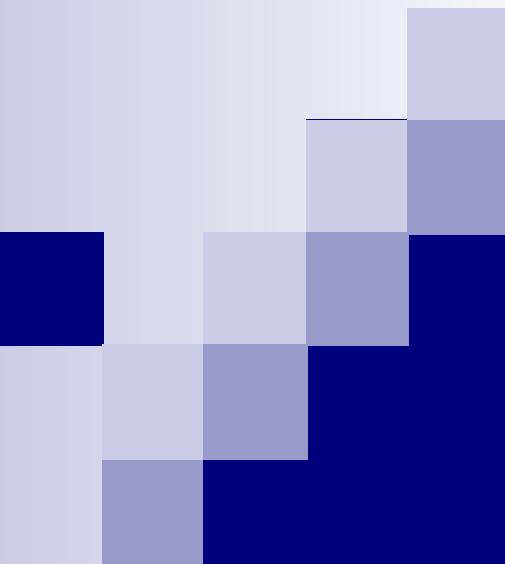


<https://courses.cs.washington.edu/446>



Machine Learning CSE446

Kevin Jamieson and Anna Karlin
University of Washington

April 1, 2019

Traditional algorithms

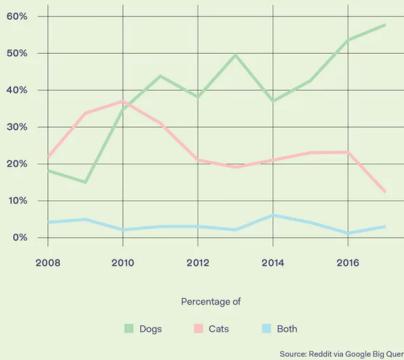
Social media mentions of Cats vs. Dogs

Reddit

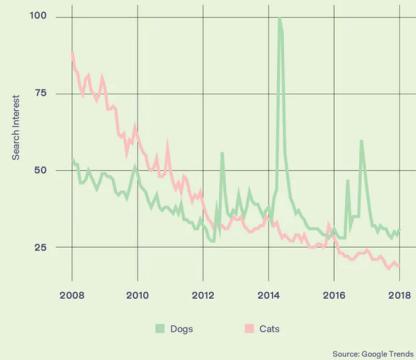
Google

Twitter?

Top 100 /r/aww Submissions About Cats and Dogs



Video Search Interest
Cats Versus Dogs



Traditional algorithms

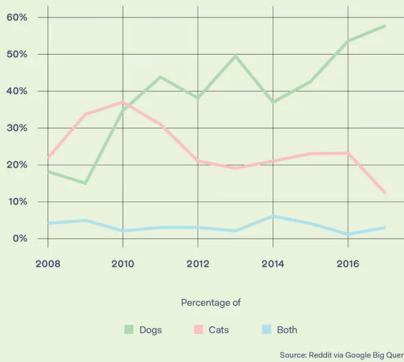
Social media mentions of Cats vs. Dogs

Reddit

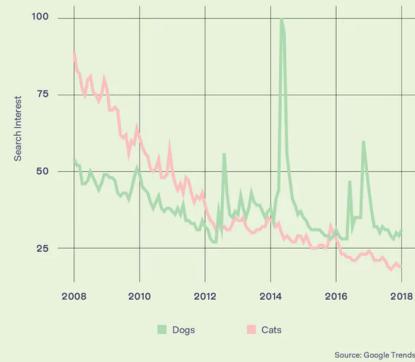
Google

Twitter?

Top 100 /r/aww Submissions About Cats and Dogs



Video Search Interest
Cats Versus Dogs



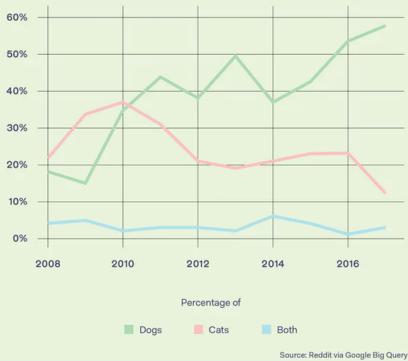
Write a program that sorts tweets into those containing “cat”, “dog”, or other

Traditional algorithms

Social media mentions of Cats vs. Dogs

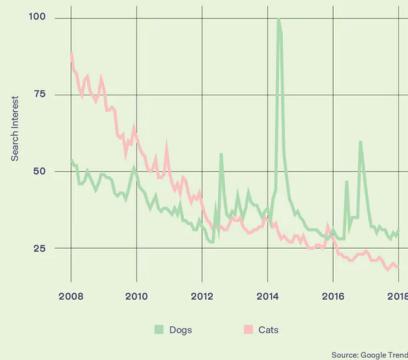
Reddit

Top 100 /r/aww Submissions About Cats and Dogs



Google

Video Search Interest
Cats Versus Dogs



Twitter?

```
cats = []
dogs = []
other = []

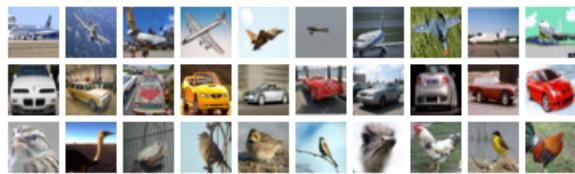
for tweet in tweets:
    if "cat" in tweet:
        cats.append(tweet)
    elif "dog" in tweet:
        dogs.append(tweet)
    else:
        other.append(tweet)

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*.**



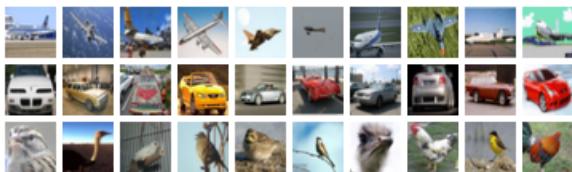
airplane

other

bird

Machine learning algorithms

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

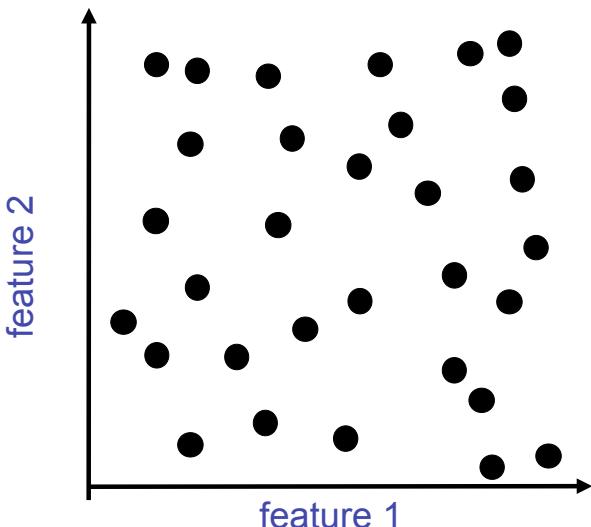
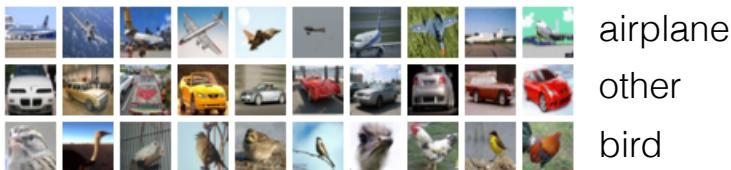


airplane
other
bird

```
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

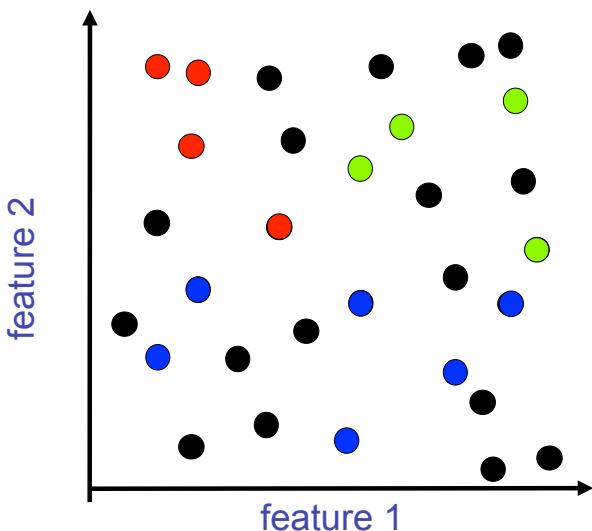
Write a program that sorts images into those containing “birds”, “airplanes”, or *other*.



```
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

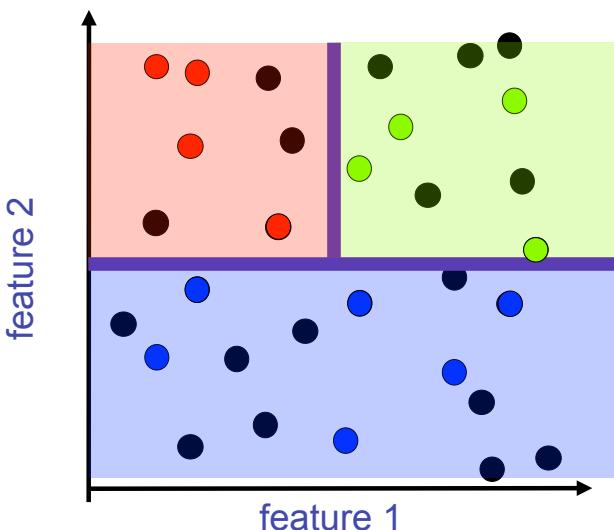
Write a program that sorts images into those containing “birds”, “airplanes”, or *other*.



```
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

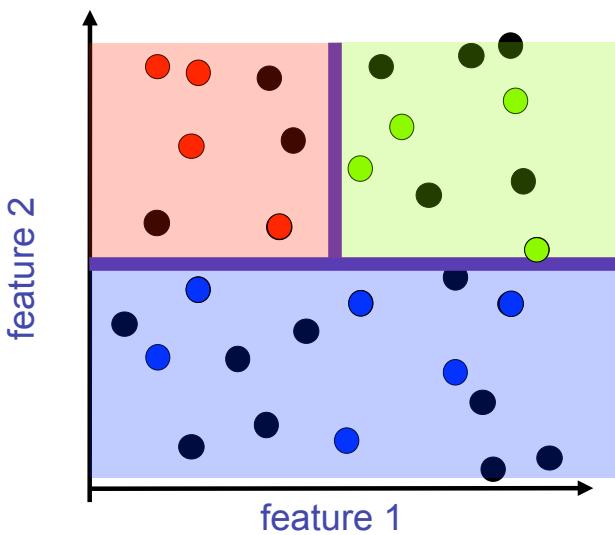
Write a program that sorts images into those containing “birds”, “airplanes”, or *other*.



```
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*.



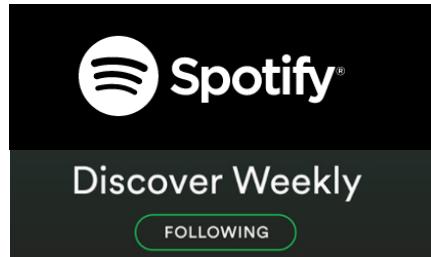
```
birds = []
planes = []
other = []
for image in images:
    if bird in image:
        birds.append(image)
    elif plane in image:
        planes.append(image)
    else:
        other.append(tweet)
return birds, planes, other
```

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

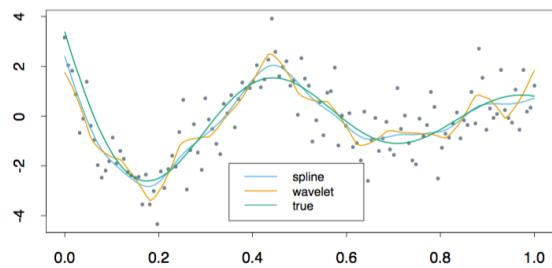
- **Data:** past observations
- **Hypotheses/Models:** devised to capture the patterns in data
- **Prediction:** apply model to forecast future observations



ML uses past data to make personalized predictions

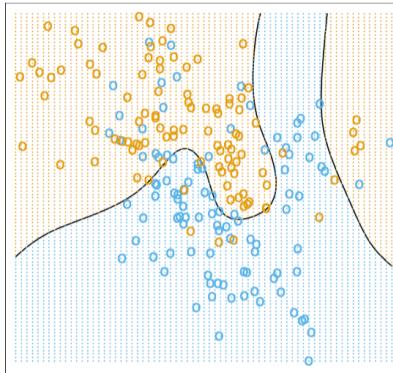


Flavors of ML



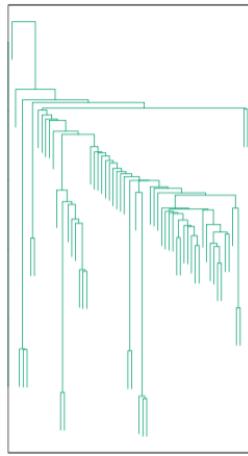
Regression

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: Machine Learning

Lecture: Monday, Wednesday 3:30-4:50 Room: [GUG 220](#)

Instructor: [Kevin Jamieson](#) and [Anna Karlin](#)

Contact: cse446-staff@cs.washington.edu

Website: <https://courses.cs.washington.edu/courses/cse446/19sp/>

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

Grading

- 5 homeworks (50%)
 - Each contains both theoretical questions and will have programming
 - Collaboration okay. You must write, submit, and understand your answers and code (which we may run)
 - Do not Google for answers.
- Midterm **May 8** (20%) and Final **June 13** (30%)

Homeworks

- HW 0 is out (**Due next Monday Midnight**, 37 points)
 - Short *review*
 - Work individually, treat as barometer for readiness
- HW 1,2,3,4 (100 points each)
 - They are not easy or short. Start early.
- Submit to Gradescope
- Regrade requests on Gradescope
- **There is no credit for late work, receives 0 points.**

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.

Communication Channels

- **Announcements, questions about class, homework help**
 - Piazza (get code on Canvas)
 - Section
 - Office hours (start tomorrow)
- **Regrade requests**
 - Directly to Gradescope
- **Personal concerns**
 - Email: cse446-staff@cs.washington.edu
- **Anonymous feedback**
 - See website for link

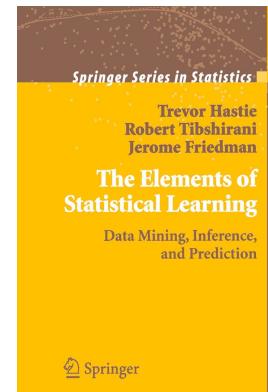
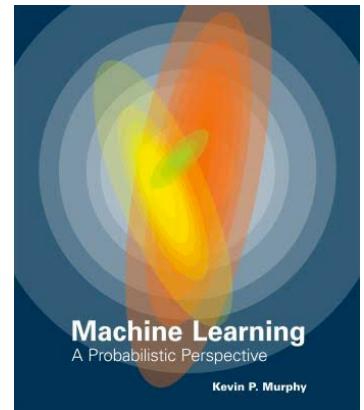
Staff

- 13 Experienced TAs, lots of office hours (see website)
- Satvik Agarwal, Kousuke Ariga, Eric Chan, Benjamin Evans, Shobhit Hathi, Zeyu Liu, Andrew Luo, Vardhman Mehta, Cheng Ni, Deric Pang, Robbie Weber, Kyle Zhang and Michael Zhang

Text Books

- 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



Add code requests

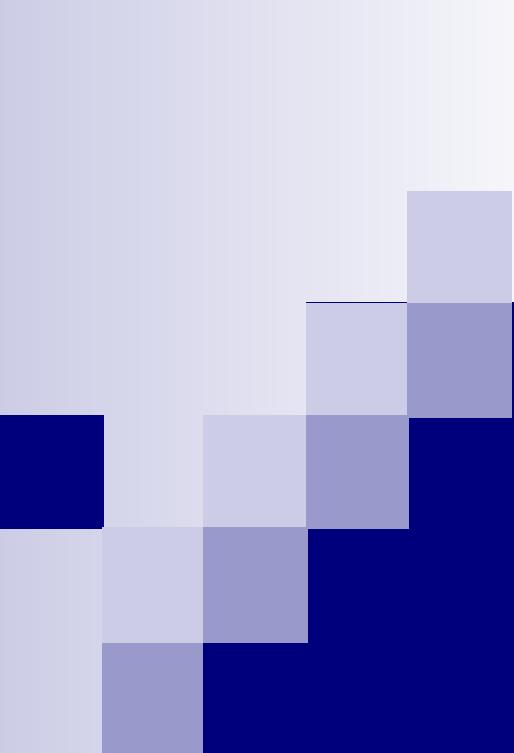
Add codes are given out according to a centralized process organized by CSE. Do **not** email instructors, we cannot help.

Resources:

- <https://www.cs.washington.edu/academics/ugrad/advising/>
- <https://www.cs.washington.edu/academics/ugrad/courses/petition>

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

Machine Learning – CSE446
Kevin Jamieson
University of Washington

April 1, 2019

Your first consulting job

- *Billionaire*: I have a special coin, if I flip it, what's the probability it will be heads?
- *You*: Please flip it a few times:

H H T H T

- *You*: The probability is: 3/5
- *Billionaire*: Why?

Coin – Binomial Distribution

- **Data:** sequence $D = (HHTHT\ldots)$, k heads out of n flips
- **Hypothesis:** $P(\text{Heads}) = \theta$, $P(\text{Tails}) = 1-\theta$
 - Flips are i.i.d.: $P(X_2 = H | X_1 = T) = P(X_2 = H)$
 - Independent events
 - Identically distributed according to Binomial distribution
- $$HHTHT$$
- $$\begin{aligned}P(D|\theta) &= \theta \cdot \theta \cdot (1-\theta) \cdot \theta \cdot (1-\theta) \cdots \\&= \theta^k (1-\theta)^{n-k}\end{aligned}$$

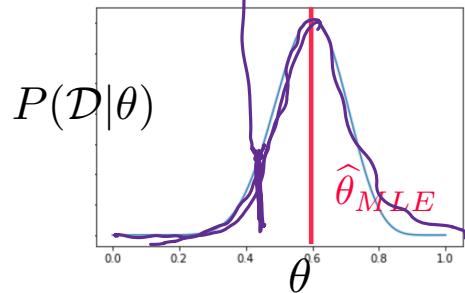
Maximum Likelihood Estimation

- **Data:** sequence $D = (HHTHT\ldots)$, **k heads** out of **n flips**
- **Hypothesis:** $P(\text{Heads}) = \theta$, $P(\text{Tails}) = 1 - \theta$

$$P(D|\theta) = \underbrace{\theta^k}_{\text{purple}} (1 - \theta)^{n-k} \underbrace{(1-\theta)^{n-k}}_{\text{purple}}$$

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

$$\begin{aligned}\widehat{\theta}_{MLE} &= \arg \max_{\theta} P(D|\theta) \\ &= \arg \max_{\theta} \log P(D|\theta)\end{aligned}$$



Your first learning algorithm

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

$$\log(ab) = \log(a) + \log(b) = \arg \max_{\theta} \log \theta^k (1-\theta)^{n-k}$$

- Set derivative to zero:

$$\log(a^b) = b \log(a)$$

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

$$\begin{aligned} &= \log(\theta^k) + \log((1-\theta)^{n-k}) \\ &= \underline{k \log(\theta)} + \underline{(n-k) \log(1-\theta)} \end{aligned}$$

$$\frac{\partial}{\partial \theta} \boxed{\quad} = \frac{k}{\theta} + \frac{n-k}{1-\theta} (-1) = 0$$

$$\underline{(1-\theta)k} + \theta \underline{(n-k)(-1)} = 0$$
$$k - \theta(n) = 0 \quad \theta = \frac{k}{n}$$

How many flips do I need?

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

- You: flip the coin 5 times. *Billionaire*: I got 3 heads.

$$\hat{\theta}_{MLE} = \frac{3}{5}$$

- You: flip the coin 50 times. *Billionaire*: I got 20 heads.

$$\hat{\theta}_{MLE} = \frac{20\%}{50} = \frac{2}{5}$$

- *Billionaire*: Which one is right? Why?

Simple bound

An estimator $\hat{\theta}: \mathcal{B} \xrightarrow{\uparrow} \mathbb{R}$ $\mathbb{E}[\hat{\theta}] = \theta^*$

(based on Hoeffding's inequality)

- For **n flips** and **k heads** the MLE is **unbiased** for true θ^* :

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

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

$k = \sum_{i=1}^n \mathbb{1}\{\text{ith coin came up heads}\}$

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

- Hoeffding's inequality says that for any $\epsilon > 0$:

$$P(|\hat{\theta}_{MLE} - \theta^*| \geq \epsilon) \leq \underline{2e^{-2n\epsilon^2}}$$

$$\epsilon = 0.1 \quad \leq 2e^{-0.02 \cdot n}$$

PAC Learning

- PAC: Probably Approximate Correct
- *Billionaire*: I want to know the parameter θ^* , within $\epsilon = 0.1$, with probability at least $1-\delta = 0.95$. How many flips?

$$P(|\hat{\theta}_{MLE} - \theta^*| \geq \epsilon) \leq 2e^{-2n\epsilon^2}$$

What about continuous variables?

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

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

Some properties of Gaussians

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

MLE for Gaussian

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

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

- Log-likelihood of data:

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

Your second learning algorithm: MLE for mean of a Gaussian

- What's MLE for mean?

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

$$= \sum \frac{(x_i - \mu)}{\sigma^2}$$

$$= 0$$

$$\Rightarrow \sum_{i=1}^n (x_i - \mu) = 0 \quad \Rightarrow \hat{\mu}_{\text{MLE}} = \frac{1}{n} \sum_{i=1}^n x_i$$

MLE for variance

- Again, set derivative to zero:

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

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

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

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

Learning Gaussian parameters

- MLE:

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

$x_i \sim N(\mu, \sigma^2)$

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

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

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

- 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” $\theta = \underline{\theta_*}$

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

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
 - 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., Hoeffding's inequality