

Machine Learning for Networking

- ① Programming in Python / Machine Learning in Python.
- ② The Internet / Computer Networking
- ③ Machine Learning Models, Algorithms

④ Putting it together

What is Machine Learning?

- Some algorithm / program that gets letter at performing a task as it gains more data.
- Supervised : Learning to predict outcomes. (^{data has labels})
 $y = f(x)$
targets/outcomes features

- Unsupervised: Learning structure from data (^{data has no labels})
 - Clusters
 - Reducing dimensions

Linear Algebra

$$\beta = [\beta_1, \beta_2, \dots, \beta_n]$$

$$x = [x_1, x_2, \dots, x_n]$$

$$\beta \cdot x = \beta^T x = x \beta^T$$

Dot Product

$$[x_1 \dots x_n] \cdot \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{bmatrix} \dots = \beta_1 x_1 + \beta_2 x_2 + \dots$$

$$y = mx + b$$

$$\approx \beta_1 x_1$$

$$\approx \beta_1 x_1 + \beta_2 x_2 + \dots$$

$$\begin{array}{c} 3 \times 2 \\ \left[\begin{array}{cc} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \end{array} \right] \end{array} \xrightarrow{\text{Transpose}} \begin{array}{c} 2 \times 3 \\ \left[\begin{array}{ccc} x_{11} & x_{21} & x_{31} \\ x_{12} & x_{22} & x_{32} \end{array} \right] \end{array}$$

Matrix Multiplication

$$A \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & \ddots & \ddots & \vdots \\ a_{m1} & & & a_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n \end{bmatrix}$$

$m \times n$ $n \times 1$ $m \times 1$

$$\begin{matrix} AB \\ \uparrow \quad \uparrow \\ m \times n \quad n \times p \end{matrix} \rightarrow M \quad m \times p$$

How would you use Machine learning in computer networking?

① Security

- Detect attacks (e.g., denial of service, spam)
- Detect unusual behavior ("anomalies")
- Detect or classify devices, applications, OSes, ...

② Performance

- Infer application performance (e.g., streaming video)

① Basic Machine Learning Model:

- Supervised vs unsupervised
- Basic Model: Linear Regression.

② The Internet : Data about the Internet

- How does data get from A to B?
 - Packets ("chunks")  "envelopes"
 - Bytes (0s & 1s) → 8 bits.

WIRESHARK { - What is in a packet? ("header")

- IP addresses, ports, protocols.

- To come: Names → IP addresses ("DNS")
google.com → 172.217.9.46

AGENDA

① Internet Background / Basics.

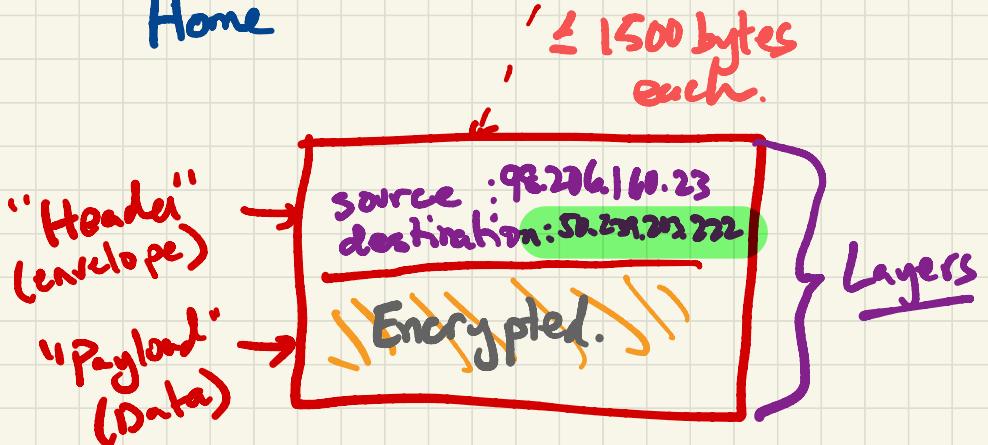
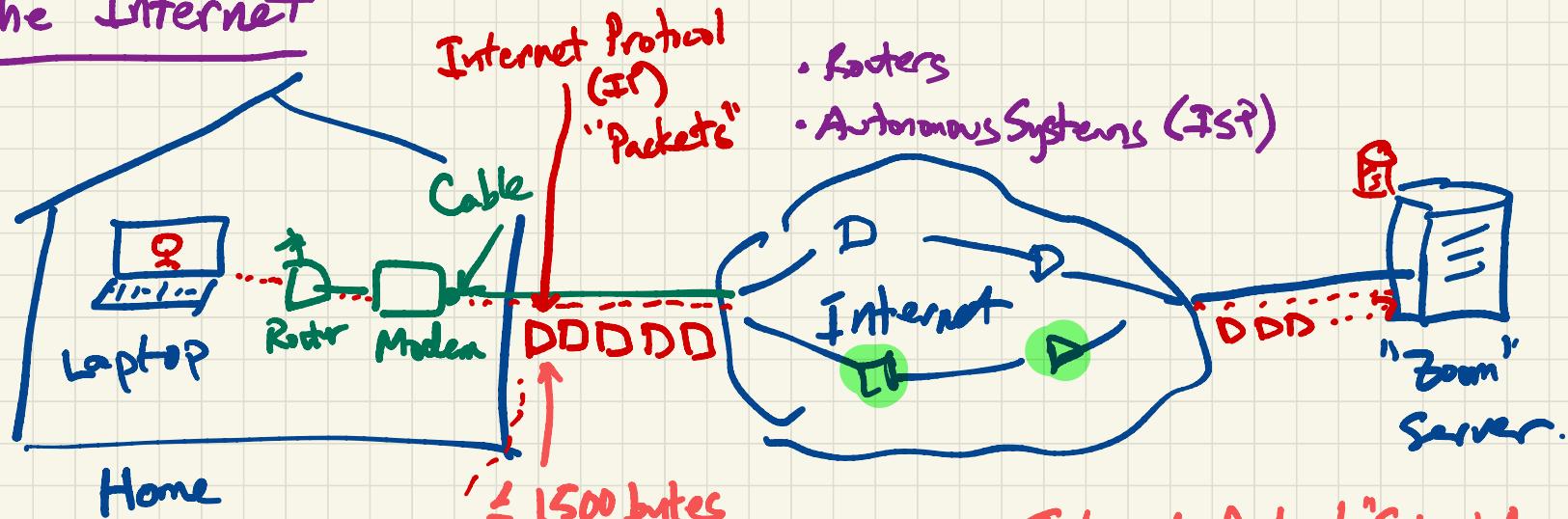
- Internet Protocol.
- Packets
- Domain Name System (DNS)

② Live Demonstration

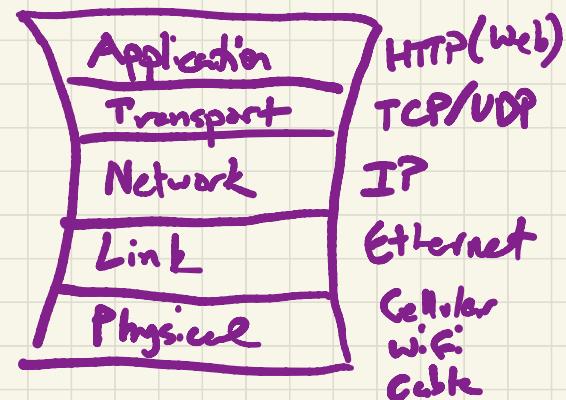
- Network Measurement / Packet Capture (Wireshark)
- Web page download.

③ Analyzing Network Traffic in Python.

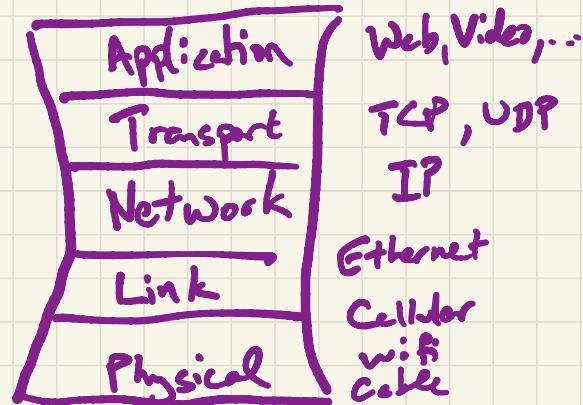
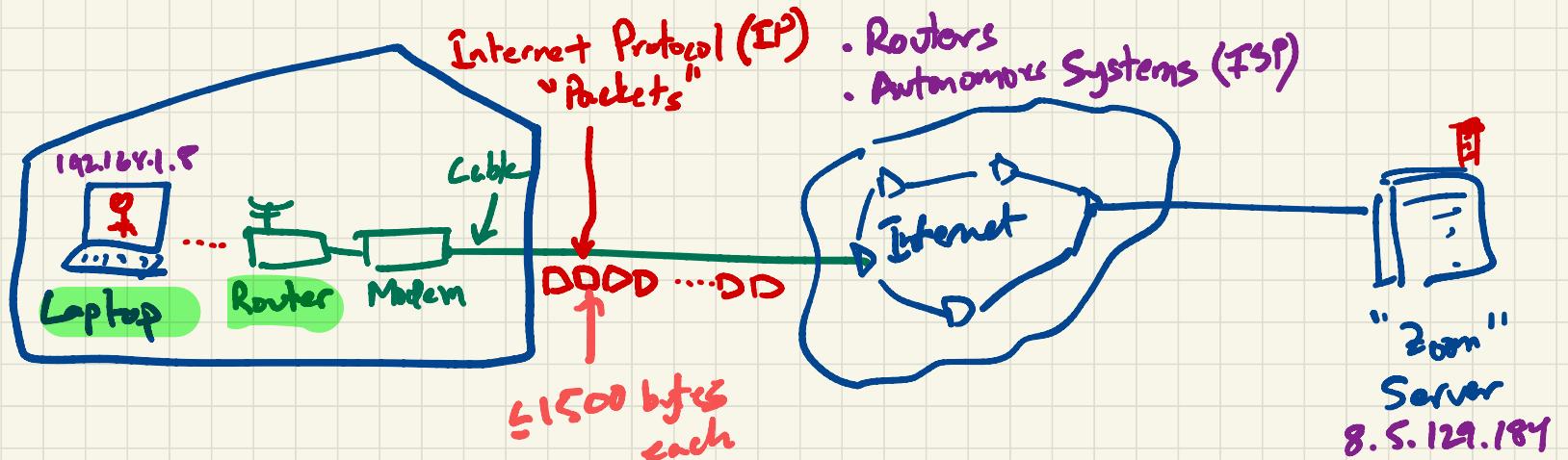
The Internet



Internet Protocol "Stack"



The Internet

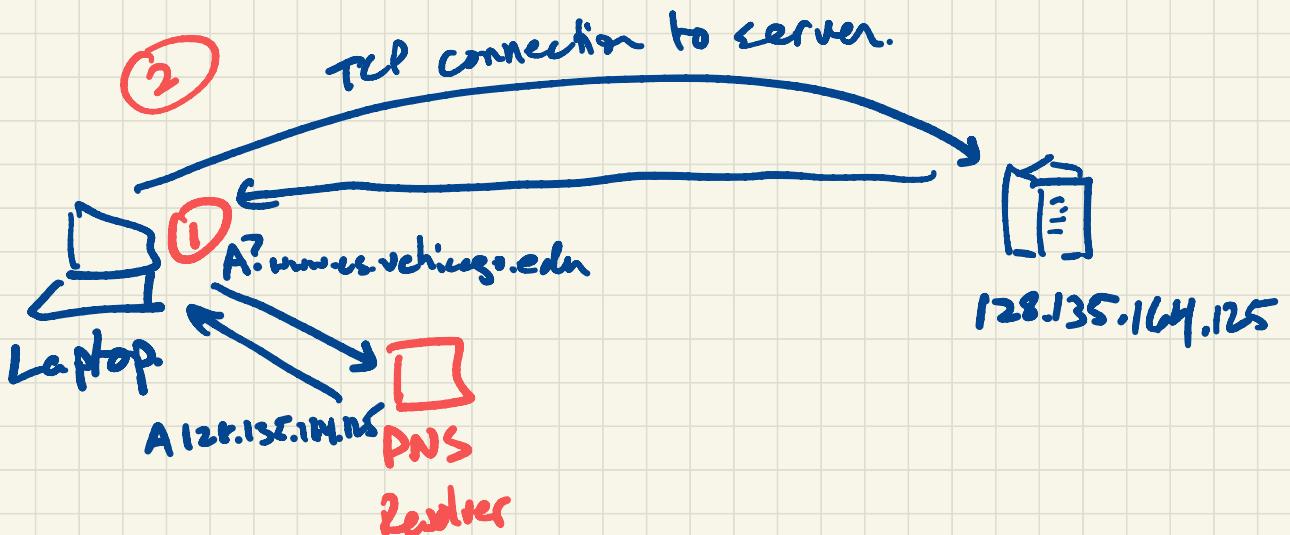


Domain Name System (DNS)

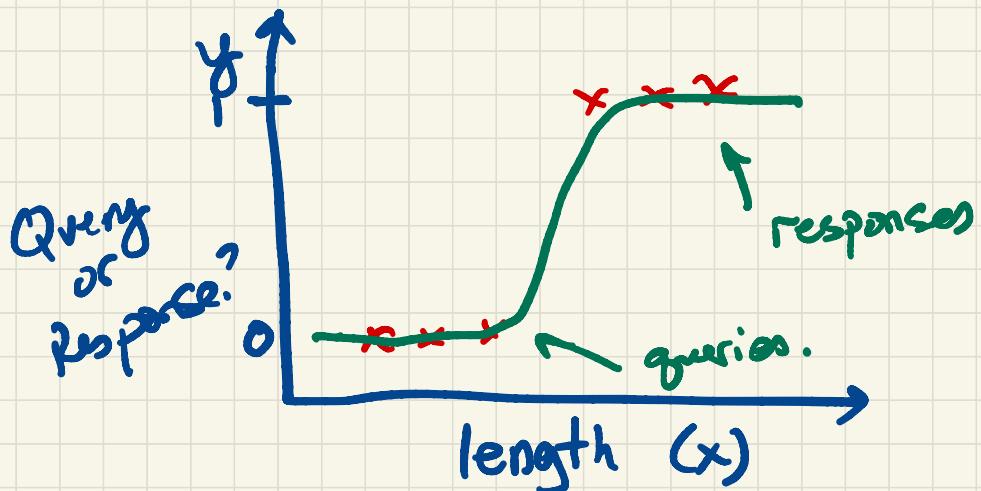
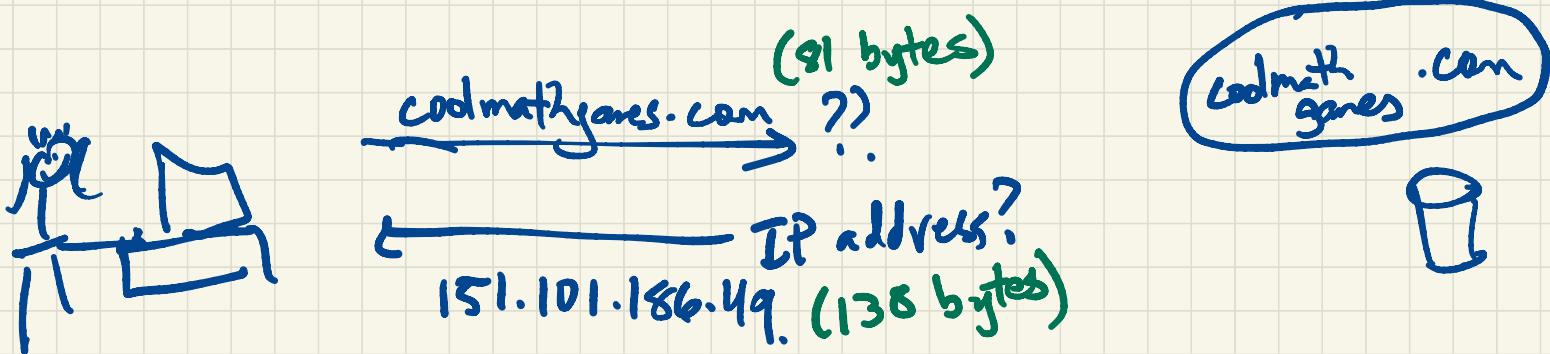
Name: www.cs.uchicago.edu

Answer: 128.135.164.125

- (1) DNS Lookup
- (2) Connect to server



Domain Name System (DNS)



HOME

- ① How does my computer connect to the Internet?
- ② What other "things" are connected to the Internet in my house?
- ③ How do so many devices in my house "share" the same Internet connection?
- ④ What places in your house have good WiFi performance/connectivity? Bad?

Applications

① What happens when I ..

- download a web page?

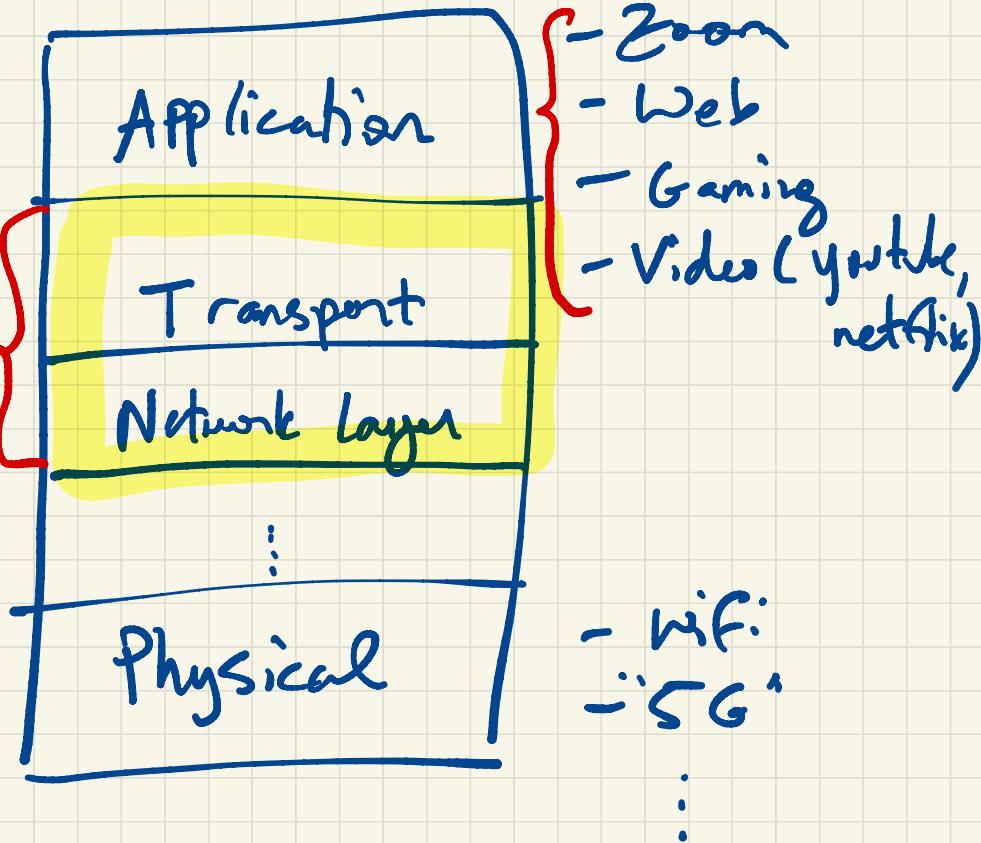
- send an email?
- stream video?

TCP / IP

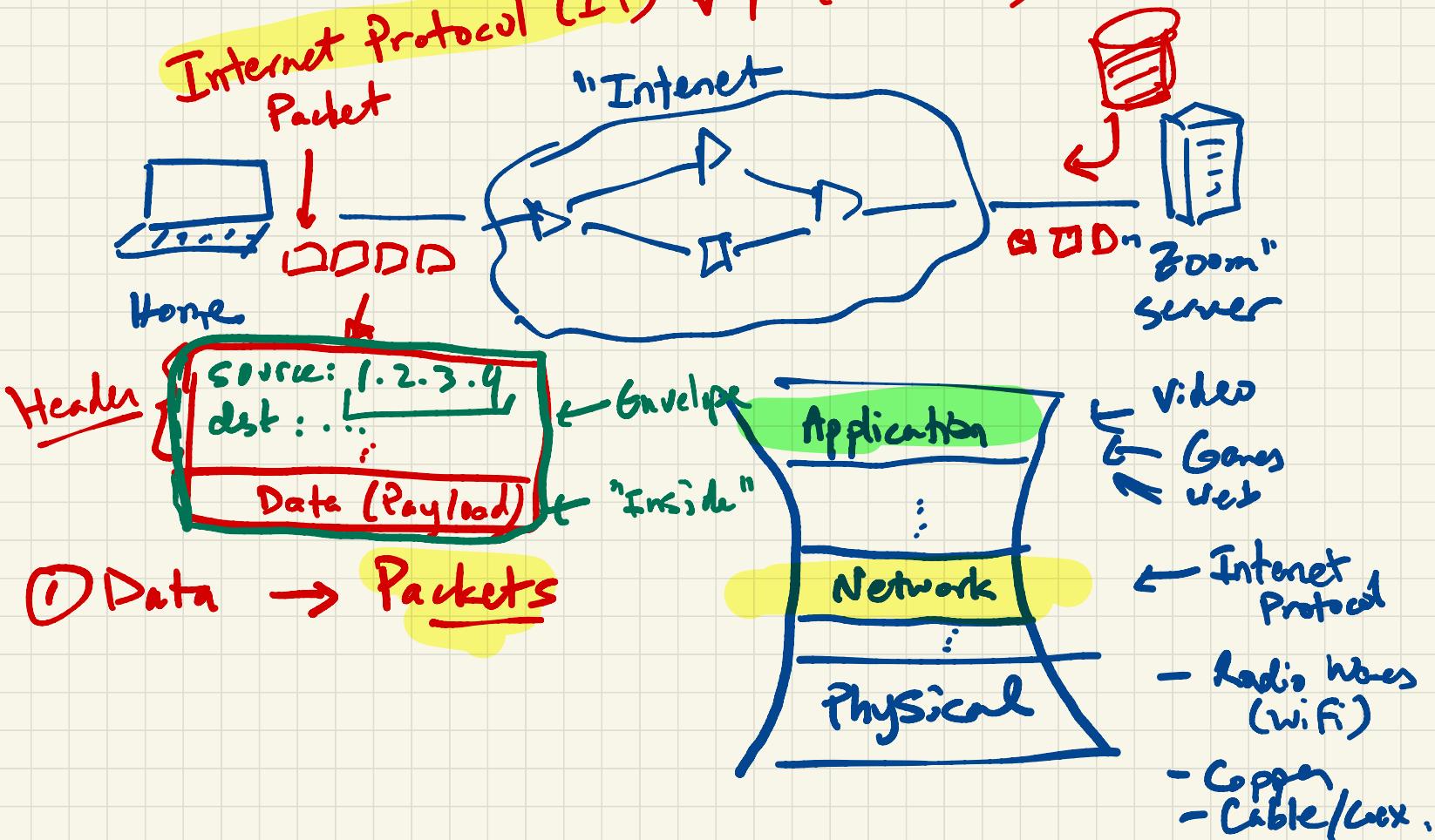
"plumbing"

① RELIABLE

② IN ORDER



Internet Protocol (IP) v4 (also v.6)

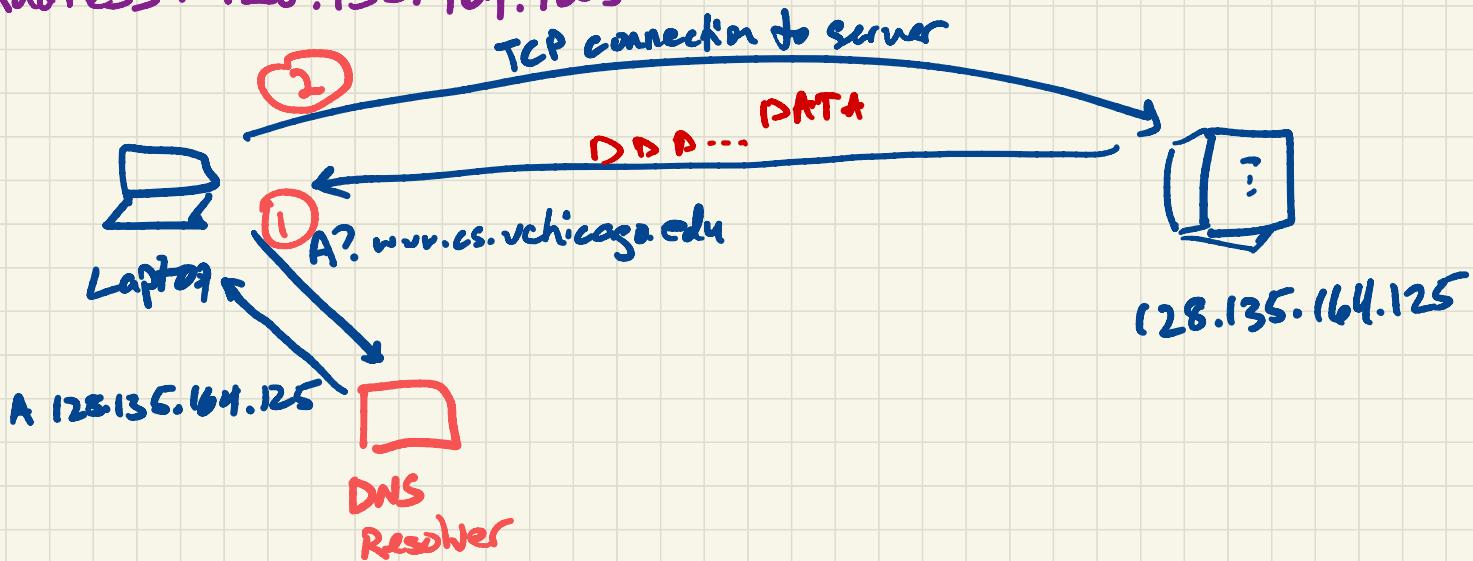


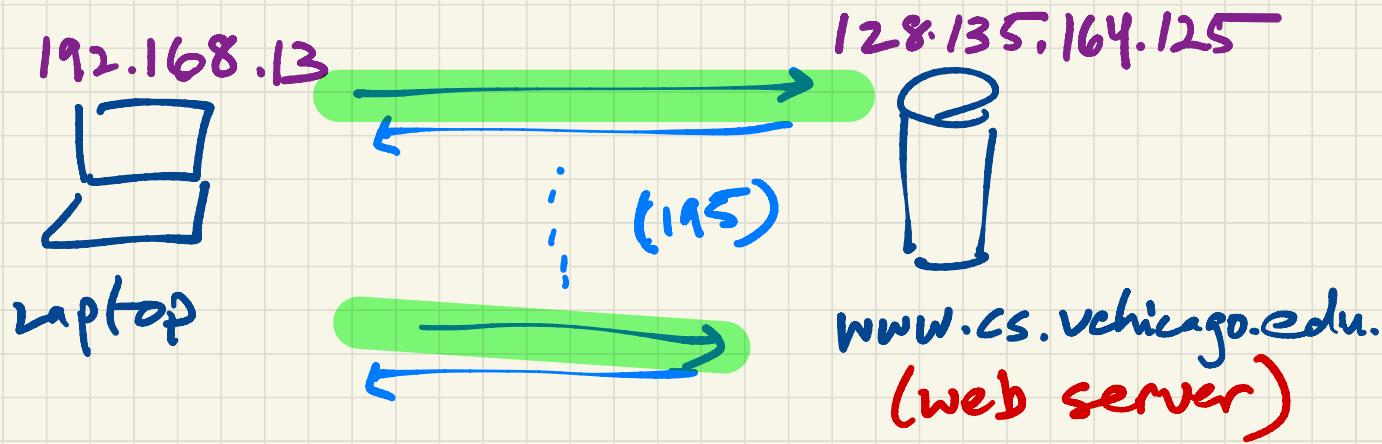
Domain Name System (DNS)

Name: www.cs.uchicago.edu

Address: 128.135.164.125

- ① DNS Lookup
- ② Connect to server
(TCP/IP connection)





MACHINE LEARNING

Input Data

X : features
 y : target

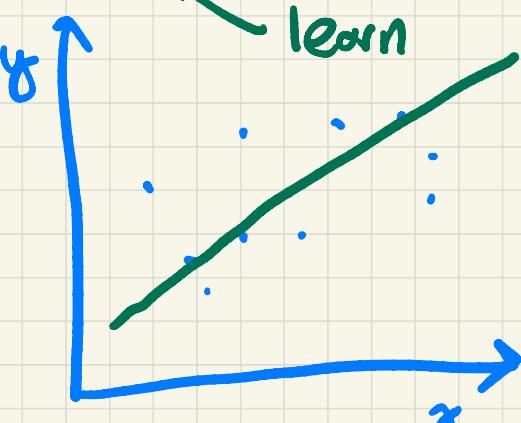
fit()

Train Model

predict()

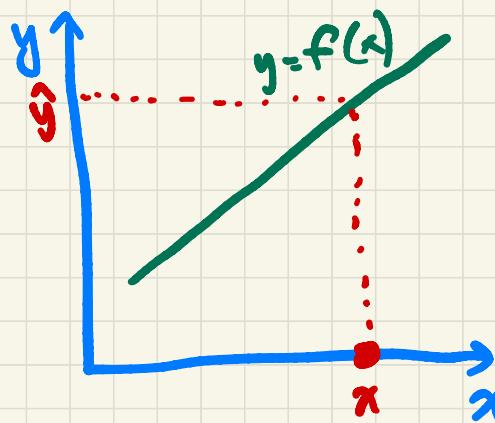
(e.g., Linear Regression)

$$y = f(x)$$

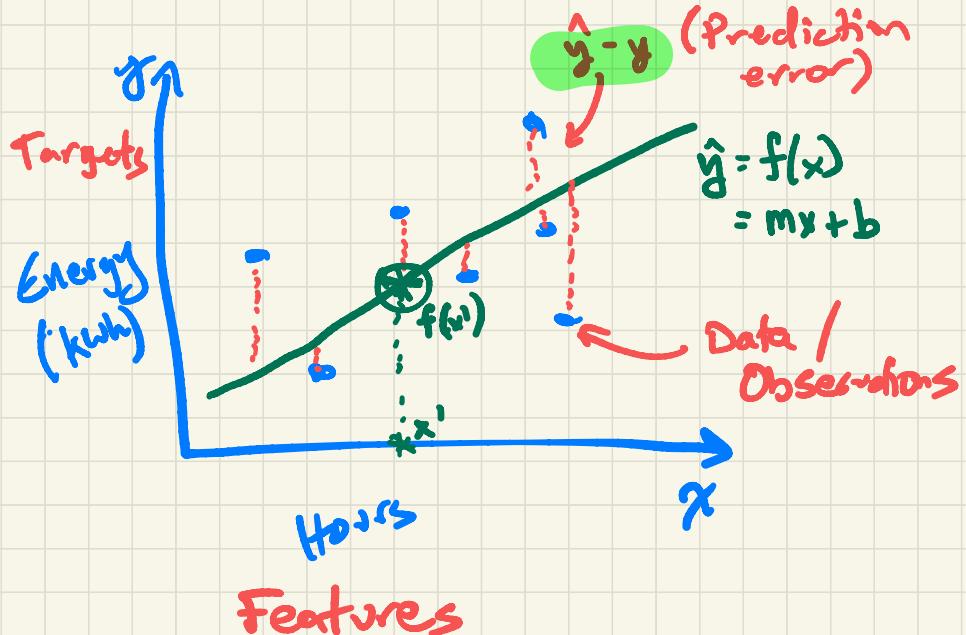


$$\hat{y} = f(x)$$

Prediction.



Linear Regression



must learn

$$\{y_1, \dots, y_n\} = f(\{x_1, \dots, x_n\})$$

Minimize Error

$$RSS = \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

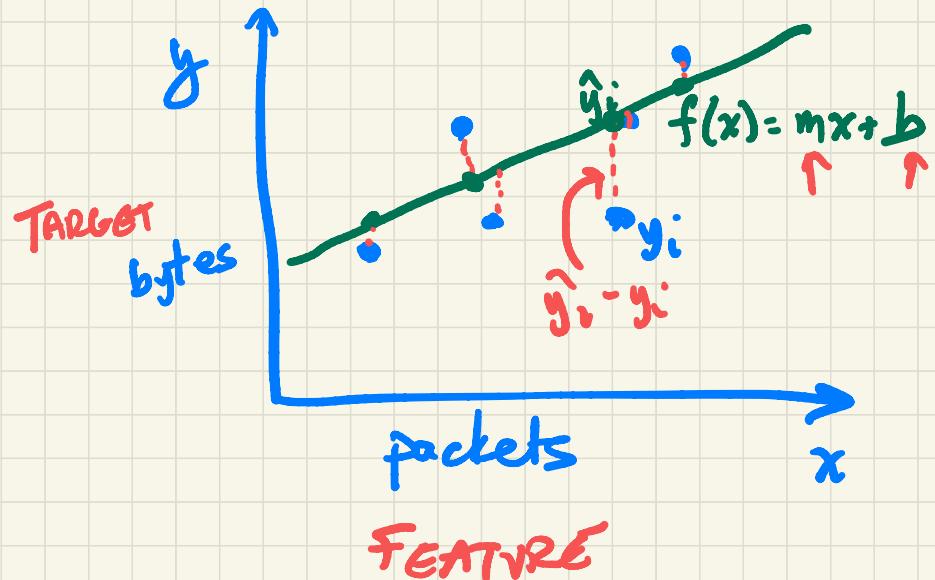
↑
prediction ↓
actual value

$$\hat{y} = f(x) = mx + b$$

↑
slope ↑
intercept

Note: Can solve in closed form. (Linear Algebra)

Linear Regression (1 feature)



$$\{y_1, \dots, y_n\} = f(\{x_1, \dots, x_n\})$$

target features
must learn
 $\hat{y}_i = f(x_i)$ linear

Goal: Minimize Error

$$RSS = \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$\hat{y} = f(x) = mx + b$$

Note: Closed form solution.

Linear Regression Optimization

$$\hat{y} = m\vec{x} + b$$

$$\hat{Y} = \hat{\beta}_0 + \sum_{i=1}^n x_i \hat{\beta}_i$$

↑ ↑ multiply
 b m x variable
 $\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots$

$$\hat{Y} = X^T \hat{\beta}$$
 (linear equations)

Solving a set of linear equations.

Error

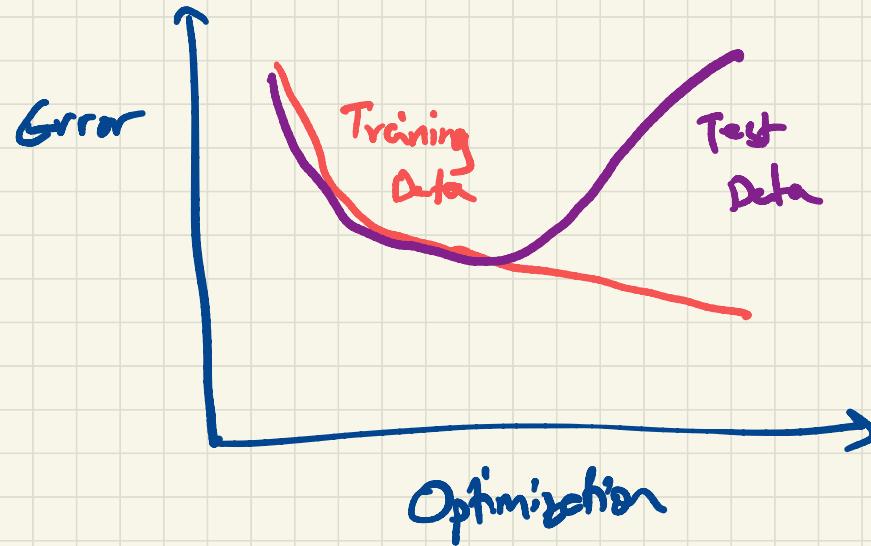
$$\begin{aligned}
 \text{RSS}(\beta) &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\
 &= \sum_{i=1}^n (y_i - x_i^T \beta)^2 \\
 &= (y - X\beta)^T (y - X\beta)
 \end{aligned}$$

Minimize:

$$X^T (y - X\beta) = 0$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

Bias - Variance Tradeoff

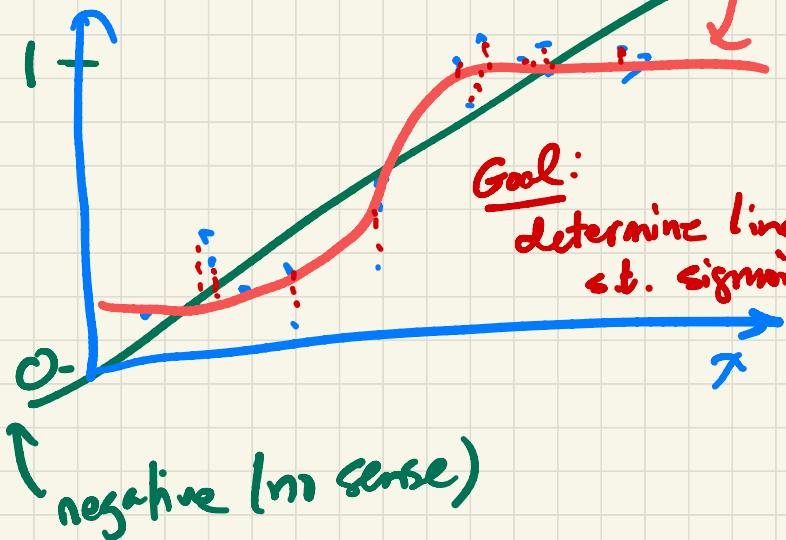


Logistic Regression

- Linear model.
- For all values of $x, y \in [0, 1]$
- Reason: Want to predict binary value

(e.g., spam/ham, malware, attack, ...) \rightarrow yes/no

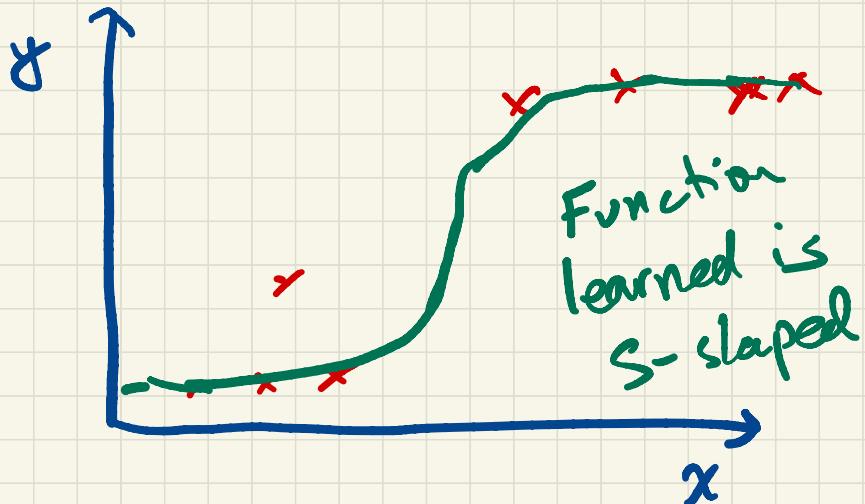
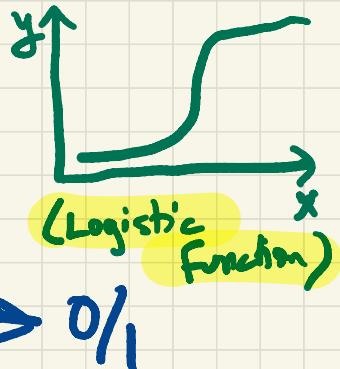
sigmoid
↓
 $y = f(dx)$



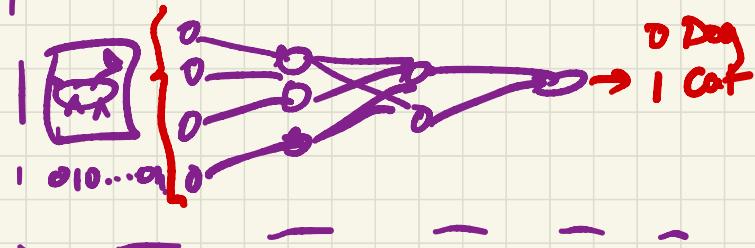
$$\hat{y} = \underbrace{\sigma(mx + b)}_{\text{output is } 0/1.}$$

Sigmoid Function: $\sigma(x) = \frac{1}{1+e^{-x}}$

Logistic Regression

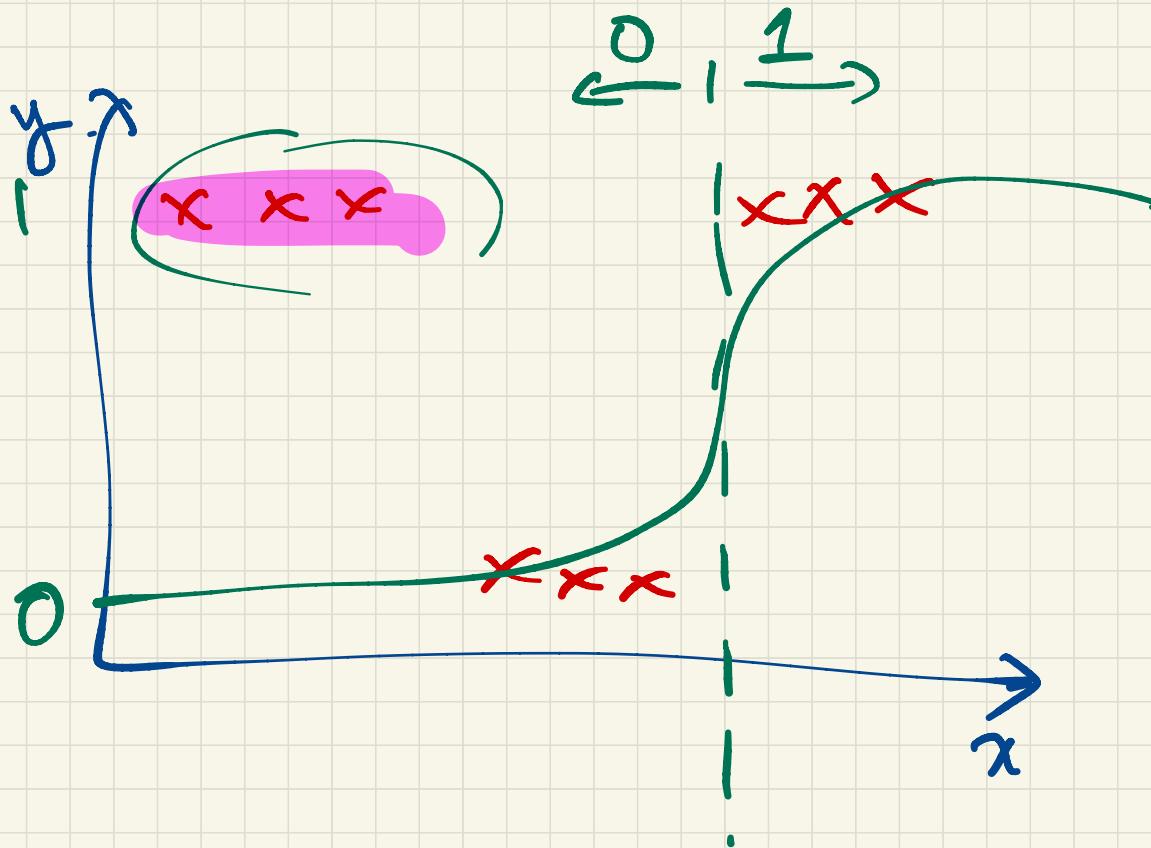


(“Deep Learning”)

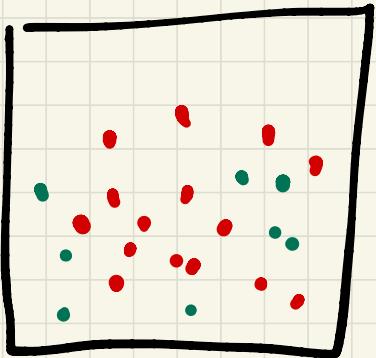


Examples: → Yes or No?
 $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$

- ① Spam email or ham/brea?
- ② Is my laptop infected?
- ③ Is the Internet under attack?



Basic Probability



100 balls total

75 red balls

25 green balls

$$P(\text{red ball}) = \frac{75}{100} = 0.75$$

$$P(\text{green ball}) = \frac{25}{100} = 0.25$$

$$P(\text{red ball, green ball}) = ?$$

① Put red ball back?

$$P(\text{red}) \cdot P(\text{green}) = \left[\frac{75}{100} \right] \times \left[0.1875 \right]$$

② Conditional Probability ② Keep red ball out?

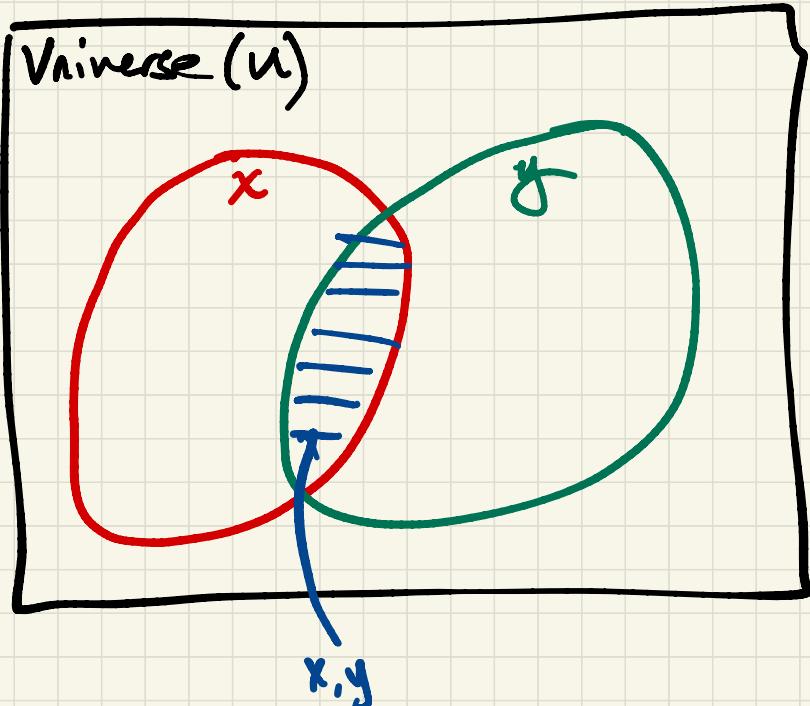
$$P(A|B)$$

$$P(B, A) = P(B|A) \cdot P(A)$$

$$P(\text{red}) \cdot P(\text{green}|\text{red}) = \left[\frac{75}{100} \right] \times \left[\frac{25}{91} \right] \times [0.19]$$

Bayes' Rule

$$P(y|x) = \frac{P(x|y) \cdot P(y)}{P(x)}$$



$$P(x) = \frac{|x|}{|U|} \quad P(y|x) = \frac{|xy|}{|x|} = \frac{\frac{|xy|}{|U|}}{\frac{|x|}{|U|}} = \frac{P(xy)}{P(x)}$$

$$P(y) = \frac{|y|}{|U|} \quad P(x,y) = \frac{|xy|}{|U|}$$
$$P(x|y) = \frac{P(xy)}{P(y)}$$

$$P(y|x) \cdot P(x) = P(x|y) \cdot P(y)$$

$$P(y|x) = \frac{P(x|y) \cdot P(y)}{P(x)}$$

predict observe

Naïve Bayes Classifier

$$P(y|x) = \frac{P(x|y) \cdot P(y)}{P(x)}$$

Predict \hat{y} , having observed x .

$$\hat{y} = \operatorname{argmax}_y \left\{ P(y=c | x=x_1, x_2, x_3, \dots, x_n) \right\}$$

(e.g., spam, ham)

observed or estimated (easy) *need to compute!*

$$= \operatorname{argmax}_y \left\{ \frac{P(x=x_1, \dots, x_n | y) \cdot P(y)}{P(x)} \right\}$$

easy

$$= \operatorname{argmax}_y \left\{ P(x_1|y) \cdot P(x_2|y) \cdots P(x_n|y) \cdot P(y) \right\}$$

Example: Spam filtering

TRAINING

SPAM ($y=SPAM$)

D_y

- MONEY, RILEX
- MONEY, PILLS
- PHARMACY, AGE

$$P(MONEY | SPAM) = \frac{2}{3}$$

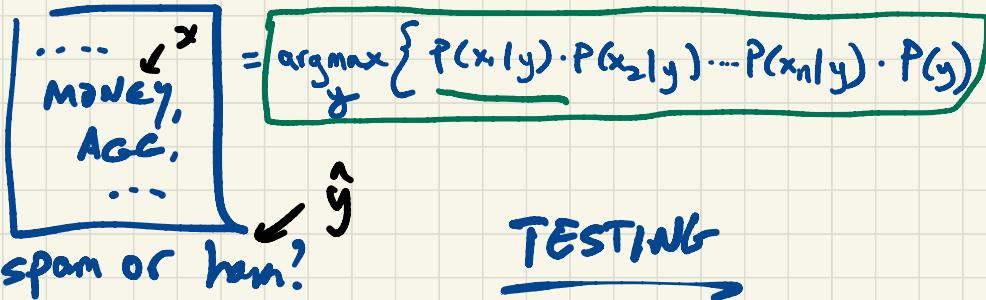
HAM ($y=HAM$)

D_y

- CLASS
- MACHINE LEARNING
- AGE

$$P(MONEY | HAM) = 0$$

$$P(AGE | HAM) = \frac{1}{3}$$



$$\textcircled{1} \quad P(y=SPAM | MONEY, AGE)$$

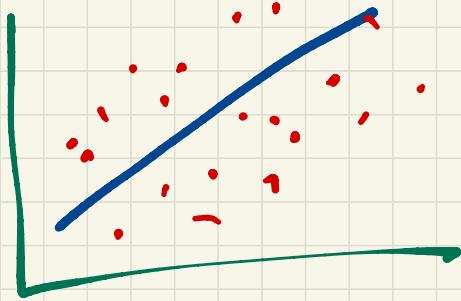
$$\approx P(MONEY | SPAM) \cdot P(AGE | SPAM) \cdot P(SPAM)$$

$$\approx \frac{2}{3} \cdot \frac{1}{3} \cdot \frac{9}{10} \approx \boxed{\frac{18}{90}}$$

$$\textcircled{2} \quad P(y=HAM | MONEY, AGE)$$

$$\times P(MONEY | HAM) \cdot P(AGE | HAM) \cdot P(HAM)$$

$$\approx 0 \cdot \frac{1}{3} \cdot \frac{1}{10} = \boxed{0}$$



Machine Learning

Supervised Learning
(data & predictions have labels)

Linear

Linear Regression
(continuous)

Probabilistic

Logistic Regression
(binary)

Naive Bayes

Unsupervised Learning
(no labels)

Clustering

Today.

Thursday

Nearest Neighbor

Example:

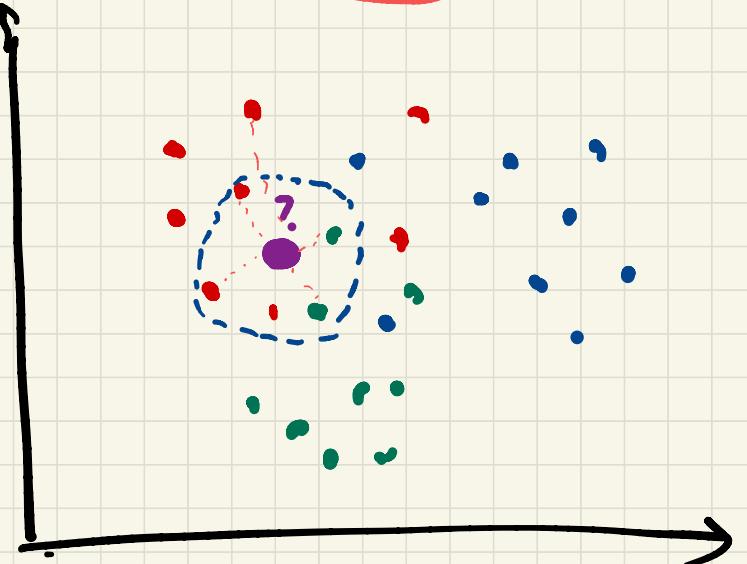
5-NN

- 3 red

- 2 green

↓
RED

(k-NN) Not known! (Explore)



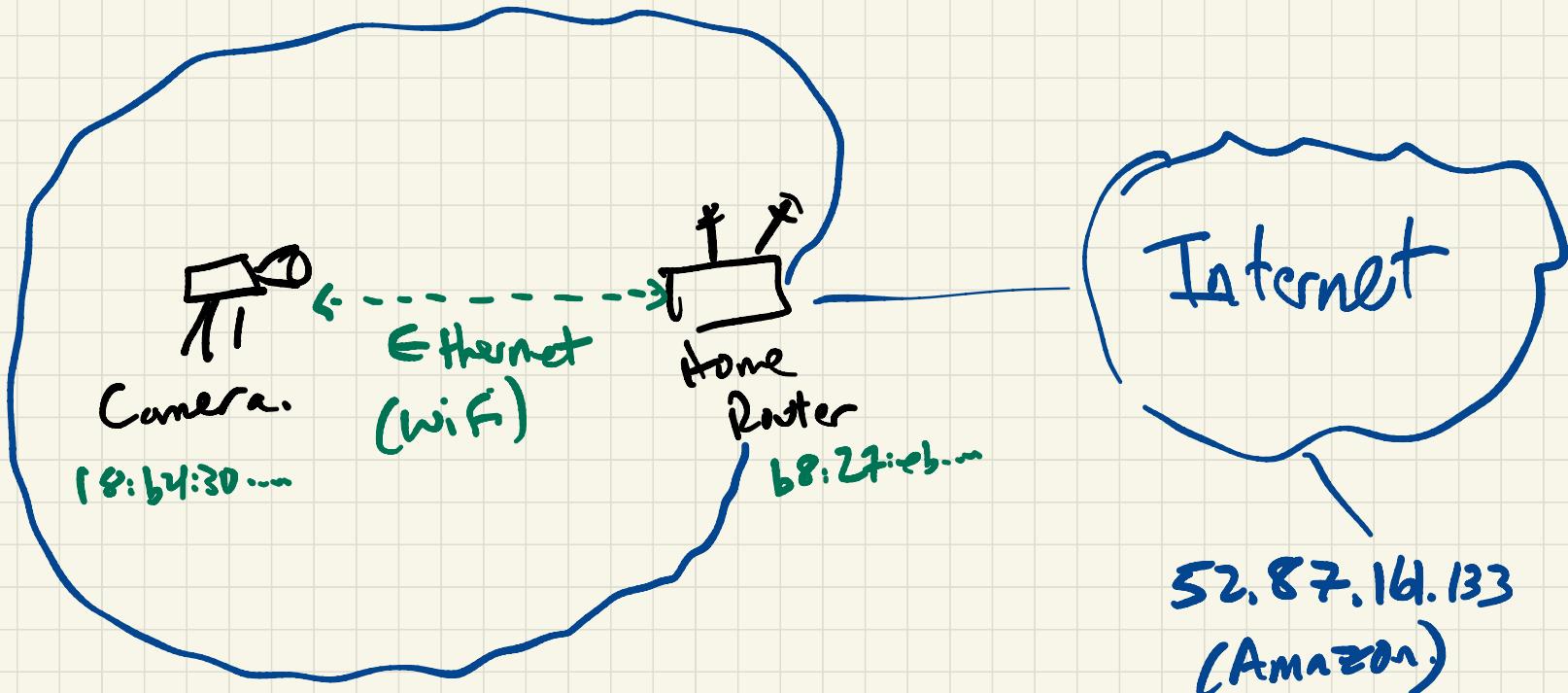
To compute:

→ Need to store
all data points

Intuition:

Point is of type
that is the same
as other "close"
points.

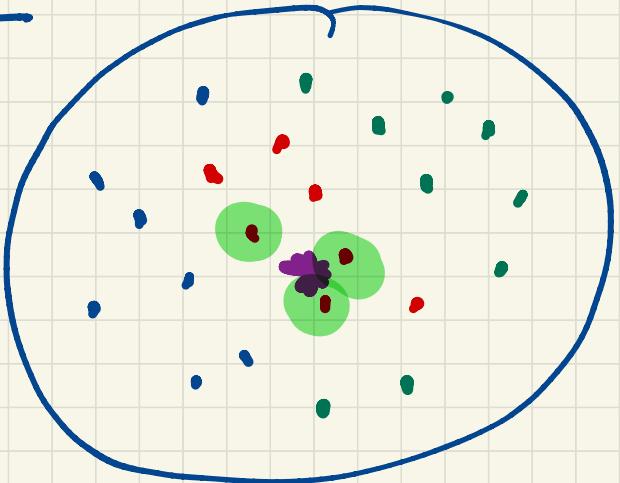
- ① Other: How many?
- ② Close?
→ Euclidean



Newest Neighbor

Parameter: k

e.g., $k=3$



$$* = \{ \cdot, \cdot, \cdot, \cdot \}$$

$\rightarrow *$

In this example, k closest points all red.

Advantages: No training \rightarrow store data points.

Disadvantages: Classification can be expensive \leftarrow computation
Must store entire dataset. \leftarrow storage

Tree-Based Methods

Idea: Split feature space into simpler regions.

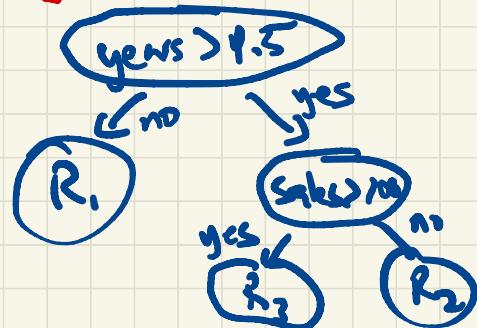
- Uses a sequence of rules
- Sequence of rules summarized in a tree

+ Easy to interpret.

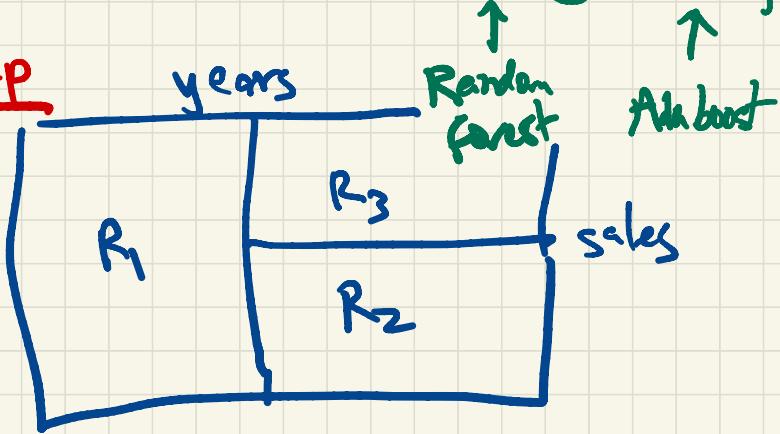
- Single tree not that accurate (Solution: Bagging, Boosting)

Example

Tree



Map



Regression Tree

- Predict Continuous Values.
- Splits in tree → Reduction in mean Squared error

Greedy Approach:

$$R_1(j, s) = \{x | x_j \leq s\}$$

$$R_2(j, s) = \{x | x_j > s\}$$

Minimize w.r.t S

$$\sum_{i \in R_1(j, s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i \in R_2(j, s)} (y_i - \hat{y}_{R_2})^2$$

↑
true values ↑
prediction

✓

Classification Trees

- Predict Qualitative/Categorical Outcome.
- Split criteria: Classification Error Rate
 - fraction of training observations that don't belong to most common class in region.

Gini Index

$$G = \sum_{k=1}^K \hat{P}_{mk} (1 - \hat{P}_{mk})$$

Entropy

$$D = - \sum_{k=1}^K \hat{P}_{mk} \log \hat{P}_{mk}$$

Ensemble Methods

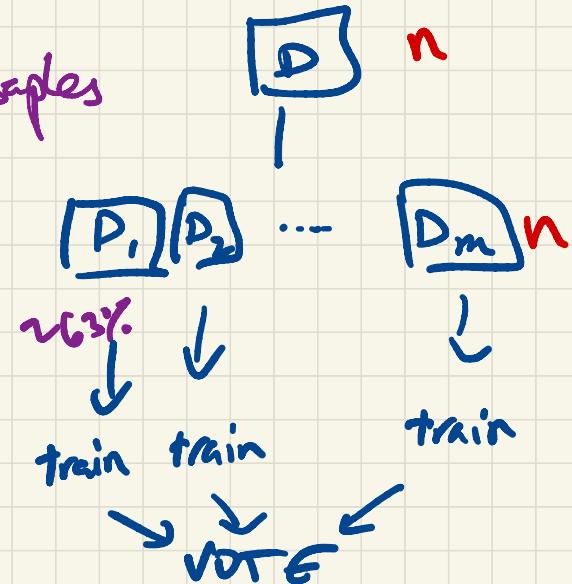
Problem: Decision trees have high variance

Bagging (Bootstrap Aggregation)

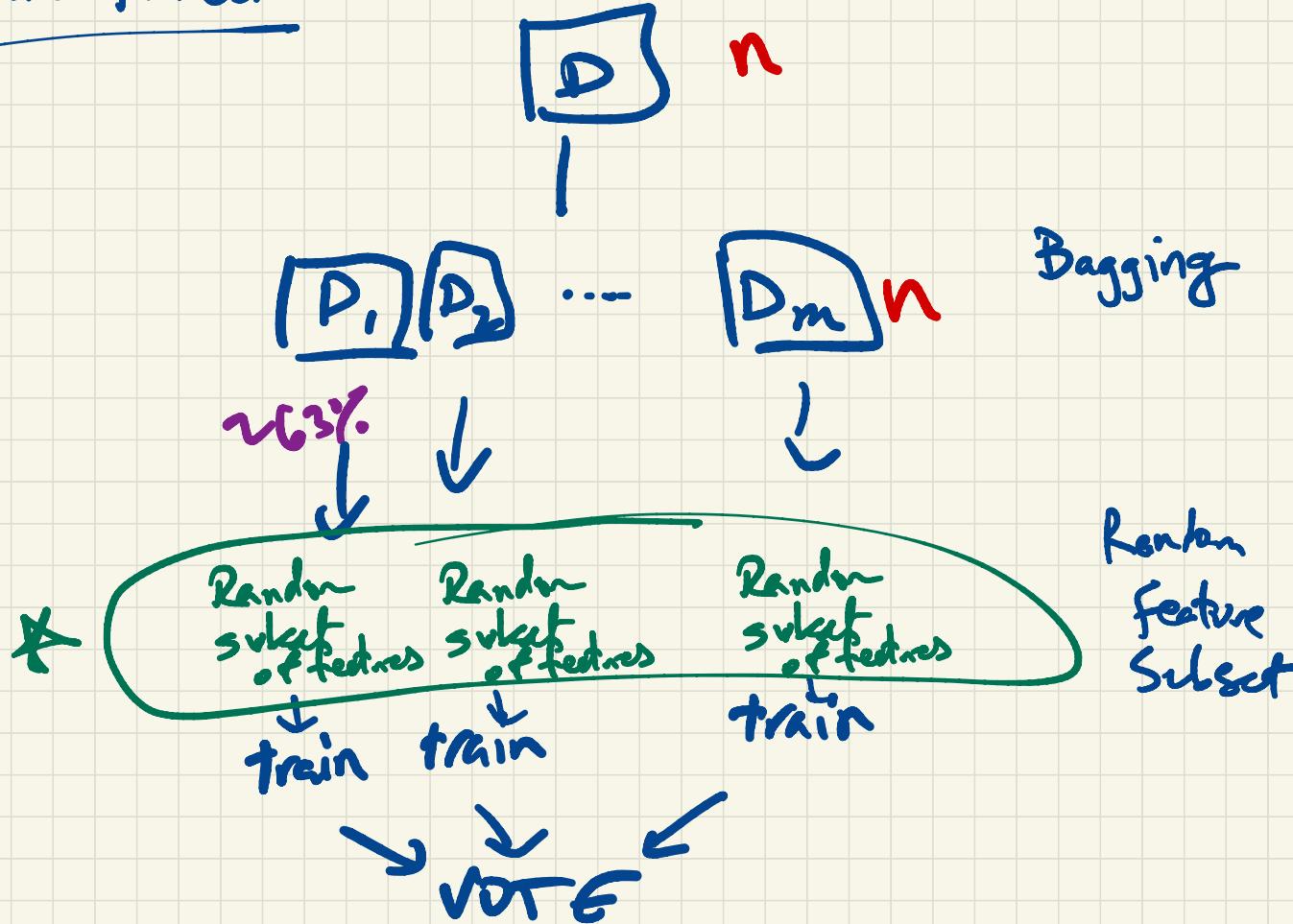
- Separate Model w/ repeated random samples
- Average predictions
- training set $\rightarrow n$

M samples of size n

✓

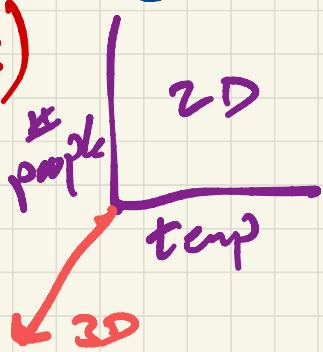


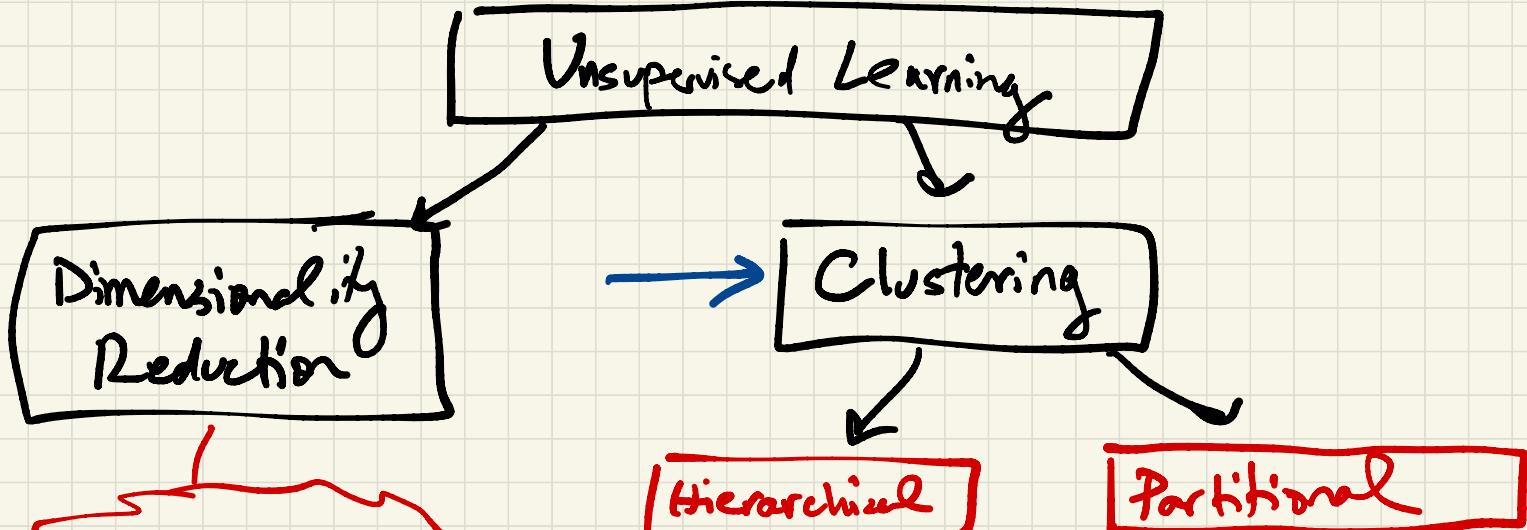
Random Forest



Why Unsupervised Learning?

1. Might be challenging / costly to get labels
2. May want to solve a problem other than making a prediction
 - Visualization (≤ 3 dimensions)
 - Understand groups / structure / similarity
 - Dimensionality Reduction.

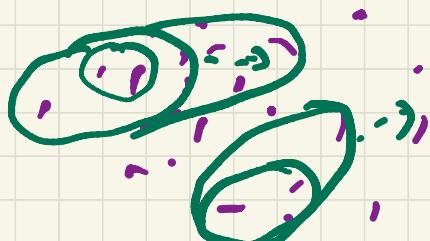




"Linear Algebra"

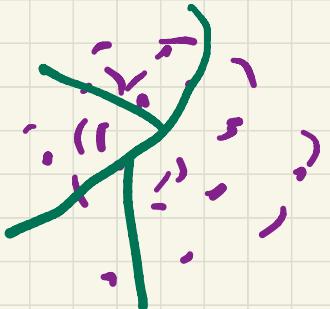
Hierarchical

- ① • linkage-based clustering



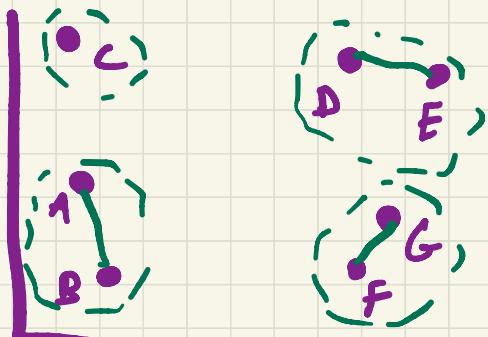
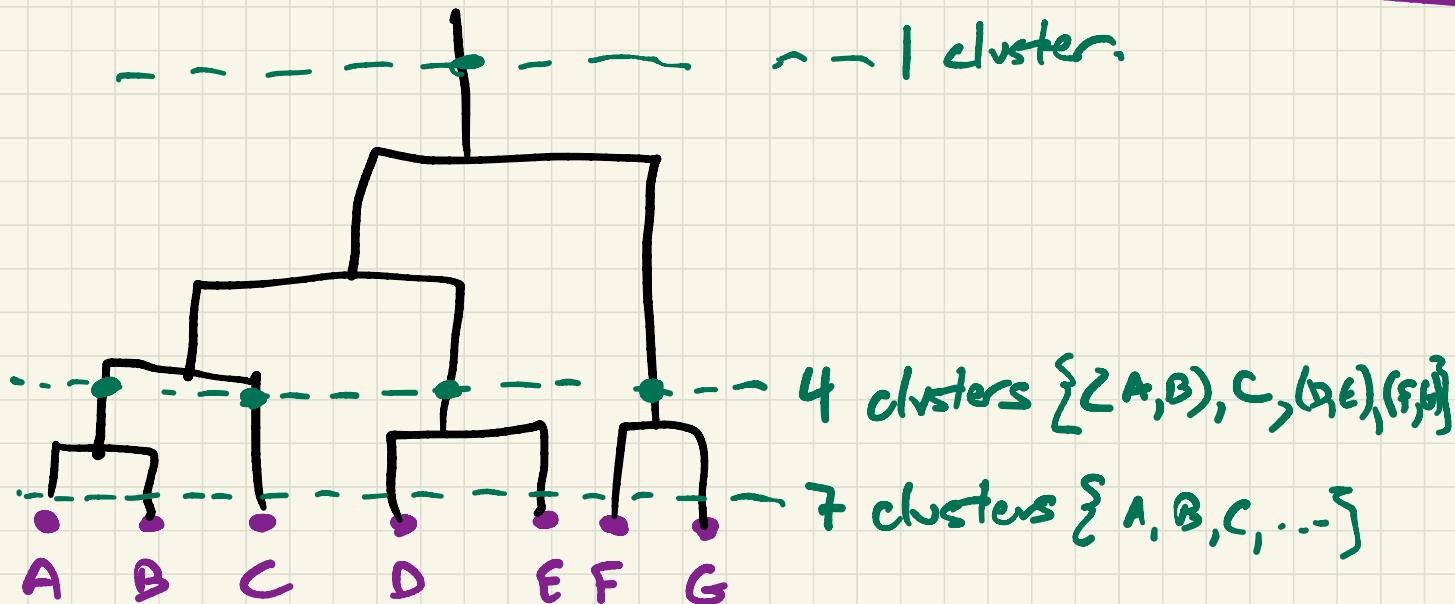
Partitional

- ② • k-means
③ • DBSCAN.



Linkage-Based Clustering / Hierarchical

Dendrogram



Not All Data is Hierarchical

Data: People. { gender, nationality }
(M, F); (US, Canada, Mexico)

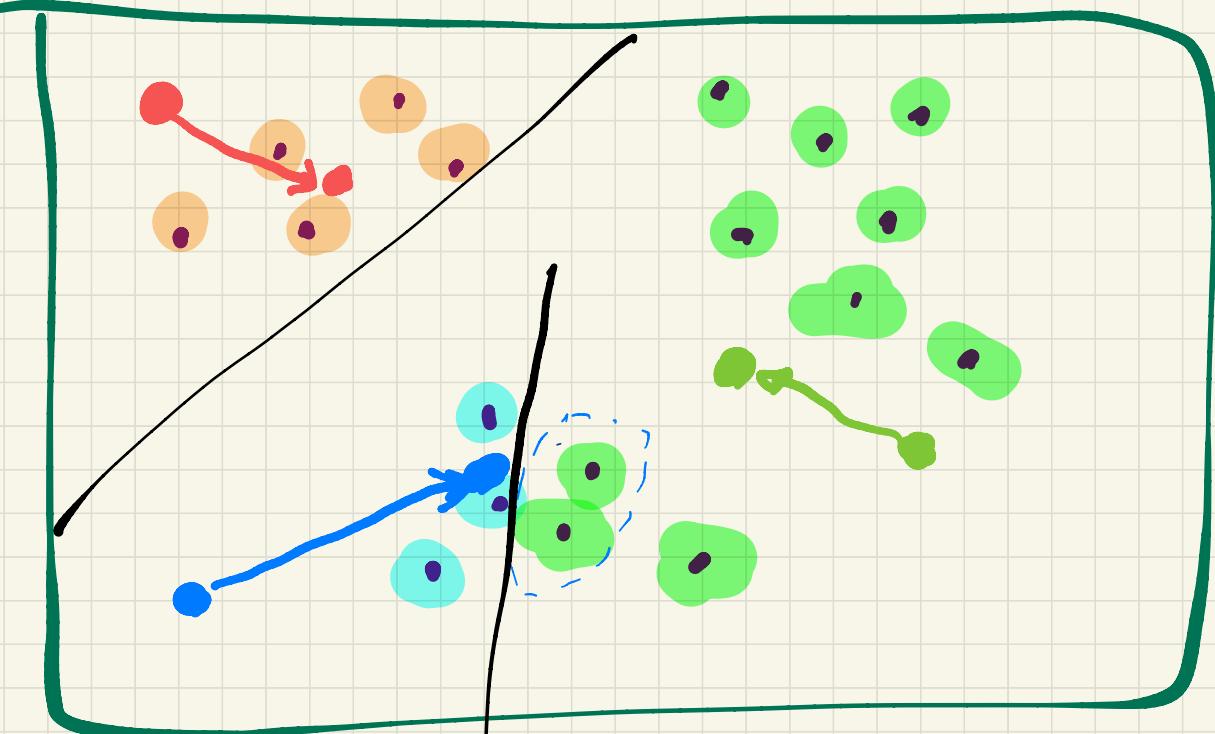
$k=2 \rightarrow$ best clustering might be based on gender

$k=3 \rightarrow$ best clustering might be based on nationality

k-means clustering

must select!

- ① e.g., $k=3$
- ② Randomly choose cluster centers
- ③ Recompute cluster centers
- ④ Repeat 2 & 3, ...



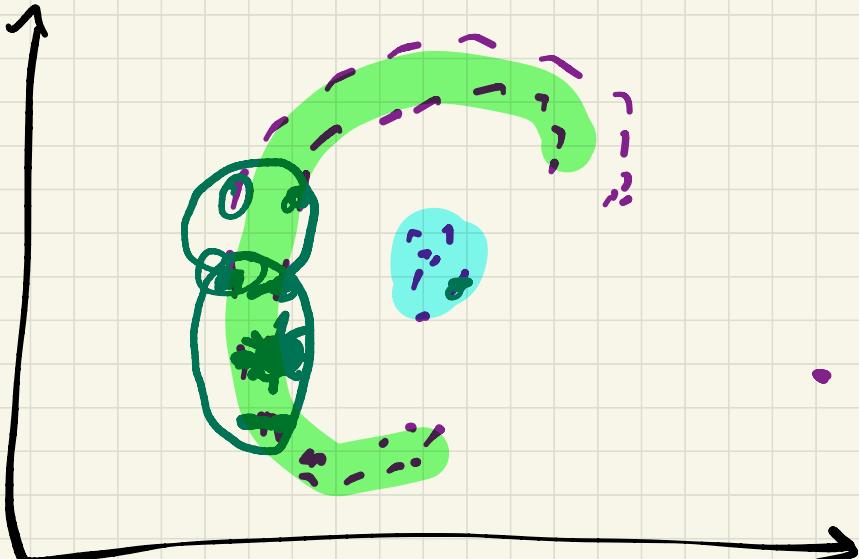
Density-Based Clustering (e.g., DBSCAN)

Concept: Regions with high density in same cluster.

① Point is in a cluster if it has a min. # of 'neighbors'

② Repeat for all point's neighbors

③ Pick new random point



k-means failure scenarios

