

CSE 446

Linear Regression

Natasha Jaques



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

Learning 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:

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

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

Under benign assumptions, as the number of observations $n \rightarrow \infty$ we have $\hat{\theta}_{MLE} \rightarrow \theta_*$

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

Regression

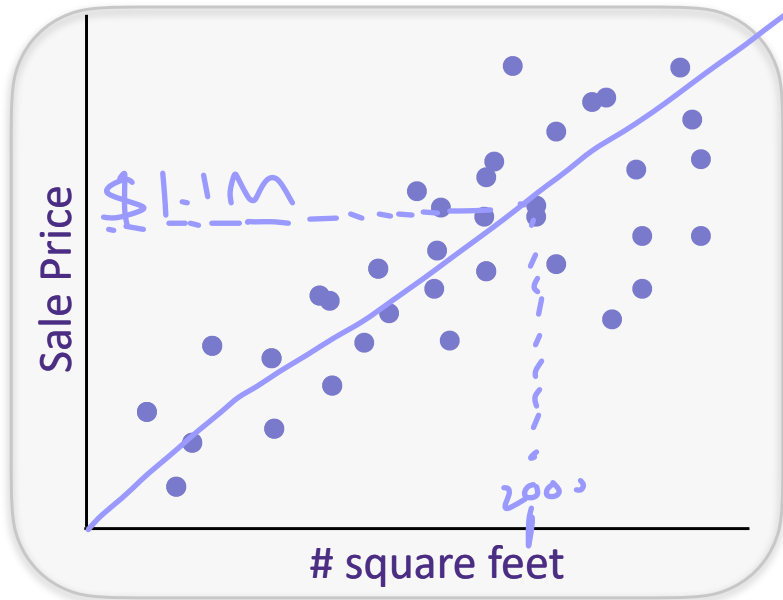


The regression problem, 1-dimensional

Given past sales data on zillow.com, predict:

$y =$ House sale price from _ ←

$x =$ {# sq. ft.}



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

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

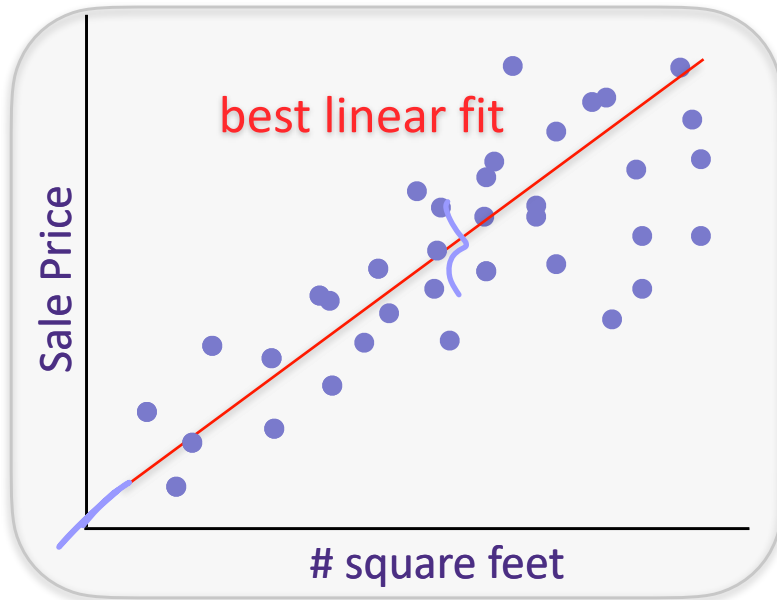
$$y_i \in \mathbb{R}$$

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

Hypothesis/Model: linear \rightarrow

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

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

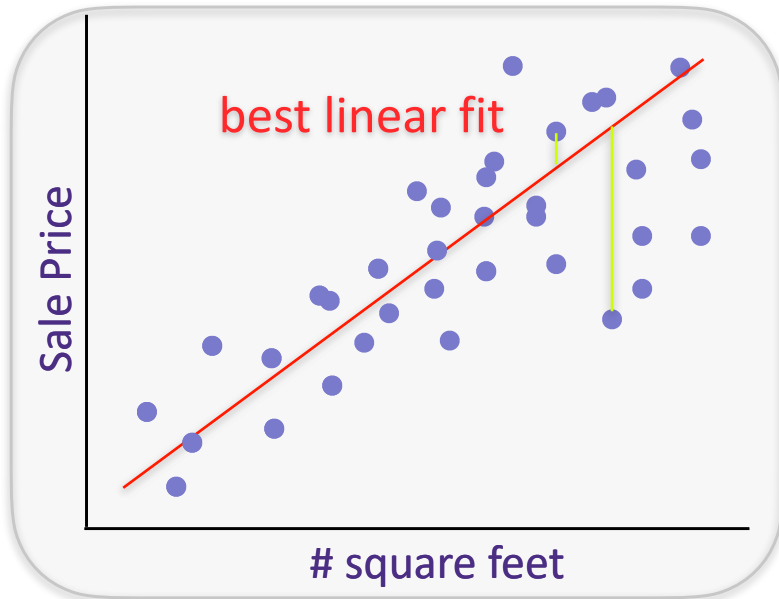
\rightarrow no intercept to make math easy

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

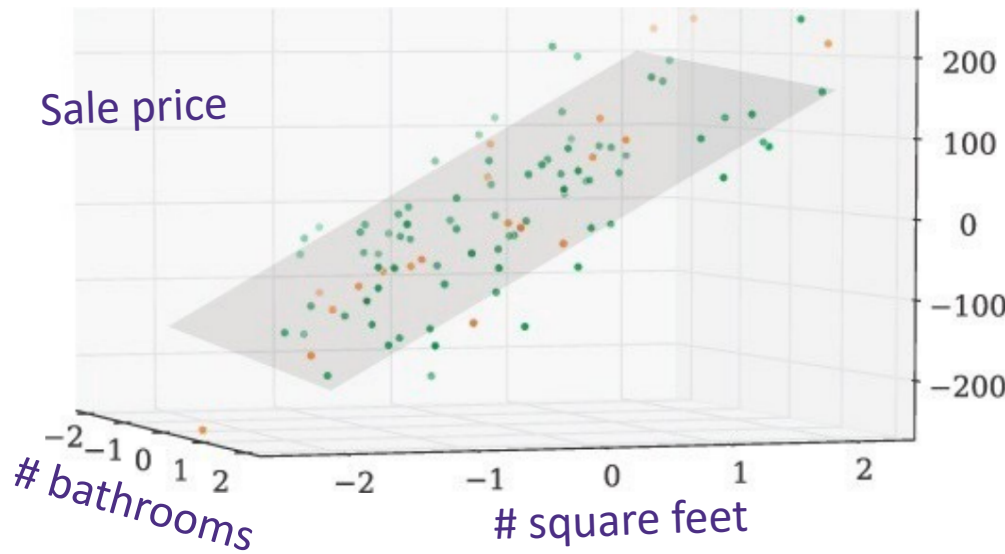
$$y_i = x_i w + \epsilon_i \quad \epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$$

The regression problem, d-dim

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

$y =$ House sale price from

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



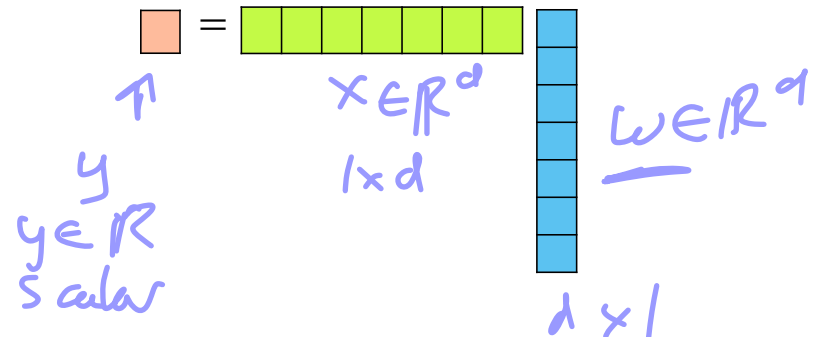
Training Data: $x_i \in \mathbb{R}^d$
 $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

Hypothesis/Model: linear

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

$\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$

$\rightarrow w \in \mathbb{R}^d$



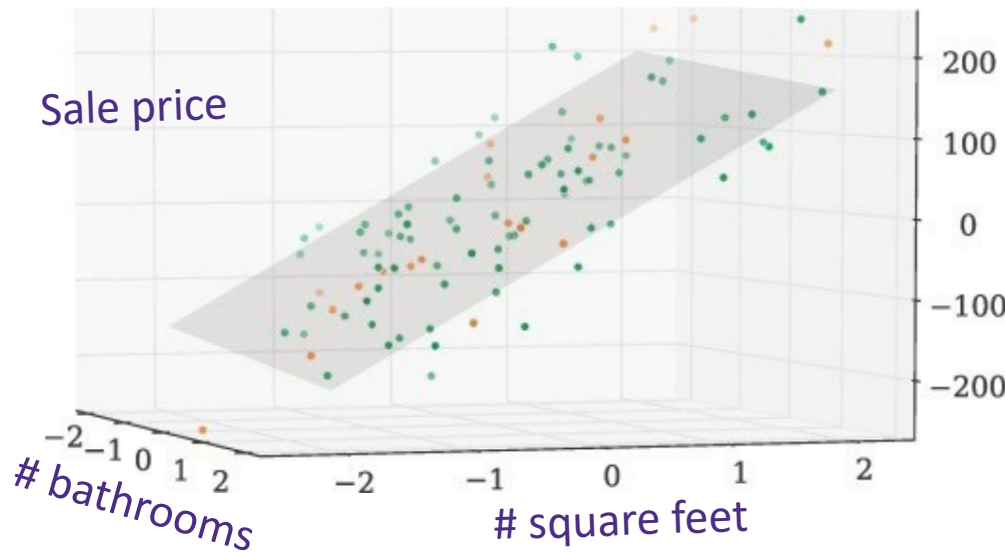


The regression problem, d-dim

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

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

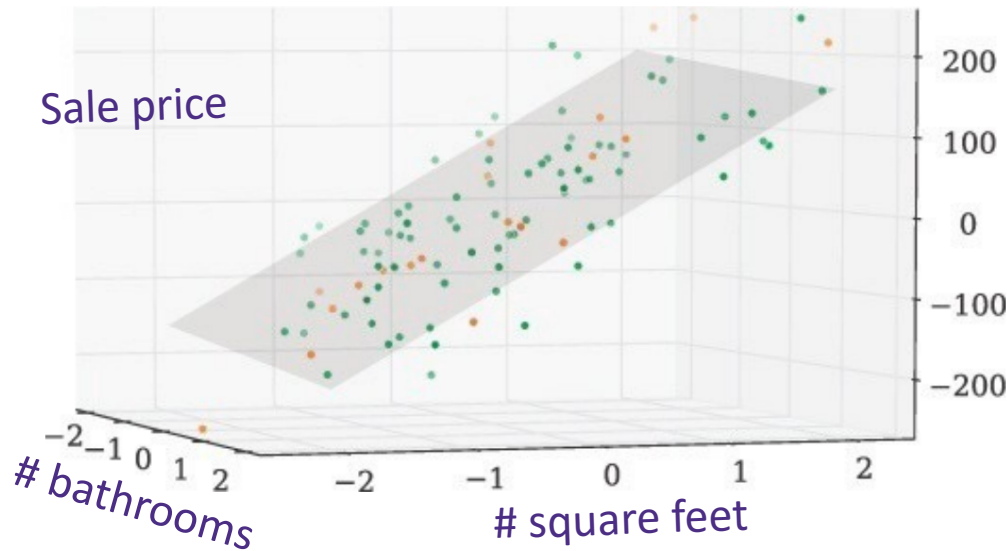
$$p(y|x, w, \sigma) = \mathcal{N}(x^T w, \sigma^2)$$

The regression problem, d-dim

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

$y =$ House sale price *from*

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



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

Hypothesis/Model: linear

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

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

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

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

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

Under benign assumptions, as the number of observations $n \rightarrow \infty$ we have $\hat{\theta}_{MLE} \rightarrow \theta_*$

Why is it useful to recover the “true” parameters θ_* of a probabilistic model?

- **Estimation** of the parameters θ_* is the goal
- Help **interpret** or summarize large datasets
- Make **predictions** about future data
- **Generate** new data $X \sim f(\cdot; \hat{\theta}_{MLE})$

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

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

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

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

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

Recap

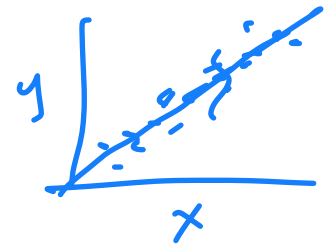
- ML "algorithms" we've learned so far...

- Maximum Likelihood Estimation (MLE)

x
 $p(x)$ $\left\{ \begin{array}{l} \rightarrow \text{fit a Bernoulli distribution (coin flips)} \\ \rightarrow \text{fit a Gaussian distribution } (\mu, \sigma) \end{array} \right.$

x, y
 $p(y|x)$ \rightarrow (currently) fit a linear predictor $x \rightarrow y$
with Gaussian noise

Maximizing log-likelihood



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

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

$$\Theta \approx w$$

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

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

↪ minimize squared error ↖

Maximizing log-likelihood

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

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

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

$$= \sum_{i=1}^n 2(y_i - x_i^T w) \begin{matrix} (-x_i) \\ \leftarrow \text{O} \\ \text{d}x \end{matrix}$$

$\begin{matrix} |x| & |x| \text{d} & \text{d}x & |x| \\ |x| & & & \end{matrix}$

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

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

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

vector of partial derivatives
 \downarrow

Set gradient=0, solve for w

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

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

// chain rule

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

"Matrix Cookbook"

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

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

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

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

matrix must be
invertible

$$n \geq d$$

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

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

we ♥ matrix notation

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

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

d : # of features

n : # of examples/datapoints



The regression problem in matrix notation

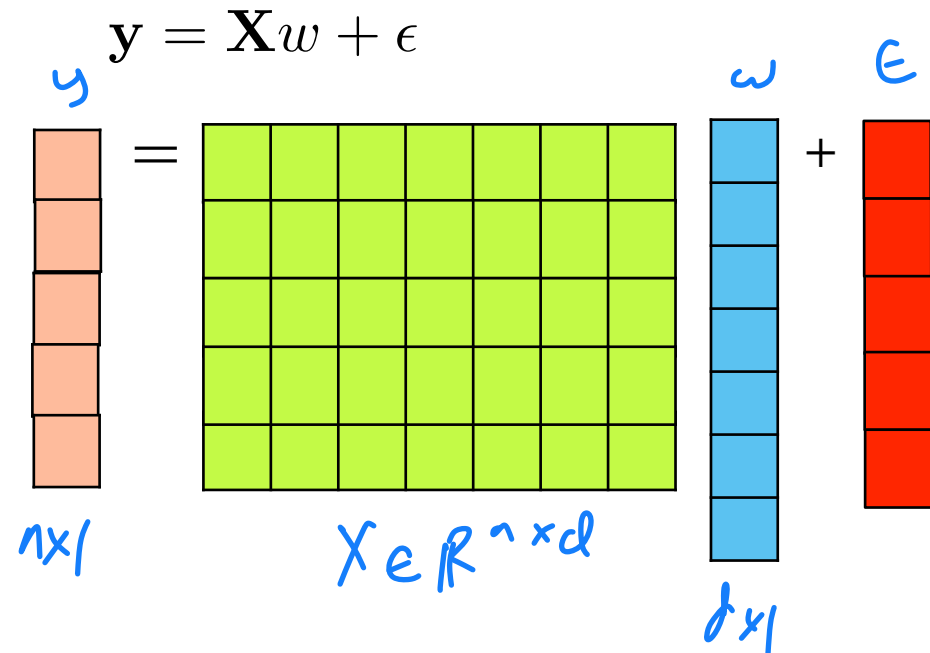
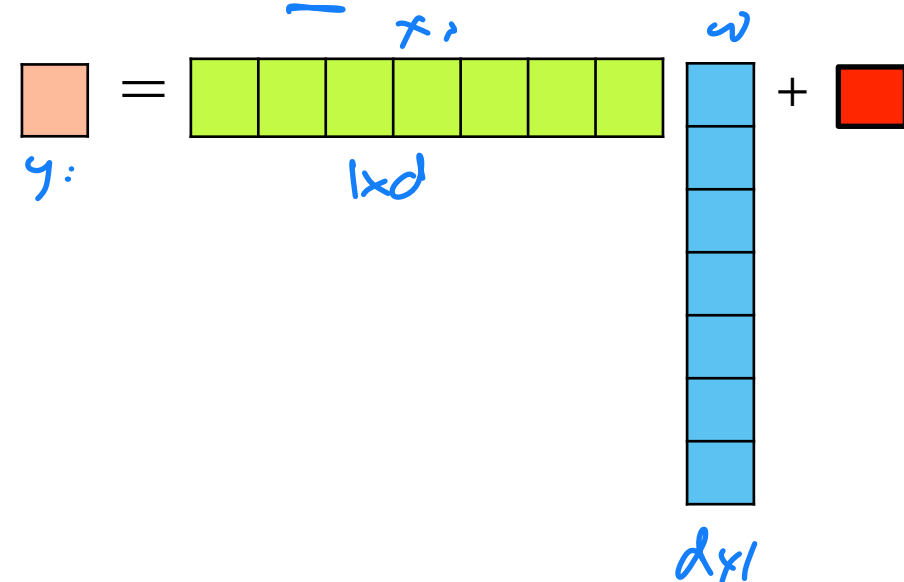
$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^\top w)^2$$

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

d : # of features

n : # of examples/datapoints

$$y_i = \underbrace{x_i^\top}_{x_i} w + \epsilon_i$$



The regression problem in matrix notation

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

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

d : # of features

n : # of examples/datapoints

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

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

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

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

$$z_i = y_i - x_i^T w$$

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

$$s \in \mathbb{R}$$

$$s = \underline{s^T}$$

$$= \arg \min_{\omega} (y - X\omega)^T (y - X\omega)$$

$$= \arg \min_{\omega} y^T y - y^T X\omega - (X\omega)^T y + (X\omega)^T (X\omega)$$

$$= \arg \min_{\omega} \cancel{y^T y} - \underbrace{y^T X\omega}_{|x_n \text{ vs } dx| \quad |x|} - \underbrace{(X\omega)^T y}_{|x \text{ vs } dx \quad |x|} + \omega^T X^T X \omega$$

$$= \arg \min_{\omega} \underline{y^T X\omega} - \underline{y^T X\omega} + \omega^T X^T X \omega$$

$$= \arg \min_{\omega} -2y^T X\omega + \omega^T X^T X \omega$$

$$\nabla_{\omega} \left[\underline{-2y^T X\omega} \right] + \underline{\omega^T X^T X \omega} = 0$$

$$= -2X^T y + 2X^T X \omega = 0$$

Useful gradients

① $\nabla_{\omega} x^T \omega = x$

② $\nabla_{\omega} x^T A \omega = \underline{A^T} x$

③ $\nabla_{\omega} \omega^T A \omega = (A + A^T) \omega$

// quadratic form

if $A = A^T$ // symmetric

then $\nabla_{\omega} \omega^T A \omega =$

$$\boxed{2A\omega}$$

$$(X^T X)^T = X^T X$$

$$= -2X^T y + 2X^T X w = 0$$

$$X^T X w = X^T y$$

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

$\frac{dxn \ n \times d}{dx \ d}$

$$= \left(\sum_{i=1}^n \underbrace{x_i x_i^T}_{\frac{dx \times d}{dx \ d}} \right)^{-1} \sum_{i=1}^n x_i y_i$$

Linear regression \equiv OLS (Ordinary Least Squares)

The regression problem in matrix notation

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

$x_i = [\text{features}]$

$w = \text{parameters}$

labels

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

d : # of features

n : # of examples/datapoints

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

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

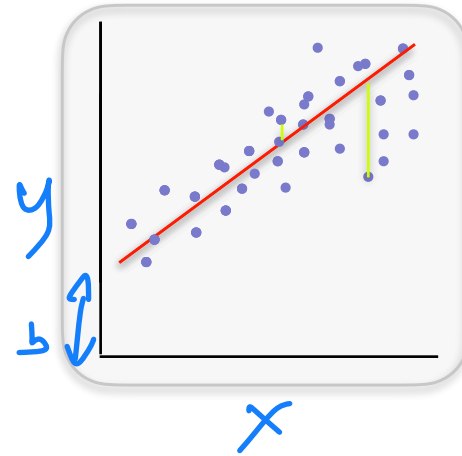
$$\ell_2 \text{ norm: } \|z\|_2 = \sqrt{\sum_{i=1}^n z_i^2} = \sqrt{z^T z}$$

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

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

The regression problem in matrix notation

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



What about an offset?

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

$n \times 1$
 $b \in \mathbb{R}$
scalar

$$\hat{y}_i = \frac{x_i^T w + b}{(y_i - \hat{y}_i)^2}$$

Dealing with an offset

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

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

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

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

$\begin{bmatrix} w \\ b \end{bmatrix} = \text{np.linalg.solve(...)}$

Dealing with an offset

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

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

$$\tilde{\mathbf{x}}_i = \begin{bmatrix} x_i \\ 1 \end{bmatrix}$$

↑
scalar $\in \mathbb{R}$

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

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

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

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

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

Dealing with an offset

pre-process
features to
have
mean zero

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

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

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

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

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

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

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?

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

→ if train features were normalized, have to normalize by subtracting the training mean

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