

CSE 446/546: Machine Learning

Simon Du and Sewoong Oh

W

Traditional algorithms

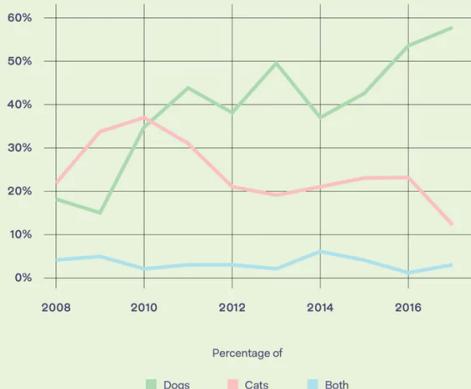
Social media mentions of Cats vs. Dogs

Reddit

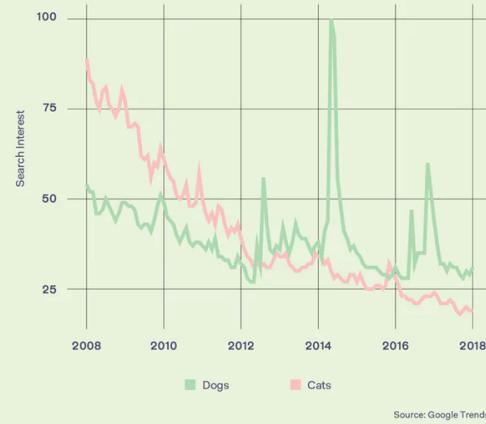
Google

Twitter?

Top 100 /r/aww Submissions About Cats and Dogs



Video Search Interest
Cats Versus Dogs



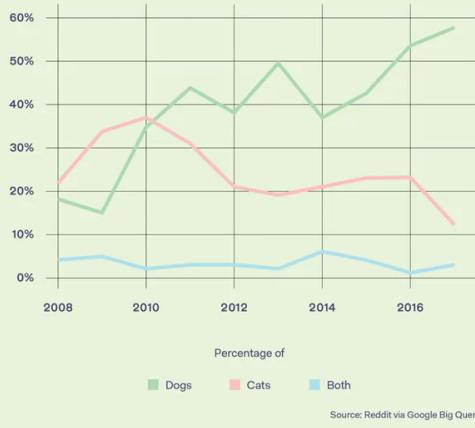
Write a program that sorts tweets into those containing “cat”, “dog”, or *other*

Traditional algorithms

Social media mentions of Cats vs. Dogs

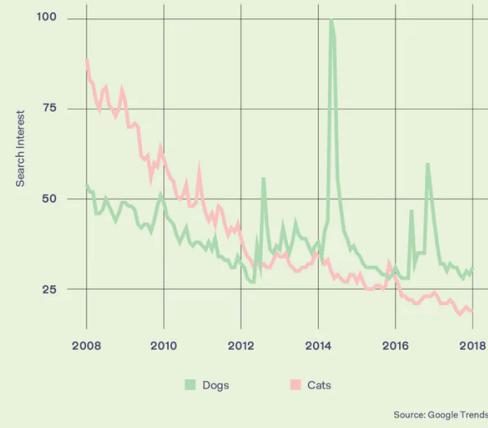
Reddit

Top 100 /r/aww Submissions About Cats and Dogs



Google

Video Search Interest
Cats Versus Dogs



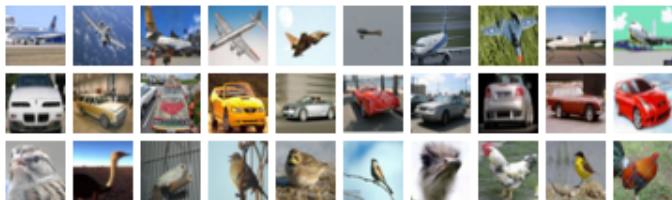
Twitter?

```
cats = []
dogs = []
other = []
for tweet in tweets:
    if "cat" in tweet:
        cats.append(tweet)
    elif "dog" in tweet:
        dogs.append(tweet)
    else:
        other.append(tweet)
return cats, dogs, other
```

Write a program that sorts tweets into those containing “cat”, “dog”, or other

Machine learning algorithms

Write a program that sorts images into those containing “birds”, “airplanes”, or *other*.



airplane
other
bird

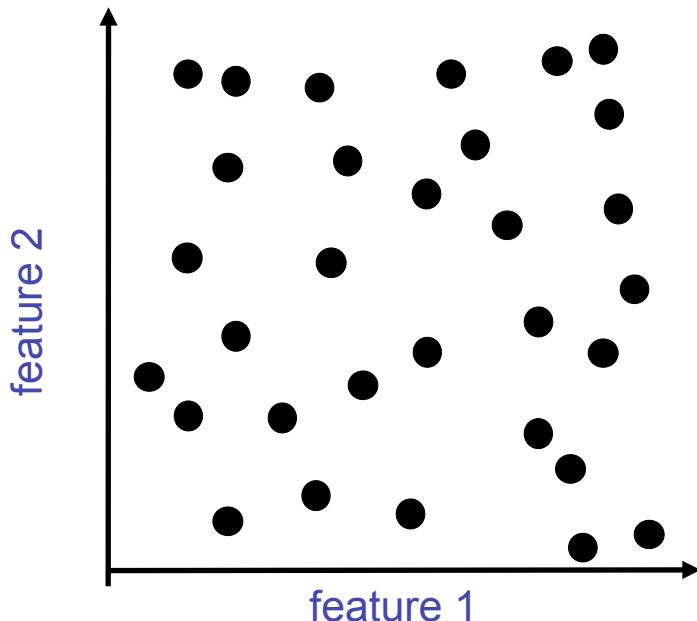
```
birds = []
planes = []
other = []

for image in images:
    if bird in image:
        birds.append(image)
    elif plane in image:
        planes.append(image)
    else:
        other.append(tweet)

return birds, planes, other
```

Machine learning algorithms

Write a program that sorts images into those containing “birds”, “airplanes”, or *other*.



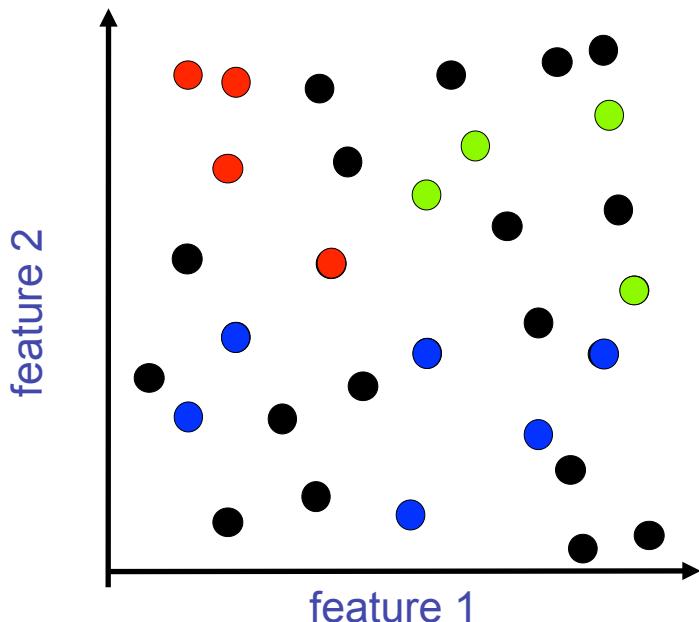
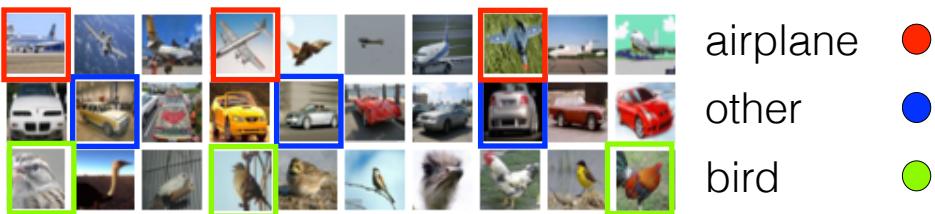
```
birds = []
planes = []
other = []

for image in images:
    if bird in image:
        birds.append(image)
    elif plane in image:
        planes.append(image)
    else:
        other.append(tweet)

return birds, planes, other
```

Machine learning algorithms

Write a program that sorts images into those containing “birds”, “airplanes”, or *other*.



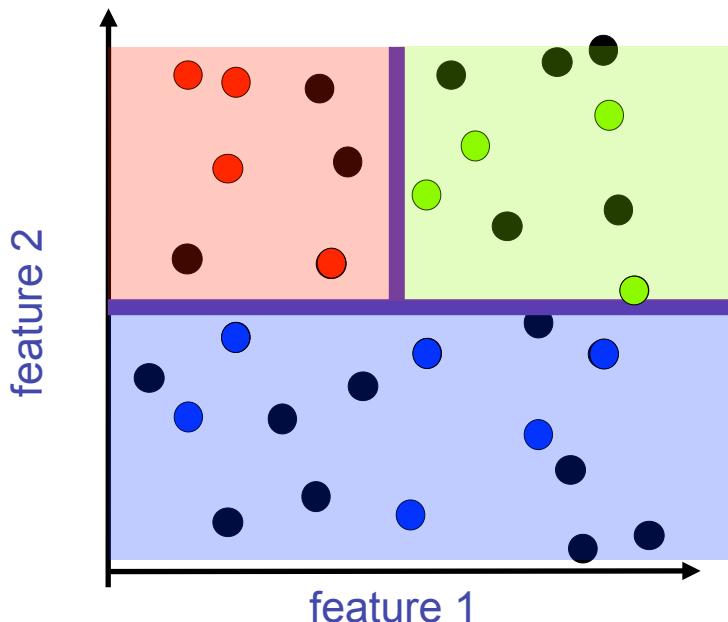
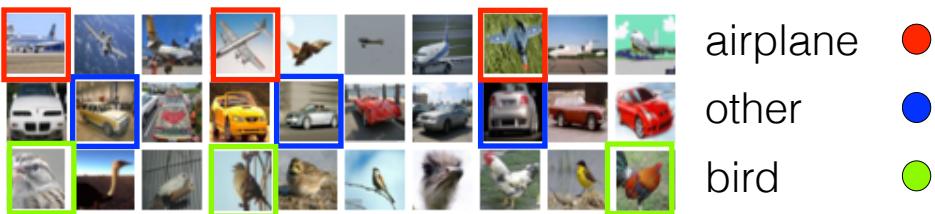
```
birds = []
planes = []
other = []

for image in images:
    if bird in image:
        birds.append(image)
    elif plane in image:
        planes.append(image)
    else:
        other.append(tweet)

return birds, planes, other
```

Machine learning algorithms

Write a program that sorts images into those containing “birds”, “airplanes”, or *other*.



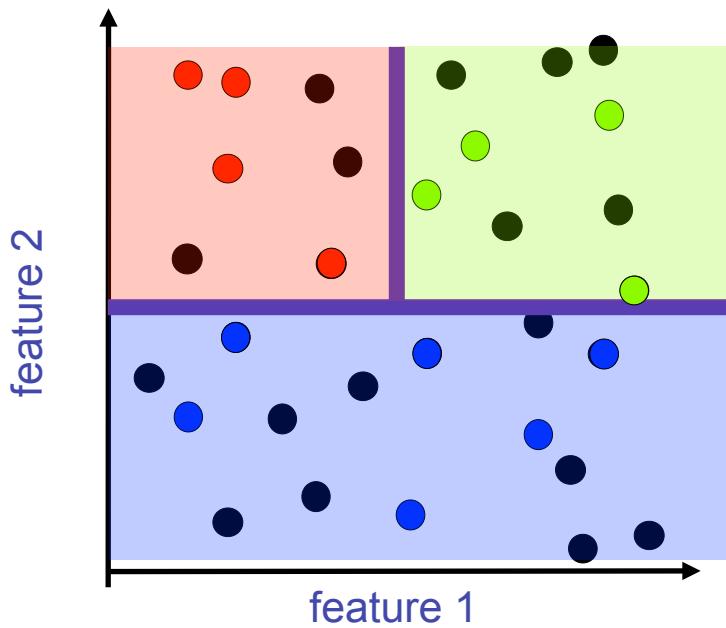
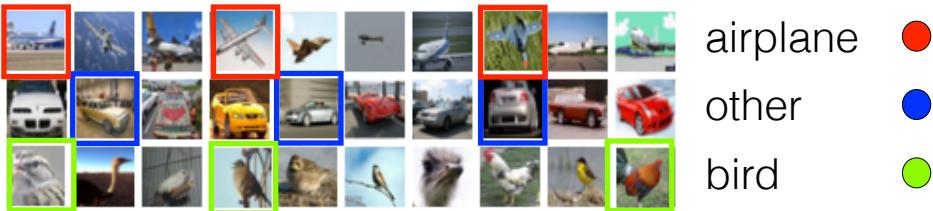
```
birds = []
planes = []
other = []

for image in images:
    if bird in image:
        birds.append(image)
    elif plane in image:
        planes.append(image)
    else:
        other.append(tweet)

return birds, planes, other
```

Machine learning algorithms

Write a program that sorts images into those containing “birds”, “airplanes”, or *other*.



```
birds = []
planes = []
other = []

for image in images:
    if bird in image:
        birds.append(image)
    elif plane in image:
        planes.append(image)
    else:
        other.append(tweet)

return birds, planes, other
```

The decision rule of
if “cat” in tweet:
is **hard coded by expert.**

The decision rule of
if bird in image:
is **LEARNED using DATA**

Machine Learning Ingredients

- **Data:** past observations
- **Hypotheses/Models:** devised to capture the patterns in data
- **Prediction:** apply model to forecast future observations

ML uses past data to make personalized predictions

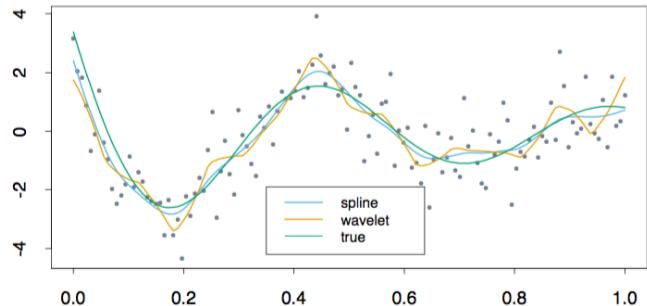


Machine learning is incredibly powerful and can have significant (unintended) negative consequences on society through targeting, excluding, and misusing.

Learning objectives of this course:

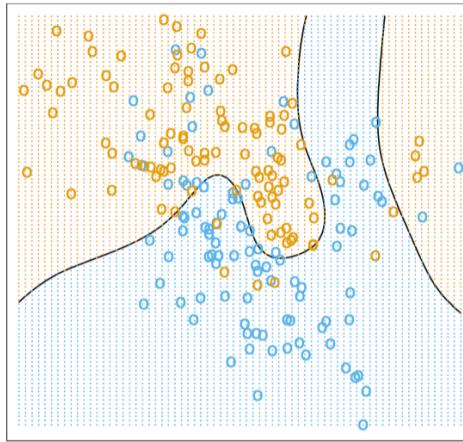
- introduction to the fundamental concepts of machine learning
- analysis and implementation of machine learning algorithms
- knowing how to use machine learning responsibly and robustly

Flavors of ML



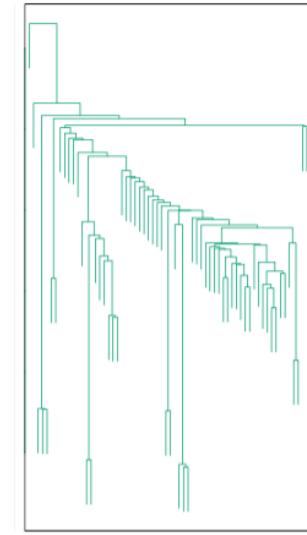
Regression

Predict continuous value:
ex: stock market, credit score,
temperature, Netflix rating



Classification

Predict categorical value:
loan or not? spam or not? what
disease is this?



Unsupervised Learning

Predict structure:
tree of life from DNA, find
similar images, community
detection

Mix of statistics (theory) and algorithms (programming)

CSE446/546: Machine Learning

Instructor: [Simon Du](#) and [Sewoong Oh](#)

Contact: cse446-staff@cs.washington.edu

Course Website: <https://courses.cs.washington.edu/courses/cse446/21sp/>

What this class is:

- **Fundamentals of ML:** bias/variance tradeoff, overfitting, optimization and computational tradeoffs, supervised learning (e.g., linear, boosting, deep learning), unsupervised models (e.g. k-means, EM, PCA)
- **Preparation for further learning:** the field is fast-moving, you will be able to apply the basics and teach yourself the latest

What this class is not:

- **Survey course:** laundry list of algorithms, how to win Kaggle
- **An easy course:** familiarity with intro linear algebra and probability are assumed, homework will be time-consuming

Prerequisites

- Formally:
 - MATH 308, CSE 312, STAT 390 or equivalent
- Familiarity with:
 - Linear algebra
 - linear dependence, rank, linear equations, SVD
 - Multivariate calculus
 - Probability and statistics
 - Distributions, marginalization, moments, conditional expectation
 - Algorithms
 - Basic data structures, complexity
- “Can I learn these topics concurrently?”
- Use HW0 to judge skills
- **See website for review materials!**

Grading

- 5 homework ($99\% = 10\% + 20\% + 20\% + 20\% + 29\%$)
 - *Each contains both theoretical questions and will have programming*
 - *Collaboration okay but must write who you collaborated with. You must write, submit, and understand your answers and code (which we may run)*
 - *Do not Google for answers.*
- NO exams
- 1% for submitting the proof of course evaluation
- We will assign random subgroups as PODs (when dust clears)

Homework

- HW 0 is out (**Due next Monday Apr 5th Midnight**)
 - Short *review*
 - Work individually, treat as barometer for readiness
- HW 1,2,3,4
 - They are not easy or short. Start early.
- Submit to Gradescope
- Regrade requests on Gradescope
- **There is no credit for late work, 5 late days**

Homework

- HW 0 is out (**Due next Monday Apr 5th Midnight**)
 - Short *review*
 - Work individually, treat as barometer for readiness
- HW 1,2,3,4
 - They are not easy or short. Start early.
- Submit to Gradescope
- Regrade requests on Gradescope
- **There is no credit for late work, 5 late days**

1. All code must be written in Python
2. All written work must be typeset (e.g., LaTeX)

See course website for tutorials and references.

Homework & Grading

- **CSE 446:**
 - Just do A problems
 - Doing B problems will not get higher grades
 - Grade is on 4.0 scale (relative to students in 446)
- **CSE 546:**
 - If just do A problems, grade is up to 3.8
 - B problems are for 0.2
 - Final grade = A grade + B grade (relative students in 546)

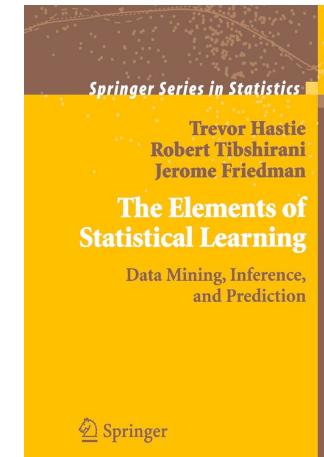
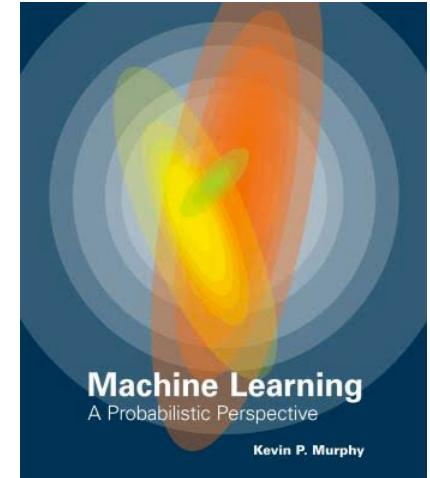
Communication Channels

- **Announcements, questions about class, homework help**
 - EdStem (invitation sent, contact TAs if you need access)
 - Weekly Section
 - Office hours (starts tomorrow)
- **Regrade requests**
 - Directly to Gradescope
- **Personal concerns**
 - Email: cse446-staff@cs.washington.edu
- **Anonymous feedback**
 - See website for link

Textbooks

- Required Textbook:
 - ***Machine Learning: a Probabilistic Perspective;***
Kevin Murphy

- Optional Books (free PDF):
 - ***The Elements of Statistical Learning: Data Mining, Inference, and Prediction;*** Trevor Hastie, Robert Tibshirani, Jerome Friedman



Addcodes

- Email: Elle Brown (ellean@cs.washington.edu) for addcodes

Enjoy!

- ML is becoming ubiquitous in science, engineering and beyond
- It's one of the hottest topics in industry today
- This class should give you the basic foundation for applying ML and developing new methods
- The fun begins...

Maximum Likelihood Estimation

W

Your first consulting job

- *Client:* I have special coin, if I flip it, what's the probability it will be heads?
 - *You:* I need to collect ***data***.
-
- *You:* The probability is:
 - *Client:* Why? What is the principle behind your prediction?

Modelling Coin Flips: Binomial Distribution

- **Data:** sequence $\mathcal{D} = (H, H, T, H, T, \dots)$
 - **k heads** out of **n flips**
- **Hypothesis:**
 - Flips are i.i.d. (independent and identically distributed):
 - Independent events
 - Identically distributed according to Bernoulli distribution
 - $P(\text{Heads}) = \theta, P(\text{Tails}) = 1 - \theta$
for some unknown **parameter** $\theta \in [0,1]$
- **Generative model:**
$$P(\mathcal{D}|\theta) =$$

Maximum Likelihood Estimation

- **Data:** sequence $\mathcal{D} = (H, H, T, H, T, \dots)$,

- **k heads out of n flips**

- **Hypothesis:** $P(\text{Heads}) = \theta$, $P(\text{Tails}) = 1 - \theta$

- **Likelihood:**

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

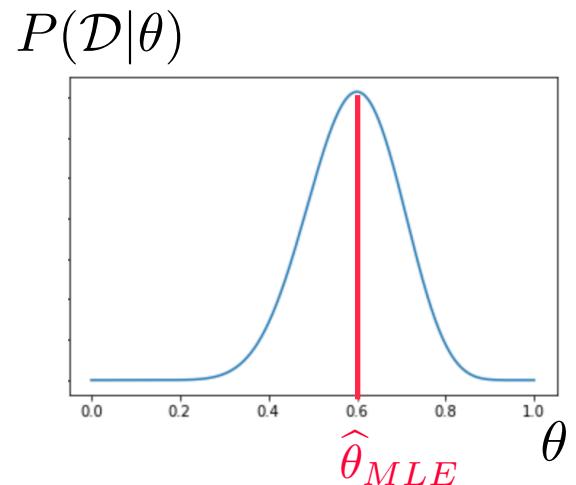
- **Maximum likelihood estimation (MLE):** Choose θ that maximizes the probability of observed data:

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(\mathcal{D}|\theta)$$

$$= \arg \max_{\theta} \log P(\mathcal{D}|\theta)$$

Your first learning algorithm

$$\begin{aligned}\hat{\theta}_{MLE} &= \arg \max_{\theta} \log P(\mathcal{D}|\theta) \\ &= \arg \max_{\theta} \log \theta^k (1-\theta)^{n-k}\end{aligned}$$



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

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

How good is MLE?

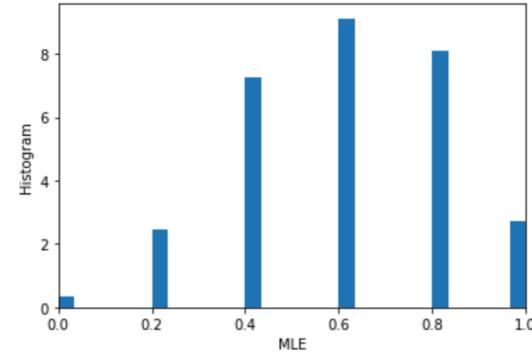
- We treat MLE $\hat{\theta}_{\text{MLE}}$ as a random variable, where there is a ground truth parameter θ^* that generates the data $\mathcal{D} = (HHTTH \dots)$ of a fixed size n
- What can we say about this random variable $\hat{\theta}_{\text{MLE}}$?
- First good property of MLE for Binomial: **unbiased**
 - Definition: **bias** of our MLE is
$$\text{Bias}(\hat{\theta}_{\text{MLE}}) := \mathbb{E}[\hat{\theta}_{\text{MLE}}] - \theta^* =$$
- **Expectation** describes how the estimator behaves *on average*

How many flips do I need?

$$\hat{\theta}_{MLE} = \frac{k}{n}$$

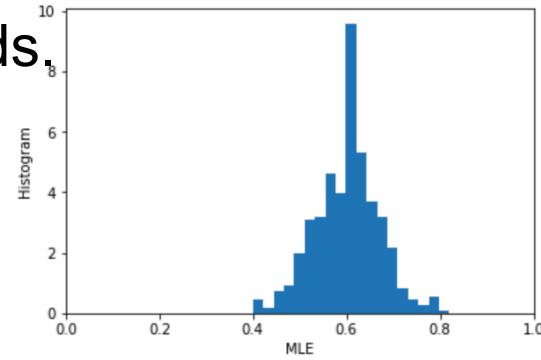
- *Client:* I flipped the coin 5 times and got 2 heads.

$$\hat{\theta}_{MLE} =$$



- *Client:* I flipped the coin 50 times and got 30 heads.

$$\hat{\theta}_{MLE} =$$



- *Client:* they are both unbiased, which one is right? Why?

Quantifying Uncertainty

- The **Variance** is the expected squared deviation from the mean:

$$\text{Variance}(\hat{\theta}_{MLE}) := \mathbb{E} \left[\left(\hat{\theta}_{MLE} - \mathbb{E}[\hat{\theta}_{MLE}] \right)^2 \right]$$

- As a rule of thumb

$$\hat{\theta}_{MLE} \simeq \mathbb{E}[\hat{\theta}_{MLE}] \pm \sqrt{\text{Variance}(\hat{\theta}_{MLE})}$$

- Second good property of MLE: **minimum (asymptotic) variance**
- **Exercise:** compute the $\text{Variance}(\hat{\theta}_{MLE})$

Expectation versus High Probability

- Tail bound of a random variable
- For any $\epsilon > 0$ can we bound $\mathbb{P}(|\hat{\theta}_{MLE} - \mathbb{E}[\hat{\theta}_{MLE}]| \geq \epsilon)$?

Markov's inequality

For any $t > 0$ and non-negative random variable X

$$\mathbb{P}(X \geq t) \leq \frac{\mathbb{E}[X]}{t}$$

- **Exercise:** Apply Markov's inequality to obtain bound.
(Hint: set $X = |\hat{\theta}_{MLE} - \theta^*|^2$)

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

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):

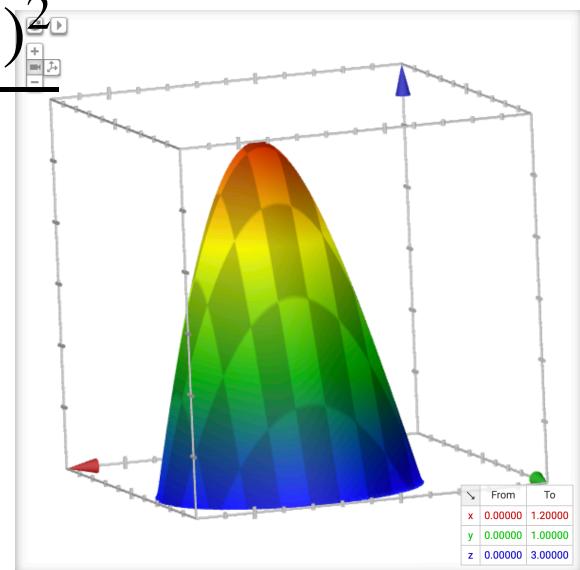
$$P(\mathcal{D}; \mu, \sigma) = P(x_1, \dots, x_n; \mu, \sigma)$$

$$= \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

- Log-likelihood of data:

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

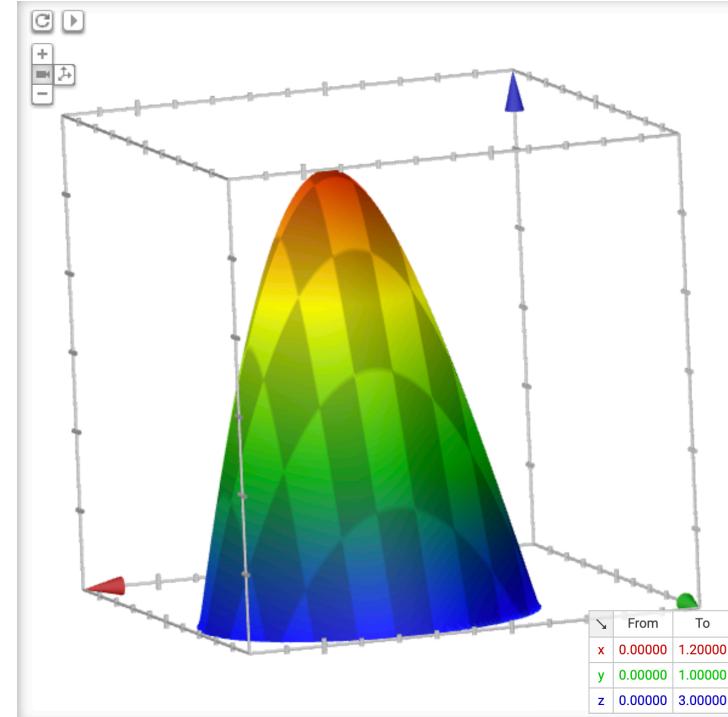
- What is $\hat{\theta}_{MLE}$ for $\theta = (\mu, \sigma^2)$?



Your second learning algorithm: MLE for mean of a Gaussian

- What's MLE for mean?

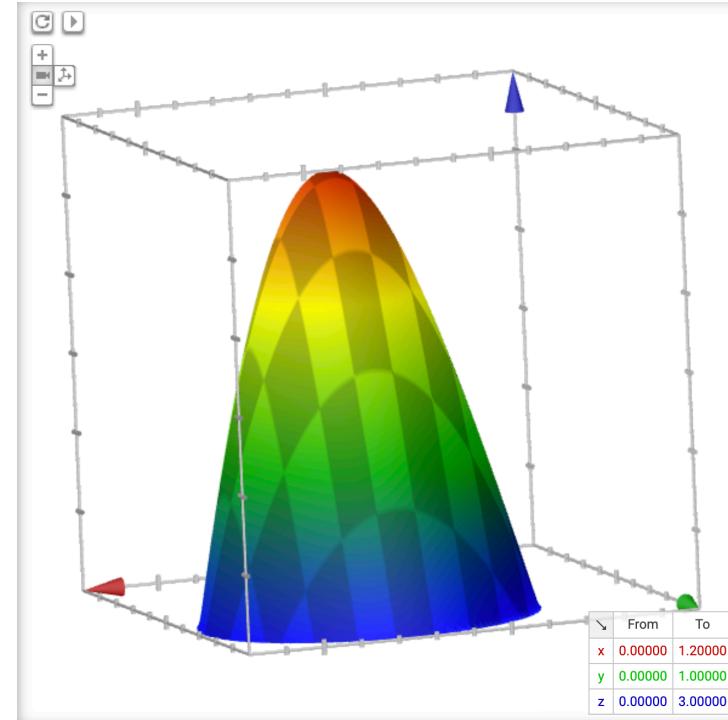
$$\frac{d}{d\mu} \log P(\mathcal{D}; \mu, \sigma) = \frac{d}{d\mu} \left[-n \log(\sigma\sqrt{2\pi}) - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$



MLE for variance

- Again, set derivative to zero:

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



What can we say about the MLE?

- MLE:

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

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

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

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

- Unbiased variance estimator:

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

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

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

- Asymptotically consistent and normal: $\frac{\hat{\theta}_{MLE} - \theta_*}{\widehat{se}} \sim \mathcal{N}(0, 1)$
- 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

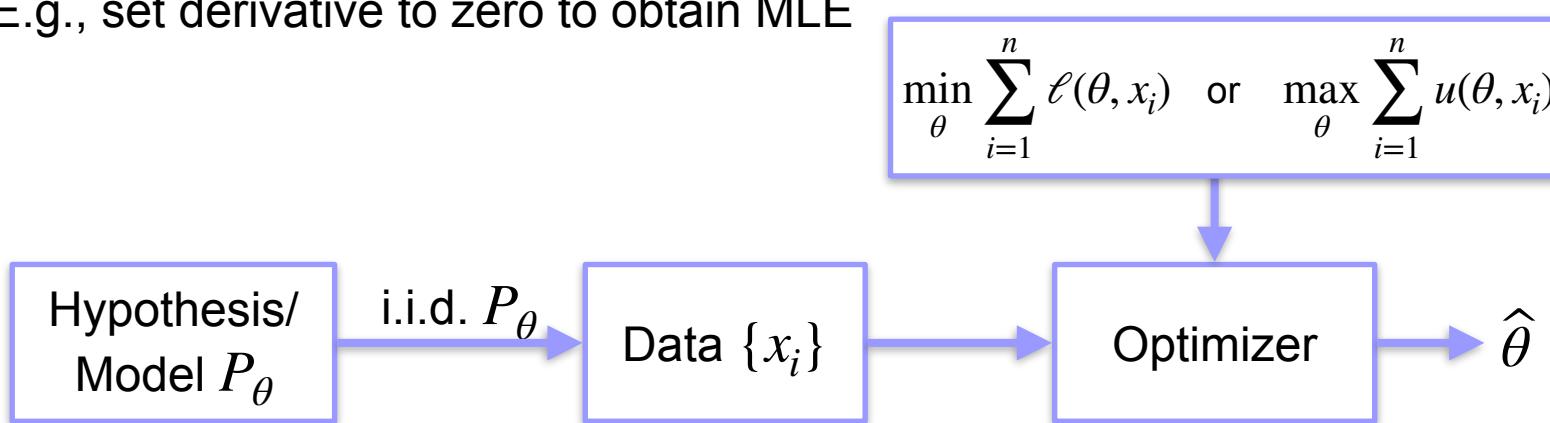
$$\min_{\theta} \sum_{i=1}^n \ell(\theta, x_i) \quad \text{or} \quad \max_{\theta} \sum_{i=1}^n u(\theta, x_i)$$



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

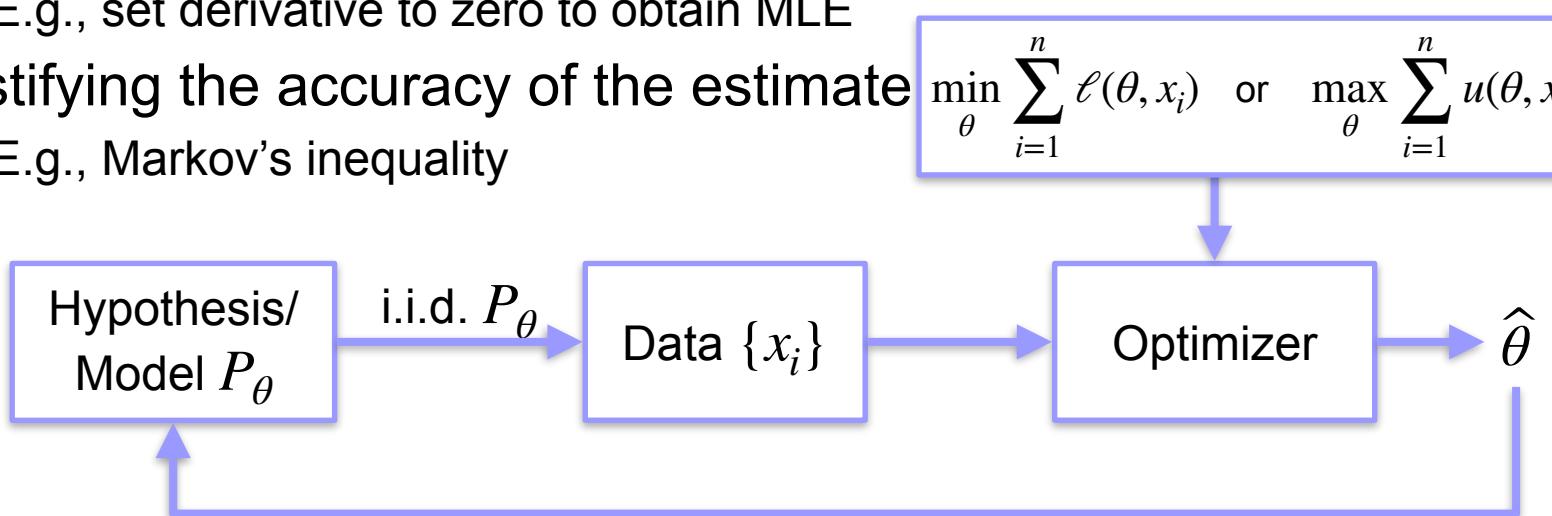
$$\min_{\theta} \sum_{i=1}^n \ell(\theta, x_i) \quad \text{or} \quad \max_{\theta} \sum_{i=1}^n u(\theta, x_i)$$



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

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



Linear Regression

UNIVERSITY *of* WASHINGTON

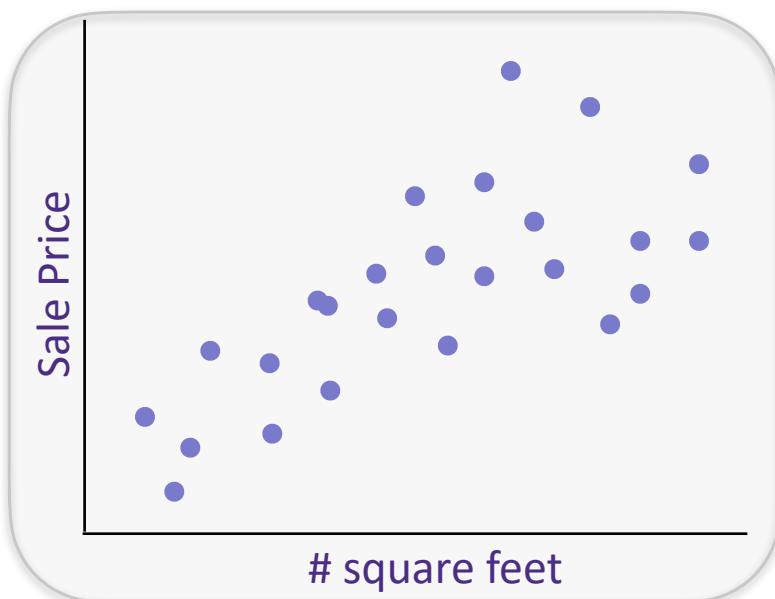
W

The regression problem, 1-dimensional

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

y = House sale price from

x = {# sq. ft.}



Training Data:

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

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

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

Process

Decide on a **model**

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

Find the function which fits the data best

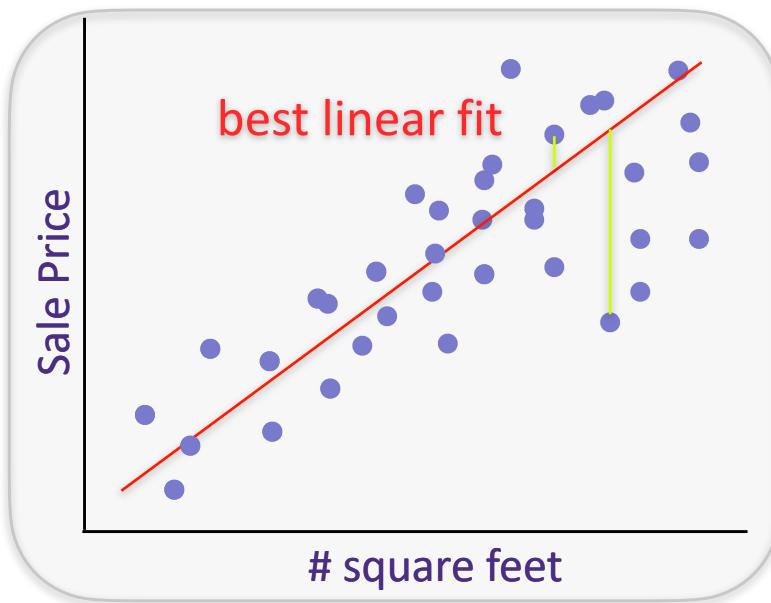
Use function to make prediction on new examples

Fit a function to our data, 1-dimension

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

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

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



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

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

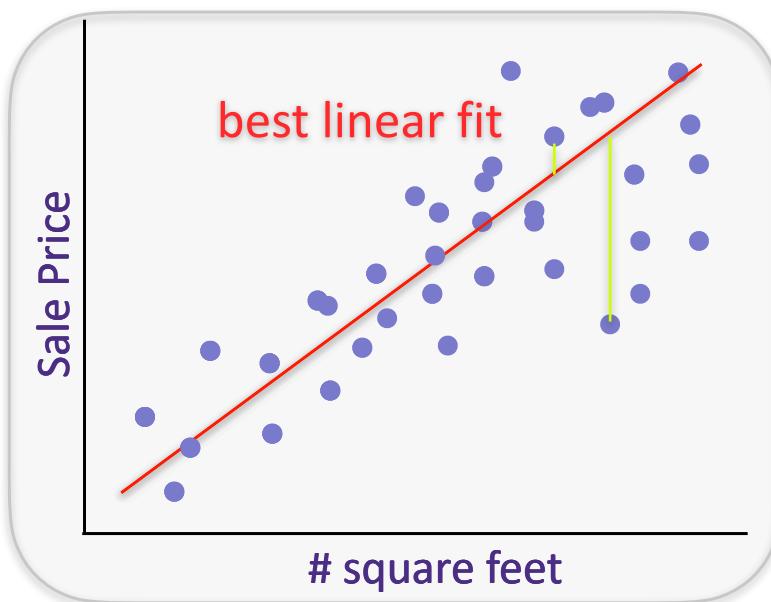
The regression problem, d-dimensions

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

Hypothesis/Model: linear

$$y_i \approx x_i^T w$$

Loss: least squares solution

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

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

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:

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

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

•

•

•

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

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

The regression problem in matrix notation

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

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

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

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

•

•

•

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

Loss: $\hat{w}_{LS} = \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 = \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2$

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

The regression problem in matrix notation

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

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

The regression problem in matrix notation

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

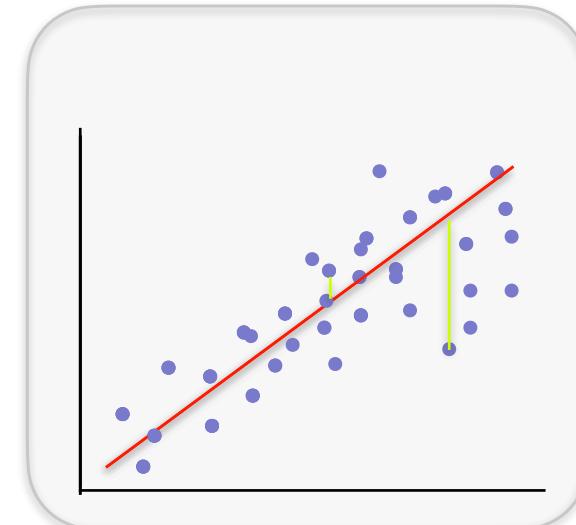
“Closed form” solution!

The regression problem in matrix notation

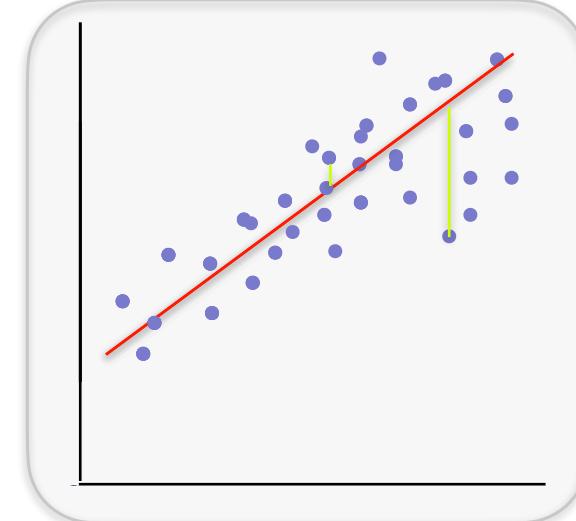
Linear model: $y_i = x_i^T w + \epsilon_i$

Least squares solution:

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



What about an offset
(a.k.a intercept)?

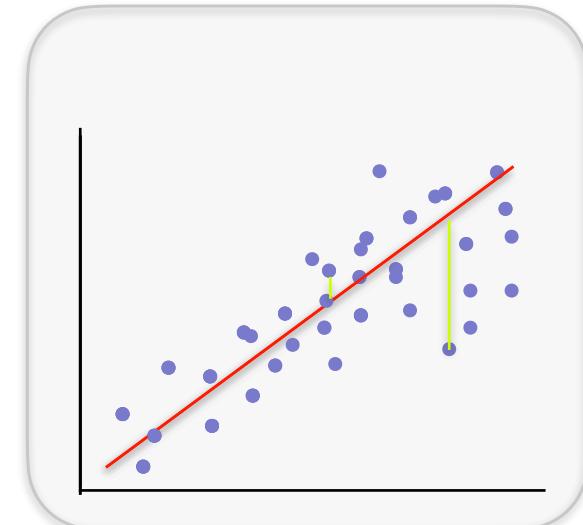


The regression problem in matrix notation

Linear model: $y_i = x_i^T w + \epsilon_i$

Least squares solution:

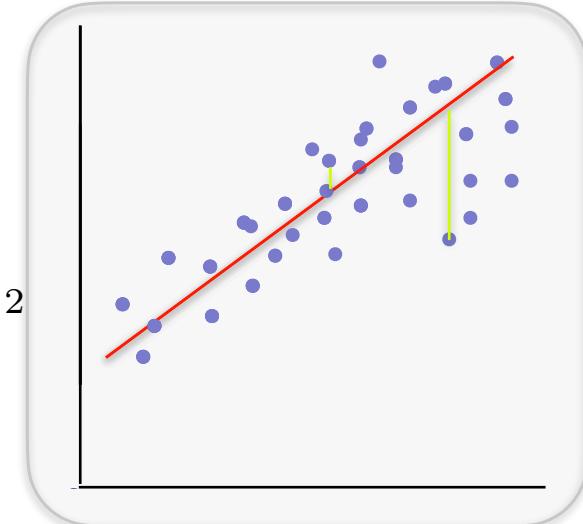
$$\begin{aligned}\hat{w}_{LS} &= \arg \min_w \|\mathbf{y} - \mathbf{X}w\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$



Affine model: $y_i = x_i^T w + b + \epsilon_i$

Least squares solution:

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



Dealing with an offset

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

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

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 $\mathbf{X}^T \mathbf{1} = 0$, if the features have zero mean,

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

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

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 $\mathbf{X}^T \mathbf{1} = 0$,

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

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

In general, when $\mathbf{X}^T \mathbf{1} \neq 0$,

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 $\mathbf{X}^T \mathbf{1} = 0$,

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

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

In general, when $\mathbf{X}^T \mathbf{1} \neq 0$,

$$\mu = \frac{1}{n} \mathbf{X}^T \mathbf{1}$$

$$\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{1}\mu^T$$

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

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

Process

Decide on a **model**: $y_i = x_i^T w + b + \epsilon_i$

Choose a loss function - least squares

Pick the function which minimizes loss on data

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w,b} \sum_{i=1}^n (y_i - (x_i^T w + b))^2$$

Use function to make prediction on new examples

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

Another way of dealing with an offset

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

reparametrize the problem as $\bar{\mathbf{X}} = [\mathbf{X}, \mathbf{1}]$ and $\bar{w} = \begin{bmatrix} w \\ b \end{bmatrix}$

$$\bar{\mathbf{X}} \bar{w} =$$

Why is least squares a good loss function?

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

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

$\implies y_i \sim$

$\implies P(y_i; x_i, w, \sigma) =$

Why is least squares a good loss function?

Maximum Likelihood Estimator:

$$\begin{aligned}\hat{w}_{\text{MLE}} &= \arg \max_w \log P(\{y_i\}_{i=1}^n; \{x_i\}_{i=1}^n, w, \sigma) \\ &= \arg \max_w -n \log(\sigma \sqrt{2\pi}) + \sum_{i=1}^n -\frac{(y_i - x_i^T w)^2}{2\sigma^2}\end{aligned}$$

Why is least squares a good loss function?

Maximum Likelihood Estimator:

$$\begin{aligned}\hat{w}_{MLE} &= \arg \max_w \log P(\{y_i\}_{i=1}^n; \{x_i\}_{i=1}^n, w, \sigma) \\ &= \arg \max_w -n \log(\sigma \sqrt{2\pi}) + \sum_{i=1}^n -\frac{(y_i - x_i^T w)^2}{2\sigma^2} \\ &= \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2 \\ \text{Recall: } \hat{w}_{LS} &= \arg \min_w \sum_{i=1}^n (y_i - x_i^T w)^2\end{aligned}$$

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

Recap of linear regression

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

**Minimize the loss
(Empirical Risk Minimization)**

Choose a loss
e.g., $(y_i - x_i^T w)^2$

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

**Maximize the likelihood
(MLE)**

Choose a Hypothesis class
e.g., $y_i = x_i^T w + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$

Maximize the likelihood,
 $\hat{w}_{\text{MLE}} = \arg \max_w \left\{ -n \log(\sigma \sqrt{2\pi}) - \sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{2\sigma^2} \right\}$

Analysis of Error under additive Gaussian noise

if $y_i = x_i^T w + \epsilon_i$ and $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$ $\mathbf{Y} = \mathbf{X}w + \epsilon$

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w + \epsilon) \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

Maximum Likelihood Estimator is unbiased:

Analysis of Error under additive Gaussian noise

if $y_i = x_i^T w + \epsilon_i$ and $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$ $\mathbf{Y} = \mathbf{X}w + \epsilon$

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w + \epsilon) \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

Covariance is:

Analysis of Error under additive Gaussian noise

if $y_i = x_i^T w + \epsilon_i$ and $\epsilon_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma^2)$ $\mathbf{Y} = \mathbf{X}w + \epsilon$

$$\begin{aligned}\hat{w}_{MLE} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}w + \epsilon) \\ &= w + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon\end{aligned}$$

$$\mathbb{E}[\hat{w}_{MLE}] = w$$

$$\text{Cov}(\hat{w}_{MLE}) = \mathbb{E}[(\hat{w} - \mathbb{E}[\hat{w}])(\hat{w} - \mathbb{E}[\hat{w}])^T] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$$

$$\hat{w}_{MLE} \sim \mathcal{N}(w, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

Questions?
