# 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

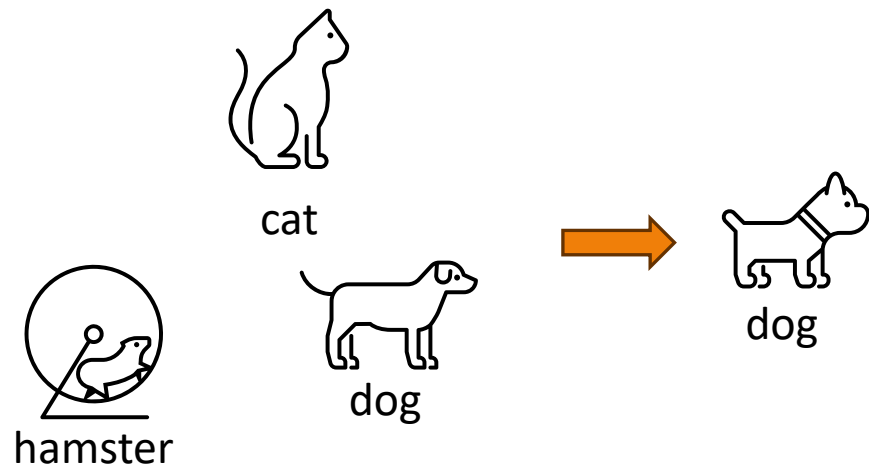# Review: Classification Evaluation: the 2-by-2 contingency table

| *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* |

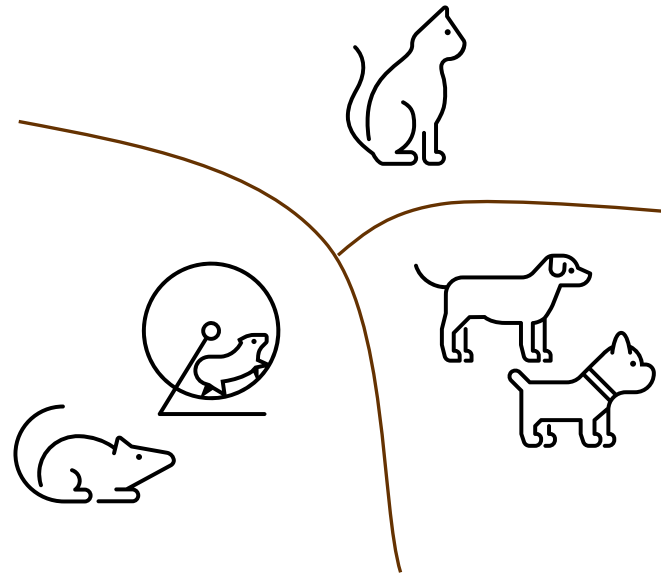Construct this table by *counting* the number of TPs, FPs, FNs, TNs

# Review: Types of Learning
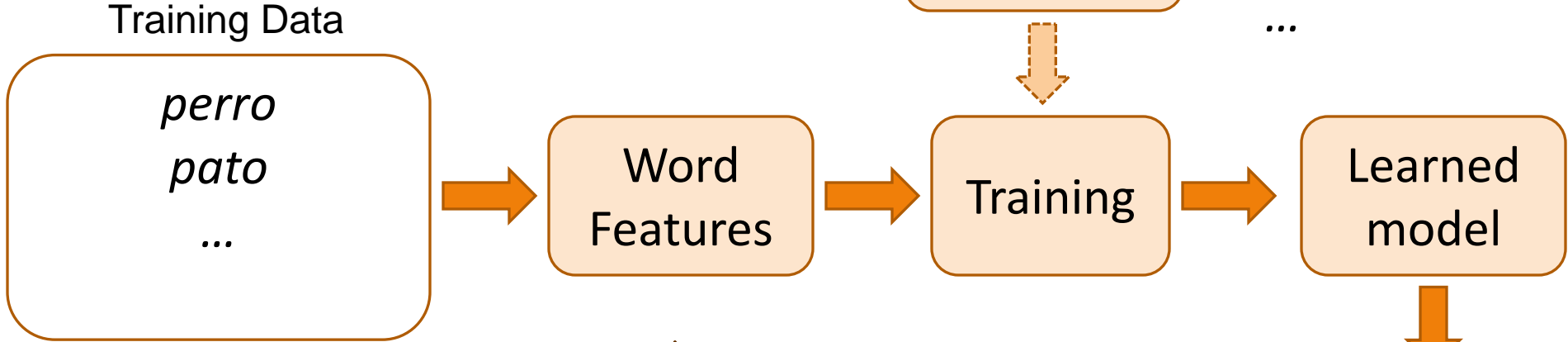
**SUPERVISED LEARNING**                    **UNSUPERVISED LEARNING**
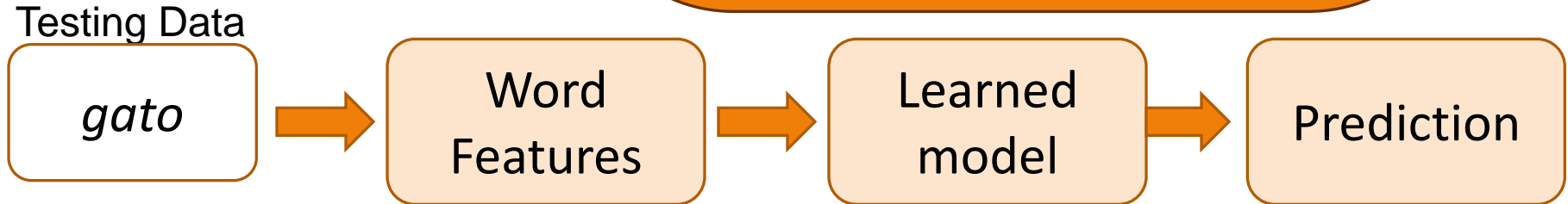


cat

dog

hamster

dog

# Review: Steps
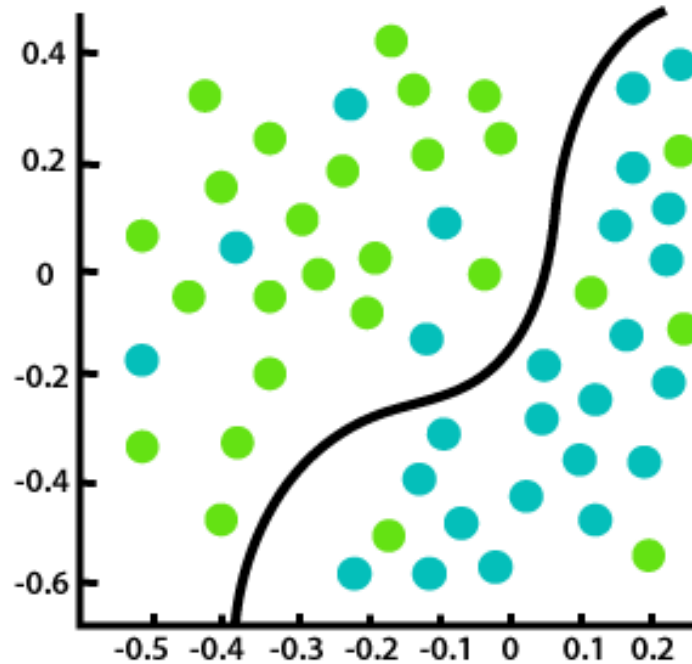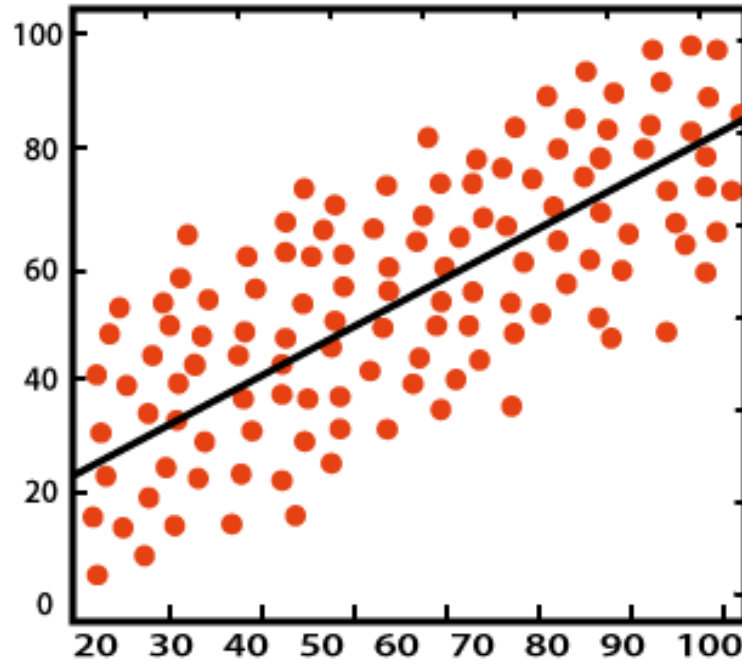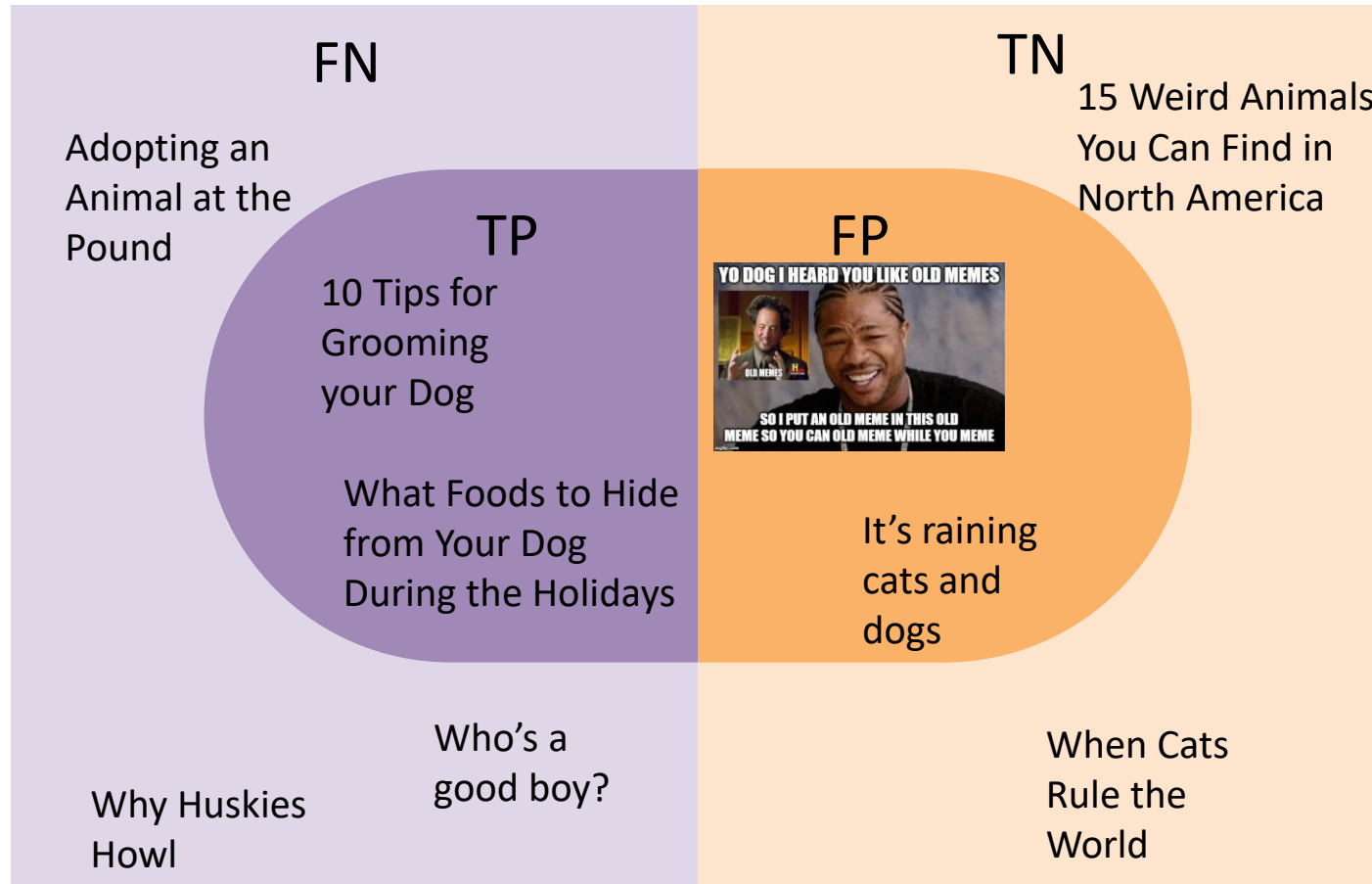
# Review: Types of models



Classification

Regression

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

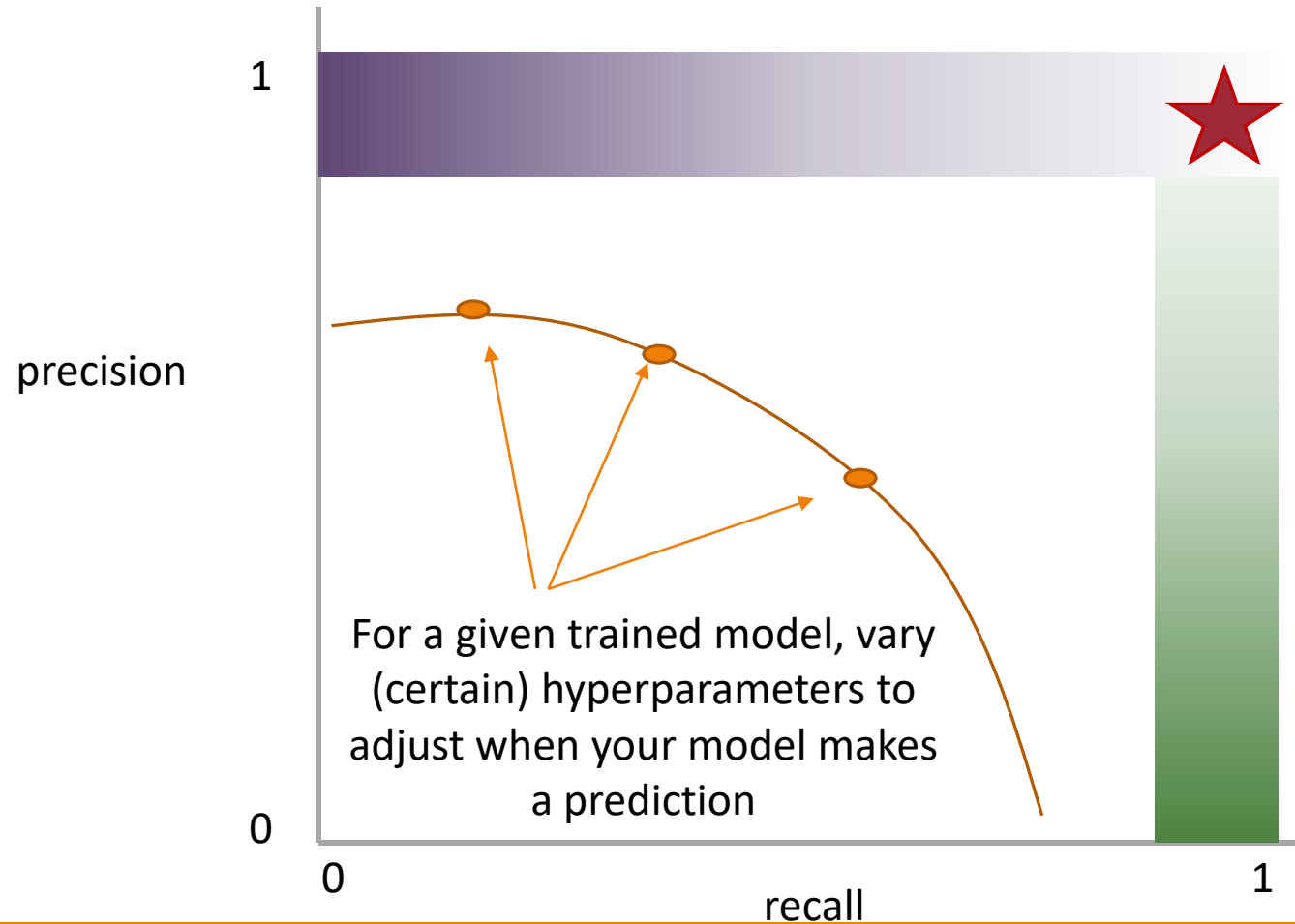| *What label does our system predict? (↓)* | *What is the actual label?* | |
| --- | --- | --- |
| | Actual Target Class ("●") | Not Target Class ("○") |
| **Selected/ Guessed ("●")** | True Positive (TP) ● *Actual* ● *Guessed* | False Positive (FP) ○ *Actual* ● *Guessed* |
| **Not selected/ not guessed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

# Contingency Table (out of table form)

Query:
Articles about dogs

**FN**

Adopting an Animal at the Pound

**TP**

10 Tips for Grooming your Dog

What Foods to Hide from Your Dog During the Holidays

Who's a good boy?

Why Huskies Howl

**TN**

15 Weird Animals You Can Find in North America

**FP**



It's raining cats and dogs

When Cats Rule the World

# Review: Precision and Recall Present a Tradeoff



precision

recall

0     1
0     1

For a given trained model, vary (certain) hyperparameters to adjust when your model makes a prediction

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

precision

recall

1

0

0

1

Improve overall model: push the curve that way

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 <u>balanced</u>

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.

$$\text{macroprecision} = \frac{1}{C}\sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FP}_c} = \frac{1}{C}\sum_c \text{precision}_c$$

$$\text{macrorecall} = \frac{1}{C}\sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FN}_c} = \frac{1}{C}\sum_c \text{recall}_c$$

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

$$\text{microprecision} = \frac{\sum_c \text{TP}_c}{\sum_c \text{TP}_c + \sum_c \text{FP}_c}$$

$$\text{microrecall} = \frac{\sum_c \text{TP}_c}{\sum_c \text{TP}_c + \sum_c \text{FN}_c}$$

when to prefer macroaveraging?

when to prefer microaveraging?

# Macro/Micro Example

Each *class* has equal weight



# Macro-Average



**Class A**

Recall: **87%.**
Precision: 72%.

True "A"

**Class B**

Recall: **33%.**
Precision: **20%.**

True "B"

**Class C**

Recall: **90%.**
Precision: **90%.**

True "C"

**Class D**

Recall: **93%.**
Precision: **100%.**

True "D"

Predicted "A"  Predicted "B"  Predicted "C"  Predicted "D"

**Macro-average**

Recall = (0.87 + 0.33 + 0.9 + 0.93)/4 = **0.76**
Precision = (0.72+0.2+0.9+1)/4=**0.71**

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

# Micro-Average

**All true positives**



True positive "A" — 13
True positive "B" — 1
True positive "C" — 9
True positive "D" — 14

**All false negatives**



False negative "A" — 2
False negative "B" — 4
False negative "C" — 1
False negative "D" — 1

**All false positives**



False positive "B" — 2
False positive "A" — 5
False positive "C" — 1
False positive "D" : 0

Total TP | Total FP | Total FN

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

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

Predicted "A"  Predicted "B"  Predicted "C"  Predicted "D"

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

# Micro- vs Macro-Average

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

$$\text{macroprecision} = \frac{1}{C}\sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FP}_c} = \frac{1}{C}\sum_c \text{precision}_c$$

$$\text{macrorecall} = \frac{1}{C}\sum_c \frac{\text{TP}_c}{\text{TP}_c + \text{FN}_c} = \frac{1}{C}\sum_c \text{recall}_c$$

$$\text{microprecision} = \frac{\sum_c \text{TP}_c}{\sum_c \text{TP}_c + \sum_c \text{FP}_c}$$

$$\text{microrecall} = \frac{\sum_c \text{TP}_c}{\sum_c \text{TP}_c + \sum_c \text{FN}_c}$$

# 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

Predicted

Actual



| Look for ● | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

| Look for ○ | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | True Positive (TP) | False Positive (FP) |
| **Not select/not guessed** | False Negative (FN) | True Negative (TN) |

| Look for ▭ | Actually Target | Actually Not Target |
|---|---|---|
| **Selected/Guessed** | 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 |
|---|---|---|
| Selected/Guessed | 2 | 1 |
| Not select/not guessed | 2 | 4 |

| Look for ○ | Actually Target | Actually Not Target |
|---|---|---|
| Selected/Guessed | 2 | 1 |
| Not select/not guessed | 1 | 5 |

| Look for ▭ | Actually Target | Actually Not Target |
|---|---|---|
| Selected/Guessed | 1 | 2 |
| Not select/not guessed | 1 | 5 |

ML EVALUATION + CLASSIFICATION

# 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 $TP$ ● ?

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

Predicted

Actual

| | Correct Value | | |
|---|---|---|---|
| | ● (orange filled circle) | ○ (orange outline circle) | ▭ (orange rounded rectangle) |
| **Guessed Value** ● | A  2 | B  0 | C  1 |
| ○ | D  1 | E  2 | F  0 |
| ▭ | G  1 | H  1 | I  1 |

How do you compute $TP$ ● ?

# 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** | 🟠 | 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** | | 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 | | |
|---|---|---|---|
|  | ● (filled circle) | ○ (open circle) | ▢ (rounded rectangle) |
| Guessed Value — ● | A 2 | B 0 | C 1 |
| Guessed Value — ○ | D 1 | E 2 | F 0 |
| Guessed Value — ▢ | G 1 | H 1 | I 1 |

How do you compute $FP_{▢}$ ?

# Generalizing the 2-by-2 contingency table

| | Correct Value | | |
|---|---|---|---|
| Q: Is this a good result? | ● (filled circle) | ○ (open circle) | ▭ (rounded rectangle) |
| **Guessed Value** ● (filled circle) | 80 | 9 | 11 |
| ○ (open circle) | 7 | 86 | 7 |
| ▭ (rounded rectangle) | 2 | 8 | 9 |

# Generalizing the 2-by-2 contingency table

| | Correct Value | | |
|---|---|---|---|
| **Guessed Value** 🔵 | 30 | 40 | 30 |
| ⭕ | 25 | 30 | 50 |
| ▭ | 30 | 35 | 35 |

Q: Is this a good result?

# Generalizing the 2-by-2 contingency table

| | Correct Value | | |
|---|---|---|---|
| Q: Is this a good result? | ⬤ | ◯ | ▢ |
| **Guessed Value** ⬤ | 7 | 3 | 90 |
| ◯ | 4 | 8 | 88 |
| ▢ | 3 | 7 | 90 |

# Classification

# 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

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

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

# 💡 Key Take-away 💡

# We will *learn* this
$$p(Y \mid X)$$

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

NLP pg. 477

# Maxent Models for Classification: Discriminatively or …

Directly model the posterior

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

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

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

Posterior (under $P(Y|X)$)

**Posterior:**
probability of event Y with <u>knowledge that X has occurred</u>
NLP pg. 478

**Likelihood:**
probability of event X given that Y <u>has occurred</u>
NLP pg. 478

**Prior:**
probability of event X occurring (regardless of what other events happen)
NLP pg. 478

# Terminology: Posterior Probability

Posterior probability:

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

Conditionally dependent probabilities:

- If $\bullet$ and $\circ$ are the only two options:

$$p(\ \bullet\ |X) + p(\ \circ\ |X) = 1$$

and

$$p(\ \bullet\ |X) \geq 0,\ p(\ \circ\ |X) \geq 0$$

# Posterior Probability with Variables

p( 🟠 |X) vs. p( ⚪ |X)

$$p(\,Y = label_1\,|X)\ \text{vs.}\ p(Y = label_0\,|X)$$

# Maxent Models for Classification: Discriminatively or Generatively Trained

Directly model the posterior

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

**Discriminatively** trained classifier

Model the posterior with Bayes rule

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

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

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

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

Modeled jointly!

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

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

We need to define:

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

- **weights** $\theta$ (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 | x)

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

◦ Denoted by a **general vector of K features**
$$f(x) = (f_1(x), \ldots, f_K(x))$$

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.

TECH

NOT TECH

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

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

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

$$f(x) = \left( f_i(x) \right)_i^V$$

| feature $f_i(x)$ | value |
|---|---|
| Amazon | 1 |
| acquired | 1 |
| behemoth | 1 |
| Bond | 2 |
| … | |
| sniffle | 0 |
| … | |

$f(x)$

$$\begin{bmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 0 \\ \ldots \end{bmatrix}$$

# 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 Tasks**

| 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**

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

# 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

# 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.

Michael Jordan and 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.

extract_and_score$_{\text{Bulls, entailed}}$(📄)      **ENTAILED**

extract_and_score$_{\text{basketball, entailed}}$(📄, ENTAILED)

extract_and_score$_{\text{Chicago, entailed}}$(📄, ENTAILED)

# Score and Combine Our Clues

$score_{1, Entailed}(\text{📄})$

$score_{2, Entailed}(\text{📄})$

$score_{3, Entailed}(\text{📄})$

...

$score_{k, Entailed}(\text{📄})$

...

**COMBINE** ➡️

posterior probability of

ENTAILED

# Scoring Our Clues

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

*(ignore the feature indexing for now)*

$$\text{score}_{1 \text{ , Entailed}}(\text{📄}) \qquad +$$

$$\text{score}_{2 \text{ , Entailed}}(\text{📄}) \qquad +$$

$$\text{score}_{3 \text{ , Entailed}}(\text{📄}) \qquad +$$

...

# Turning Scores into Probabilities

score( [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 ) > score( [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 )

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.] ) > p( 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.] )

KEY IDEA

# Turning Scores into Probabilities (More Generally)

$$score(x, y_1) > score(x, y_2)$$

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

KEY IDEA

ML EVALUATION + CLASSIFICATION

# Maxent Modeling



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

*This must be a probability*

*This could be any real number*

*Convert through function G? What is this function?*

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

# What function G…

operates on any real number?

is never less than 0?

is monotonic? (a < b ➔ G(a) < G(b))

# What function G…

operates on any real number?

is never less than 0?

is monotonic? (a < b ➔ G(a) < G(b))

**G(x) = exp(x)**

# Maxent Modeling

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

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

# Maxent Modeling

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

$$\exp(\ \begin{array}{l} \text{score}_{1,\text{Entailed}}(📄)\ + \\ \text{score}_{2,\text{Entailed}}(📄)\ + \\ \text{score}_{3,\text{Entailed}}(📄)\ + \\ \quad\ldots \end{array}\ ))$$

# Maxent Modeling

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

$$\exp(\ \begin{array}{l}\text{weight}_{1,\,\text{Entailed}} * \text{applies}_1(\text{📄})\ \textbf{+}\\[4pt]\text{weight}_{2,\,\text{Entailed}} * \text{applies}_2(\text{📄})\ \textbf{+}\\[4pt]\text{weight}_{3,\,\text{Entailed}} * \text{applies}_3(\text{📄})\ \textbf{+}\\[4pt]\qquad\qquad\ldots\end{array}\ ))$$

# Maxent Modeling

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

$$\exp\Big(\ \begin{array}{l} \text{weight}_{1,\ \text{Entailed}} * \text{applies}_1(\text{📄})\ \textbf{+} \\ \text{weight}_{2,\ \text{Entailed}} * \text{applies}_2(\text{📄})\ \textbf{+} \\ \text{weight}_{3,\ \text{Entailed}} * \text{applies}_3(\text{📄})\ \textbf{+} \\ \qquad\qquad\qquad\qquad ... \end{array}\ \Big)\Big)$$

| K different weights… | for K different features |
|---|---|

$$\theta \qquad\qquad f(x)$$

$$\begin{bmatrix} .31 \\ -.5 \\ .1 \\ .002 \\ .522 \\ ... \end{bmatrix} \qquad \begin{bmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 0 \\ ... \end{bmatrix}$$

# Maxent Modeling

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

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

...

K different              for K different
weights...                  features

multiplied and then summed

# Maxent Modeling

$$p\left( \text{ENTAILED} \;\middle|\; \boxed{\begin{array}{l} \text{s: Michael Jordan, coach Phil} \\ \text{Jackson and the star cast,} \\ \text{including Scottie Pippen, took} \\ \text{the Chicago Bulls to six} \\ \text{National Basketball Association} \\ \text{championships.} \\ \text{h: The Bulls basketball team is} \\ \text{based in Chicago.} \end{array}} \right) \propto$$

$$\exp\Big( \text{Dot\_product of Entailed weight\_vec feature\_vec(📄)} \Big)$$

K different weights…　　for K different features　　multiplied and then summed
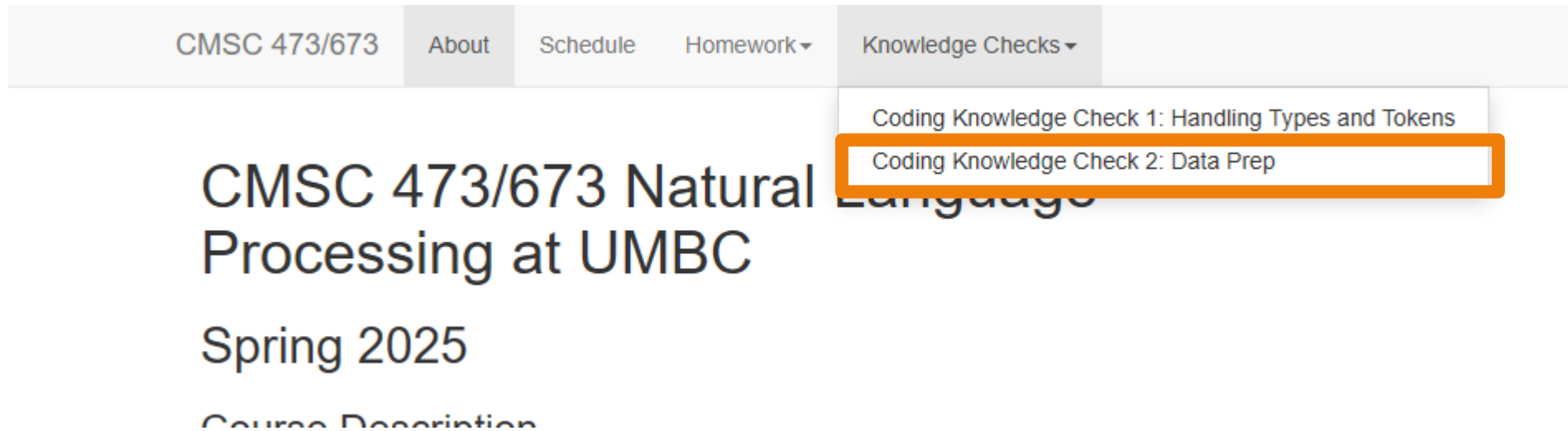
# Maxent Modeling

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

$$\begin{bmatrix} .31\ -.5\ .1\ .002\ .522\ \dots \end{bmatrix} \times \begin{bmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 0 \\ \dots \end{bmatrix}$$

$$\exp(\quad \theta^{T}_{\text{ENTAILED}}\, f(\;📄\;) \quad )$$

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

# Knowledge Check: Data Prep

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