

# Natural Language Processing

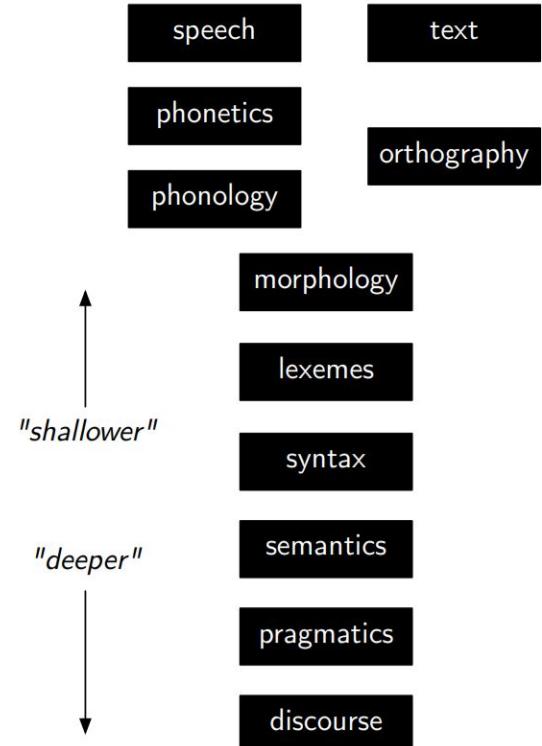
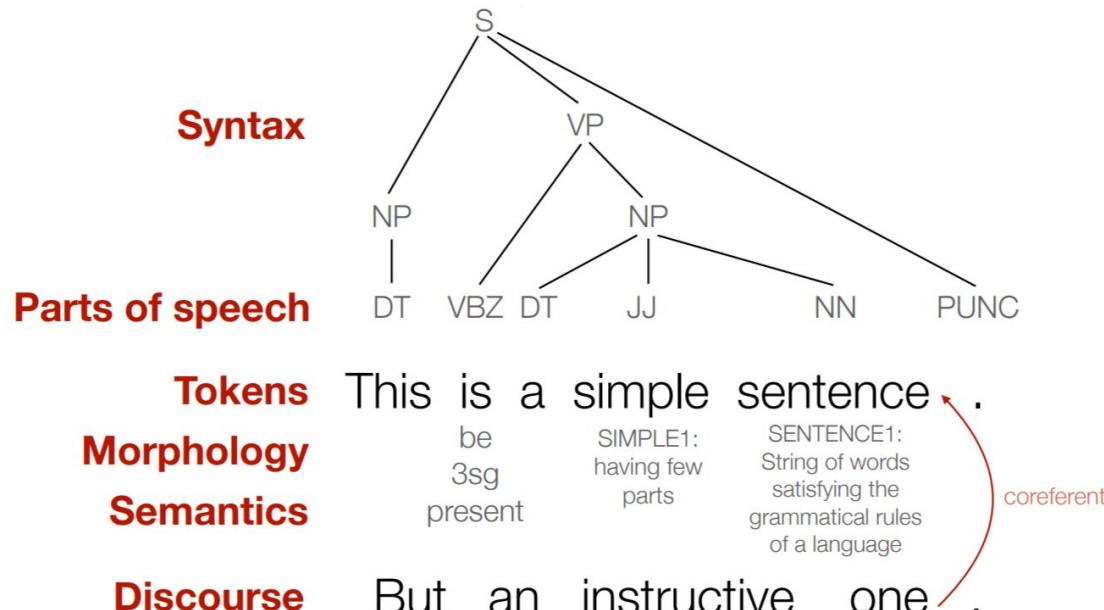
## Introduction, Text classification

Yulia Tsvetkov

[yuliats@cs.washington.edu](mailto:yuliats@cs.washington.edu)

# Announcements

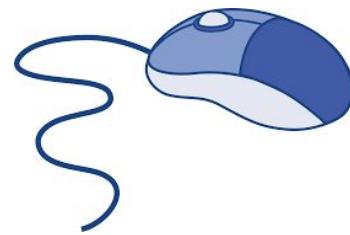
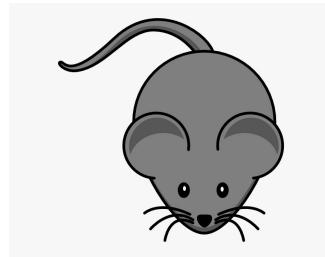
# Language structure & corresponding linguistic subfields



# Why is language interpretation hard?

1. Ambiguity
2. Variation
3. Sparsity
4. Expressivity
5. Unmodeled variables
6. Unknown representation  $\mathcal{R}$

# Ambiguity: word sense disambiguation



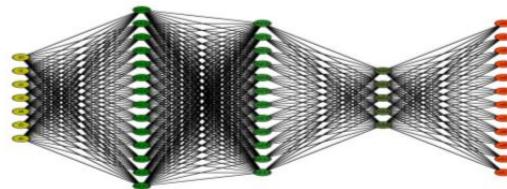
# Ambiguity

- Ambiguity at multiple levels:
  - Word senses: **bank** (finance or river?)
  - Part of speech: **chair** (noun or verb?)
  - Syntactic structure: **I can see a man with a telescope**
  - Multiple: **I saw her duck**



# Dealing with ambiguity

- How can we model ambiguity and choose the correct analysis in context?
  - non-probabilistic methods (FSMs for morphology, CKY parsers for syntax) return **all possible analyses**.
  - probabilistic models (HMMs for part-of-speech tagging, PCFGs for syntax) and algorithms (Viterbi, probabilistic CKY) return **the best possible analysis**, i.e., the most probable one according to the model
  - Neural networks, pretrained language models now provide end-to-end solutions



- But the “best” analysis is only good if our probabilities are accurate. Where do they come from?

# Corpora

- A corpus is a collection of text
  - Often annotated in some way
  - Sometimes just lots of text
- Examples
  - Penn Treebank: 1M words of parsed WSJ
  - Canadian Hansards: 10M+ words of aligned French / English sentences
  - Yelp reviews
  - The Web: billions of words of who knows what



# Why is language interpretation hard?

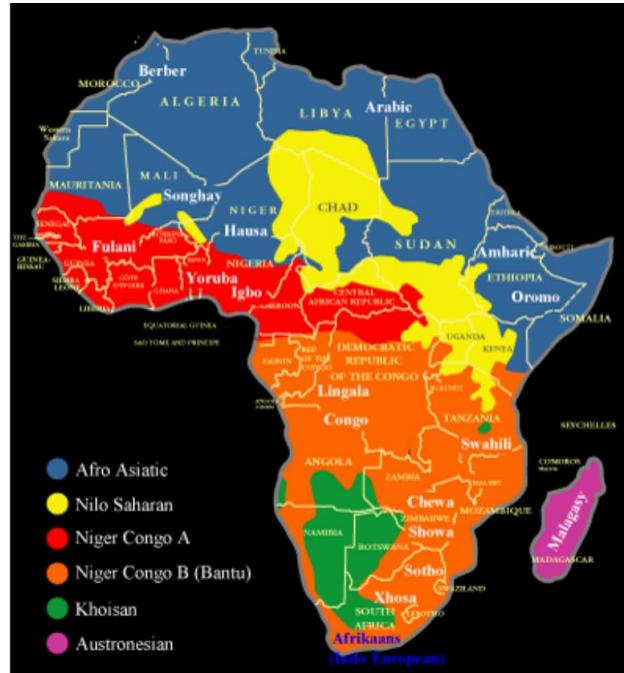
1. Ambiguity
2. Variation
3. Sparsity
4. Expressivity
5. Unmodeled variables
6. Unknown representation  $\mathcal{R}$

# Variation

- ~7K languages
- Thousands of language varieties



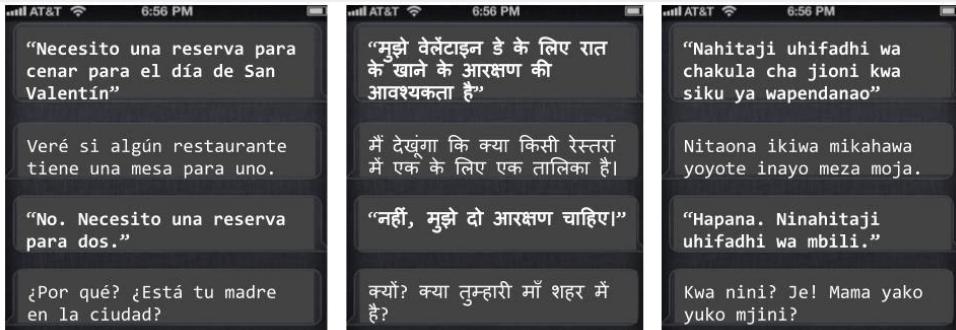
Engishes



Africa is a continent with a very high linguistic diversity: there are an estimated **1.5-2K African languages** from 6 language families. **1.33 billion people**

# NLP beyond English

- ~7,000 languages
- thousands of language varieties



Spanish

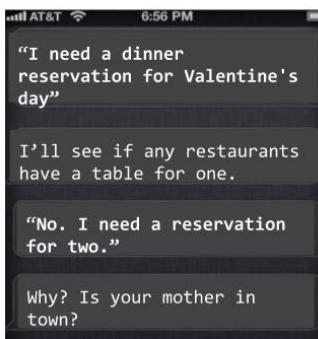
534 million speakers

Hindi

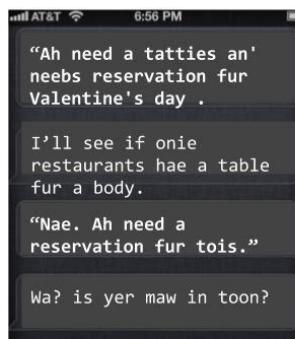
615 million speakers

Swahili

100 million speakers



American English

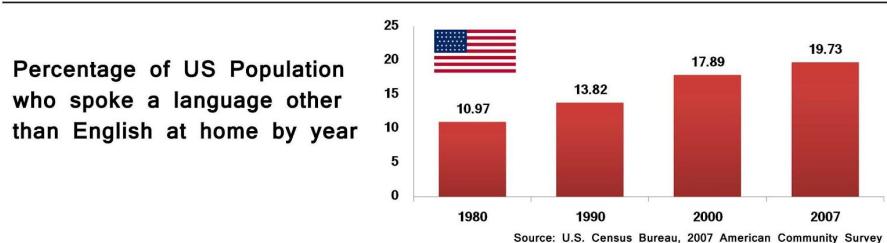
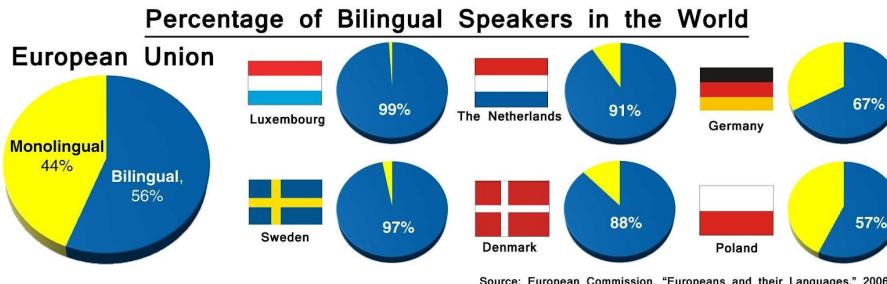


Scottish English



Hinglish

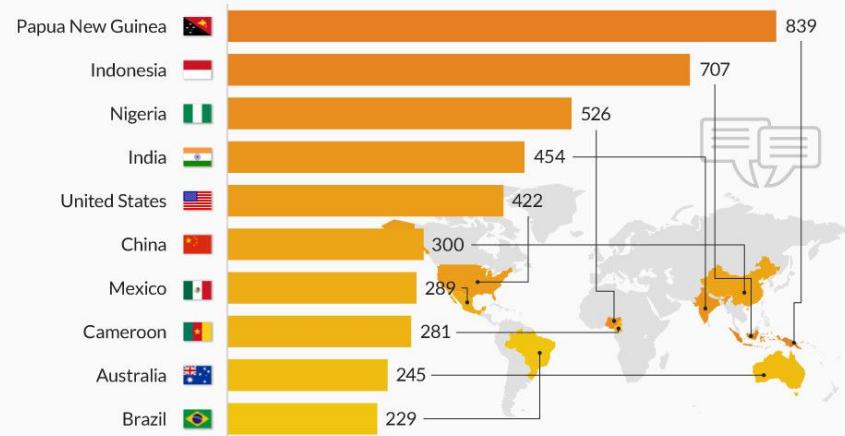
# Most of the world today is multilingual



Source: US Census Bureau

## The Countries With The Most Spoken Languages

Number of living languages spoken per country in 2015



# Tokenization

这是一个简单的句子

WORDS

This is a simple sentence

זה משפט פשוט

# Tokenization + disambiguation

in tea  
her daughter

בתה

in tea	בתה
in the tea	בהתה
that in tea	שבתה
that in the tea	שבחתה
and that in the tea	ושבחתה

- most of the vowels unspecified

and her saturday	+שבת+ה
and that in tea	+ש+בת+ה
and that her daughter	+ש+בת+ה

- most of the vowels unspecified
- particles, prepositions, the definite article, conjunctions attach to the words which follow them
- tokenization is highly ambiguous

# Tokenization + morphological analysis

- Quechua

Much'ananayakapushasqakupuniñataqsunamá

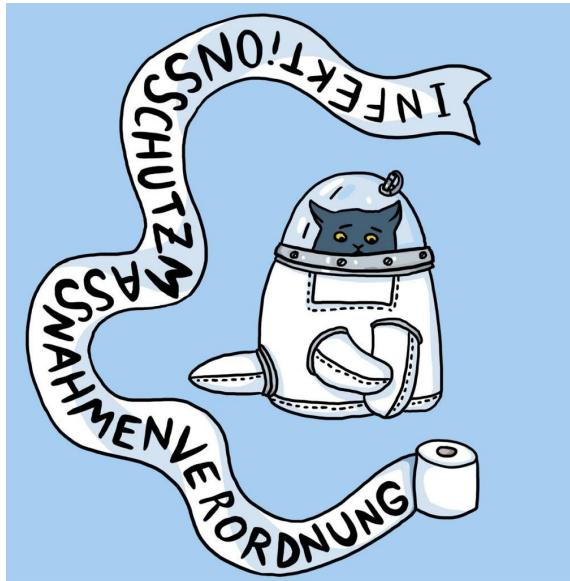
Much'a -na -naya -ka -pu -sha -sqa -ku -puni -ña -taq -suna -má

*"So they really always have been kissing each other then"*

Much'a	to kiss
-na	expresses obligation, lost in translation
-naya	expresses desire
-ka	diminutive
-pu	reflexive (kiss *eachother*)
-sha	progressive (kiss*ing*)
-sqa	declaring something the speaker has not personally witnessed
-ku	3rd person plural (they kiss)
-puni	definitive (really*)
-ña	always
-taq	statement of contrast (...then)
-suna	expressing uncertainty (So...)
-má	expressing that the speaker is surprised

# Tokenization + morphological analysis

- German



Infektionsschutzmaßnahmenverordnung

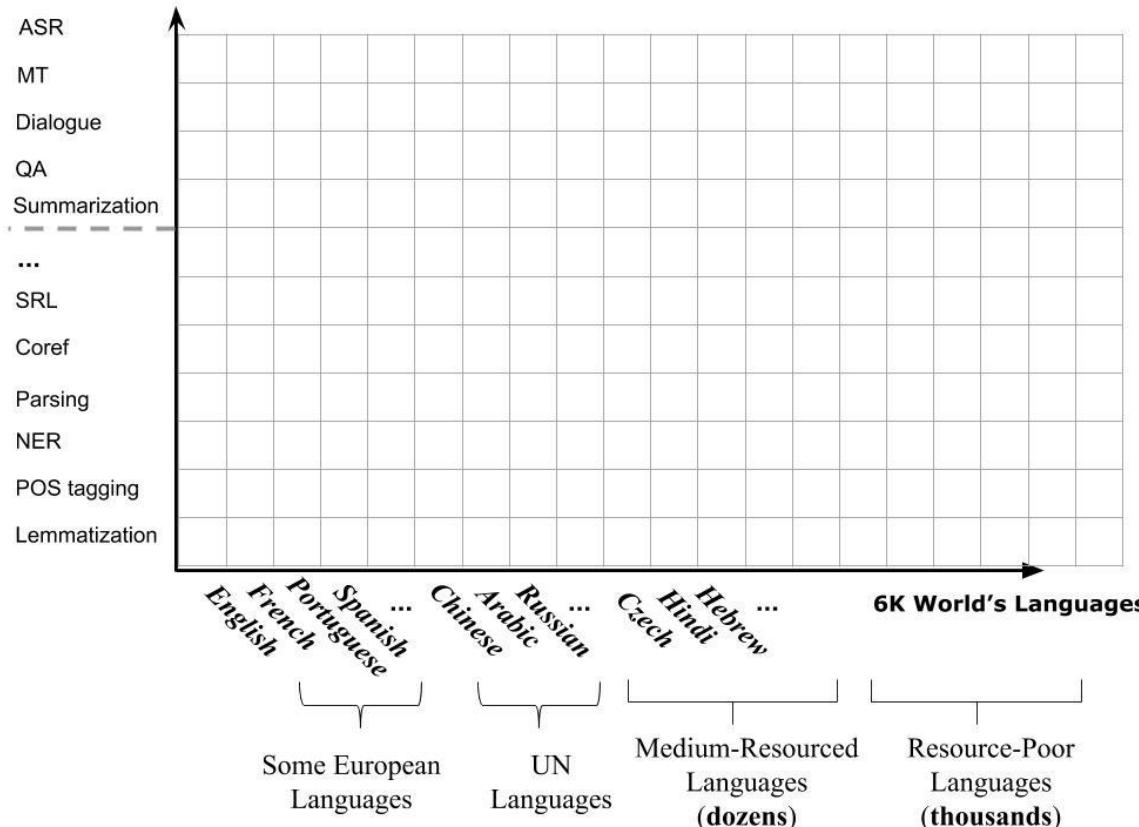
# Semantic analysis

- Every language sees the world in a different way
  - For example, it could depend on cultural or historical conditions



- Russian has very few words for colors, Japanese has hundreds
- Multiword expressions, e.g. [happy as a clam](#), [it's raining cats and dogs](#) or [wake up](#) and metaphors, e.g. [love is a journey](#) are very different across languages

### NLP Technologies/Applications



# Linguistic variation

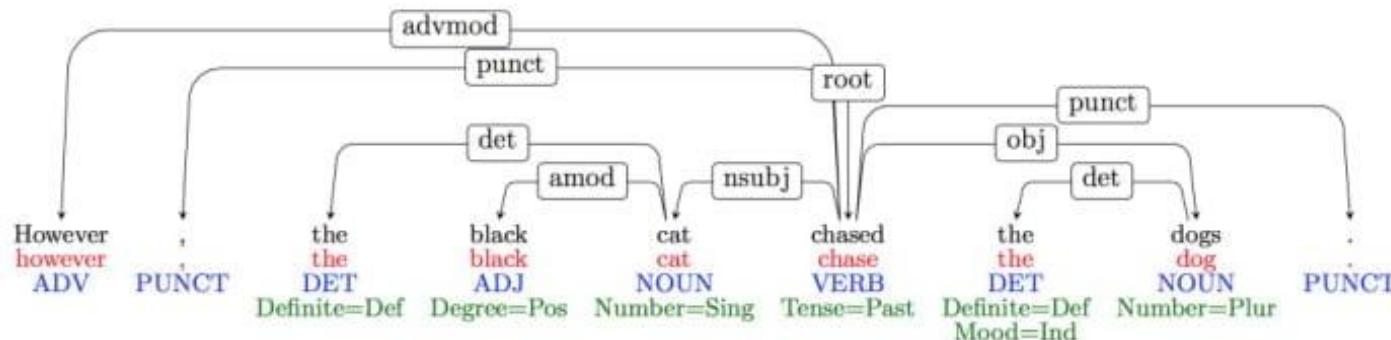
- Non-standard language, emojis, hashtags, names



**chowdownwithchan** #crab and #pork #xiaolongbao at  
@dintaifungusa... where else? 😂🤷‍♀️ Note the cute little  
crab indicator in the 2nd pic 🦀💕

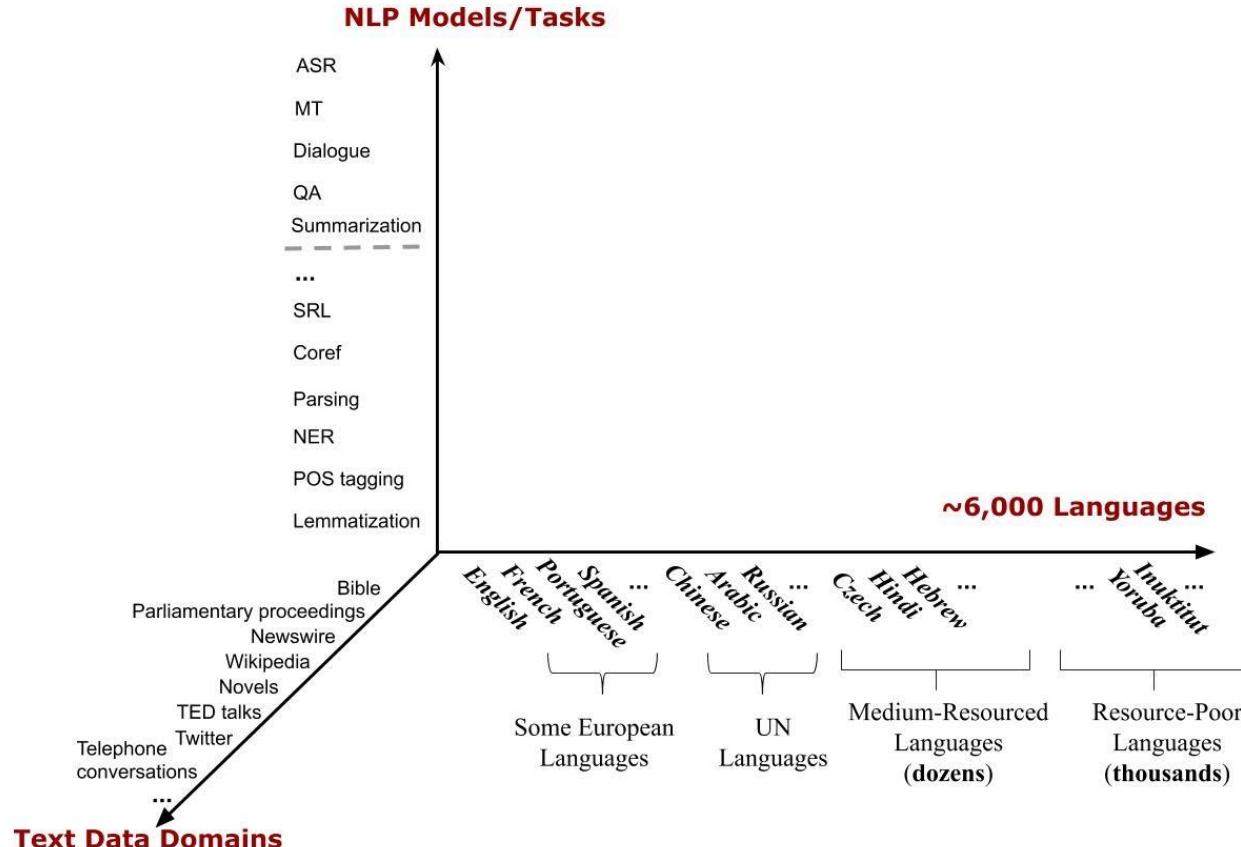
# Variation

- Suppose we train a part of speech tagger or a parser on the Wall Street Journal



- What will happen if we try to use this tagger/parser for social media??

@\_rkptrnte hindi ko alam babe eh, absent ako  
 kanina I'm sick rn hahaha 😊🙌



# Why is language interpretation hard?

1. Ambiguity
2. Scale
3. Variation
4. Sparsity
5. Expressivity
6. Unmodeled variables
7. Unknown representation  $\mathcal{R}$

# Sparsity

Sparse data due to **Zipf's Law**

- To illustrate, let's look at the frequencies of different words in a large text corpus
- Assume “word” is a string of letters separated by spaces

# Word Counts

Most frequent words in the English Europarl corpus (out of 24m word tokens)

any word		nouns	
Frequency	Token	Frequency	Token
1,698,599	the	124,598	European
849,256	of	104,325	Mr
793,731	to	92,195	Commission
640,257	and	66,781	President
508,560	in	62,867	Parliament
407,638	that	57,804	Union
400,467	is	53,683	report
394,778	a	53,547	Council
263,040	I	45,842	States

# Word Counts

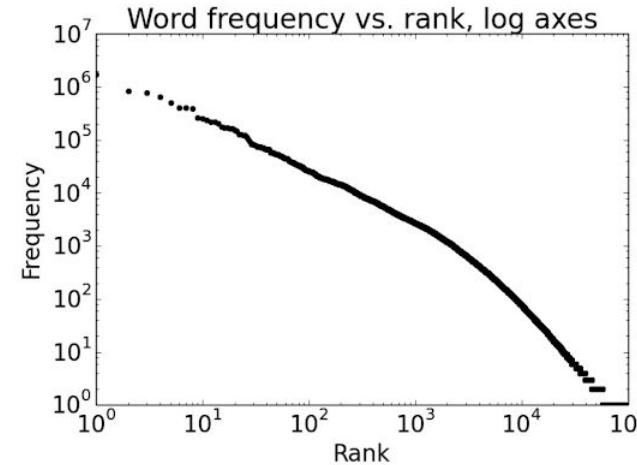
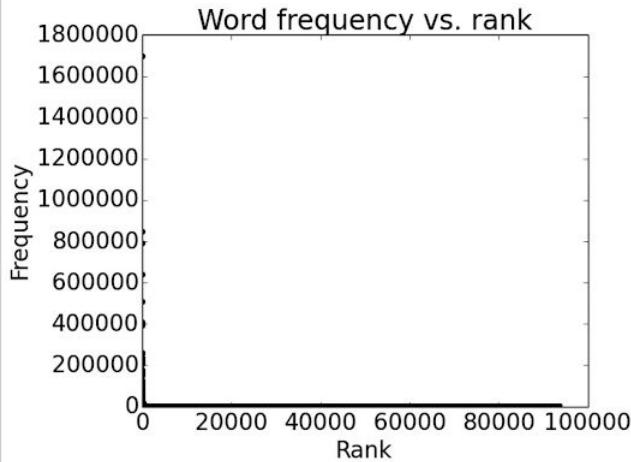
But also, out of 93,638 distinct words (word types), 36,231 occur only once.

Examples:

- cornflakes, mathematicians, fuzziness, jumbling
- pseudo-rapporteur, lobby-ridden, perfunctorily,
- Lycketoft, UNCITRAL, H-0695
- policyfor, Commissioneris, 145.95, 27a

# Plotting word frequencies

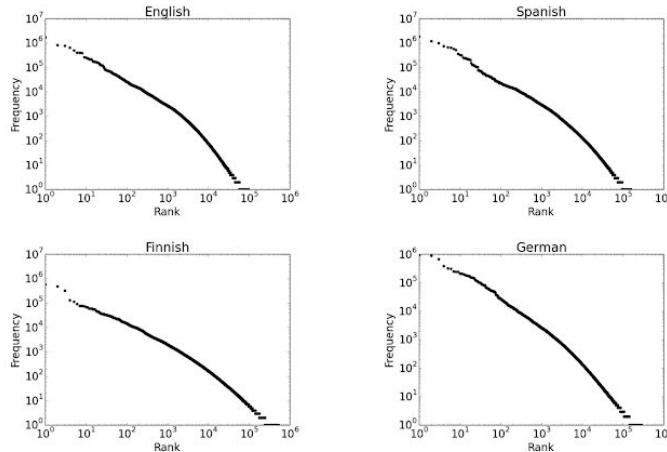
Order words by frequency. What is the frequency of nth ranked word?



# Zipf's Law

## Implications

- Regardless of how large our corpus is, there will be a lot of infrequent (and zero-frequency!) words
- This means we need to find clever ways to estimate probabilities for things we have rarely or never seen



# Why is language interpretation hard?

1. Ambiguity
2. Scale
3. Variation
4. Sparsity
5. Expressivity
6. Unmodeled variables
7. Unknown representation  $\mathcal{R}$

# Expressivity

Not only can one form have different meanings (ambiguity) but the same meaning can be expressed with different forms:

She gave the book to Tom      vs.      She gave Tom the book

Some kids popped by      vs.      A few children visited

Is that window still open?      vs.      Please close the window

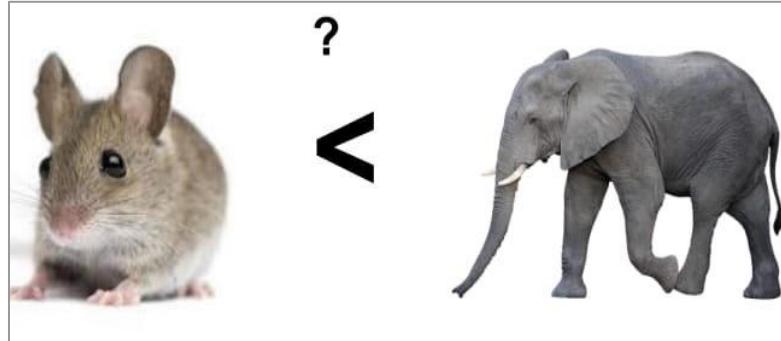
# Why is language interpretation hard?

1. Ambiguity
2. Scale
3. Variation
4. Sparsity
5. Expressivity
6. **Unmodeled variables**
7. Unknown representation  $\mathcal{R}$

# Unmodeled variables



“Drink this milk”



## World knowledge

- I dropped the glass on the floor and it broke
- I dropped the hammer on the glass and it broke

# Why is language interpretation hard?

1. Ambiguity
2. Scale
3. Variation
4. Sparsity
5. Expressivity
6. Unmodeled variables
7. Unknown representation  $\mathcal{R}$

# Unknown representation

- Very difficult to capture what is  $\mathcal{R}$ , since we don't even know how to represent the knowledge a human has/needs:
  - What is the “meaning” of a word or sentence?
  - How to model context?
  - Other general knowledge?

# Desiderata for NLP models

- Sensitivity to a wide range of phenomena and constraints in human language
- Generality across languages, modalities, genres, styles
- Strong formal guarantees (e.g., convergence, statistical efficiency, consistency)
- High accuracy when judged against expert annotations or test data
- Ethical

# Text Classification

# Is this spam?

from: ECRES 2022 <2022@ecres.net> via amazoneses.com  
reply-to: 2022@ecres.net  
to: yuliats@cs.washington.edu  
date: Feb 22, 2022, 7:21 AM  
subject: The Best Renewable Energy Conference ( Last chance ! )  
signed-by: amazoneses.com  
security: Standard encryption (TLS) [Learn more](#)

Dear Colleague,

Account: [yuliats@cs.washington.edu](mailto:yuliats@cs.washington.edu)

Good news: Due to many requests, the submission deadline has been extended to **10 March 2022** (It is firm date).

We would like to invite you to submit a paper to 10. European Conference on Renewable Energy Systems (ECRES). **ECRES 2022 will be held hybrid mode, the participants can present their papers physically or online.** The event is going to be organized in Istanbul/Turkey under the technical sponsorship of Istanbul Medeniyet University and many international institutions. The conference is highly international with the participants from all continents and more than 40 countries.

**The submission deadline and special and regular issue journals can be seen in [ecres.net](http://ecres.net)**

There will be keynote speakers who will address specific topics of energy as you would see at [ecres.net/keynotes.html](http://ecres.net/keynotes.html)

#### CLICK FOR PAPER SUBMISSION

All accepted papers will be published in a special Conference Proceedings under a specific ISBN. Besides, the extended versions will be delivered to reputable journals **indexed in SCI, E-SCI, SCOPUS, and EBSCO**. You can check our previous journal publications from [ecres.net](http://ecres.net). Please note that the official journal of the event, **Journal of Energy Systems** ([dergipark.org.tr/jes](http://dergipark.org.tr/jes)) is also indexed in SCOPUS.

# Spam classification

Dear Colleague,

Account: [yuliats@cs.washington.edu](mailto:yuliats@cs.washington.edu)

Good news: Due to many requests, the submission deadline has been extended to **10 March 2022** (It is firm date).

We would like to invite you to submit a paper to our conference on Renewable Energy Systems (ECRES). **ECRES 2022** will be held online, the participants can present their papers physically or online. The conference is going to be organized in Istanbul/Turkey under the technical support of Marmara Medeniyet University and many international institutions. The conference is international with the participants from all continents and more than 40 countries.

**The submission deadline and specific and regular issue journals can be seen in [ecres.net](http://ecres.net)**

There will be keynote speakers who will address specific topics of energy as you would see at [ecres.net/keynotes.html](http://ecres.net/keynotes.html)

**[CLICK FOR PAPER SUBMISSION](#)**

All accepted papers will be published in a special Conference Proceedings under a specific ISBN. Besides, the extended versions will be delivered to reputable journals **indexed in SCI, E-SCI, SCOPUS, and EBSCO**. You can check our previous journal publications from [ecres.net](http://ecres.net). **Please note that the official journal of the event, Journal of Energy Systems ([dergipark.org.tr/jes](http://dergipark.org.tr/jes)) is also indexed in SCOPUS.**

spam

Invitation to present at the February 2022 Wikimedia Research Showcase



Emily Lescak <[lescak@wikimedia.org](mailto:lescak@wikimedia.org)>

to [yuliats@cs.washington.edu](mailto:yuliats@cs.washington.edu) ▾

Hi Yulia,

My name is Emily Lescak and I am a member of the [Research team](#) at the Wikimedia Foundation. On behalf of the [Wikimedia Research team](#), I would like to invite you to present your research on social biases on Wikipedia at our [Research Showcase](#) in February 2022. This topic fits into our theme for this showcase, which is gaps and biases on Wikipedia.

The [Wikimedia Research Showcase](#) is a monthly, public lecture series where Foundation, academic, and independent researchers share their work related to Wikipedia, Wikimedia, peer production, and open-source software. We focus on topics and projects that we think our audience—a global community of academic researchers, Wikimedians, and Wikimedia Foundation staff—would find interesting and relevant to their work.

Research Showcase presentations are generally 20 minutes long, with an additional 10 minutes for questions and discussion. We invite two presenters to every showcase. Most presenters choose to use slides to present their work.

The February showcase takes place on the 16th at 9:15AM Pacific / 17:15 UTC. If you are unable to attend the live presentation, your presentation will be recorded and also archived for later viewing on the [Wikimedia Foundation's YouTube channel](#).

If this date does not work for you, but you are still interested in giving a showcase presentation, please let me know and we can discuss other options.

I hope to get a chance to see your work presented at the Research Showcase!

Sincerely,

Emily

not spam

# Language ID

Аяны замд түр зогсон тэнгэрийн байдлыг ажиглаад хөдлөх зуур гутал дор шинэхэн орсон цас шаржигнан дуугарч байв. Цасны тухай бодол сонин юм. Хот хүрээ тийш цас орвол орно л биз гэсэн хэнэггүй бодол маань хөдөө талд,.govийн ээрэм хөндийд, малын бэлчээрт, малчдын хотонд болохоор солигдож эргэцүүлэн бodoх нь хачин. Цас хэр орсон бол?

Београд, 16. јун 2013. године – Председник Владе Републике Србије Ивица Дачић честитao је кајакашици златне медаље у олимпијској дисциплини K-1, 500 метара, као и у двоструко дужој стази освојене на првенству Европе у Португалији.

Beograd, 16. jun 2013. godine – Predsednik Vlade Republike Srbije Ivica Dačić čestitao je kajakašici zlatne medalje u olimpijskoj disciplini K-1, 500 metara, kao i u dvostruko dužoj stazi osvojene na prvenstvu Evrope u Portugaliji.

Nestrankarski Urad za vladno odgovornost ZDA je objavil eksplozivno mnenje, da je vlada predsednika Donalda Trumpa kršila zvezno zakonodajo, ko je zadrževala izplačilo kongresno potrjene vojaške pomoći Ukrajini zaradi političnih razlogov. Predstavniški dom kongresa je prav zaradi tega sprožil ustavno obtožbo proti Trumpu.

# Language ID

Аяны замд түр зогсон тэнгэрийн байдлыг ажиглаад хөдлөх зуур гутал дор шинэхэн орсон цас шаржигнан дуугарч байв. Цасны тухай бодол сонин юм. Хот **mongolian** рвол орно л биз гэсэн хэнэггүй бодол маань хөдөө тал **mongolian** өндийд, малын бэлчээрт, малчдын хотонд болохоор солигдож эргэцүүлэн бodoх нь хачин. Цас хэр орсон бол?

Београд, 16. јун 2013. године – Председник Владе Републике Србије **serbian** честитао је кајакашици златне медаље у или **serbian** ини K-1, 500 метара, као и у двоструко дужој стази освојене на првенству Европе у Португалији.

Beograd, 16. jun 2013. godine – Predsednik Vlade Republike Srbije **serbian** je čestitao je kajakašici zlatne medalje u oru **serbian** K-1, 500 metara, kao i u dvostruko dužoj stazi osvojene na prvenstvu Evrope u Portugaliji.

Nestrankarski Urad za vladno odgovornost ZDA je objavil eksplozivno mnenje, да је **vlast predsednika** Donalda Trumpa krшила зvezno zakonodajo, ко је задрžевала izplačilo **k** **slovenian** vojaške помоћи Украјини заради политичних razlogov. Представниšки **d** **slovenian** в заради тega sprožil ustavno obtožbo proti Trumpu.

# Sentiment analysis



By [John Neal](#)

This review is from: Accoutrements Horse Head Mask (Toy)

When I turned State's Witness, they didn't have enough money to put me in the Witness Protection Program, so they bought me this mask and gave me a list of suggested places to move. Since then I've lived my life in peace and safety knowing that my old identity is forever obscured by this life-saving item.



By [Christine E. Tork](#)

Verified Purchase ([What's this?](#))

First of all, for taste I would rate these a 5. So good. Soft, true-to-taste fruit flavors like the sugar variety...I was a happy camper.

BUT (or should I say BUTT), not long after eating about 20 of these all hell broke loose. I had a gastrointestinal experience like nothing I've ever imagined. Cramps, sweating, bloating beyond my worst nightmare. I've had food poisoning from some bad shellfish and that was almost like a skip in the park compared to what was going on inside

# Sentiment analysis



By [John Neal](#)

This review is from: Accoutrements Horse Head Mask (Toy)

When I turned State's Witness, they didn't have enough money to put me in the Witness Protection Program, so they bought me this... and gave me a list of suggested places to move. Since then, I've lived my life in peace and safety knowing that my old identity is forever obscured by this life-saving item.



By [Christine E. Torok](#)

Verified Purchase ([What's this?](#))

First of all, for taste I would rate these a 5. So good. Soft, true-to-taste fruit flavors like the sugar variety...I was a happy camper.

BUT (or should I say BUTT), not long after eating about 20 of these all hell broke loose. I had a gastrointestinal experience like nothing I've ever imagined. Cramps, sweating, floating beyond my worst nightmare. I've had food poisoning from some bad shellfish and that was almost like a skip in the park compared to what was going on inside



# Topic classification

## MEDLINE Article



## MeSH Subject Category Hierarchy

- Antagonists and Inhibitors
- Blood Supply
- Chemistry
- Drug Therapy
- Embryology
- Epidemiology
- ...

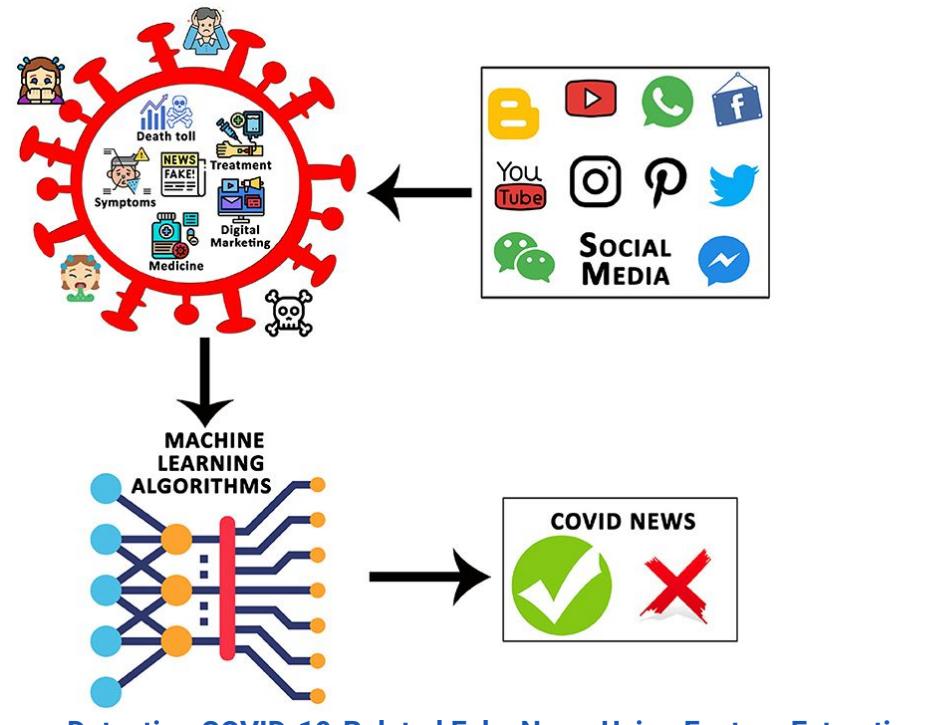
# Authorship attribution: is the author male or female?

By 1925 Vietnam was divided into three parts under French colonial rule. The southern region embracing Saigon and the Mekong delta was the colony Cochin-China; the central area with its imperial capital at Hue was the protectorate of Annam.

Clara never failed to be astonished by the extraordinary felicity of her own name. She found it hard to trust herself to the mercy of fate, which had managed over the years to convert her greatest shame into one of the greatest assets...

S. Argamon, M. Koppel, J. Fine, A. R. Shimoni, 2003. "Gender, Genre, and Writing Style in Formal Written Texts," *Text*, volume 23, number 3, pp. 321–346

# Fact verification: trustworthy or fake?



# Text classification

- We might want to categorize the **content** of the text:
  - Spam detection (binary classification: spam/not spam)
  - Sentiment analysis (binary or multiway)
    - movie, restaurant, product reviews (pos/neg, or 1-5 stars)
    - political argument (pro/con, or pro/con/neutral)
    - Topic classification (multiway: sport/finance/travel/etc)
  - Language Identification (multiway: languages, language families)
  - ...
- Or we might want to categorize the **author** of the text (authorship attribution)
  - Human- or machine generated?
  - Native language identification (e.g., to tailor language tutoring)
  - Diagnosis of disease (psychiatric or cognitive impairments)
  - Identification of gender, dialect, educational background, political orientation (e.g., in forensics [legal matters], advertising/marketing, campaigning, disinformation)
  - ...

# Text classification



Goal: create a function  $f$  that makes a prediction  $\hat{y}$  given an input  $x$

# Over the next couple of classes, we'll investigate:

1. How do we “digest” text into a form usable by a function?

(Keywords for this section: features, feature extraction, feature selection, representations)

2. What kinds of strategies might we use to create our function  $f$ ?

(Keyword for this section: models)

3. How do we evaluate our function  $f$ ?

(Keyword for this section: ... evaluation)



# How do we “digest” text into a form usable by a function?

# Classification: features (measurements)

- Perform measurements and obtain features



4.2, 212, 3.4, 1332  
↓      ↓      ↓      ↓  
diameter, weight, softness, color



5.2, 315, 5.7, 4567  
↓      ↓      ↓      ↓  
diameter, weight, softness, color

# Text classification – feature extraction

What can we measure over text? Consider this movie review:

I love this movie! It's sweet, but with satirical humor. The dialogue is great, and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it just to about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it before.

# Text classification – feature extraction

What can we measure over text? Consider this movie review:

I **love** this movie! It's **sweet**, but with **satirical humor**. The dialogue is **great**, and the adventure scenes are **fun**... It manages to be **whimsical** and **romantic** while **laughing** at the conventions of the fairy tale genre. I would **recommend** it just to about anyone. I've seen it **several** times, and I'm always happy to see it **again** whenever I have a friend who hasn't seen it before.

# Text classification – feature extraction

I **love** this movie! It's **sweet**, but with **satirical humor**. The dialogue is **great**, and the adventure scenes are **fun**... It manages to be **whimsical** and **romantic** while **laughing** at the conventions of the fairy tale genre. I would **recommend** it just to about anyone. I've seen it **several** times, and I'm always happy to see it **again** whenever I have a friend who hasn't seen it before.

(almost) the entire lexicon

word	count	relative frequency
love	10	0.0007
great	...	
recommend		
laugh		
happy		
...		
several		
boring		
...		

# Types of textual features

- Words
  - content words, stop-words
  - punctuation? tokenization? lemmatization? lowercase?
- Word sequences
  - bigrams, trigrams, n-grams
- Grammatical structure, sentence parse tree
- Words' part-of-speech
- Word vectors
- ...

# Possible representations for text

- Bag-of-Words (BOW)
  - Easy, no effort required
  - Variable size, ignores sentential structure
- Hand-crafted features
  - Full control, can use NLP pipeline, class-specific features
  - Over-specific, incomplete, makes use of NLP pipeline
- Learned feature representations
  - Can learn to contain all relevant information
  - Needs to be learned

# Bag-of-Words (BOW)

- Given a document  $d$  (e.g., a movie review) – how to represent  $d$  ?

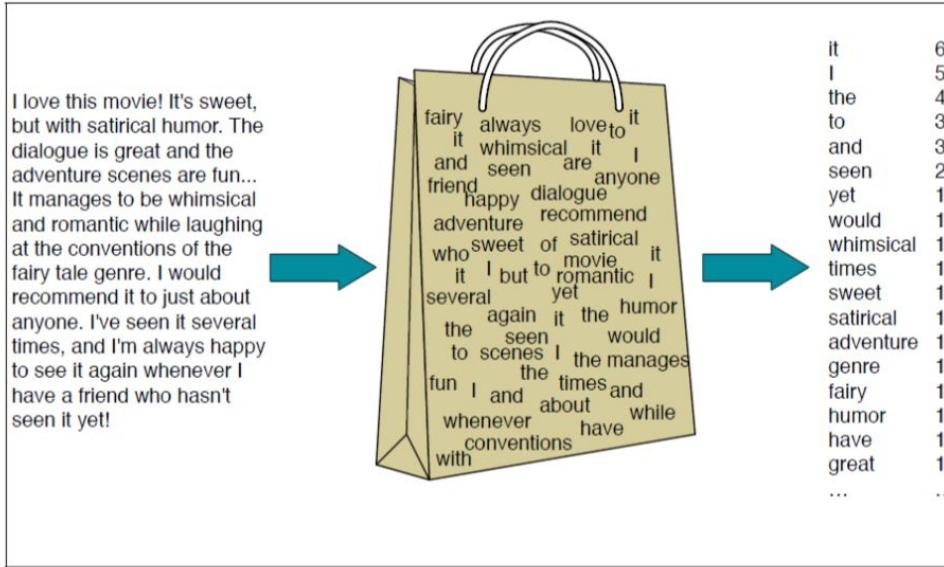
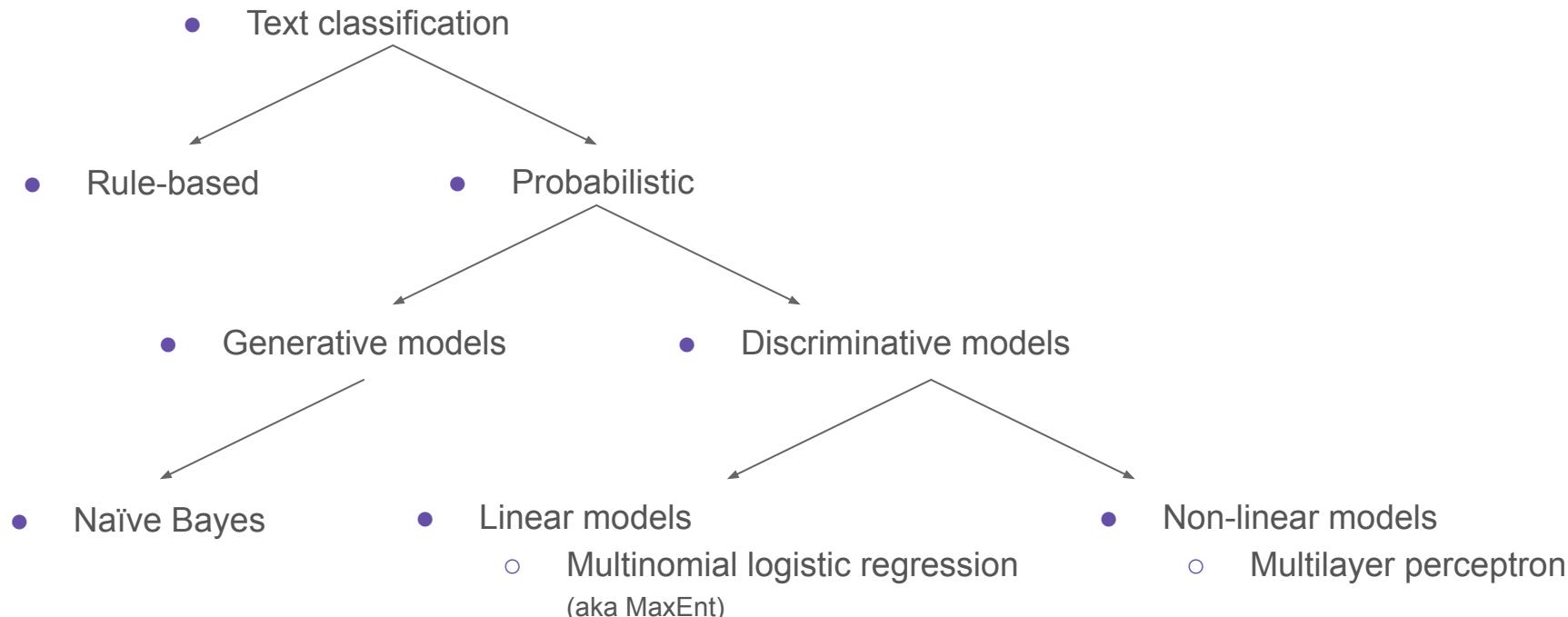


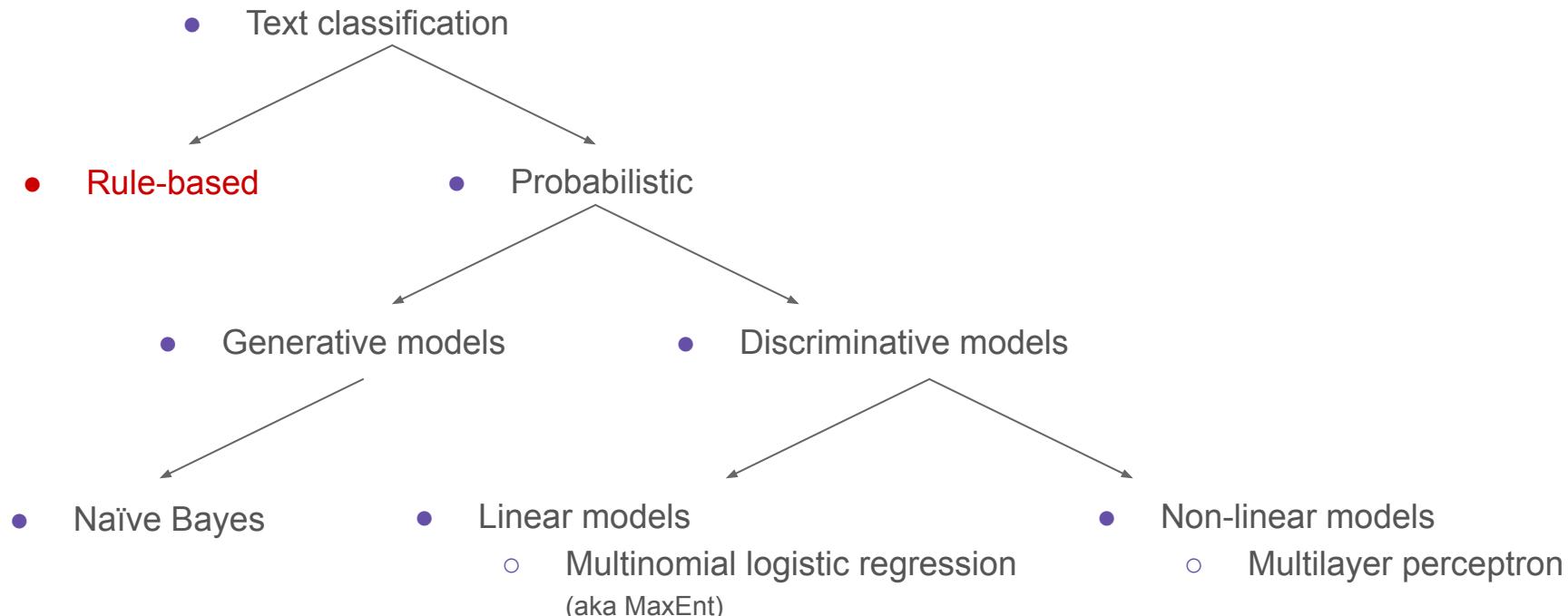
Figure from J&M 3rd ed. draft, sec 7.1

What kinds of strategies might we use  
to create our function  $f$ ?

# We'll consider alternative models for classification



# We'll consider alternative models for classification



# Rule-based classifier

```
def classify_sentiment(document):
    for word in document:
        if word in {"good", "wonderful", "excellent"}:
            return 5
        if word in {"bad", "awful", "terrible"}:
            return 1
```

# Rule-based classification: challenges

Sentiment: Half submarine flick, half ghost story, all in one a criminally neglected film.

# Rule-based classification: challenges

Sentiment: Half submarine flick, half ghost story, all in one a criminally neglected film.

→ hard to identify a priori which words are informative (and what information they carry!)

# Rule-based classification: challenges

Sentiment: Half submarine flick, half ghost story, all in one a criminally neglected film.

→ hard to identify *a priori* which words are informative (and what information they carry!)

Sentiment: It's not life-affirming, it's vulgar, it's mean, but I liked it.

# Rule-based classification: challenges

Sentiment: Half submarine flick, half ghost story, all in one a criminally neglected film.

→ hard to identify *a priori* which words are informative (and what information they carry!)

Sentiment: It's not life-affirming, it's vulgar, it's mean, but I liked it.

→ language pragmatics is complex to model at word level, word order (syntax) matters, but hard to encode in rules!

# Rule-based classification: challenges

Sentiment: Half submarine flick, half ghost story, all in one a criminally neglected film.

→ hard to identify a priori which words are informative (and what information they carry!)

Sentiment: It's not life-affirming, it's vulgar, it's mean, but I liked it.

→ language pragmatics is complex to model at word level, word order (syntax) matters, but hard to encode in rules!

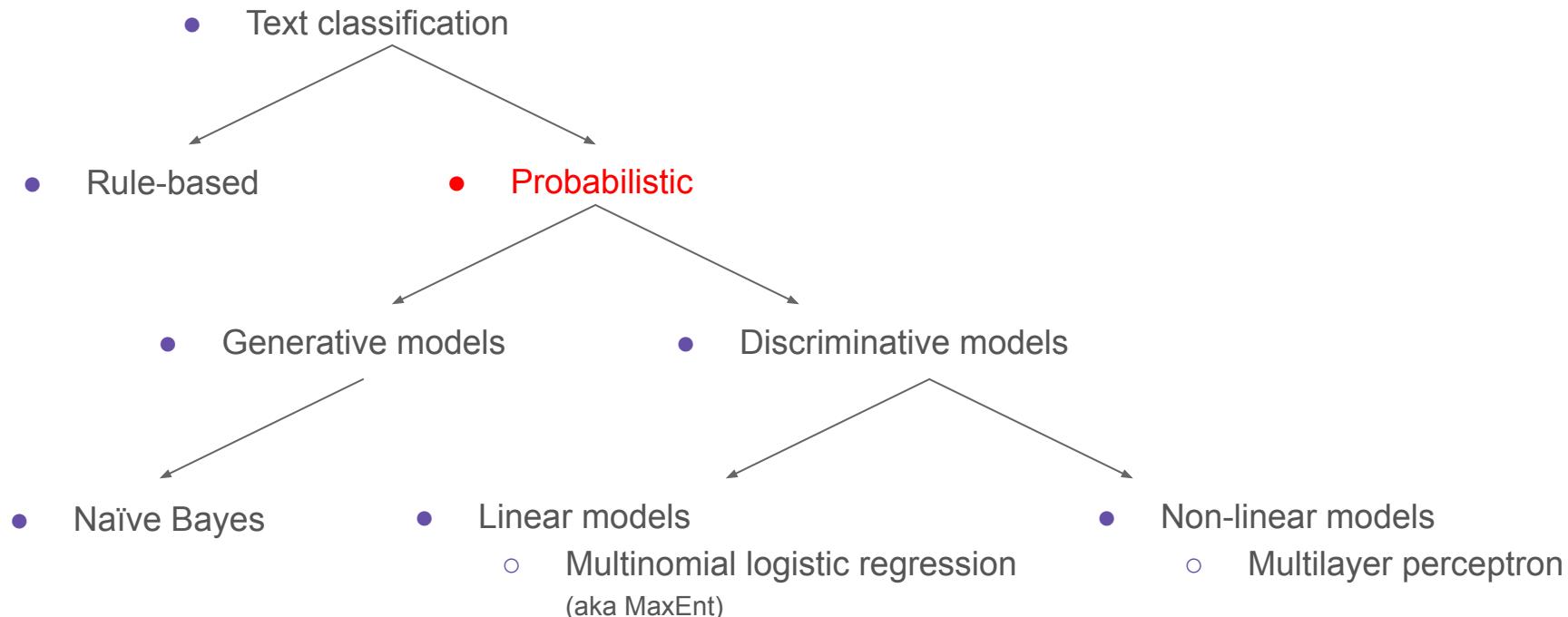
Language ID: All falter, stricken in kind.

→ simple features can be misleading!

# Rule-based classification

But don't forget: if you don't have access to data, speaker intuition and a bit of coding get you pretty far!

# We'll consider alternative models for classification



# Learning-based classification



pick the function  $f$  that does “best” on training data

Goal: ~~create a function  $f$  that makes a prediction  $\hat{y}$  given an input  $x$~~

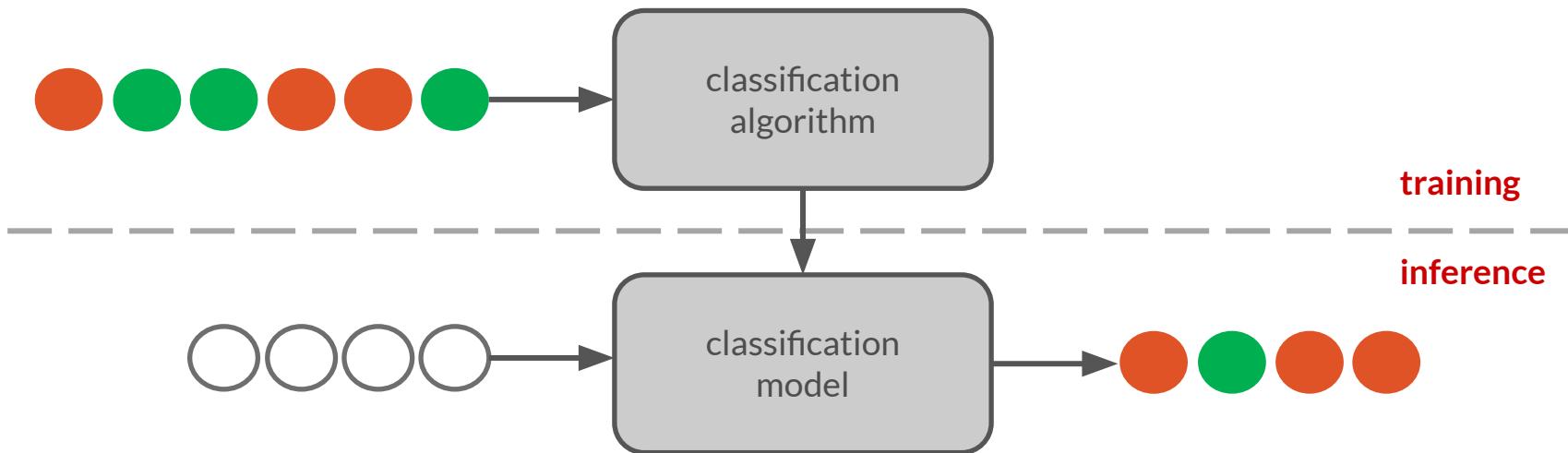
# Classification: learning from data

- Supervised
  - labeled examples
    - Binary (true, false)
    - Multi-class classification (politics, sports, gossip)
    - Multi-label classification (#party #FRIDAY #fail)
- Unsupervised
  - no labeled examples
- Semi-supervised
  - labeled examples + non-labeled examples
- Weakly supervised
  - heuristically-labeled examples

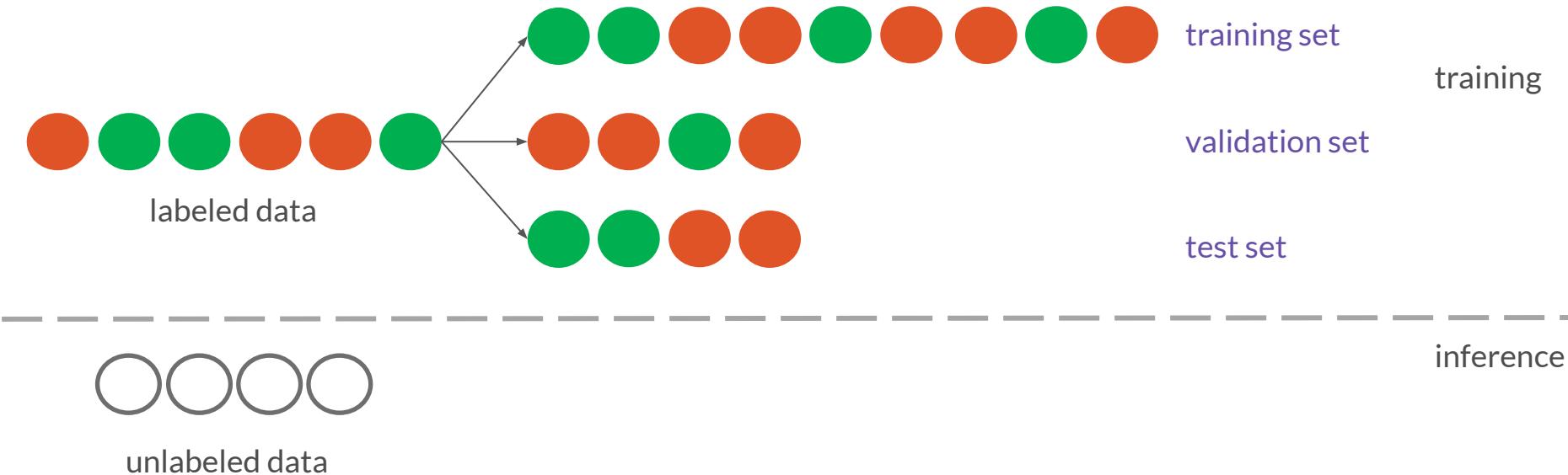
# Where do datasets come from?

Human institutions	Noisy labels	Expert annotation	Crowd workers
Government proceedings	Domain names	Treebanks	Question answering
Product reviews	Link text	Biomedical corpora	Image captions

# Supervised classification



# Training, validation, and test sets



# Supervised classification: formal setting

- Learn a **classification model** from labeled data on
  - properties (“**features**”) and their importance (“**weights**”)
- **X**: set of attributes or features  $\{x_1, x_2, \dots, x_n\}$ 
  - e.g. fruit measurements, or word counts extracted from an input documents
- **y**: a “class” label from the label set  $Y = \{y_1, y_2, \dots, y_k\}$ 
  - e.g., fruit type, or spam/not spam, positive/negative/neutral

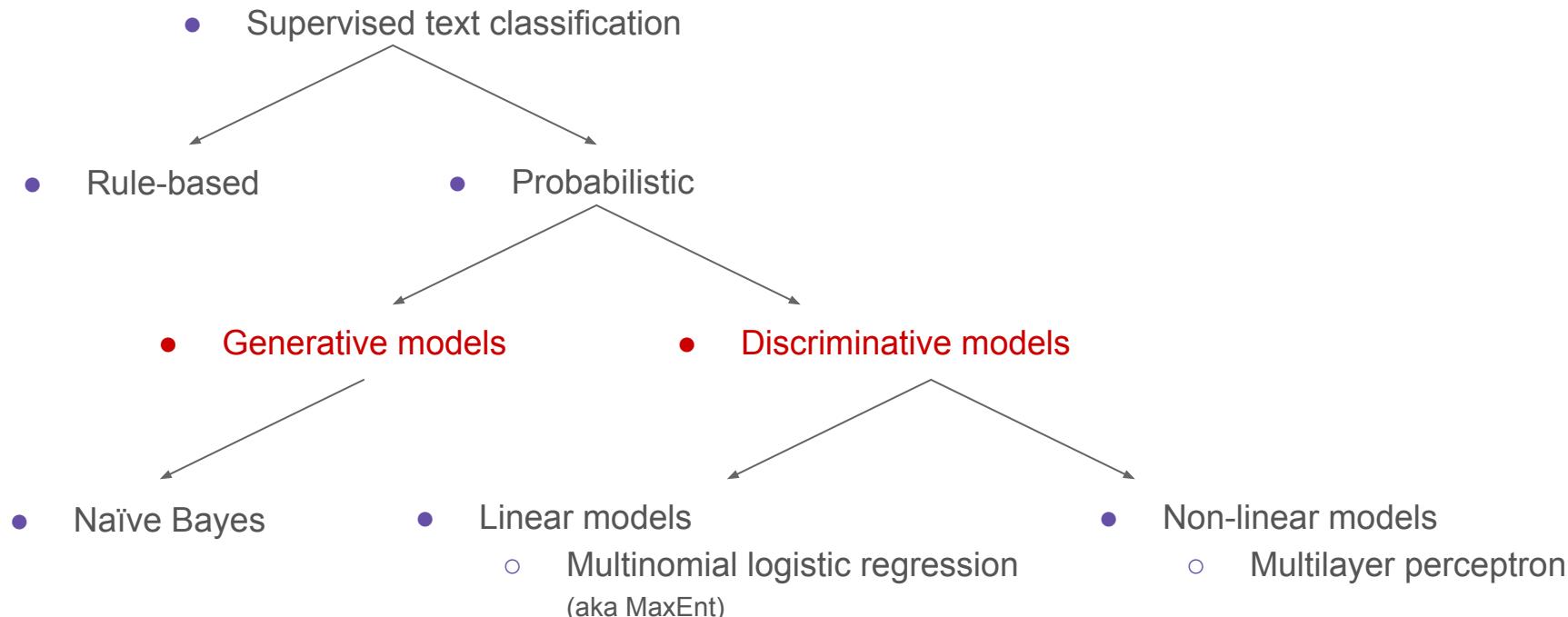
# Supervised classification: formal setting

- Learn a **classification model** from labeled data on
  - properties (“features”) and their importance (“weights”)
- **X**: set of attributes or features  $\{x_1, x_2, \dots, x_n\}$ 
  - e.g. fruit measurements, or word counts extracted from an input documents
- **y**: a “class” label from the label set  $Y = \{y_1, y_2, \dots, y_k\}$ 
  - e.g., fruit type, or spam/not spam, positive/negative/neutral
- Given data samples  $\{x_1, x_2, \dots, x_n\}$  and corresponding labels  $Y = \{y_1, y_2, \dots, y_k\}$
- We **train** a function  $f: x \in X \rightarrow y \in Y$  (the model)

# Supervised classification: formal setting

- Learn a **classification model** from labeled data on
  - properties (“features”) and their importance (“weights”)
- **X**: set of attributes or features  $\{x_1, x_2, \dots, x_n\}$ 
  - e.g. fruit measurements, or word counts extracted from an input documents
- **y**: a “class” label from the label set  $Y = \{y_1, y_2, \dots, y_k\}$ 
  - e.g., fruit type, or spam/not spam, positive/negative/neutral
- At inference time, apply the model on new instances to **predict the label**

# We'll consider alternative models for classification



# Generative and discriminative models

- Generative model: a model that calculates the probability of the input data itself

$P(X, Y)$

joint

- Discriminative model: a model that calculates the probability of a latent trait given the data

$P(Y | X)$

conditional

# Generative and discriminative models

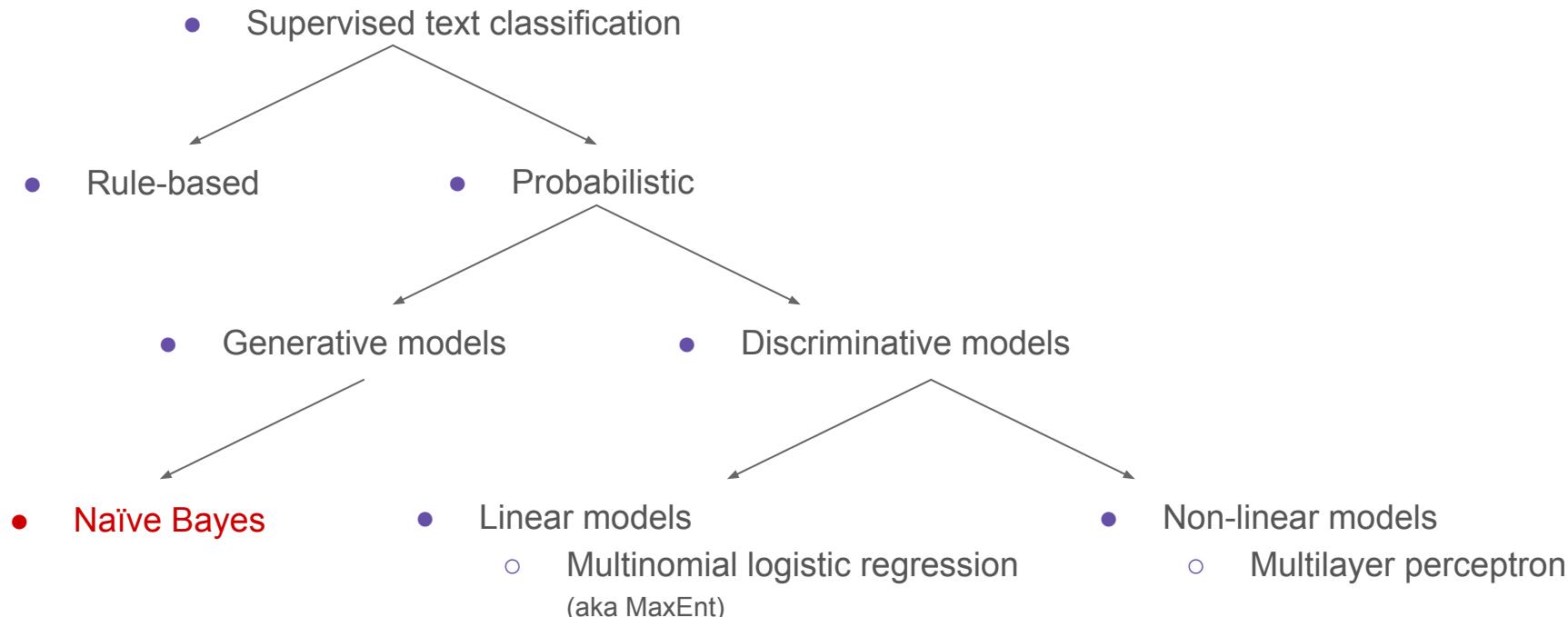
- Generative text classification: Learn a model of the joint  $P(X, y)$ , and find

$$\hat{y} = \operatorname{argmax}_{\tilde{y}} P(X, \tilde{y})$$

- Discriminative text classification: Learn a model of the conditional  $P(y | X)$ , and find

$$\hat{y} = \operatorname{argmax}_{\tilde{y}} P(\tilde{y}|X)$$

# We'll consider alternative models for classification



# Generative text classification: naïve Bayes

- Simple (naïve) classification method
  - based on the Bayes rule
- Relies on very simple representation of documents
  - bag-of-words, no relative order
- A good baseline for more sophisticated models

Andrew Y. Ng and Michael I. Jordan, On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes, Advances in Neural Information Processing Systems 14 (NIPS), 2001.

# Naïve Bayes

## Sentiment analysis: movie reviews

- Given a document  $d$  (e.g., a movie review)
- Decide which class  $c$  it belongs to: positive, negative, neutral
- Compute  $P(c | d)$  for each  $c$ 
  - $P(\text{positive} | d), P(\text{negative} | d), P(\text{neutral} | d)$
  - select the one with max  $P$

# Bag-of-Words (BOW) (I told you it'd be back soon!)

- Given a document  $d$  (e.g., a movie review) – how to represent  $d$  ?



**Figure 7.1** Intuition of the multinomial naive Bayes classifier applied to a movie review. The position of the words is ignored (the *bag of words* assumption) and we make use of the frequency of each word.

Figure from J&M 3rd ed. draft, sec 7.1

# Naïve Bayes

- Given a document  $d$  and a class  $c$ , Bayes' rule:

$$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

# Naïve Bayes

- Given a document  $d$  and a class  $c$ , Bayes' rule:

$$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

$$P(\text{'positive'}|d) \propto P(d|\text{'positive'})P(\text{'positive'})$$

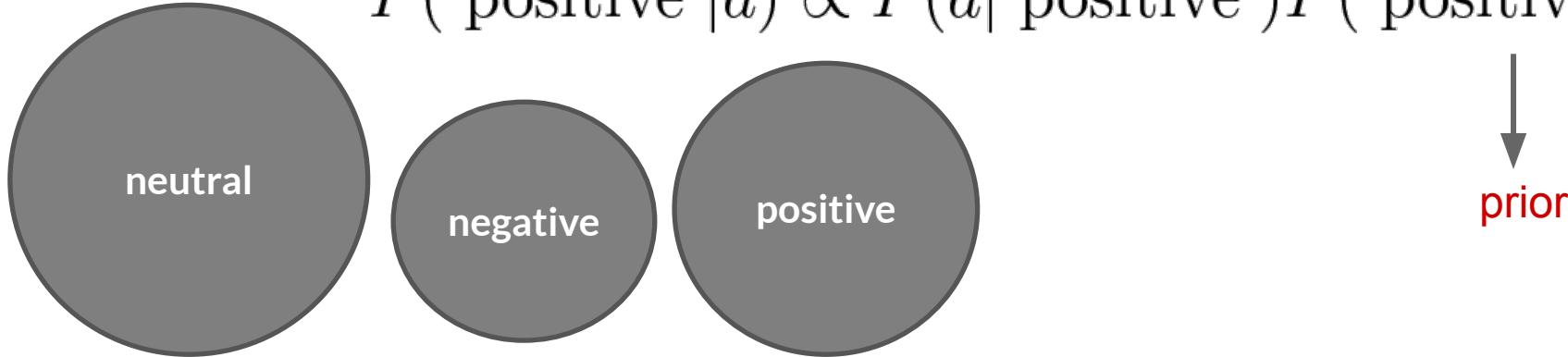
↓                                    ↓  
likelihood                              prior

# Naïve Bayes

- Given a document  $d$  and a class  $c$ , Bayes' rule:

$$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

$$P(\text{'positive'}|d) \propto P(d|\text{'positive'})P(\text{'positive'})$$



# Naïve Bayes

- Given a document  $d$  and a class  $c$ , Bayes' rule:

$$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

$$P(\text{'positive'}|d) \propto P(d|\text{'positive'})P(\text{'positive'})$$



likelihood

# Naïve Bayes independence assumptions

$$P(w_1, w_2, \dots, w_n | c)$$

- **Bag of Words assumption:** Assume position doesn't matter
- **Conditional Independence:** Assume the feature probabilities  $P(w_i | c_j)$  are independent given the class  $c$

$$P(w_1, w_2, \dots, w_n | c) = P(w_1 | c) \times P(w_2 | c) \times P(w_3 | c) \times \dots \times P(w_n | c)$$

# Document representation

I love this movie. It's sweet but with satirical humor. The dialogue is great and the adventure scenes are fun... it manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



**bag of words  
(BOW)**



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

# Document representation

I love this movie. It's sweet but with satirical humor. The dialogue is great and the adventure scenes are fun... it manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



bag of words  
(BOW)

it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

$$P(d|c) = P(w_1, w_2, \dots, w_n | c) = \prod_i P(w_i | c)$$

# Generative text classification: Naïve Bayes

$$\text{C}_{NB} = \operatorname{argmax}_c P(c|d) = \operatorname{argmax}_c \frac{P(d|c)P(c)}{P(d)} \propto \text{Bayes rule}$$

$$\operatorname{argmax}_c P(d|c)P(c) = \text{same denominator}$$

$$\operatorname{argmax}_c P(w_1, w_2, \dots, w_n|c)P(c) = \text{representation}$$

$$\operatorname{argmax}_{c_j} P(c_j) \prod_i P(w_i|c) \text{conditional independence}$$

# Underflow prevention: log space

- Multiplying lots of probabilities can result in floating-point underflow
- Since  $\log(xy) = \log(x) + \log(y)$ 
  - better to sum logs of probabilities instead of multiplying probabilities
- Class with highest un-normalized log probability score is still most probable

$$C_{NB} = \operatorname{argmax}_{c_j} P(c_j) \prod_i P(w_i|c)$$

$$C_{NB} = \operatorname{argmax}_{c_j} \log(P(c_j)) + \sum_i \log(P(w_i|c))$$

- Model is now just max of sum of weights

# Learning the multinomial naïve Bayes

- How do we learn (train) the NB model?

# Learning the multinomial naïve Bayes

- How do we learn (train) the NB model?
- We learn  $P(c)$  and  $P(w_i|c)$  from training (labeled) data

$$C_{NB} = \operatorname{argmax}_{c_j} \log(\underline{P(c_j)}) + \sum_i \log(\underline{P(w_i|c)})$$

# Parameter estimation

- Parameter estimation during training
- Concatenate all documents with category  $c$  into one mega-document
- Use the frequency of  $w_i$  in the mega-document to estimate the word probability

$$\text{C}_{NB} = \operatorname{argmax}_{c_j} \log(P(c_j)) + \sum_i \log(P(w_i|c))$$

$$\hat{P}(c_j) = \frac{\text{doccount}(C = c_j)}{N_{\text{doc}}}$$

$$\hat{P}(w_i|c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

# Parameter estimation

$$\hat{P}(w_i|c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

- fraction of times word  $w_i$  appears among all words in documents of topic  $c_j$
- Create mega-document for topic  $j$  by concatenating all docs in this topic
  - Use frequency of  $w$  in mega-document

# Problem with Maximum Likelihood

- What if we have seen no training documents with the word “fantastic” and classified in the topic **positive**?

# Problem with Maximum Likelihood

- What if we have seen no training documents with the word “fantastic” and classified in the topic **positive**?

$$\hat{P}(\text{“fantastic”} | c = \text{positive}) = \frac{\text{count(“fantastic”, positive)}}{\sum_{w \in V} \text{count}(w, \text{positive})} = 0$$

- Zero probabilities cannot be conditioned away, no matter the other evidence!

$$\operatorname{argmax}_{c_j} P(c_j) \prod_i P(w_i | c)$$

# Laplace (add-1) smoothing for naïve Bayes

$$\hat{P}(w_i|c_j) = \frac{\text{count}(w_i, c_j) + 1}{\sum_{w \in V} (\text{count}(w, c_j) + 1)}$$

# Laplace (add-1) smoothing for naïve Bayes

$$\hat{P}(w_i|c_j) = \frac{\text{count}(w_i, c_j) + 1}{\sum_{w \in V} (\text{count}(w, c_j) + 1)}$$

$$= \frac{\text{count}(w_i, c_j) + 1}{(\sum_{w \in V} (\text{count}(w, c_j))) + |V|}$$

# Multinomial naïve Bayes : learning

- From training corpus, extract *Vocabulary*
- Calculate  $P(c_j)$  terms
  - For each  $c_j$  do
    - $docs_j \leftarrow$  all docs with class =  $c_j$
    - $P(c_j) \leftarrow \frac{|docs_j|}{total \# documents}$

# Multinomial naïve Bayes : learning

- From training corpus, extract *Vocabulary*
- Calculate  $P(c_j)$  terms
  - For each  $c_j$  do
    - $docs_j \leftarrow$  all docs with class =  $c_j$
    - $P(c_j) \leftarrow \frac{|docs_j|}{total \# documents}$
- Calculate  $P(w_i | c_j)$  terms
  - $Text_j \leftarrow$  single doc containing all docs<sub>j</sub>
  - For each word  $w_i$  in *Vocabulary*
    - $n_i \leftarrow$  # of occurrences of  $w_i$  in  $Text_j$
    - $P(w_j | c_j) \leftarrow \frac{n_i + \alpha}{n + \alpha |Vocabulary|}$

# Example

	Doc	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shanghai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

# Example

$$\hat{P}(c) = \frac{N_c}{N}$$

Priors:

$$P(c) = \frac{3}{4}$$

$$P(j) = \frac{1}{4}$$

	Doc	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shanghai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

# Example

$$\hat{P}(c) = \frac{N_c}{N}$$

$$\hat{P}(w|c) = \frac{\text{count}(w,c)+1}{\text{count}(c)+|V|}$$

Priors:

$$P(c) = \frac{3}{4}$$

$$P(j) = \frac{1}{4}$$

Conditional Probabilities:

$$P(\text{Chinese}|c) = (5+1) / (8+6) = 6/14 = 3/7$$

$$P(\text{Tokyo}|c) = (0+1) / (8+6) = 1/14$$

$$P(\text{Japan}|c) = (0+1) / (8+6) = 1/14$$

$$P(\text{Chinese}|j) = (1+1) / (3+6) = 2/9$$

$$P(\text{Tokyo}|j) = (1+1) / (3+6) = 2/9$$

$$P(\text{Japan}|j) = (1+1) / (3+6) = 2/9$$

	Doc	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shanghai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

# Example

$$\hat{P}(c) = \frac{N_c}{N} \quad \hat{P}(w|c) = \frac{\text{count}(w,c)+1}{\text{count}(c)+|V|}$$

Priors:

$$P(c) = \frac{3}{4} \quad P(j) = \frac{1}{4}$$

Conditional Probabilities:

$$P(\text{Chinese}|c) = (5+1) / (8+6) = 6/14 = 3/7$$

$$P(\text{Tokyo}|c) = (0+1) / (8+6) = 1/14$$

$$P(\text{Japan}|c) = (0+1) / (8+6) = 1/14$$

$$P(\text{Chinese}|j) = (1+1) / (3+6) = 2/9$$

$$P(\text{Tokyo}|j) = (1+1) / (3+6) = 2/9$$

$$P(\text{Japan}|j) = (1+1) / (3+6) = 2/9$$

	Doc	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shanghai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

Choosing a class:

$$P(c|d5) \propto 3/4 * (3/7)^3 * 1/14 * 1/14 \\ \approx 0.0003$$

$$P(j|d5) \propto 1/4 * (2/9)^3 * 2/9 * 2/9 \\ \approx 0.0001$$

# Summary: naïve Bayes is not so naïve

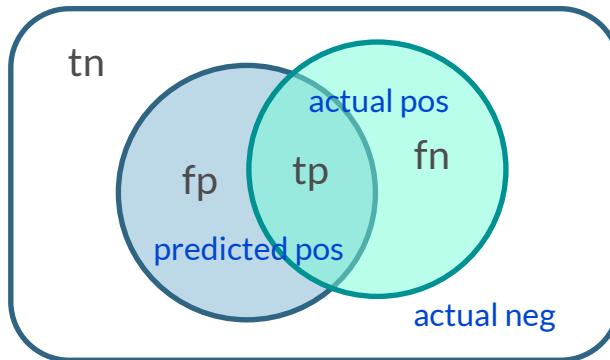
- Naïve Bayes is a probabilistic model
- Naïve because it assumes features are independent of each other for a class
- Very fast, low storage requirements
- Robust to Irrelevant Features
  - Irrelevant Features cancel each other without affecting results
- Very good in domains with many equally important features
  - Decision Trees suffer from fragmentation in such cases – especially if little data
- Optimal if the independence assumptions hold: If assumed independence is correct, then it is the Bayes Optimal Classifier for problem
- A good dependable baseline for text classification
  - But we will see other classifiers that give better accuracy

# How do we evaluate our function $f$ ?

# Classification evaluation

- Contingency table: model's predictions are compared to the correct results
  - a.k.a. confusion matrix

	actual pos	actual neg
predicted pos	true positive (tp)	false positive (fp)
predicted neg	false negative (fn)	true negative (tn)



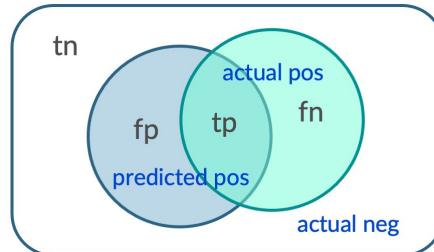
# Classification evaluation

- Borrowing from Information Retrieval, empirical NLP systems are usually evaluated using the notions of **precision** and **recall**

# Classification evaluation

- Precision (P) is the proportion of the selected items that the system got right in the case of text categorization
  - it is the % of documents classified as “positive” by the system which are indeed “positive” documents
- Reported per class or average

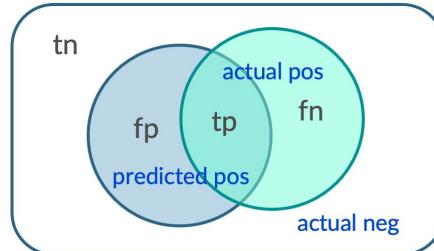
$$\text{precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}} = \frac{tp}{tp + fp}$$



# Classification evaluation

- Recall (R) is the proportion of actual items that the system selected in the case of text categorization
  - it is the % of the “positive” documents which were actually classified as “positive” by the system
- Reported per class or average

$$\text{recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}} = \frac{tp}{tp + fn}$$



# Classification evaluation

- We often want to trade-off precision and recall
  - typically: the higher the precision the lower the recall
  - can be plotted in a precision-recall curve
- It is convenient to combine P and R into a single measure
  - one possible way to do that is F measure

$$F_{\beta} = \frac{(\beta^2 + 1)PR}{\beta^2P + R} \quad \text{for } \beta=1, F_1 = \frac{2PR}{P+R}$$

# Classification evaluation

- Additional measures of performance: accuracy and error
  - accuracy is the proportion of items the system got right
  - error is its complement

$$\text{accuracy} = \frac{tp + tn}{tp + fp + tn + fn}$$

# Micro- vs. macro-averaging

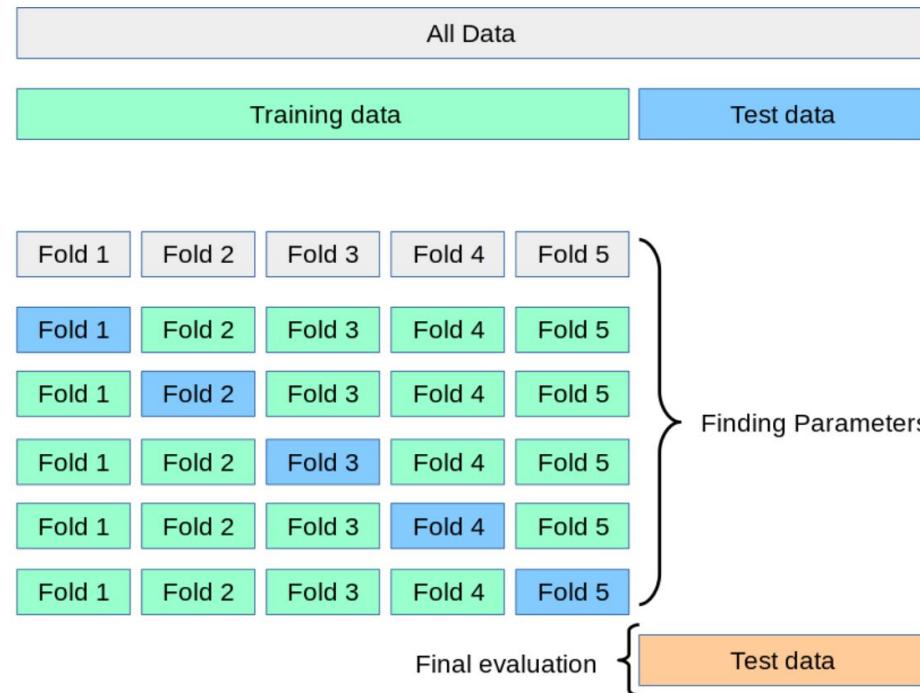
If we have more than one class, how do we combine multiple performance measures into one quantity?

- Macroaveraging
  - Compute performance for each class, then average.
- Microaveraging
  - Collect decisions for all classes, compute contingency table, evaluate.

# Classification common practices

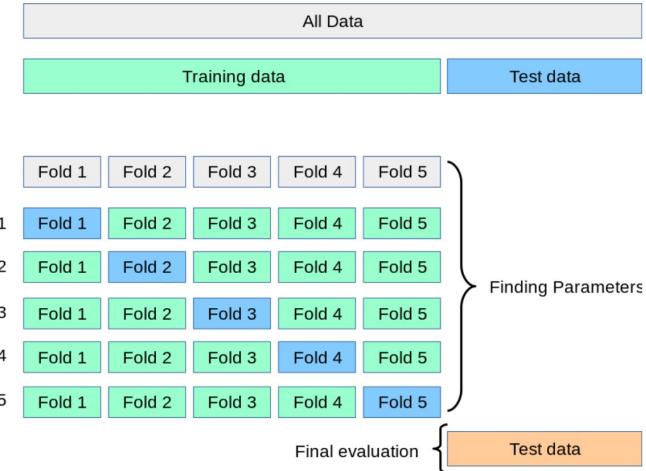
- Divide the training data into  $k$  folds (e.g.,  $k=10$ )
- Repeat  $k$  times: train on  $k-1$  folds and test on the holdout fold, cyclically
- Average over the  $k$  folds' results

# K-fold cross-validation

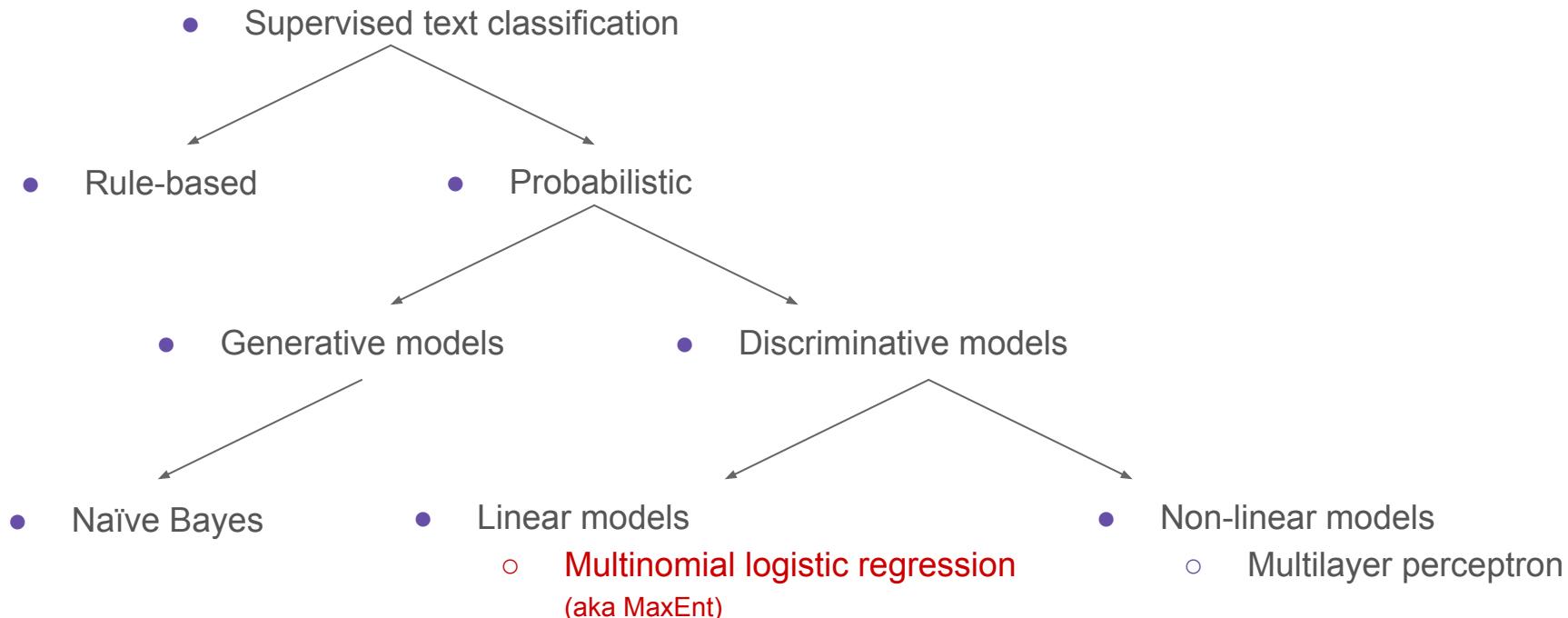


# K-fold cross-validation

- Metric: P/R/F1 or Accuracy
- Unseen test set
  - avoid overfitting ('tuning to the test set')
  - more conservative estimate of performance
- Cross-validation over multiple splits
  - Handles sampling errors from different datasets
  - Pool results over each split
  - Compute pooled dev set performance



# Next class



# Readings

- Eis 2
- J&M III 4
- Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. In Proceedings of EMNLP, 2002
- Andrew Y. Ng and Michael I. Jordan, On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes, In Proceedings of NeurIPS, 2001.