tilde{X} = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_n^T \end{bmatrix} = \begin{bmatrix} X & \mathbf{1} \end{bmatrix}$$

$$\hat{\mathbf{y}} = \tilde{X}^T \tilde{\omega}$$

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

Same equation on modified input matrix containing all 1s feature

Dealing with an offset

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

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

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

It's good practice to pre-process your features to have a mean of zero

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

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

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

Make Predictions

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

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

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

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

How to normalize features to have mean 0?

For each feature, calc the mean **of the training data** and subtract it from each data point

So what do I need to do to x_{new} to be able to make a prediction?

Subtract the mean **of the training data** for each feature

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