



CMPT 413/713: Natural Language Processing

Classification - Evaluation

Spring 2024
2024-01-17

Adapted from slides from Danqi Chen, Karthik Narasimhan

Evaluation

Evaluation

- Consider binary classification
- Table of predictions

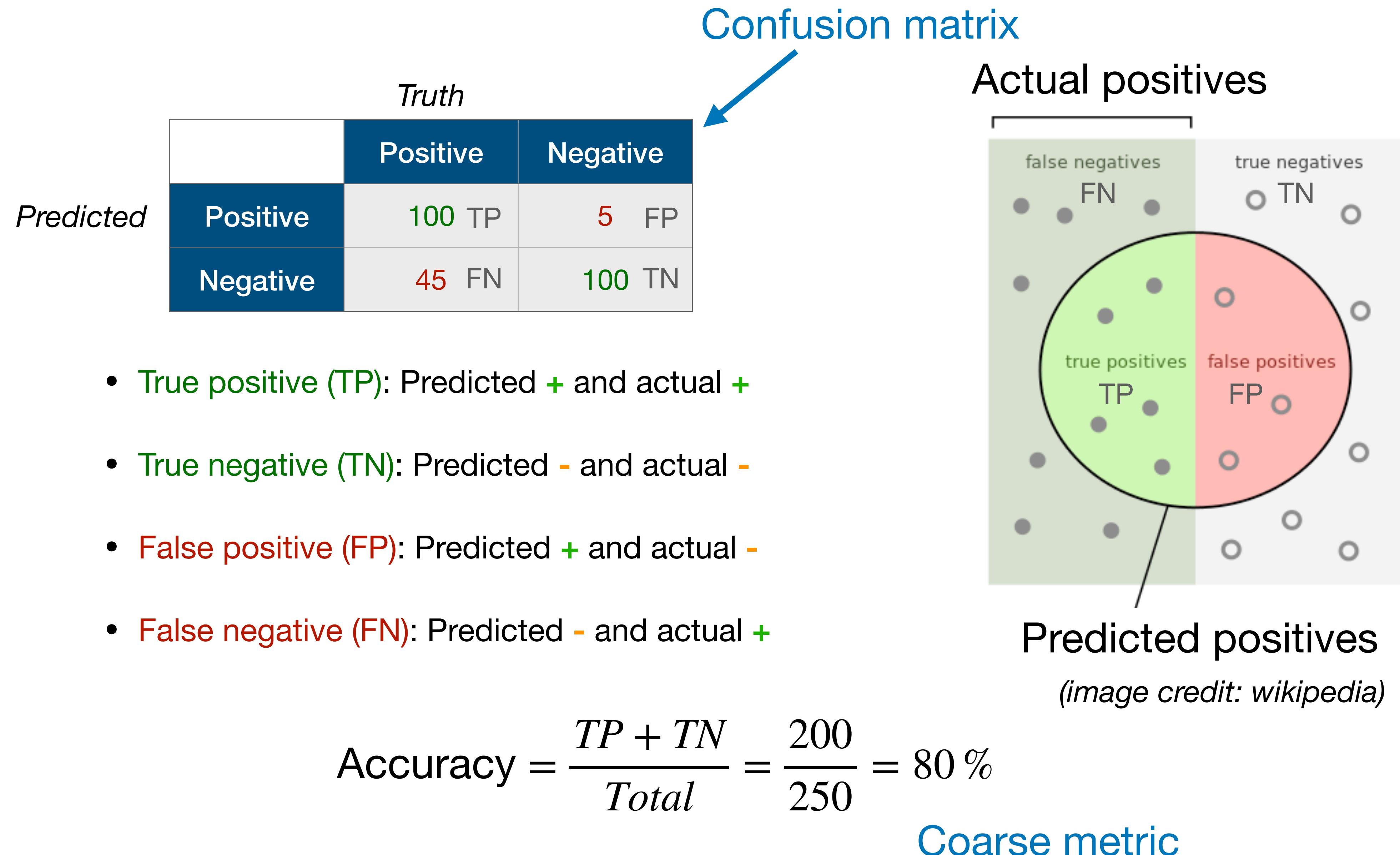
		<i>Truth</i>	
		Positive	Negative
<i>Predicted</i>	Positive	100	5
	Negative	45	100

Confusion matrix

- Ideally, we want:

		<i>Truth</i>	
		Positive	Negative
<i>Predicted</i>	Positive	145	0
	Negative	0	105

Evaluation Metrics



Evaluation Metrics

		Truth	
		Positive	Negative
Predicted	Positive	100	5
	Negative	45	100
		Positive	Negative
	Positive	50	25
	Negative	25	150

- True positive (TP): Predicted + and actual +
- True negative (TN): Predicted - and actual -
- False positive (FP): Predicted + and actual -
- False negative (FN): Predicted - and actual +

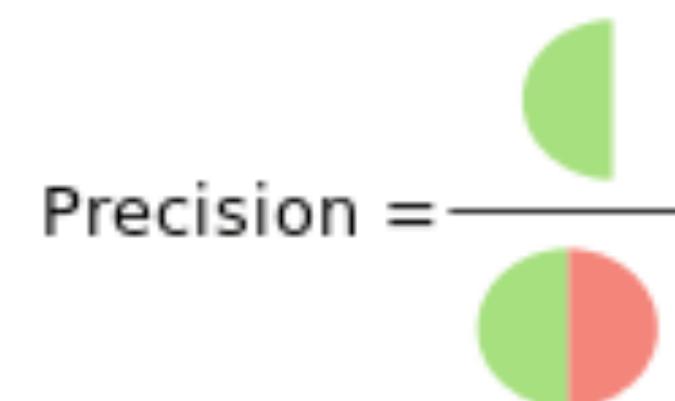
Accuracy cannot distinguish
between the two models!

$$\text{Accuracy} = \frac{TP + TN}{Total} = \frac{200}{250} = 80\%$$

Precision and Recall

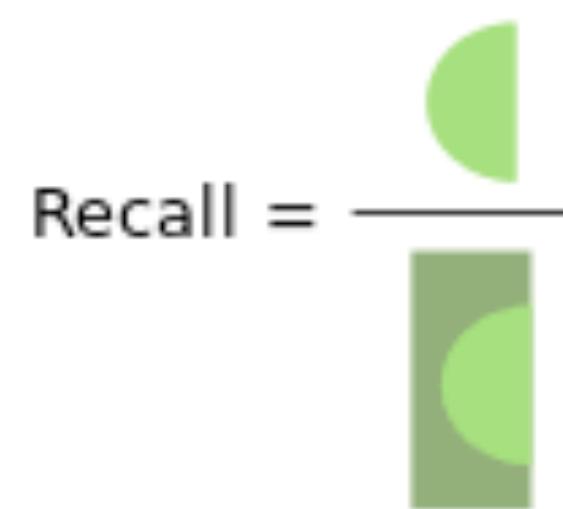
- Precision: % of selected classes that are correct

$$\text{Precision}(+) = \frac{TP}{TP + FP}$$

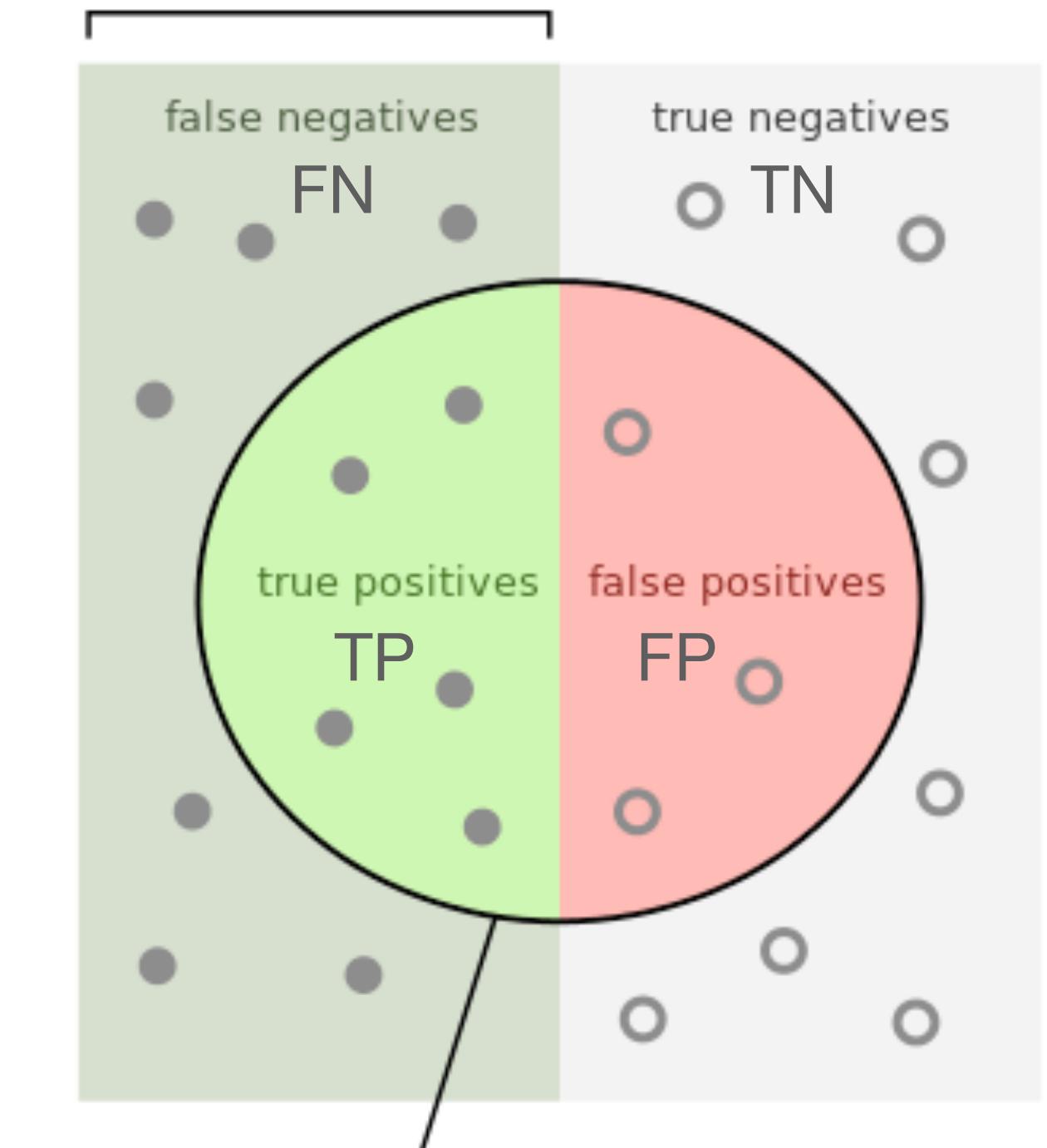


- Recall: % of correct items selected

$$\text{Recall}(+) = \frac{TP}{TP + FN}$$



Actual positives (relevant)



Predicted positives
(selected/retrieved)

(image credit: wikipedia)

Evaluation Metrics

		Truth	
		Positive	Negative
Predicted	Positive	100	5
	Negative	45	100
		Positive	Negative
	Positive	50	25
	Negative	25	150

$$\text{Precision}(+) = \frac{TP}{TP + FP}$$

$$\frac{100}{100 + 5} = 0.95$$

$$\frac{50}{50 + 25} = 0.75$$

$$\text{Recall}(+) = \frac{TP}{TP + FN}$$

$$\frac{100}{100 + 45} = 0.69$$

$$\frac{50}{50 + 25} = 0.75$$

Two metrics - which one to use when comparing the two models?

F-Score

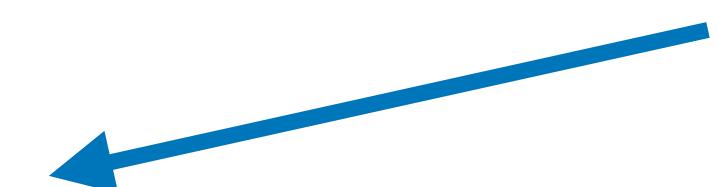
- Combined measure
- Harmonic mean of Precision and Recall

$$F_1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

- Or more generally,

$$F_\beta = \frac{(1 + \beta^2) \cdot \text{Precision} \cdot \text{Recall}}{\beta^2 \cdot \text{Precision} + \text{Recall}}$$

Use β to control importance of
Precision vs Recall



Evaluation Metrics

		Truth	
		Positive	Negative
Predicted	Positive	100	5
	Negative	45	100
		Positive	Negative
	Positive	50	25
	Negative	25	150

$$\text{Precision}(+) = \frac{TP}{TP + FP}$$

$$\frac{100}{100 + 5} = 0.95$$

$$\frac{50}{50 + 25} = 0.75$$

$$\text{Recall}(+) = \frac{TP}{TP + FN}$$

$$\frac{100}{100 + 45} = 0.69$$

$$\frac{50}{50 + 25} = 0.75$$

$$F_1(+) = \frac{2 \cdot P(+)R(+) }{P(+) + R(+)}$$

$$0.8$$

$$0.75$$

Evaluation Metrics

		<i>Truth</i>	
		Positive	Negative
<i>Predicted</i>	Positive	TP	FP
	Negative	FN	TN

Use a simple rule, can you
design a classifier with

$$\text{Precision}(+) = \frac{TP}{TP + FP}$$

Q. perfect precision?

$$\text{Recall}(+) = \frac{TP}{TP + FN}$$

Q. perfect recall?

Aggregating scores

- How to handle more than 2 classes?
- We have Precision, Recall, F1 for each class

		gold labels		
		urgent	normal	spam
system output	urgent	8	10	1
	normal	5	60	50
	spam	3	30	200

precision_u= $\frac{8}{8+10+1}$

precision_n= $\frac{60}{5+60+50}$

precision_s= $\frac{200}{3+30+200}$

recall_u= $\frac{8}{8+5+3}$

recall_n= $\frac{60}{10+60+30}$

recall_s= $\frac{200}{1+50+200}$

(Credits: Dan Jurafsky)

Aggregating scores

- How to handle more than 2 classes?
- We have Precision, Recall, F1 for each class
- How to combine them for an overall score?
 - **Macro-average:** Compute for each class, then average
 - **Micro-average:** Collect predictions for all classes and jointly evaluate

Macro vs Micro average

- Micro-averaged score is dominated by score on **common classes**

Class 1: Urgent		
	true true	
urgent	true urgent	not
system		
urgent	8	11
not	8	340

Class 2: Normal		
	true true	
normal	true normal	not
system		
normal	60	55
not	40	212

Class 3: Spam		
	true true	
spam	true spam	not
system		
spam	200	33
not	51	83

Pooled		
	true true	
yes	true yes	no
system		
yes	268	99
no	99	635

$$\text{precision} = \frac{8}{8+11} = .42$$

$$\text{precision} = \frac{60}{60+55} = .52$$

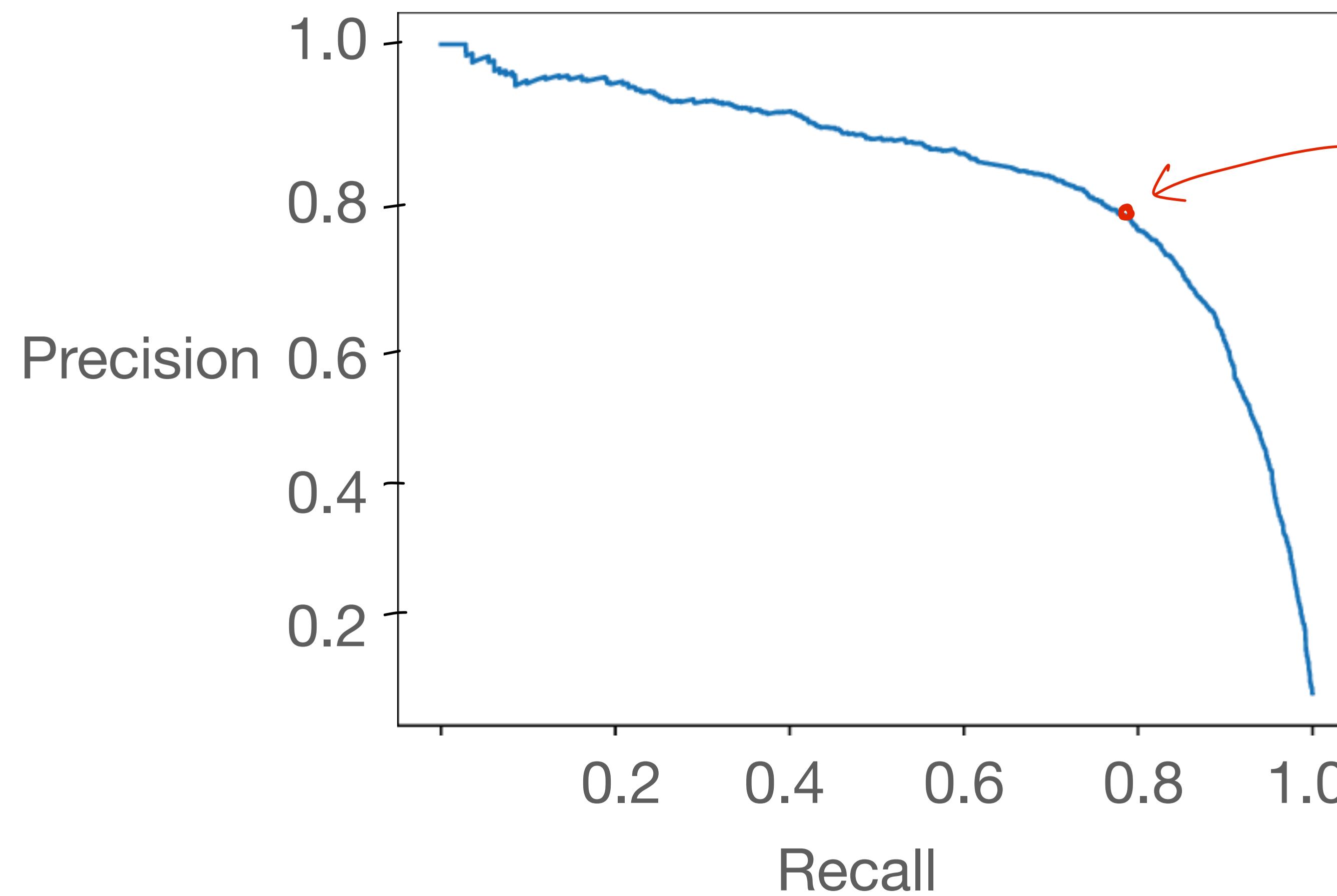
$$\text{precision} = \frac{200}{200+33} = .86$$

$$\text{microaverage precision} = \frac{268}{268+99} = .73$$

$$\text{macroaverage precision} = \frac{.42+.52+.86}{3} = .60$$

(Credits: Dan Jurafsky)

Precision Recall tradeoff



Maximum F1

Vary hyperparameters

- Smoothing α
- Threshold T

$$\frac{P(+ | d)}{P(- | d)} > T$$

Tune on validation set

Train, val, test split

- Train model on **training** set
- Tune hyperparameters on **validation** set
- Evaluate performance on unseen **test** set



Why do we do this?

Summary

- Evaluation Metrics
 - Accuracy - coarse metric
 - Precision, Recall, F1 for each class
- Aggregated scores
 - Macro-average: Compute for each class, then average
 - Micro-average: Collect predictions for all classes and jointly evaluate (dominated by common classes)
- Precision-Recall curve: pick threshold for maximum F1
 - Use validation set to tune hyperparameters, test set should remain “unseen”