

# ML Evaluation + Logistic Regression Models

---

Instructor: Lara J. Martin (she/they)

TA: Omkar Kulkarni (he)

<https://laramartin.net/NLP-class/>

*Slides modified from Dr. Frank Ferraro*

# Learning Objectives

---

Extend P/R to multi-class problems

Identify when you might want certain evaluation metrics over others

Model classification problems using logistic regression

Define appropriate features for a logistic regression problem

# Review

---

Argmax:

Returning the argument corresponding to the maximum probability of a distribution

Precision:

% of selected items that are correct

Recall:

% of correct items that are selected

Accuracy:

% of items that are correct

# Review: Contingency Table

<i>What label does our system predict? (↓)</i>	<i>What is the actual label?</i>	
	Actual Target Class ("●")	Not Target Class ("○")
Selected/ Guessed ("●")		
Not selected/ not guessed ("○")		

# Review: Contingency Table

<i>What label does our system predict? (↓)</i>	<i>What is the actual label?</i>	
	Actual Target Class ("●")	Not Target Class ("○")
Selected/ Guessed ("●")	True Positive (TP)	
Not selected/ not guessed ("○")		

# Review: Contingency Table

<i>What label does our system predict? (↓)</i>	<i>What is the actual label?</i>	
	Actual Target Class ("●")	Not Target Class ("○")
Selected/ Guessed ("●")	True Positive (TP)	False Positive (FP)
Not selected/ not guessed ("○")		

# Review: Contingency Table

<i>What label does our system predict? (↓)</i>	<i>What is the actual label?</i>	
	Actual Target Class ("●")	Not Target Class ("○")
Selected/ Guessed ("●")	True Positive (TP)	False Positive (FP)
Not selected/ not guessed ("○")	False Negative (FN)	

# Review: Contingency Table

<i>What label does our system predict? (↓)</i>	<i>What is the actual label?</i>	
	Actual Target Class ("●")	Not Target Class ("○")
Selected/ Guessed ("●")	True Positive (TP)	False Positive (FP)
Not selected/ not guessed ("○")	False Negative (FN)	True Negative (TN)





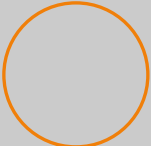
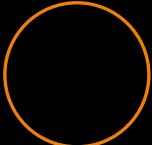
# The Importance of “Polarity” in Binary Classification

---

What are you trying to “identify” in your classification?

That is, are you trying to find  or ?


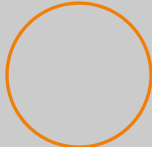



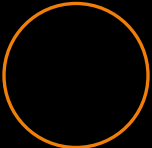


If  is our target

		Correct Value	
Guessed Value			
		?	?
		?	?







Where do  
TP / FP / FN / FN go?

If  is our target







---

		Correct Value	
			
Guessed Value		<i>TP</i> 	<i>FP</i> 
		<i>FN</i> 	<i>TN</i> 

If  is our target

Predicted:      


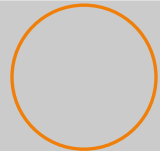



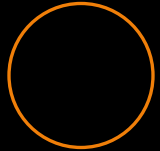


---

Actual:      

$$acc = \frac{TP+TN}{TP+FP+FN+TN}$$

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


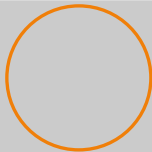
$$R = \frac{TP}{TP+FN}$$

		Correct Value	
			
Guessed Value		$TP = 2$ 	$FP = 2$ 
		$FN = 1$ 	$TN = 1$ 

What are the  
accuracy, recall, and  
precision values?

Accuracy: 50%  
Recall: 66.67%  
Precision: 50%






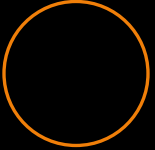


If ○ is our target

		Correct Value	
Guessed Value	●		
	○	?	?
		?	?

Where do  
TP / FP / FN / FN go?

If  is our target

---

		Correct Value	
			
Guessed Value		<i>TN</i> 	<i>FN</i> 
		<i>FP</i> 	<i>TP</i> 

If ○ is our target

Predicted: ○ ● ● ● ○ ●

Actual: ● ● ● ○ ○ ○

$$acc = \frac{TP+TN}{TP+FP+FN+TN}$$

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




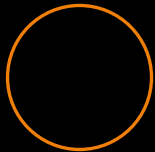
$$R = \frac{TP}{TP+FN}$$

		Correct Value	
		●	○
Guessed Value	●	<div></div> <div>●</div> <div><math>TN</math> ○ = 2</div>	<div></div> <div>○</div> <div><math>FN</math> ○ = 2</div>
	○	<div></div> <div>○</div> <div><math>FP</math> ○ = 1</div>	<div></div> <div>○</div> <div><math>TP</math> ○ = 1</div>

What are the  
accuracy, recall, and  
precision values?

Accuracy: 50%  
Recall: 33.34%  
Precision: 50%

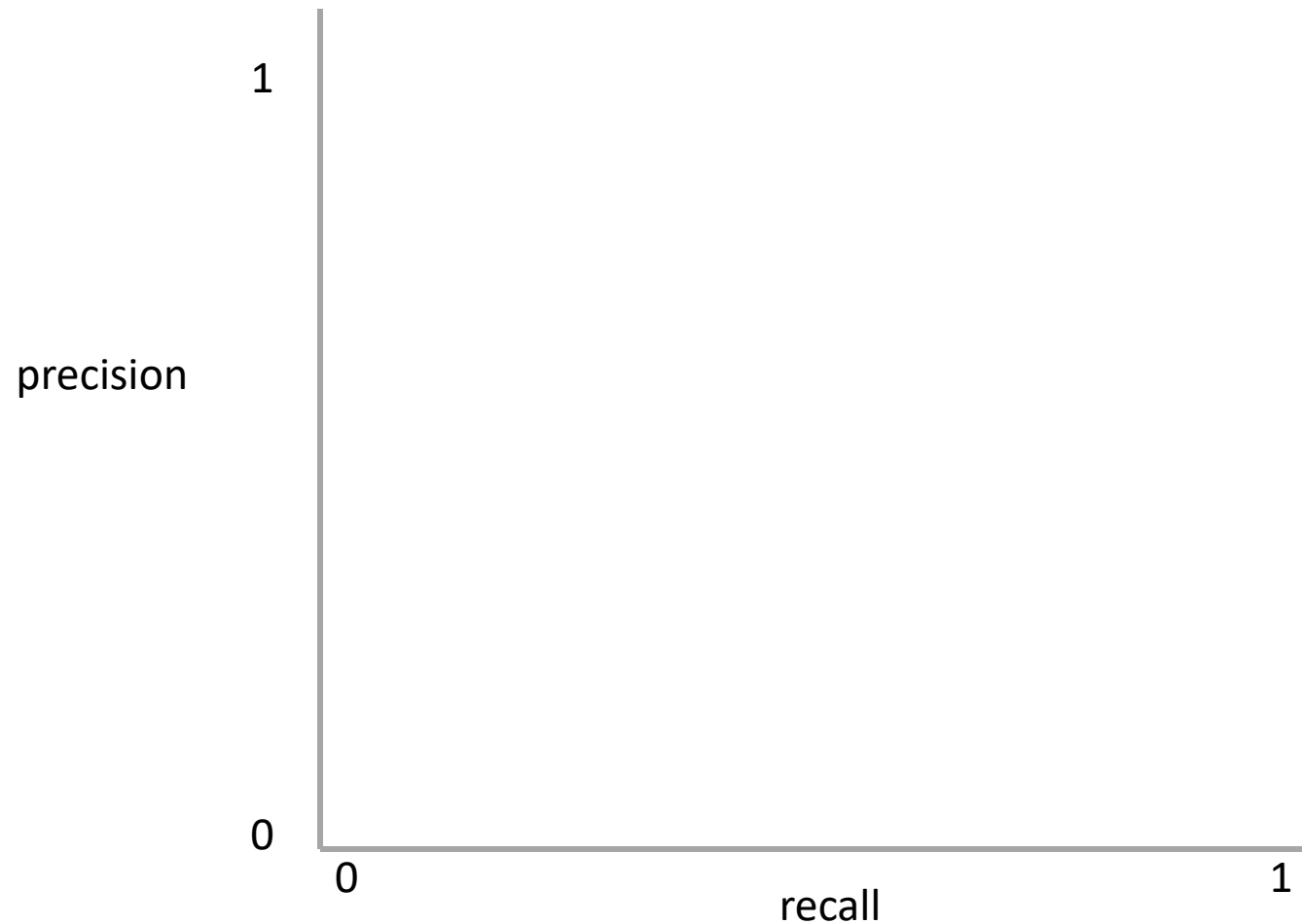
# When there are two classes, TP/TN & FP/FN are symmetrical

		Correct Value	
			
Guessed Value		$TP \text{ }  = TN \text{ } $	$FP \text{ }  = FN \text{ } $
		$FN \text{ }  = FP \text{ } $	$TN \text{ }  = TP \text{ } $

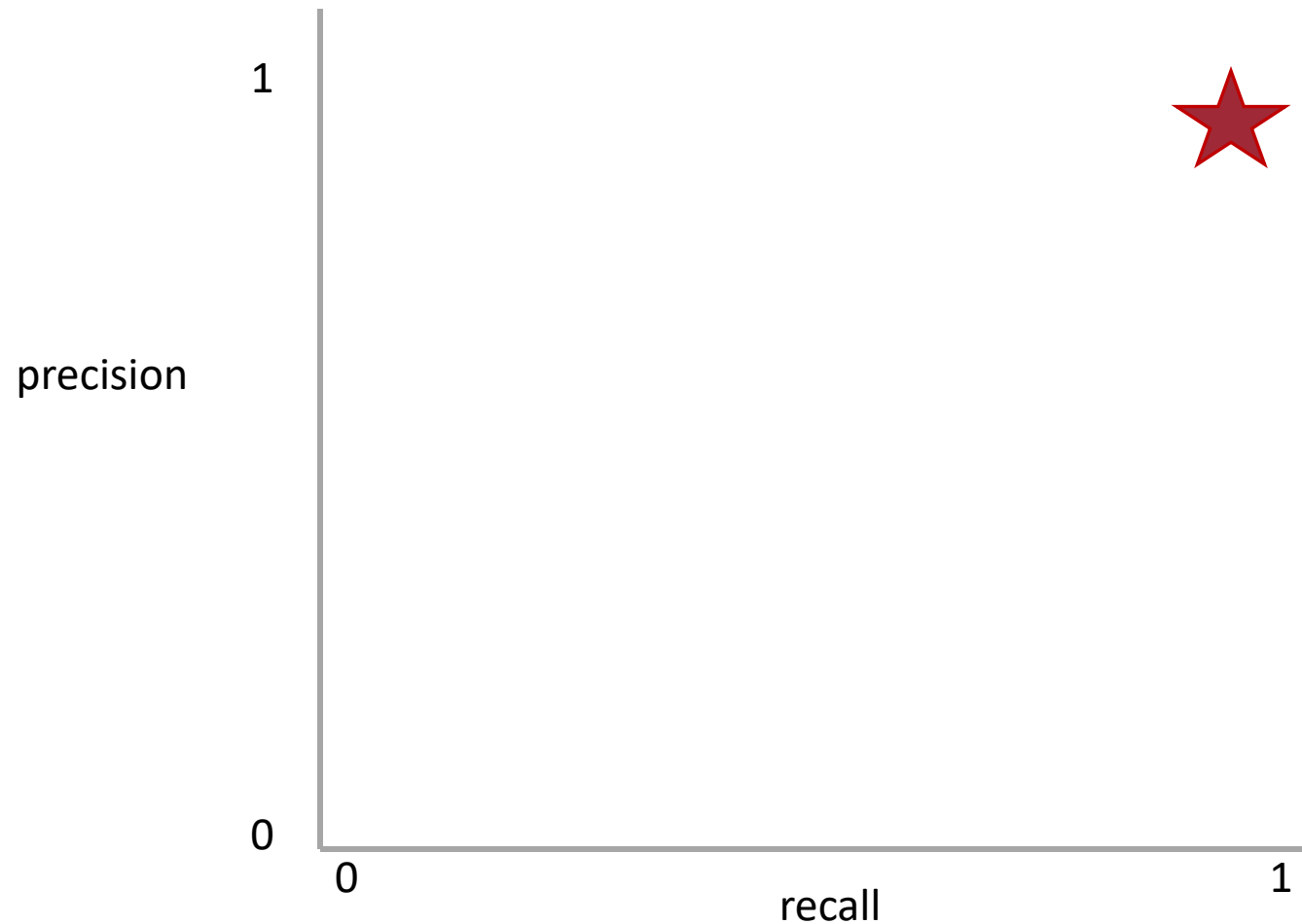


# Precision and Recall Present a Tradeoff

Q: Where do you  
want your ideal  
model ?



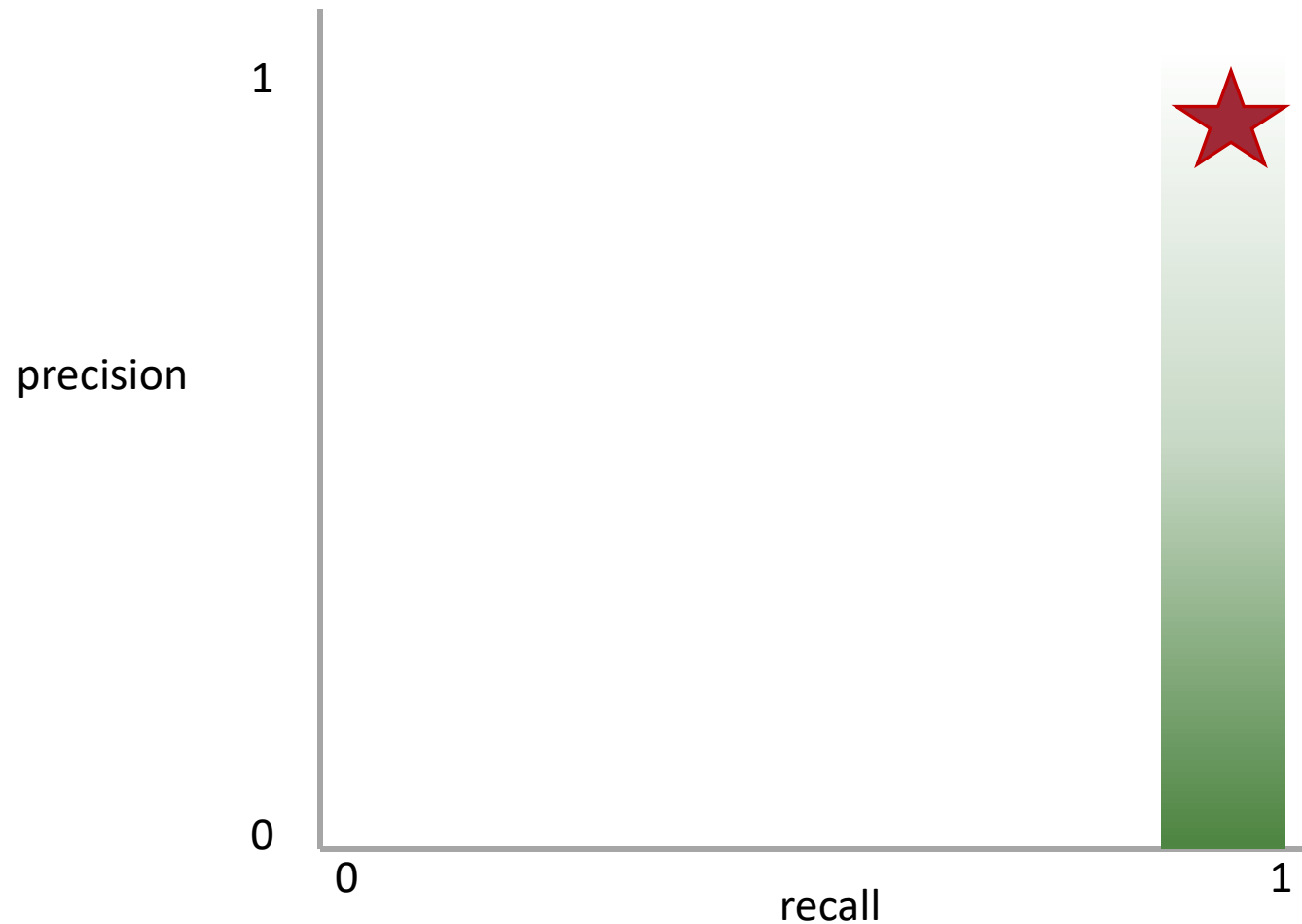
# Precision and Recall Present a Tradeoff



Q: Where do you want your ideal model?

Q: You have a model that always identifies correct instances. Where on this graph is it?

# Precision and Recall Present a Tradeoff

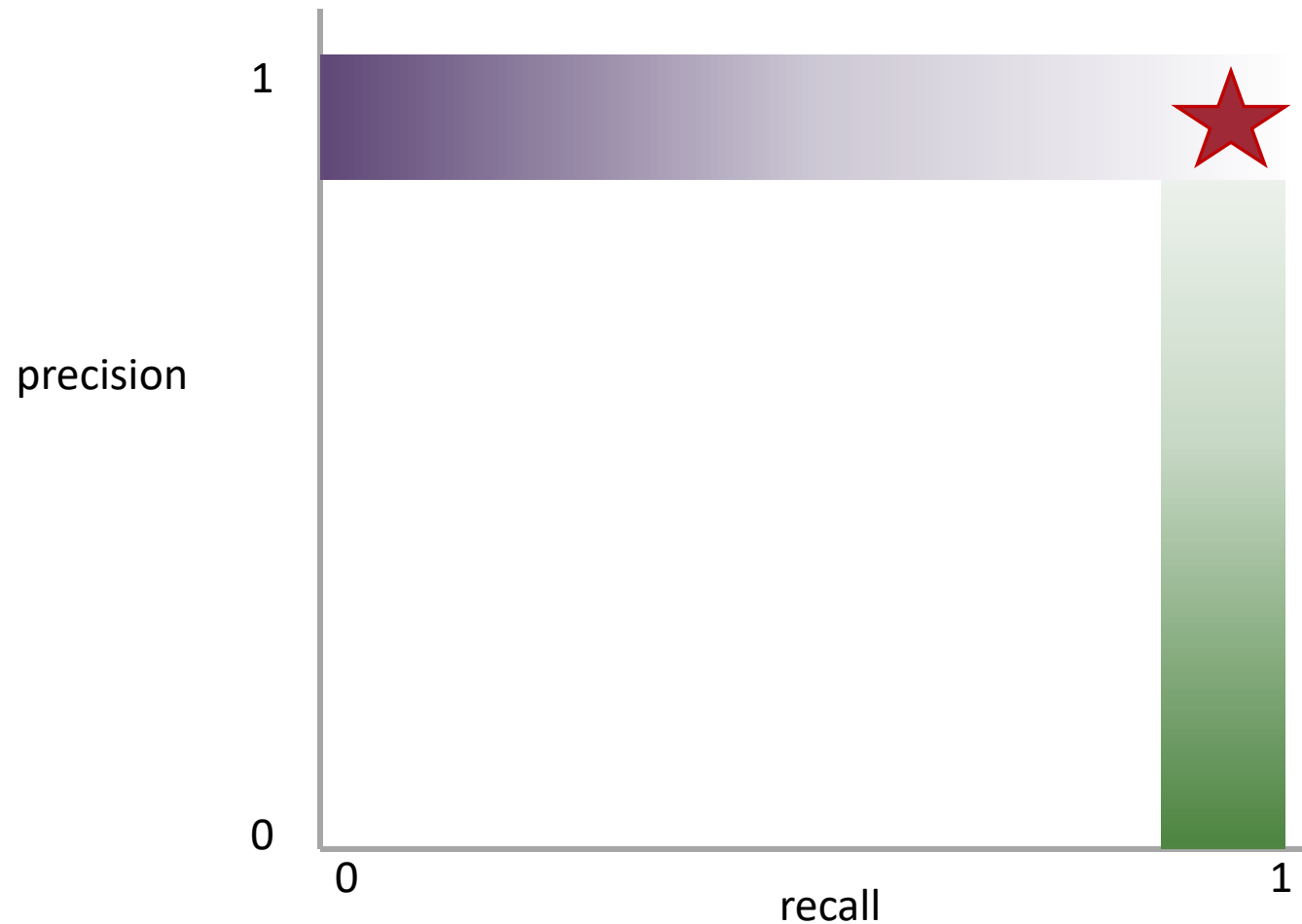


Q: Where do you want your ideal model ?

Q: You have a model that always identifies correct instances. Where on this graph is it?

Q: You have a model that only make correct predictions. Where on this graph is it?

# Precision and Recall Present a Tradeoff

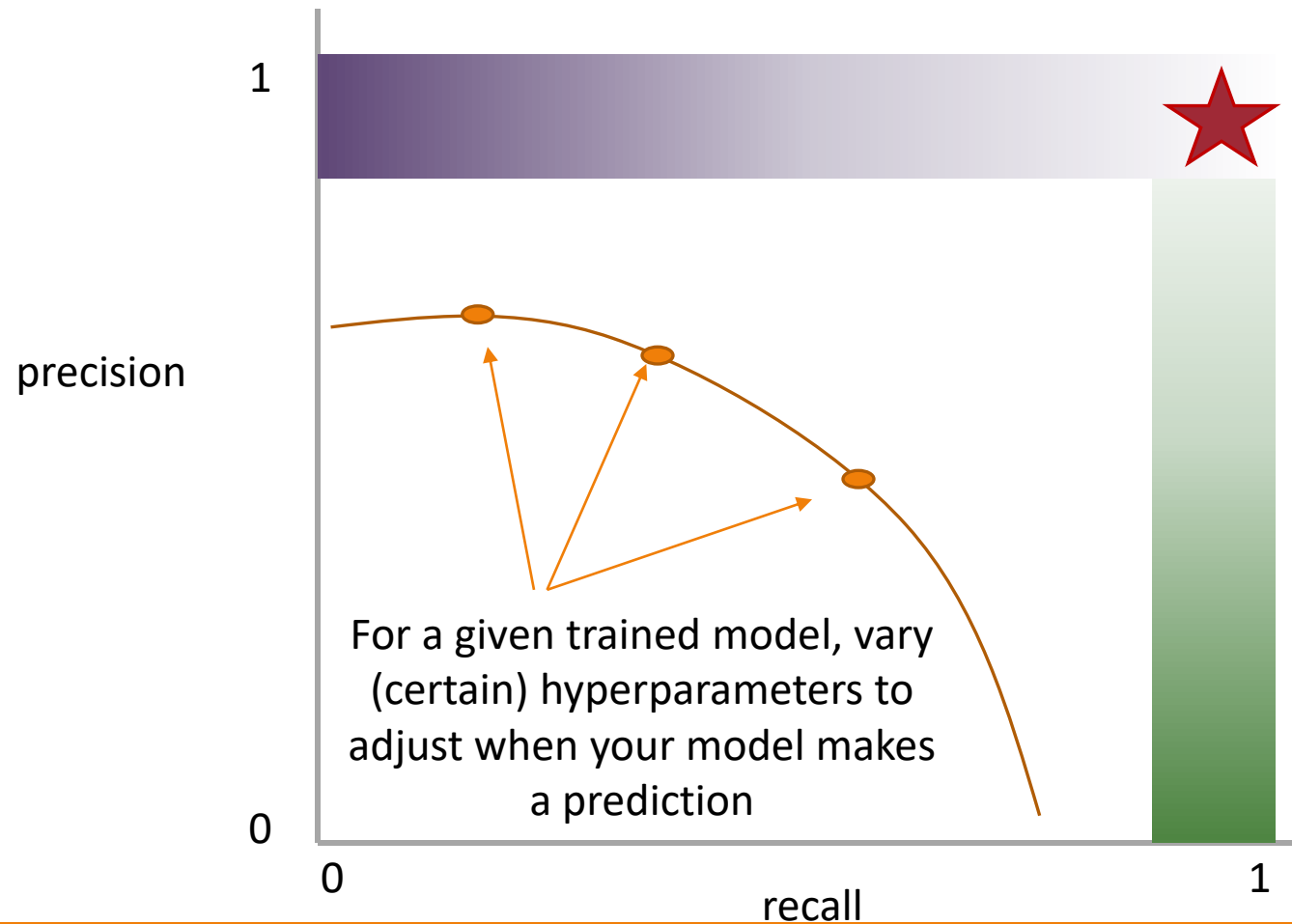


Q: Where do you want your ideal model ?

Q: You have a model that always identifies correct instances. Where on this graph is it?

Q: You have a model that only make correct predictions. Where on this graph is it?

# Precision and Recall Present a Tradeoff



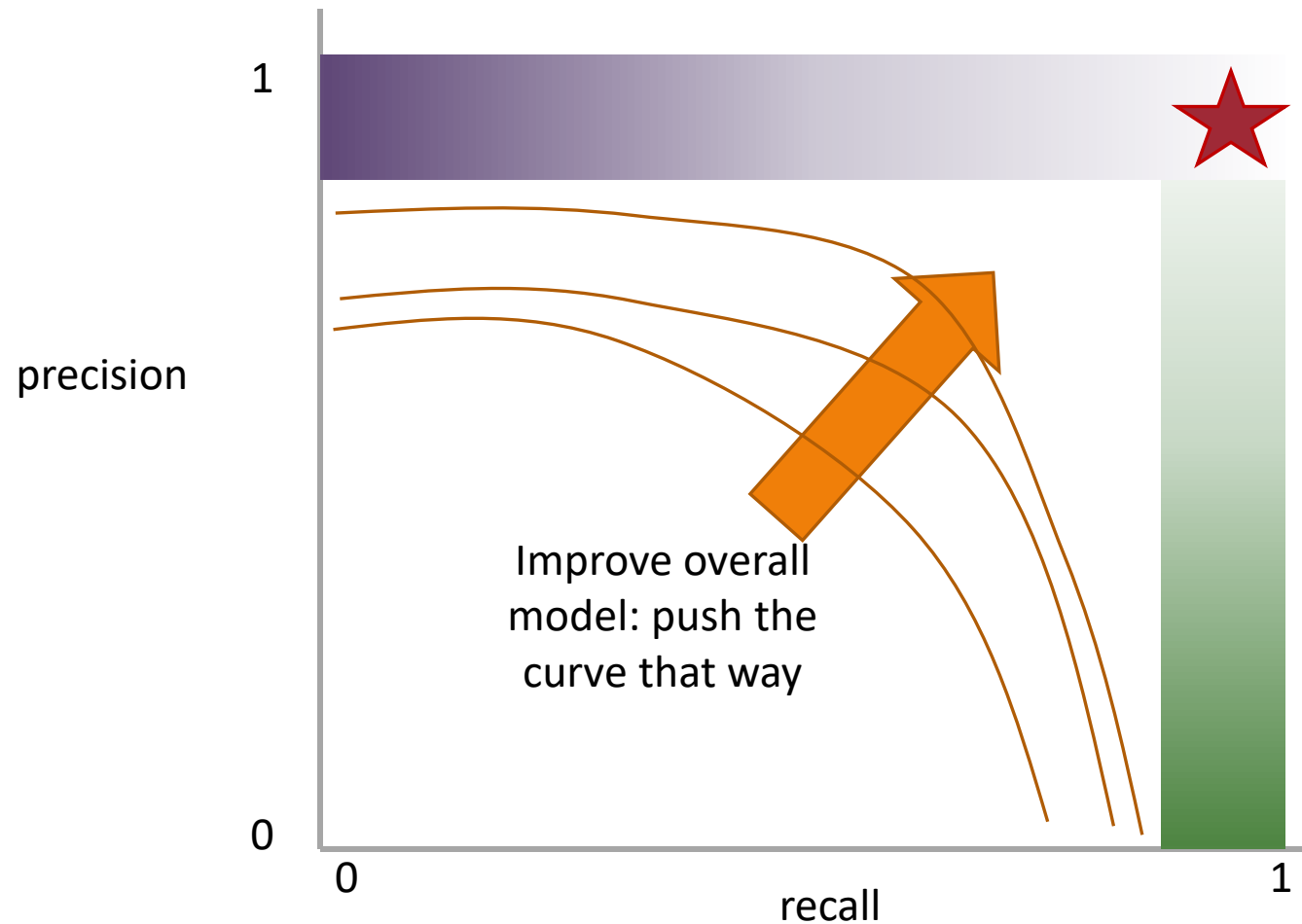
Q: Where do you want your ideal model ?

Q: You have a model that always identifies correct instances. Where on this graph is it?

Q: You have a model that only make correct predictions. Where on this graph is it?

Idea: measure the tradeoff between precision and recall

# Precision and Recall Present a Tradeoff



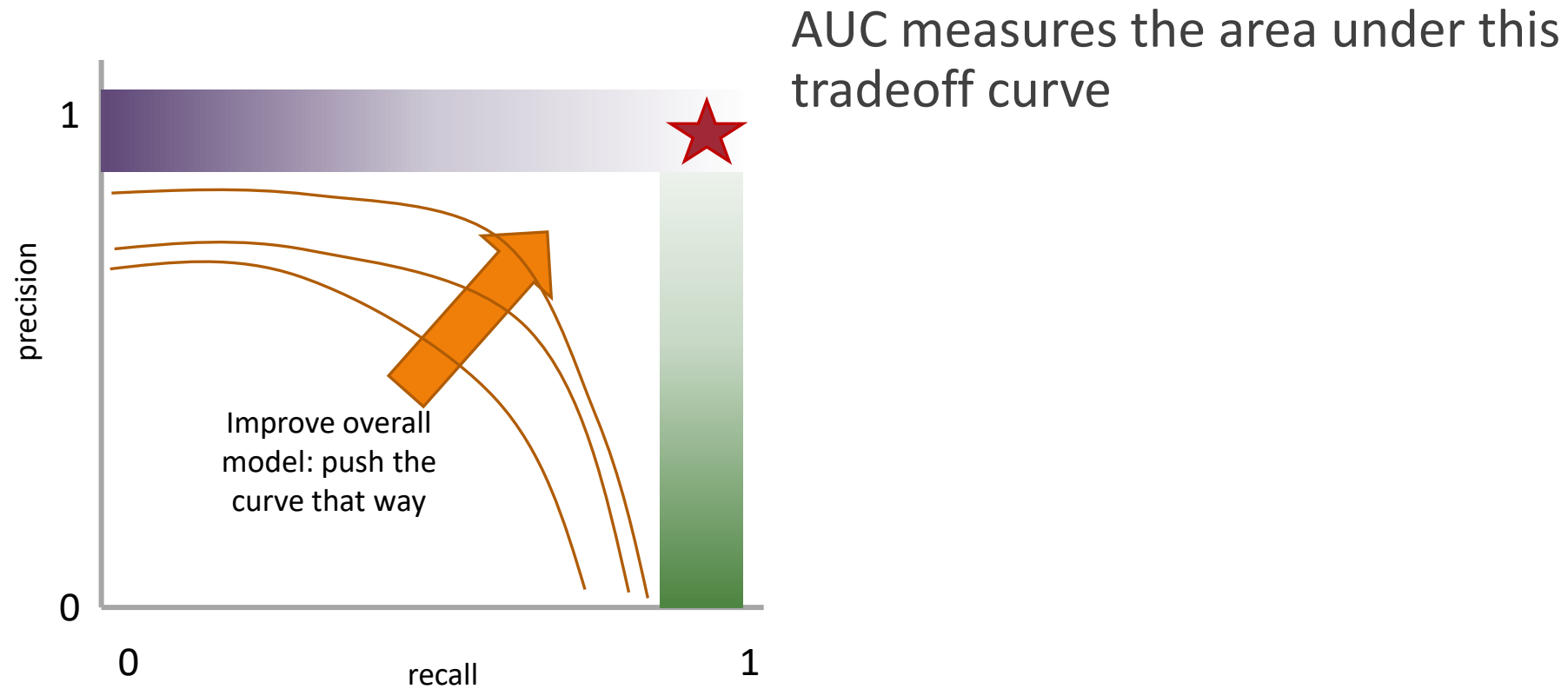
Q: Where do you want your ideal model ?

Q: You have a model that always identifies correct instances. Where on this graph is it?

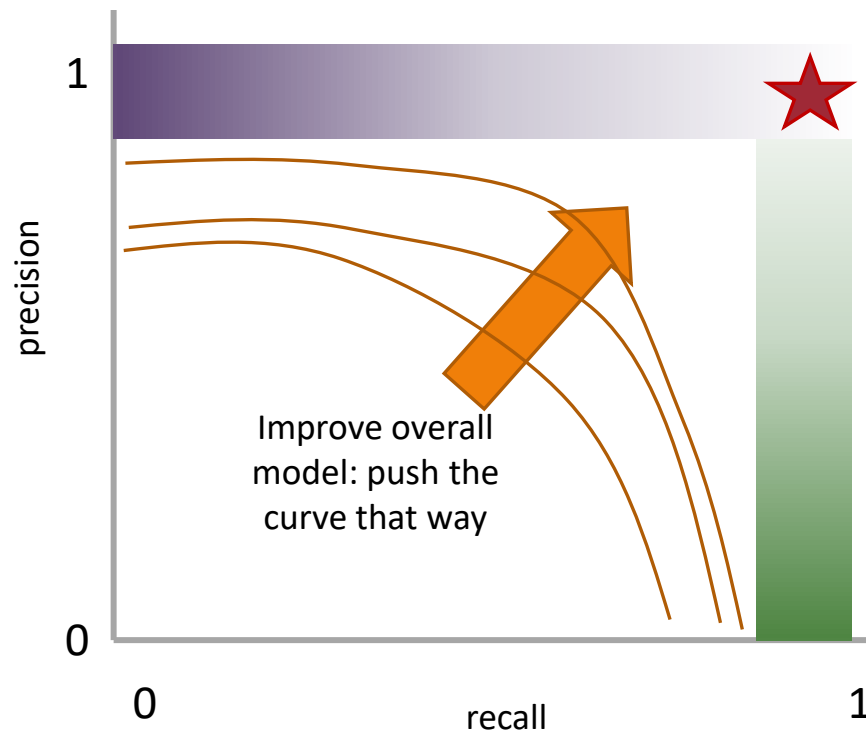
Q: You have a model that only make correct predictions. Where on this graph is it?

Idea: measure the tradeoff between precision and recall

# Measure this Tradeoff: Area Under the Curve (AUC)



# Measure this Tradeoff: Area Under the Curve (AUC)



AUC measures the area under this tradeoff curve

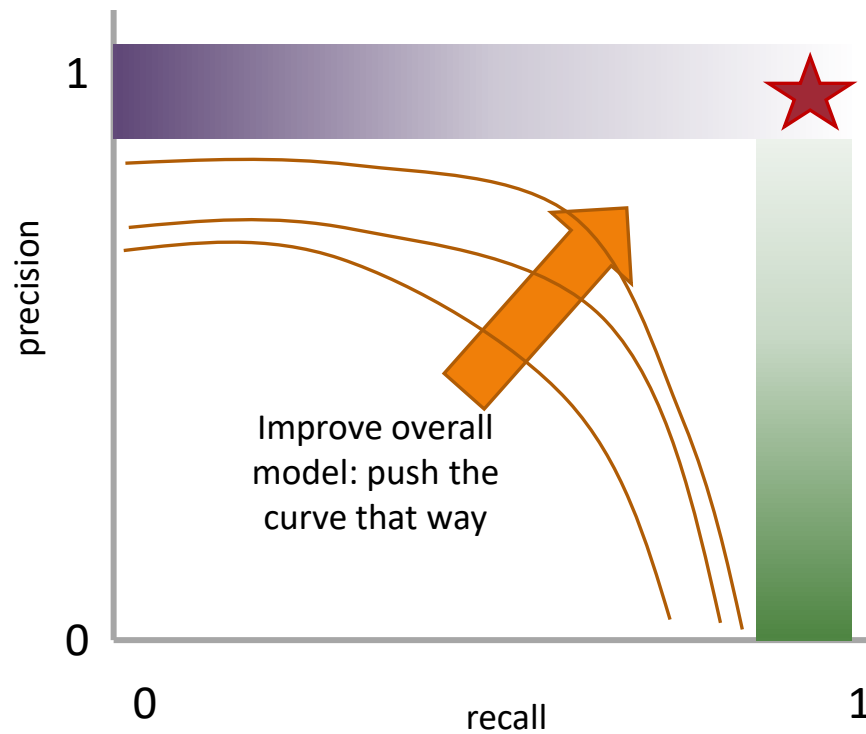
## 1. Computing the curve

You need true labels & predicted labels with some score/confidence estimate

Threshold the scores and for each threshold compute precision and recall



# Measure this Tradeoff: Area Under the Curve (AUC)



AUC measures the area under this tradeoff curve

1. Computing the curve  
You need true labels & predicted labels with some score/confidence estimate  
Threshold the scores and for each threshold compute precision and recall
2. Finding the area  
How to implement: trapezoidal rule (& others)

**In practice:** external library like the `sklearn.metrics` module

# A combined measure: F1 (or F-score)

---

Weighted (harmonic) average of **P**recision & **R**ecall

F1 measure: equal weighting between precision and recall

$$F_1 = \frac{2 * P * R}{P + R}$$

# A combined measure: F1 (or F-score)

---

Weighted (harmonic) average of **P**recision & **R**ecall

F1 measure: equal weighting between precision and recall

$$F_1 = \frac{2 * P * R}{P + R} = \frac{2 * TP}{2 * TP + FP + FN}$$

(useful when  $P = R = 0$ )

# Comparing Accuracy & F1

**Accuracy:** % of items correct

$$\frac{TP + TN}{TP + FP + FN + TN}$$

$$F_1 = \frac{2 * P * R}{P + R} = \frac{2 * TP}{2 * TP + FP + FN}$$

When would you want to use accuracy vs F1?

Accuracy works better if the dataset is balanced

Accuracy takes everything in consideration

F-Score is focused on TP

	Actually Target	Actually Not Target
Selected/Guessed	True Positive (TP)	False Positive (FP)
Not select/not guessed	False Negative (FN)	True Negative (TN)

# P/R/F in a Multi-class Setting: Micro- vs. Macro-Averaging

---

*If we have more than one class, how do we combine multiple performance measures into one quantity?*

# P/R/F in a Multi-class Setting: Micro- vs. Macro-Averaging

**Macroaveraging:** Compute performance for each class, then average.

$$\text{macroprecision} = \frac{1}{C} \sum_c \frac{TP_c}{TP_c + FP_c} = \frac{1}{C} \sum_c \text{precision}_c$$

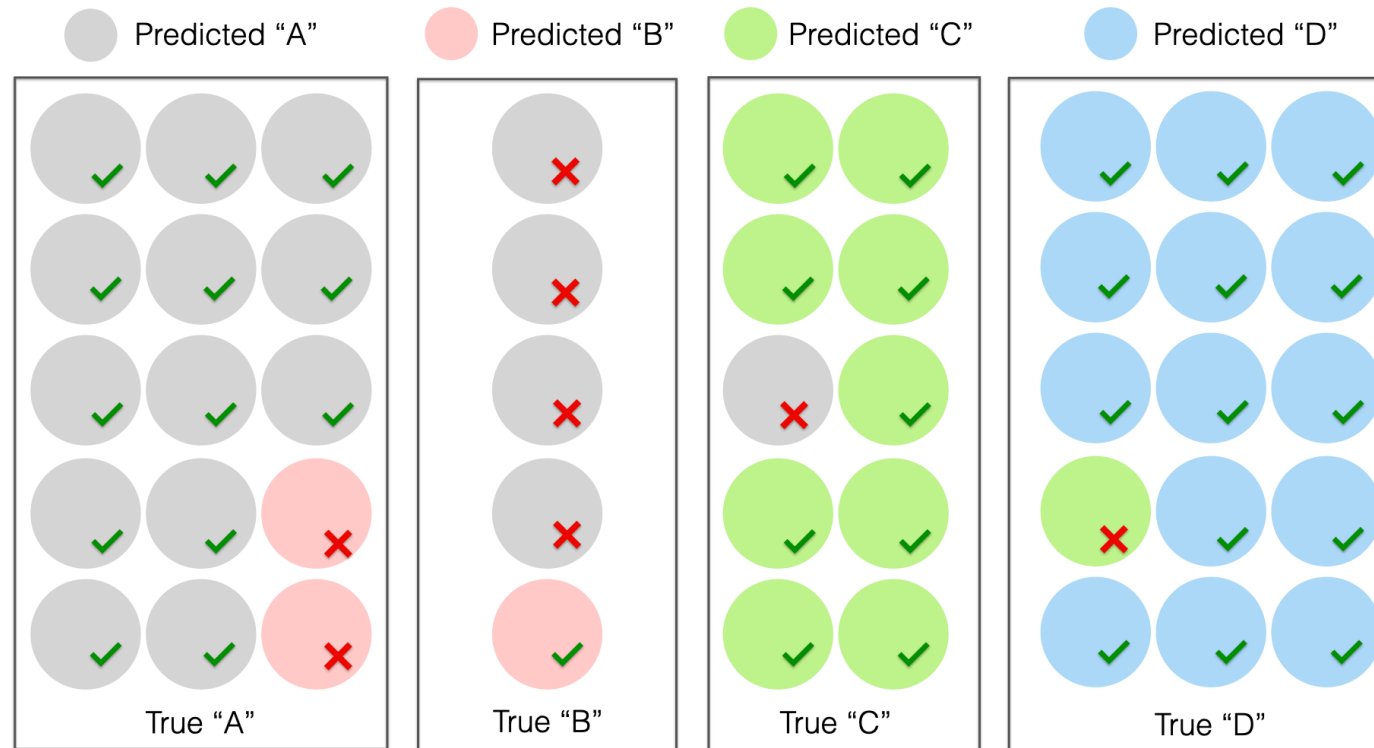
$$\text{macrorecall} = \frac{1}{C} \sum_c \frac{TP_c}{TP_c + FN_c} = \frac{1}{C} \sum_c \text{recall}_c$$

**Microaveraging:** Collect decisions for all classes, compute contingency table, evaluate.

$$\text{microprecision} = \frac{\sum_c TP_c}{\sum_c TP_c + \sum_c FP_c}$$

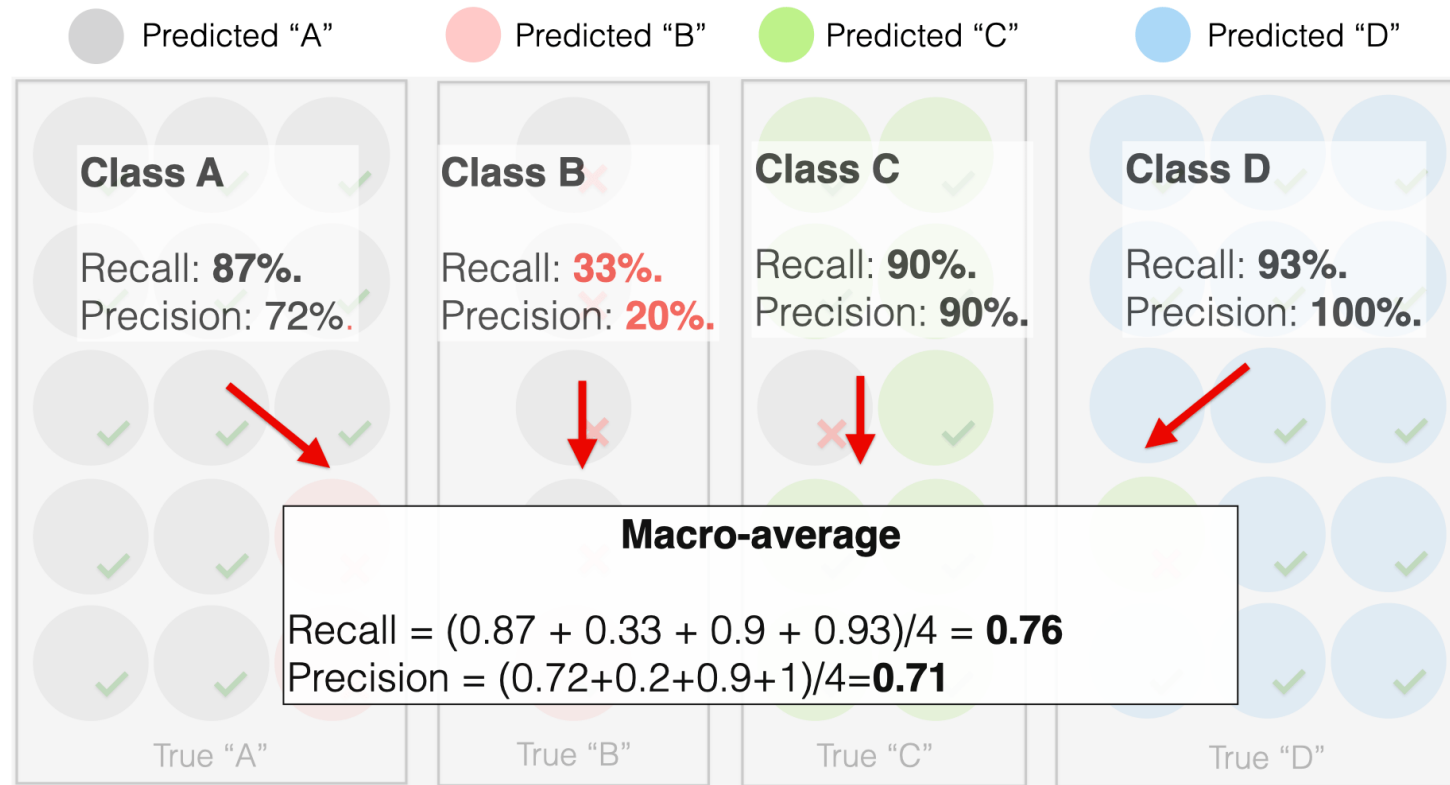
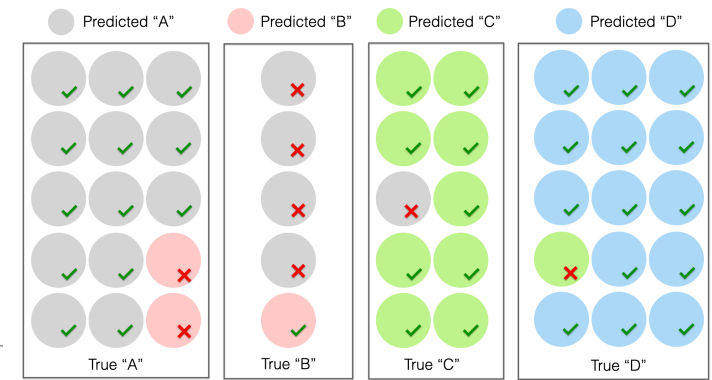
$$\text{microrecall} = \frac{\sum_c TP_c}{\sum_c TP_c + \sum_c FN_c}$$

# Macro/Micro Example



Each *class* has equal weight

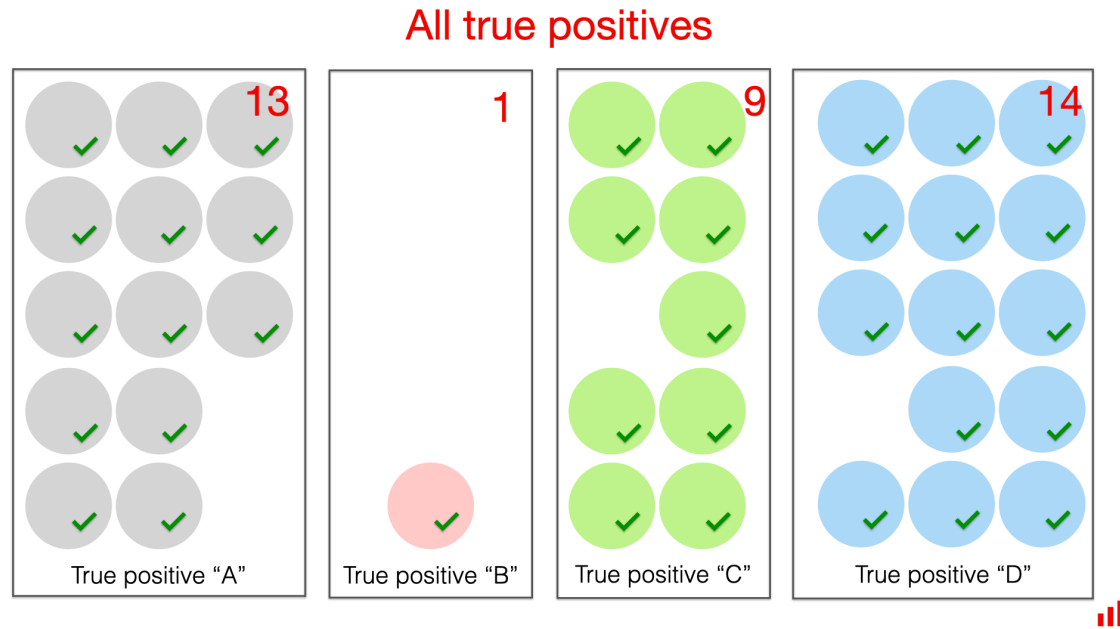
# Macro-Average





Each *instance* has equal weight

# Micro-Average

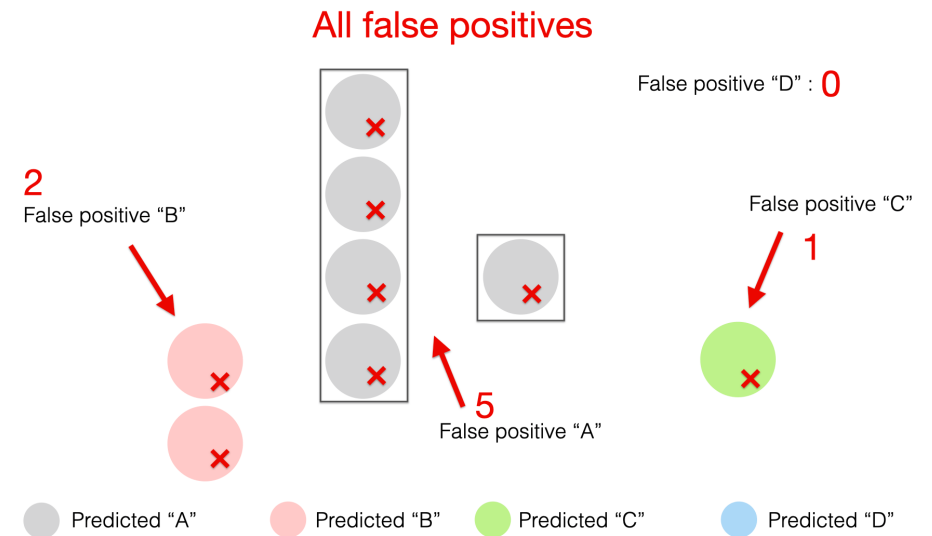
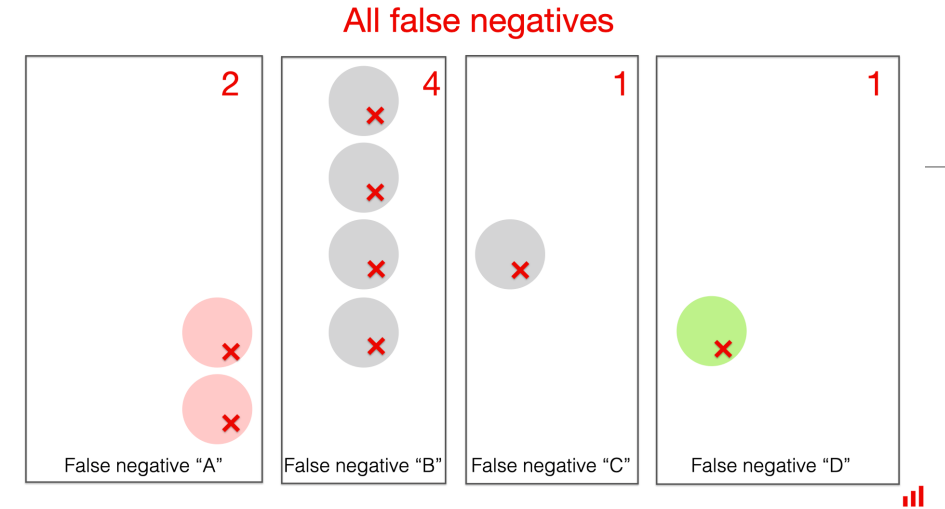


Total TP: 13 + 1 + 9 + 14 = 37

Total FP: 2 + 5 + 1 + 0 = 8

Total FN: 2 + 4 + 1 + 1 = 8

$$\text{Precision}_{\text{Micro-average}} = \frac{13 + 1 + 9 + 14}{13 + 1 + 9 + 14 + 2 + 5 + 1 + 0} = 0.82$$
$$\text{Recall}_{\text{Micro-average}} = \frac{13 + 1 + 9 + 14}{13 + 1 + 9 + 14 + 2 + 4 + 1 + 1} = 0.82$$



# P/R/F in a Multi-class Setting: Micro- vs. Macro-Averaging

**Macroaveraging:** Compute performance for each class, then average.

$$\text{macroprecision} = \frac{1}{C} \sum_c \frac{TP_c}{TP_c + FP_c} = \frac{1}{C} \sum_c \text{precision}_c$$

$$\text{macrorecall} = \frac{1}{C} \sum_c \frac{TP_c}{TP_c + FN_c} = \frac{1}{C} \sum_c \text{recall}_c$$

**Microaveraging:** Collect decisions for all classes, compute contingency table, evaluate.

$$\text{microprecision} = \frac{\sum_c TP_c}{\sum_c TP_c + \sum_c FP_c}$$

$$\text{microrecall} = \frac{\sum_c TP_c}{\sum_c TP_c + \sum_c FN_c}$$

When would we want to prefer micro-averaging vs macro-averaging?

# But how do we compute stats for multiple classes?

---

We already saw how the “polarity” affects the stats we compute...

Two main approaches. Either:

1. Compute “one-vs-all” 2x2 tables. OR
2. Generalize the 2x2 tables and compute per-class TP / FP / FN based on the diagonals and off-diagonals



# 1. Compute “one-vs-all” 2x2 tables


Predicted



Actual



Look for 	Actually Target	Actually Not Target	Look for 	Actually Target	Actually Not Target
Selected/G uessed	True Positive (TP)	False Positive (FP)	Selected/G uessed	True Positive (TP)	False Positive (FP)
Not select/not guessed	False Negative (FN)	True Negative (TN)	Not select/not guessed	False Negative (FN)	True Negative (TN)

Look for 	Actually Target	Actually Not Target
Selected/G uessed	True Positive (TP)	False Positive (FP)
Not select/not guessed	False Negative (FN)	True Negative (TN)

# 1. Compute “one-vs-all” 2x2 tables

Predicted




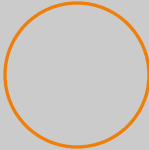


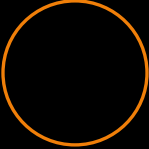

Actual



Look for 	Actually Target	Actually Not Target	Look for 	Actually Target	Actually Not Target
Selected/G uessed	2	1	Selected/G uessed	2	1
Not select/not guessed	2	4	Not select/not guessed	1	5

Look for 	Actually Target	Actually Not Target
Selected/G uessed	1	2
Not select/not guessed	1	5

## 2. Generalizing the 2-by-2 contingency table

		Correct Value		
				
Guessed Value		#	#	#
		#	#	#
		#	#	#

This is also called a **Confusion Matrix**


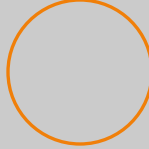


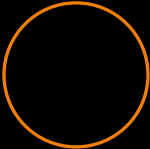

## 2. Generalizing the 2-by-2 contingency table

Predicted



Actual



		Correct Value		
				
Guessed Value		a #	b #	c #
		d #	e #	f #
		g #	h #	i #


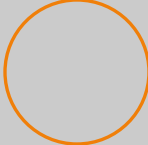


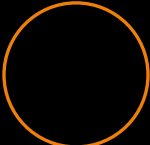

## 2. Generalizing the 2-by-2 contingency table

Predicted



Actual



		Correct Value		
				
Guessed Value		a 2	b 0	c 1
		d 1	e 2	f 0
		g 1	h 1	i 1




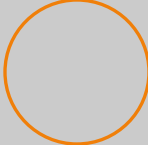


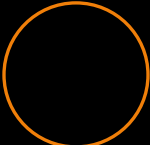

## 2. Generalizing the 2-by-2 contingency table

Predicted



Actual



		Correct Value		
				
Guessed Value		a 2	b 0	c 1
		d 1	e 2	f 0
		g 1	h 1	i 1

How do you compute  $TP_{\bullet}$ ?


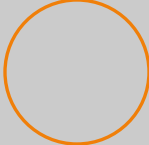


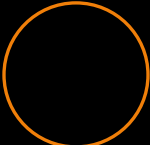

## 2. Generalizing the 2-by-2 contingency table

Predicted



Actual



		Correct Value		
				
Guessed Value		2	0	1
		1	2	0
		1	1	1

How do you compute  $TP_{\bullet}$ ?





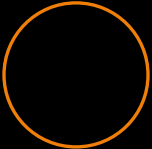

## 2. Generalizing the 2-by-2 contingency table

Predicted



Actual



		Correct Value		
				
Guessed Value		a 2	b 0	c 1
		d 1	e 2	f 0
		g 1	h 1	i 1

How do you compute  $FN$ ?


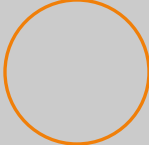


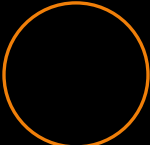

## 2. Generalizing the 2-by-2 contingency table

Predicted



Actual



		Correct Value		
				
Guessed Value		2	0	1
		1	2	0
		1	1	1

How do you compute  $FN$ ?


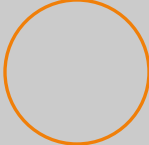


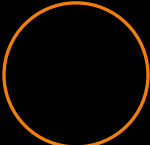

## 2. Generalizing the 2-by-2 contingency table

Predicted



Actual



		Correct Value		
				
Guessed Value		a 2	b 0	c 1
		d 1	e 2	f 0
		g 1	h 1	i 1

How do you compute  $FP_{\square}$ ?

## 2. Generalizing the 2-by-2 contingency table

Predicted




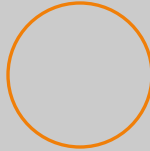


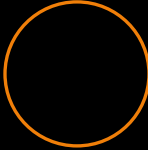

Actual



		Correct Value		
Guessed Value		2	0	1
		1	2	0
		1	1	1


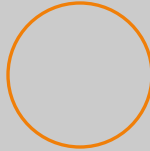


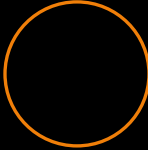

How do you compute  $FP_{\square}$ ?

# Performance of a Classifier using a Confusion Matrix

		Correct Value		
				
Guessed Value		80	9	11
		7	86	7
		2	8	9

Q: Is this a good result?


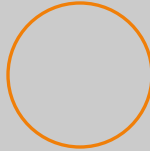


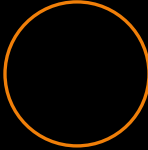

# Performance of a Classifier using a Confusion Matrix

		Correct Value		
				
Guessed Value		30	40	30
		25	30	50
		30	35	35

Q: Is this a good result?



# Performance of a Classifier using a Confusion Matrix

		Correct Value		
				
Guessed Value		7	3	90
		4	8	88
		3	7	90

Q: Is this a good result?

# Max Entropy / Logistic Regression Models

---

# Outline

---

Maximum Entropy classifiers

Defining the model

Defining the objective

Learning: Optimizing the objective

# Outline

---

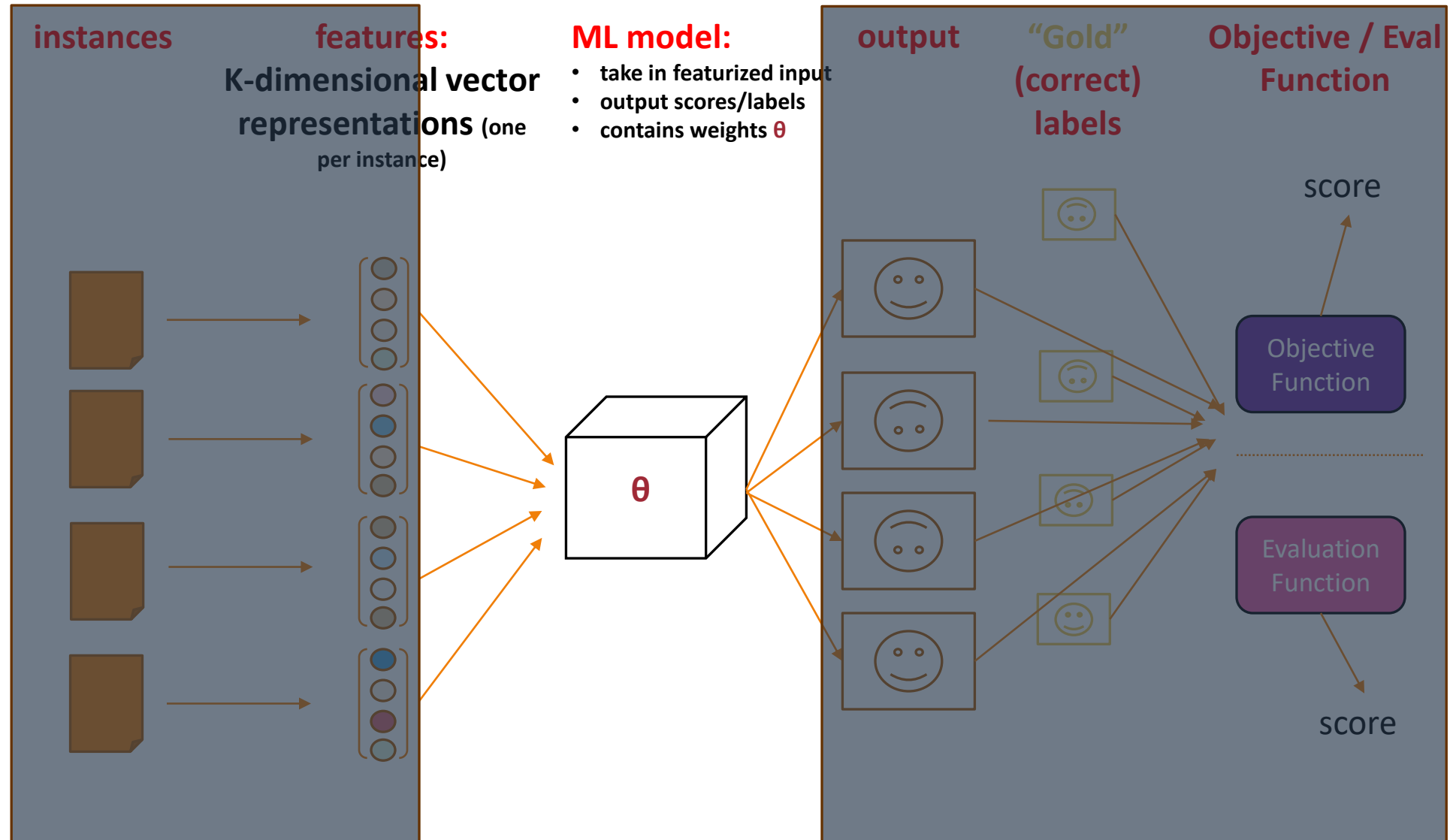
Maximum Entropy classifiers

**Defining the model**

Defining the objective

Learning: Optimizing the objective

# Defining the Model



# Terminology

---

common NLP term	Log-Linear Models
as statistical regression	(Multinomial) logistic regression Softmax regression
based in information theory	Maximum Entropy models (MaxEnt)
a form of	Generalized Linear Models
viewed as	Discriminative Naïve Bayes
to be cool today	Very shallow (sigmoidal) neural nets

# Maxent Models are Flexible

---

Maxent models can be used:

- to design discriminatively trained classifiers, or
- to create featureful language models

(among other approaches in NLP and ML more broadly)

# Examining Assumption 3 Made for Classification Evaluation

---

Given  $X$ , our classifier produces a score for each possible label

$$p(\text{●} | X) \text{ vs. } p(\text{○} | X)$$

$$\text{best label} = \arg \max_{\text{label}} P(\text{label} | \text{example})$$





## Key Take-away



We will *learn* this  
 $p(Y | X)$

**Conditional probability:**  
probability of event Y,  
assuming event X  
happens too

NLP pg. 477

# Maxent Models for Classification: Discriminatively or ...

---

Directly model  
the posterior

$$p(Y | X) = \textbf{maxent}(X; Y)$$

Discriminatively trained classifier

“Discriminative classifiers like logistic regression instead learn what features from the input are most useful to discriminate between the different possible classes.”

SLP, ch. 4

# Bayes' Rule

$$\underbrace{P(Y|X)}_{\text{Posterior}} = \frac{\overbrace{P(X|Y)}^{\text{Likelihood}} \cdot \overbrace{P(Y)}^{\text{Prior}}}{P(X)}$$

## Posterior:

probability of event Y  
with knowledge that X  
has occurred

NLP pg. 478

## Likelihood:

probability of event X  
given that Y has occurred

NLP pg. 478

## Prior:

probability of event X  
occurring (regardless of  
what other events  
happen)

NLP pg. 478


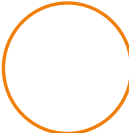
# Terminology: Posterior Probability

---

Posterior probability:

$$p(\text{●} | X) \text{ vs. } p(\text{○} | X)$$

Conditionally dependent probabilities:

- If  and  are the only two options:

$$p(\text{●} | X) + p(\text{○} | X) = 1$$

and

$$p(\text{●} | X) \geq 0, p(\text{○} | X) \geq 0$$

# Posterior Probability with Variables

---

$$p(\text{●} | X) \text{ vs. } p(\text{○} | X)$$



$$p(Y = \text{label}_1 | X) \text{ vs. } p(Y = \text{label}_0 | X)$$

# Maxent Models for Classification: Discriminatively or Generatively Trained

Directly model  
the posterior

$$p(Y | X) = \text{maxent}(X; Y)$$

**Discriminatively** trained classifier

Model the  
posterior with  
Bayes rule


$$p(Y | X) \propto \text{maxent}(X | Y)p(Y)$$

**Generatively** trained classifier with  
maxent-based language model

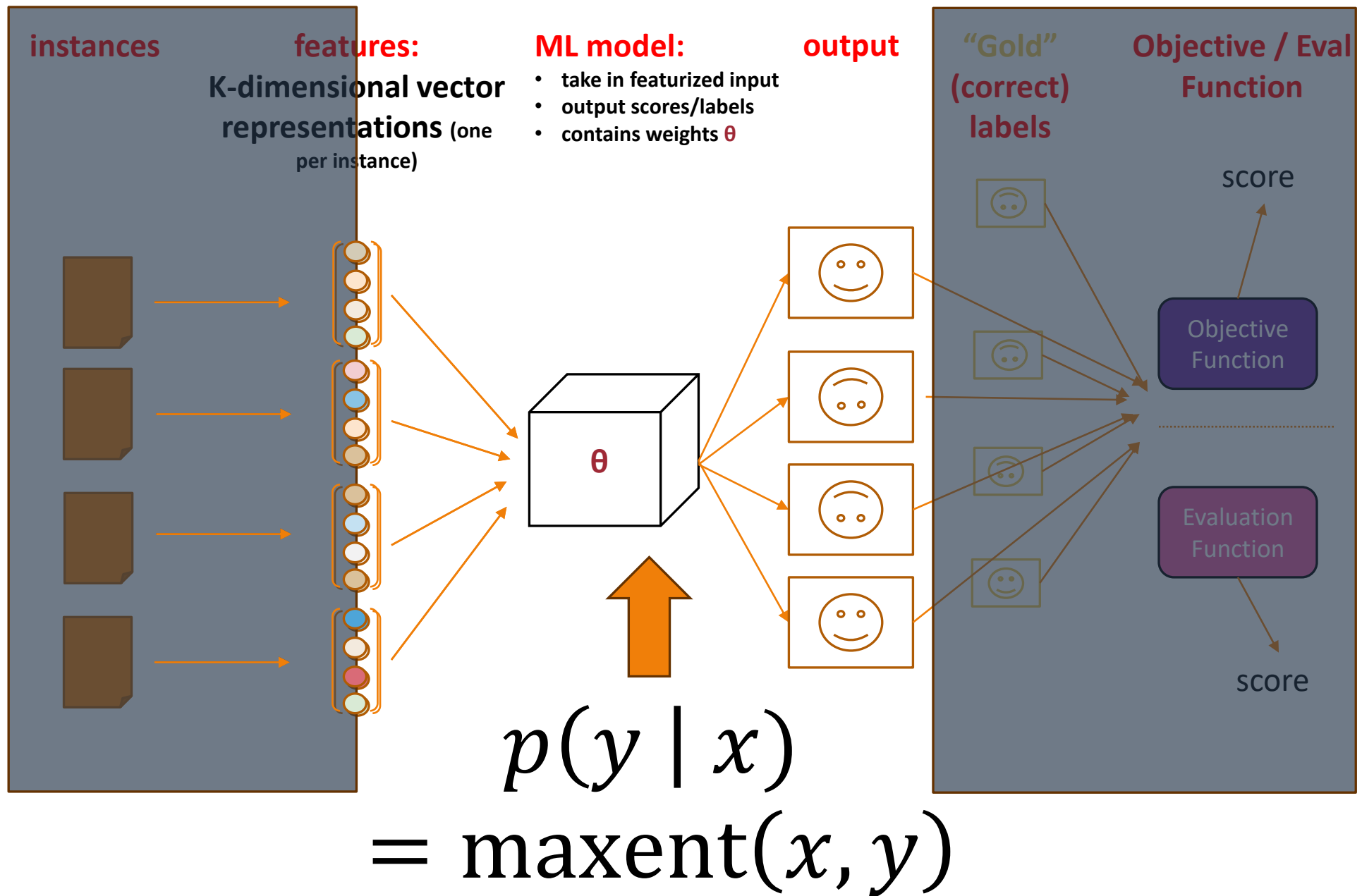
# Maximum Entropy (Log-linear) Models For Discriminatively Trained Classifiers

---

$$p(y | x) = \text{maxent}(x, y)$$



Modeled  
jointly!





# Core Aspects to Maxent Classifier $p(y|x)$

---

We need to define:

- **features**  $f(x)$  from  $x$  that are meaningful;
- **weights**  $\theta$  (at least one per feature, often one per feature/label combination) to say how important each feature is; and
- a way to **form probabilities** from  $f$  and  $\theta$

# Overview of Featurization

---

Common goal: probabilistic classifier  $p(y \mid x)$

Often done by defining **features** between  $x$  and  $y$  that are meaningful

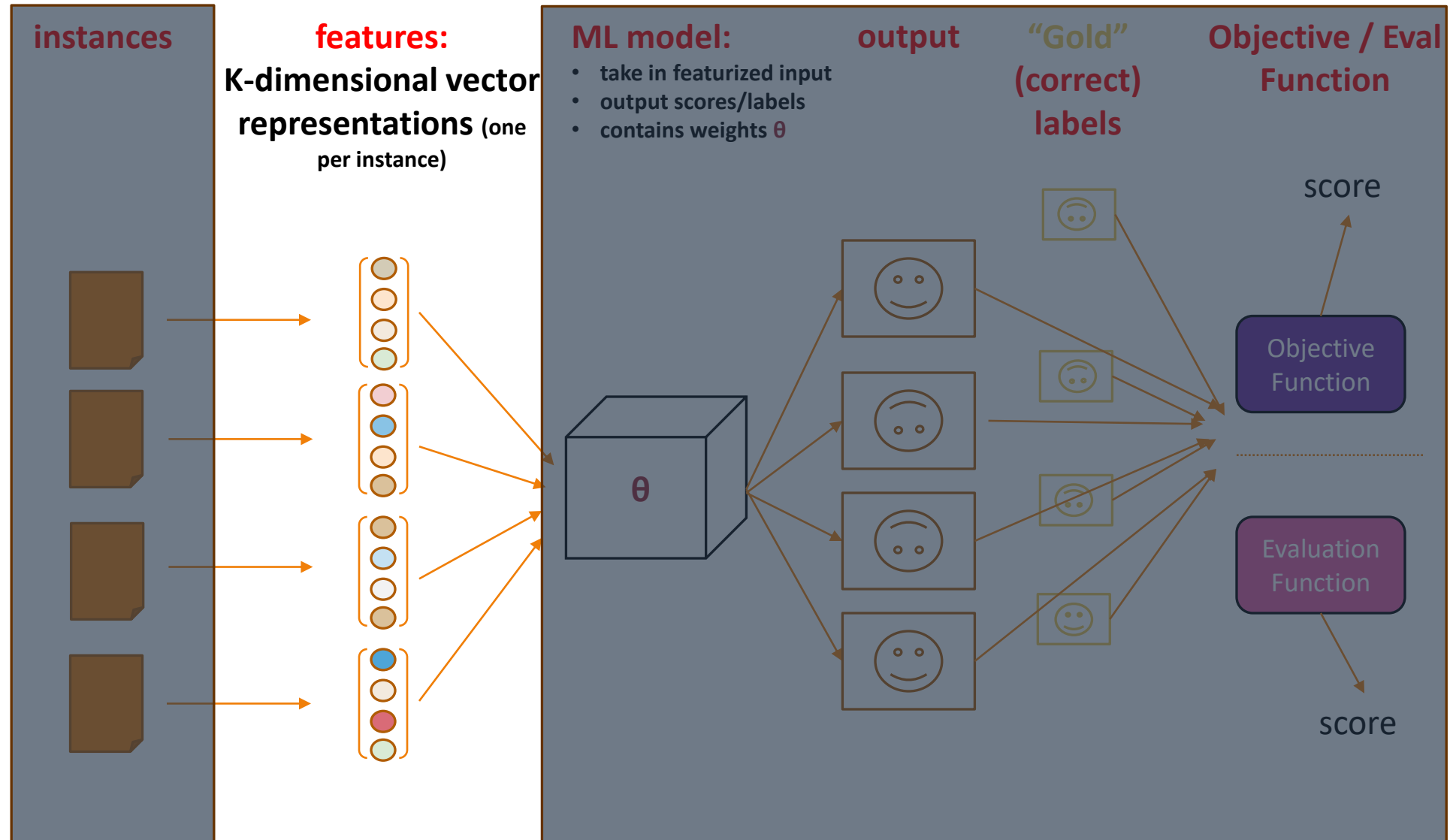
- Denoted by a **general vector of  $K$  features**

$$f(x) = (f_1(x), \dots, f_K(x))$$

Features can be thought of as “soft” rules

- E.g., POSITIVE sentiments tweets *may* be more likely to have the word “happy”

# Defining the Model



# Review: Document Classification via Bag-of-Words Features (Example)

Amazon acquired MGM in 2022, taking over a sprawling library that includes more than 4,000 feature films and 17,000 television shows. The tech behemoth also earned the rights to distribute all the Bond movies, but the new deal solidifies the company's oversight of Bond's big-screen future.

With  $V$  word types, define  $V$  feature functions  $f_i(x)$  as

$f_i(x)$  = # of times word type  $i$  appears in document  $x$

$$f(x) = (f_i(x))_i^V$$

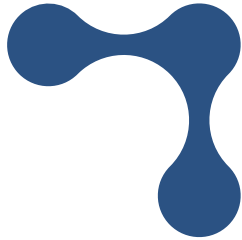
TECH  
NOT TECH

Core assumption:  
the label can be  
predicted from  
counts of individual  
word types

feature $f_i(x)$	value
Amazon	1
acquired	1
behemoth	1
Bond	2
...	
sniffle	0
...	

$$f(x) = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 0 \\ \dots \end{bmatrix}$$

# Example Classification Tasks



GLUE

<https://gluebenchmark.com/>

🤖 datasets: glue

## GLUE Tasks

Name	Download
The Corpus of Linguistic Acceptability	<a href="#">Download</a>
The Stanford Sentiment Treebank	<a href="#">Download</a>
Microsoft Research Paraphrase Corpus	<a href="#">Download</a>
Semantic Textual Similarity Benchmark	<a href="#">Download</a>
Quora Question Pairs	<a href="#">Download</a>
MultiNLI Matched	<a href="#">Download</a>
MultiNLI Mismatched	<a href="#">Download</a>
Question NLI	<a href="#">Download</a>
Recognizing Textual Entailment	<a href="#">Download</a>
Winograd NLI	<a href="#">Download</a>
Diagnostics Main	<a href="#">Download</a>

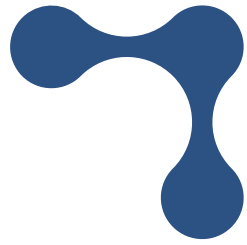
## SuperGLUE 1

Name	Identifier
Broadcoverage Diagnostics	AX-b
CommitmentBank	CB
Choice of Plausible Alternatives	COPA
Multi-Sentence Reading Comprehension	MultiRC
Recognizing Textual Entailment	RTE
Words in Context	WiC
The Winograd Schema Challenge	WSC
BoolQ	BoolQ
Reading Comprehension with Commonsense Reasoning	ReCoRD
Winogender Schema Diagnostics	AX-g

 **SuperGLUE**

<https://super.gluebenchmark.com/>

🤖 datasets: super\_glue



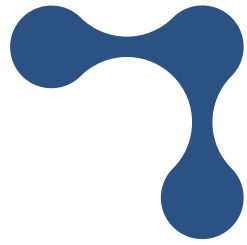
# Recognizing Textual Entailment (RTE)

---

Given a premise sentence  $s$  and hypothesis sentence  $h$ ,  
determine if  $h$  “follows from”  $s$

ENTAILMENT (yes):

NOT ENTAILED (no):



# Recognizing Textual Entailment (RTE)

---

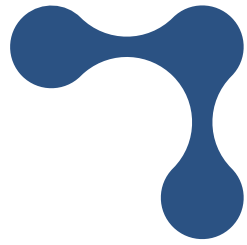
Given a premise sentence  $s$  and hypothesis sentence  $h$ ,  
determine if  $h$  “follows from”  $s$

ENTAILMENT (yes):

$s$ : Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the Chicago Bulls to six National Basketball Association championships.

$h$ : The Bulls basketball team is based in Chicago.

NOT ENTAILED (no):



# Recognizing Textual Entailment (RTE)

---

Given a premise sentence  $s$  and hypothesis sentence  $h$ , determine if  $h$  “follows from”  $s$

ENTAILMENT (yes):

$s$ : Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the Chicago Bulls to six National Basketball Association championships.

$h$ : The Bulls basketball team is based in Chicago.

NOT ENTAILED (no):

$s$ : Based on a worldwide study of smoking-related fire and disaster data, UC Davis epidemiologists show smoking is a leading cause of fires and death from fires globally.

$h$ : Domestic fires are the major cause of fire death.



# RTE

---

s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the Chicago Bulls to six National Basketball Association championships.

h: The Bulls basketball team is based in Chicago.

ENTAILED

p(

ENTAILED

|

s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the Chicago Bulls to six National Basketball Association championships.

h: The Bulls basketball team is based in Chicago.

)

# Discriminative Document Classification

---

s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the Chicago Bulls to six National Basketball Association championships.

ENTAILED

h: The Bulls basketball team is based in Chicago.

# Discriminative Document Classification

---

s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the **Chicago** Bulls to six National Basketball Association championships.

h: The Bulls basketball team is based in **Chicago**.

ENTAILED

These extractions are all **features** that have **fired** (likely have some significance)

# Discriminative Document Classification

---

s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the **Chicago Bulls** to six National Basketball Association championships.

h: The **Bulls** basketball team is based in **Chicago**.

ENTAILED

These extractions are all **features** that have **fired** (likely have some significance)

# Discriminative Document Classification

---

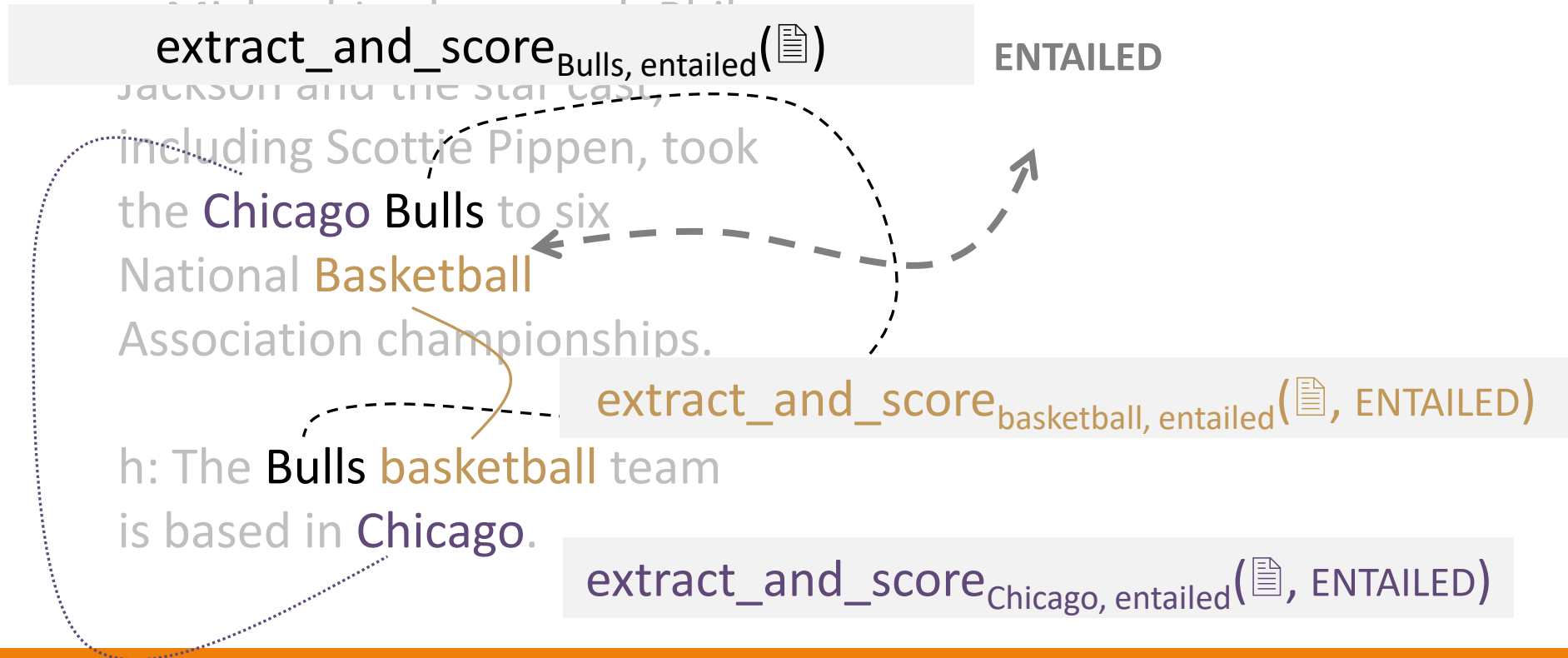
s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the **Chicago Bulls** to six National **Basketball** Association championships.

h: The **Bulls basketball** team is based in **Chicago**.

ENTAILED

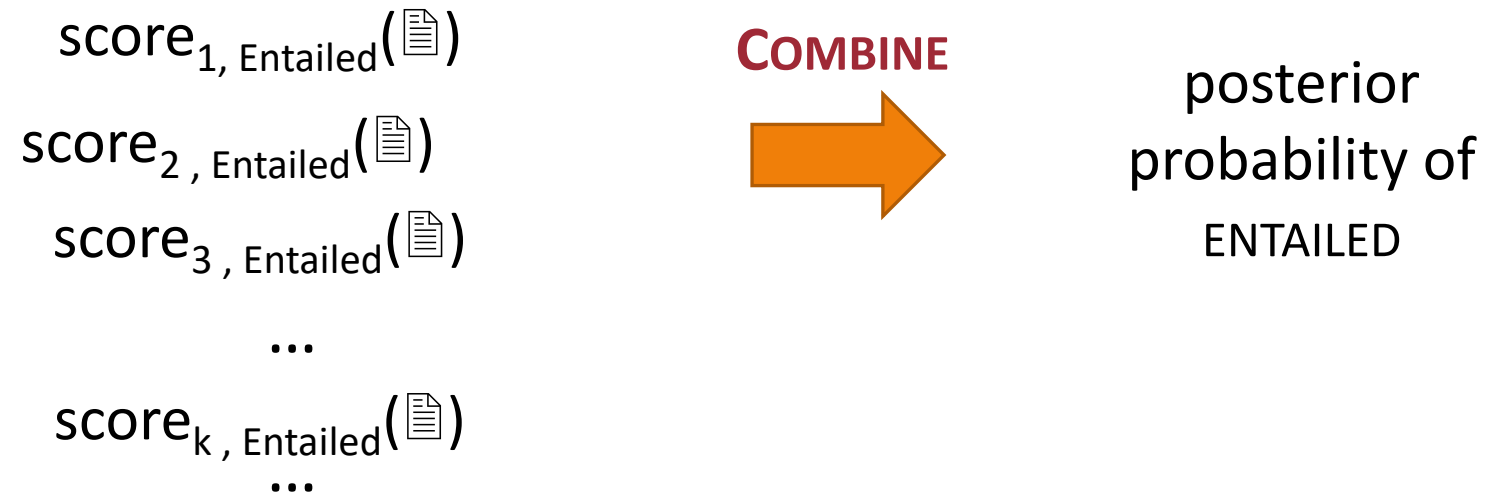
These extractions are all **features** that have **fired** (likely have some significance)

We need to *score* the different extracted clues.



# Score and Combine Our Clues

---



# Scoring Our Clues

score( , ENTAILED ) =

s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the Chicago Bulls to six National Basketball Association championships.

h: The Bulls basketball team is based in Chicago.

*(ignore the  
feature indexing  
for now)*

score<sub>1</sub> , Entailed (📄)

+

score<sub>2</sub> , Entailed (📄)

+

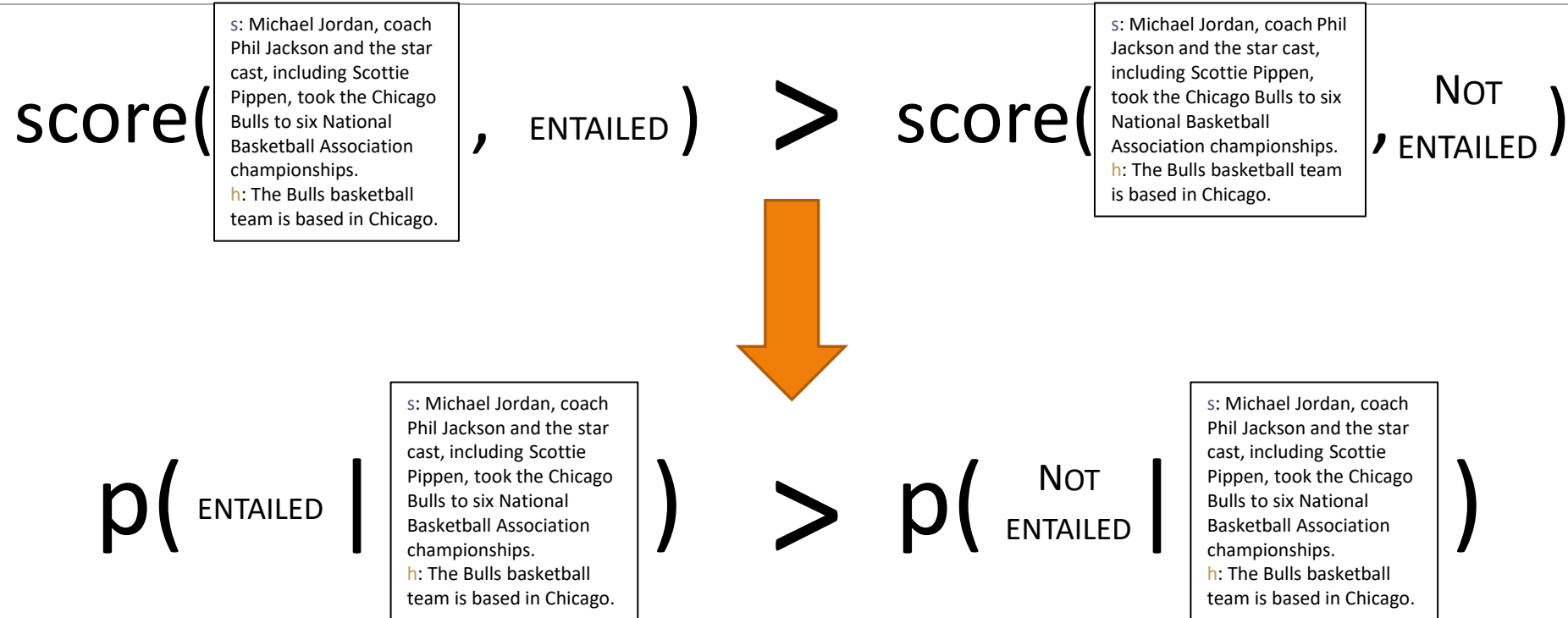
score<sub>3</sub> , Entailed (📄)

+

...



# Turning Scores into Probabilities



KEY IDEA

# Turning Scores into Probabilities (More Generally)

---

$$\text{score}(x, y_1) > \text{score}(x, y_2)$$



$$p(y_1 | x) > p(y_2 | x)$$

KEY IDEA