

ML Evaluation → Classification

CMSC 473/673 - NATURAL LANGUAGE PROCESSING

Slides modified from Dr. Frank Ferraro

Learning Objectives

Develop an intuition about precision & recall









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

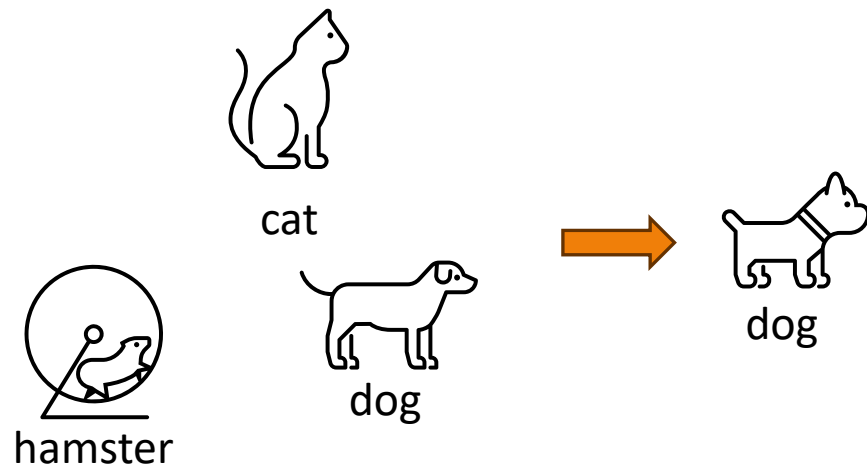
Review: Classification Evaluation: the 2-by-2 contingency table

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

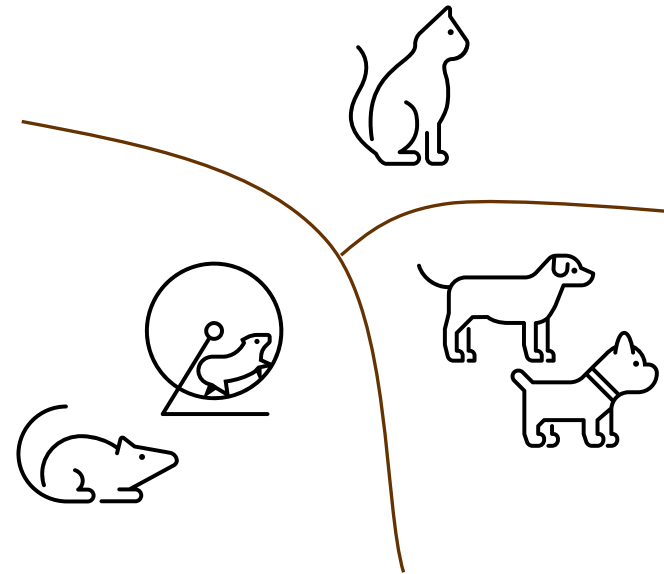
Construct this table by *counting*
the number of TPs, FPs, FNs, TNs

Review: Types of Learning

SUPERVISED LEARNING



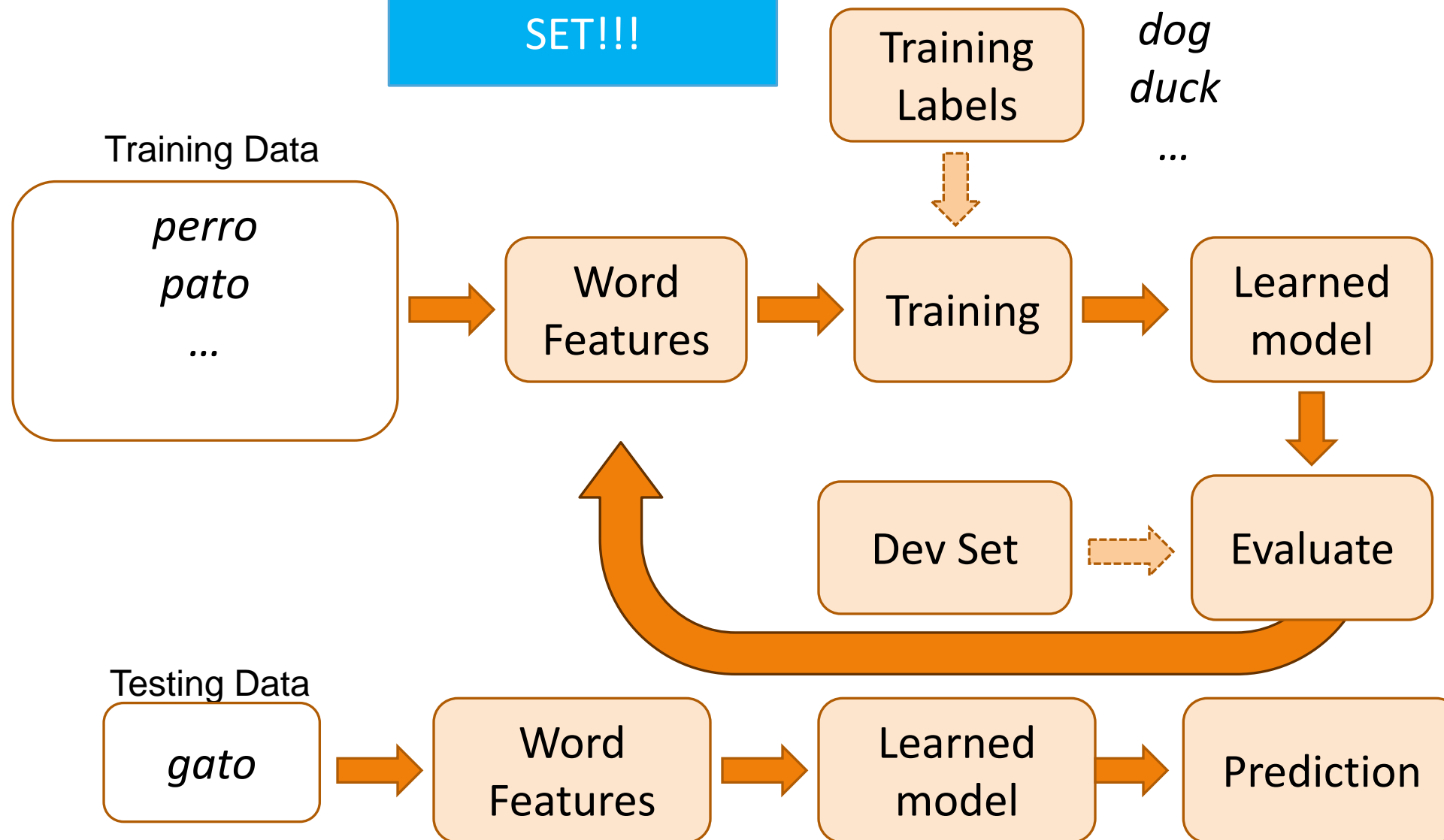
UNSUPERVISED LEARNING



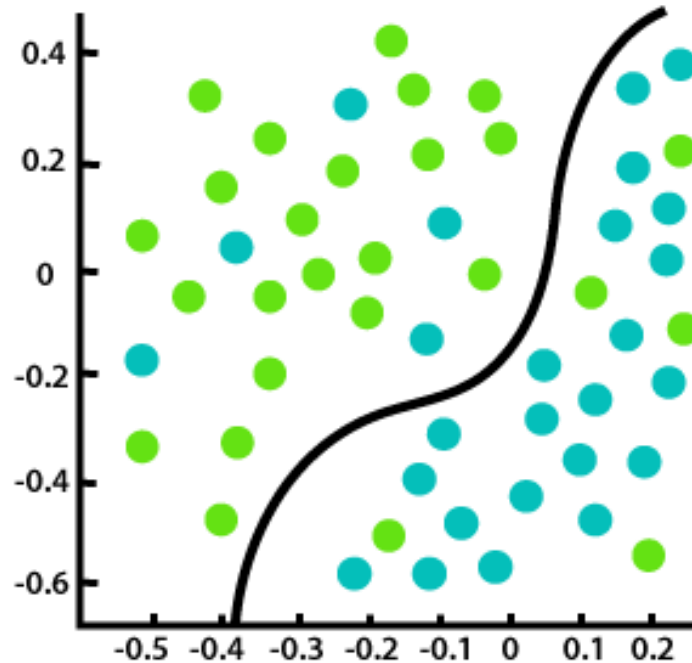
Review: Steps

DO NOT ITERATE
ON THE TESTING
SET!!!

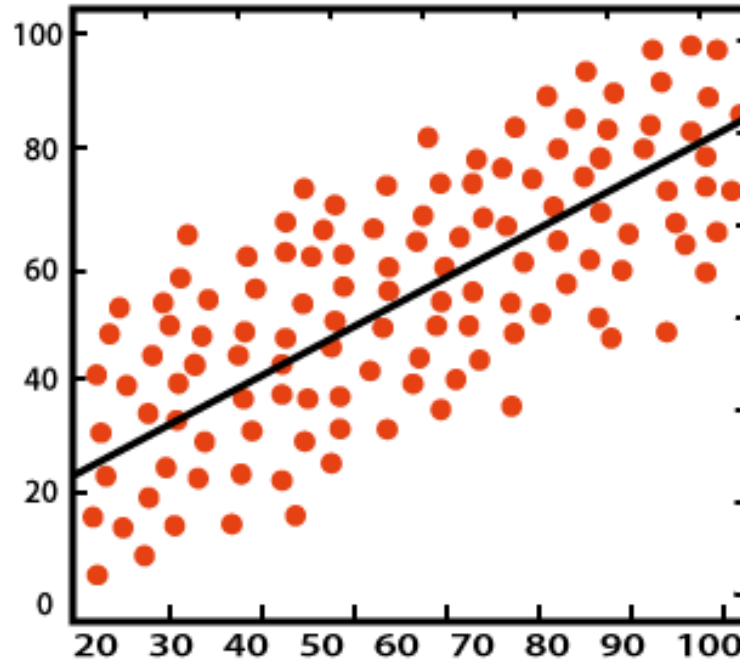
Training



Review: Types of models











Classification



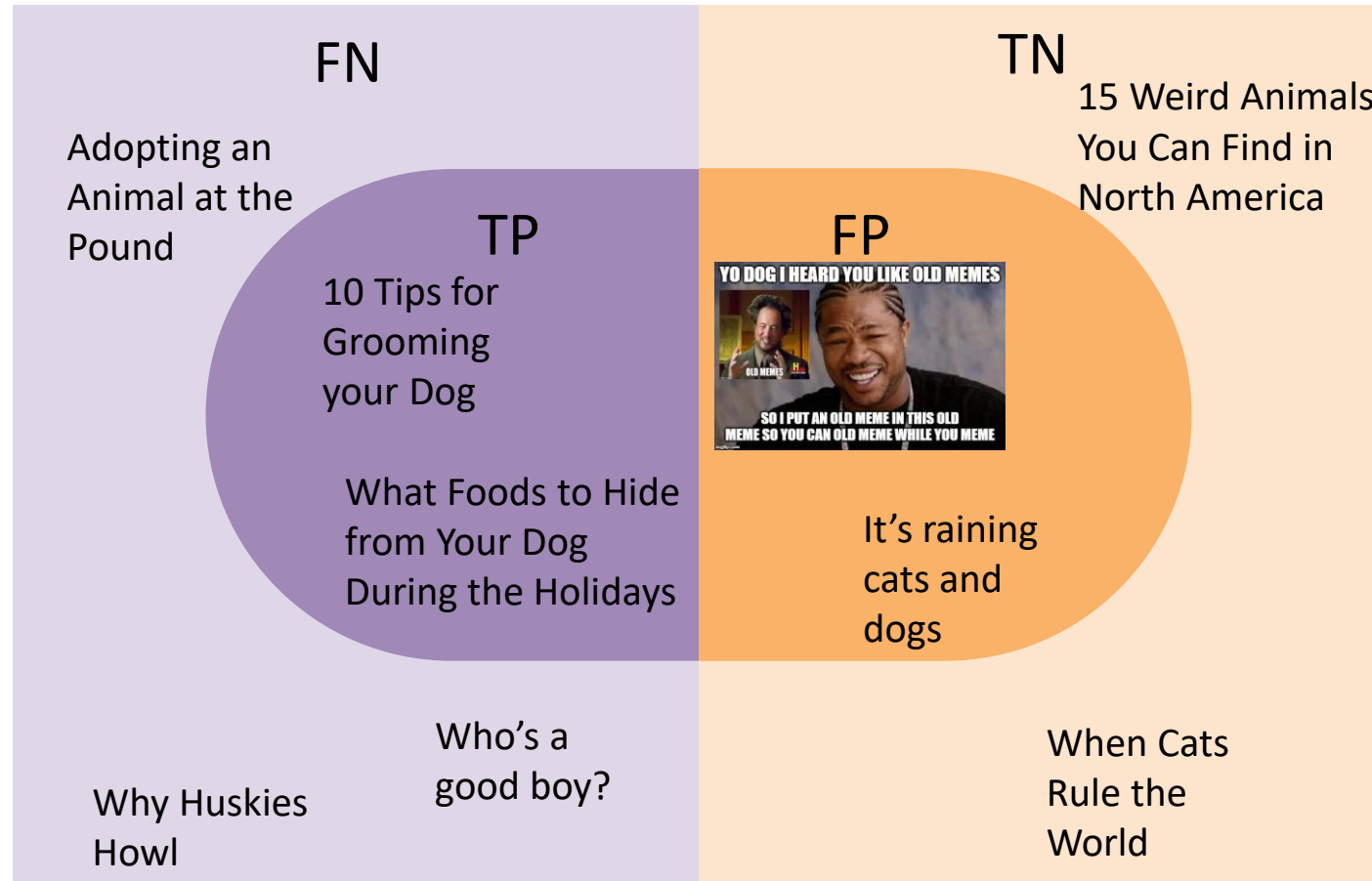
Regression

Review: Classification Evaluation: the 2-by-2 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)  <i>Actual</i>  <i>Guessed</i>	False Positive (FP)  <i>Actual</i>  <i>Guessed</i>
Not selected/ not guessed (“○”)	False Negative (FN)  <i>Actual</i>  <i>Guessed</i>	True Negative (TN)  <i>Actual</i>  <i>Guessed</i>

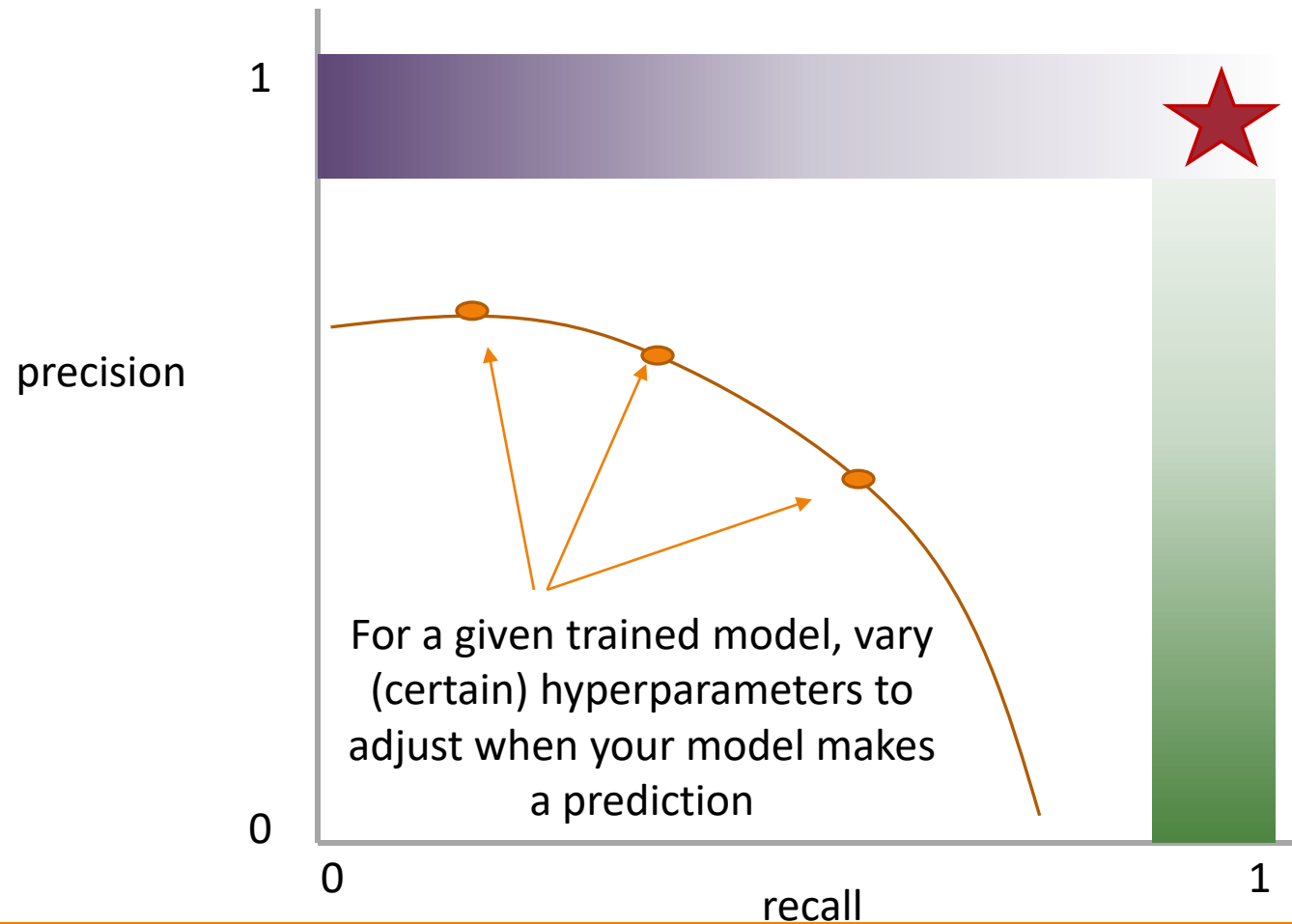
Contingency Table (out of table form)

Query:
Articles about
dogs



Meme from: https://www.reddit.com/r/AdviceAnimals/comments/ck8xh0/yo_dawg_i_heard_you_like_old_memes/

Review: Precision and Recall Present a Tradeoff



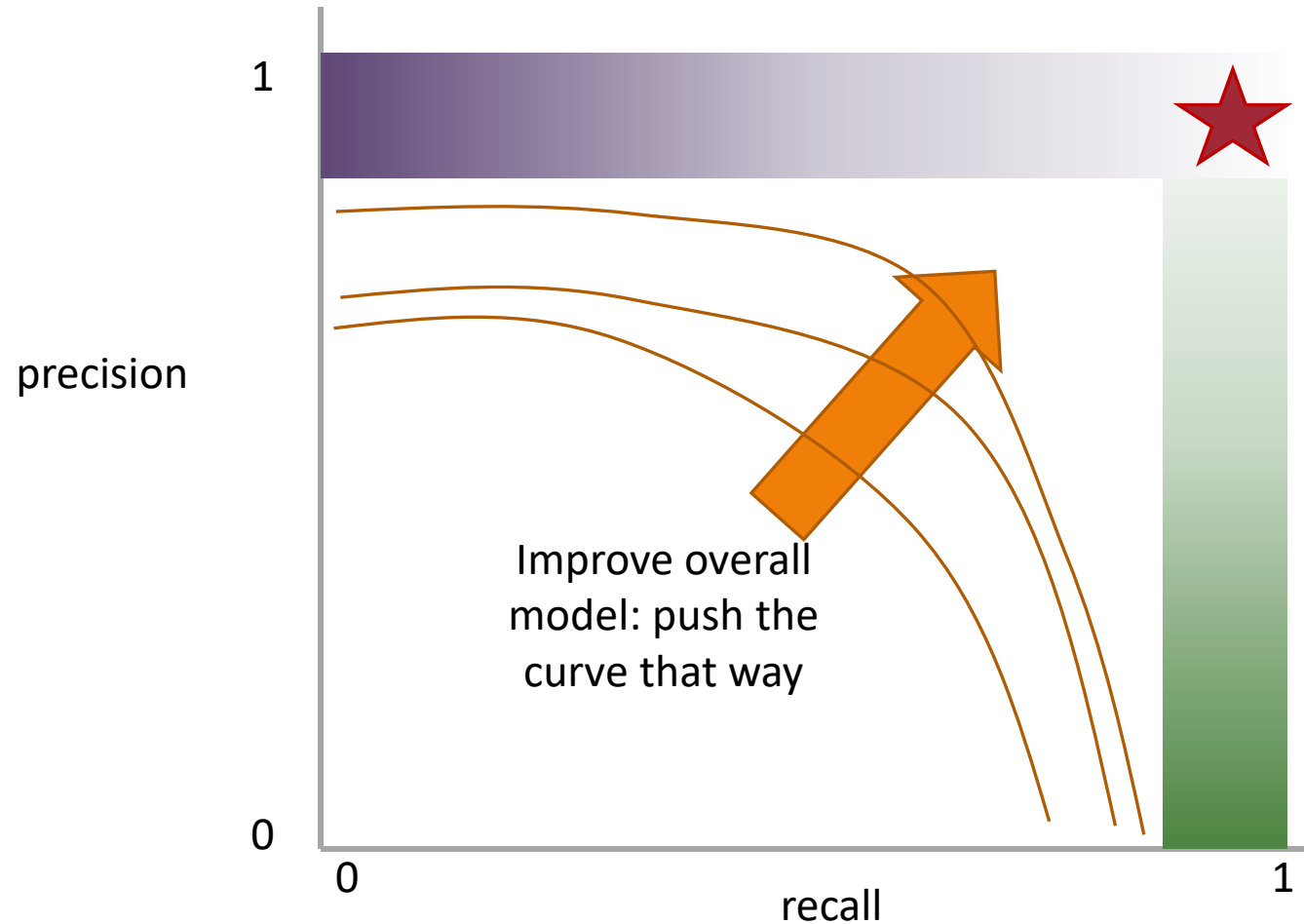
Q: Where do you want your ideal model ?

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

Review: 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

Review: A combined measure: 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 $P = R = 0$)

Classification Evaluation: Accuracy, Precision, and Recall

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

Macroaveraging: Compute performance for each class, then average.

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

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

when to prefer
macroaveraging?

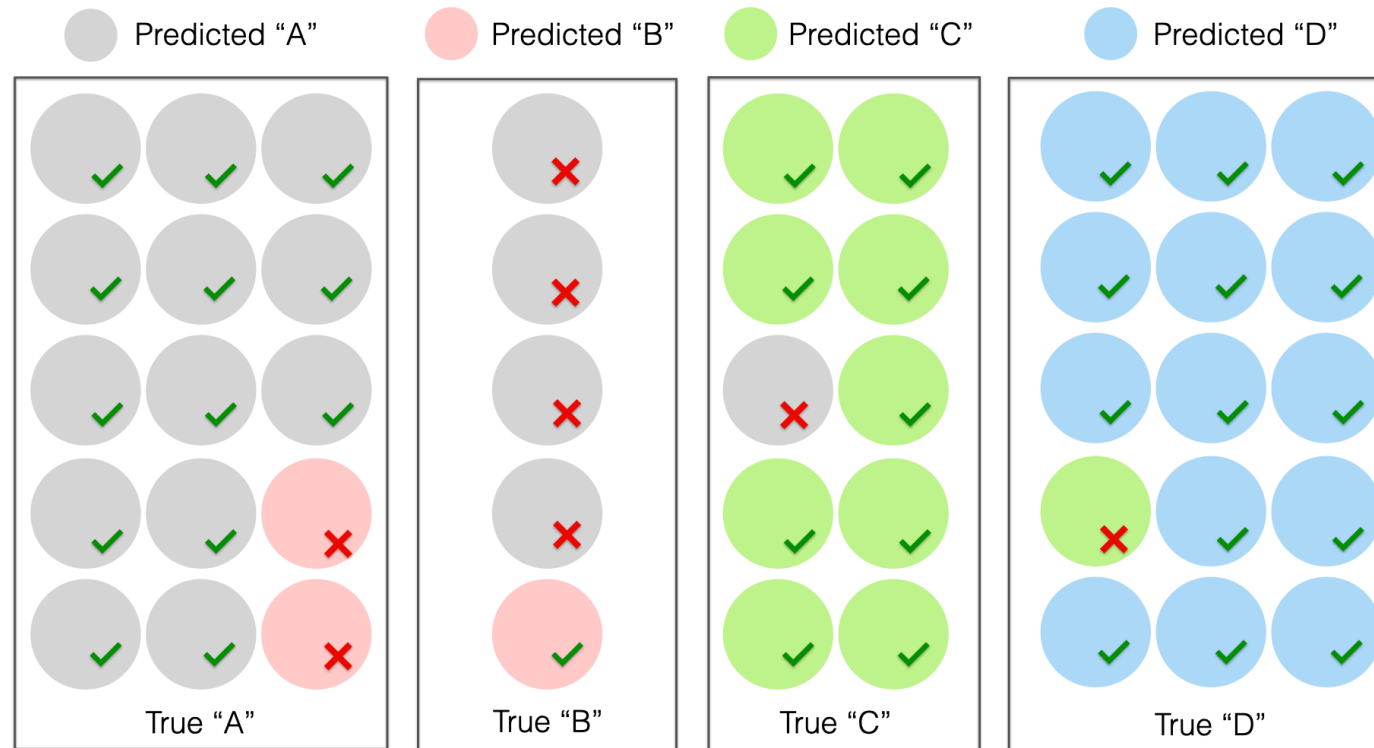
Microaveraging: Collect decisions for all classes, compute contingency table, evaluate.

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

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

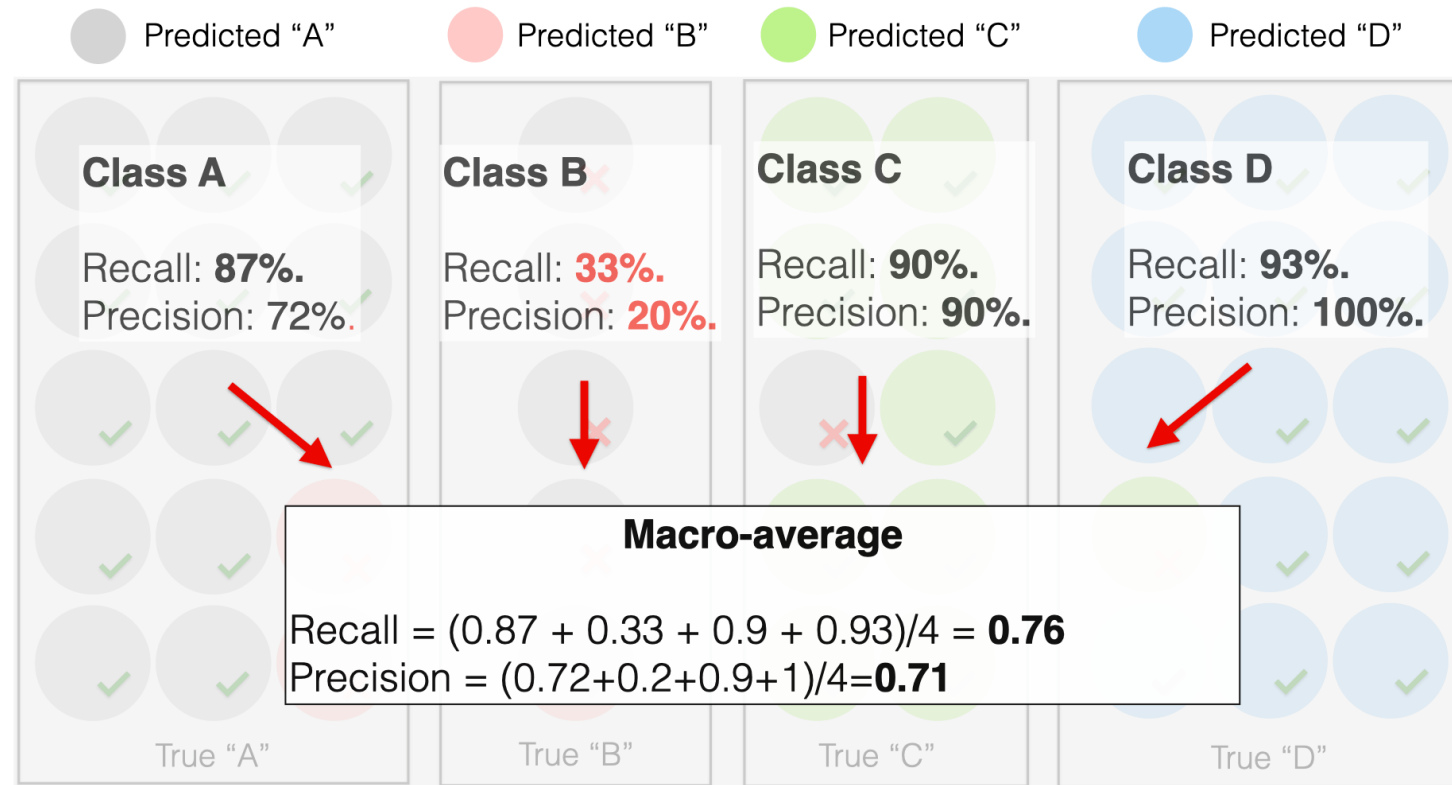
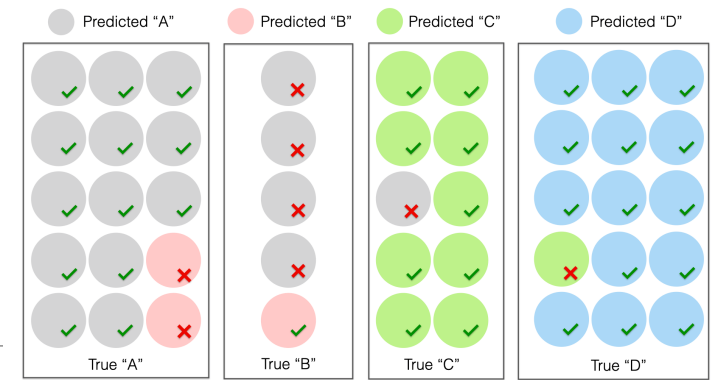
when to prefer
microaveraging?

Macro/Micro Example



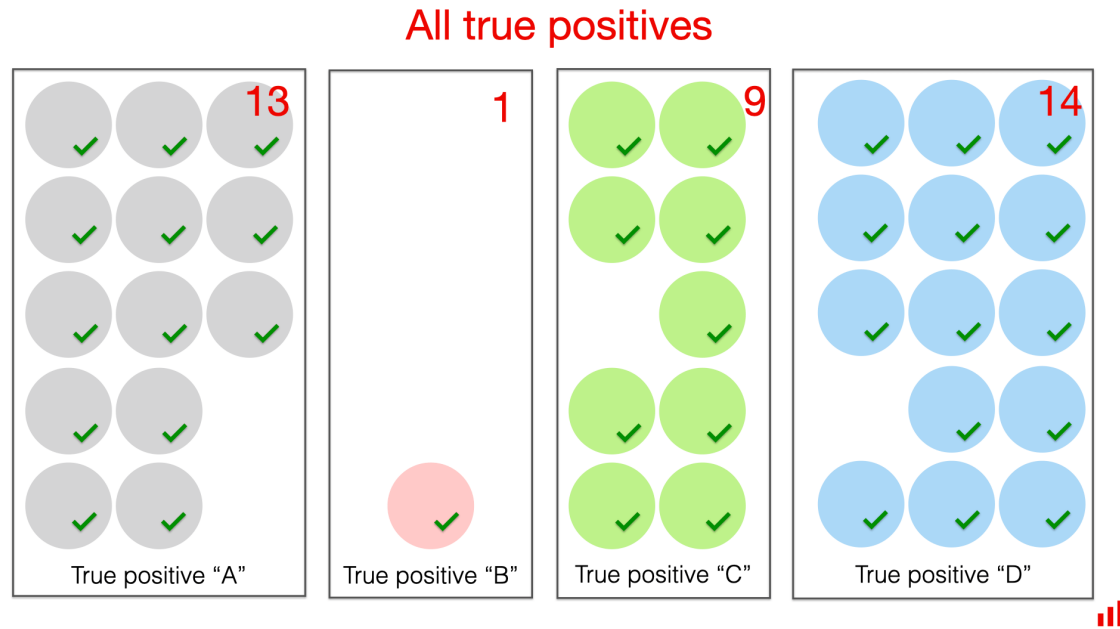
Each *class* has equal weight

Macro-Average



Each *instance* has equal weight

Micro-Average

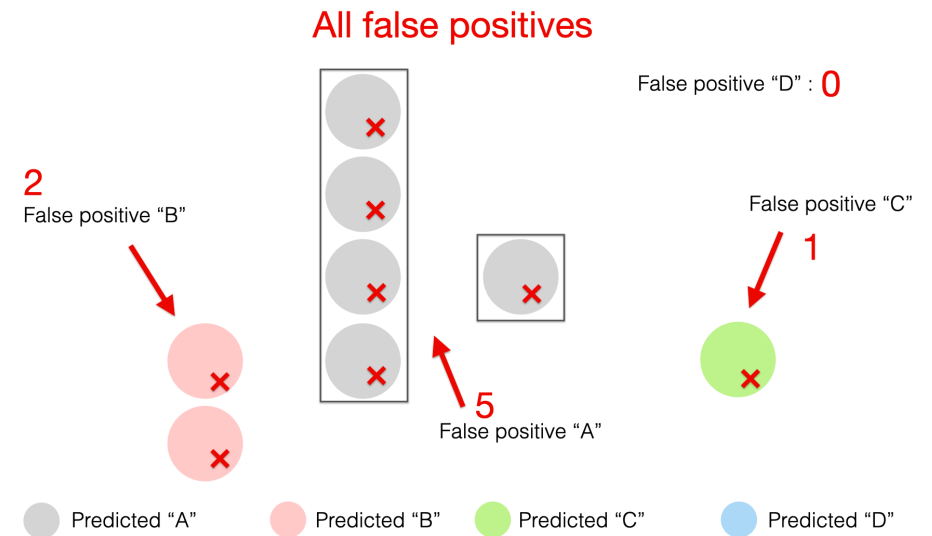
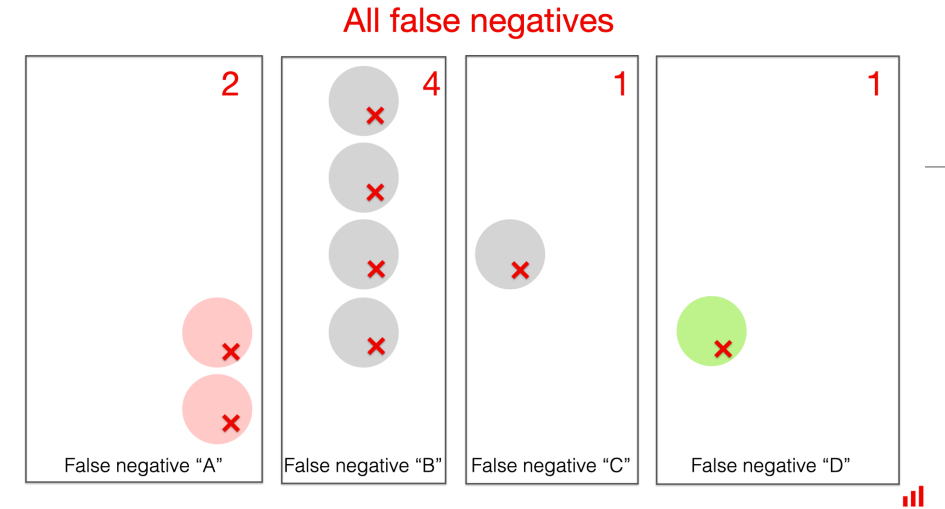


Total TP: 13 + 1 + 9 + 14 = 37

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

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

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



Micro- vs Macro-Average

So when would we want to prefer micro-averaging vs macro-averaging?

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

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

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

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

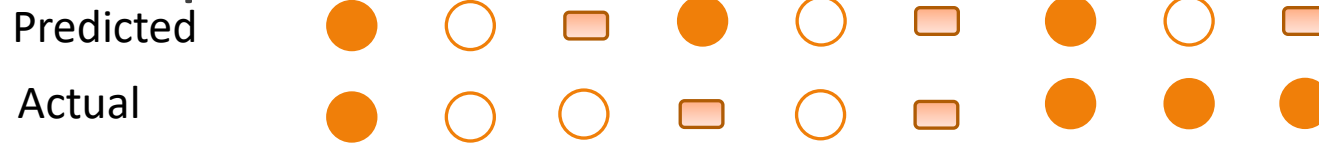
But how do we compute stats for multiple classes?



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



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




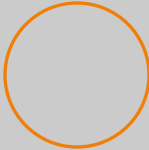


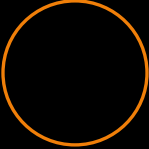

Actual



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

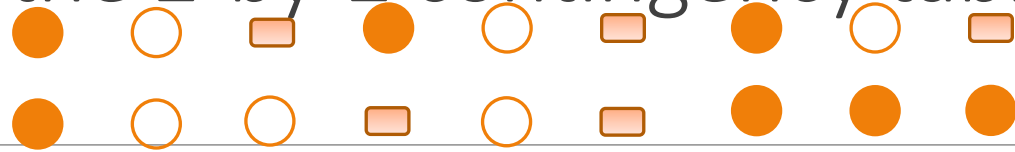
		Correct Value		
				
Guessed Value		#	#	#
		#	#	#
		#	#	#

This is also called a **Confusion Matrix**

2. Generalizing the 2-by-2 contingency table

Predicted

Actual



		Correct Value		
		●	○	□
Guessed Value	●	#	#	#
	○	#	#	#
	□	#	#	#


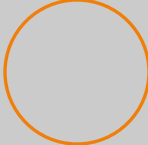


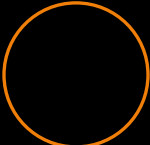

2. Generalizing the 2-by-2 contingency table

Predicted



Actual



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





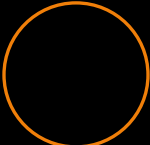

2. Generalizing the 2-by-2 contingency table

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} ?


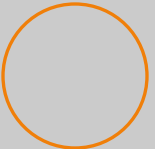


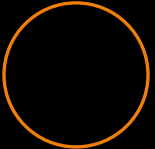

2. Generalizing the 2-by-2 contingency table

Predicted



Actual



		Correct Value		
				
Guessed Value		A 2	B 0	C 1
		D 1	E 2	F 0
		G 1	H 1	I 1

How do you compute TP_{\bullet} ?


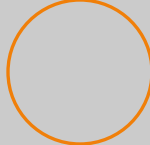


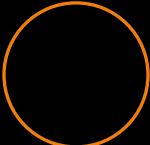

2. Generalizing the 2-by-2 contingency table

Predicted



Actual



		Correct Value		
				
Guessed Value		A 2	B 0	C 1
		D 1	E 2	F 0
		G 1	H 1	I 1

How do you compute FN_{\bullet} ?


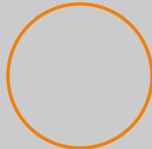


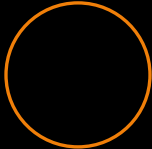

2. Generalizing the 2-by-2 contingency table

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_{\bullet} ?


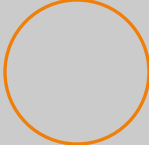


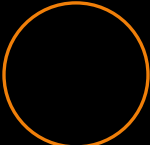

2. Generalizing the 2-by-2 contingency table

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 FP_{\square} ?


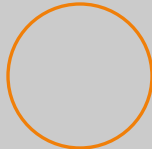


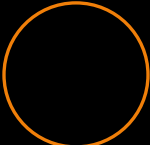

2. Generalizing the 2-by-2 contingency table

Predicted




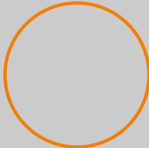


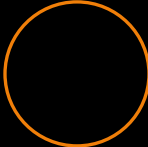

Actual




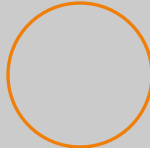


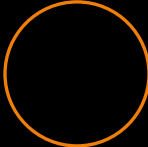

		Correct Value		
				
Guessed Value		A 2	B 0	C 1
		D 1	E 2	F 0
		G 1	H 1	I 1

How do you compute $FP_{\text{orange rectangle}}$?


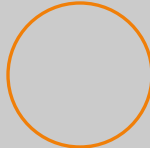


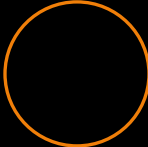

Generalizing the 2-by-2 contingency table

Q: Is this a good result?		Correct Value		
				
Guessed Value		80	9	11
		7	86	7
		2	8	9

Generalizing the 2-by-2 contingency table

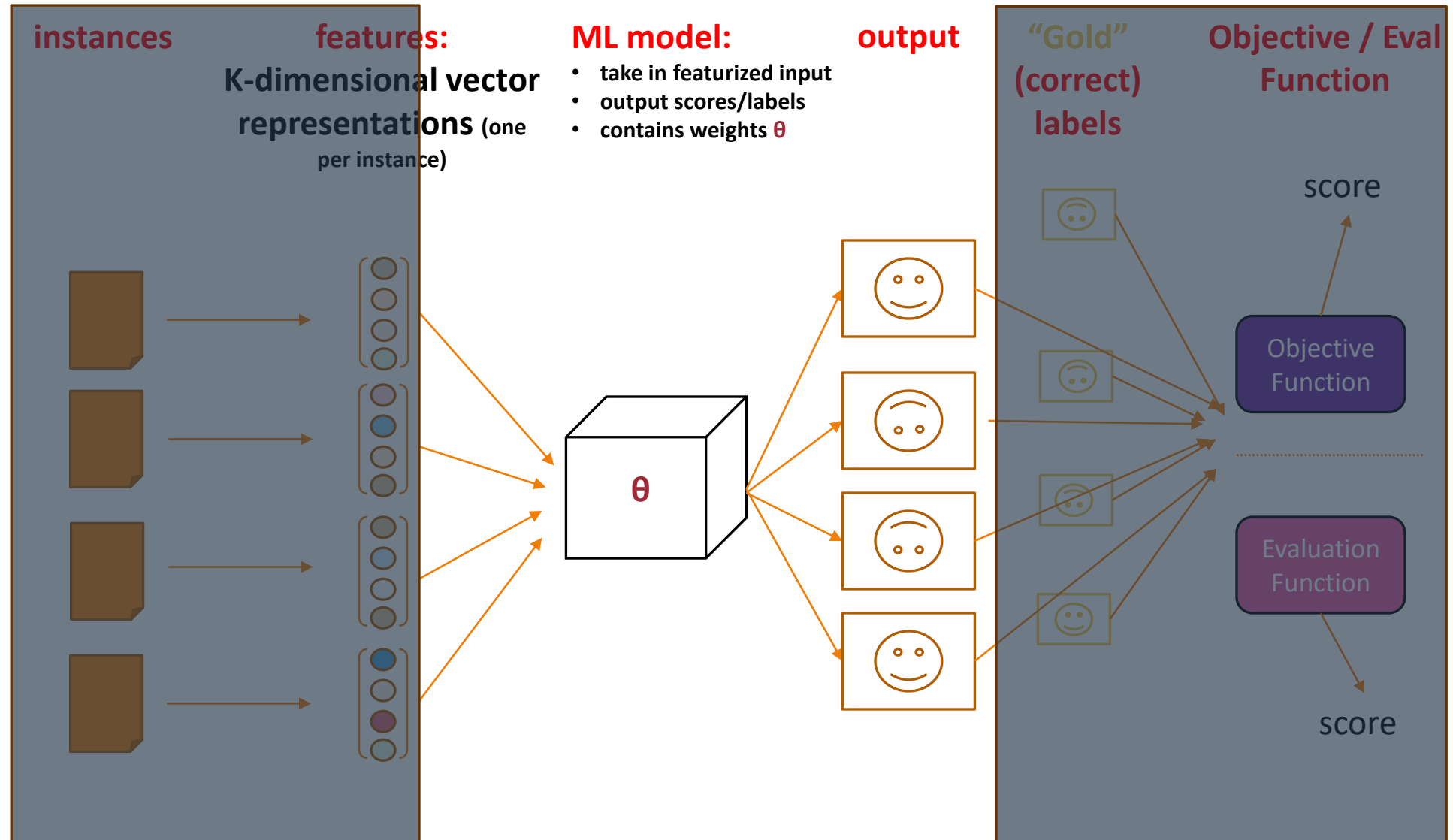
Q: Is this a good result?		Correct Value		
				
Guessed Value		30	40	30
		25	30	50
		30	35	35

Generalizing the 2-by-2 contingency table

Q: Is this a good result?		Correct Value		
				
Guessed Value		7	3	90
		4	8	88
		3	7	90

Classification

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

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

Normally (*but this can be adjusted!)


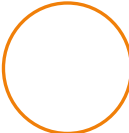
$$\text{best label} = \arg \max_{\text{label}} P(\text{label} | \text{example})$$

Terminology: Posterior Probability

Posterior probability:

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

These are conditional probabilities

- If  and  are the only two options:

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

and

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

Bayes' Rule

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

Posterior: $P(Y|X)$

Posterior probability:
probability of event Y
with knowledge that X
has occurred

NLP pg. 450

Terminology (with variables)

Posterior probability:

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

Conditional probabilities:

$$p(Y = \text{label}_1 | X) + p(Y = \text{label}_0 | X) = 1$$

$$\begin{aligned} p(Y = \text{label}_1 | X) &\geq 0, \\ p(Y = \text{label}_0 | X) &\geq 0 \end{aligned}$$

Conditional probability:

probability of event Y,
assuming event X
happens too

NLP pg. 449



Key Take-away



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

Maxent Models for Classification: Discriminatively or ...

Directly model
the posterior

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

Discriminatively trained classifier

Maxent Models for Classification: Discriminatively or Generatively Trained

Directly model
the posterior

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

Discriminatively trained classifier



Model the
posterior with
Bayes rule

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

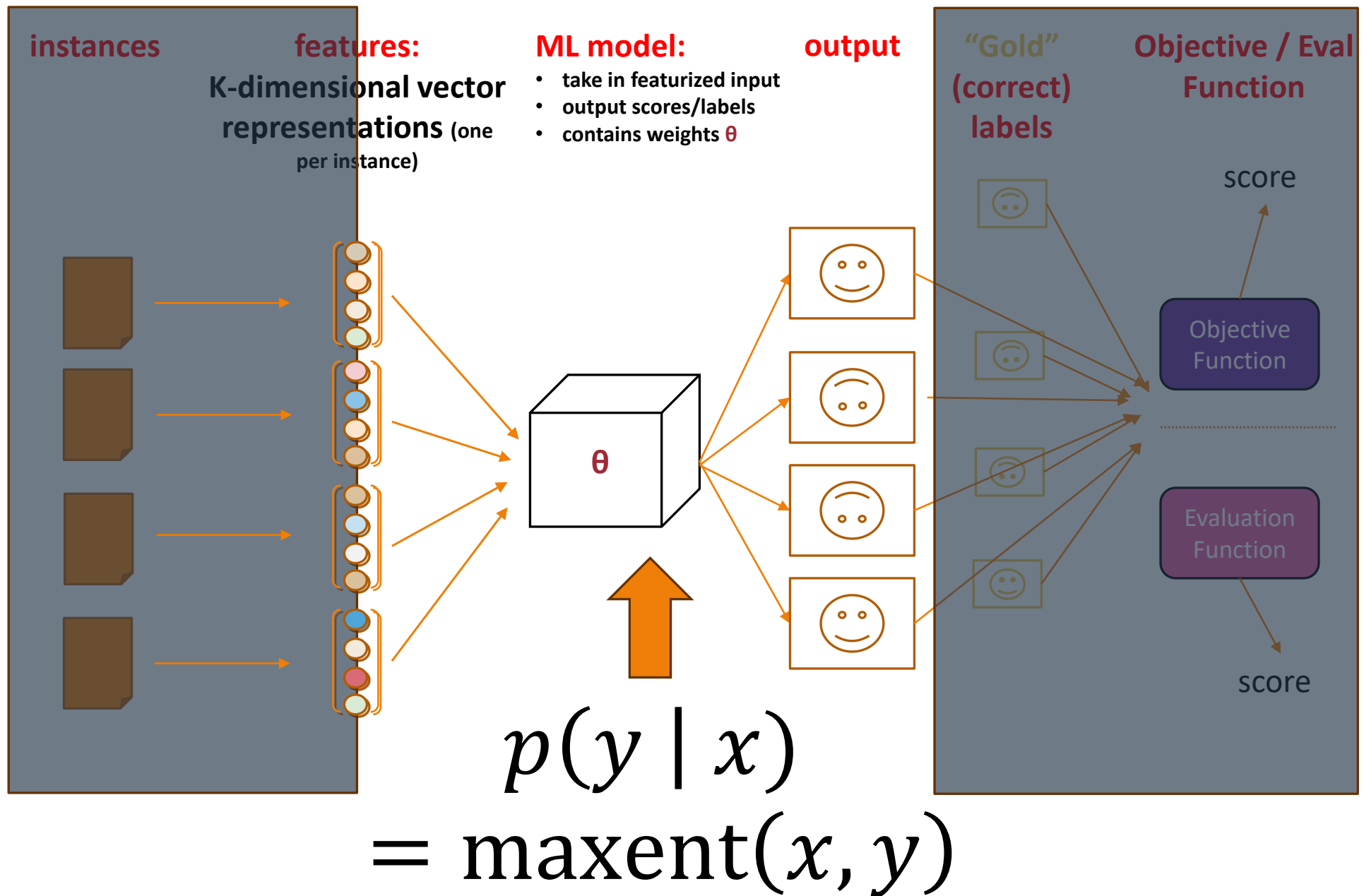
Generatively trained classifier with
maxent-based language model

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

(we'll start with this one)

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

discriminatively trained:
classify in one go



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

We need to define:

- **features** $f(x)$ from x that are meaningful;
- **weights** θ (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 θ

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”

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

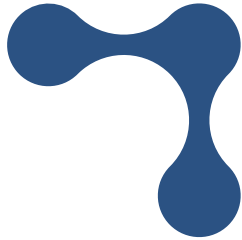
$$f(x) = (f_i(x))_i^V$$

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
...	

Example Classification Tasks



GLUE

<https://gluebenchmark.com/>

🤖 datasets: glue

GLUE Tasks

Name	Download
The Corpus of Linguistic Acceptability	Download
The Stanford Sentiment Treebank	Download
Microsoft Research Paraphrase Corpus	Download
Semantic Textual Similarity Benchmark	Download
Quora Question Pairs	Download
MultiNLI Matched	Download
MultiNLI Mismatched	Download
Question NLI	Download
Recognizing Textual Entailment	Download
Winograd NLI	Download
Diagnostics Main	Download

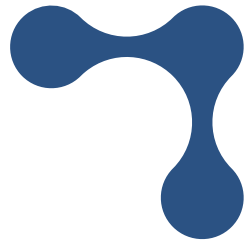
SuperGLUE 1

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



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

🤖 datasets: super_glue

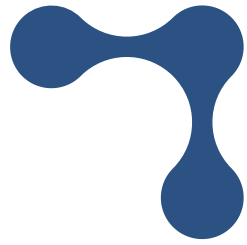


Recognizing Textual Entailment (RTE)

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

ENTAILMENT (yes):

NOT ENTAILED (no):



Recognizing Textual Entailment (RTE)

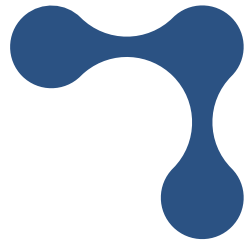
Given a premise sentence s and hypothesis sentence h ,
determine if h “follows from” s

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

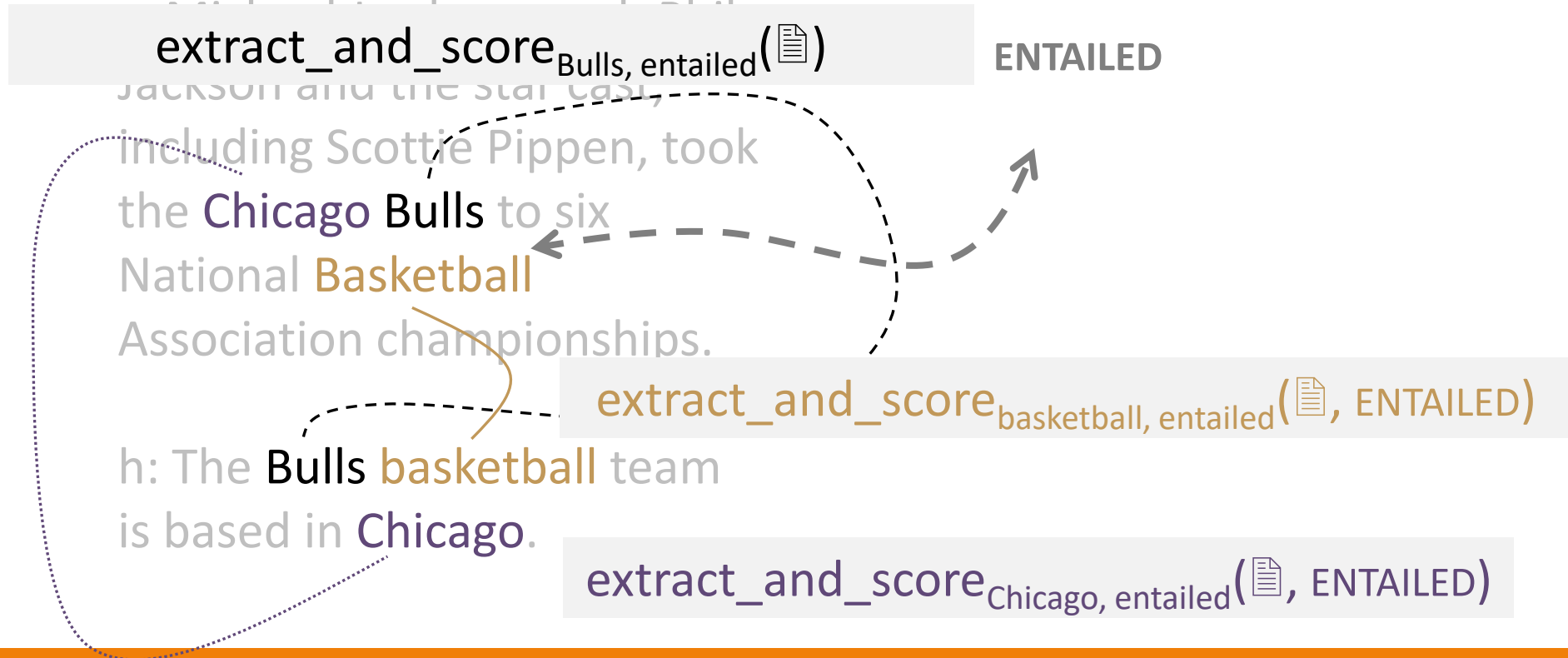
s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the **Chicago Bulls** to six National **Basketball** Association championships.

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

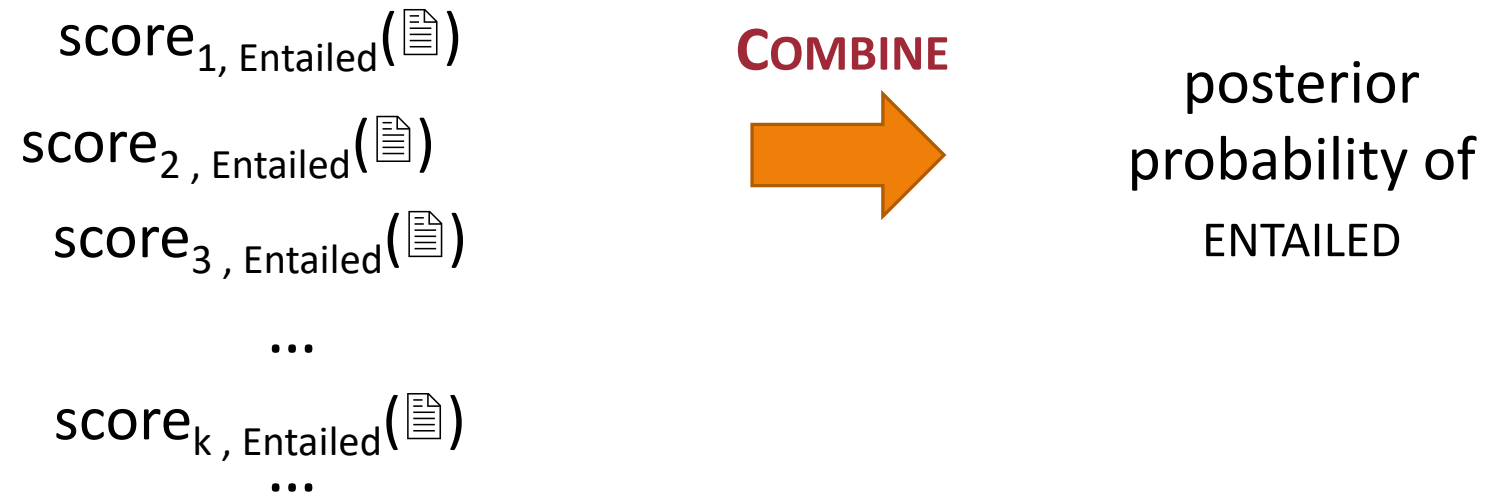
ENTAILED

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)*

score₁ , Entailed (📄)

+

score₂ , Entailed (📄)

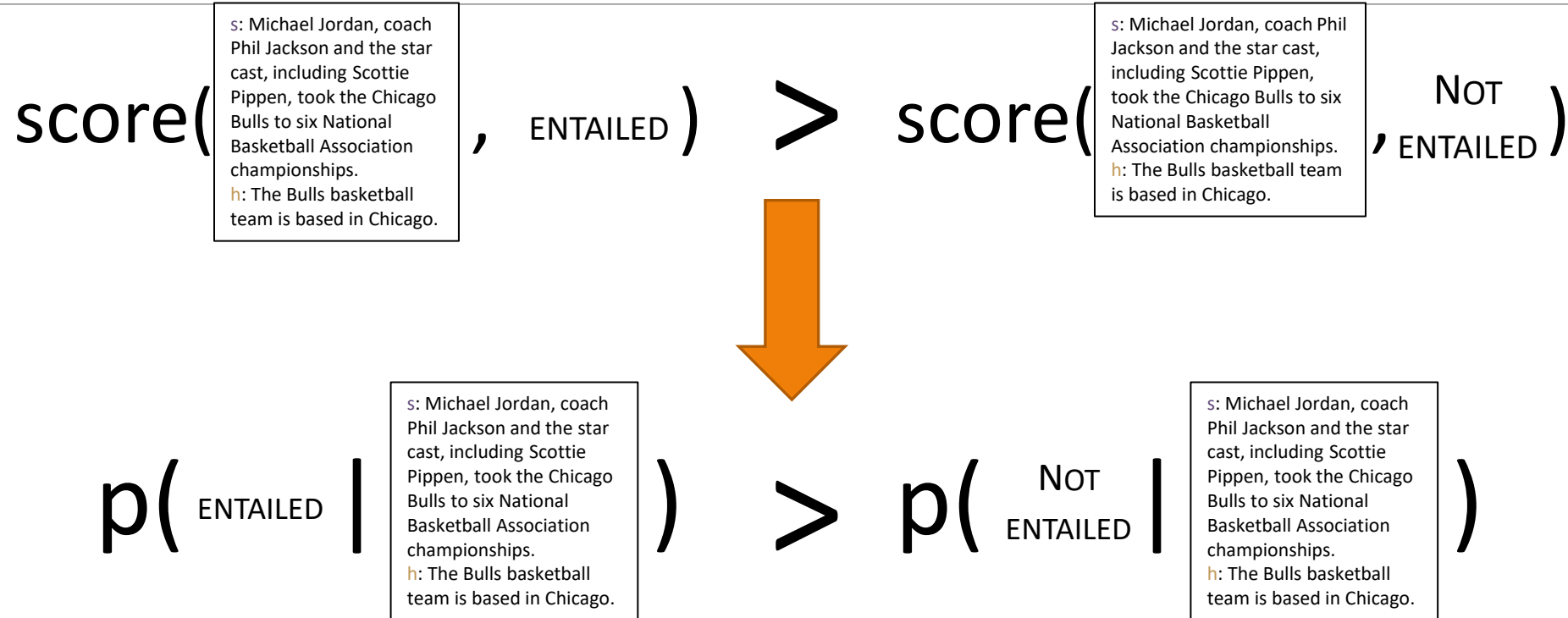
+

score₃ , 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

Maxent Modeling

$p(\text{ENTAILED} \mid \text{s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the Chicago Bulls to six National Basketball Association championships. h: The Bulls basketball team is based in Chicago.}) \propto$

Convert through
function G ?

What is this
function?

This must be a probability

This could be any real number

$G(\text{score}(\text{s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the Chicago Bulls to six National Basketball Association championships. h: The Bulls basketball team is based in Chicago.}, \text{ENTAILED}))$

What function G...

operates on any real number?

is never less than 0?

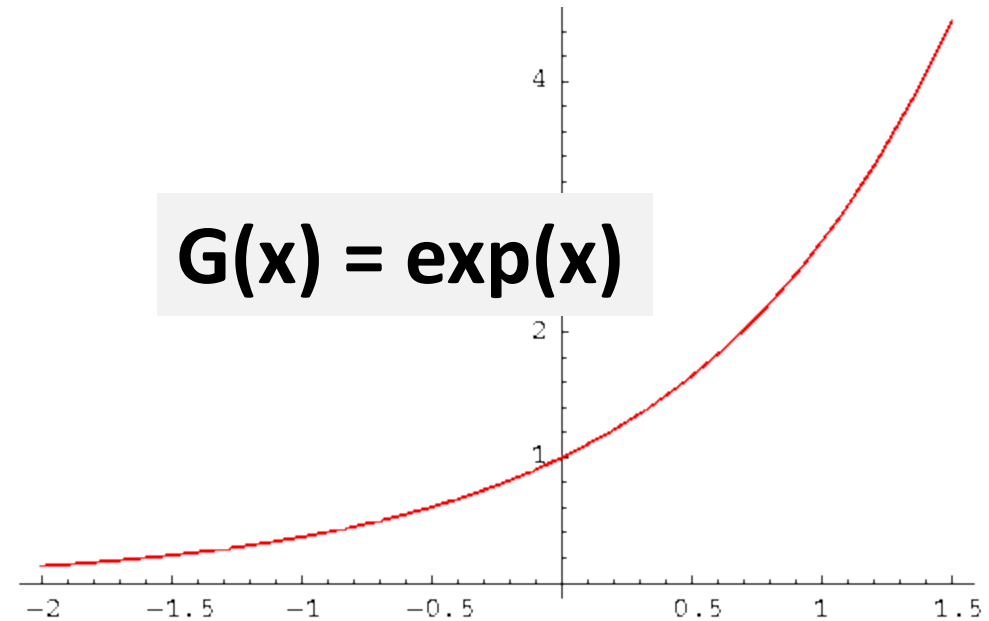
monotonic? ($a < b \rightarrow G(a) < G(b)$)

What function G...

operates on any real number?

is never less than 0?

monotonic? ($a < b \rightarrow G(a) < G(b)$)



Maxent Modeling

$$p(\text{ENTAILED} \mid \text{s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the Chicago Bulls to six National Basketball Association championships. h: The Bulls basketball team is based in Chicago.}) \propto \exp(\text{score}(\text{ENTAILED}))$$

Maxent Modeling

$$p(\text{ENTAILED} \mid \text{...}) \propto$$

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.

$$\exp(\text{score}_{1, \text{Entailed}} + \text{score}_{2, \text{Entailed}} + \text{score}_{3, \text{Entailed}} + \dots)$$

Maxent Modeling

$p(\text{ENTAILED} \mid \text{...}) \propto$

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.

$\exp(\text{weight}_{1, \text{Entailed}} * \text{applies}_1(\text{...}) + \text{weight}_{2, \text{Entailed}} * \text{applies}_2(\text{...}) + \text{weight}_{3, \text{Entailed}} * \text{applies}_3(\text{...}) + \dots)$

Maxent Modeling

$$p(\text{ENTAILED} \mid \text{...}) \propto$$

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.

$$\exp\left(\begin{array}{l} \text{weight}_{1, \text{Entailed}} * \text{applies}_1(\text{...}) + \\ \text{weight}_{2, \text{Entailed}} * \text{applies}_2(\text{...}) + \\ \text{weight}_{3, \text{Entailed}} * \text{applies}_3(\text{...}) + \\ \vdots \end{array}\right)$$

K different
weights...

for K different
features

Maxent Modeling

$p(\text{ENTAILED} \mid \text{...}) \propto$

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.

$\exp(\text{weight}_{1, \text{Entailed}} * \text{applies}_1(\text{...}) + \text{weight}_{2, \text{Entailed}} * \text{applies}_2(\text{...}) + \text{weight}_{3, \text{Entailed}} * \text{applies}_3(\text{...}) + \dots)$

K different
weights...

for K different
features

multiplied and then summed

Maxent Modeling

$p(\text{ENTAILED} | \text{...}) \propto$

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.

$\exp(\text{Dot_product of Entailed weight_vec feature_vec}(\text{📄}))$

K different
weights...

for K different
features

multiplied and
then summed

Maxent Modeling

$p(\text{ENTAILED} | \text{ }) \propto$

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.

$\exp(\theta_{\text{ENTAILED}}^T f(\text{document}))$

K different weights... for K different features multiplied and then summed



Implementation Check

$$\exp\left(\theta_{\text{ENTAILED}}^T f(\text{document})\right)$$

Assume we can compute $\mathbf{x} = f(\text{document})$
(a one-dimensional tensor)

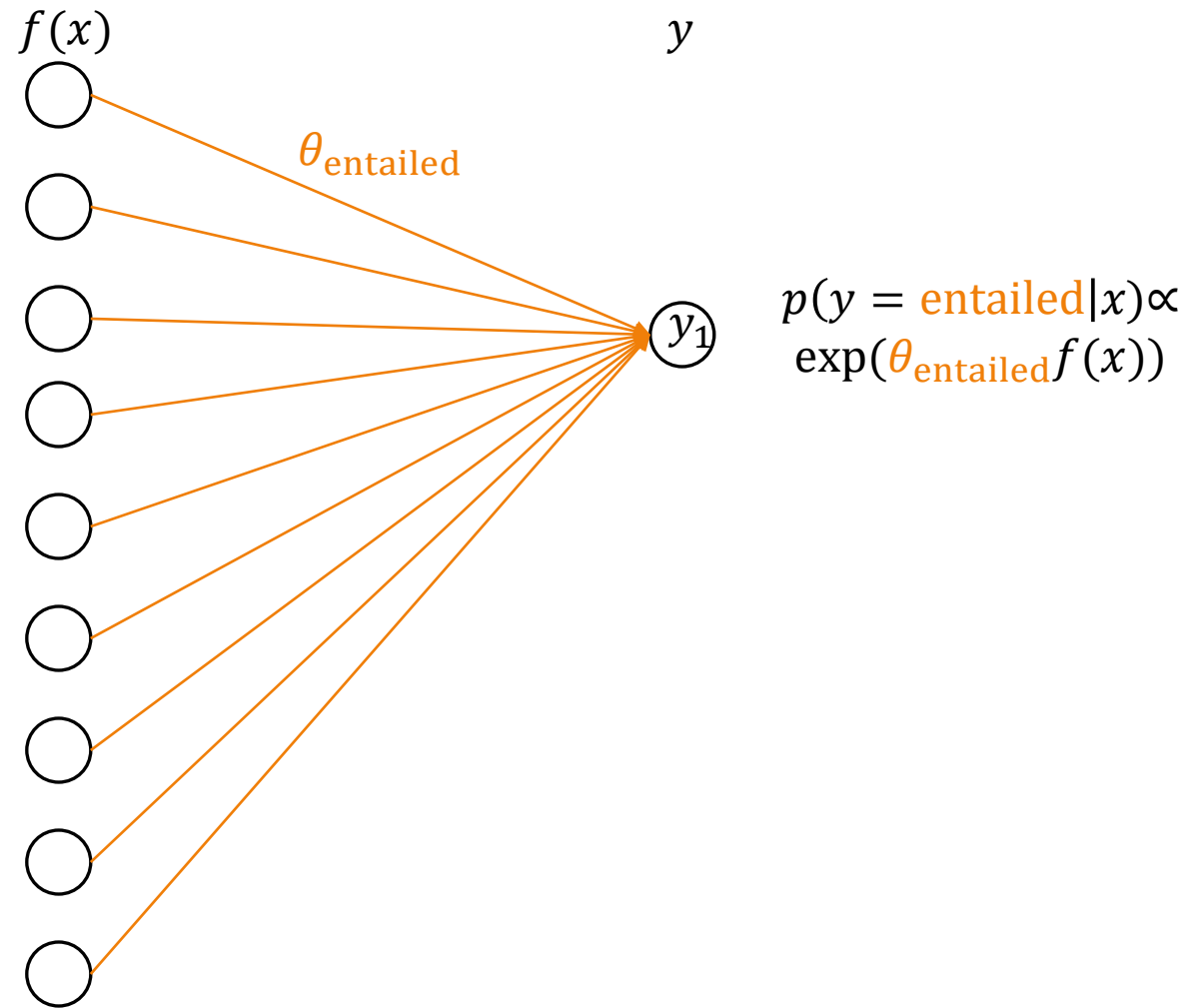
How can we implement the above computation in Pytorch?

Knowledge Check: Data Prep

<https://colab.research.google.com/drive/19yg0EUXQtHozBiSuO6cKOBhoSPzQHgug?usp=sharing>



Maxent Classifier, schematically



Maxent Modeling

$p(\text{ENTAILED} | \text{...}) =$

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.

Q: How do we define Z?

$\frac{1}{Z} \exp(\theta \text{ENTAILED} f(\text{document}))$

K different
weights...

for K different
features...

multiplied and
then

Normalization for Classification

$$Z = \sum_{\text{label } j} \exp(\theta_j^T f(\text{document icon}))$$

$$p(y | x) \propto \exp(\theta_y^T f(x))$$

classify doc x with label y in one go

Normalization for Classification (long form)

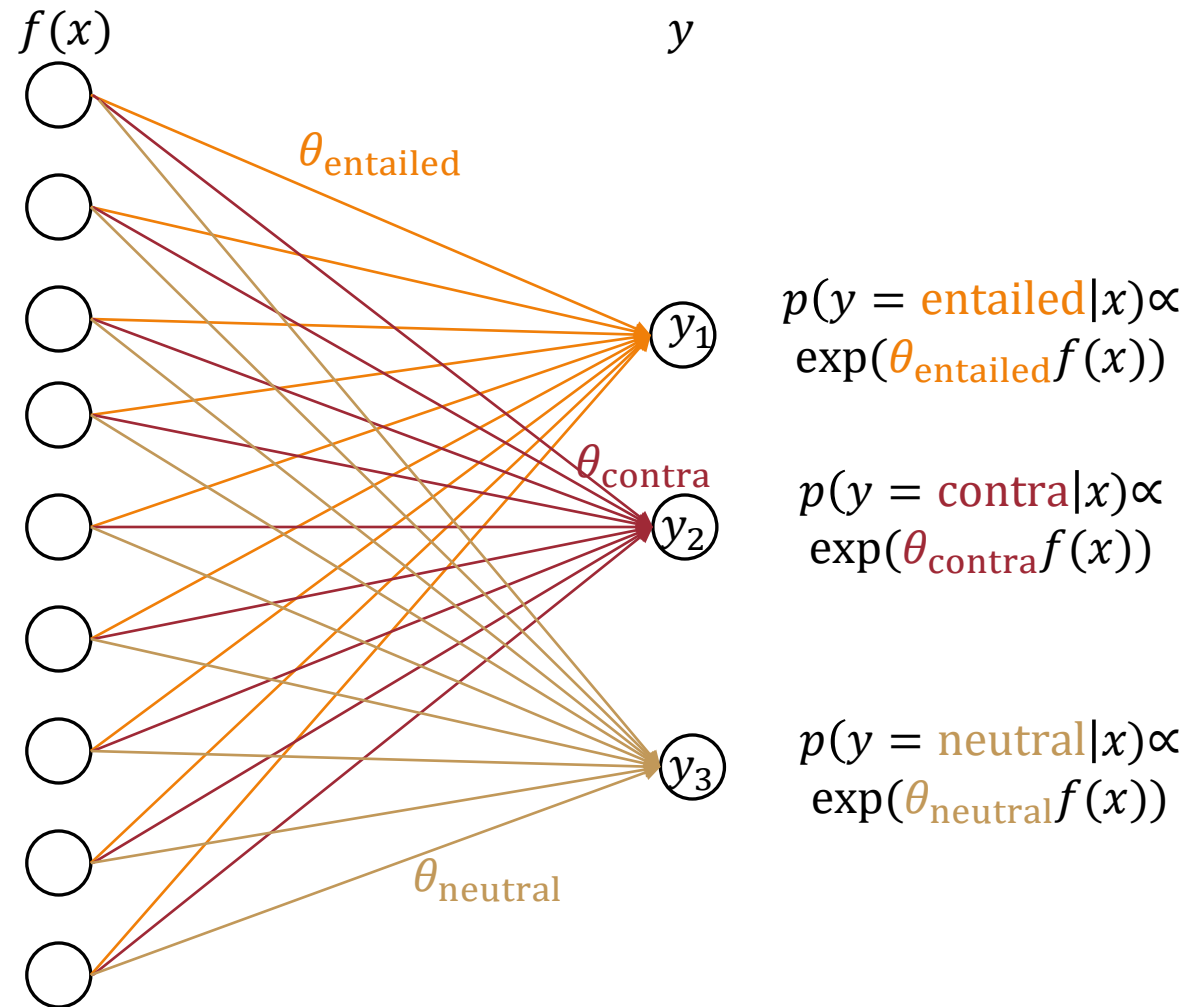
$$Z = \sum_{\text{label } j} \exp(\text{weight}_{1,j} * \text{applies}_1(\text{📄}) + \text{weight}_{2,j} * \text{applies}_2(\text{📄}) + \text{weight}_{3,j} * \text{applies}_3(\text{📄}) + \dots)$$

$$p(y | x) \propto \exp(\theta_y^T f(x))$$

classify doc x with label y in one go

Maxent Classifier, schematically

Why would we want to normalize the weights?



output:
 $i = \text{argmax score}_i$
class i

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

features $f(x)$ from x that are meaningful;

weights θ (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 θ

$$p(\mathbf{y} | x) = \frac{\exp(\theta_{\mathbf{y}}^T f(x))}{\sum_{y'} \exp(\theta_{y'}^T f(x))}$$

Different Notation, Same Meaning

$$p(Y = y | x) = \frac{\exp(\theta_y^T f(x))}{\sum_{y'} \exp(\theta_{y'}^T f(x))}$$

Different Notation, Same Meaning

$$p(Y = y | x) = \frac{\exp(\theta_y^T f(x))}{\sum_{y'} \exp(\theta_{y'}^T f(x))}$$

$$p(Y = y | x) \propto \exp(\theta_y^T f(x))$$

Different Notation, Same Meaning

$$p(Y = y | x) = \frac{\exp(\theta_y^T f(x))}{\sum_{y'} \exp(\theta_{y'}^T f(x))}$$

$$p(Y = y | x) \propto \exp(\theta_y^T f(x))$$

$$p(Y | x) = \text{softmax}(\theta f(x))$$

Defining Appropriate Features in a Maxent Model

Feature functions help extract useful features (characteristics) of the data

They turn *data* into *numbers*

Features that are not 0 are said to have fired

Generally *templated*

Binary-valued (0 or 1) or real-valued

Representing a Linguistic “Blob”

User-defined

Integer representation/on e-hot encoding

Assign each word to some index i , where $0 \leq i < V$

Represent each word w with a V -dimensional **binary** vector e_w , where $e_{w,i} = 1$ and 0 otherwise

Model-produced

Dense embedding

Let E be some *embedding size* (often 100, 200, 300, etc.)

Represent each word w with an E -dimensional **real-valued** vector e_w

Featurization is Similar but...

Vocab types (V) / embedding dimension (E) → number of features (number of “clues”)

“Linguistic blob” → Instances to represent

Features are extracted on each instance

Review: Bag-of-words as a Function

Based on some tokenization, turn an input document into an array (or dictionary or set) of its unique vocab items

Think of getting a BOW rep. as a function f

input: Document

output: Container of size E , indexable by

each vocab type v

Some Bag-of-words Functions

Kind	Type of f_v	Interpretation
Binary	0, 1	Did v appear in the document?
Count-based	Natural number (int ≥ 0)	How often did v occur in the document?
Averaged	Real number (≥ 0 , ≤ 1)	How often did v occur in the document, normalized by doc length?
TF-IDF (term frequency, inverse document frequency)	Real number (≥ 0)	How frequent is a word, tempered by how prevalent it is across the corpus (to be covered later!)
...		

Q: Is this a reasonable representation?

Q: What are some tradeoffs (benefits vs. costs)?

Useful Terminology: n-gram

Within a larger string (e.g., sentence),
a contiguous sequence of n items (e.g., words)

Colorless green ideas sleep furiously

n	Commonly called	History Size (Markov order)	Example n-gram ending in “furiously”
1	unigram	0	furiously
2	bigram	1	sleep furiously
3	trigram (3-gram)	2	ideas sleep furiously
4	4-gram	3	green ideas sleep furiously
n	n-gram	$n-1$	$w_{i-n+1} \dots w_{i-1} w_i$

Templated Features

Define a feature `fclue()` for each clue you want to consider

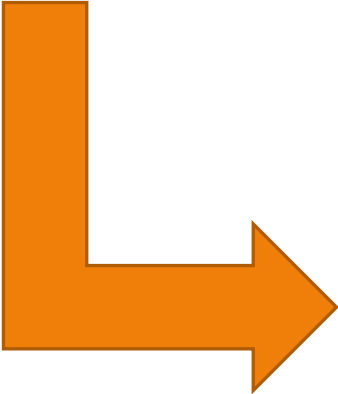
The feature `fclue` fires if the clue applies to/can be found in 

Clue is often a target phrase (an n-gram)

Maxent Modeling: Templated Binary Feature Functions

$$p(\text{ENTAILED} \mid \text{...}) \propto \exp\left(\text{weight}_{1, \text{Entailed}} * \text{applies}_1(\text{...}) + \text{weight}_{1, \text{Entailed}} * \text{applies}_2(\text{...}) + \text{weight}_{1, \text{Entailed}} * \text{applies}_3(\text{...}) + \dots\right)$$

$\text{applies}_{\text{target}}(\text{...}) = \begin{cases} 1, & \text{target matches } \text{...} \\ 0, & \text{otherwise} \end{cases}$
binary



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.

Example of a Templated Binary Feature Functions

$$\text{applies}_{\text{target}}(\text{document}) = \begin{cases} 1, & \text{target matches document} \\ 0, & \text{otherwise} \end{cases}$$



$$\text{applies}_{\text{ball}}(\text{document}) = \begin{cases} 1, & \text{ball in both s and h of document} \\ 0, & \text{otherwise} \end{cases}$$

Example of a Templated Binary Feature Functions

$$\text{applies}_{\text{target}}(\text{doc}) = \begin{cases} 1, & \text{target matches doc} \\ 0, & \text{otherwise} \end{cases}$$



$$\text{applies}_{\text{ball}}(\text{doc}) = \begin{cases} 1, & \text{ball in both s and h of doc} \\ 0, & \text{otherwise} \end{cases}$$

Q: If there are V vocab types and L label types:

1. How many features are defined if unigram targets are used (w/ each label)?

Example of a Templated Binary Feature Functions

$$\text{applies}_{\text{target}}(\text{doc}) = \begin{cases} 1, & \text{target matches doc} \\ 0, & \text{otherwise} \end{cases}$$



$$\text{applies}_{\text{ball}}(\text{doc}) = \begin{cases} 1, & \text{ball in both s and h of doc} \\ 0, & \text{otherwise} \end{cases}$$

Q: If there are V vocab types and L label types:

1. How many features are defined if unigram targets are used (w/ each label)?

A1: VL

Example of a Templated Binary Feature Functions

$$\text{applies}_{\text{target}}(\text{document}) = \begin{cases} 1, & \text{target matches document} \\ 0, & \text{otherwise} \end{cases}$$



$$\text{applies}_{\text{ball}}(\text{document}) = \begin{cases} 1, & \text{ball in both s and h of document} \\ 0, & \text{otherwise} \end{cases}$$

Q: If there are V vocab types and L label types:

1. How many features are defined if unigram targets are used (w/ each label)?

A1: VL

2. How many features are defined if bigram targets are used?

Example of a Templated Binary Feature Functions

$$\text{applies}_{\text{target}}(\text{doc}) = \begin{cases} 1, & \text{target matches doc} \\ 0, & \text{otherwise} \end{cases}$$



$$\text{applies}_{\text{ball}}(\text{doc}) = \begin{cases} 1, & \text{ball in both s and h of doc} \\ 0, & \text{otherwise} \end{cases}$$

Q: If there are V vocab types and L label types:

1. How many features are defined if unigram targets are used (w/ each label)?

$$A1: VL$$

2. How many features are defined if bigram targets are used (w/ each label)?

$$A2: V^2L$$

Example of a Templated Binary Feature Functions

$$\text{applies}_{\text{target}}(\text{document}) = \begin{cases} 1, & \text{target matches document} \\ 0, & \text{otherwise} \end{cases}$$



$$\text{applies}_{\text{ball}}(\text{document}) = \begin{cases} 1, & \text{ball in both s and h of document} \\ 0, & \text{otherwise} \end{cases}$$

Q: If there are V vocab types and L label types:

1. How many features are defined if unigram targets are used (w/ each label)?

$$A1: VL$$

2. How many features are defined if bigram targets are used (w/ each label)?

$$A2: V^2L$$

3. How many features are defined if unigram and bigram targets are used (w/ each label)?

Example of a Templated Binary Feature Functions

$$\text{applies}_{\text{target}}(\text{doc}) = \begin{cases} 1, & \text{target matches doc} \\ 0, & \text{otherwise} \end{cases}$$



$$\text{applies}_{\text{ball}}(\text{doc}) = \begin{cases} 1, & \text{ball in both s and h of doc} \\ 0, & \text{otherwise} \end{cases}$$

Q: If there are V vocab types and L label types:

1. How many features are defined if unigram targets are used (w/ each label)?

$$A1: VL$$

2. How many features are defined if bigram targets are used (w/ each label)?

$$A2: V^2L$$

3. How many features are defined if unigram and bigram targets are used (w/ each label)?

$$A2: (V + V^2)L$$