

# Lecture 2: MLE for Gaussian and linear regression

*Seung Oh*

- HW0 due Wednesday October 8th midnight
- Good idea to use the LaTeX source we provide, and use overleaf
  - we now have the full.png on the web also
- Some office hours can get cancelled due to travel, etc. Check Ed/Logistics.



# Recap: Maximum Likelihood Estimation

- **Observe**  $\mathcal{D} = \{X_1, X_2, \dots, X_n\}$  drawn i.i.d. from  $P(X_i; \theta)$  for some ground truth  $\theta = \theta^*$ , unknown to us

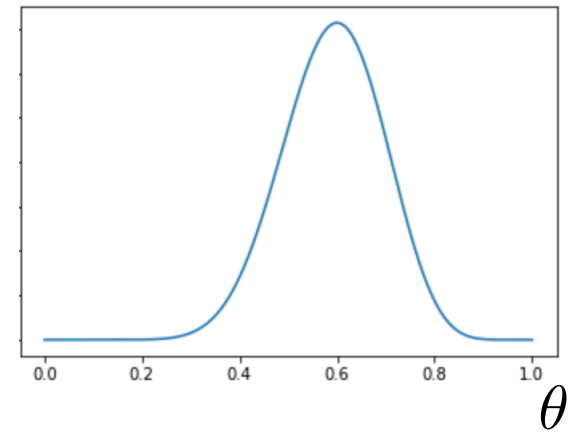
*Binomial (n, θ)*

- **Maximize log-likelihood** when we observe  $k$  heads in  $n$  flips

$$\log P(\mathcal{D}; \theta) = \log \{ \theta^k (1-\theta)^{n-k} \}$$

$$= k \log(\theta) + (n-k) \log(1-\theta)$$

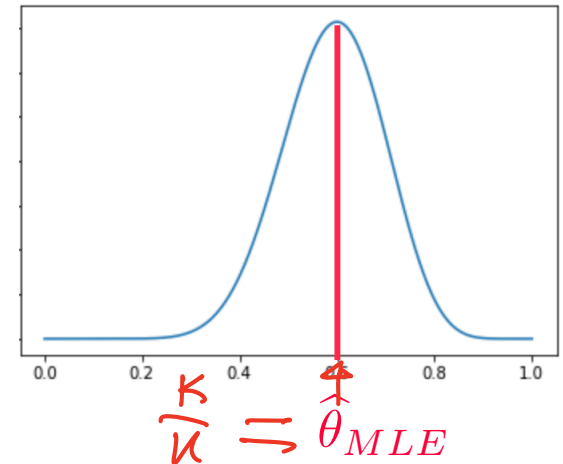
$P(\mathcal{D}; \theta)$  when  $k$  heads observed in  $n$  flips



# Recap: Maximum Likelihood Estimation

- **Observe**  $\mathcal{D} = X_1, X_2, \dots, X_n$  drawn i.i.d. from  $P(X_i; \theta)$  for some ground truth  $\theta = \theta^*$ , unknown to us
- **Maximize log-likelihood** when we observe  $k$  heads in  $n$  flips  
 $\log P(\mathcal{D}; \theta) = k \log(\theta) + (n - k) \log(1 - \theta)$

$P(\mathcal{D}; \theta)$  when  $k$  heads observed in  $n$  flips



- Use the fact that derivative is zero at maxima (and also minima)
- Set derivative to zero,

and find  $\theta$  satisfying:

$$\frac{d}{d\theta} \log P(\mathcal{D}; \theta) = 0$$

# Maximum Likelihood Estimation

- **Observe**  $X_1, X_2, \dots, X_n$  drawn i.i.d. from  $P(X_i; \theta)$  for some true  $\theta = \theta^*$

- **Likelihood function:**  $L_n(\theta) = \prod_{i=1}^n P(X_i; \theta)$

- **Log-likelihood function:**  $\ell_n(\theta) = \log L_n(\theta) = \sum_{i=1}^n \log P(X_i; \theta)$

- **Maximum Likelihood Estimator (MLE):**  $\hat{\theta}_{\text{MLE}} \stackrel{\text{①}}{=} \arg \max_{\theta} \ell_n(\theta)$   
②

- Warning when setting the derivative to zero to find the MLE:

- The solution includes maxima, minima, and stationary points  $\Rightarrow$  needs to be checked
- It does not always lead to an explicit expression in a closed form  $\Rightarrow$  alternative methods

"optimization"

# What about continuous variables?

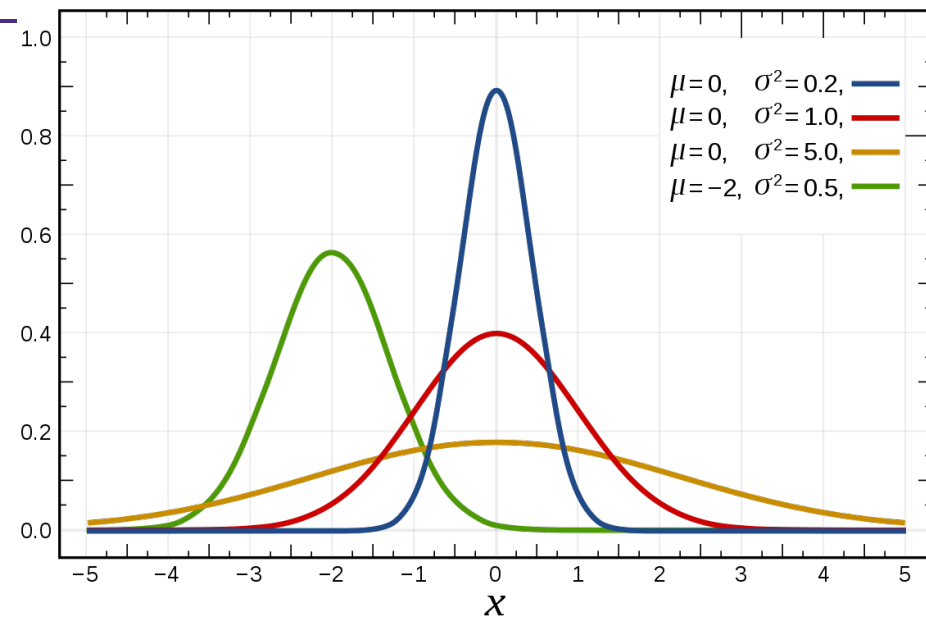
- *Client*: What if I am measuring a **continuous variable**?

- *You*: Let me tell you about **Gaussians**...

- A Gaussian random variable is written as  $X \sim \mathcal{N}(\mu, \sigma^2)$  with mean  $\mu \triangleq \mathbb{E}[X]$  and variance  $\sigma^2 \triangleq \mathbb{E}[(X - \mathbb{E}[X])^2]$

- The p.d.f. (Probability Density Function) of  $X$  is

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



Binom  
← Gaussian

[EdDiscussion Question: What distributions do we need to memorize?]

# Some useful properties of Gaussians

- Affine transformation  
(multiplying by scalar and adding a constant)

- $X \sim \mathcal{N}(\mu, \sigma^2)$

- $Y = aX + b \implies Y \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$

Gaussian  
↓

- Sum of Gaussians  $\Rightarrow$  Set of Gaussian distributions is "closed" under summation.

- $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$

- $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$

- $Z = X + Y \implies Z \sim \mathcal{N}(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$

Gaussian  
↓

- [HW0 Questions A3 and A4]

# MLE for Gaussian

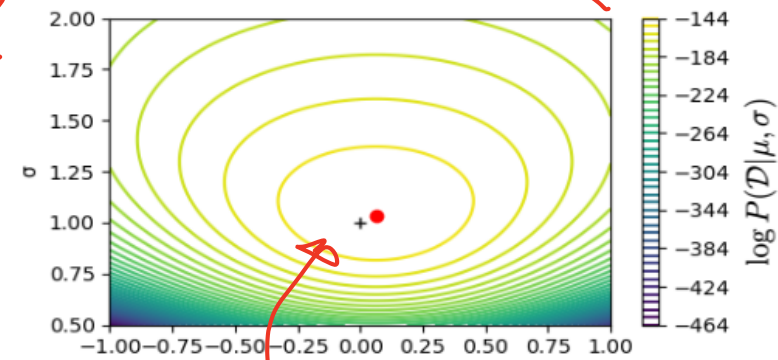
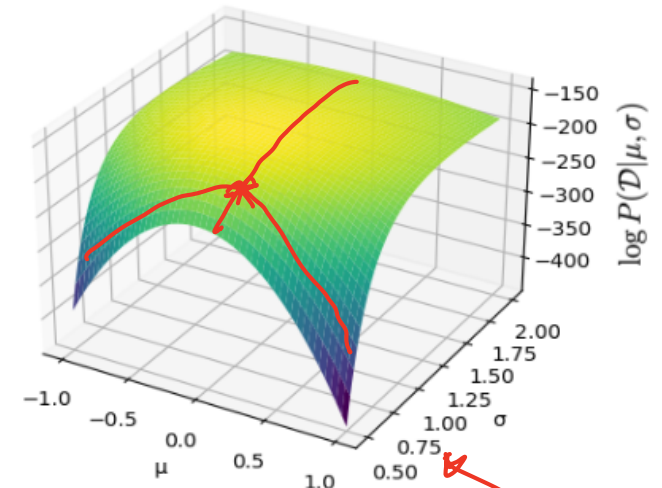
- Hypothesis:** i.i.d. samples  $\mathcal{D} = \{x_1, x_2, \dots, x_n\}$  from  $\mathcal{N}(\mu, \sigma^2)$

$$\begin{aligned}
 P(\mathcal{D}; \underbrace{(\mu, \sigma^2)}_{\theta}) &= P(x_1, \dots, x_n; \mu, \sigma^2) \\
 &= P(x_1; \mu, \sigma^2) \times P(x_2; \mu, \sigma^2) \times \dots \times P(x_n; \mu, \sigma^2) \\
 &= \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}
 \end{aligned}$$

- Log-likelihood of data:**

$$\begin{aligned}
 \log P(\mathcal{D}; \mu, \sigma^2) &= \sum_{i=1}^n \left\{ -\log(6\sqrt{2\pi}) - \frac{(x_i - \mu)^2}{2\sigma^2} \right\} \\
 &= -n \log(6\sqrt{2\pi}) - \frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}
 \end{aligned}$$

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



level set

+  $(\mu_{\text{True}}, \sigma_{\text{True}})$   
 •  $(\hat{\mu}_{\text{MLE}}, \hat{\sigma}_{\text{MLE}})$

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

- What's MLE for mean? Set partial derivative to zero:

$$\frac{d}{d\mu} \log P(\mathcal{D}; \mu, \sigma^2) = \frac{d}{d\mu} \left[ \underbrace{-n \log(\sigma\sqrt{2\pi})}_0 - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$= \frac{+ \sum_{i=1}^n \{+ 2(x_i - \mu)\}}{2\sigma^2} = 0$$

$$\sum_{i=1}^n x_i = n \cdot \mu$$

$$\hat{\mu}_{MLE} = \frac{1}{n} \cdot \sum_{i=1}^n x_i \quad \text{Empirical Mean}$$

# MLE for variance of a Gaussian distribution

- Again, set partial derivative to zero:

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

$$= -n \frac{1}{\sigma} + \sum_{i=1}^n \left\{ \frac{(x_i - \mu)^2}{\sigma^3} \right\} = 0$$

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

# What can we say about the MLE?

- MLE for the mean of a Gaussian is **unbiased**

$$\hat{\mu}_{\text{MLE}} = \frac{1}{n} \sum_{i=1}^n x_i \sim N(\mu^*, \sigma^2)$$

$$\mathbb{E}[\hat{\theta}] = \theta^*$$

$$\mathbb{E}[\hat{\mu}_{\text{MLE}}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[x_i] = \mu^*$$

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

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

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

# Maximum Likelihood Estimation

---

- **Observe**  $X_1, X_2, \dots, X_n$  drawn i.i.d. from  $P(X_i; \theta)$  for some true  $\theta = \theta^*$
- **Likelihood function:**  $L_n(\theta) = \prod_{i=1}^n P(X_i; \theta)$
- **Log-likelihood function:**  $\ell_n(\theta) = \log L_n(\theta) = \sum_{i=1}^n \log P(X_i; \theta)$
- **Maximum Likelihood Estimator (MLE):**  $\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \ell_n(\theta)$
  
- Properties (under benign regularity conditions—smoothness, identifiability, etc.):
  - MLE converges to the ground truths  $\theta^*$  as the number of samples  $n \rightarrow \infty$

# Linear Regression

---

# Maximum Likelihood Estimation

- **Observe**  $X_1, X_2, \dots, X_n$  drawn i.i.d. from  $P(X_i; \theta)$  for some true  $\theta = \theta^*$
- **Likelihood function:**  $L_n(\theta) = \prod_{i=1}^n P(X_i; \theta)$
- **Log-likelihood function:**  $\ell_n(\theta) = \log L_n(\theta) = \sum_{i=1}^n \log P(X_i; \theta)$
- **Maximum Likelihood Estimator (MLE):**  $\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \ell_n(\theta)$
- Why do we care about recovering the “true” parameter  $\theta^*$ ?
  - **Estimation** of the parameter  $\theta^*$  can be a goal.
  - Help **Interpret** or summarize large datasets.
  - Make **predictions** about future data.
  - **Generate** new data  $X \sim f(\cdot; \hat{\theta}_{\text{MLE}})$

# Estimation

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

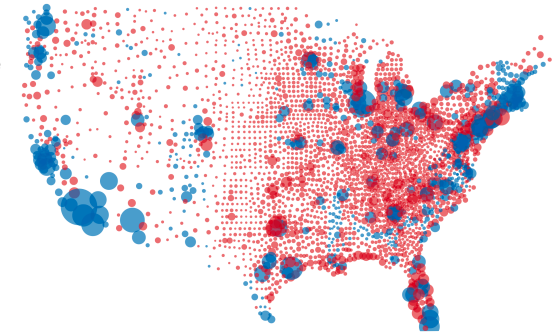
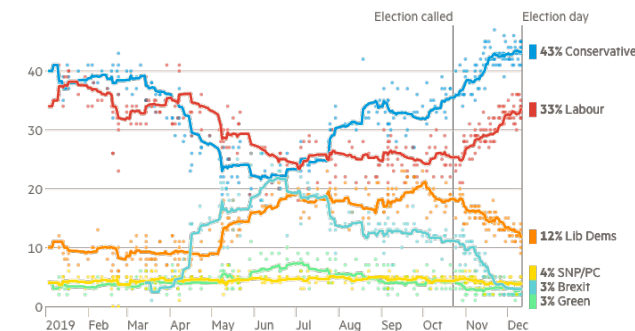
## Opinion polls

How does the greater population feel about an issue?  
Correct for over-sampling?

- $\theta_*$  is “true” average opinion
- $X_1, X_2, \dots$  are sample calls

UK poll tracker

Lines represent weighted averages, points represent polls (%)



## A/B testing

How do we figure out which ad results in more click-through?

- $\theta_*$  are the “true” average rates
- $X_1, X_2, \dots$  are binary “clicks”

The image shows two advertisement banners for Humana Medicare plans. The top banner, labeled 'Control', features a woman smiling and text that reads: 'Save on prescription drugs - over \$3,637\* a year!'. Below this, it says: 'Last year, Humana's Medicare Advantage plan members saved, on average, \$3,637\* on prescription drugs! Choose your Humana Medicare Advantage plan and you could enjoy savings on prescription drugs, plus:'. The bottom banner, labeled 'Treatment', features a man and a woman smiling and text that reads: 'Explore Humana's Medicare plans'. Below this, it says: 'Let us help you determine the Humana plan that's best for your needs.' and a 'Get started now' button.

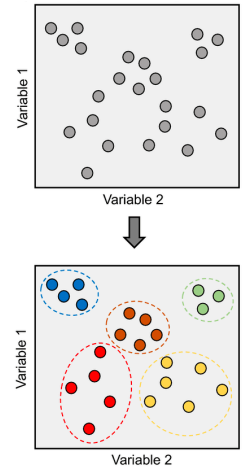
# Interpret

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

## Customer segmentation / clustering

Can we identify distinct groups of customers by their behavior?

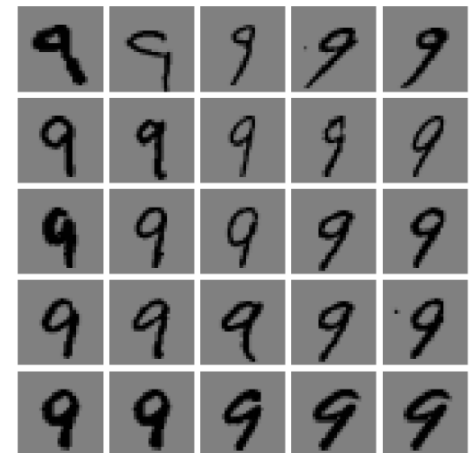
- $\theta_*$  describes “center” of distinct groups
- $X_1, X_2, \dots$  are individual customers



## Data exploration

What are the degrees of freedom of the dataset?

- $\theta_*$  describes the principle directions of variation
- $X_1, X_2, \dots$  are the individual images



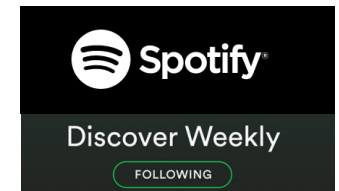
# Predict

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

## Content recommendation

Can we predict how much someone will like a movie based on past ratings?

- $\theta_*$  describes user’s preferences
- $X_1, X_2, \dots$  are (movie, rating) pairs



## Object recognition / classification

Identify a flower given just its picture?

- $\theta_*$  describes the characteristics of each kind of flower
- $X_1, X_2, \dots$  are the (image, label) pairs



(a)



(b)



(c)

Figure 1.1: Three types of Iris flowers: Setosa, Versicolor and Virginica. Used with kind permission of Dennis Krumb and SIGNA.

index	sl	sw	pl	pw	label
0	5.1	3.5	1.4	0.2	Setosa
1	4.9	3.0	1.4	0.2	Setosa
...					
50	7.0	3.2	4.7	1.4	Versicolor
...					
149	5.9	3.0	5.1	1.8	Virginica

# Generate

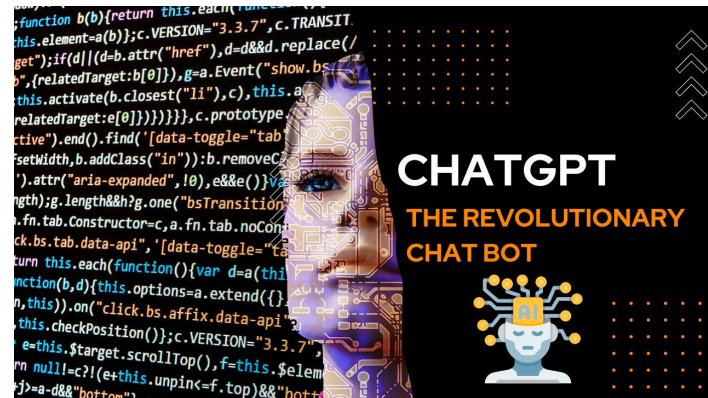
Observe  $X_1, X_2, \dots, X_n$  drawn IID from  $f(x; \theta)$  for some “true”  $\theta = \theta_*$

## Text generation

Can AI generate text that could have been written like a human?

- $\theta_*$  describes language structure
- $X_1, X_2, \dots$  are text snippets found online

“Kaia the dog wasn't a natural pick to go to mars.  
No one could have predicted she would...”



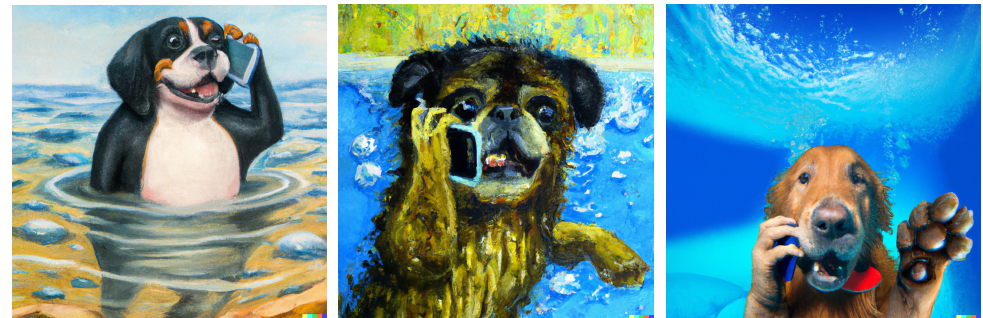
<https://chat.openai.com/chat>

## Image to text generation

Can AI generate an image from a prompt?

- $\theta_*$  describes the coupled structure of images and text
- $X_1, X_2, \dots$  are the (image, caption) pairs found online

“dog talking on cell phone under water, oil painting”



<https://labs.openai.com/>

# CSE 446/546

---

- week 1: **Estimation**
  - Maximum Likelihood Estimation
- week 1~8: **Prediction**
  - week 1~4: Linear regression models
  - week 4~5: Linear classification models (also called Logistic regression)
  - Midterm Exam
  - week 6~7: Non-linear models
- week 8~9: **Interpretation**
- week 10: **Generation... (?)** → 4935/599 Spring 26.

# Linear regression model, 1-dimensional

You want to sell your house that is 2,500 sq.ft.

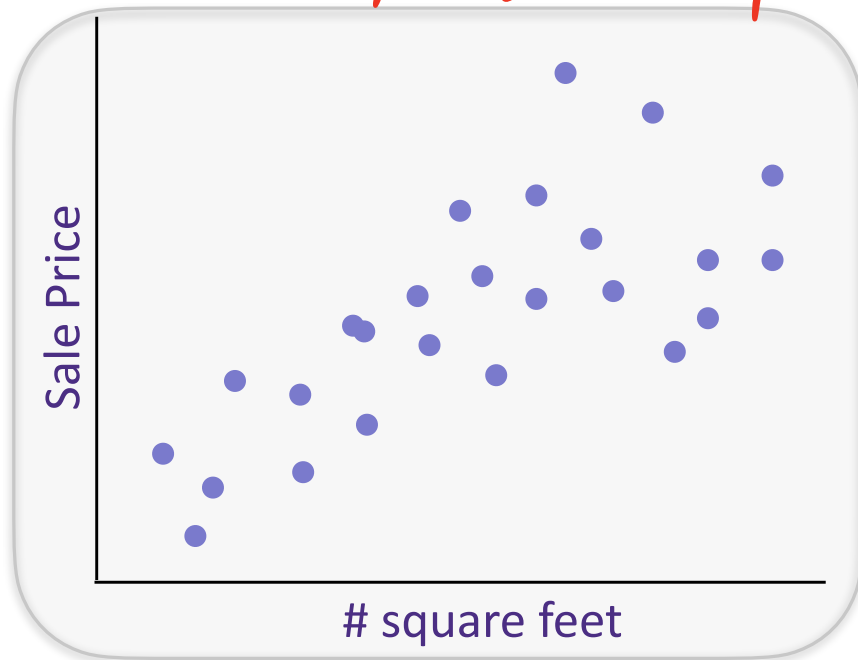
Q. What is the right price?

Collect past sales data on [zillow.com](https://www.zillow.com):

$y =$  House sale price and  $x = \{\# \text{ sq. ft.}\}$

*Label, dependant var, response*

*independant var, input,*



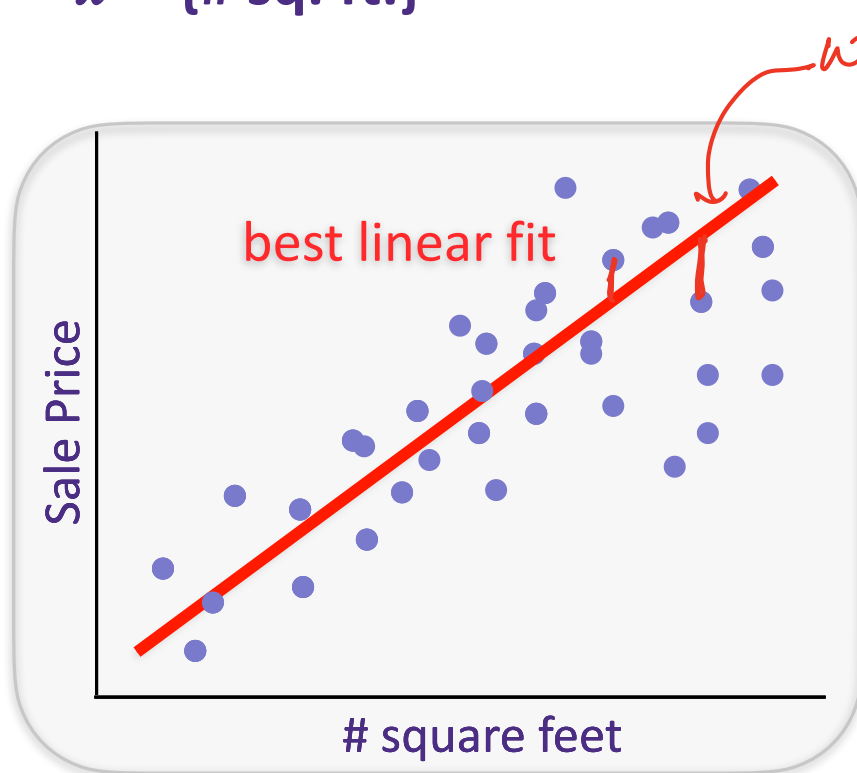
Training Data:  $x_i \in \mathbb{R}$   $y_i \in \mathbb{R}$   
 $\{(x_i, y_i)\}_{i=1}^n$

# Linear regression model, 1-dimension

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

$y =$  House sale price

$x =$  {# sq. ft.}



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

2. Hypothesis/Model: linear *model*

$$y_i = w \cdot x_i + \epsilon_i$$

$\uparrow$   
 $\mathbb{R}$   
model  
param.

$\uparrow$   
 $\mathbb{R}$   
noise

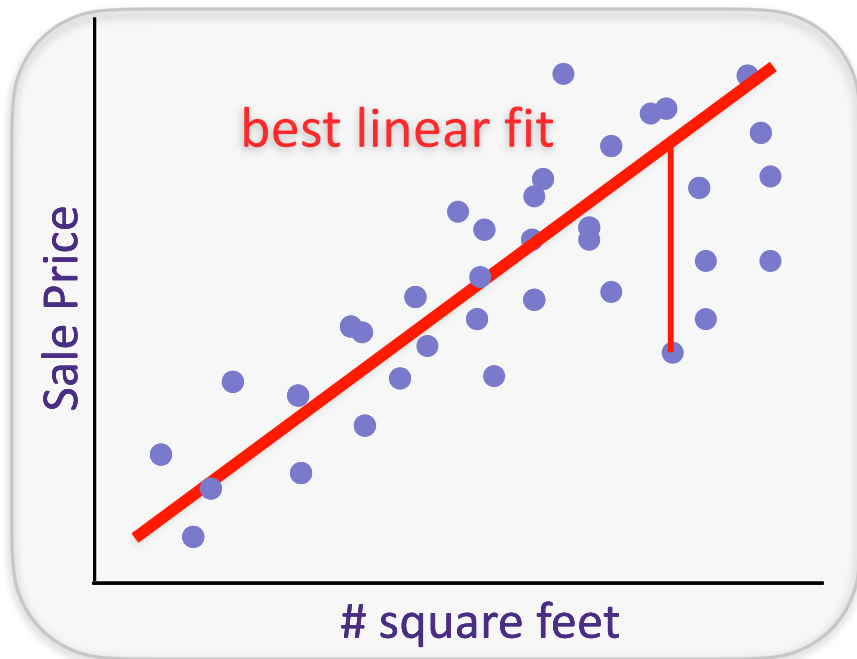
For now we assume there is no y-intercept in the model, and will handle it later

# Linear regression model, 1-dimension

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

$y =$  House sale price

$x =$  {# sq. ft.}



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

2. Hypothesis/Model: linear

$$y_i = w \cdot x_i + \epsilon_i$$

3. Noise: i.i.d. Gaussian with

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

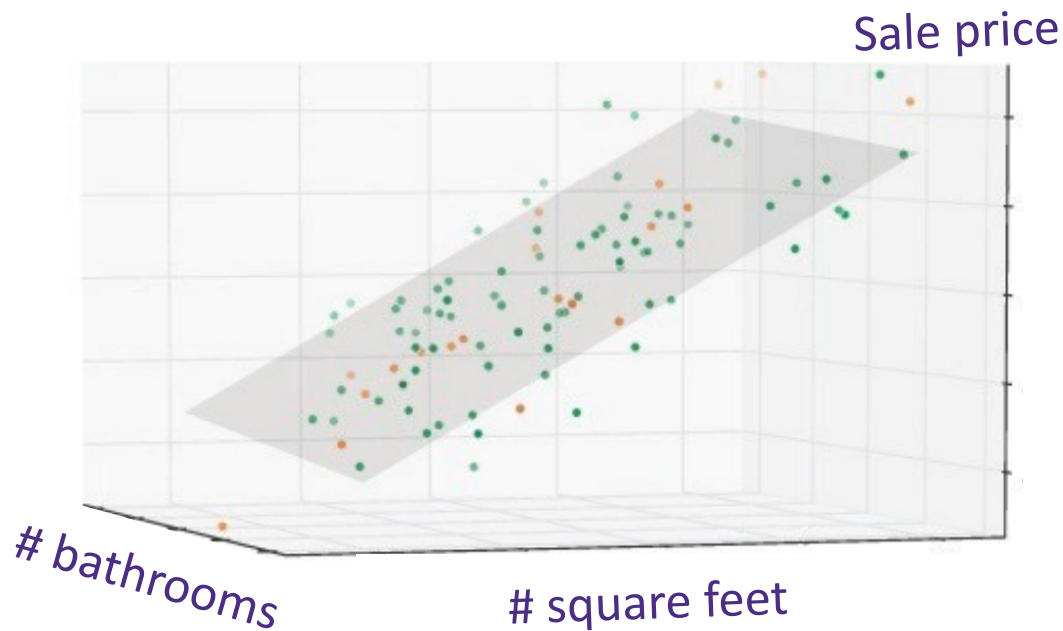
For now we assume there is no  $y$ -intercept in the model, and will handle it later

# Linear regression model, d-dim

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

$y =$  House sale price

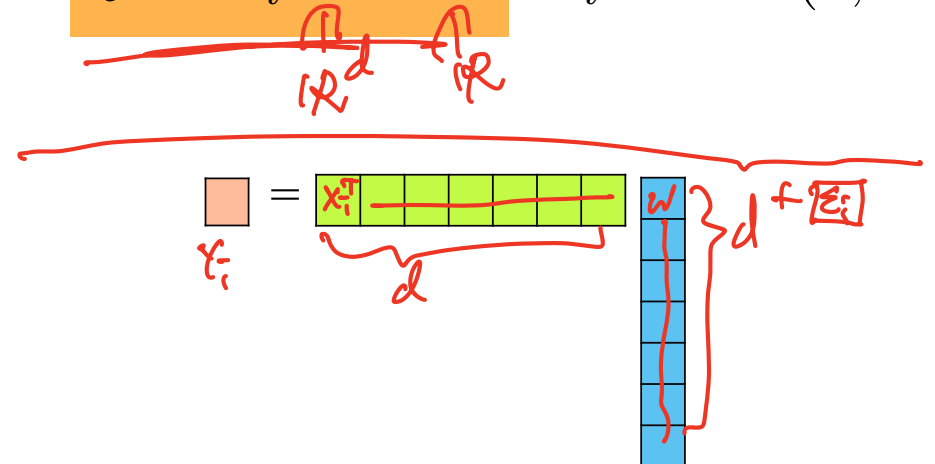
$x = \{\# \text{ sq. ft.}, \text{zip code}, \text{date of sale}, \text{etc.}\}$



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

Hypothesis/Model: linear

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

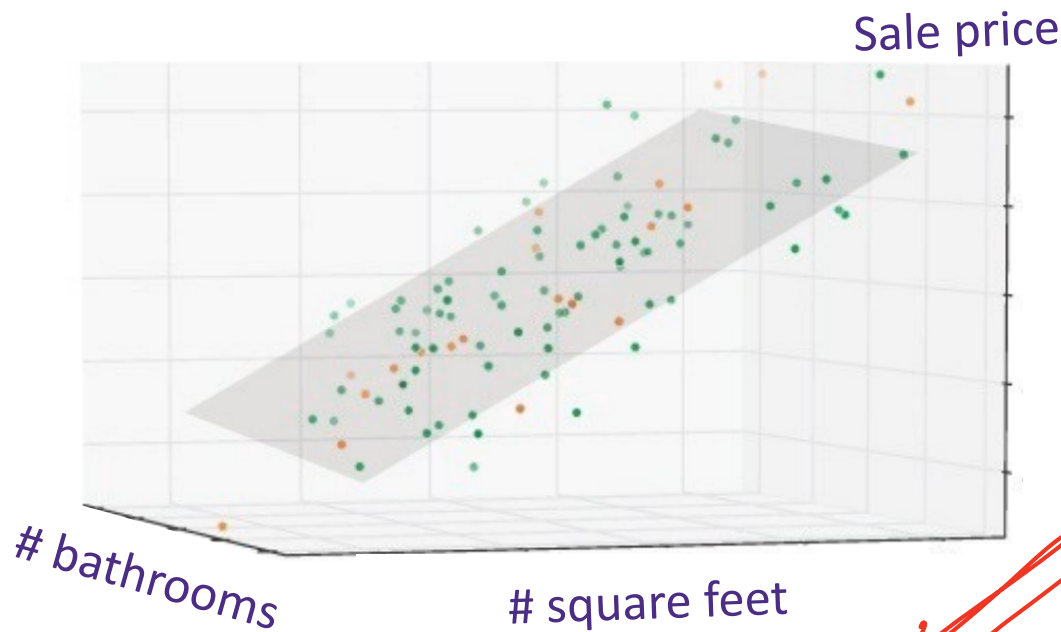


# Linear regression model, d-dim

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

$y =$  House sale price

$x =$  {# sq. ft., zip code, 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

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

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

$$P(y; x, w, \sigma)$$

$$P(Y, X, w, \sigma)$$

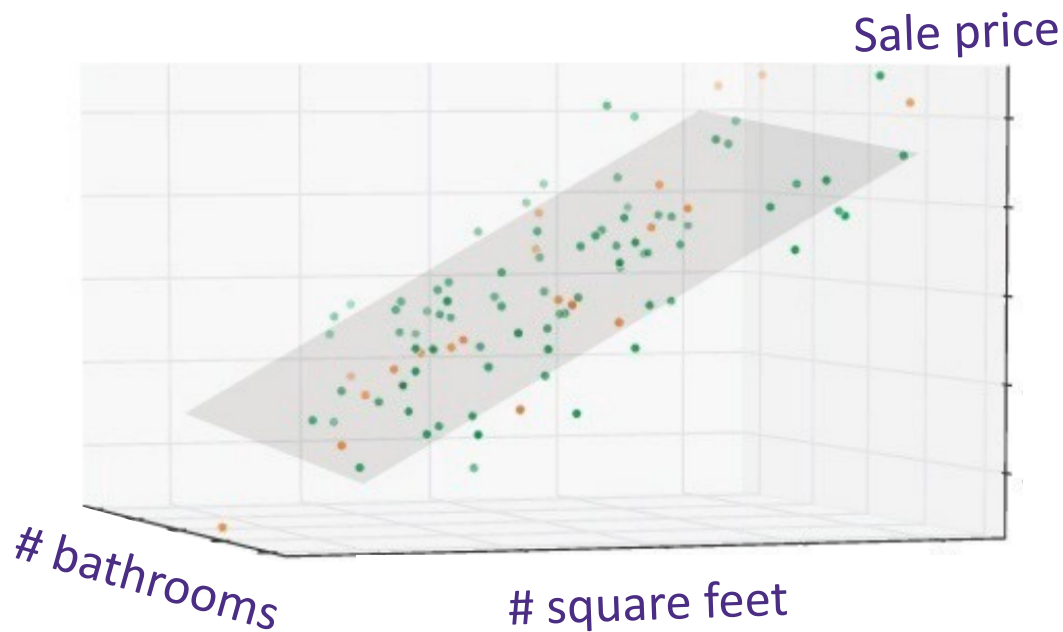
$$\frac{P(Y, X, w, \sigma)}{P(X, w, \sigma)} \leftarrow \text{Bayes' Rule}$$

# Linear regression model, d-dim

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

$y =$  House sale price

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



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

Hypothesis/Model: linear

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

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

# Maximizing log-likelihood

**Training Data:**  $x_i \in \mathbb{R}^d$   
 $\{(x_i, y_i)\}_{i=1}^n$   $y_i \in \mathbb{R}$  *Random*  $\downarrow$   $p(y|x, w, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-x^\top w)^2/2\sigma^2}$   
*noise*

$$y_i = x_i^\top w + \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

*↳ model parameter.*

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

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

# Maximum Likelihood Estimation

- **Observe**  $X_1, X_2, \dots, X_n$  drawn i.i.d. from  $P(X_i; \theta)$  for some true  $\theta = \theta^*$
- **Likelihood function:**  $L_n(\theta) = \prod_{i=1}^n P(X_i; \theta)$
- **Log-likelihood function:**  $\ell_n(\theta) = \log L_n(\theta) = \sum_{i=1}^n \log P(X_i; \theta)$
- **Maximum Likelihood Estimator (MLE):**  $\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \ell_n(\theta)$
- Properties (under benign regularity conditions—smoothness, identifiability, etc.):
  - MLE converges to the ground truths  $\theta^*$  as the number of samples  $n \rightarrow \infty$

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

$\hat{w}_{MLE} = \underset{w}{\text{arg max}} \sum_{i=1}^n -\log(\sqrt{2\pi\sigma^2}) \Rightarrow \frac{(y_i - x_i^\top w)^2}{2\sigma^2}$   
 ignore ignore

$\hat{w}_{MLE} = \underset{w}{\text{arg min}} \sum_{i=1}^n (y_i - x_i^\top w)^2$  } changes max/min  
 Mean Squared Error

# 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

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

$$\nabla_w l(w)$$

$$d \left[ \sum_{i=1}^n \frac{\partial l}{\partial w_i} \right]$$

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

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

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

$$w = A^{-1} \cdot b$$



Set gradient=0, solve for w

$$b = A \cdot w$$

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

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

$$b = A^{-1} \cdot w$$

$$A^T \cdot A = A \cdot A^T = I$$

$$A^{-1} b = A^{-1} A w$$

$$= w$$

$n \geq d$

# 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

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

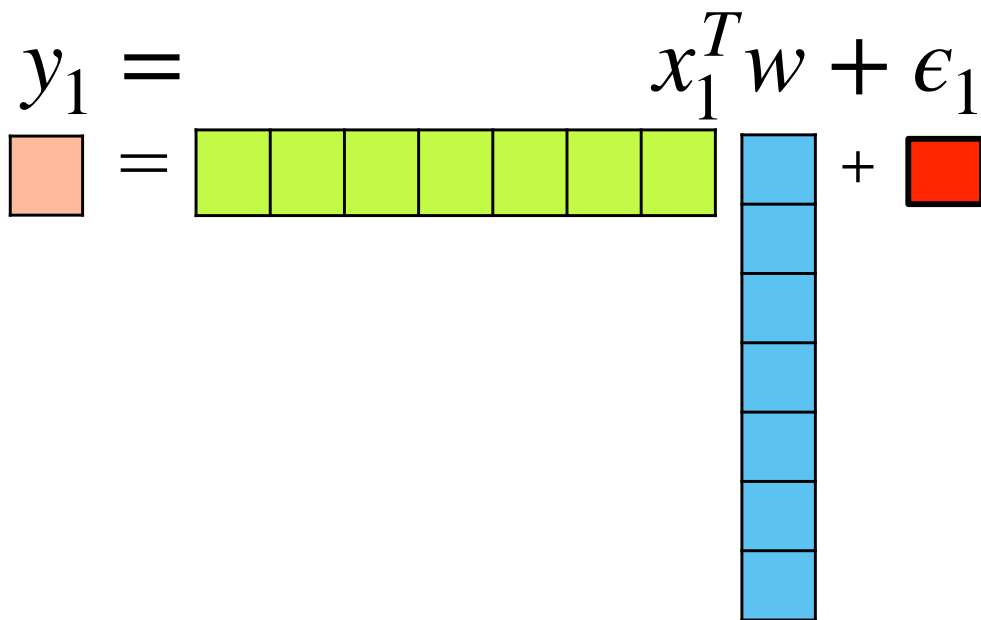
$$\boxed{\quad}^{-1} \cdot \boxed{\quad}$$

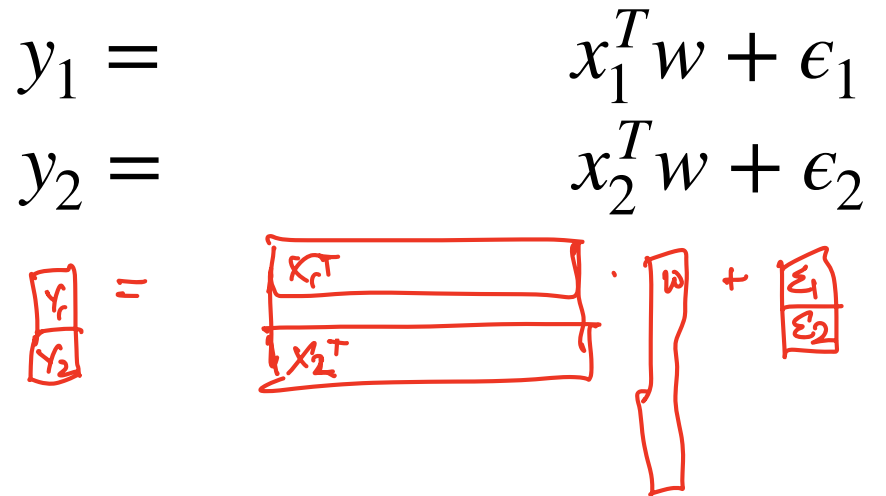
# The regression problem in matrix notation

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/size of the input  
 $n$  : # of examples/datapoints

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


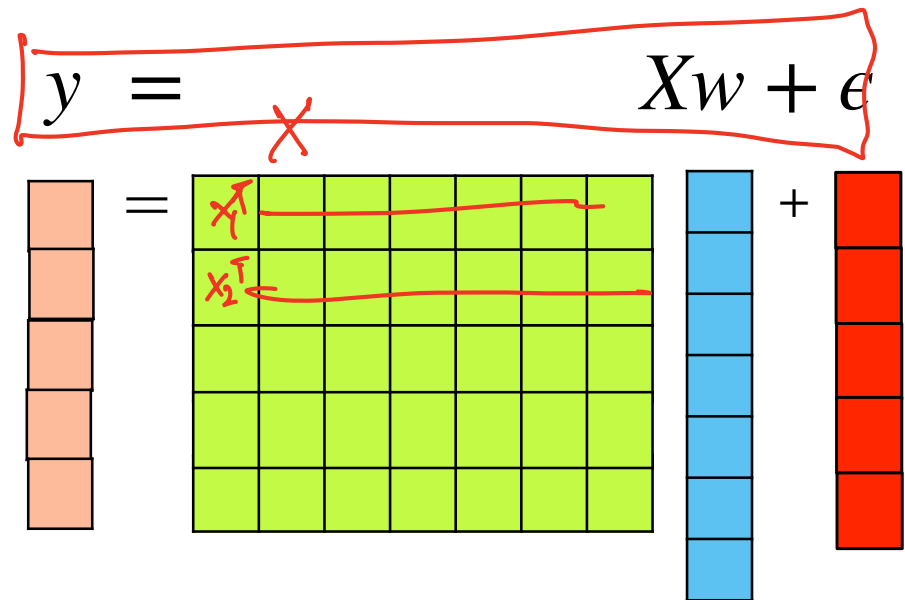
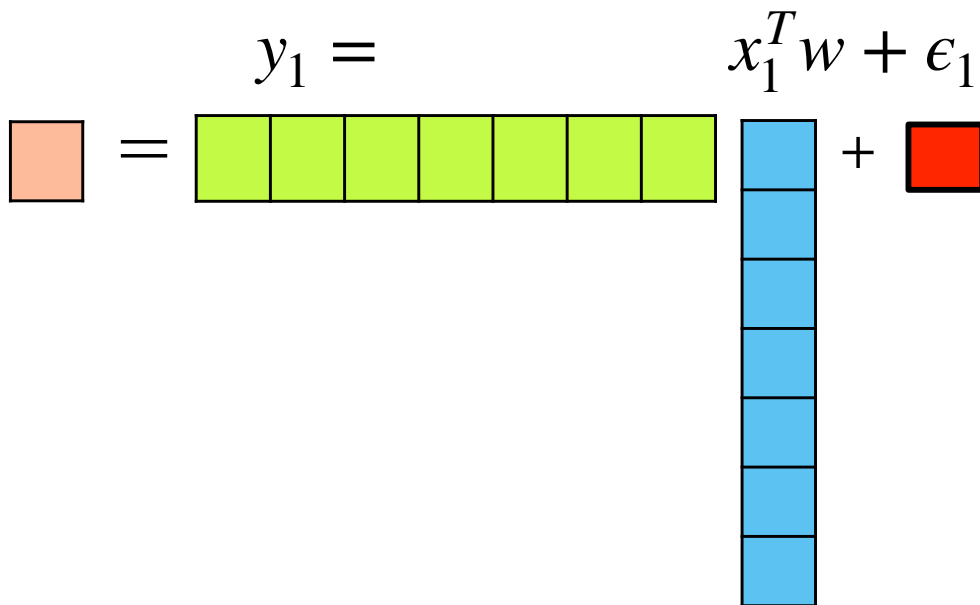
$$\begin{aligned} y_1 &= x_1^T w + \epsilon_1 \\ y_2 &= x_2^T w + \epsilon_2 \end{aligned}$$


# 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/size of the input  
 $n$  : # of examples/datapoints



# The regression problem in matrix notation

$$\hat{w}_{MLE} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 \quad \leftarrow (y - Xw)_{i\text{-th}}$$

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

$d$  : # of features

$n$  : # of examples/datapoints

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

*error vector*

*Square*

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

*$\ell_2$ -norm /  $\ell_p$ -norm*

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

$$\|z\|_2^2 = z^T z$$

# The regression problem in matrix notation

---

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

# The regression problem in matrix notation

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

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

$d$  : # of features

$n$  : # of examples/datapoints

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

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

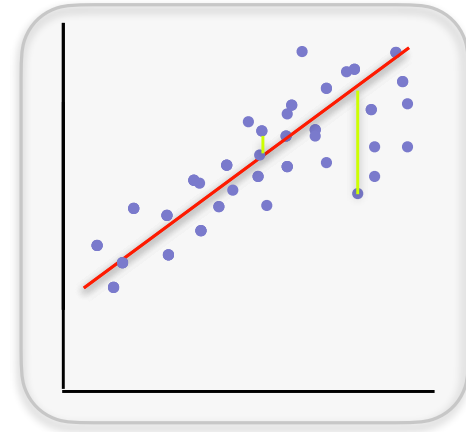
$$\hat{w}_{LS} = \hat{w}_{MLE} = \left( \mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

# The regression problem in matrix notation

Recall that we start with a linear model with no offset

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

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



We can add the offset to the linear model, with a new parameter  $b$

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

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

# Dealing with an offset

---

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

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

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

# 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 or we pre-processed the data have zero-mean, 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?

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

# Questions?

---