

Lecture 2: MLE for Gaussian and linear regression

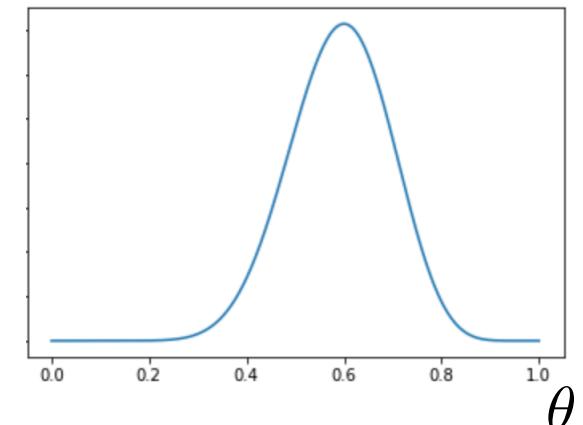
Sewoong Oh

Ed Discussion:

W

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
- $$\log P(\mathcal{D}; \theta) = \log \left(\theta^k (1-\theta)^{n-k} \right)$$
- $$= k \cdot \log \theta + (n-k) \log (1-\theta)$$
- $X_i \sim \text{Binomial}(\theta)$
- $P(\mathcal{D}; \theta)$ when k heads observed in n flips



Recap: Maximum Likelihood Estimation

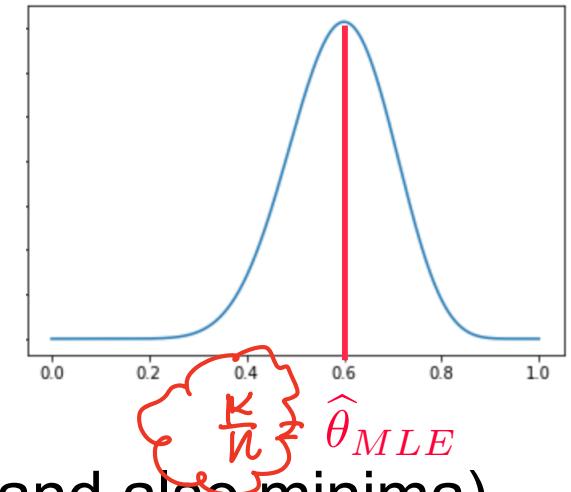
- **Observe** $\mathcal{D} = X_1, X_2, \dots, X_n$ drawn i.i.d. from $P(X_i; \theta)$ for some ground truth $\theta = \theta^*$, unknown to us
- **Maximize log-likelihood** when we observe k heads in n flips

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

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

- Maximum likelihood estimate (MLE)

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \{ k \lg \theta + (n-k) \lg(1-\theta) \}$$



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

and find θ 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 $\theta = \theta^*$

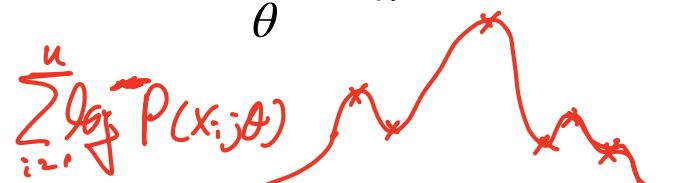
- Likelihood function: $L_n(\theta) = \prod_{i=1}^n P(X_i; \theta)$

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

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

- 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



What about continuous variables?

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

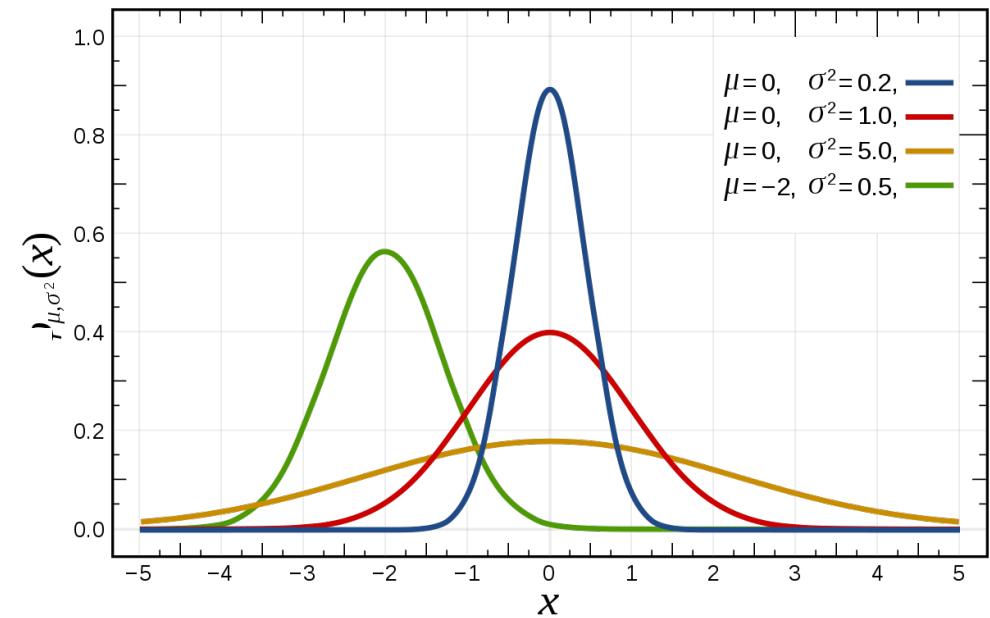
- A Gaussian random variable is written as $X \sim \mathcal{N}(\mu, \sigma^2)$ with mean $\mu \triangleq \mathbb{E}[X]$ and variance $\sigma^2 \triangleq \mathbb{E}[(X - \mathbb{E}[X])^2]$

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

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

$$P(x | \mu, \sigma^2) = \frac{P(x, \mu, \sigma^2)}{P(\mu, \sigma^2)}$$

• Conditional on Bayes Rule
• Given



Binomial (θ)
Gaussian (μ, σ^2)

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

Some useful properties of Gaussians

- Affine transformation
(multiplying by scalar and adding a constant)
 - $X \sim \mathcal{N}(\mu, \sigma^2)$
 - $Y = aX + b \implies Y \sim \mathcal{N}(\underbrace{a\mu + b}_{\text{mean}}, \underbrace{a^2\sigma^2}_{\text{variance}})$
- Sum of Gaussians , Set of Gaussian Distributions is "closed under summation"
 - $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$
 - $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$
 - $Z = X + Y \implies Z \sim \mathcal{N}(\underbrace{\mu_X + \mu_Y}_{\text{mean}}, \underbrace{\sigma_X^2 + \sigma_Y^2}_{\text{variance}})$
- [HW0 Questions A3 and A4]

MLE for Gaussian

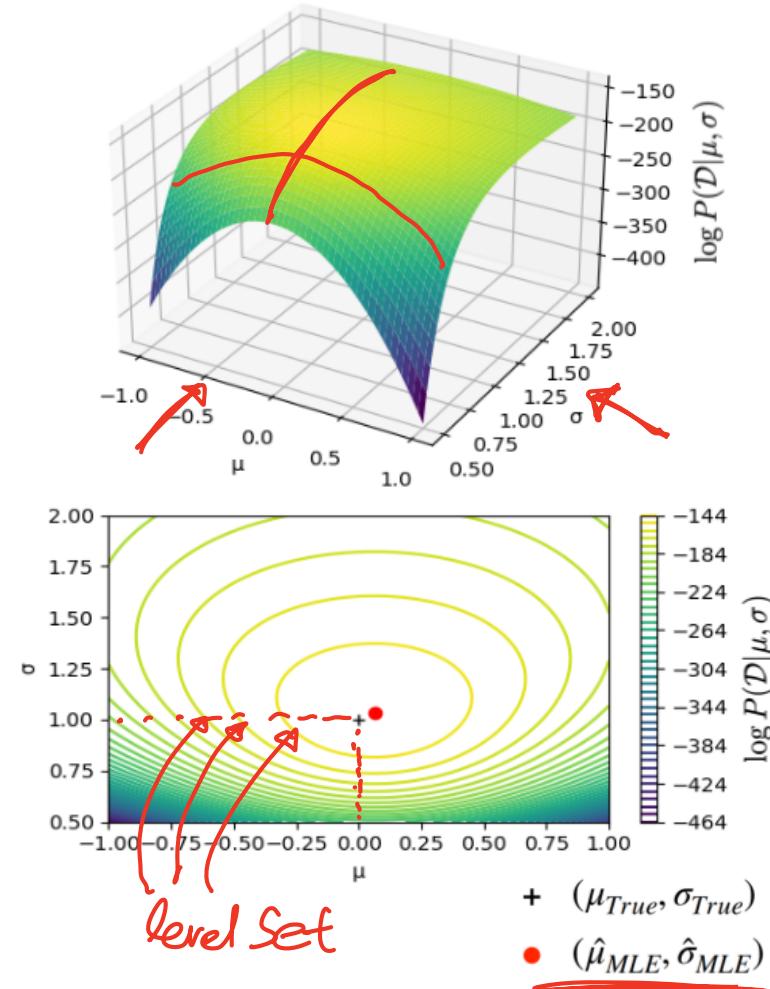
- **Hypothesis:** i.i.d. samples $\mathcal{D} = \{x_1, x_2, \dots, x_n\}$ from $\mathcal{N}(\mu, \sigma^2)$

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

- **Log-likelihood** of data:

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

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



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

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

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

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

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

: Empirical mean

Does not depend on σ .

MLE for variance of a Gaussian distribution

- Again, set partial derivative to zero:

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

$$= \frac{-n}{6} + \sum_{i=1}^n \frac{(x_i - \mu)^2}{6^3} = 0$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 : \text{Empirical Variance.}$$

↓
plug in $\hat{\mu}_{MLE}$

What can we say about the MLE?

- MLE for the mean of a Gaussian is unbiased

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

$$\begin{aligned}\mathbb{E}[\hat{\theta}] &= \theta^* \\ x_i &\sim N(\mu^*, \sigma^2) \\ \mathbb{E}[\hat{\mu}_{\text{MLE}}] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[x_i] = \mu^*\end{aligned}$$

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

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

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

$$\begin{aligned}\mathbb{E}[\hat{\sigma}_{\text{MLE}}^2] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(x_i - \hat{\mu})^2] \\ &= \frac{1}{n} \sum_{i=1}^n \left\{ \mathbb{E}[x_i^2] - 2\mathbb{E}[x_i]\hat{\mu} + \hat{\mu}^2 \right\} \\ &= \sigma^2 + (\mu - \hat{\mu})^2\end{aligned}$$

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 θ^* as the number of samples $n \rightarrow \infty$

Linear Regression

UNIVERSITY *of* WASHINGTON



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 θ^* ?
 - **Estimation** of the parameter θ^* can be a goal.
 - Help **Interpret** or summarize large datasets.
 - Make **predictions** about future data.
 - Generate new data $X \sim f(\cdot; \hat{\theta}_{\text{MLE}})$

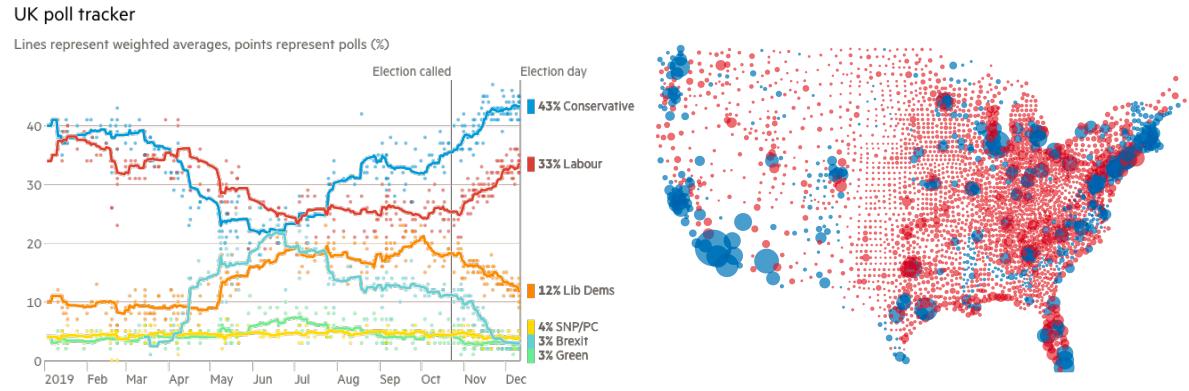
Estimation

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?

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



A/B testing

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

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



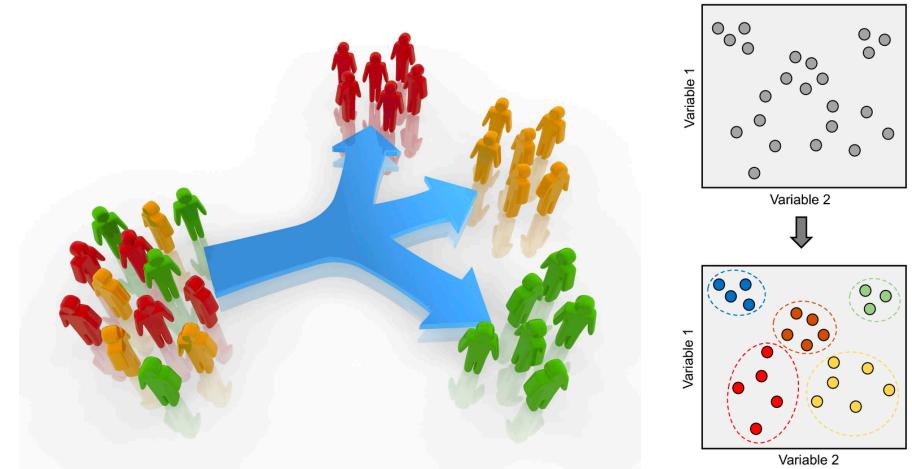
Interpret

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?

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



Data exploration

What are the degrees of freedom of the dataset?

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

9	9	9	9	9
9	9	9	9	9
9	9	9	9	9
9	9	9	9	9
9	9	9	9	9

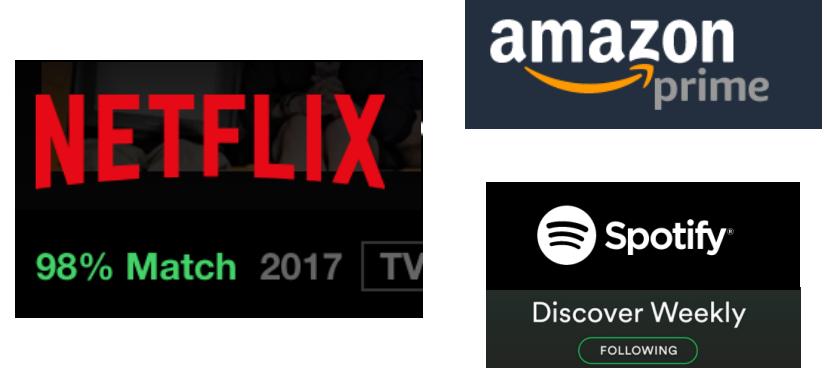
Predict

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?

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



Object recognition / classification

Identify a flower given just its picture?

- θ_* 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 Kramb 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?

- θ_* 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...”

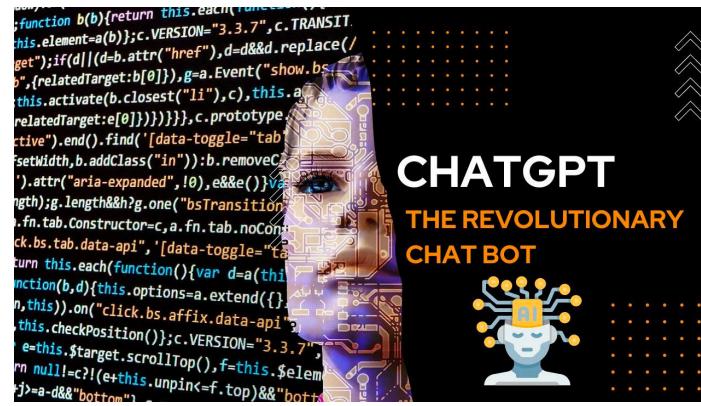


Image to text generation

Can AI generate an image from a prompt?

- θ_* 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
- week 6~7: Non-linear models
- week 8~9: **Interpretation**
- week 10: **Generation...?** → CSE 493S/599 2026sp

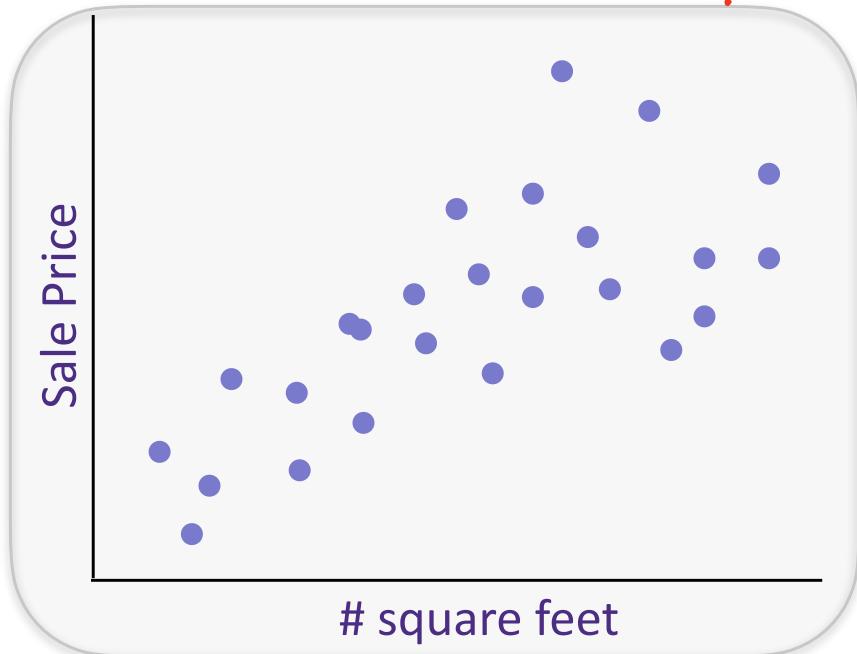
Linear regression model, 1-dimensional

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, response, dependent variable / Covariate, input variable, independent variable



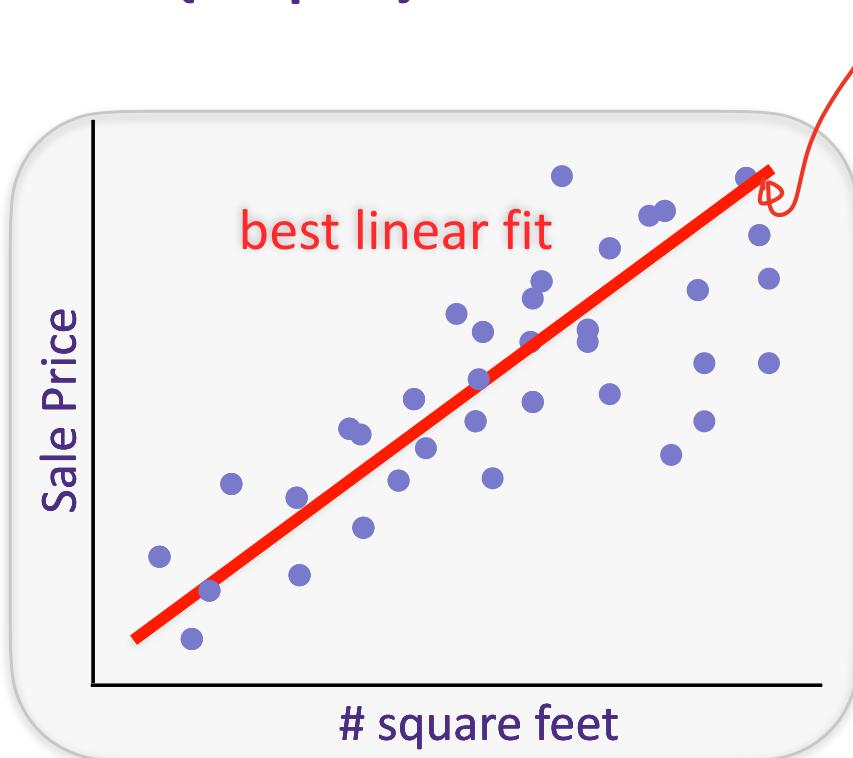
Training Data: $x_i \in \mathbb{R}$ $y_i \in \mathbb{R}$
 $\{(x_i, y_i)\}_{i=1}^n$

Linear regression model, 1-dimension

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

y = House sale price from

x = {# sq. ft.}



w : slope

1. Training Data:

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

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

2. Hypothesis/Model: linear *model*

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

\uparrow
IR

Parameter
of the model

\uparrow
IR

Noise iid

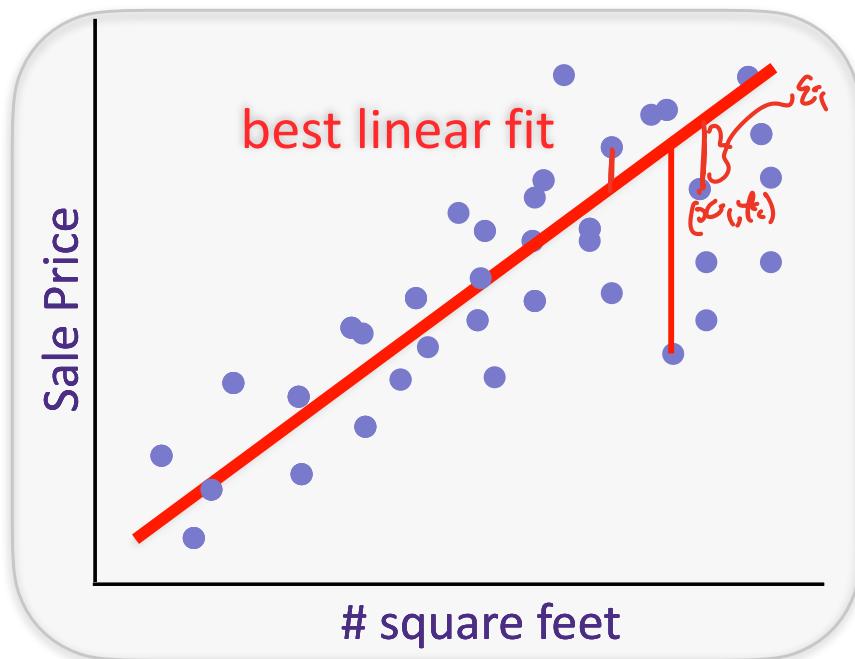
For now we assume there is no y -intercept in the model, and will handle it later

Linear regression model, 1-dimension

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

y = House sale price from

x = {# sq. ft.}



1. Training Data:

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

2. Hypothesis/Model: linear

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

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

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

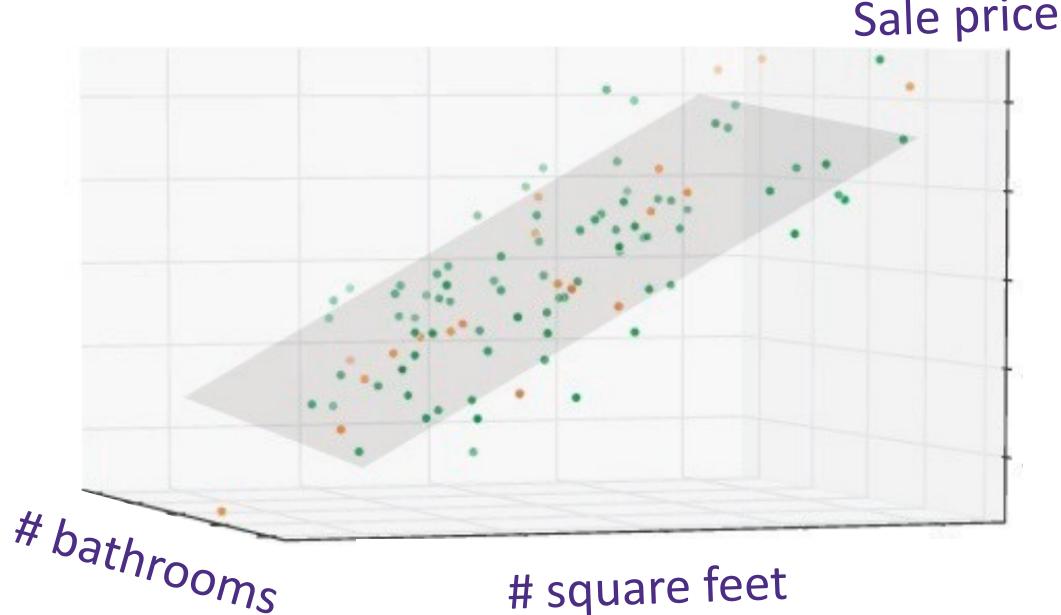
For now we assume there is no y -intercept in the model, and will handle it later

Linear regression model, d-dim

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

y = House sale price from

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



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

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

Hypothesis/Model: linear

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

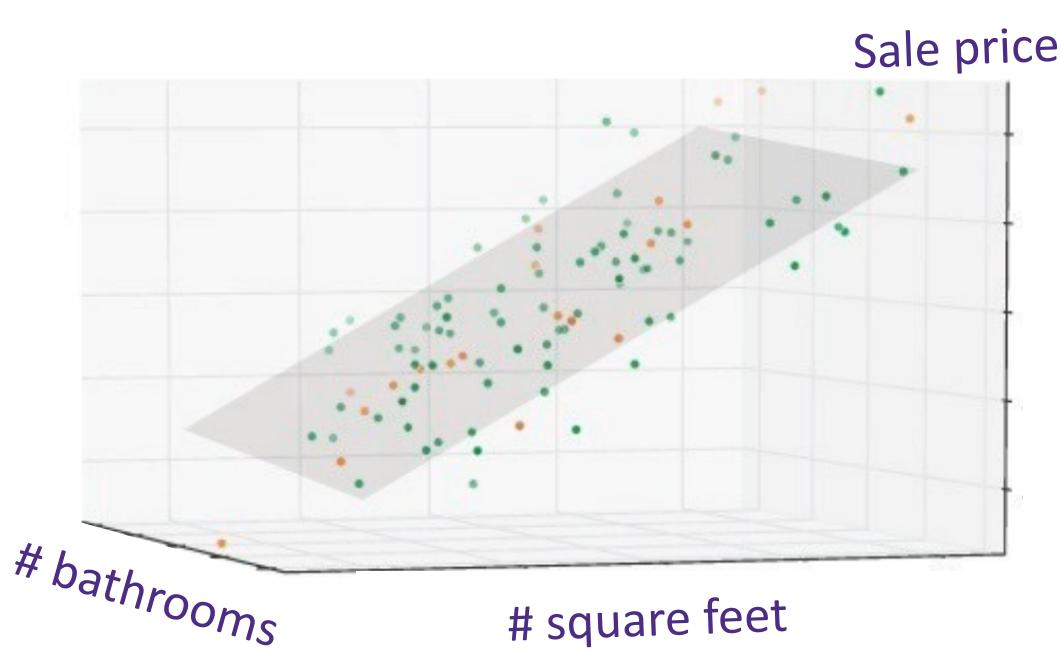
$$\begin{aligned} y_i &= \underbrace{\text{[} \text{]}}_{\mathbb{R}^d} \xrightarrow{x_i^T} \underbrace{\text{[} \text{]}}_{\mathbb{R}^d} + \underbrace{\epsilon_i}_{\mathbb{R}} \\ &= \underbrace{\text{[} \text{]}}_{\mathbb{R}^d} \xrightarrow{w^T} \underbrace{\text{[} \text{]}}_{\mathbb{R}^d} + \underbrace{\epsilon_i}_{\mathbb{R}} \end{aligned}$$

Linear regression model, d-dim

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

y = House sale price from

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

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

$\sim \mathcal{N}(x_i^T w, \sigma^2)$

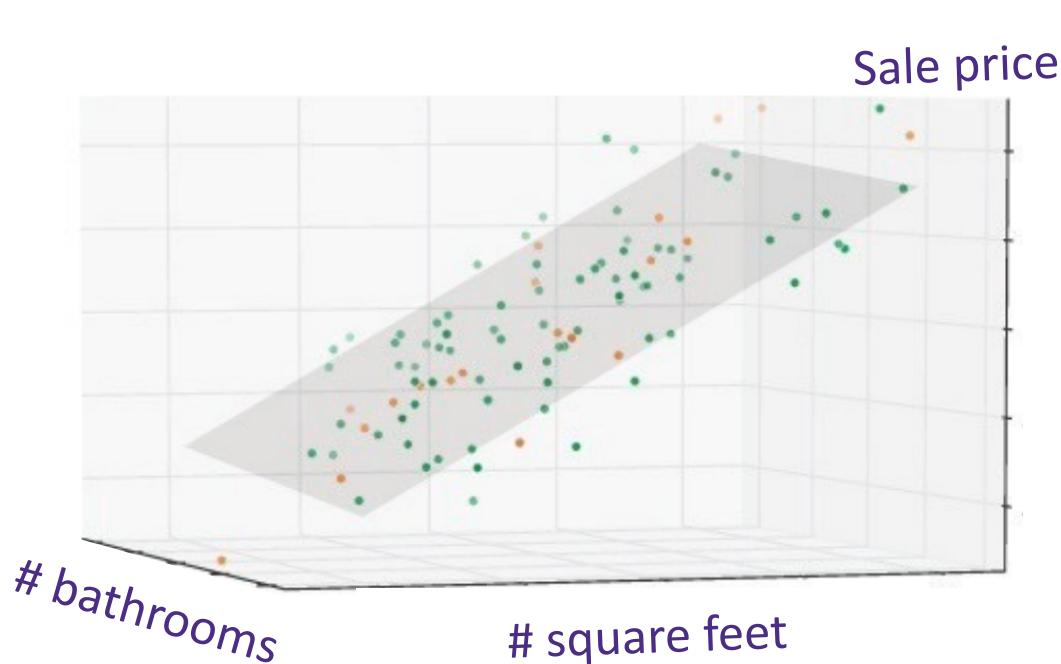
$$\begin{aligned} p(y|x, w, \sigma) &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{\epsilon_i^2}{2\sigma^2}} \\ &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - x_i^T w)^2}{2\sigma^2}} \end{aligned}$$

Linear regression model, d-dim

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

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

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



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

Hypothesis/Model: linear

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

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

Maximizing log-likelihood

Training Data: $x_i \in \mathbb{R}^d$ $y_i \in \mathbb{R}$

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

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

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

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

ignore $\log \sqrt{2\pi\sigma^2}$

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

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

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

* MLE is minimizing squared error

* Different noise results in different loss

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

$$\frac{\partial f}{\partial w_i} = \underbrace{-\sum_{i=1}^n (y_i - x_i^\top w) \cdot x_i}_{\text{i-th row}}$$

if $\sum_{i=1}^n x_i x_i^\top$ is invertible.

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

$$\boxed{\sum x_i x_i^\top} \cdot \boxed{w} = \sum \boxed{y_i x_i}$$

$$w = \left(\sum_{i=1}^n x_i x_i^\top \right)^{-1} \cdot \sum_{j=1}^n y_j x_j$$

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

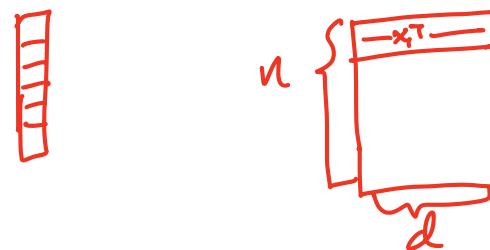
$$\hat{w}_{MLE} = \left(\sum_{i=1}^n x_i x_i^\top \right)^{-1} \sum_{i=1}^n x_i y_i$$

The regression problem in matrix notation

Data:

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

d : # of features/size of the input
n : # of examples/datapoints



$$y_1 =$$

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

A diagram illustrating the regression equation $y_1 = \mathbf{x}_1^T \mathbf{w} + \epsilon_1$. On the left, a red square labeled y_1 is followed by a green horizontal bar divided into segments. To its right is a blue vertical bar representing \mathbf{w} . An arrow points from the green bar to the blue bar, indicating multiplication. A red square labeled ϵ_1 is shown being added to the result.

$$y_2 =$$

$$= \begin{bmatrix} y_1 \\ y_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \end{bmatrix}$$

$$\begin{aligned} & \mathbf{x}_1^T \mathbf{w} + \epsilon_1 \\ & \mathbf{x}_2^T \mathbf{w} + \epsilon_2 \\ & \vdots \end{aligned}$$

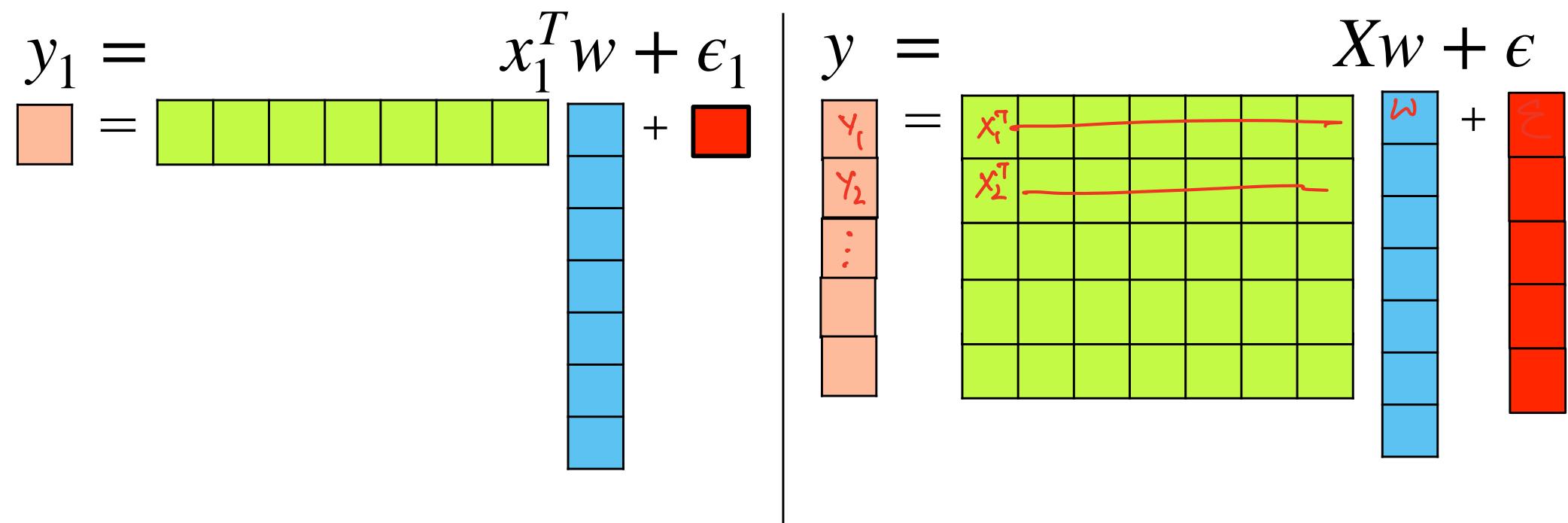
A diagram showing the stacked form of the regression equations. On the left, two equations are shown: $y_1 = \mathbf{x}_1^T \mathbf{w} + \epsilon_1$ and $y_2 = \mathbf{x}_2^T \mathbf{w} + \epsilon_2$. To the right, a vertical stack of red squares labeled y is shown, with a bracket indicating it is a vector. Below it, a vertical stack of red rectangles labeled \mathbf{x}^T is shown, with a bracket indicating it is a matrix. Red arrows point from the terms in the equations to the corresponding parts in the vector and matrix forms.

The regression problem in matrix notation

Data:

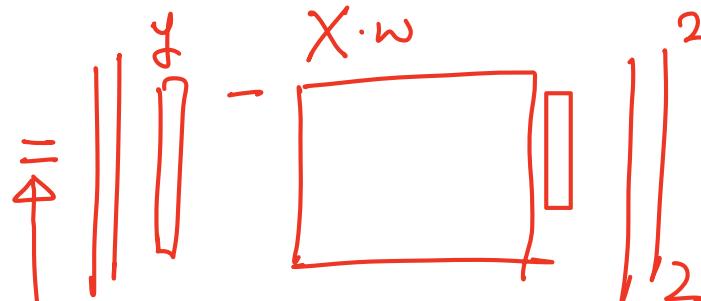
$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}$$

d : # of features/size of the input
n : # of examples/datapoints



The regression problem in matrix notation

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



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

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

$$\|z\|_2 = \sqrt{\sum_{i=1}^d z_i^2}$$

Square of a scalar
we use ℓ_2 -norm
 $\|z\|_p^2$: ℓ_p -norm.

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

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

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

[related to HW0 questions A6 and A7]

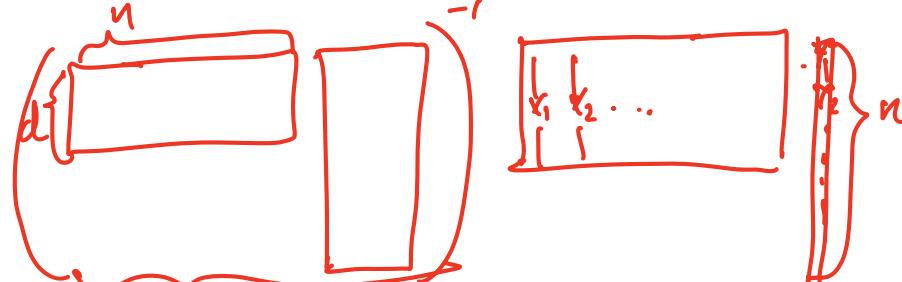
The regression problem in matrix notation

$$= \arg \min_w (\mathbf{y} - \mathbf{X}w)^T (\mathbf{y} - \mathbf{X}w)$$
$$= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}w + w^T \mathbf{X}^T \mathbf{X}w.$$

$$\nabla_w f(w) = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}w = 0$$

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

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

$$\left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \left(\sum_{i=1}^n \mathbf{x}_i^T \mathbf{y}_i \right)$$


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

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

$$\nabla_w w^T A w = 2A \cdot w$$

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

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

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

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

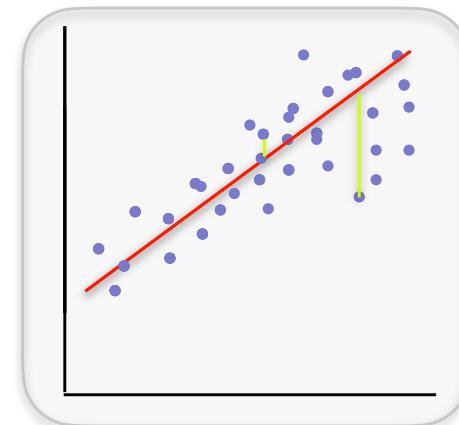
$$\hat{w}_{LS} = \hat{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

The regression problem in matrix notation

Recall that we start with a linear model with no offset

$$y_i = x_i^T w + \epsilon_1$$

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



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

$$\theta = (w \in \mathbb{R}^d, b \in \mathbb{R})$$

$$y_i = x_i^T w + b + \epsilon_1$$

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

all 1's vector.

Dealing with an offset

$$\mathbf{y}^T \mathbf{y} - 2 \mathbf{y}^T \mathbf{X} w + 2 \mathbf{y}^T \mathbf{1} b + w^T \mathbf{X}^T \mathbf{X} w + b^T \mathbf{1}^T \mathbf{1} = f(w, b)$$

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

$$\nabla_w f(w, b) = 0 \rightarrow \mathbf{X}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{X}^T \mathbf{1} = \mathbf{X}^T \mathbf{y}$$

$$\nabla_b f(w, b) = 0 \rightarrow \mathbf{1}^T \mathbf{X} \hat{w}_{LS} + \hat{b}_{LS} \mathbf{1}^T \mathbf{1} = \mathbf{1}^T \mathbf{y}$$

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

$$\tilde{\mathbf{X}} = [\mathbf{X} \ \mathbf{1}]$$

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

$$\tilde{w} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

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

$$\begin{aligned}\hat{w}_{LS} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ \hat{b}_{LS} &= \frac{1}{n} \sum_{i=1}^n y_i \quad \leftarrow \text{Avg Response.}\end{aligned}$$

Make Predictions

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

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

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

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

Questions?
