# Features and classification

CSC401/2511 – Natural Language Computing – Spring 2026
Lecture 3 Ken Shi and Gerald Penn
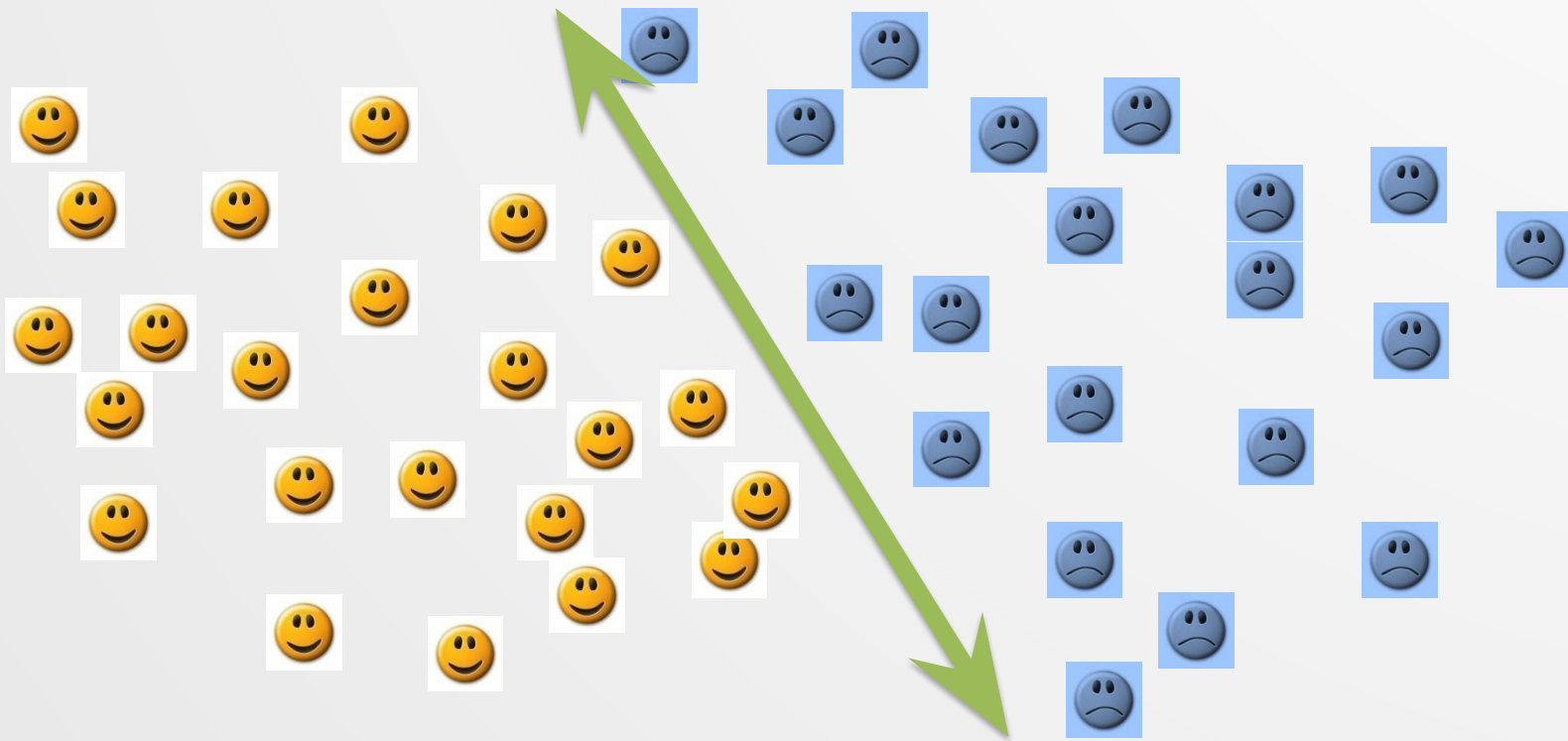University of Toronto

1

# Lecture 3 (topic-wise) overview

- Today:

- **Classification** overview

- Quick introduction to Text Classification

- **Feature extraction** from text.

  - Deep Learning Manifesto: We don't need this

  - Practice: You sure?

- Some slides *may* be based on content from Bob Carpenter, Dan Klein, Roger Levy, Josh Goodman, Dan Jurafsky, and Christopher Manning.

UNIVERSITY OF TORONTO
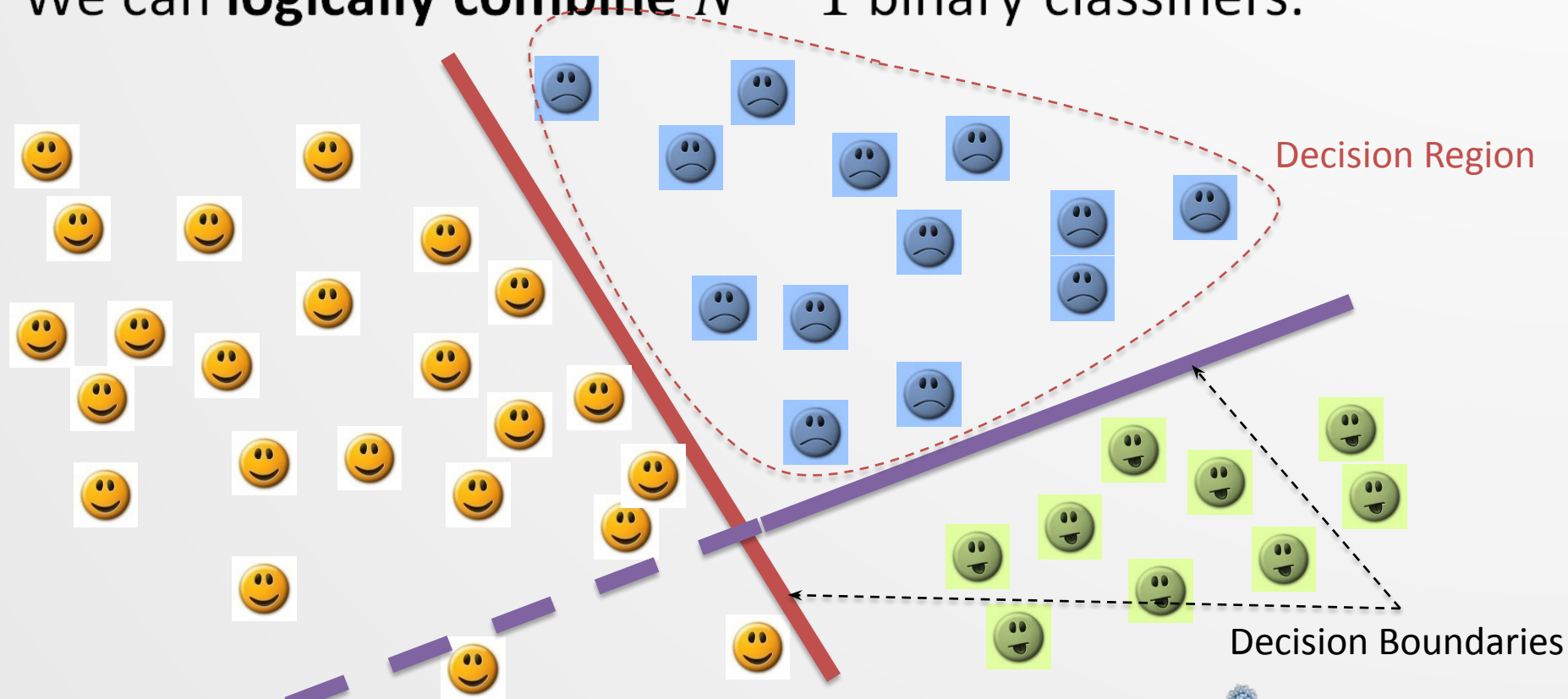
# Classification and Classifier Accuracy

UNIVERSITY OF TORONTO

# Binary and linearly separable

- Perhaps the easiest case.
  - Extends to dimensions $d \geq 3$, line becomes (hyper-)plane.
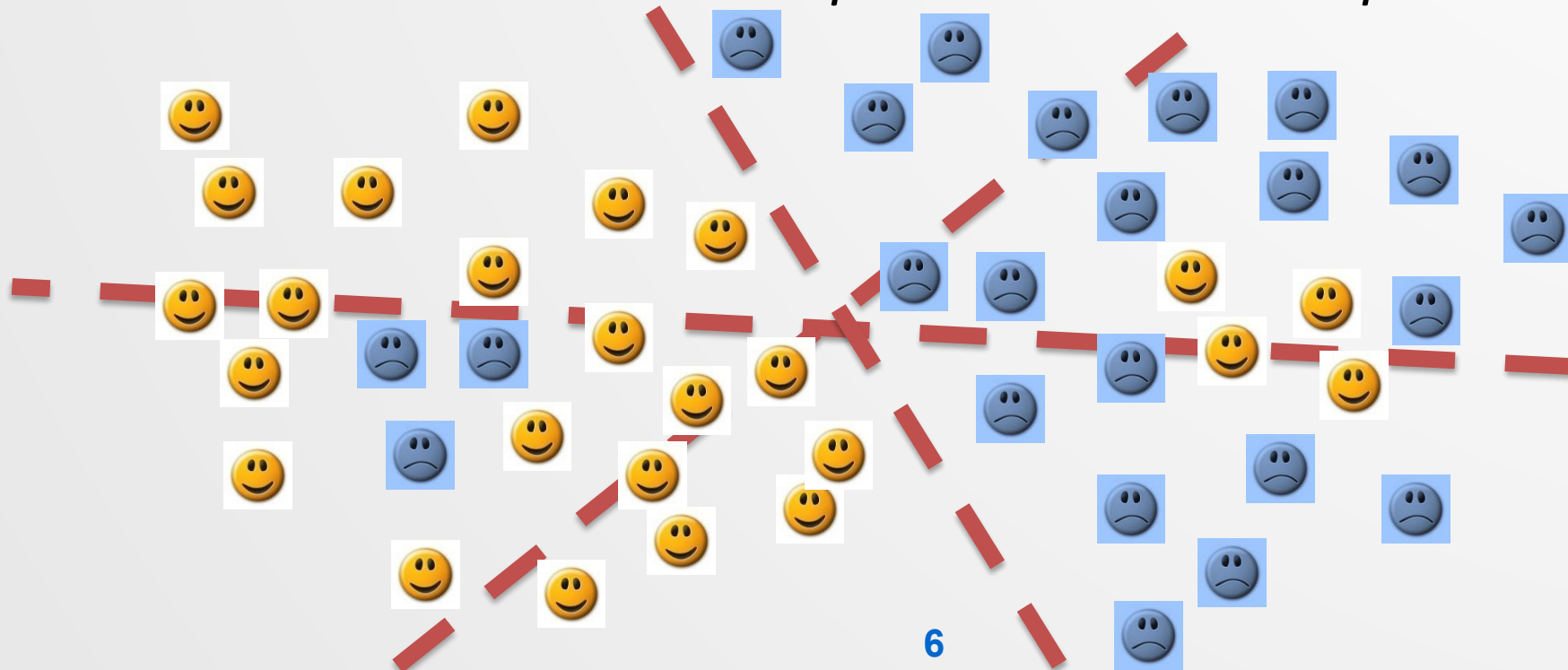
UNIVERSITY OF TORONTO

# *N*-ary and linearly separable

- A bit harder – random guessing gives $\frac{1}{N}$ accuracy (given equally likely classes).
  - We can **logically combine** $N - 1$ binary classifiers.



Decision Region

Decision Boundaries

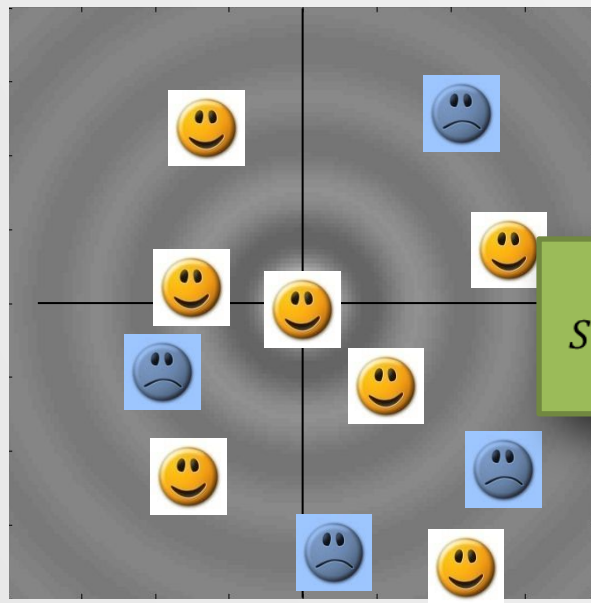UNIVERSITY OF
TORONTO

# Class holes

- Sometimes it can be impossible to draw *any* lines through the data to separate the classes.
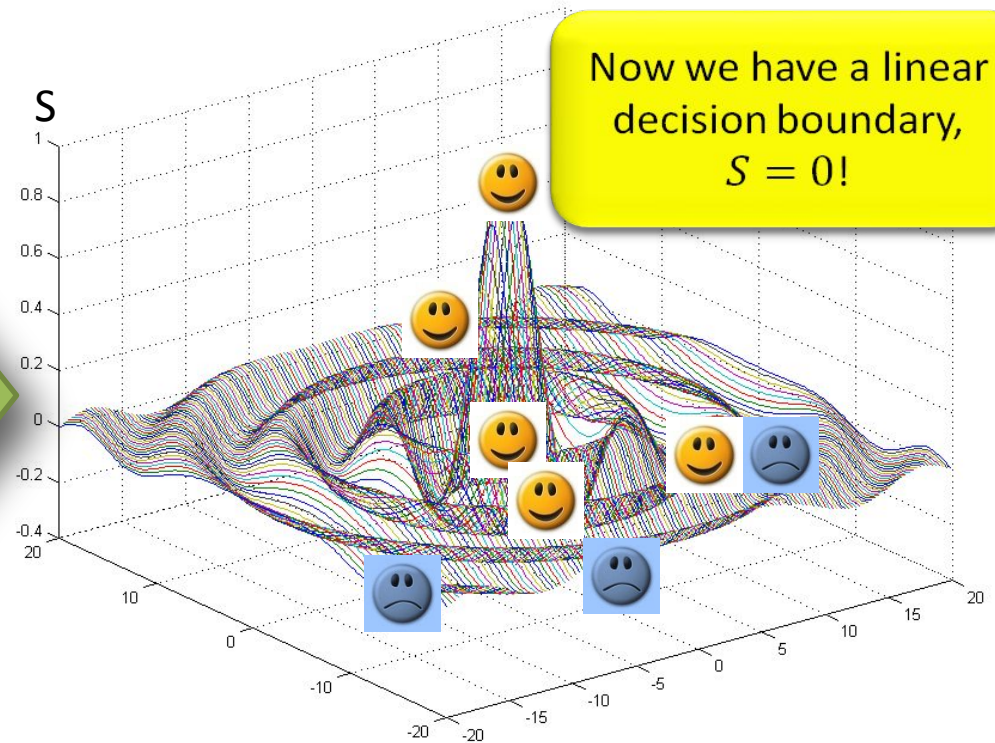  - *Are those troublesome points noise or real phenomena?*

# The kernel trick

- We can sometimes linearize a non-linear case by moving the data into a higher dimension with a **kernel function**. E.g.,



$$S = \frac{\sin\left(\sqrt{x^2 + y^2}\right)}{\sqrt{x^2 + y^2}}$$

Now we have a linear decision boundary, $S = 0$!

UNIVERSITY OF TORONTO

# Precision and Recall

- **Precision**: $\dfrac{N_{\text{relevant \& retrieved}}}{N_{\text{retrieved}}}$
  - Among all retrieved documents, how many are relevant?
  - Precision in machine learning: $\dfrac{TP}{P}$

- **Recall**: $\dfrac{N_{\text{relevant \& retrieved}}}{N_{\text{relevant}}}$
  - Among all relevant documents, how many are retrieved?
  - Recall in machine learning: $\dfrac{TP}{T}$

- Note: Precision and recall has some tradeoff.

# F-measure
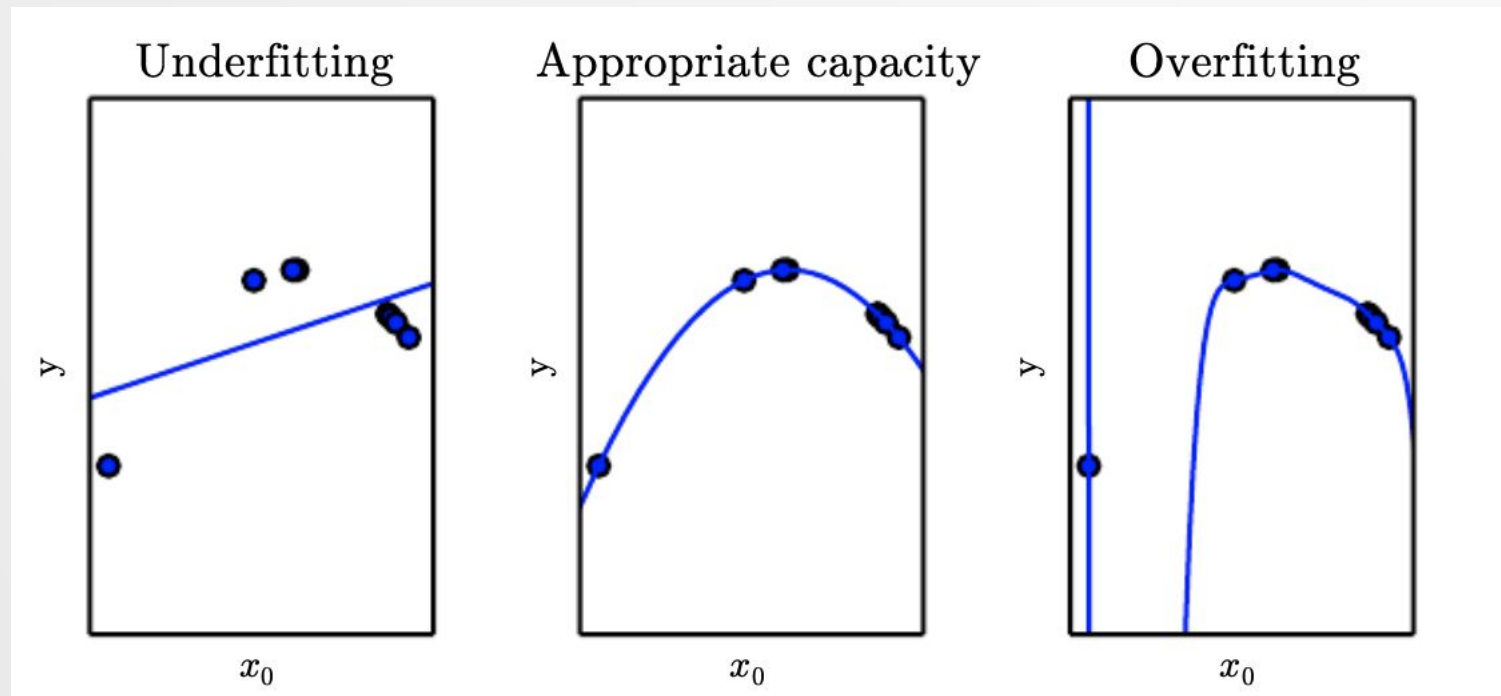
- **F-measure** is the weighted harmonic mean of precision and recall:
  - $F = \dfrac{1}{\alpha\frac{1}{p} + (1-\alpha)\frac{1}{r}}$

- Where $p$ is precision, $r$ is recall, and $\alpha \in [0,1]$.

- Notes:
  - When $\alpha = \frac{1}{2}$, we have $F_1 = \dfrac{2pr}{p+r}$
  - If either of precision or recall is 0 (i.e., true positive count $TP = 0$), then $F$ is arbitrarily set to 0.

# Capacity and over/under-fitting

- A central challenge in machine learning is that our models should **generalize** to unseen data, so we need to set our (hyper-)parameters appropriately.



From Goodfellow

UNIVERSITY OF
TORONTO

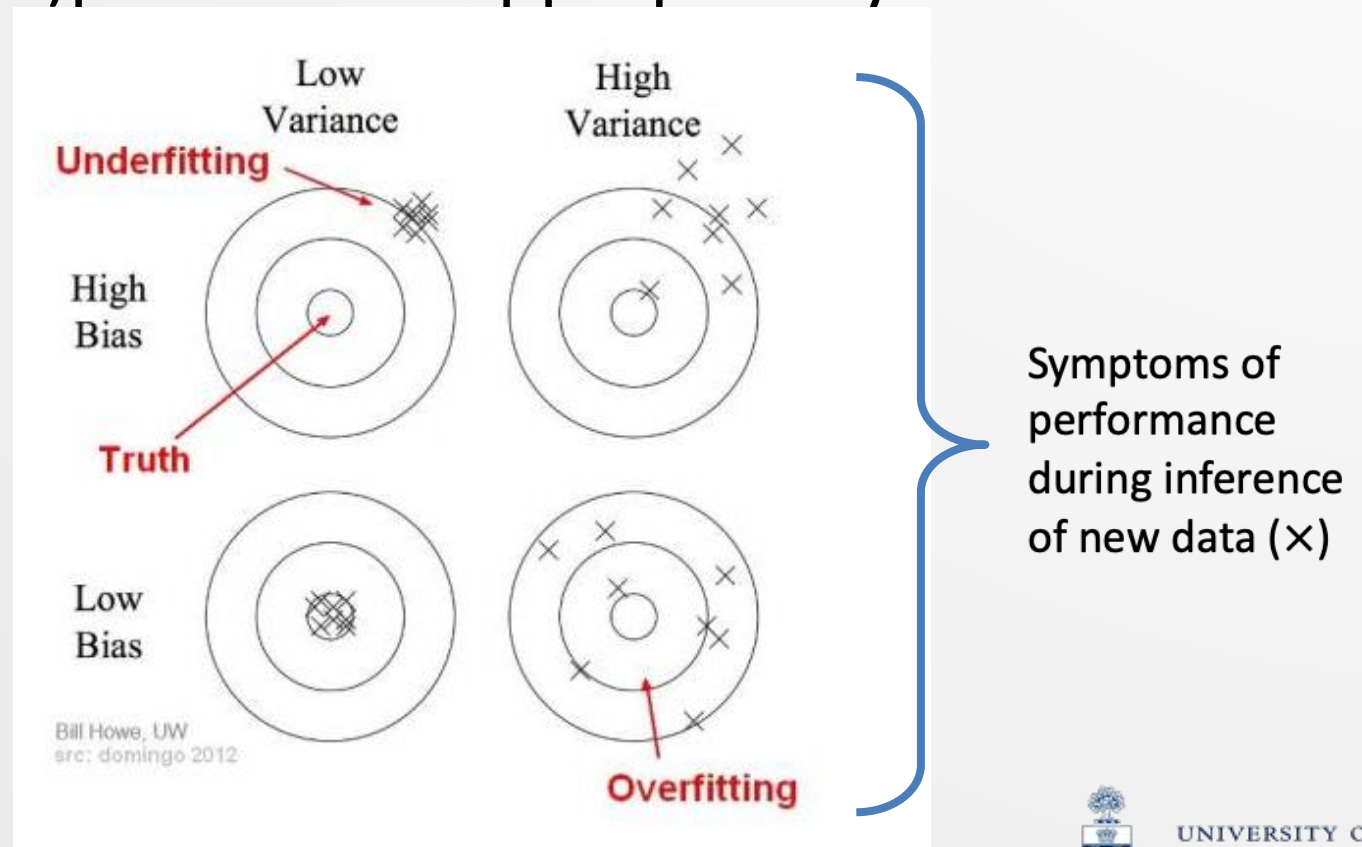# Capacity and over/under-fitting

- A central challenge in machine learning is that our models should **generalize** to unseen data, so we need to set our (hyper-)parameters appropriately.

# Bias and Variance

▸ Even though this analysis only applies to squared error, we often loosely use "bias" and "variance" as synonyms for "underfitting" and "overfitting".

  ○ **Bias**: how wrong the expected prediction is (corresponds to underfitting).

  ○ **Variance**: the amount of variability in the predictions (corresponds to overfitting).

UNIVERSITY OF
TORONTO

# Bias and Variance



High bias



High variance

From Technology Upskilling ML Software Foundations  by Juhan Bae and En-Shiun Annie Lee

UNIVERSITY OF TORONTO

# General process

1. We gather a big and relevant **training** corpus.
2. We learn our **parameters** (e.g., probabilities) from that corpus to build our **model**.
3. Once that model is fixed, we use those probabilities to evaluate **testing** data.

UNIVERSITY OF TORONTO

# General process

- Often, **training data** consist of 80% to 90% of the available data.
  - Often, some subset of *this* is used as a **validation/development set**.

- **Testing data** are **_not_** used for training but often come from the same *corpus*.
  - It often consists of the remaining available data.
  - Sometimes, it's important to **partition** speakers/writers so they **don't** appear in both training and testing.
  - *But what if we just partitioned (un)luckily??*

UNIVERSITY OF TORONTO

# Better process: *K*-fold cross-validation

- ***K*-fold cross validation**: *n.* splitting all data into *K* **partitions** and iteratively testing on each after training on the rest (report means and variances).

| | Part 1 | Part 2 | Part 3 | Part 4 | Part 5 | |
|---|---|---|---|---|---|---|
| **Iteration 1** | ■ | | | | | : Err1 % |
| **Iteration 2** | | ■ | | | | : Err2 % |
| **Iteration 3** | | | ■ | | | : Err3 % |
| **Iteration 4** | | | | ■ | | : Err4 % |
| **Iteration 5** | | | | | ■ | : Err5 % |

**5-fold cross-validation**

| | |
|---|---|
| ■ | **Testing Set** |
| | **Training Set** |

UNIVERSITY OF TORONTO

# (Some) Types of classifiers



- **Generative** classifiers model the data.
  - Parameters set to maximize likelihood of training data.
  - We can *generate* new observations from these.
    - e.g., hidden Markov models

**Vs.**

- **Discriminative** classifiers emphasize **class boundaries**.
  - Parameters set to minimize error on training data.
    - e.g., support vector machines, decision trees.

- …*What do class boundaries look like in the data?*

UNIVERSITY OF TORONTO

# Classification vs. Regression



Classification

Regression

UNIVERSITY OF
TORONTO

# Quick Intro to Text Classification

From Technology Upskilling Machine Learning Software Foundations by En-Shiun Annie Lee

UNIVERSITY OF TORONTO

# Features

- **Feature**: *n.* A measurable **variable** that is (or *should be*) **distinctive** of something we want to model.

- We often choose features to **classify** something.

  - e.g., an emotional, whiny **tone** is likely to indicate that the speaker is not professional, scientific, nor political.

  - Note that in neural networks, e.g., '**features**' refer to something distinctive but often not *nameable.*

- We often need **various, heterogeneous** features to adequately model something,

  e.g. tone plus aspects of grammar.

# Example: Feature vectors

- **Values** for **several** features of an **observation** can be put into a single **vector**.



| | # proper nouns | # 1st person pronouns | # commas |
|---|---|---|---|
| Damien Fahey — Rush Limbaugh looks like if someone put a normal human being in landscape mode. | 2 | 0 | 0 |
| Faux John Madden — BREAKING: Apple Maps projecting Barack Obama to win Brazil. | 5 | 0 | 0 |
| Jim Gaffigan — If there was an award for most pessimistic, I probably wouldn't even be nominated. | 0 | 1 | 1 |

UNIVERSITY OF TORONTO

# Feature vectors

- Features should be useful in **discriminating** between categories.

Table 3: Features to be computed for each text

- Counts:
  — First person pronouns
  — Second person pronouns
  — Third person pronouns
  — Coordinating conjunctions
  — Past-tense verbs
  — Future-tense verbs
  — Commas
  — Colons and semi-colons
  — Dashes
  — Parentheses
  — Ellipses
  — Common nouns
  — Proper nouns
  — Adverbs
  — *wh*-words
  — Modern slang acroynms
  — Words all in upper case (at least 2 letters long)
- Average length of sentences (in tokens)
- Average length of tokens, excluding punctuation tokens (in characters)
- Number of sentences

**Higher** values → this person is referring to themselves (to their opinion, too?)

**Higher** values → looking forward to (or dreading) some future event?

**Lower** values → this tweet is more formal. Perhaps not overly sentimental?

UNIVERSITY OF TORONTO

# Different features for different tasks

- **Alzheimer's disease** involves atrophy in the brain.
  - Excessive **pauses** (acoustic disfluencies),
  - Excessive **word type repetition**, and
  - Simplistic or **short** sentences.
    - '**function words**' like *the* and *an* are often **dropped**.
- To **diagnose** Alzheimer's disease, one might measure:
  - **Proportion** of utterance spent in **silence**.
  - **Entropy** of **word type** usage.
  - **Number** of word **tokens** in a sentence.
    - **Number** of **prepositions** and **determiners** (explained shortly).
      **Explainability/Interpretability!**

UNIVERSITY OF TORONTO

# Features in Sentiment Analysis

- **Sentiment analysis** can involve detecting:

  - **Stress** or **frustration** in a conversation.

  - **Interest**, **confusion**, or **preferences**. Useful to marketers.

    - e.g., *'got socks for xmas wanted #ps5 fml'* 👍

  - **Deceipt**. e.g., *'Let's watch Netflix and chill.'*

- Complicating factors include **sarcasm**, **implicitness**, and a **subtle** spectrum from **negative** to **positive** opinions.

- **Useful features** for sentiment analyzers include:

  - Trigrams.

  - First-person <u>pronouns</u>?

  - Passive voice.

> What does this mean?

> **Pronouns? Voice?**

UNIVERSITY OF
TORONTO

# Features = Dimensions?

- In modern NLP, features don't necessarily correspond to specific dimensions.

- However, this doesn't mean that features no longer exist in them. (more on these models later)

# Pre-processing

- **Pre-processing** involves **preparing** your data to make feature extraction easier or more valid.

  - E.g., **punctuation** likes to press up against words. The sequence " *example,* " should be counted as **two** tokens – not one.
    - We separate the punctuation, as in " *example  ,* ".


- ***There is no perfect pre-processor***. Mutually exclusive approaches can often **both** be justified.

  - E.g., Is *Newfoundland-Labrador* **one** word type or **two**? Each answer has a unique implication for splitting the dash.
  - Often, **noise-reduction** removes *some* information.
  - Being **consistent** is important.

UNIVERSITY OF TORONTO

# Parts of Speech

UNIVERSITY OF TORONTO

# Parts-of-speech (PoS)

- Linguists like to group words according to their **syntactic function** in building sentences.
  - This is similar to grouping Lego by their shapes.

- **Part-of-speech**: *n.* lexical category or morphological class.

> Nouns collectively constitute a part-of-speech (called *Noun*)

UNIVERSITY OF TORONTO

# Example parts-of-speech

| Part of Speech | Description | Examples |
| --- | --- | --- |
| Noun | is usually a **person**, **place**, **event**, or **entity**. | *chair, pacing, monkey, breath.* |
| Verb | is usually an **action** or **predicate**. | *run, debate, explicate.* |
| Adjective | modifies a **noun** to further describe it. | *orange, obscene, disgusting.* |
| Adverb | modifies a **verb** to further describe it. | *lovingly, horrifyingly, often* |

# Example parts of speech

| Part of Speech | Description | Examples |
|---|---|---|
| Preposition | Often specifies aspects of **space**, **time**, or **means**. | *around, over, under, after, before, with* |
| Pronoun | Substitutes for nouns; referent typically understood in context. | *I, we, they* |
| Determiner | logically **quantify** words, usually nouns. | *the, an, both, either* |
| Conjunction | **combines** words or phrases. | *and, or, although* |

UNIVERSITY OF
TORONTO

# Content categories

- Some PoSs convey content labels more than function or linguistic structure.

  - Usually nouns, verbs, adjectives, adverbs.

  - **Content** categories are usually multifarious.

    - e.g., there are more **nouns** than **prepositions**.

  - **New** content words are continually **added**

    e.g., an *app*, *to google, to misunderestimate*.

  - Some **archaic** content words go **extinct**.

    e.g., *fumificate, v., (1721-1792),*
    *frenigerent, adj., (1656-1681),*
    *melanochalcographer, n., (c. 1697).*

UNIVERSITY OF
TORONTO

# Functional categories

- Some PoS are '**glue**' that holds others together.
  - E.g., prepositions, determiners, conjunctions.
  - **Functional** PoS usually cover a **small** and **fixed** number of word types (i.e., a '**closed class**').

  - Their **semantics** depend on the contentful words with which they're used.
    - E.g., *I'm **on** time* vs. *I'm **on** a boat*

UNIVERSITY OF TORONTO

# Grammatical features

- There are several **grammatical features** that can be associated with words:
  - **Case**
  - **Person**
  - **Number**
  - **Gender**

- These features can **restrict** other words in a sentence.

# Other features of nouns

- **Proper noun**:   **named** things (e.g., *"they've killed **Bill**!"*)
- **Common noun**:  **unnamed** things
  (e.g., *"they've killed the **bill**!"*)


- **Mass noun**:   **divisible** and **uncountable**
  (e.g., *"butter"* split in two gives two piles of butter – not two *'butters'*)
- **Count noun**:  **indivisible** and **countable**.
  (e.g., a *"pig"* split in two does not give two pigs)

UNIVERSITY OF TORONTO

# Agreement

- Parts-of-speech **should** match (i.e., **agree**) in certain ways.

- **Articles** 'have' to **agree** with the **number** of their **noun**
  - e.g., "_these pretzels are making me thirsty_" 🙂
  - e.g., "_a winters are coming_" 😛

- **Verbs** 'have' to **agree** (at least) with their **subject** (in English)
  - e.g., "_the dogs eats the gravy_" 😛 **no number** agreement
  - e.g., "_Yesterday, all my troubles seem so far away_" 😛
    **bad tense** – should be past tense _seemed_
  - e.g., "_Can you handle me the way I are?_" 😛

UNIVERSITY OF TORONTO

# Lecture Review Slide

- What are some examples of Text Classification
- **What are features?**
  - What are unique features for the specific tasks of sentiment analysis versus spam detection?
  - What are some words with multiple POS tags?

UNIVERSITY OF
TORONTO

# Tagging

UNIVERSITY OF
TORONTO

# PoS tagging

- **Tagging**: *v.g.* the process of **assigning** a **part-of-speech** to each word in a sequence.

- E.g., using the '**Penn treebank**' tag set (see appendix):

| Word | The | nurse | put | the | sick | patient | to | sleep |
|------|-----|-------|-----|-----|------|---------|----|-------|
| Tag | DT | NN | VBD | DT | JJ | NN | IN | NN |

UNIVERSITY OF TORONTO

# Ambiguities in parts-of-speech

- Word types can have many parts-of-speech.
    - E.g., *back:*
        - *The back/JJ door*            (adjective)
        - *On its back/NN*            (noun)
        - *Win the voters back/RB*      (adverb)
        - *Promise to back/VB you in a fight*    (verb)

- We want to determine the **appropriate** tag for a given *token* in its context.

# Why is tagging useful?

- First step towards many practical purposes.
  - **Speech synthesis**: how to pronounce text
    - *I'm con**TENT**/JJ*      vs. *the **CON**tent/NN*
    - *I ob**JECT**/VBP*      vs. *the **OBJ**ect/NN*
    - *I **lead**/VBP ("l iy d")*    vs. *it's **lead**/NN ("l eh d")*
  - **Information extraction**:
    - Help to find names and relations.
  - **Machine translation**:
    - Help to identify phrase boundaries
  - **Explainability?**

UNIVERSITY OF TORONTO

# Tagging as classification

- We have access to a **sequence of observations** and are expected to decide on the best assignment of a **hidden variable**, i.e., the PoS

| | | | NN | | |
|---|---|---|---|---|---|
| | | | **VB** | | |
| | VBN | | JJ | | **NN** |
| **PRP** | **VBD** | **TO** | RB | **DT** | VB |
| **she** | **promised** | **to** | **back** | **the** | **bill** |

Hidden variable (rows 1–4)

Observation (last row)

UNIVERSITY OF TORONTO

# Reminder: Bayes' Rule



P(X,Y)

P(X)

P(Y)

$$P\left(track \mid football\right) = \frac{P\left(track \cap football\right)}{P\left(football\right)} = \frac{\left(\frac{4}{24}\right)}{\left(\frac{15}{24}\right)} = \frac{4}{15}$$

$$P(X,Y) = P(X)P(Y|X)$$
$$P(X,Y) = P(Y)P(X|Y)$$

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

UNIVERSITY OF TORONTO

# Statistical PoS tagging

- Determine the **most likely** tag sequence $t_{1:n}$ by:

$$\underset{t_{1:n}}{\operatorname{argmax}} P(t_{1:n}|w_{1:n}) = \underset{t_{1:n}}{\operatorname{argmax}} \frac{P(w_{1:n}|t_{1:n})P(t_{1:n})}{P(w_{1:n})}$$

By Bayes' Rule

$$= \underset{t_{1:n}}{\operatorname{argmax}} \frac{P(w_{1:n}|t_{1:n})P(t_{1:n})}{P(w_{1:n})}$$

Only maximize numerator

$$\approx \underset{t_{1:n}}{\operatorname{argmax}} \prod_{i}^{n} P(w_i|t_i)P(t_i|t_{i-1})$$

Assuming independence

Assuming Markov

UNIVERSITY OF
TORONTO

# Those are hidden Markov models!

- We'll see these soon…



HMMs

Image sort of from *2001:A Space Odyssey by MGM pictures*

UNIVERSITY OF
TORONTO

# Word likelihood probability $P(w_i|t_i)$

- **VBZ** (verb, 3$^{rd}$ person singular present) is likely *is*.
- Compute $P(is|VBZ)$ by **counting** in a corpus that has **already** been **tagged**:

$$P(w_i|t_i) = \frac{Count(w_i \text{ tagged as } t_i)}{Count(t_i)}$$

e.g.,

$$P(is|VBZ) = \frac{Count(is \text{ tagged as } VBZ)}{Count(VBZ)} = \frac{10,073}{21,627} = 0.47$$

UNIVERSITY OF TORONTO

# Tag-transition probability $P(t_i | t_{i-1})$

- *Will/MD the/DT **chair/NN** **chair/??** the/DT meeting/NN from/IN that/DT **chair/NN**?*

a)

MD → DT → NN → VB → …
↓ Will   ↓ the   ↓ chair   ↓ chair

b)

MD → DT → NN → NN → …
↓ Will   ↓ the   ↓ chair   ↓ chair

UNIVERSITY OF TORONTO

Let's summarize a few of the classifiers from Assignment 1

UNIVERSITY OF
TORONTO

# Naïve Bayes and SoftMax

- Broadly, Bayesian probability conceives of probability *not* as frequency of some phenomenon occurring, but rather as an expectation related to our own certainty.
- Given an observation $x$, **Naïve Bayes** simply chooses the class $c \in C$ that maximizes $P(c \mid x)$.
  - This can be done in many ways.

$$\underset{c}{\operatorname{argmax}} P(c|x) = \frac{P(c)}{P(x)} P(x|c)$$

Estimate the $P(\cdot)$ using Gaussians, or...

UNIVERSITY OF
TORONTO

# Bayesian Classifier



Given features $\mathbf{x} = [x_1, x_2, \cdots, x_D]^T$
want to compute class probabilities using Bayes Rule:

$$\underbrace{p(c|\mathbf{x})}_{\text{Pr. class given feature}} = \frac{\overbrace{p(\mathbf{x}|c)}^{\text{Pr. feature given class}} \ p(c)}{p(\mathbf{x})}$$

In words,

$$\text{Posterior for class} = \frac{\text{Pr. of feature given class} \times \text{Prior for class}}{\text{Pr. of feature}}$$

To compute $p(c|\mathbf{x})$ we need: $p(\mathbf{x}|c)$ and $p(c)$.

From Technology Upskilling ML Software Foundations by Juhan Bae and En-Shiun Annie Lee

# Independence Assumption

▸ Naive assumption: The features $x_i$ are conditionally independent given the class $c$.
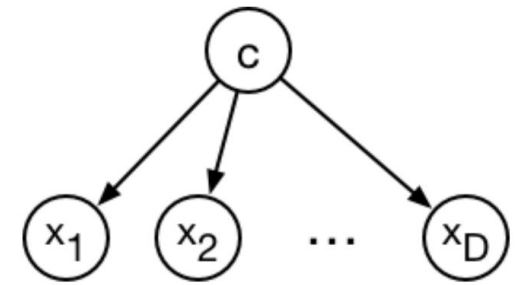
▸ Allows us to decompose the joint distribution:

$$p(c, x_1, \ldots, x_D) = p(c)\, p(x_1|c) \cdots p(x_D|c).$$

▸ Compact representation of the joint distribution.

○ Prior probability of class:

$$p(c = 1) = \pi$$

○ Conditional probability of feature given class:

$$p(x_j = 1|c) = \theta_{jc}$$

UNIVERSITY OF
TORONTO

# Naïve Bayes and SoftMax



- Assume $x \in \mathbb{R}^d$, learning a linear decision boundary is tantamount to learning $W \in \mathbb{R}^{C \times d}$.

P(Class|features) = P(features|Class)*P(Class)

$$\forall c \in C: \boldsymbol{f_c} = W[c, \cdots] \cdot x = \sum_{i=1}^{d} W[c, i] \cdot x[i]$$

Uh oh − $\boldsymbol{f_c}$ can be negative and we want something on $[0,1]$, to be a probability.

Solution: Just raise it with an exponent

**Softmax:**

$$P(y|x) = \frac{\exp(f_y)}{\sum_{c=1}^{C} \exp(f_c)}$$

Naive Bayes:   https://www.youtube.com/watch?v=O2L2Uv9pdDA
SoftMax:   https://www.youtube.com/watch?v=8ps_JEW42xs
Example on Text:   https://www.youtube.com/watch?v=temQ8mHpe3k
Naive Bayes on Spam: https://youtu.be/M59h7CFUwPU
Why Naive Bayes are Cool:   https://www.youtube.com/watch?v=8NEfN3JbINA

UNIVERSITY OF TORONTO

# Naive Bayes Properties

- ▸ An amazingly cheap learning algorithm!

- ▸ **Training time**: Estimate parameters using maximum likelihood.

  - ○ Compute co-occurrence counts of each feature with the labels. I Requires only one pass through the data!

- ▸ **Test time**: Apply Bayes' Rule.

  - ○ Cheap because of the model structure. For more general models, Bayesian inference can be very expensive and/or complicated.

- ▸ Analysis easily extends to prob. distributions other than Bernoulli.

- ▸ Less accurate in practice compared to discriminative models due to its "naive" independence assumption.

UNIVERSITY OF
TORONTO

# Readings

- J&M: 5.1-5.5 (2$^{nd}$ edition)
- M&S: 16.1, 16.4

UNIVERSITY OF TORONTO

# Appendix – prepositions from CELEX

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| of | 540,085 | through | 14,964 | worth | 1,563 | pace | 12 |
| in | 331,235 | after | 13,670 | toward | 1,390 | nigh | 9 |
| for | 142,421 | between | 13,275 | plus | 750 | re | 4 |
| to | 125,691 | under | 9,525 | till | 686 | mid | 3 |
| with | 124,965 | per | 6,515 | amongst | 525 | o'er | 2 |
| on | 109,129 | among | 5,090 | via | 351 | but | 0 |
| at | 100,169 | within | 5,030 | amid | 222 | ere | 0 |
| by | 77,794 | towards | 4,700 | underneath | 164 | less | 0 |
| from | 74,843 | above | 3,056 | versus | 113 | midst | 0 |
| about | 38,428 | near | 2,026 | amidst | 67 | o' | 0 |
| than | 20,210 | off | 1,695 | sans | 20 | thru | 0 |
| over | 18,071 | past | 1,575 | circa | 14 | vice | 0 |

UNIVERSITY OF TORONTO

# Appendix – particles

| | | | | | |
|---|---|---|---|---|---|
| aboard | aside | besides | forward(s) | opposite | through |
| about | astray | between | home | out | throughout |
| above | away | beyond | in | outside | together |
| across | back | by | inside | over | under |
| ahead | before | close | instead | overhead | underneath |
| alongside | behind | down | near | past | up |
| apart | below | east, etc. | off | round | within |
| around | beneath | eastward(s),etc. | on | since | without |

UNIVERSITY OF TORONTO

# Appendix – conjunctions

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| and | 514,946 | yet | 5,040 | considering | 174 | forasmuch as | 0 |
| that | 134,773 | since | 4,843 | lest | 131 | however | 0 |
| but | 96,889 | where | 3,952 | albeit | 104 | immediately | 0 |
| or | 76,563 | nor | 3,078 | providing | 96 | in as far as | 0 |
| as | 54,608 | once | 2,826 | whereupon | 85 | in so far as | 0 |
| if | 53,917 | unless | 2,205 | seeing | 63 | inasmuch as | 0 |
| when | 37,975 | why | 1,333 | directly | 26 | insomuch as | 0 |
| because | 23,626 | now | 1,290 | ere | 12 | insomuch that | 0 |
| so | 12,933 | neither | 1,120 | notwithstanding | 3 | like | 0 |
| before | 10,720 | whenever | 913 | according as | 0 | neither nor | 0 |
| though | 10,329 | whereas | 867 | as if | 0 | now that | 0 |
| than | 9,511 | except | 864 | as long as | 0 | only | 0 |
| while | 8,144 | till | 686 | as though | 0 | provided that | 0 |
| after | 7,042 | provided | 594 | both and | 0 | providing that | 0 |
| whether | 5,978 | whilst | 351 | but that | 0 | seeing as | 0 |
| for | 5,935 | suppose | 281 | but then | 0 | seeing as how | 0 |
| although | 5,424 | cos | 188 | but then again | 0 | seeing that | 0 |
| until | 5,072 | supposing | 185 | either or | 0 | without | 0 |

# Appendix – Penn TreeBank PoS tags

| Tag | Description | Example | Tag | Description | Example |
|-----|-------------|---------|-----|-------------|---------|
| CC | coordin. conjunction | *and, but, or* | SYM | symbol | *+,%, &* |
| CD | cardinal number | *one, two, three* | TO | "to" | *to* |
| DT | determiner | *a, the* | UH | interjection | *ah, oops* |
| EX | existential 'there' | *there* | VB | verb, base form | *eat* |
| FW | foreign word | *mea culpa* | VBD | verb, past tense | *ate* |
| IN | preposition/sub-conj | *of, in, by* | VBG | verb, gerund | *eating* |
| JJ | adjective | *yellow* | VBN | verb, past participle | *eaten* |
| JJR | adj., comparative | *bigger* | VBP | verb, non-3sg pres | *eat* |
| JJS | adj., superlative | *wildest* | VBZ | verb, 3sg pres | *eats* |
| LS | list item marker | *1, 2, One* | WDT | wh-determiner | *which, that* |
| MD | modal | *can, should* | WP | wh-pronoun | *what, who* |
| NN | noun, sing. or mass | *llama* | WP$ | possessive wh- | *whose* |
| NNS | noun, plural | *llamas* | WRB | wh-adverb | *how, where* |
| NNP | proper noun, singular | *IBM* | $ | dollar sign | *$* |
| NNPS | proper noun, plural | *Carolinas* | # | pound sign | *#* |
| PDT | predeterminer | *all, both* | " | left quote | *' or "* |
| POS | possessive ending | *'s* | " | right quote | *' or "* |
| PRP | personal pronoun | *I, you, he* | ( | left parenthesis | *[, (, {, <* |
| PRP$ | possessive pronoun | *your, one's* | ) | right parenthesis | *], ), }, >* |
| RB | adverb | *quickly, never* | , | comma | *,* |
| RBR | adverb, comparative | *faster* | . | sentence-final punc | *. ! ?* |
| RBS | adverb, superlative | *fastest* | : | mid-sentence punc | *: ; ... – -* |
| RP | particle | *up, off* | | | |

# Example – Hero classification

| Hero | Hair length | Height | Age | Hero Type |
|---|---|---|---|---|
| Aquaman | 2" | 6'2" | 35 | Hero |
| Batman | 1" | 5'11" | 32 | Hero |
| Catwoman | 7" | 5'9" | 29 | Villain |
| Deathstroke | 0" | 6'4" | 28 | Villain |
| Harley Quinn | 5" | 5'0" | 27 | Villain |
| Martian Manhunter | 0" | 8'2" | 128 | Hero |
| Poison Ivy | 6" | 5'2" | 24 | Villain |
| Wonder Woman | 6" | 6'1" | 108 | Hero |
| Zatanna | 10" | 5'8" | 26 | Hero |

**Training data** brackets the rows above.

| Test data | Red Hood | 2" | 6'0" | 22 | ? |
|---|---|---|---|---|---|

Characters © DC

UNIVERSITY OF TORONTO