

# Lecture 2: MLE for Gaussian and linear regression

2nd MLE example

↑  
1st ML algorithm.

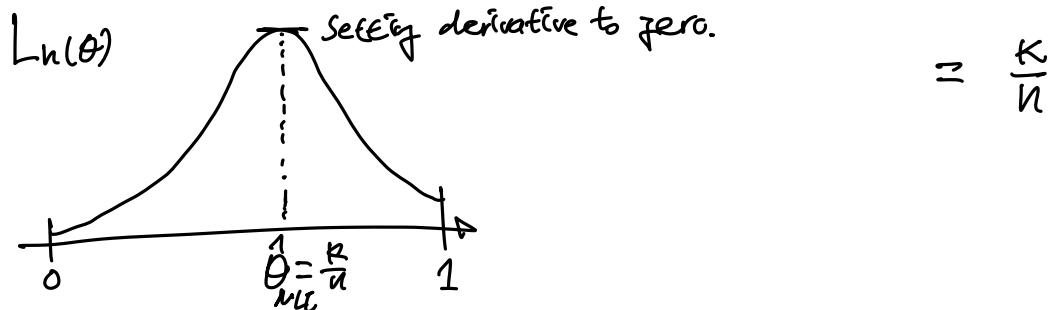
W

# Recap: Maximum Likelihood Estimation

$$\text{Bern}(\theta) \rightarrow P(H|\theta) = \theta$$

$$P(T|\theta) = 1-\theta$$

- **Observe**  $X_1, X_2, \dots, X_n$  drawn i.i.d. from  $P(X_i; \theta)$  for some true  $\theta = \theta^*$
- **Likelihood function:**  $L_n(\theta) = \prod_{i=1}^n P(X_i; \theta) = \theta^K (1-\theta)^{n-K}$  if  $K$  heads.
- **Log-likelihood function:**  $\ell_n(\theta) = \log L_n(\theta) = \sum_{i=1}^n \log P(X_i; \theta) = K \cdot \log \theta + (n-K) \log(1-\theta)$
- **Maximum Likelihood Estimator (MLE):**  $\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \ell_n(\theta)$



# What about continuous variables?

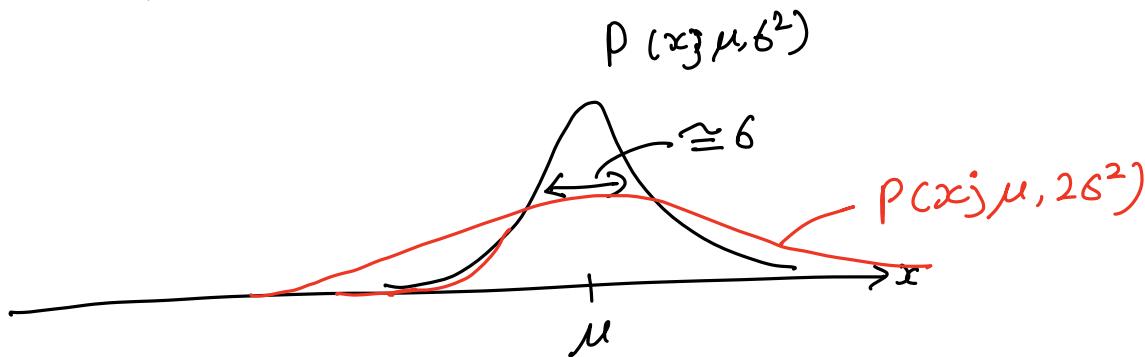
- *Client:* What if I am measuring a **continuous variable**?
- *You:* Let me tell you about **Gaussians...**  $X \sim N(\mu, \sigma^2)$

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

$\mu \triangleq \mathbb{E}[x]$        $\sigma^2 \triangleq \mathbb{E}[(x - \mathbb{E}[x])^2]$

Defined as

← P.d.f  
Probability  
density  
function



# Some properties of Gaussians

---

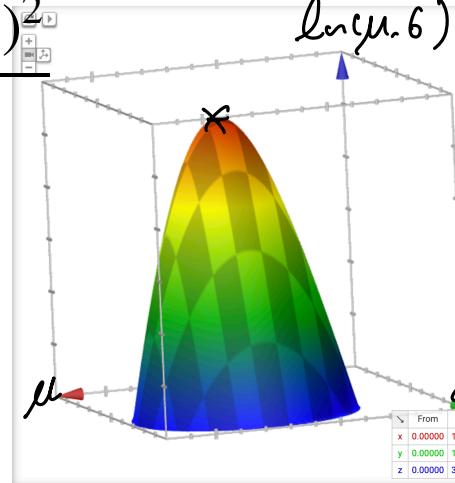
- affine transformation (multiplying by scalar and adding a constant)
  - $X \sim N(\mu, \sigma^2)$
  - $Y = aX + b \rightarrow Y \sim N(a\mu + b, a^2\sigma^2)$
- Sum of Gaussians
  - $X \sim N(\mu_X, \sigma^2_X)$
  - $Y \sim N(\mu_Y, \sigma^2_Y)$
  - $Z = X+Y \rightarrow Z \sim N(\mu_X + \mu_Y, \sigma^2_X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

# MLE for Gaussian

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

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

- Log-likelihood of data:

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


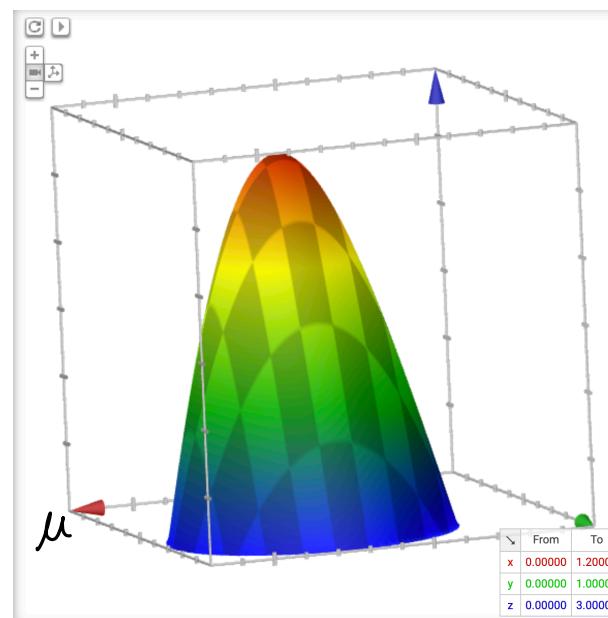
- What is  $\hat{\theta}_{MLE}$  for  $\theta = (\mu, \sigma^2)$  ?

# Your second learning algorithm: MLE for mean of a Gaussian

- What's MLE for mean?

$$\begin{aligned}\frac{d}{d\mu} \log P(\mathcal{D}; \mu, \sigma) &= \frac{d}{d\mu} \left[ -n \log(\sigma\sqrt{2\pi}) - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right] \\&= -\sum_{i=1}^n \frac{2(x_i - \mu)}{2\sigma^2} \quad \text{for } \hat{\mu}_{\text{MLE}} \\&= -\frac{1}{\sigma^2} \left\{ \sum_{i=1}^n x_i - n \cdot \mu \right\} \stackrel{!}{=} 0 \\&\Rightarrow \hat{\mu}_{\text{MLE}} = \frac{1}{n} \sum_{i=1}^n x_i : \text{empirical mean}\end{aligned}$$

\* Note that  $\hat{\mu}_{\text{MLE}}$  does not depend on  $\hat{\sigma}_{\text{MLE}}$ .  
for Gaussian.



# MLE for variance

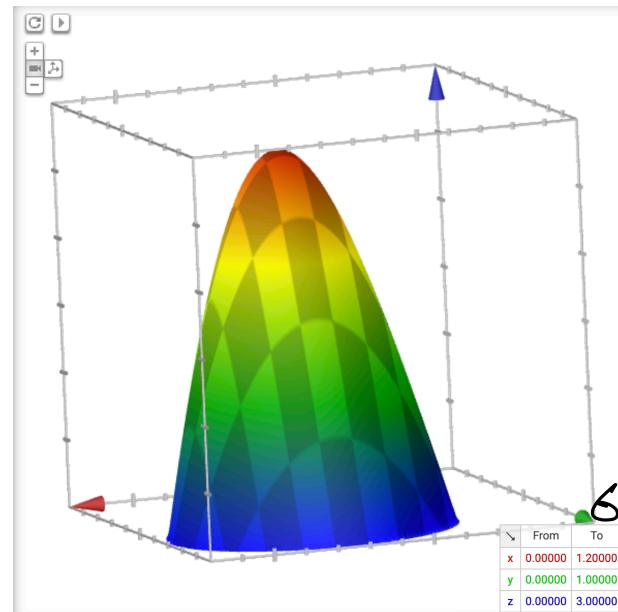
- Again, set derivative to zero:

$$\frac{d}{d\sigma} \log P(\mathcal{D}; \mu, \sigma) = \frac{d}{d\sigma} \left[ -n \log(\sigma\sqrt{2\pi}) - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$
$$= -\frac{n}{6} - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^3} (-2)$$

$$= -\frac{n}{6\sigma^3} \left\{ 6^2 - \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \right\} \Big|_{\mu = \hat{\mu}_{MLE}} = 0$$

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

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



# What can we say about the MLE?

- 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**

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

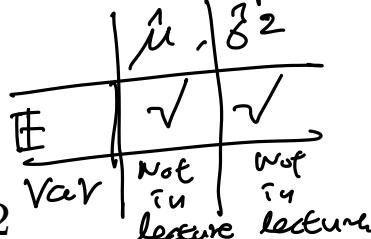
$$\mathbb{E}[\hat{\sigma}^2_{MLE}] = (1 - \frac{1}{n}) \sigma^2$$

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

- Unbiased variance estimator:

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

\* Treat  $(\hat{\mu}_{MLE}, \hat{\sigma}^2_{MLE})$  as Random Variables.  
Imagine we run many experiments.



$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \mathbb{E}\left[\frac{1}{n} \cdot n \cdot \mu\right] = \mu$$

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

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

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

# Recap

---

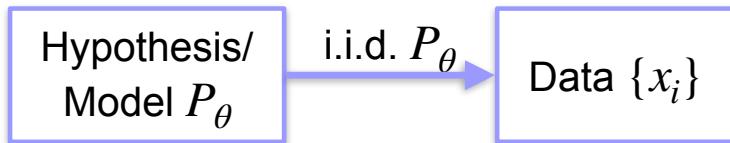
- Learning is...
  - Collect some data
    - E.g., coin flips

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

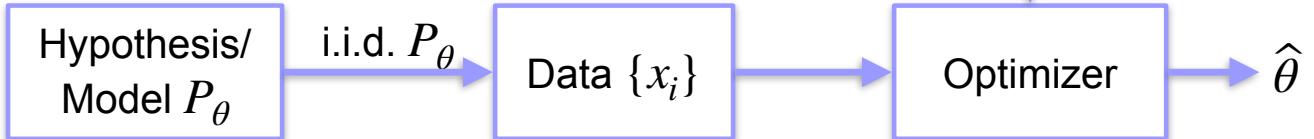
$$\min_{\theta} \sum_{i=1}^n \ell(\theta, x_i) \quad \text{or} \quad \max_{\theta} \sum_{i=1}^n u(\theta, x_i)$$



# 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

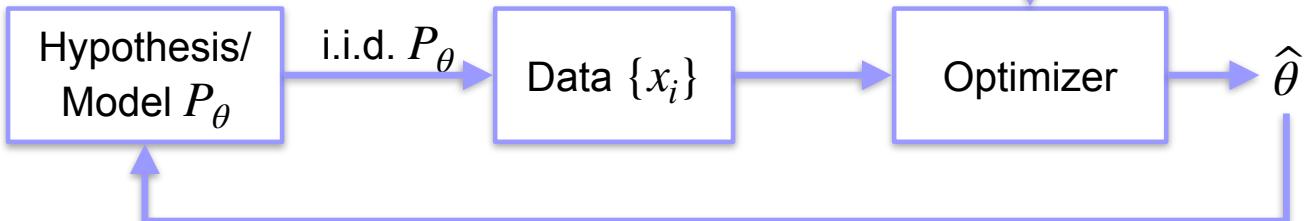
$$\min_{\theta} \sum_{i=1}^n \ell(\theta, x_i) \quad \text{or} \quad \max_{\theta} \sum_{i=1}^n u(\theta, x_i)$$



# 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

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



hypothesis class is linear functions

↓  
Supervised Learning with continuous labels.

# Linear Regression

---

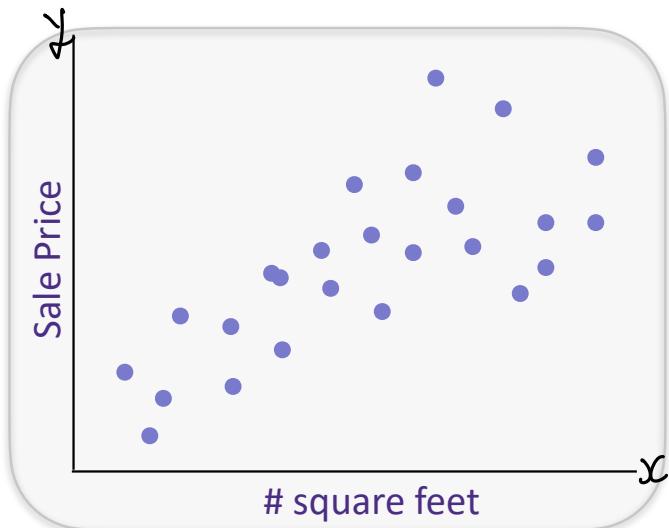
W

# The regression problem, 1-dimensional

You are trying to sell your house, what is the right price?  
Given past sales data on zillow.com, predict:

$y$  = House sale price from

$x$  = {# sq. ft.}



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

↑  
input      ↑  
label

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

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

# Process

---

Decide on a **model** / hypothesis explaining Data, i.e. how price is related to sqft.

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

$$y_i$$
$$x_i$$

Find the function which fits the data best  
or model

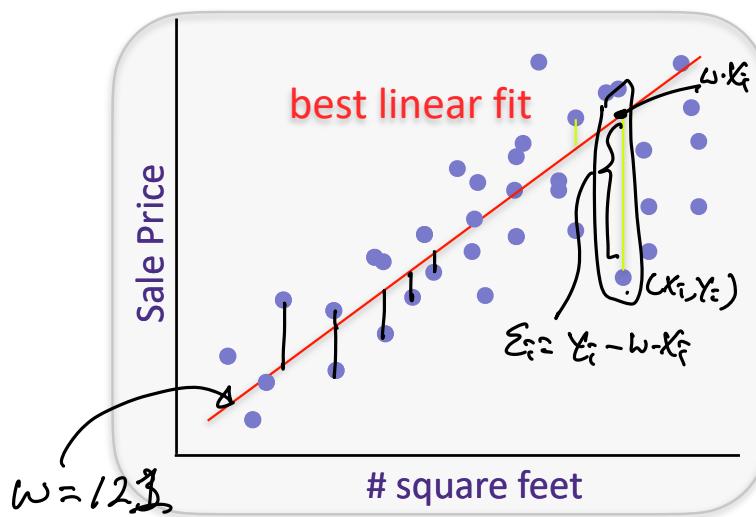
Use function to make prediction on new examples

# 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

$y_i \approx \underbrace{w \cdot x_i}_{\text{model}} + \underbrace{\epsilon_i}_{\text{error/noise}}$   
label input  
model parameter

Loss: least squares solution

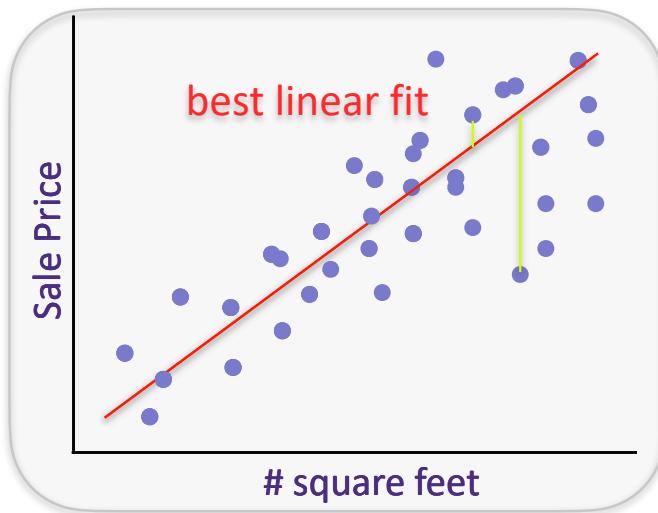
$$\min_w \sum_{i=1}^n (y_i - x_i w)^2$$

# 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:

$$y_i = x_i w + \epsilon_i$$

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

$$x_i \in \mathbb{R}^d$$
$$y_i \in \mathbb{R}$$

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix} \quad \left. \right\}^d$$

Hypothesis/Model: linear

$$y_i \approx x_i^T w = [x_{i1} \ x_{i2} \ \dots \ x_{id}] \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} = x_{i1} w_1 + x_{i2} w_2 + \dots + x_{id} w_d$$

Loss: least squares solution

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

$\overbrace{\quad}^{n \text{ rows}}$

# The regression problem in matrix notation

Data:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$
$$n = \left\{ \begin{bmatrix} x_{11} & \cdots & x_{1d} \\ \hline x_2^T \\ \vdots \\ \hline x_n^T \end{bmatrix} \right\}_d$$

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

\* it is really important  
to keep track of dimensions.

# 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:

$$y_1 = \mathbf{x}_1^T \mathbf{w} + \epsilon_1$$
$$y_2 = \mathbf{x}_2^T \mathbf{w} + \epsilon_2$$
$$\vdots$$
$$y_n = \mathbf{x}_n^T \mathbf{w} + \epsilon_n$$

*label*    *input · model*    *noise*

*equivalent*

$$\xleftrightarrow{n} \mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$
$$\left[ \begin{array}{c} \mathbf{y} \\ \vdots \\ \mathbf{y}_n \end{array} \right] = \left[ \begin{array}{c} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{array} \right] \left[ \begin{array}{c} \mathbf{w} \\ \vdots \\ \mathbf{w}_d \end{array} \right] + \left[ \begin{array}{c} \epsilon_1 \\ \vdots \\ \epsilon_n \end{array} \right]$$

# 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:**

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

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

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

norm

•

•

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

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

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

# 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:**

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

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

•

⋮

•

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

$$\begin{aligned} \textbf{Loss: } \hat{w}_{LS} &= \arg \min_w \sum_{i=1}^n \underbrace{(y_i - x_i^T w)}_{\varepsilon_i}_2^2 = \arg \min_w \underbrace{\|\mathbf{y} - \mathbf{X}w\|_2^2}_{\varepsilon} \\ &= \arg \min_w (\underbrace{\mathbf{y} - \mathbf{X}w}_\varepsilon)^T (\underbrace{\mathbf{y} - \mathbf{X}w}_\varepsilon) \end{aligned}$$

# The regression problem in matrix notation

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

model parameter  $\theta$

Set gradient w.r.t.  $w$  to zero to find the minima:

$$\begin{aligned}\nabla_w \mathcal{L}(w) &= -\mathbf{X}^T (\mathbf{y} - \mathbf{X}w) \\ &= -\mathbf{X}^T \mathbf{y} + \mathbf{X}^T \mathbf{X} w \Big|_{w=\hat{w}_{LS}} = 0\end{aligned}$$

$$\Rightarrow \mathbf{X}^T \mathbf{X} \cdot w = \mathbf{X}^T \mathbf{y} \in \mathbb{R}^d \Rightarrow w = (\mathbf{X}^T \mathbf{X})^{-1} \cdot \mathbf{X}^T \mathbf{y}$$

$d \left[ \begin{array}{c} \\ \vdots \\ n \end{array} \right] \left[ \begin{array}{c} \\ \vdots \\ n \end{array} \right]^T \left[ \begin{array}{c} \\ \vdots \\ d \end{array} \right] = d \left[ \begin{array}{c} \\ \vdots \\ n \end{array} \right] \left[ \begin{array}{c} \\ \vdots \\ n \end{array} \right] \quad \begin{matrix} d \square \\ \square d \end{matrix}$

\* We will assume that  $\mathbf{X}^T \mathbf{X}$  is full rank and invertible

# 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!

# Questions?

---