# CS 360: Machine Learning

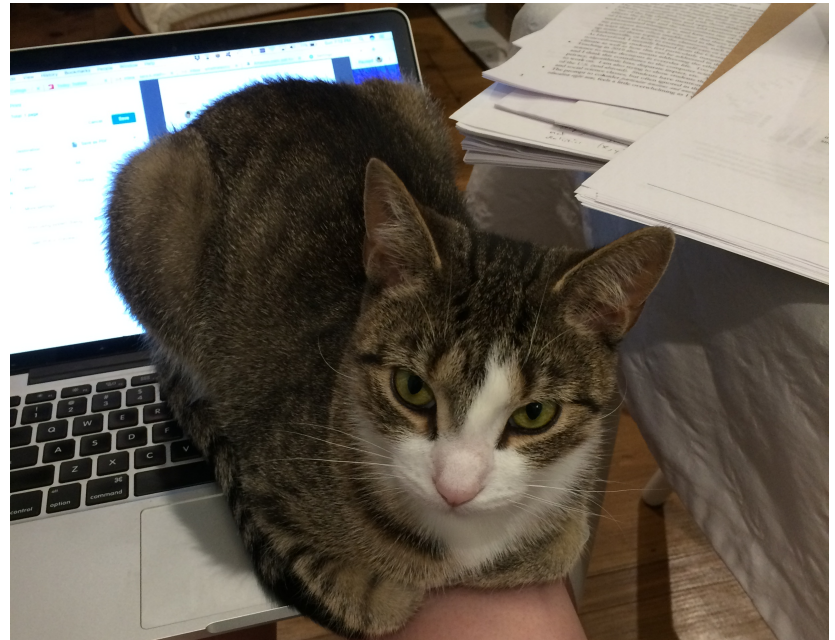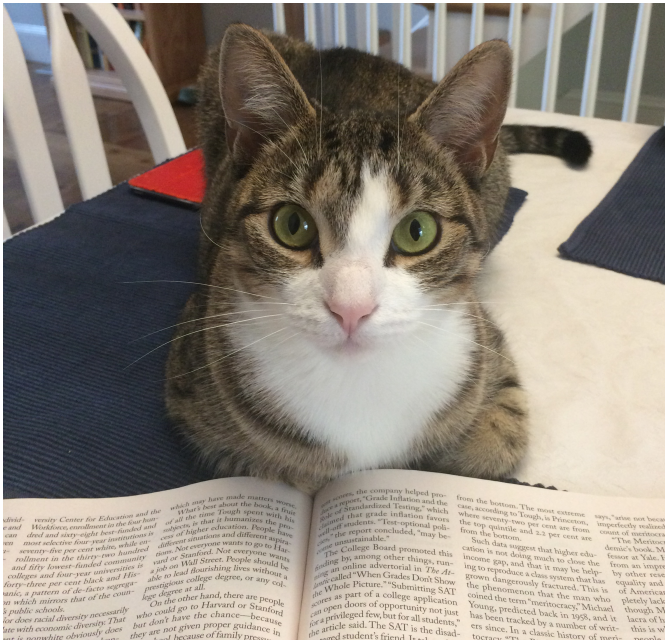## Prof. Sara Mathieson

## Fall 2019

# Admin

- **Lab 5 TODAY!**
  - Office hours today 12:30—1:30pm (H110)

- **Reading Quiz Thursday** (Duame Section 13.1)

- Lab 6 due Friday Nov 1
  - Checkpoint during lab on Thursday Oct 31 (Part 1 and 2)

# In lab Thursday

- Hand back the midterm
- Go over common issues
- Start Lab 6

# Outline for October 22

- Evaluation metrics
  - Confusion matrices revisited
  - ROC curves
  - Relationship to probabilistic methods

- Ensemble methods
  - Bagging
  - Random forests

# Outline for October 22

- Evaluation metrics
  - Confusion matrices revisited
  - ROC curves
  - Relationship to probabilistic methods

- Ensemble methods
  - Bagging
  - Random forests

# For now: assume binary classification task

- Transactions that indicate credit card fraud
- Detecting which scans show tumors
- Prenatal test for Down's Syndrome
- Finding genes under natural selection
- Finding regions of the genome with high recombination rate ("hotspots")

# For now: assume binary classification task

- Transactions that indicate credit card fraud

- Detecting which scans show tumors

- Prenatal test for Down's Syndrome

- Finding genes under natural selection

- Finding regions of the genome with high recombination rate ("hotspots")

In all these examples, we are trying to find unusual items ("needle in a haystack") -- we call these *positives*

# Goals of Evaluation

- Think about what metrics are important for the problem at hand

- Compare different methods on the same problem

- Common set of tools that other researchers/users can understand

# Back to Confusion Matrices...

pred class

|  | negative | positive |
|---|---|---|
| negative | 70 | 20 |
| positive | 8 | 15 |

true class

- false positive
- FP
- "false alarm"
- Type I error
  - want low

- TP
- true positive
- want high.
- "flagged"

- TN
- high (want)
- true negative

- false negatives
- "miss"
- FN
- type II error
- want low

**Recall**: how many positives were found?

true positive rate

$$TPR = \frac{TP}{FN + TP}$$

$$= \frac{15}{8 + 15} \approx 0.65$$

**Precision**:

$$= \frac{TP}{FP + TP}$$

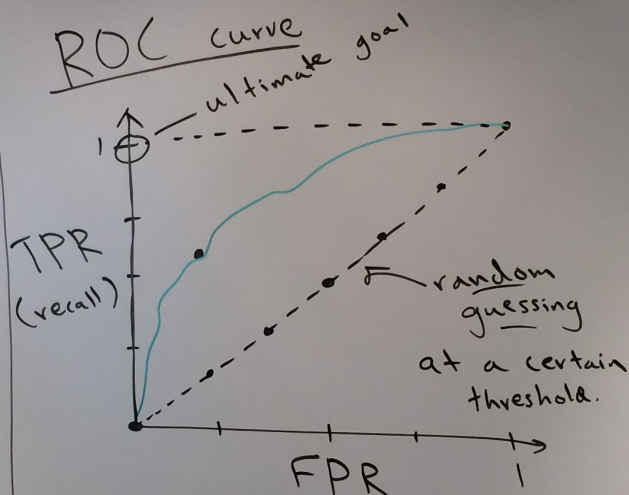$$= \frac{15}{20 + 15} \approx 0.43$$

false positive rate

$$FPR = \frac{FP}{FP + TN}$$

$\underline{\text{want}}$
$\underline{\text{low}}$

$$= \frac{20}{70 + 20} \approx 0.22$$

ROC curve



ultimate goal

TPR (recall)

random guessing

at a certain threshold.

FPR

1

| 0 | 90 | 90 = N |
|---|----|--------|
| 0 | 23 | 23 = P |

$N^* = 0 \quad P^* = 113$

predict all positive

$$FPR = \frac{90}{0 + 90} = 1$$

$$TPR = \frac{23}{0 + 23} = 1$$

| 90 | 0 |
|----|---|
| 23 | 0 |

$FPR = 0$

$TPR = 0$

predict all negative.

$$\begin{array}{|c|c|} \hline 45 & 45 \\ \hline 12 & 11 \\ \hline \end{array}$$

$$TPR = \frac{11}{23} \approx 0.5$$

$$FPR = \frac{45}{90} \approx 0.5$$

## Probalistic Model

} only at test time

threshold: 0.25

$$p(y=1 \mid x) \begin{cases} > 0.25 \Rightarrow \hat{y}=1 \\ \leq 0.25 \Rightarrow \hat{y}=0 \end{cases}$$

threshold: 0.75

Want to be confident

# Handout 12

**(1)**

|     | N  | P |
|-----|----|---|
| N   | 77 | 3 |
| P   | 13 | 7 |

$N^* = 90 \quad P^* = 10$

**(2)** $N = 80$

$P = 20$

precision $= \dfrac{7}{3+7}$

$= 0.70$

recall (TPR) $= \dfrac{7}{13+7}$

$= 0.35$

FPR $= \dfrac{3}{80} = 0.04$

**(3)**



| 68 | 12 |
|----|----|
| 2  | 18 |

$FPR = \dfrac{12}{80} \approx .15$

$TPR = \dfrac{18}{20} = 0.9$

# Precision and Recall

- <u>Precision</u>: of all the "flagged" examples, which ones are actually relevant (i.e. positive)?


- <u>Recall</u>: of all the relevant results, which ones did I actually return?

# Precision and Recall

- <u>Precision</u>: of all the "flagged" examples, which ones are actually relevant (i.e. positive)?

  (Purity)

- <u>Recall</u>: of all the relevant results, which ones did I actually return?
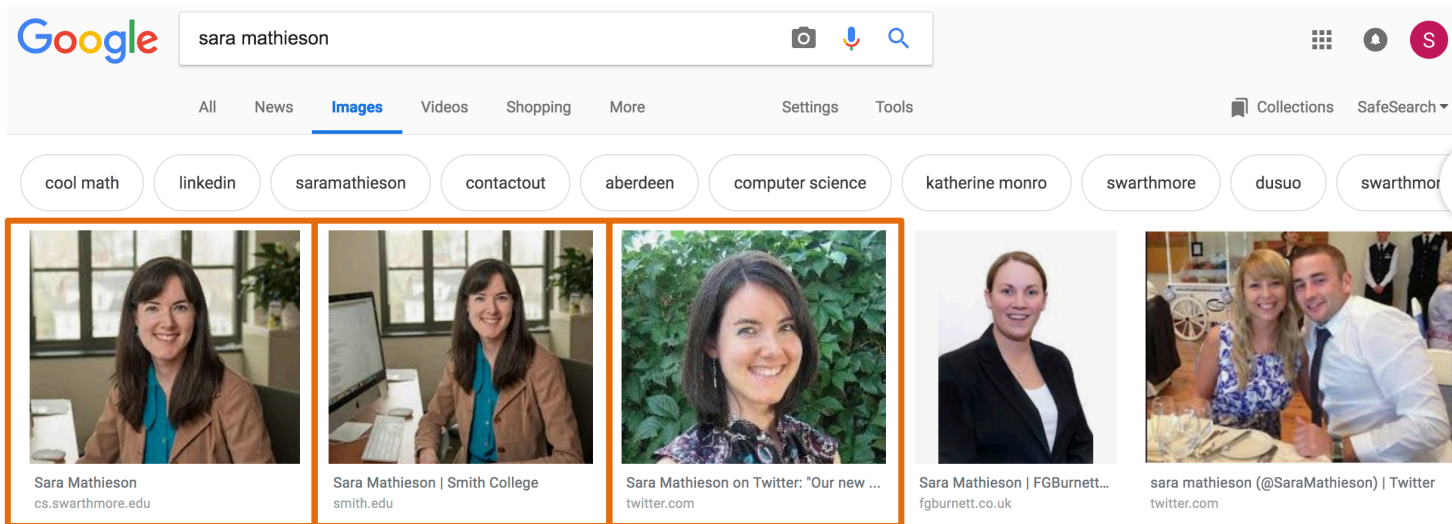
  (Completeness)

# Precision and Recall



*P*=6 (number of images that are actually me)

- Precision?

- Recall?

# Precision and Recall



*P*=6 (number of images that are actually me)
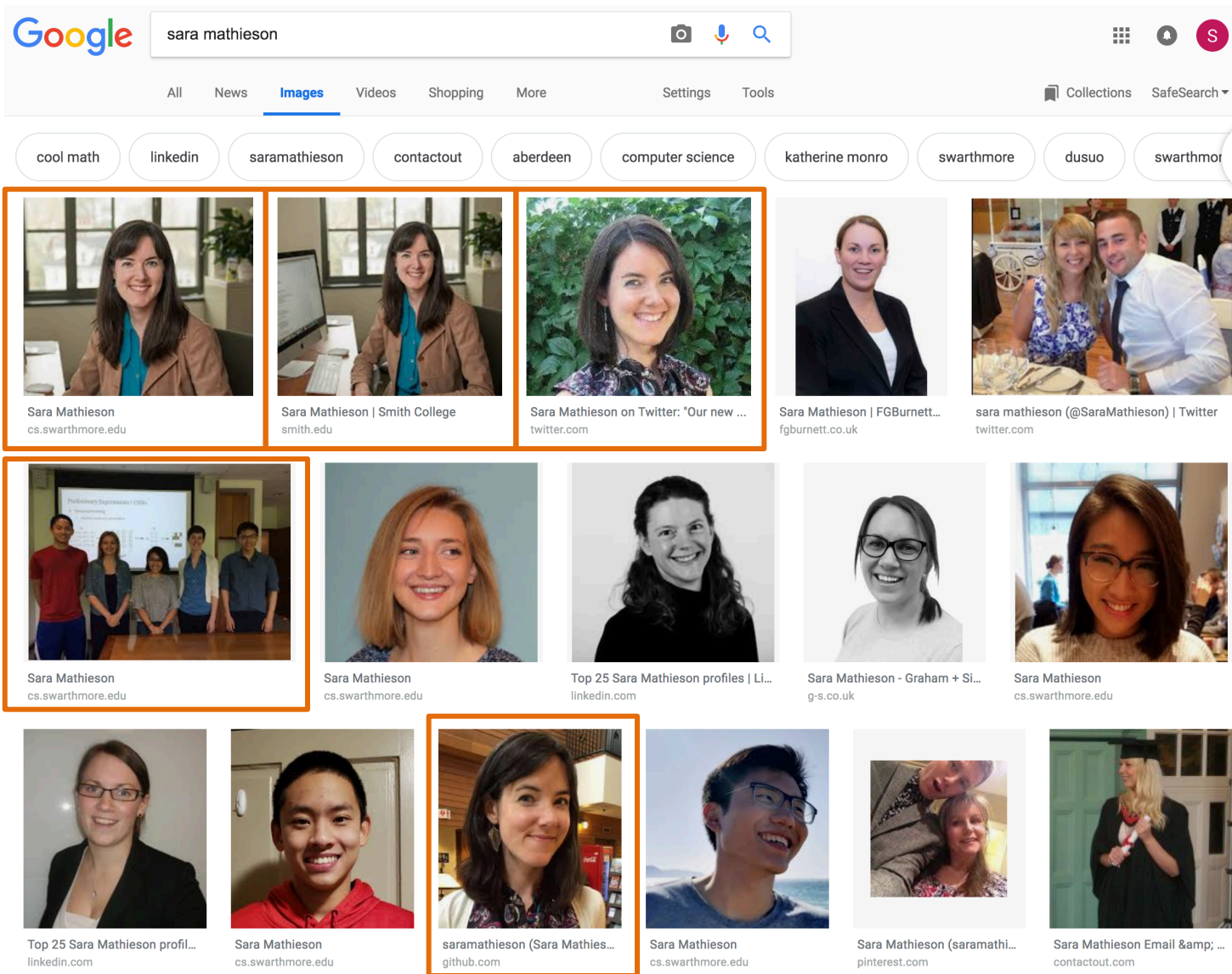
- Precision = TP/(FP+TP) = 3/5

- Recall?

# Precision and Recall



*P*=6 (number of images that are actually me)

- Precision = TP/(FP+TP) = 3/5

- Recall = TP/(FN+TP) = 3/6

# Precision and Recall



*P*=6 (number of images that are actually me)

- Precision = 5/16

- Recall = 5/6

# Recap Confusion Matrices

Predicted class

|  | Negative | Positive |
|---|---|---|
| **Negative** | True negative (TN) | False positive (FP) |
| **Positive** | False negative (FN) | True positive (TP) |

True class

# Recap Confusion Matrices

Predicted class

|  | Negative | Positive |  |
|---|---|---|---|
| **Negative** | True negative (TN) | False positive (FP) "false alarm" | N (total number of true negatives) |
| **Positive** | False negative (FN) "miss" | True positive (TP) | P (total number of true positives) |

True class

N* (what we said was negative)    P* (what we said was positive "flagged")

# Recap Confusion Matrices

Predicted class

| | Negative | Positive | |
|---|---|---|---|
| **Negative** | True negative (TN) ✔ | False positive (FP) "false alarm" ✖ | N |
| **Positive** | False negative (FN) "miss" ✖ | True positive (TP) ✔ | P |
| | N* | P* | |

True class

# Recap Confusion Matrices

Predicted class

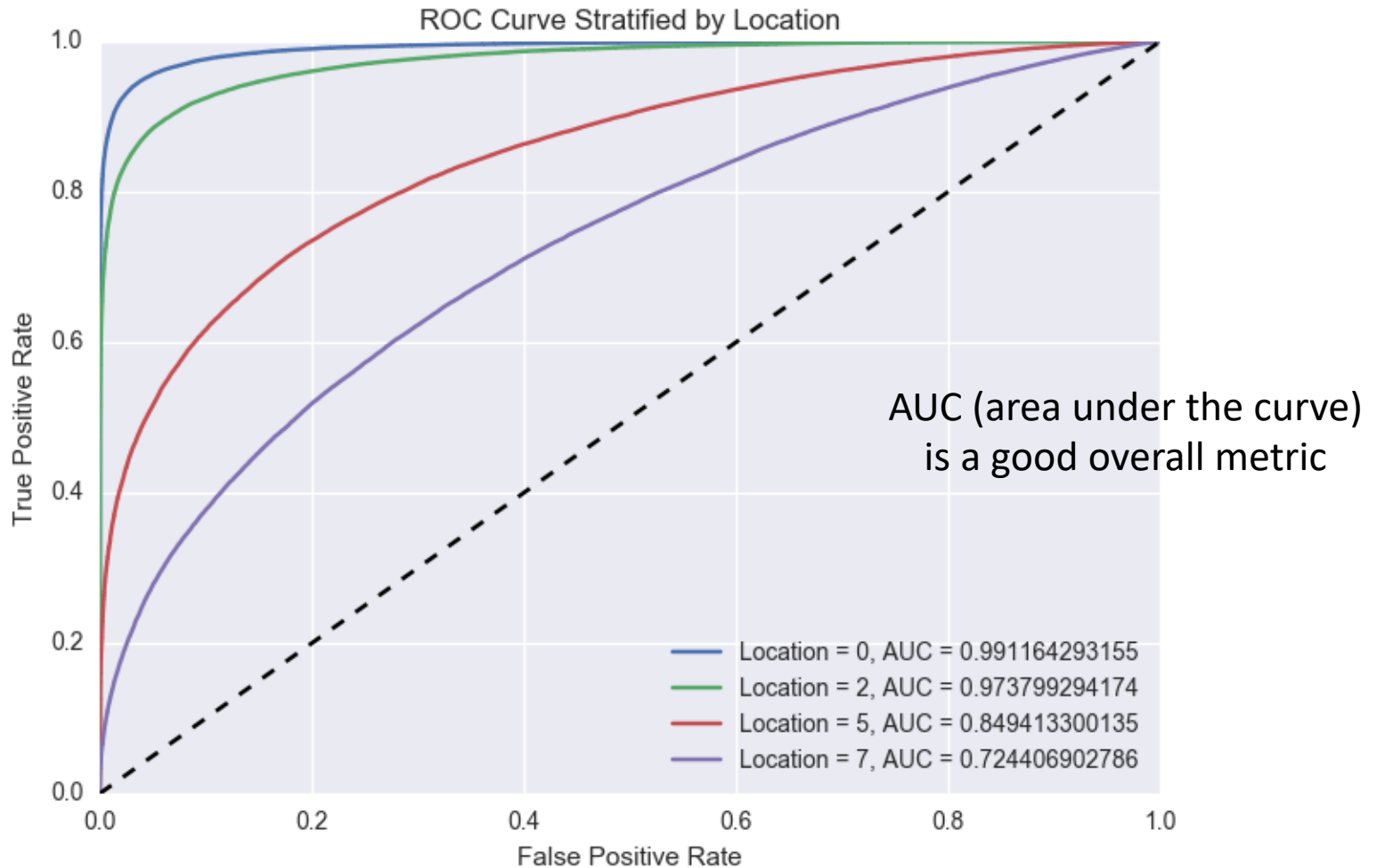|  | Negative | Positive |  |
|---|---|---|---|
| Negative | True negative (TN) | False positive (FP) "false alarm" | N |
| Positive | False negative (FN) "miss" | True positive (TP) | P |
|  | N* | P* |  |

True class

Error:

(FN+FP)/(TN+FP+FN+TP)

= (FN+FP)/(N+P)

# Recap Confusion Matrices

Predicted class

|  | Negative | Positive |
|---|---|---|
| Negative | True negative (TN) | False positive (FP) "false alarm" | N |
| Positive | False negative (FN) "miss" | True positive (TP) | P |
|  | N* | P* |  |

True class

Accuracy = 1-Error:

$$(TN+TP)/(TN+FP+FN+TP)$$

$$= (TN+TP)/(N+P)$$

# Recap Confusion Matrices

Predicted class

|  | Negative | Positive |  |
|---|---|---|---|
| Negative | True negative (TN) | False positive (FP) "false alarm" | N |
| Positive | False negative (FN) "miss" | True positive (TP) | P |
|  | N* | P* |  |

True class

Precision:

$$TP/(FP+TP) = TP/P^*$$

# Recap Confusion Matrices

Predicted class

|  | Negative | Positive |  |
|---|---|---|---|
| Negative | True negative (TN) | False positive (FP) "false alarm" | N |
| Positive | False negative (FN) "miss" | True positive (TP) | P |
|  | N* | P* |  |

True class

Recall
(True Positive Rate):

$$TP/(FN+TP) = TP/P$$

# Recap Confusion Matrices

Predicted class

|  | Negative | Positive |  |
|---|---|---|---|
| Negative | True negative (TN) | False positive (FP) "false alarm" | N |
| Positive | False negative (FN) "miss" | True positive (TP) | P |
|  | N* | P* |  |

True class

False Positive Rate:

FP/(TN+FP) = FP/N

# ROC curve (Receiver Operating Characteristic)

# ROC curve example: comparing methods



Example of a ROC curve from my research
Chan, Perrone, Spence, Jenkins, Mathieson, Song

# How to get a ROC curve for probabilistic methods?

- Usually we use 0.5 as a threshold for binary classification


- Vary the threshold!  (i.e. choose 0.25)


    - $P(y=1 \mid x) > 0.25$       => classify as 1 (positive)
    - $P(y=1 \mid x) <= 0.25$     => classify as 0 (negative)

# Outline for October 22

- Evaluation metrics
  - Confusion matrices revisited
  - ROC curves
  - Relationship to probabilistic methods

- Ensemble methods
  - Bagging
  - Random forests

# Quiz: recap bias and variance



A                    B                    C

Label each picture with variance (high or low) and bias (high or low)

# Quiz: recap bias and variance



Variance:     low
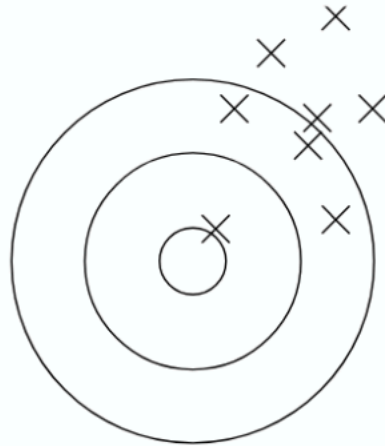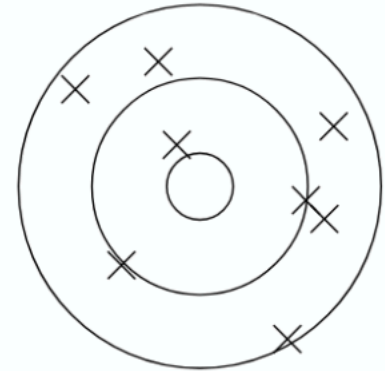Bias:         high

Label each picture with variance (high or low) and bias (high or low)

# Quiz: recap bias and variance
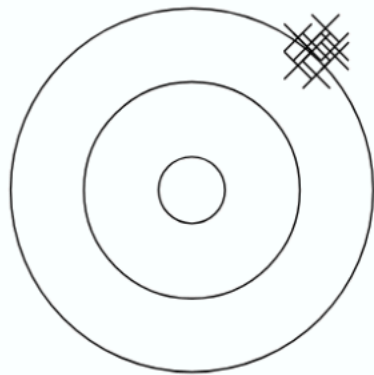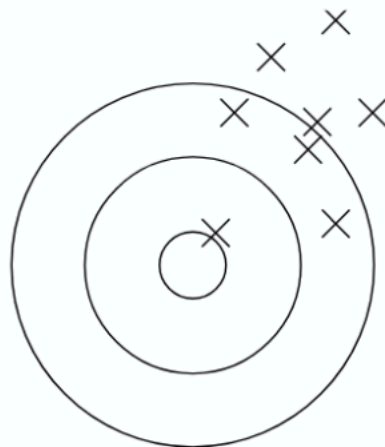


A

B

C

Variance: low
Bias: high

Variance: high
Bias: high

Label each picture with variance (high or low) and bias (high or low)
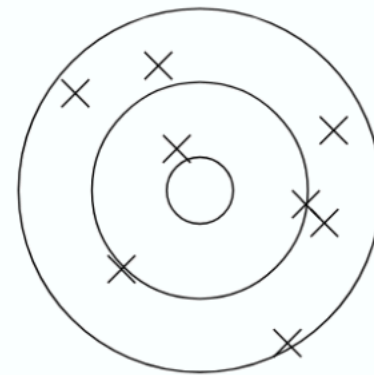
# Quiz: recap bias and variance



| A | B | C |
|---|---|---|
| Variance: low | Variance: high | Variance: high |
| Bias: high | Bias: high | Bias: low |

Label each picture with variance (high or low) and bias (high or low)

# Quiz: recap bias and variance



**A**

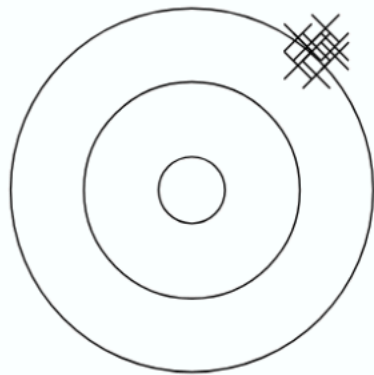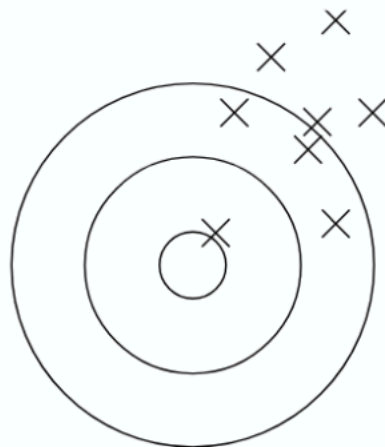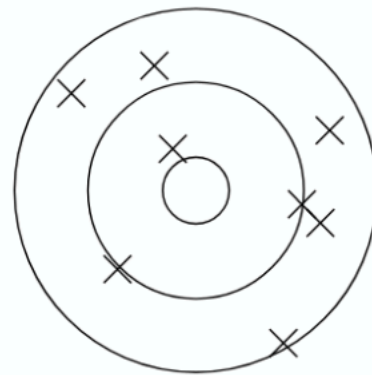Variance: low
Bias: high

**B**

Variance: high
Bias: high

**C**

Variance: high
Bias: low

This is the type of classifier we want to average!

Label each picture with variance (high or low) and bias (high or low)

# Ensemble Idea

- Average the results from several models with high variance and low bias
  - Important that models be diverse (don't want them to be wrong in the same ways)

- If $n$ observations each have variance $s^2$, then the mean of the observations has variance $s^2/n$ (reduce variance by averaging!)

# Learning Theory

Let $H$ be the hypothesis space

Three sources of limitations for traditional classifiers:

❖ Statistical - $H$ is too large relative to size of data

    ❖ Many hypotheses can fit the data by chance

❖ Computational - $H$ is too large to completely search for "best" model

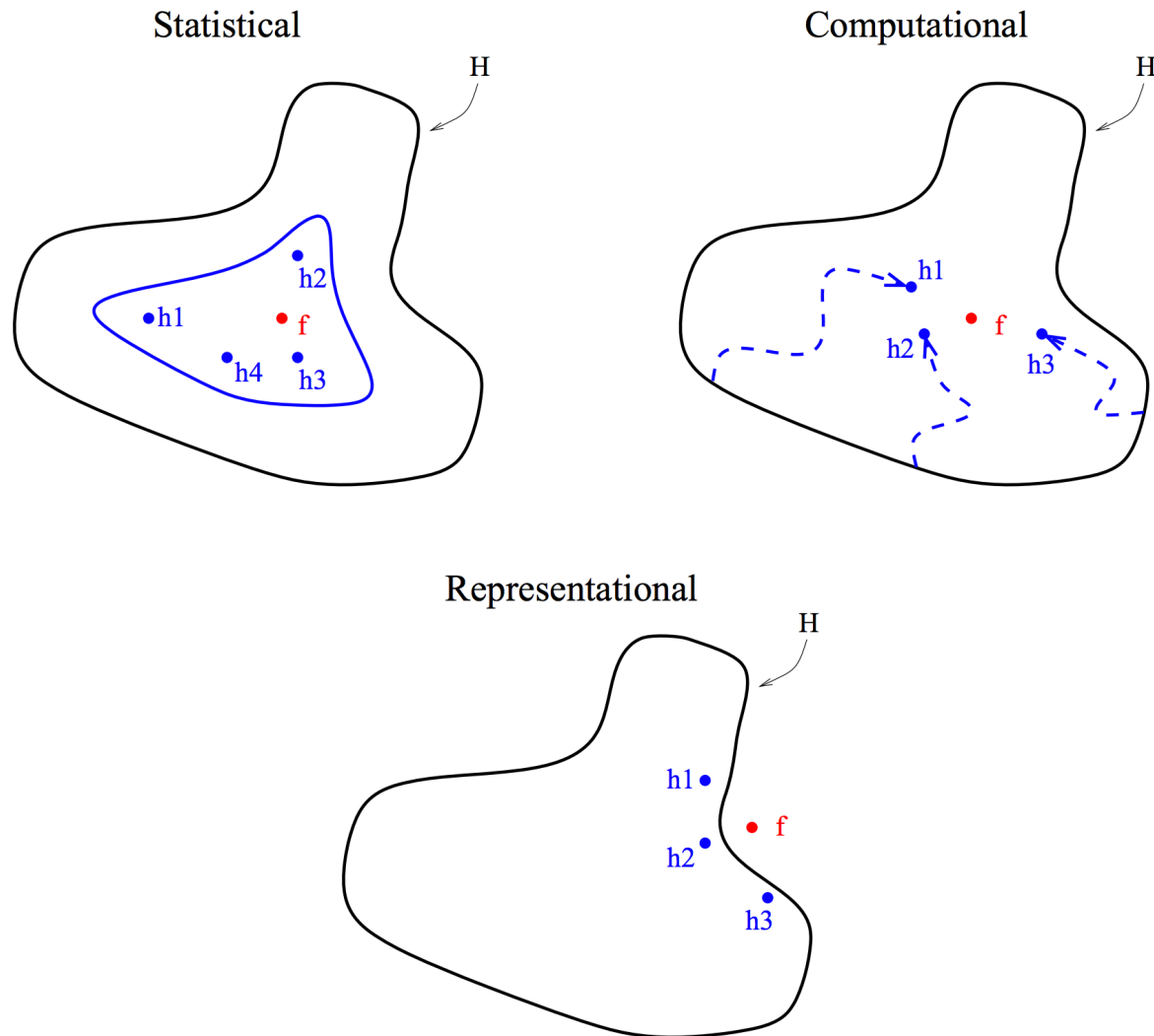❖ Representational - $H$ is not expressive enough

# Learning Theory

❖ <u>Statistical</u>: Average of unstable models (high variance) has more stability

❖ <u>Computational</u>: searching from multiple starting points is better approximation than one starting point

❖ <u>Representational</u>: sum of many models can represent more hypotheses than an individual model

# Learning Theory

❖ <u>Statistical</u>: Average of unstable models (high variance) has more stability

❖ <u>Computational</u>: searching from multiple starting points is better approximation than one starting point

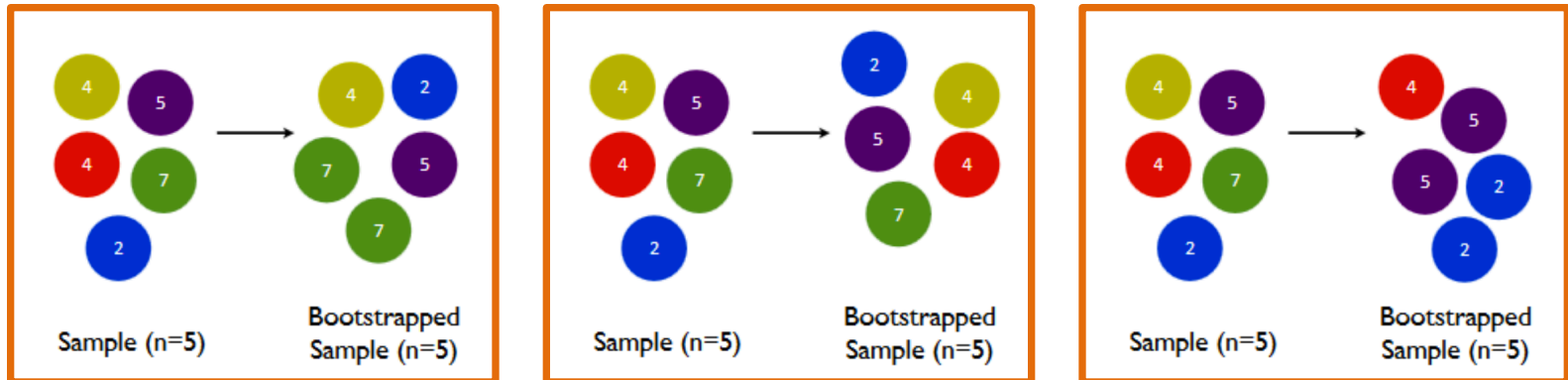❖ <u>Representational</u>: sum of many models can represent more hypotheses than an individual model

Ensembles can address all 3!

# Learning Theory



Figure from Tom Dietterich

# Bagging Algorithm

❖ Bagging = Bootstrap Aggregation [Brieman, 1996]

❖ *Bootstrap* (randomly sample <u>with replacement</u>) original data to create many different training sets

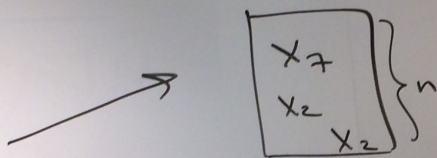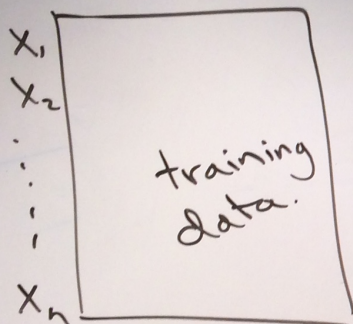❖ Run base learning algorithm on each new data set independently
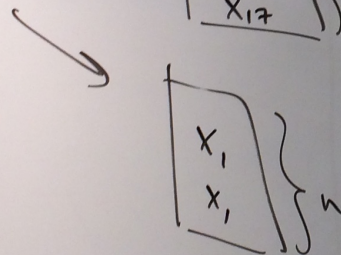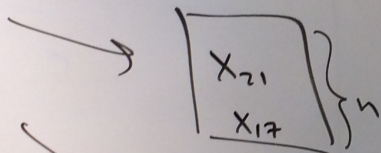


*Desmond Ong, Stanford*

# Bootstrap

Sample <u>with</u> replacement

$X_1$
$X_2$
$\cdot$
$\cdot$
$\cdot$
$X_n$

training data.

$X_7$
$X_2$
$X_2$ $\Big\}$ n

prob we don't choose an example

$X_{21}$
$X_{17}$ $\Big\}$ n

$\left(\dfrac{n-1}{n}\right)^n = \left(1 - \dfrac{1}{n}\right)^n \xrightarrow[\substack{\text{lim} \\ \text{as} \\ n \to \infty}]{} e^{-1} \approx \dfrac{0.38}{\text{large}}$

$X_1$
$X_1$ $\Big\}$ n

$\lim\limits_{n \to \infty} \left(1 + \dfrac{x}{n}\right)^n = e^x$

For Ensembles

$T$ = #models/classifiers (index $t$)

$x$ = test example (could be vector)

$x^{(t)}$ = bootstrap training set $t$

$h^{(t)}(x)$ = hypothesis about $x$
from model $t$

$r$ = prob. of error of individual
model

$R$ = # votes for wrong class.

# Bagging (Bootstrap Aggregation) , $y \in \{0,1\}$

## Train   for $t$ in range($T$):

- create bootstrap dataset $X^{(t)}_{(n \times p)}$

- train on $X^{(t)}$ to get model $h^{(t)}$

also in $\{0,1\}$

threshold already applied

## Test   for $x$ in test data:

$$h(x) = \arg\max_{y \in \{0,1\}} \sum_{t=1}^{T} \mathbb{1}(h^{(t)}(x) = y)$$