

CS 66: Machine Learning

Prof. Sara Mathieson

Spring 2019



Outline for April 22

- Lab 5 Analysis Questions
 - Midterm 2 Review
 - Practice Problems
 - Questions for Wednesday
-
- Final Project Proposals: all feedback and repos finished
 - Wednesday in class “office hours” (**submit a question today!**)
 - Fri: Guest lecture by Prof. Matt Zucker
 - **Office hours today: 12:30-2pm & 3-4pm**
 - I will not be on campus tomorrow, so make sure to come to office hours today (**and also post on Piazza!**)

Outline for April 22

- Lab 5 Analysis Questions
- Midterm 2 Review
- Practice Problems
- Questions for Wednesday

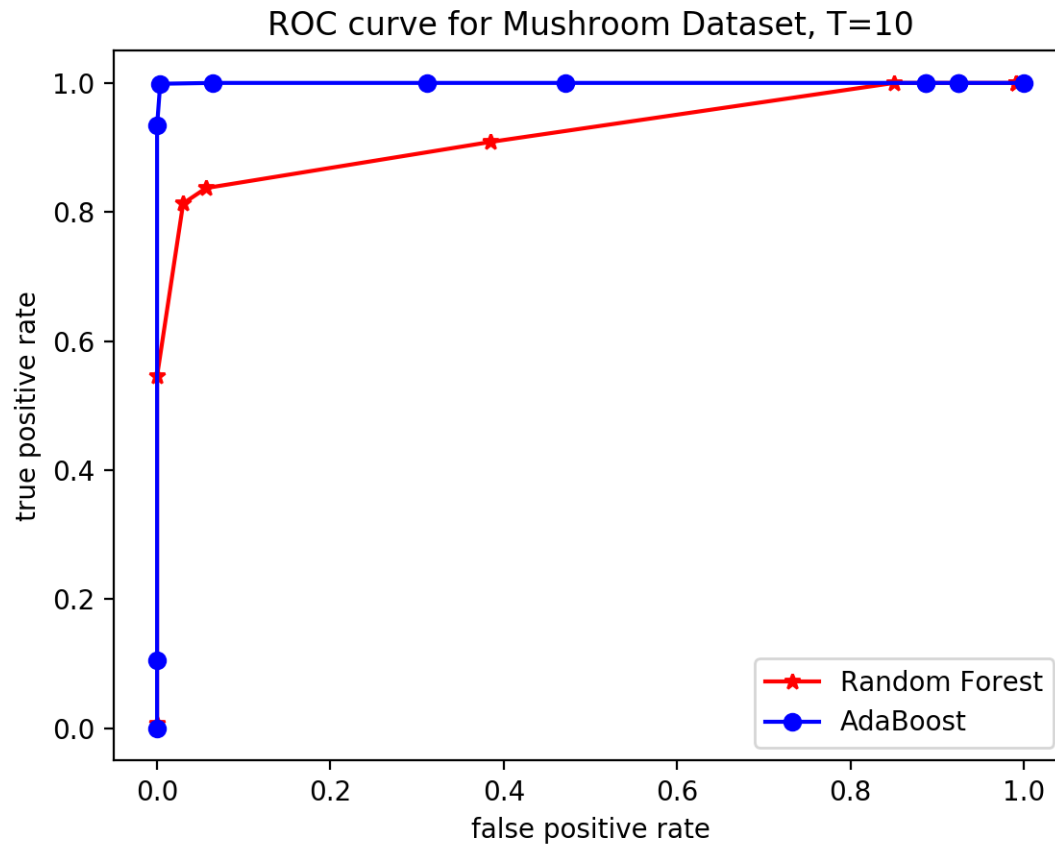
Lab 5: Runtimes

- n =number of examples, p =number of features, T =number of classifiers
- Random Forests: $O(\sqrt{p}nT)$
- AdaBoost: $O(pnT)$
- Random forests is better even given this analysis, but it is also very parallelizable! AdaBoost is not

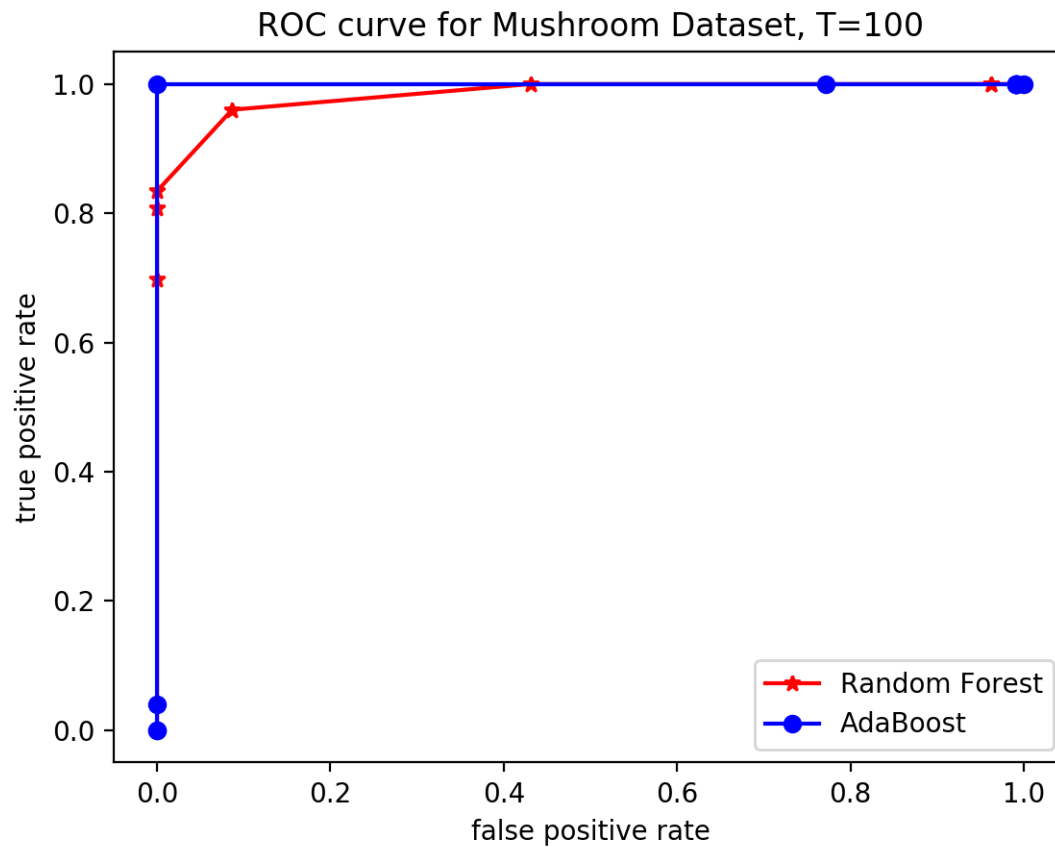
Ensemble Methods: overfitting

- Ensemble methods are very robust to overfitting!
- If all classifiers in the ensemble are “weak”, then nothing about the overall model is fit to the “noise” in the data
- Very powerful idea, and one reason why ensemble methods are still widely used

Lab 5: ROC curve (T=10)



Lab 5: ROC curve (T=100)



Outline for April 22

- Lab 5 Analysis Questions
- **Midterm 2 Review**
- Practice Problems
- Questions for Wednesday

Naïve Bayes

* pickup notecard
& handout

* write Question
on card for
Wednesday

* can use double
sided cheat
Sheet for
exam!

Naive Bayes

$$\underbrace{P(y=k|\vec{x})}_{\text{posterior}} =$$

$$\frac{\overbrace{P(y=k)}^{\text{prior}} \overbrace{P(\vec{x}|y=k)}^{\text{likelihood}}}{\underbrace{P(\vec{x})}_{\text{evidence}}}$$

$$k = 1, 2, \dots, K$$

↑
of
classes

$$P(A, B | C)$$
$$= P(A | B, C) P(B | C)$$

Naive Bayes

$$P(y=k|\vec{x}) = \frac{\overbrace{P(y=k)}^{\text{prior}} \overbrace{P(\vec{x}|y=k)}^{\text{likelihood}}}{\underbrace{P(\vec{x})}_{\text{evidence}}}$$

posterior

$k = 1, 2, \dots, K$
 \uparrow
 # of classes

$$P(A, B | C)$$

$$= P(A | B, C) P(B | C)$$

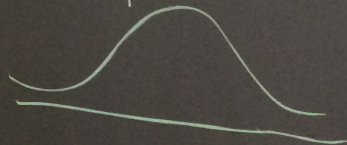
$$P(\vec{x} | y=k) = P(\underbrace{x_1}_A, \underbrace{x_2}_{B}, \dots, \underbrace{x_p}_C | y=k)$$

$$= P(x_1 | x_2, \dots, x_p, y=k) P(x_2 \dots x_p | y=k)$$

$$\approx \prod_{j=1}^p P(x_j | y=k)$$

if cont. fit a Gaussian

NB assumption.



Naïve Bayes Assumption

- Feature j is independent of all other features, given the class label

$$\begin{array}{c|c} x_1 \dots x_5 \dots x_p & y \\ \hline C \dots A \dots B & 4 \end{array}$$

estimat

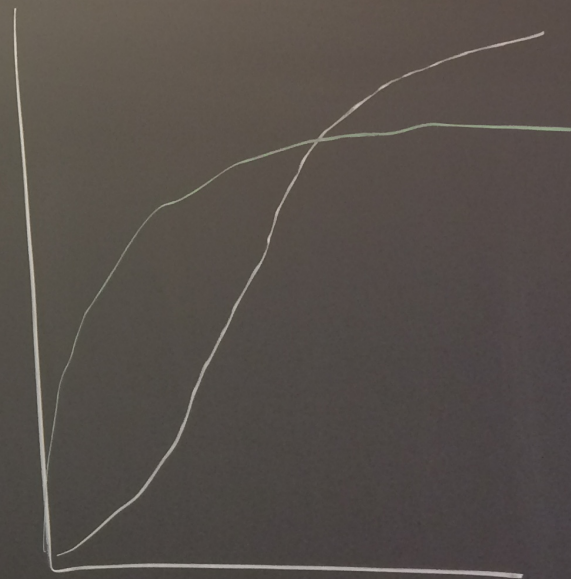
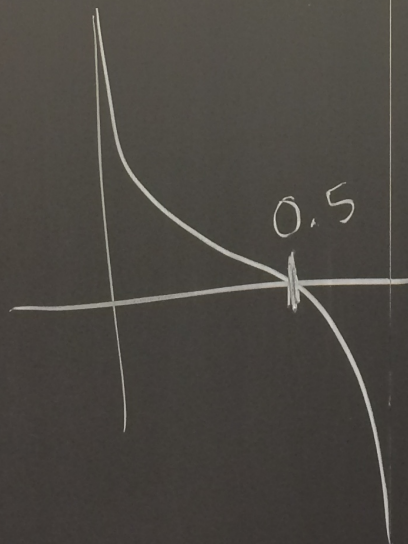
train

no

estimate $p(x_j | y=k) = \frac{p(x_j, y=k)}{p(y=k)}$

train $x_p \quad y$
no $\textcircled{B \quad 4}$

$$\approx \frac{N_{j,v,k} + 1}{N_k + |f_j|}$$



Evaluation Metrics

Recap Precision and Recall

- Precision: of all the “flagged” examples, which ones are actually relevant (i.e. positive)?

(Purity)

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

(Completeness)

Recap Confusion Matrices

Predicted class

Negative

Positive

Negative

True negative
(TN)

False positive
(FP)
“false alarm”

N (total number of true negatives)

True
class

Positive

False negative
(FN)
“miss”

True positive
(TP)

P (total number of true positives)

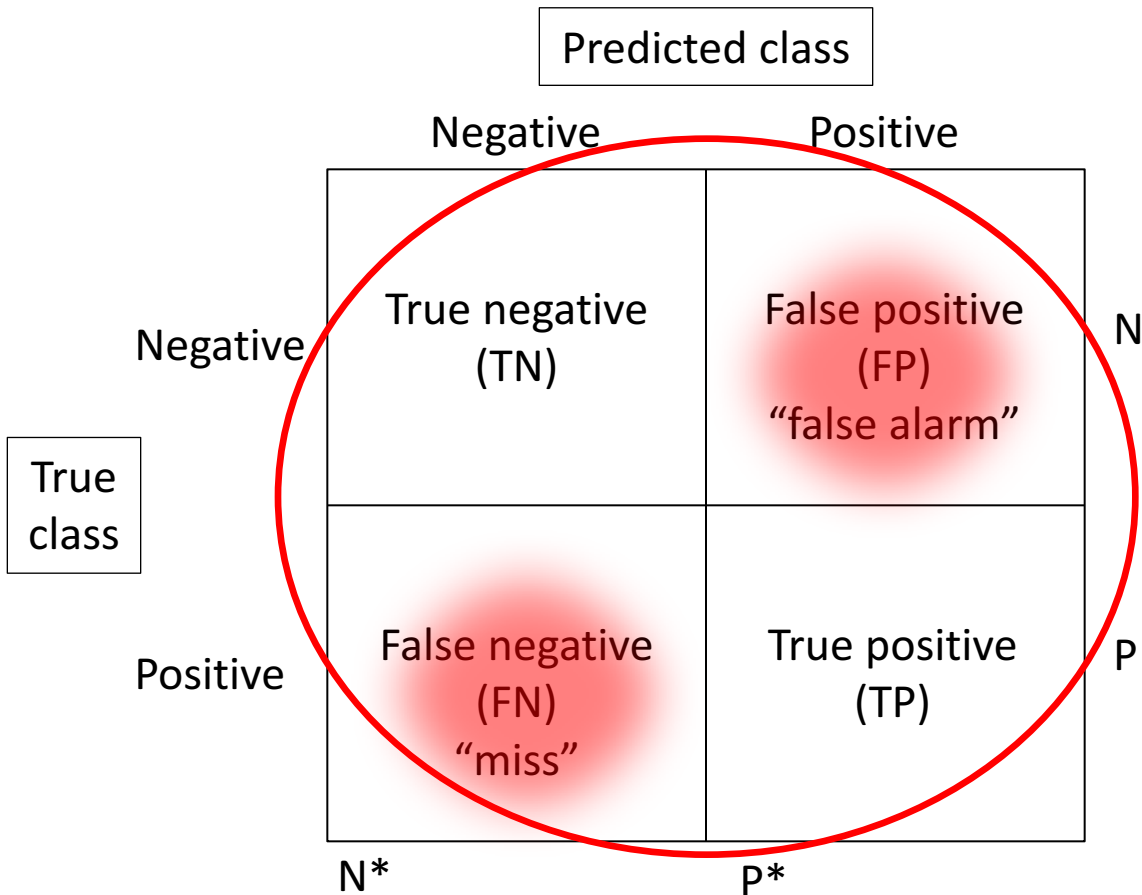
N* (what we said
was negative)

P* (what we said was
positive “flagged”)

Recap Confusion Matrices

		Predicted class		
		Negative	Positive	
True class	Negative	True negative (TN) ✓	False positive (FP) "false alarm" ✗	N
	Positive	False negative (FN) "miss" ✗	True positive (TP) ✓	P
		N*	p*	

Recap Confusion Matrices

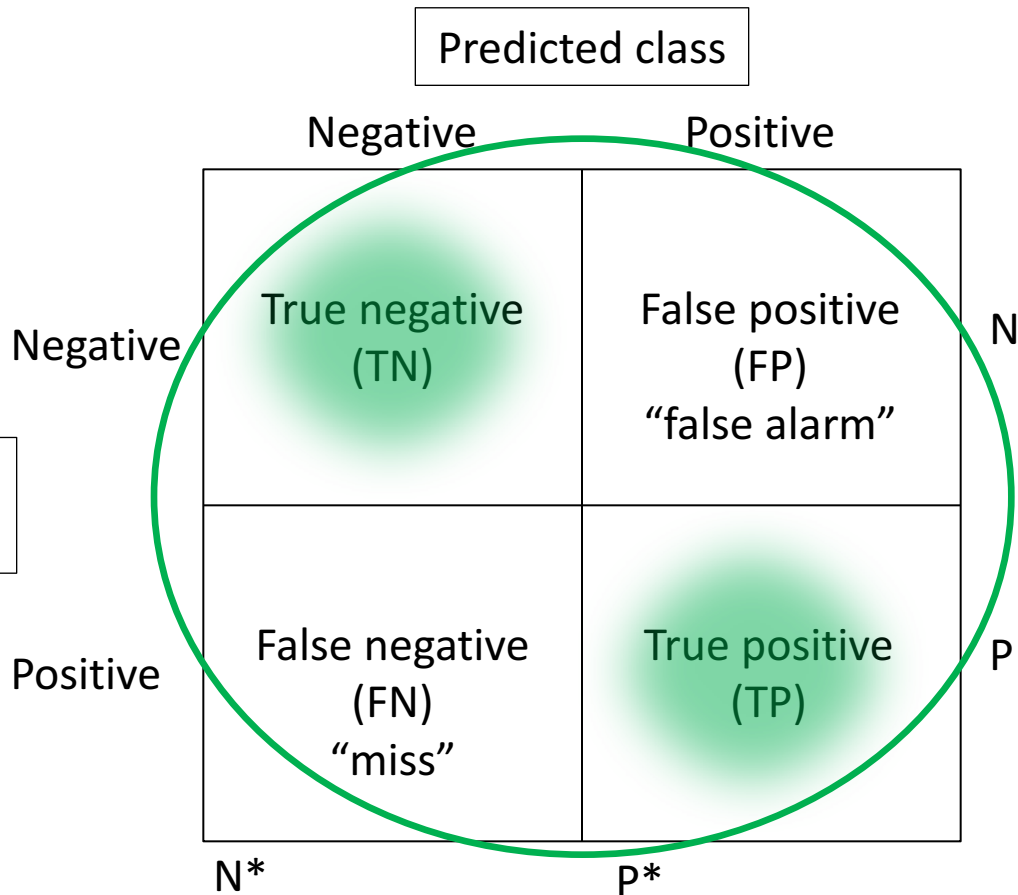


Error:

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

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

Recap Confusion Matrices



Accuracy = 1-Error:

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

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

Recap Confusion Matrices

		Predicted class		
		Negative	Positive	
True class	Negative	True negative (TN)	False positive (FP) "false alarm"	N
	Positive	False negative (FN) "miss"	True positive (TP)	P
		N*	p*	

Precision:

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

Recap Confusion Matrices

		Predicted class		
		Negative	Positive	
True class	Negative	True negative (TN)	False positive (FP) “false alarm”	N
	Positive	False negative (FN) “miss”	True positive (TP)	P
		N*	p*	

Recall
(True Positive Rate):

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

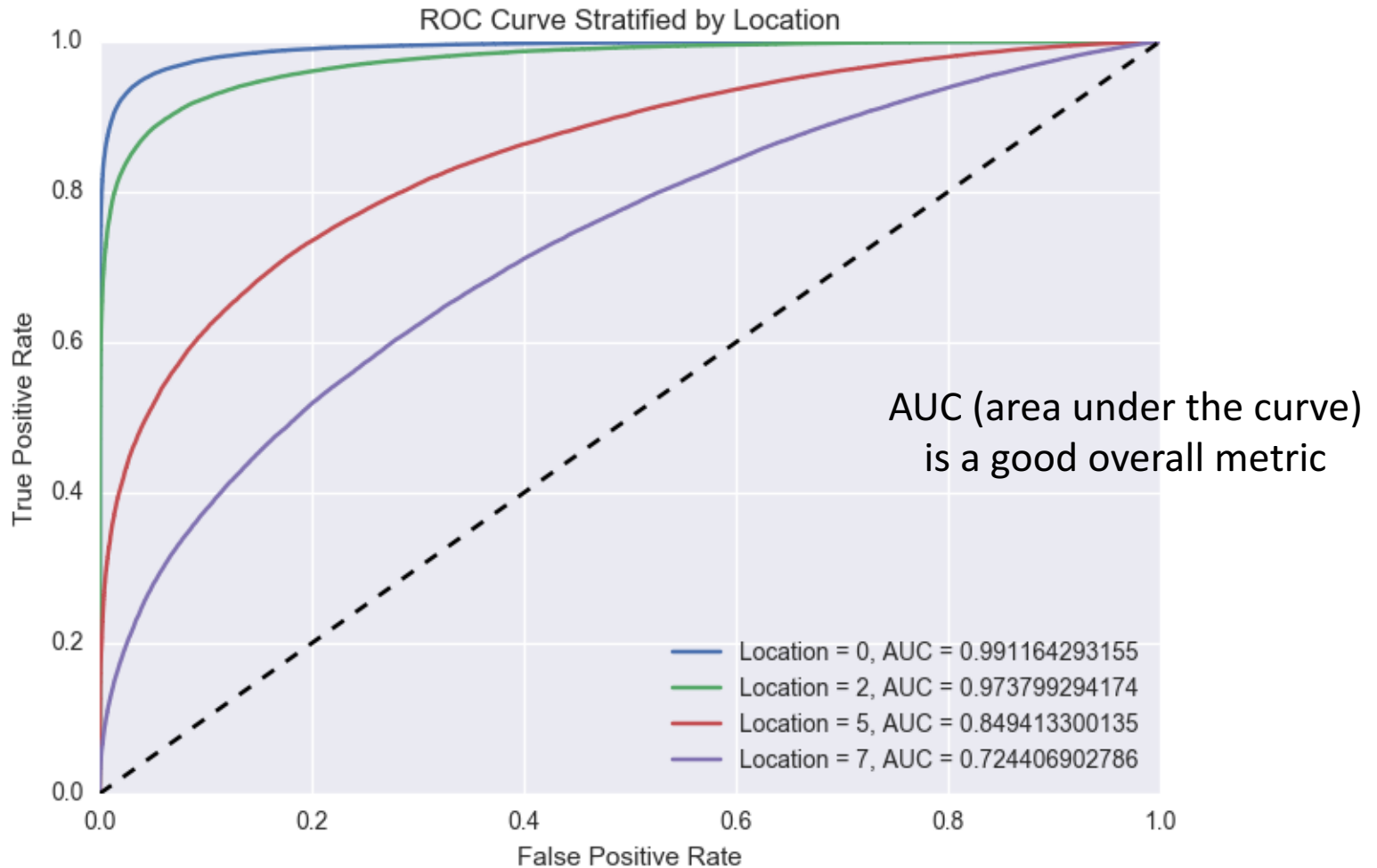
Recap Confusion Matrices

		Predicted class		
		Negative	Positive	
True class	Negative	True negative (TN)	False positive (FP) "false alarm"	N
	Positive	False negative (FN) "miss"	True positive (TP)	P
		N*	p*	

False Positive Rate:

$$FP/(TN+FP) = FP/N$$

ROC curve example: comparing methods



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

Cross Validation

- Allows us to choose best hyper-parameters
- Allows us to return multiple independent accuracy results
- We can use this distribution of accuracy numbers in statistical frameworks (find mean/variance, compare with other methods, etc)

Ensemble Methods

Learning Theory

Let H be the hypothesis space

Three sources of limitations for traditional classifiers:

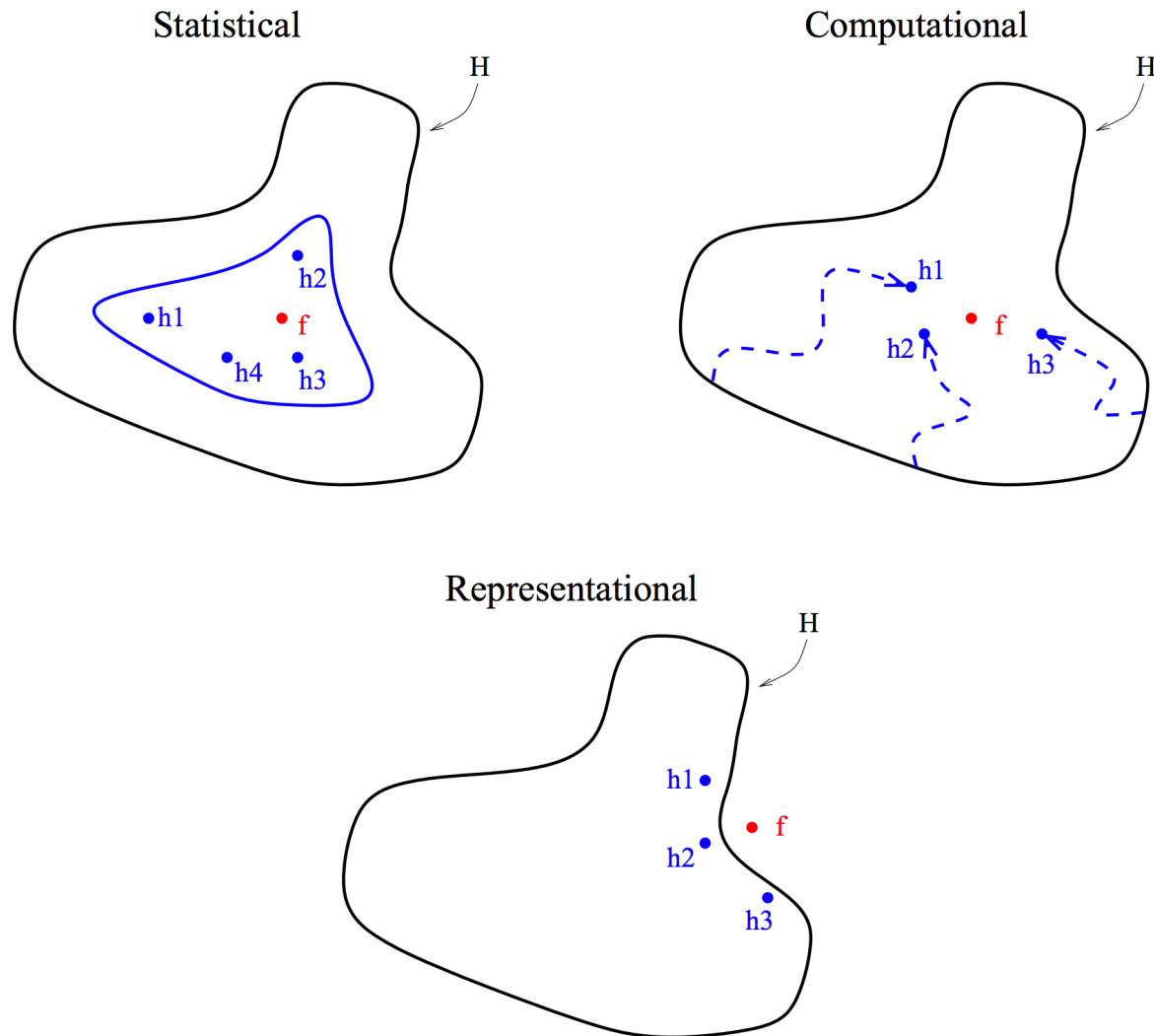
- ❖ Statistical - H is too large relative to size of data
 - ❖ Many hypotheses can fit the data by chance
- ❖ Computational - H is too large to completely search for “best” model
- ❖ Representational - H is not expressive enough

Learning Theory

- ❖ Statistical: Average of unstable models (high variance) has more stability
- ❖ Computational: searching from multiple starting points is better approximation than one starting point
- ❖ Representational: sum of many models can represent more hypotheses than an individual model

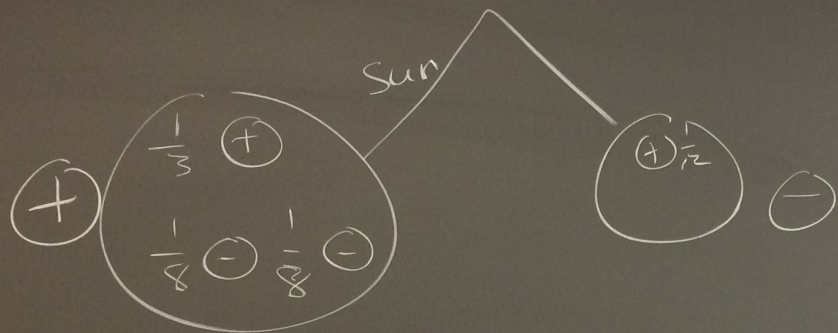
Ensembles can address all 3!

Learning Theory



Outline for April 22

- Lab 5 Analysis Questions
- Midterm 2 Review
- **Practice Problems**
- Questions for Wednesday



$$\epsilon_t = \frac{1}{8} + \frac{1}{8} + \frac{1}{12} = \frac{1}{3} \quad \star$$

$$P(+)=\frac{\frac{1}{3}}{\frac{1}{3}+\frac{1}{8}+\frac{1}{8}}=\frac{12}{21}=\frac{4}{7}$$

$$\frac{4}{7} \geq 0.5$$

=k)

r. p. 7
Gaussian

Handout 19, Question 1, parts (a) and (b)

Handout 19, Question 2

$$r = \frac{1}{3}, T = 5 \quad T = 4$$

$R = \#$ votes for wrong class

$$R = 0, 1, 2, \underbrace{(3, 4, 5)}_{\text{pt classified incorrectly overall}}$$

$$P(R=k)$$

$$= \binom{T}{k} \underbrace{r^k}_{\text{wrong}} \underbrace{(1-r)^{T-k}}_{\text{right}}$$

$$P(R=5) = \binom{5}{5} \left(\frac{1}{3}\right)^5 \quad 0! = 1$$

$$P(R=4) = \binom{5}{4} \left(\frac{1}{3}\right)^4 \left(\frac{2}{3}\right) \quad \left(\binom{5}{0} = 1\right)$$

$$P(R=3) = \binom{5}{3} \left(\frac{1}{3}\right)^3 \left(\frac{2}{3}\right)^2$$

$$\sum_{k=0}^T P(R=k) \approx \underbrace{0.21}$$

$$k = \left\lceil \frac{T+1}{2} \right\rceil$$

$$\frac{4+1}{2} = 2.5$$

Overall prob of being wrong.
→ 3

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$