

CS 260: Foundations of Data Science

Prof. Sara Mathieson

Fall 2021



HVERFORD
COLLEGE

Admin

- **Lab 1** grades posted on Moodle
 - Note we did not take off for style, but will in the future
- **Lab 3** due TONIGHT
- Office hours TODAY **3:30-5pm in H204**
- TA hours: **6:30-8:30pm in H110**
- **Lab 4** posted today
- **Note-taker:** Julia

Outline for September 21

- Introduction to classification
 - Decision tree models
 - Probabilistic interpretation
- Evaluation Metrics
 - Confusion matrices
 - ROC curves
 - Precision and recall

Outline for September 21

- Introduction to classification
 - Decision tree models
 - Probabilistic interpretation
- Evaluation Metrics
 - Confusion matrices
 - ROC curves
 - Precision and recall

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*

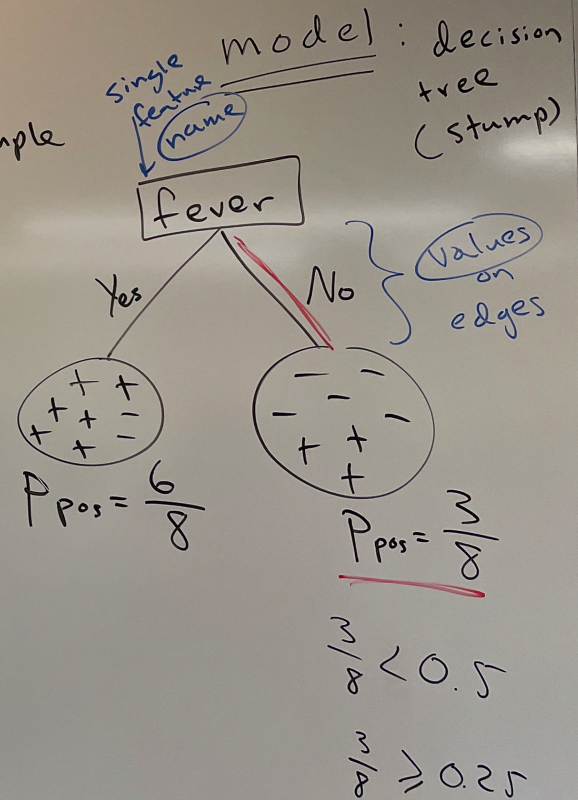
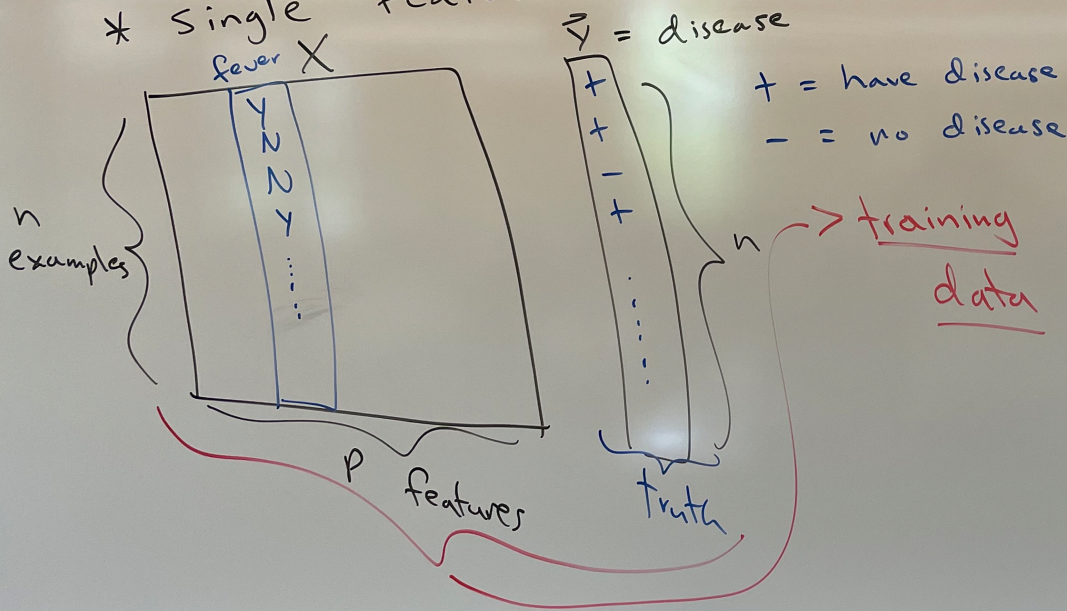
Goals of Evaluation

- Think about what metrics are important for the problem at hand
- Compare different methods or models on the same problem
- Common set of tools that other researchers/users can understand

PLEASE LEAVE COMPUTERS ON

Introduction to Classification

* single feature model \rightarrow medical example



classify a test example

$$\vec{X}_{\text{test}} = \left[\dots \overset{\text{fever}}{N} \dots \right]^T$$

typically: threshold = 0.5 \Rightarrow

threshold = 0.25 \Rightarrow

$$\hat{Y}_{\text{test}} = \text{negative}$$

$$\hat{Y}_{\text{test}} = \text{positive}$$

$P_{\text{pos}} \geq \text{threshold} \Rightarrow \text{positive}$
else: $\Rightarrow \text{negative}$

5

Handout 7

model 1

Slope

up

flat

down

$P_{pos} = \frac{0}{4}$

(-)

(-)

$\frac{5}{7}$

(+)

(+)

$\frac{3}{5}$

(+)

(+)

Which is better?

model 2

thal

fixed

normal

reverse

$\frac{1}{3}$

(-)

(+)

$\frac{3}{8}$

(-)

(+)

$\frac{4}{5}$

(+)

(+)

threshold = 0.5

threshold = $\frac{1}{12}$

Outline for September 21

- Introduction to classification
 - Decision tree models
 - Probabilistic interpretation
- **Evaluation Metrics**
 - Confusion matrices
 - ROC curves
 - Precision and recall

PLEASE LEAVE COMPUTERS ON

Evaluation Metrics

①

Confusion matrix

threshold = 0.5

	-	+	
-	65	15	N=80
+	7	13	P=20

FP (False Positive) points to 15

TP (True Positive) points to 13

low (red arrow) points to 7

want high (blue arrow) points to 13

②

$$\text{accuracy} = \frac{\# \text{ correct}}{\text{total}} = \frac{65 + 13}{100} = \boxed{0.78}$$

test data: m = 100 examples, 80 negatives, 20 positives

threshold = 0.25

	-	+
-	50	30
+	1	19

$$\text{accuracy} = \frac{50 + 19}{100} = \boxed{0.69}$$

threshold = 0.75

	-	+
-	76	4
+	11	9

false negatives "miss" (red arrow) points to 11

$$\text{accuracy} = \frac{76 + 9}{100} = \boxed{0.85}$$

③ ROC curve

$$\text{FPR} = \frac{\text{FP}}{N} \leftarrow \begin{array}{l} \text{false} \\ \text{positives} \\ \text{all} \\ \text{negatives} \end{array}$$

$$\text{TPR} = \frac{\text{TP}}{P} \leftarrow \begin{array}{l} \text{true} \\ \text{positives} \\ \text{all} \\ \text{positives} \end{array}$$

$$\text{thresh} = 0.5 \Rightarrow \text{FPR} = \frac{15}{80} \approx 19\%$$

$$\text{TPR} = \frac{13}{20} \approx 65\%$$

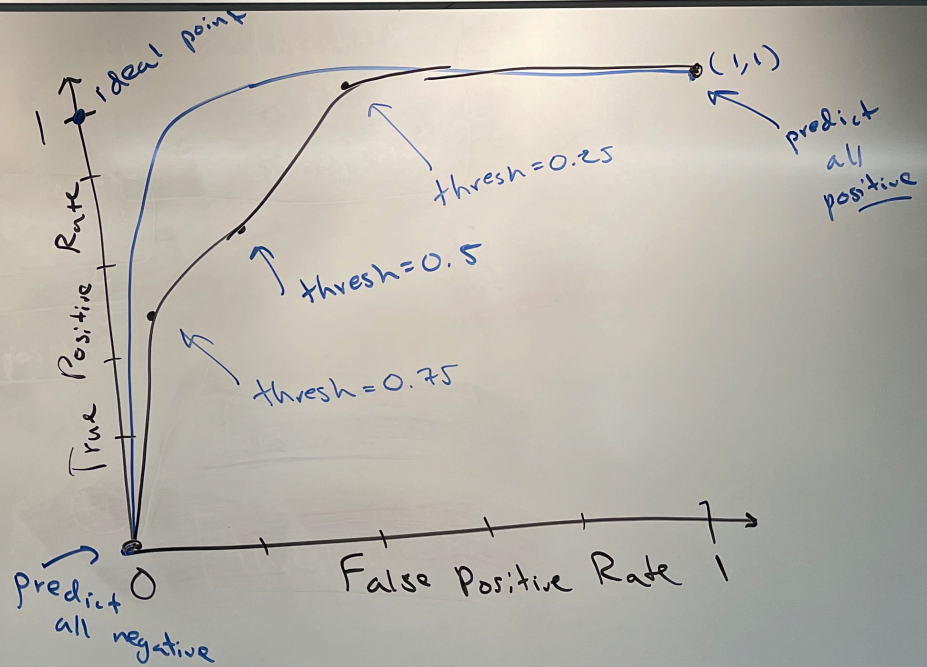
$$\text{thresh} = 0.25 \Rightarrow \text{FPR} = \frac{30}{80} \approx 38\% -$$

$$\text{TPR} = \frac{19}{20} \approx 95\% \star$$

$$\text{thresh} = 0.75$$

$$\text{FPR} = \frac{4}{80} = 5\% \star$$

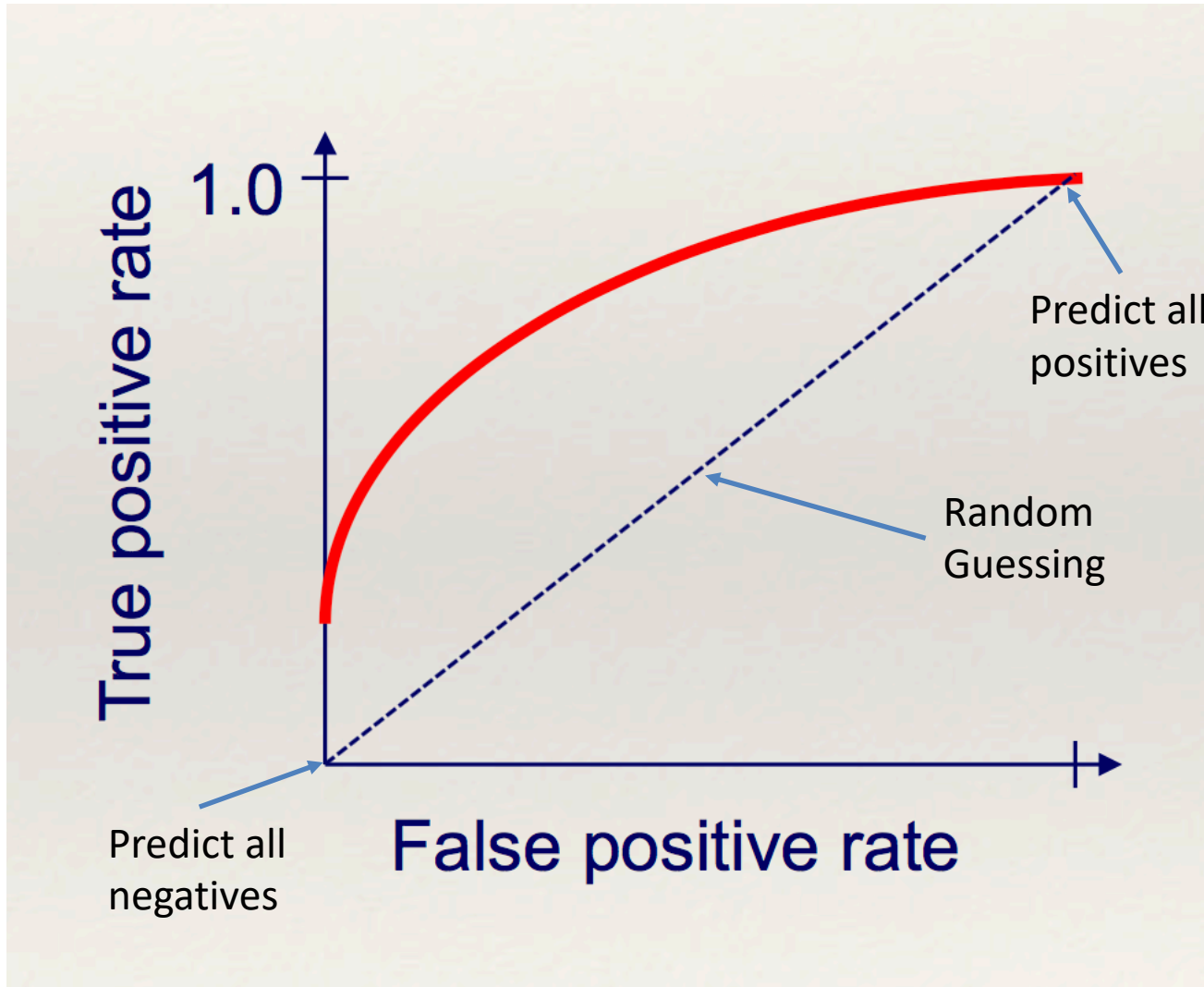
$$\text{TPR} = \frac{9}{20} = 45\% -$$



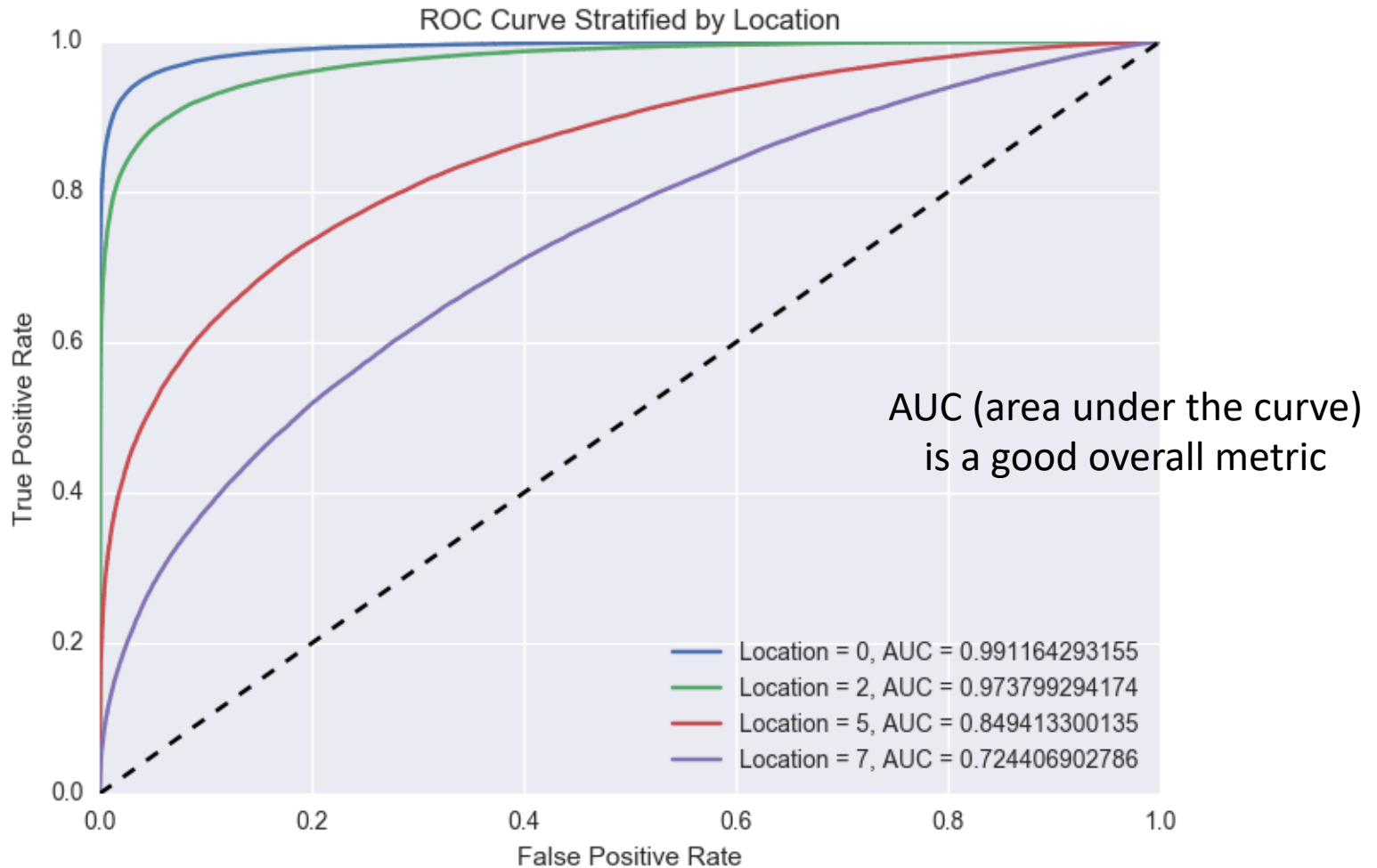
ROC curve (Receiver Operating Characteristic)

More history here!

https://en.wikipedia.org/wiki/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)
 - $P(y=1 \mid x) \geq 0.25$ \Rightarrow classify as 1 (positive)
 - $P(y=1 \mid x) < 0.25$ \Rightarrow classify as 0 (negative)