

CMSC 473/673

Natural Language Processing

Instructor: Lara J. Martin (she/they)

TA: Duong Ta (he)

Slides modified from Dr. Frank Ferraro

Learning Objectives

Model classification problems using logistic regression

Define appropriate features for a logistic regression problem

Create a method to prepare data for a BoW model

Define an objective for LR modeling

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

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

Review: Classification Evaluation

Accuracy: % of items correct

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

F-score: Weighted (harmonic) average of
Precision & Recall

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

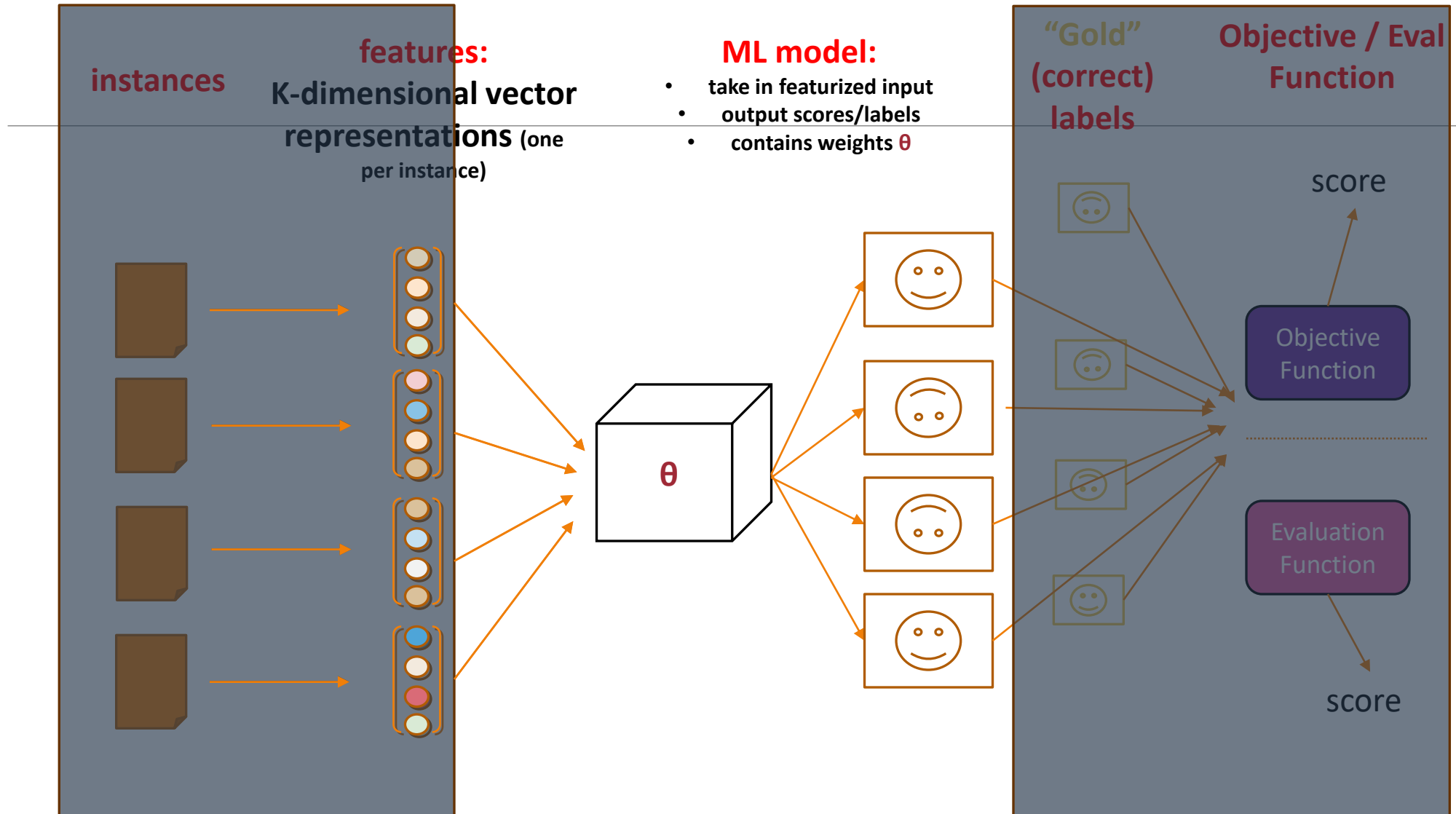
When would you want to use
accuracy vs F1?

Accuracy works better if
the dataset is balanced

Accuracy takes
everything in
consideration

F-Score is
focused on TP

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

$$p(\bullet | X) \text{ vs. } p(\circ | X)$$

Normally (*but this can be adjusted!)


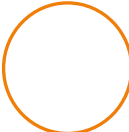
$$\text{best label} = \arg \max_{\text{label}} P(\text{label} | \text{example})$$

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

Likelihood Prior

$$\frac{P(X|Y) \cdot P(Y)}{P(X)}$$

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:

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

$$p(Y = \text{label}_1 | X) \geq 0,$$

$$p(Y = \text{label}_0 | X) \geq 0$$

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

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

Discriminatively trained classifier

Maxent Models for Classification: Discriminatively or Generatively Trained

Directly model
the posterior

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

Discriminatively trained classifier



Model the
posterior with
Bayes rule

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

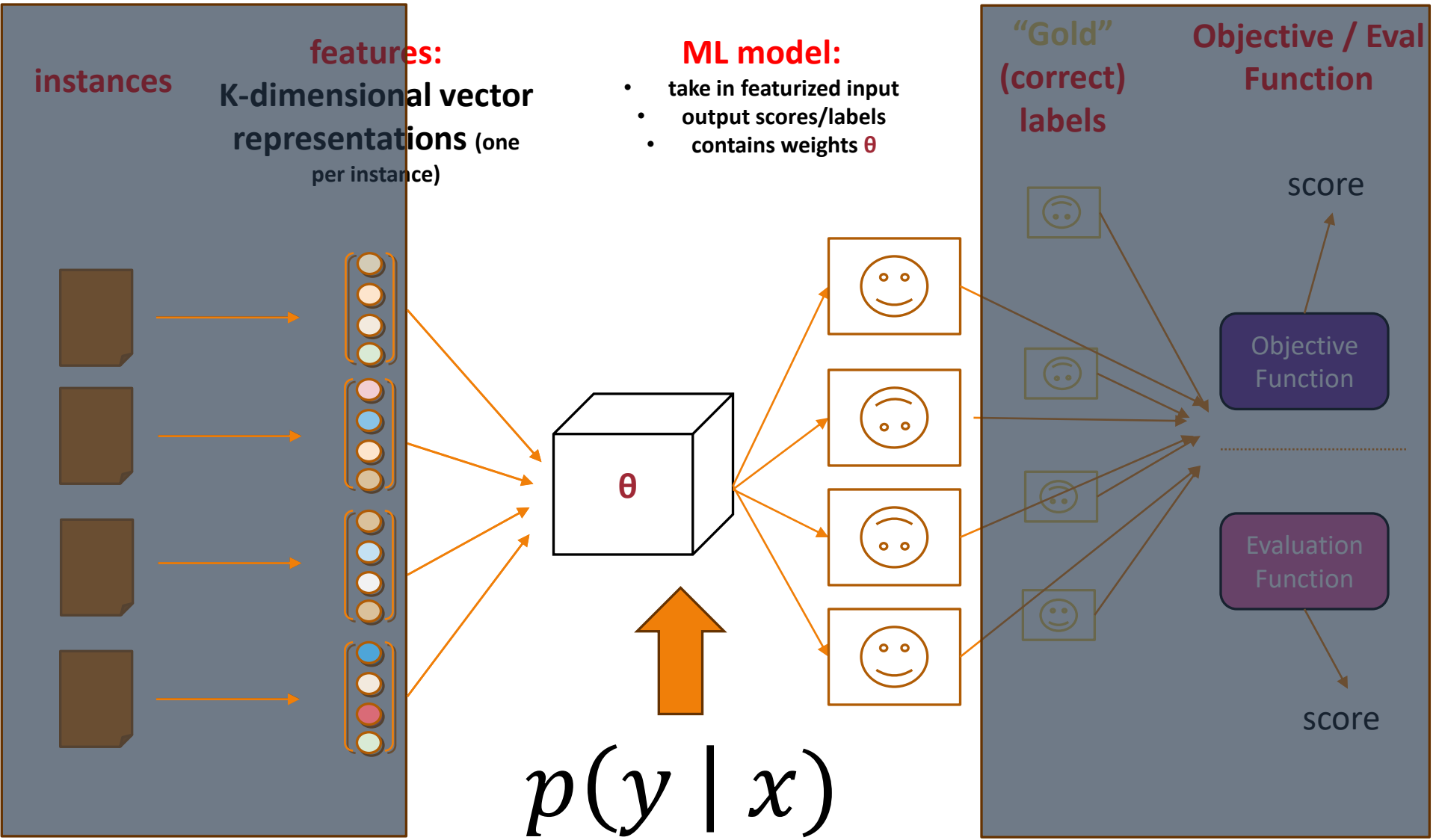
Generatively trained classifier with
maxent-based language model

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

(we'll start with this one)

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

discriminatively trained:
classify in one go



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

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 $p(y | x)$

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)

Electronic alerts have been used to assist the authorities in moments of chaos and potential danger: after the Boston bombing in 2013, when the Boston suspects were still at large, and last month in Los Angeles, during an active shooter scare at the airport.

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

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

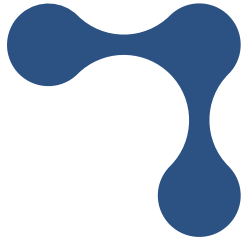
$$f(x) = (f_i(x))_i^V$$

TECH
NOT TECH

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

feature $f_i(x)$	value
alerts	1
assist	1
bombing	1
Boston	2
...	
sniffle	0
...	

Example Classification Tasks



GLUE

<https://gluebenchmark.com/>

🤖 datasets: glue

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

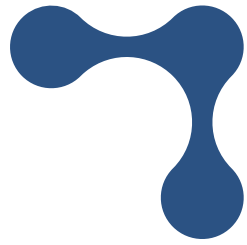
SuperGLUE 1

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

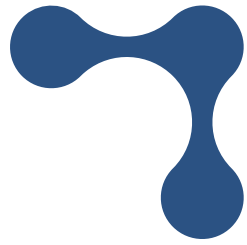


Recognizing Textual Entailment (RTE)

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

ENTAILMENT (yes):

NOT ENTAILED (no):



Recognizing Textual Entailment (RTE)

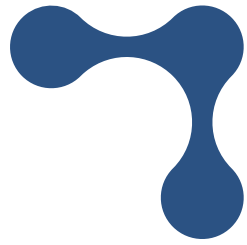
Given a premise sentence s and hypothesis sentence h , determine if h “follows from” s

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

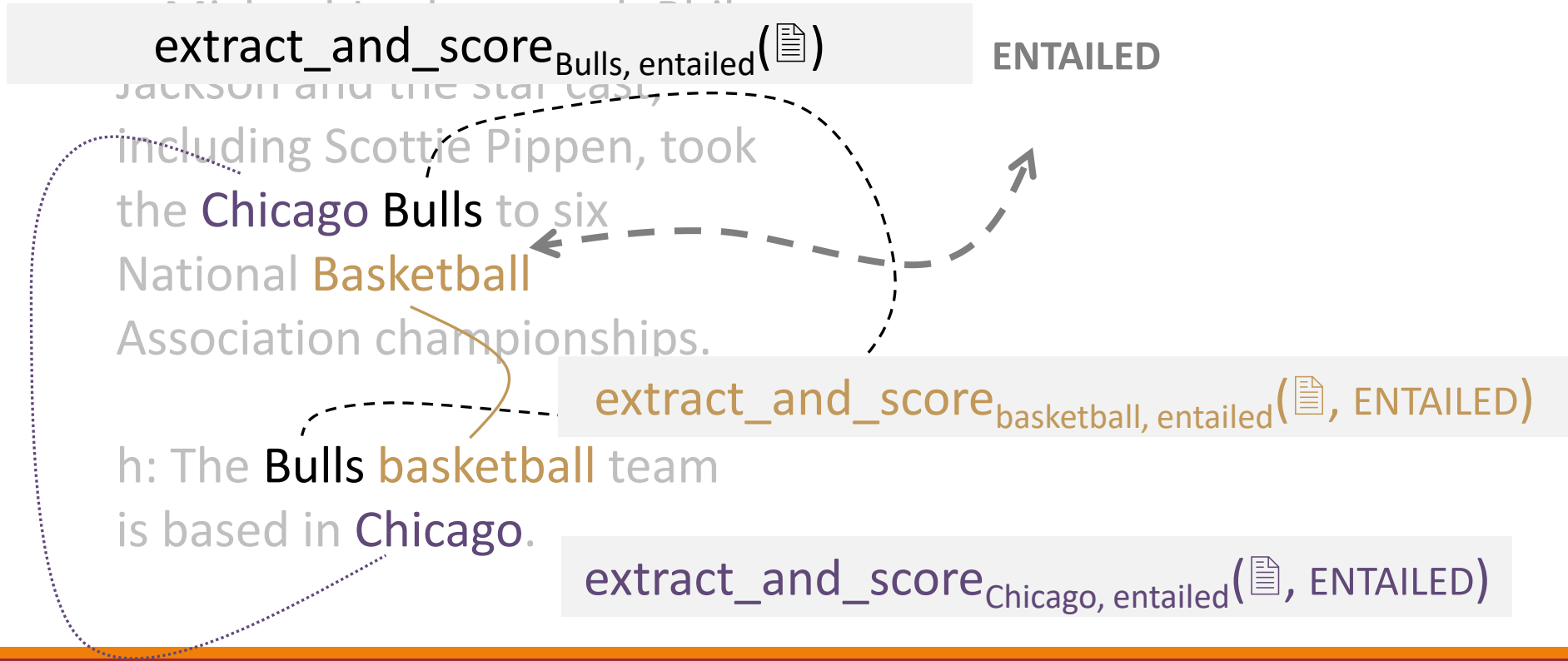
s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the **Chicago Bulls** to six National **Basketball** Association championships.

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

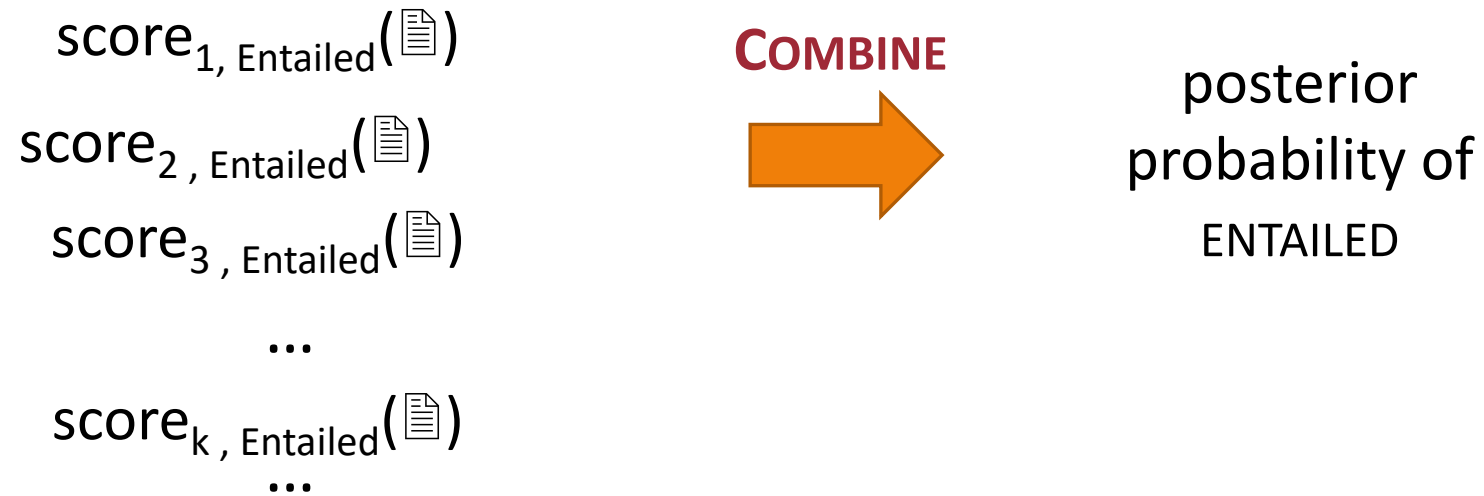
ENTAILED

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)*

score₁, Entailed (📄)

+

score₂, Entailed (📄)

+

score₃, Entailed (📄)

+

...

Turning Scores into Probabilities

$$\text{score}(s, \text{ENTAILED}) > \text{score}(s, \text{NOT ENTAILED})$$

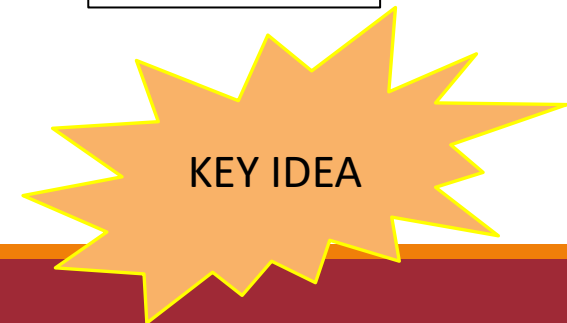
s: Michael Jordan, coach Phil Jackson and the star cast, including Scottie Pippen, took the Chicago Bulls to six National Basketball Association championships.
h: The Bulls basketball team is based in Chicago.

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.

$$p(\text{ENTAILED} | s) > p(\text{NOT ENTAILED} | s)$$

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.

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.



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?

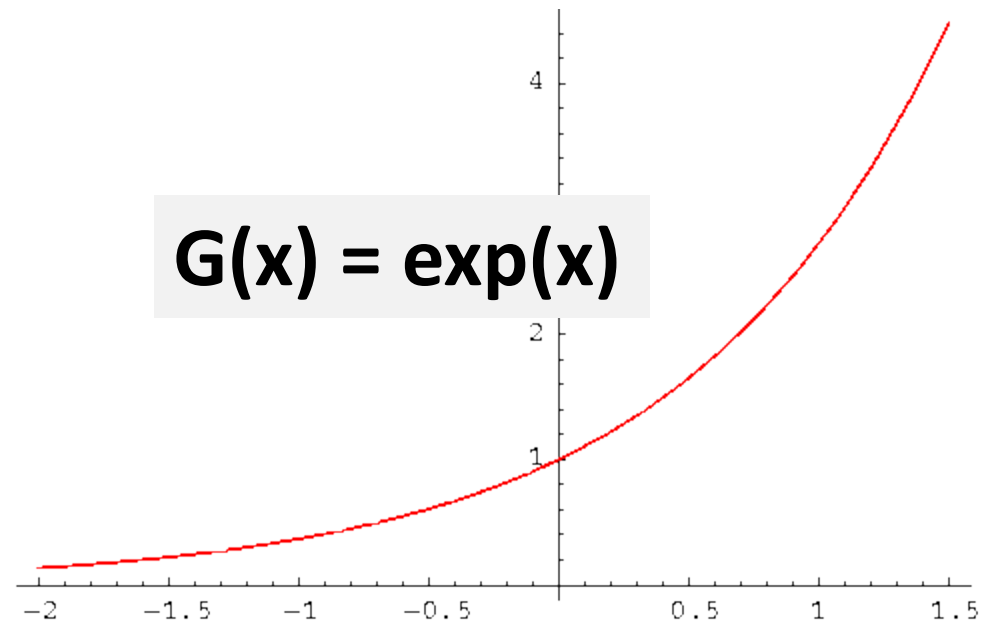
monotonic? ($a < b \rightarrow G(a) < G(b)$)

What function G...

operates on any real number?

is never less than 0?

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



Maxent Modeling

$p(\text{ENTAILED} |$

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

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

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.

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

Maxent Modeling

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

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.

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

K different
weights...

for K different
features

Maxent Modeling

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

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.

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

K different weights...

for K different features

multiplied and then summed

Maxent Modeling

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

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

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

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.

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

K different weights...

for K different features

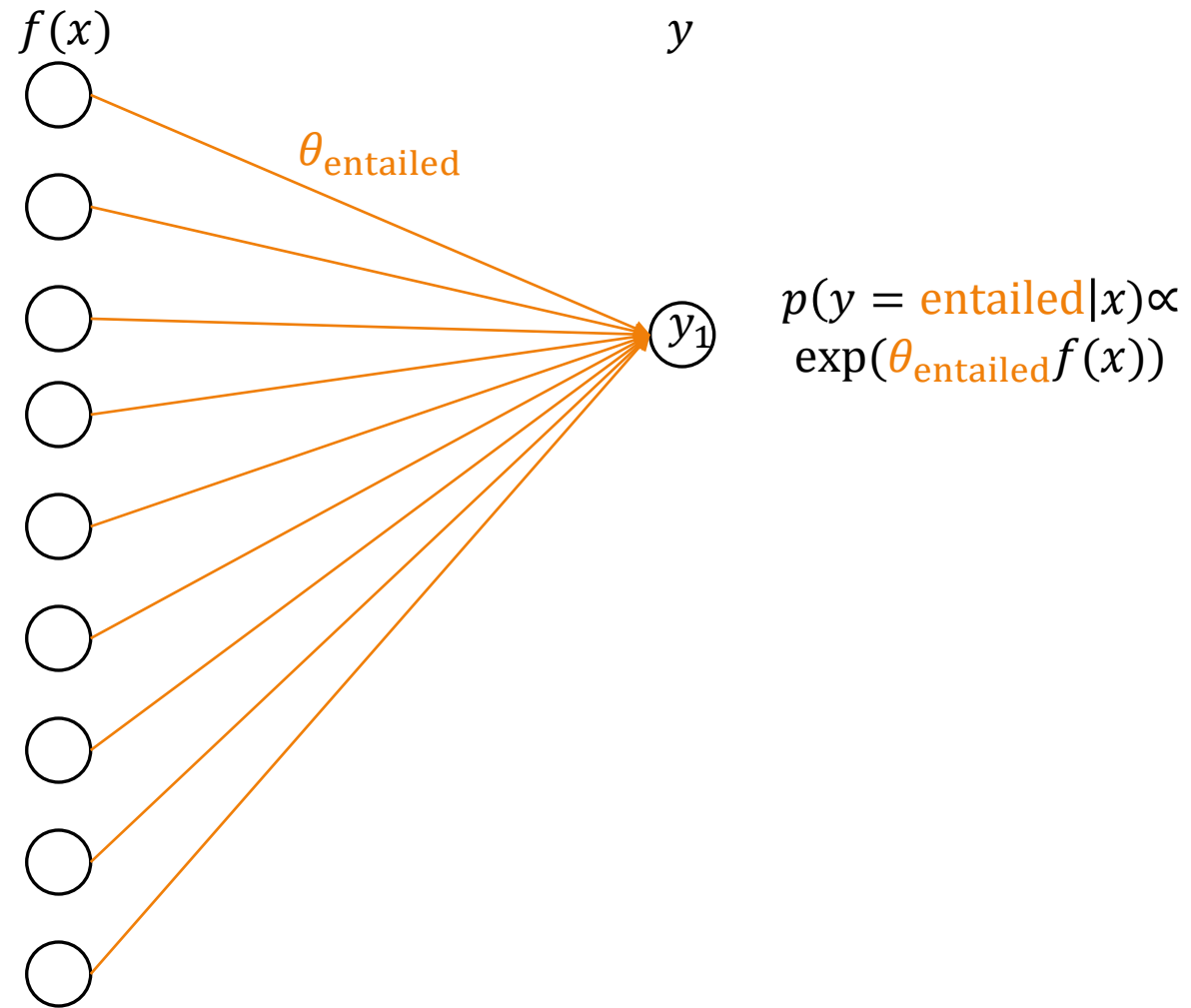
multiplied and then summed

Knowledge Check: Data Prep

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



Maxent Classifier, schematically



Maxent Modeling

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

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.

Q: How do we define Z?

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

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)

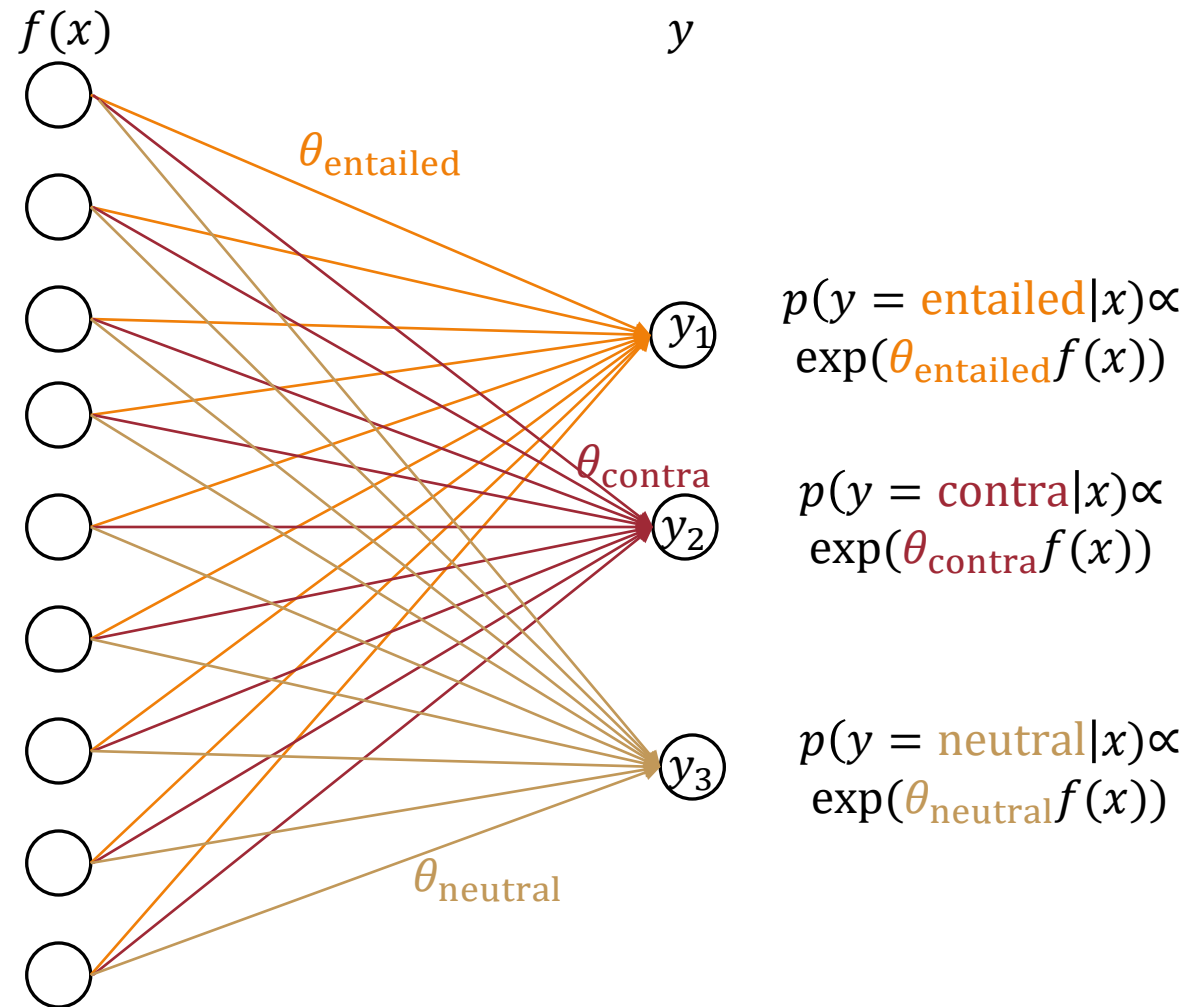
$$Z = \sum_{\text{label } j} \exp(\text{weight}_{1,j} * \text{applies}_1(\text{📄}) + \text{weight}_{2,j} * \text{applies}_2(\text{📄}) + \text{weight}_{3,j} * \text{applies}_3(\text{📄}) + \dots)$$

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

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

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

Different Notation, Same Meaning

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

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

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 ($\geq 0, \leq 1$)	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!)
...		