

Logistic Regression Models

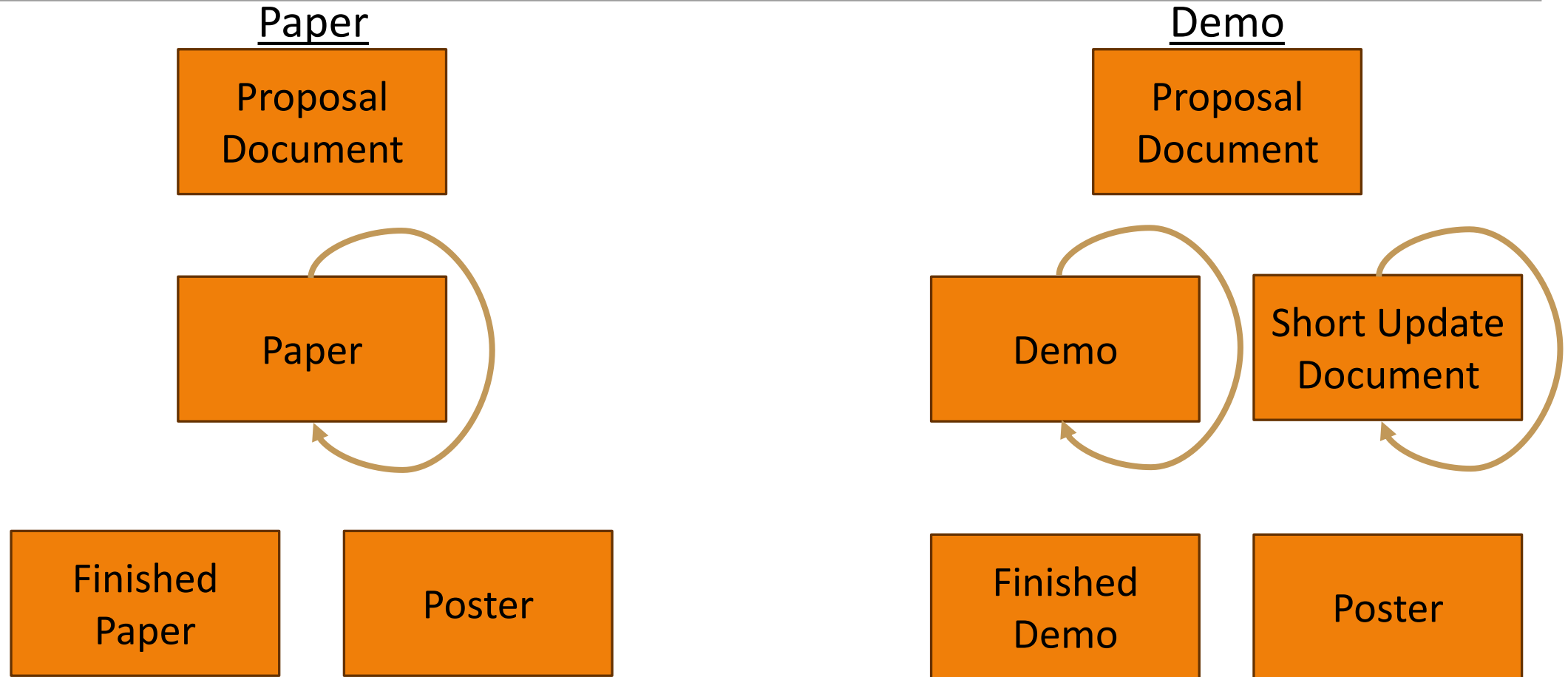
Instructor: Lara J. Martin (she/they)

TA: Omkar Kulkarni (he)

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

Slides modified from Dr. Frank Ferraro

Project Flow



Learning Objectives

Model classification problems using logistic regression

Define appropriate features for a logistic regression problem

Review: F1 (or F-score)

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

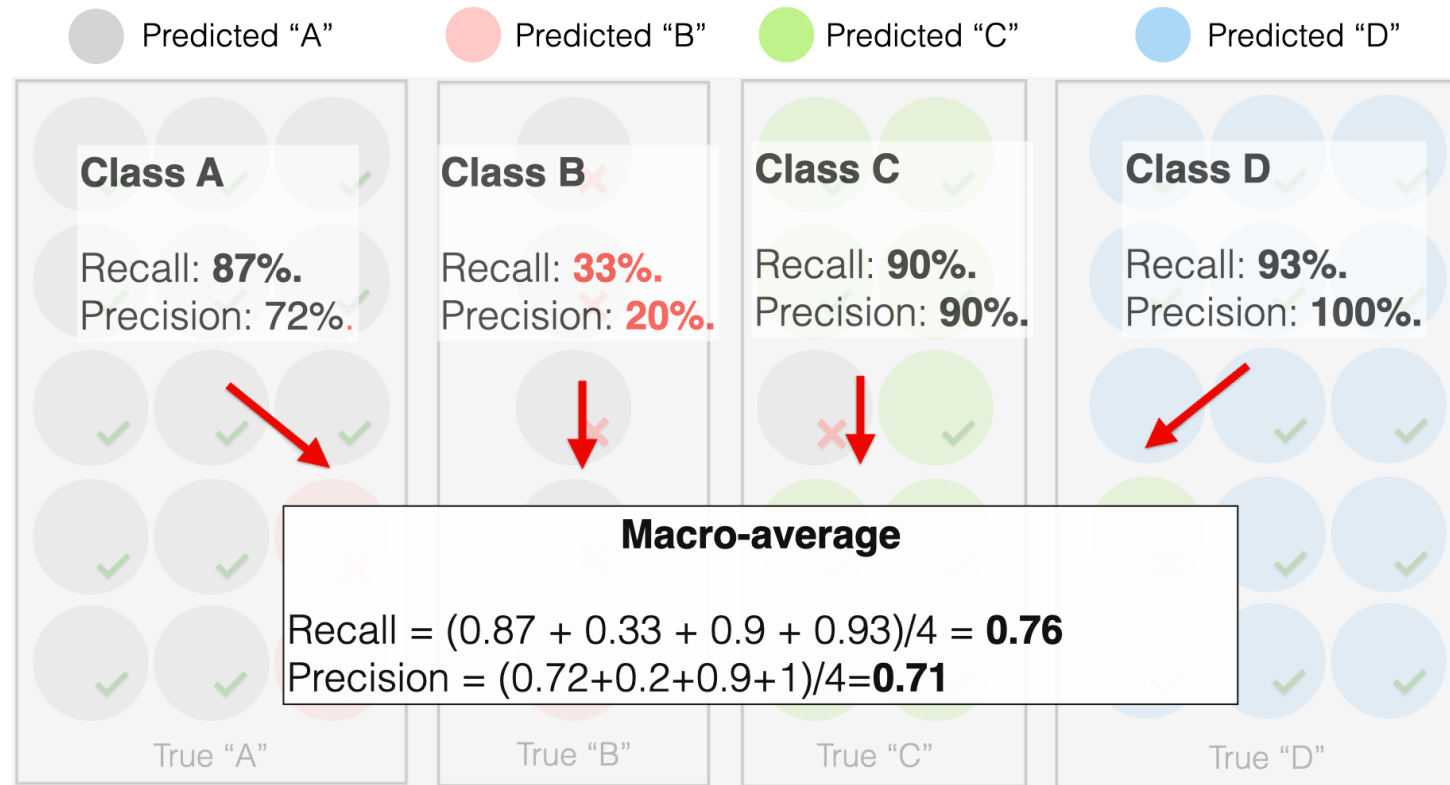
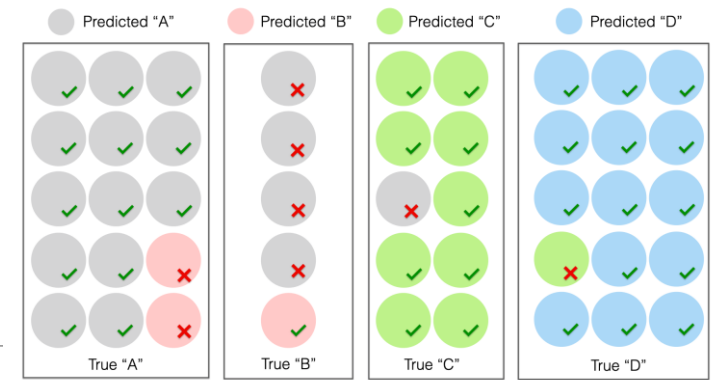
F1 measure: equal weighting between precision and recall

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

(useful when $P = R = 0$)

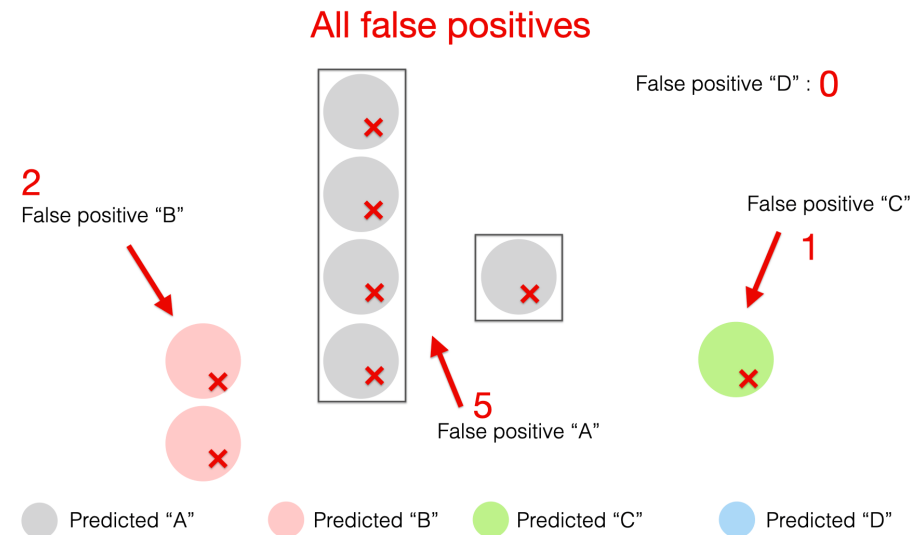
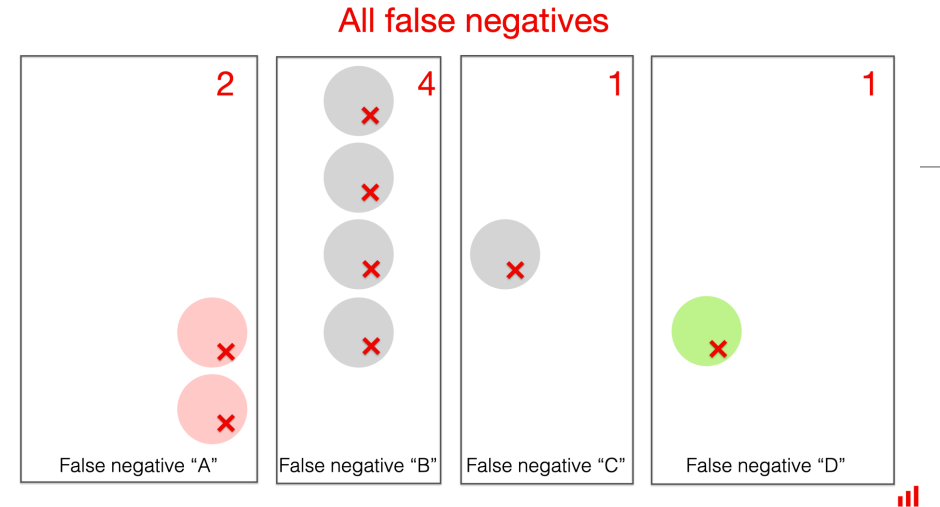
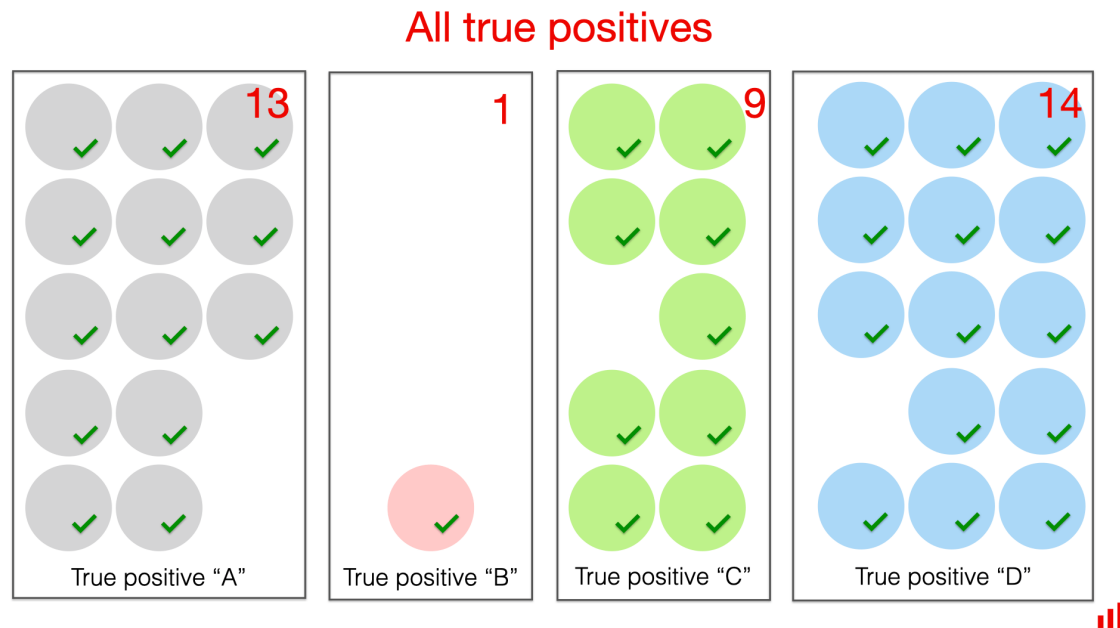
Each *class* has equal weight

Review: Macro-Average



Each *instance* has equal weight

Review: Micro-Average



	Total TP	Total FP	Total FN
Precision	13 + 1 + 9 + 14	2 + 5 + 1 + 0	
Recall	13 + 1 + 9 + 14		2 + 4 + 1 + 1

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

Micro-average Recall = $\frac{13 + 1 + 9 + 14}{13 + 1 + 9 + 14 + 2 + 4 + 1 + 1} = 0.82$

<https://www.evidentlyai.com/classification-metrics/multi-class-metrics>

Types of Classification Metrics

AUC <https://scikit-learn.org/stable/modules/generated/sklearn.metrics.auc.html>

F-score https://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1_score.html

Accuracy https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html

Confusion matrix https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html

Precision (can specify macro/micro average) https://scikit-learn.org/stable/modules/generated/sklearn.metrics.precision_score.html

Recall (can specify macro/micro average) https://scikit-learn.org/stable/modules/generated/sklearn.metrics.recall_score.html

Outline

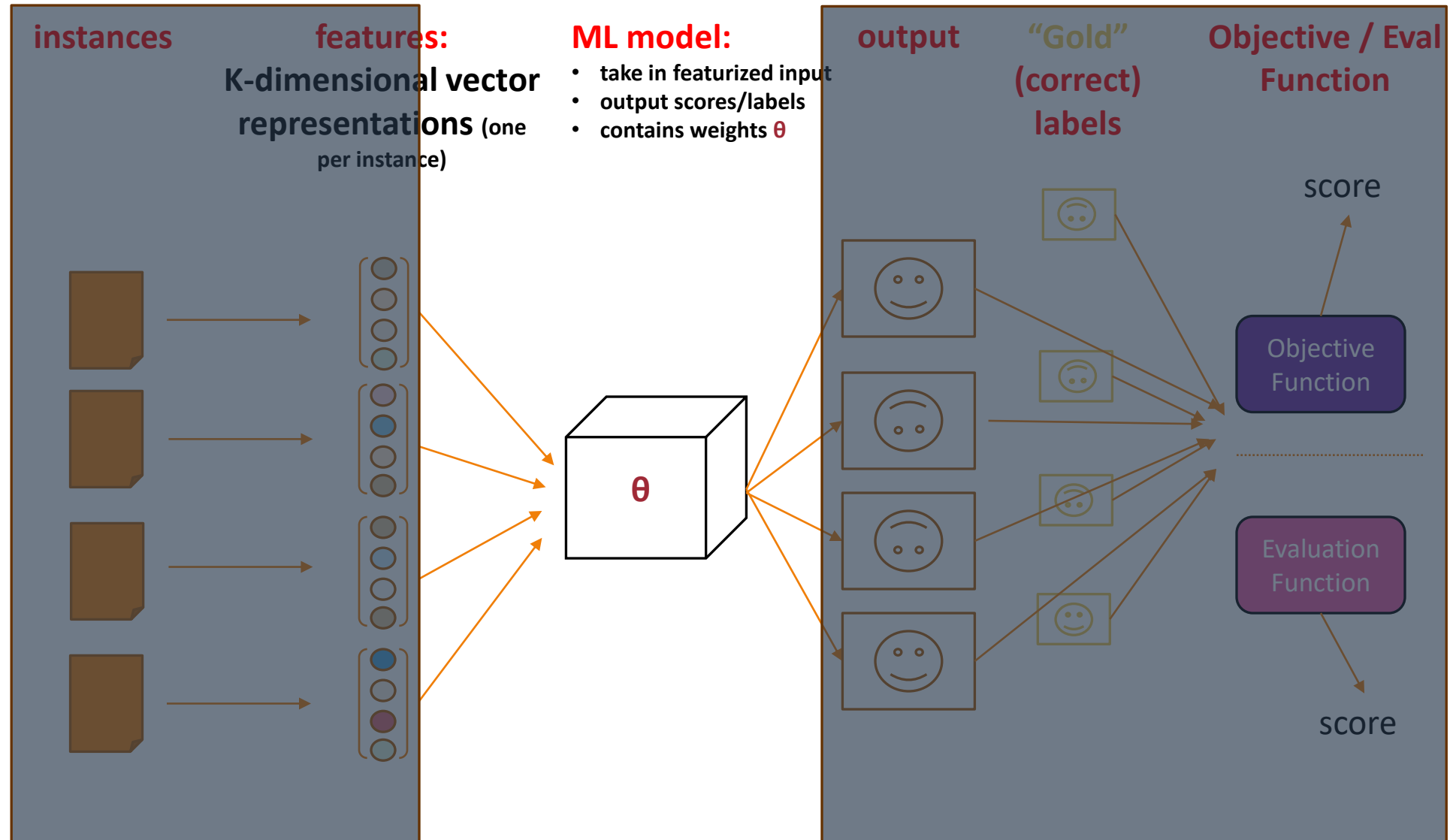
Maximum Entropy classifiers

Defining the model

Defining the objective

Learning: Optimizing the objective

Defining the Model



Maxent Models for Classification: Discriminatively or Generatively Trained

Directly model
the posterior

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

Discriminatively trained classifier

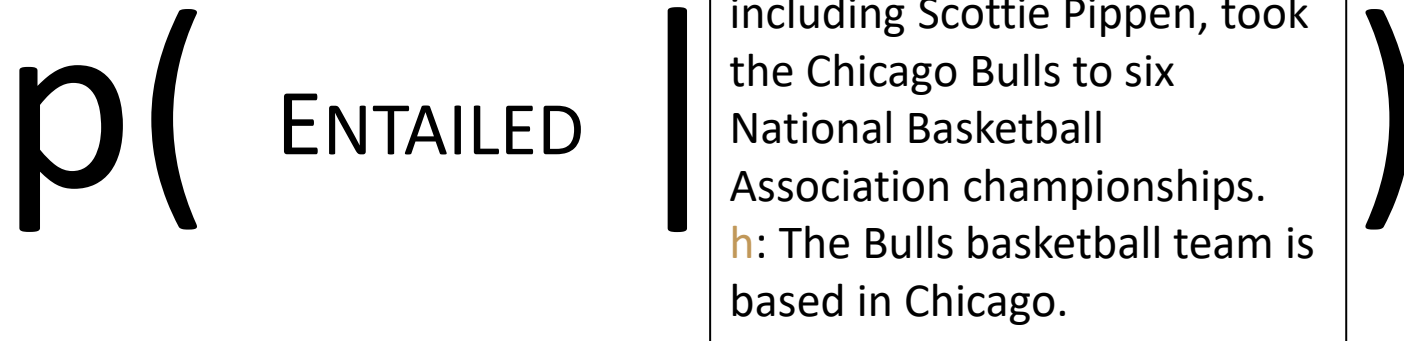
Model the
posterior with
Bayes rule

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

Generatively trained classifier with
maxent-based language model

Review:

Discriminative Model using Document Classification Example



Review: Extracting Features

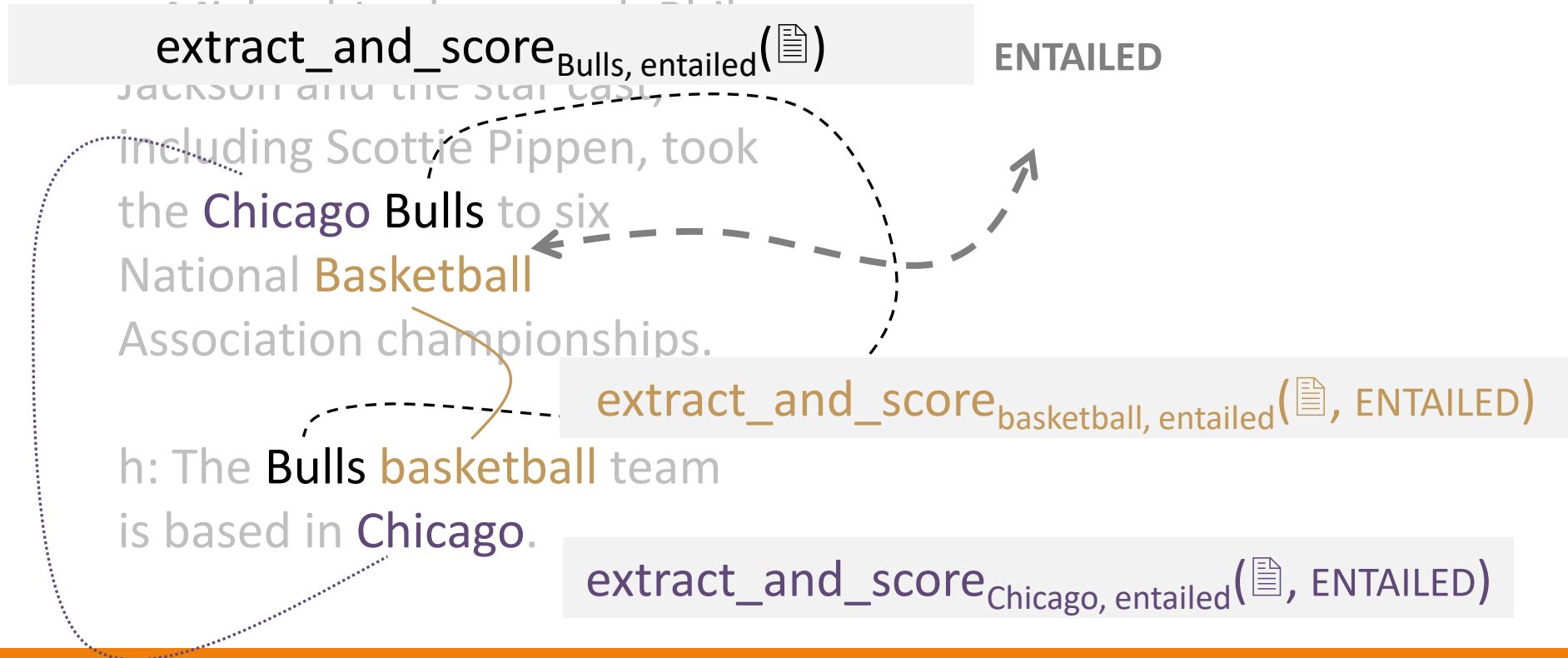
s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the **Chicago Bulls** to six National **Basketball** Association championships.

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.



Review: 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 (📄)

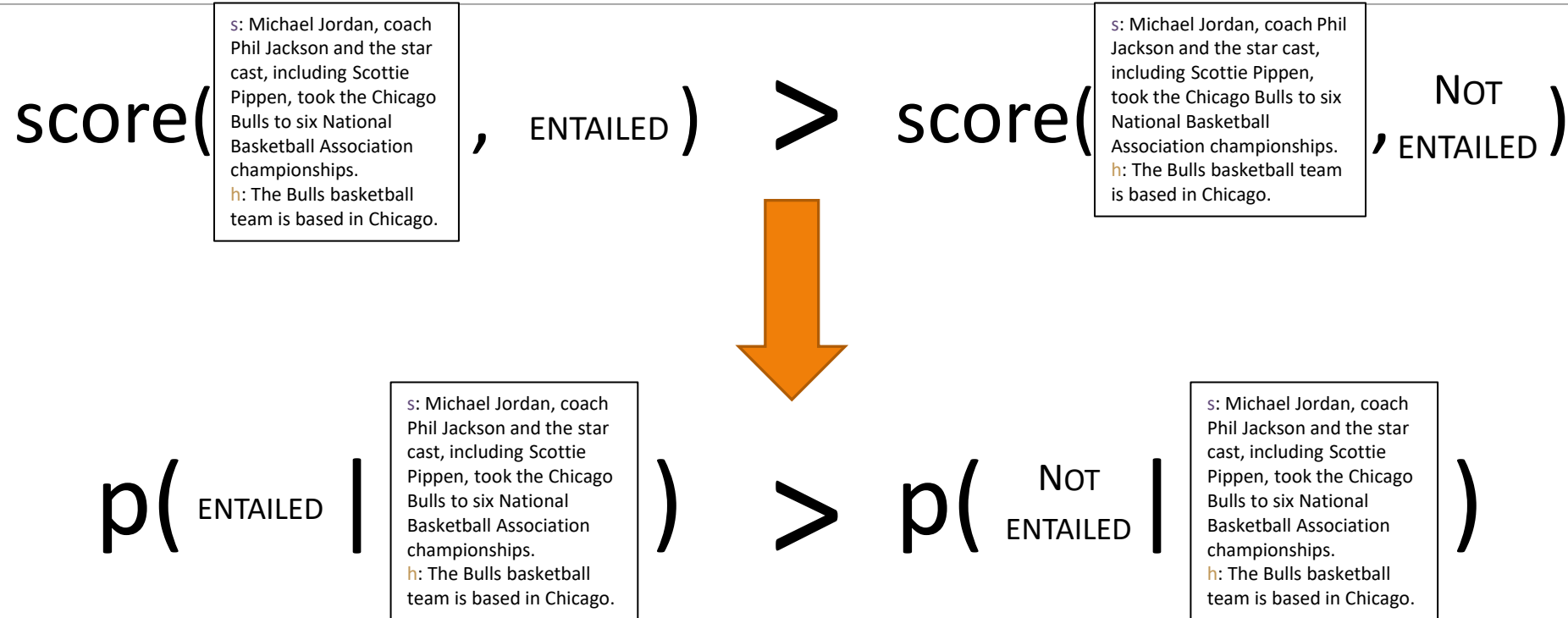
+

score₃ , Entailed (📄)

+

...

Review: Turning Scores into Probabilities



KEY IDEA

Maxent Modeling

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

Proportional to

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?

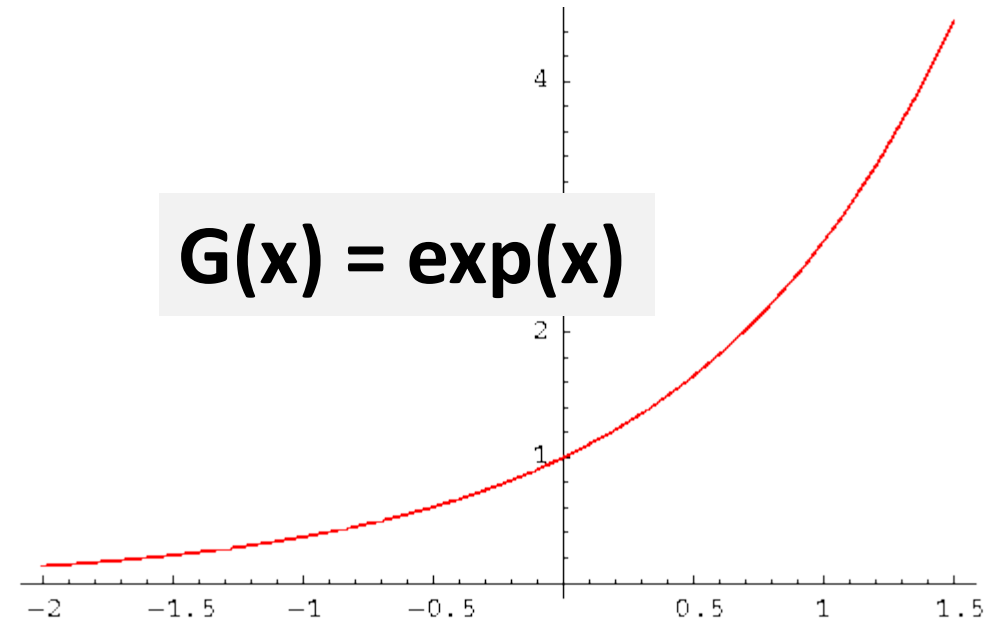
is monotonic? (if $a < b$, then $G(a) < G(b)$)

What function G...

operates on any real number?

is never less than 0?

is monotonic? (if $a < b$, then $G(a) < G(b)$)



Maxent Modeling

$p(\text{ENTAILED} \mid$

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.

$) \propto$

$\exp(\text{score}(\text{ENTAILED}))$

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.

Maxent Modeling

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

$$\exp\left(\begin{array}{l} \text{score}_{1, \text{Entailed}}(\text{document icon}) + \\ \text{score}_{2, \text{Entailed}}(\text{document icon}) + \\ \text{score}_{3, \text{Entailed}}(\text{document icon}) + \\ \dots \end{array}\right)$$

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.

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

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.

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

K different
weights...

for K different
features

$$\begin{bmatrix} \theta \\ .31 \\ -.5 \\ .1 \\ .002 \\ .522 \\ \dots \end{bmatrix}$$

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

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.

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

K different
weights...

for K different
features

multiplied and then summed

Maxent Modeling

$p(\text{ENTAILED} \mid$

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.

$) \propto$

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

K different
weights...

for K different
features

multiplied and
then summed

Maxent Modeling

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

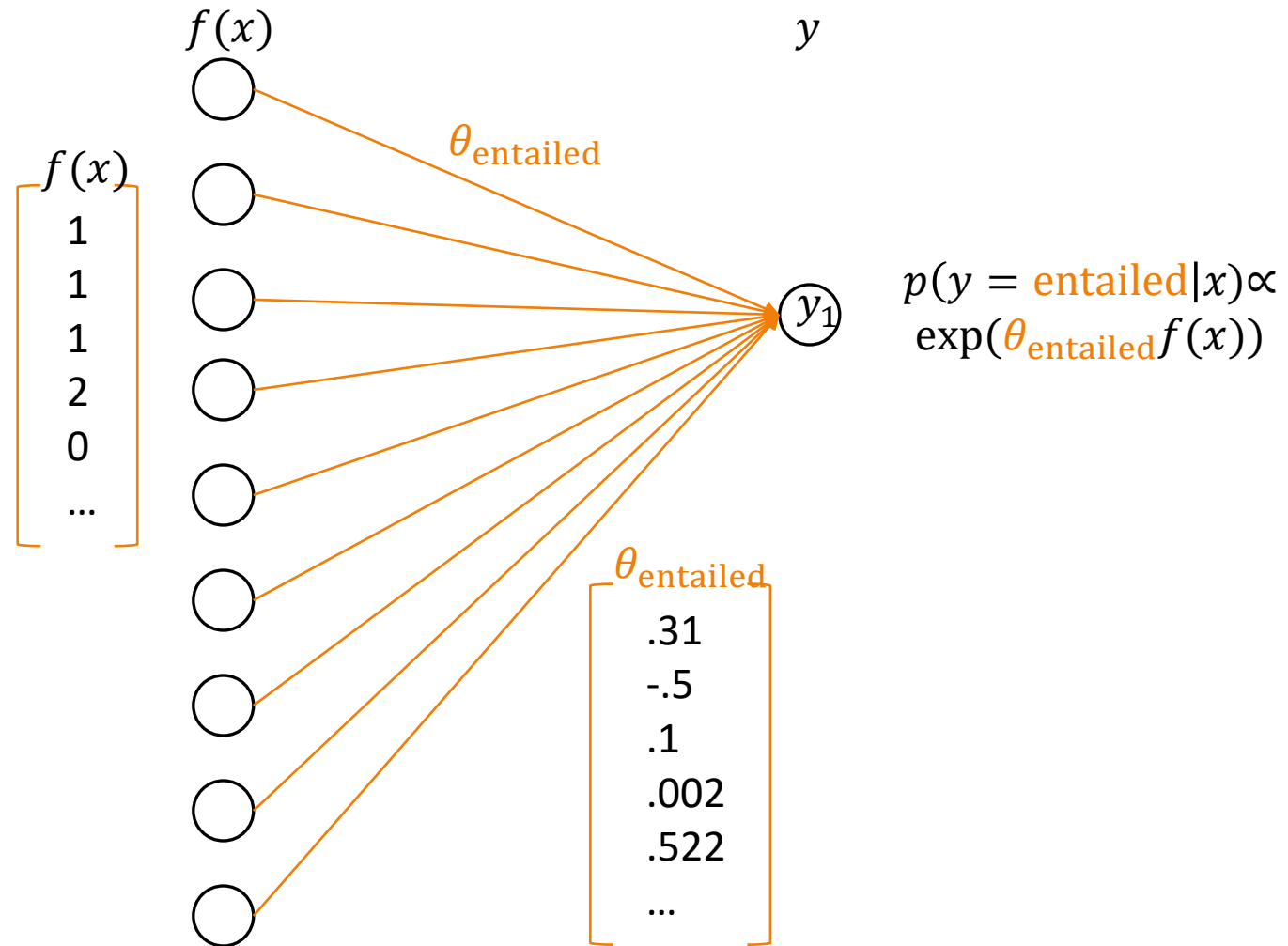
s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the Chicago Bulls to six National Basketball Association championships.

h: The Bulls basketball team is based in Chicago.

$$\exp\left(\theta_{\text{ENTAILED}}^T f(\text{document})\right) \times \begin{bmatrix} .31 & -.5 & .1 & .002 & .522 & \dots \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 0 \\ \dots \end{bmatrix}$$

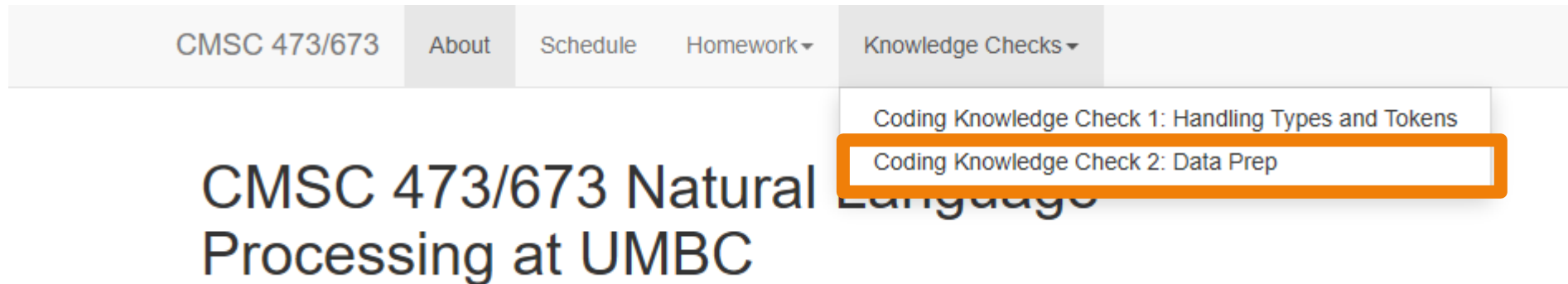
K different weights... for K different features multiplied and then summed

Maxent Classifier, schematically



Knowledge Check: Data Prep

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



The image shows a navigation bar for the course 'CMSC 473/673'. The bar contains links for 'About', 'Schedule', 'Homework', and 'Knowledge Checks'. The 'Knowledge Checks' link is expanded, showing two options: 'Coding Knowledge Check 1: Handling Types and Tokens' and 'Coding Knowledge Check 2: Data Prep'. The second option is highlighted with an orange border. Below the navigation bar, the course title 'CMSC 473/673 Natural Language Processing at UMBC' is visible.

CMSC 473/673

About Schedule Homework Knowledge Checks

Coding Knowledge Check 1: Handling Types and Tokens

Coding Knowledge Check 2: Data Prep

CMSC 473/673 Natural Language Processing at UMBC

Maxent Modeling

$$p(\text{ENTAILED} \mid \text{[document]}) \propto$$

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.

$$\frac{1}{Z} \exp(\theta_{\text{ENTAILED}}^T f(\text{[document]}))$$

Maxent Modeling

$p(\text{ENTAILED} \mid$

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.

$) =$

How do we define Z?

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

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)

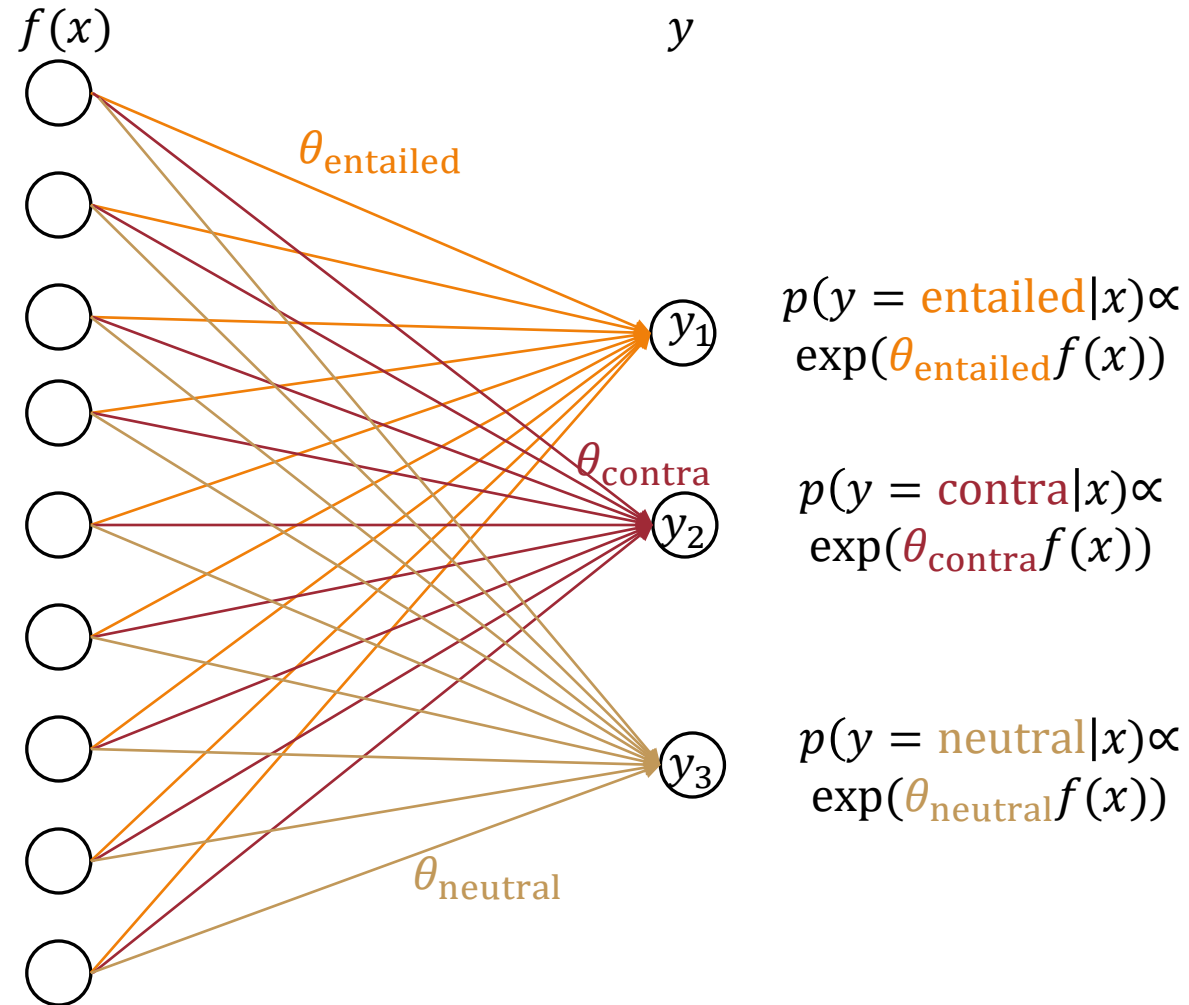
$$Z = \sum_{\text{label } j} \exp(\text{weight}_{1,j} * \text{applies}_1(\text{📄}) + \text{weight}_{2,j} * \text{applies}_2(\text{📄}) + \text{weight}_{3,j} * \text{applies}_3(\text{📄}) + \dots)$$

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

classify doc x with label y in one go

Multiclass Maxent Classifier, schematically

Why would we want to normalize the weights?



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

Final Equation for Logistic Regression

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

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

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

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

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)?