

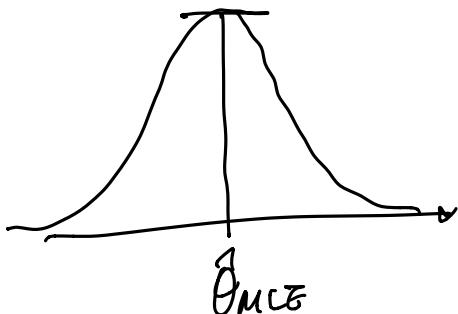
Maximum Likelihood Estimation

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

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

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

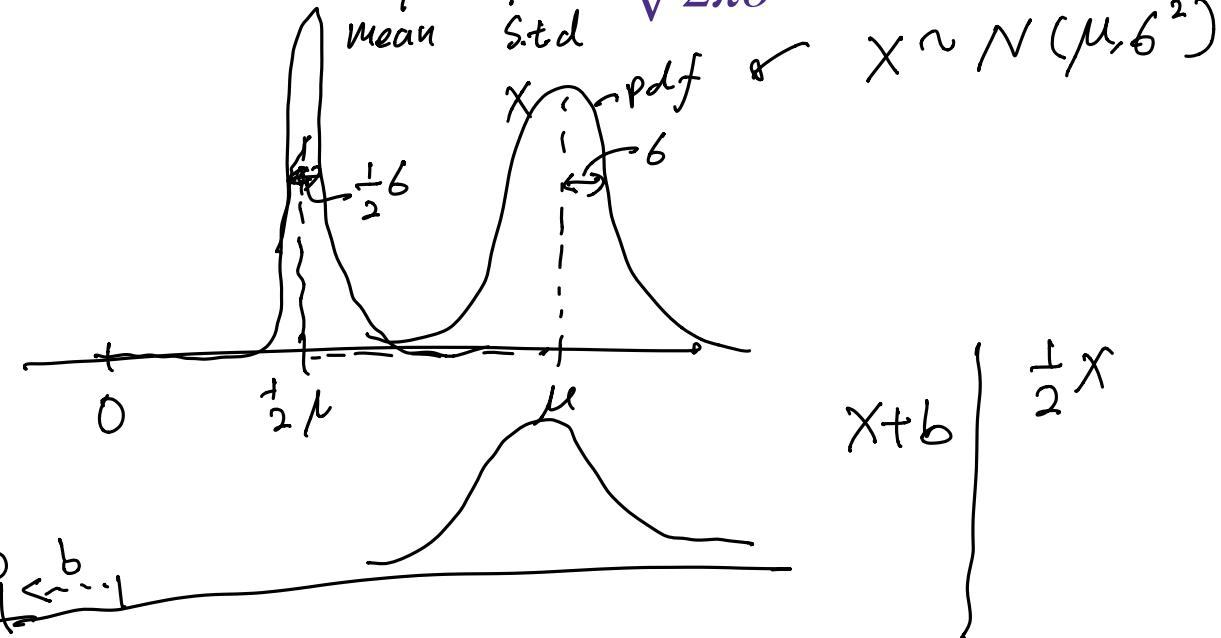
Maximum Likelihood Estimator (MLE) $\hat{\theta}_{MLE} = \arg \max_{\theta} L_n(\theta) = \frac{k}{n}$



What about continuous variables?

- **Client:** What if I am measuring a **continuous variable**?
- **You:** Let me tell you about **Gaussians**...

$$P(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Some properties of Gaussians

- affine transformation (multiplying by scalar and adding a constant)
 - $X \sim N(\mu, \sigma^2)$
 - $Y = aX + b \rightarrow Y \sim N(a\mu + b, a^2\sigma^2)$
- Sum of Gaussians
 - $X \sim N(\mu_X, \sigma^2_X)$
 - $Y \sim N(\mu_Y, \sigma^2_Y)$
 - $Z = X+Y \rightarrow Z \sim N(\mu_X + \mu_Y, \sigma^2_X + \sigma^2_Y)$

MLE for Gaussian

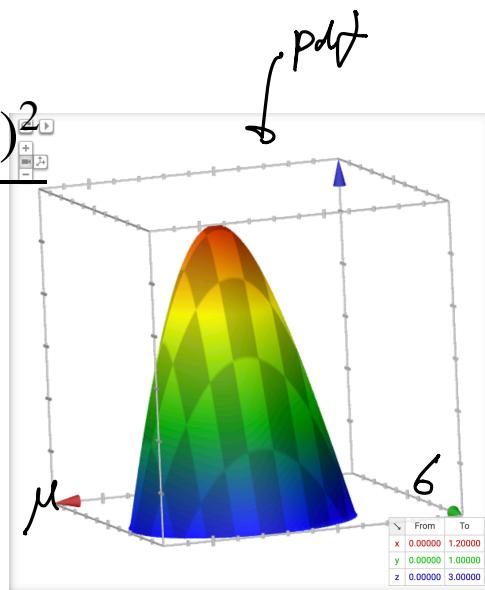
- Prob. of i.i.d. samples $D=\{x_1, \dots, x_n\}$ (e.g., temperature):

$$\begin{aligned} P(\mathcal{D}; \mu, \sigma) &= P(x_1, \dots, x_n; \mu, \sigma) = P(x_1 | \mu, \sigma) \times \dots \times P(x_n | \mu, \sigma) \\ &= \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} = \underline{f(x_i | \mu, \sigma)} \end{aligned}$$

- Log-likelihood of data:

$$\log P(\mathcal{D}; \mu, \sigma) = -n \log(\sigma \sqrt{2\pi}) - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}$$

- What is $\hat{\theta}_{MLE}$ for $\theta = (\mu, \sigma^2)$?



Your second learning algorithm: MLE for mean of a Gaussian

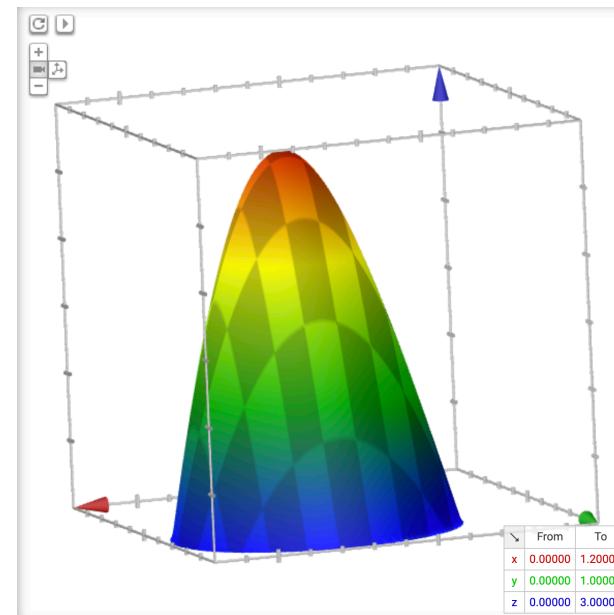
- What's MLE for mean?

$$\frac{d}{d\mu} \log P(\mathcal{D}; \mu, \sigma) = \frac{d}{d\mu} \left[-n \log(\sigma\sqrt{2\pi}) - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

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

$$= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = \frac{1}{\sigma^2} \left(\sum_{i=1}^n x_i - n\mu \right) \Big|_{\mu=\hat{\mu}} = 0$$

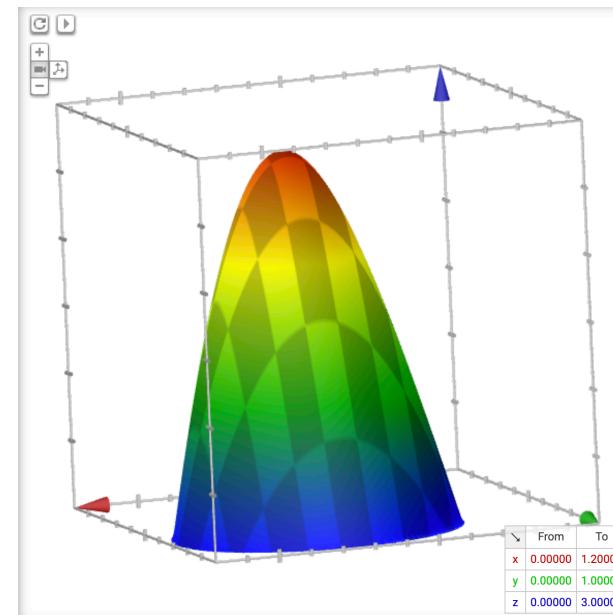
$$\frac{1}{n} \sum_{i=1}^n x_i = \hat{\mu}_{LS}$$



MLE for variance

- Again, set derivative to zero:

$$\begin{aligned} \frac{d}{d\sigma} \log P(\mathcal{D}; \mu, \sigma) &= \frac{d}{d\sigma} \left[-n \log(\sigma\sqrt{2\pi}) - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right] \\ &= \frac{-n}{6} - \sum_{i=1}^n \frac{(x_i - \mu)^2 (-2)}{2\sigma^3} \\ &= \frac{1}{6\sigma^3} \left(-n\sigma^2 + \sum_{i=1}^n (x_i - \mu)^2 \right) \Bigg|_{\mu = \hat{\mu}, \sigma = \hat{\sigma}} \\ &\quad - n\hat{\sigma}^2 + \sum_{i=1}^n (x_i - \hat{\mu})^2 = 0 \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 \end{aligned}$$



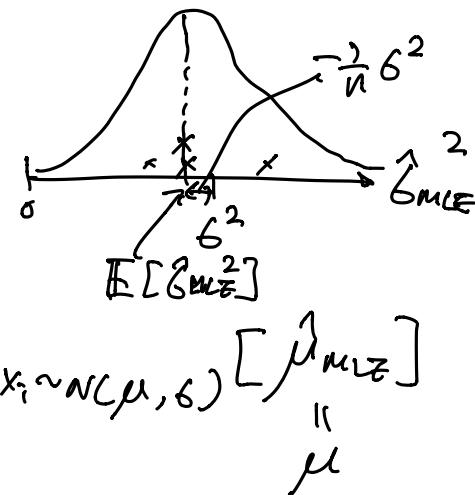
What can we say about the MLE?

- MLE:

	mean	variance
μ	\checkmark	\checkmark
VAR	?	?

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

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



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

$$\frac{n-1}{n} \cdot \sigma^2 = \mathbb{E}[\hat{\sigma}^2_{MLE}] \neq \sigma^2$$

$$(1 - \frac{1}{n})\sigma^2 = \mathbb{E}[\hat{\sigma}^2_{MLE}]$$

bias = $-\frac{1}{n}\sigma^2$

- Unbiased variance estimator:

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

Maximum Likelihood Estimation

$$\text{Normal Distribution: } f(x; \theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad \theta = (\mu, \sigma)$$

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

Properties (under benign regularity conditions—smoothness, identifiability, etc.):

- Asymptotically consistent and normal: $\frac{\hat{\theta}_{MLE} - \theta_*}{\hat{s}_\theta} \sim \mathcal{N}(0, 1)$
- Asymptotic Optimality, minimum variance (see Cramer-Rao lower bound)

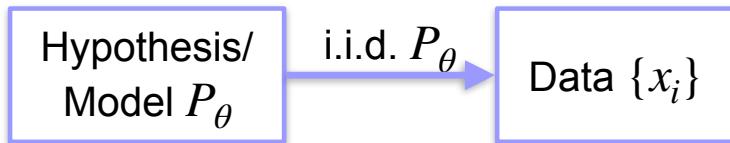
Recap

- Learning is...
 - Collect some data
 - E.g., coin flips

Data $\{x_i\}$

Recap

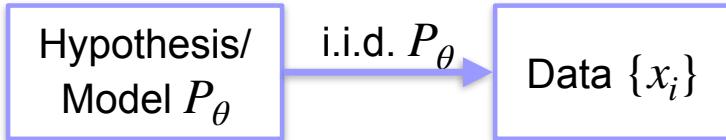
- Learning is...
 - Collect some data
 - E.g., coin flips
 - Choose a hypothesis class or model
 - E.g., binomial



Recap

- Learning is...
 - Collect some data
 - E.g., coin flips
 - Choose a hypothesis class or model
 - E.g., binomial
 - Choose a loss function
 - E.g., data likelihood

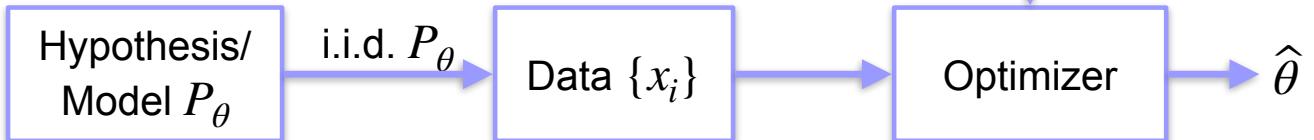
$$\min_{\theta} \sum_{i=1}^n \ell(\theta, x_i) \quad \text{or} \quad \max_{\theta} \sum_{i=1}^n u(\theta, x_i)$$



Recap

- Learning is...
 - Collect some data
 - E.g., coin flips
 - Choose a hypothesis class or model
 - E.g., binomial
 - Choose a loss function
 - E.g., data likelihood
 - Choose an optimization procedure
 - E.g., set derivative to zero to obtain MLE

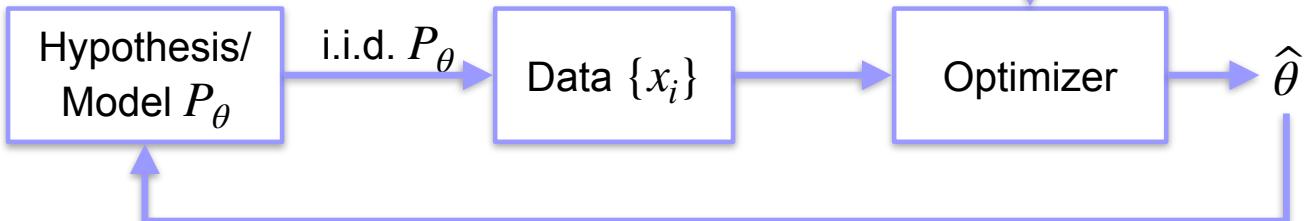
$$\min_{\theta} \sum_{i=1}^n \ell(\theta, x_i) \quad \text{or} \quad \max_{\theta} \sum_{i=1}^n u(\theta, x_i)$$



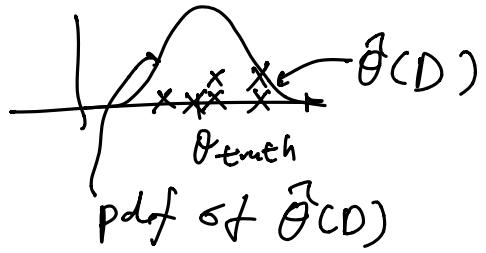
Recap

- Learning is...
 - Collect some data
 - E.g., coin flips
 - Choose a hypothesis class or model
 - E.g., binomial
 - Choose a loss function
 - E.g., data likelihood
 - Choose an optimization procedure
 - E.g., set derivative to zero to obtain MLE
 - Justifying the accuracy of the estimate
 - E.g., Markov's inequality

$$\min_{\theta} \sum_{i=1}^n \ell(\theta, x_i) \quad \text{or} \quad \max_{\theta} \sum_{i=1}^n u(\theta, x_i)$$



Usually, $D = \{x_i\}_{i=1}^n \rightarrow \hat{\theta}(D)$
 Random Variable



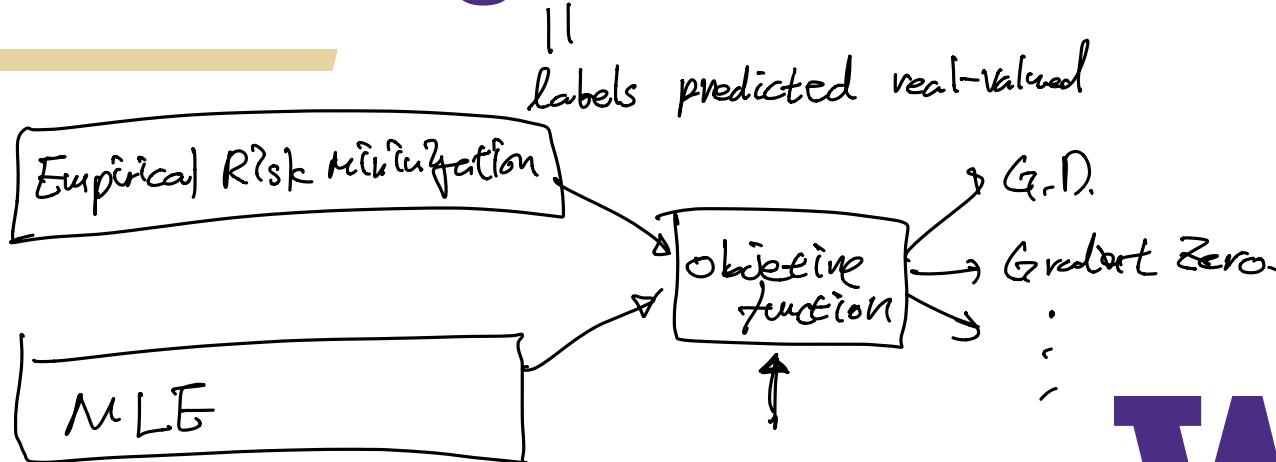
Imagine, again,

Hypothesis class is linear function

||

$$\sum_{i=1}^d w_i \cdot x_{i,f} = r_i$$

Linear Regression



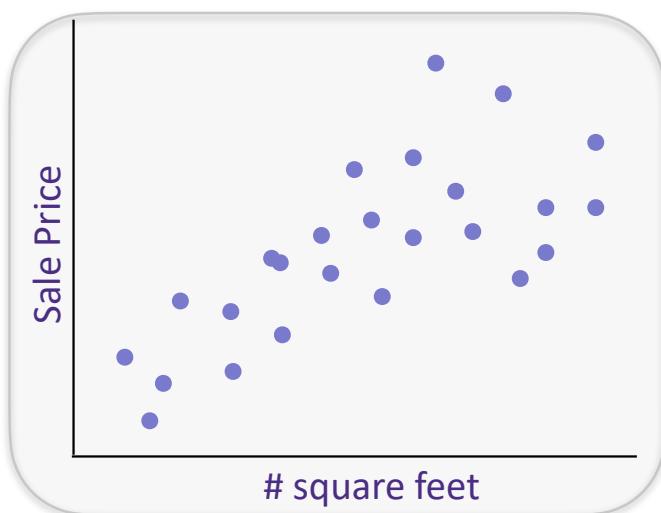
W

The regression problem, 1-dimensional

Given past sales data on [zillow.com](#), predict:

$y = \text{House sale price from}$

$x = \{\# \text{ sq. ft.}\}$



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

input labels

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

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

Process

Decide on a model

y_i
assume house sale price is a linear function of
square feet.
 x_i

Find the function which fits the data best

Use function to make prediction on new examples

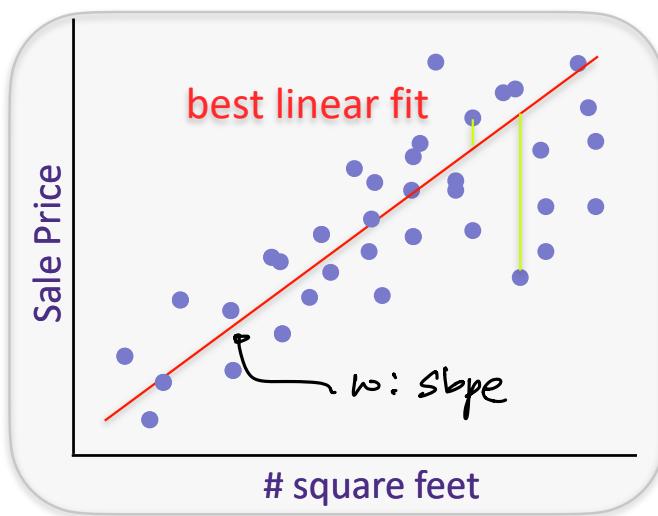
$$y_{\text{new}} = f(x_{\text{new}})$$

Fit a function to our data, 1-dimension

Given past sales data on [zillow.com](#), predict:

y = House sale price from

x = {# sq. ft.}



Error

$$y_i = x_i w + \epsilon_i$$

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

Hypothesis/Model: linear

$$y_i \approx x_i w$$

Loss: least squares solution

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

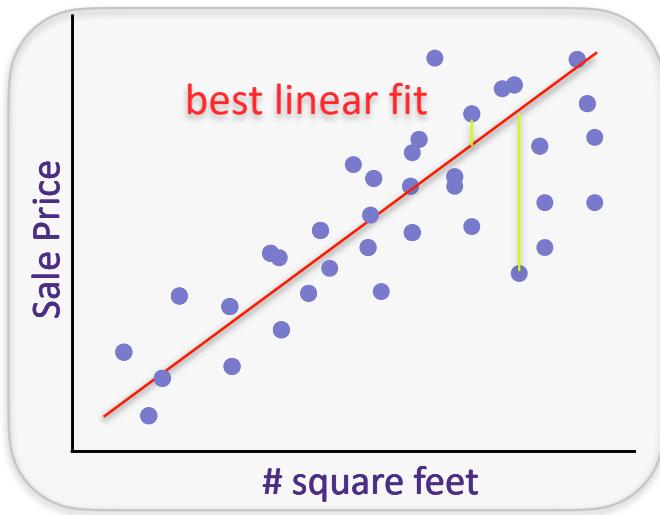
$\underbrace{(y_i - x_i w)}_{\epsilon_i : \text{error}}$

The regression problem, d-dimensions

Given past sales data on [zillow.com](#), predict:

y = House sale price from

x = {# sq. ft., zip code, date of sale, etc.}



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

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

Hypothesis/Model: linear

$$y_i \approx x_i^T w \approx x_{i,1}w_1 + x_{i,2}w_2 + \dots + x_{i,d}w_d$$

Loss: least squares solution

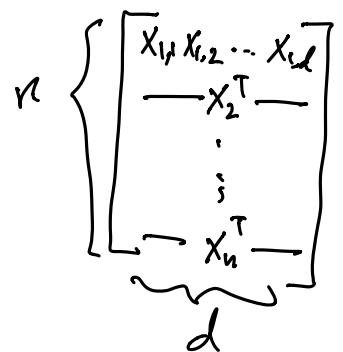
$$\min_w \sum_{i=1}^n \underbrace{(y_i - x_i^T w)^2}_{\text{error}}$$

The regression problem in matrix notation

Data:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$

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



The regression problem in matrix notation

Data:

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

d : # of features

n : # of examples/datapoints

Model:

$$y_1 = x_1^T w + \epsilon_1$$

$$y_2 = x_2^T w + \epsilon_2$$

•

•

$$y_n = x_n^T w + \epsilon_n$$

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

A hand-drawn diagram illustrating the regression model. It shows a vertical vector \mathbf{y} on the left, labeled 'n' below it, followed by an equals sign. To its right is a large bracket labeled 'd' above it, enclosing a vertical rectangle labeled ' \mathbf{X} '. To the right of ' \mathbf{X} ' is a plus sign, followed by another large bracket labeled 'n' below it, enclosing a vertical rectangle labeled ' w '. To the right of ' w ' is a plus sign, followed by a final large bracket labeled 'n' below it, enclosing a vertical rectangle labeled ' $\boldsymbol{\epsilon}$ '.

The regression problem in matrix notation

Data:

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

d : # of features

n : # of examples/datapoints

Model:

$$y_1 = x_1^T w + \epsilon_1 \quad \mathbf{y} = \mathbf{X}w + \epsilon$$

$$y_2 = x_2^T w + \epsilon_2$$

•

⋮

•

$$y_n = x_n^T w + \epsilon_n$$

Loss: $\hat{w}_{LS} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 \leftarrow \sum_i \|v_i\|_2^2 = \|v\|_2^2 = v^T v$

$$\|v\|_2 = \sqrt{v_1^2 + v_2^2 + \dots}$$

The regression problem in matrix notation

Data:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$

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

Model:

$$y_1 = \mathbf{x}_1^T \mathbf{w} + \epsilon_1$$

$$y_2 = \mathbf{x}_2^T \mathbf{w} + \epsilon_2$$

•

•

$$y_n = \mathbf{x}_n^T \mathbf{w} + \epsilon_n$$

Loss: $\hat{\mathbf{w}}_{LS} = \arg \min_w \sum_{i=1}^n \underbrace{(y_i - \mathbf{x}_i^T \mathbf{w})^2}_{\hat{\epsilon}_i}$

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

$$\begin{bmatrix} y_1 - \mathbf{x}_1^T \mathbf{w} \\ \vdots \\ y_n - \mathbf{x}_n^T \mathbf{w} \end{bmatrix} = \mathbf{y} - \mathbf{X}\mathbf{w}$$

$$\begin{aligned} \|\mathbf{v}\|_2 &\triangleq \sqrt{v_1^2 + \dots + v_d^2} \\ \|\mathbf{v}\|_2^2 &= (\|\mathbf{v}\|_2)^2 \\ &= \sqrt{v_1^2 + \dots + v_d^2}^2 \\ &= \sqrt{\mathbf{v}^T \mathbf{v}} \end{aligned}$$

The regression problem in matrix notation

parameters $\theta = w$

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

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

$$\mathcal{L}(w) \parallel$$

$$f(w) = w^T w = w_1^2 + w_2^2 + \dots$$

$$\nabla_w f(w) = 2 \cdot w$$

$$\nabla_w f(w) = 2(\mathbf{A}w)^T \mathbf{A}w$$

$$\nabla_w f(w) = 2 \mathbf{A}^T \mathbf{A}w$$

$$\nabla_w f(w) = (\mathbf{A}w + b)^T (\mathbf{A}w + b)$$

$$\nabla_w f(w) = 2 \mathbf{A}^T (\mathbf{A}w + b)$$

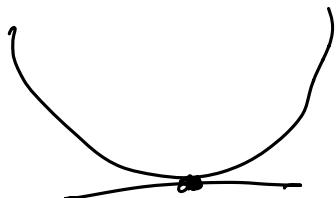
$$\|\mathbf{v}\|_p = (\sqrt[p]{v_1^p + \dots + v_n^p})^{1/p}$$

$$\boxed{p=0, 1, 2, \infty}$$

NNZ

$$\boxed{\frac{1}{2} \nabla_w \mathcal{L}(w) = -\mathbf{X}^T (\mathbf{y} - \mathbf{X}w)} \quad ?$$

$$= -\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X}w \Big|_{w=\hat{w}_{LS}} = 0$$



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

$$\underbrace{d \square}_d \quad \underbrace{\square d}_n = \underbrace{d \square}_n \quad \underbrace{\square}_n$$

invertible.

$$\square \cdot \square$$

Closed-form

The regression problem in matrix notation

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

“Closed form” solution!