

Natural Language Processing

Text classification

Yulia Tsvetkov

yuliats@cs.washington.edu

Some questions from the previous lecture

- NLP for sign languages

Announcements

- HW1 is out today, please start early!

Is this spam?

from: **ECRES 2022 <2022@ecres.net>** [via](#) amazonses.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: amazonses.com
security: Standard encryption (TLS) [Learn more](#)

Dear Colleague,

Account: 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](#)

There will be keynote speakers who will address specific topics of energy as you would see at [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](#) . **Please note that the official journal of the event, Journal of Energy Systems ([dergipark.org.tr/jes](#)) is also indexed in SCOPUS.**

Spam classification

Dear Colleague,

Account: 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 the 2022 Conference on Renewable Energy Systems (ECRES). **ECRES 2022** will be held in Istanbul, Turkey. **ECRES 2022 is a hybrid conference, the participants can present their papers physically or virtually.** The conference is being organized in Istanbul/Turkey under the technical sponsorship of 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 special and regular issue journals can be seen in ecres.net

There will be keynote speakers who will address specific topics of energy as you would see at 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. **Please note that the official journal of the event, [Journal of Energy Systems \(dergipark.org.tr/jes\)](http://Journal of Energy Systems (dergipark.org.tr/jes)) is also indexed in SCOPUS.**

Invitation to present at the February 2022 Wikimedia Research Showcase



Emily Lescak <elescak@wikimedia.org>
to 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 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 community researchers present 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, Wikipedia editors, and Wikimedia staff—would find interesting and relevant to their work.

Research Showcase presentations are generally 20 minutes long, with an additional 10 minutes for questions. 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. You can watch past showcases on our [YouTube](#) 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, please let us know so we can discuss other options.

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

Sincerely,

Emily

--



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, говийн ээрэм хөндийд, малын бэлчээрт, малчдын хотонд болохоор солигдож эргэцүүлэн бодох нь хачин. Цас хэр орсон бол?

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

Beograd, 16. јun 2013. године – Председник Владе Републике Србије Ivica Dačić čestitao је кајакашиси златне медаље у олимпијској дисциплини К-1, 500 метара, као и у двоstrуко дуžој стази освојене на првенству Европе у Portugaliji.

Nestranski Urad за vladno odgovornost ZDA је objavil eksplozivno mnenje, da је vlada predsednika Donalda Trumpa kršila zvezno zakonodajo, ko је zadrževala izplačilo kongresno potrjene vojaške pomoči Ukrajini zaradi političnih razlogov. Predstavniški dom kongresa је 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 бодох нь хачин. Цас хэр орсон бол?

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

Beograd, 16. јun 2013. године – Председник Владе Републике Ср **serbian** неститао је кајакашици златне медаље у о **serbian** K-1, 500 метара, као и у двојструко дуђој стази освојене на првенству Европе у Португалији.

Nestranski Urad za vladno odgovornost ZDA je objavil eksplozivno mnenje, da je vlada predsednika Donalda Trumpa kršila zvezno zakonodajo, ko je zadrževala izplačilo k **slovenian** vojaške pomoči Ukrajini zaradi političnih razlogov. Predstavniški d **slovenian** av zaradi tega 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. 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, 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 mask and I came all over sugar and 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

Syntactic frame and verb bias in aphasia: Plausibility judgments of underdog-subject sentences
 Susanna Galli,¹ Lisa Mann,² Carl Ramscar,³ David R. Just,⁴ Elizabeth Bates,⁵ Moly Ravegg,⁶ and L. Holland Aulaby,⁷

¹University of California, Berkeley, Berkeley, CA, USA
²University of California, Berkeley, Berkeley, CA, USA
³University of California, Berkeley, Berkeley, CA, USA
⁴University of California, Berkeley, Berkeley, CA, USA
⁵University of California, Berkeley, Berkeley, CA, USA
⁶University of California, Berkeley, Berkeley, CA, USA
⁷University of California, Berkeley, Berkeley, CA, USA

Abstract
 The study investigates how factors that have been argued to define "lexical bias" in sentence comprehension interact with syntactic frame and degree of ambiguity to influence plausibility judgments in aphasia. The study examines the effects of syntactic frame and degree of ambiguity on plausibility judgments in aphasia. The study examines the effects of syntactic frame and degree of ambiguity on plausibility judgments in aphasia. The study examines the effects of syntactic frame and degree of ambiguity on plausibility judgments in aphasia.

1. Introduction
 The concept of "lexical bias" or "lexical word bias" in normal and aphasic comprehension has often been taken as evidence in the normal comprehension literature that the lexicon is not just a list of words (Marslen-Wilson & Levy, 2008). The original source of the term "lexical bias" was the observation that the comprehension of "underdog-subject" sentences is faster than the comprehension of "dog-subject" sentences (Just & Carpenter, 1982). The original source of the term "lexical bias" was the observation that the comprehension of "underdog-subject" sentences is faster than the comprehension of "dog-subject" sentences (Just & Carpenter, 1982).



MeSH Subject Category Hierarchy

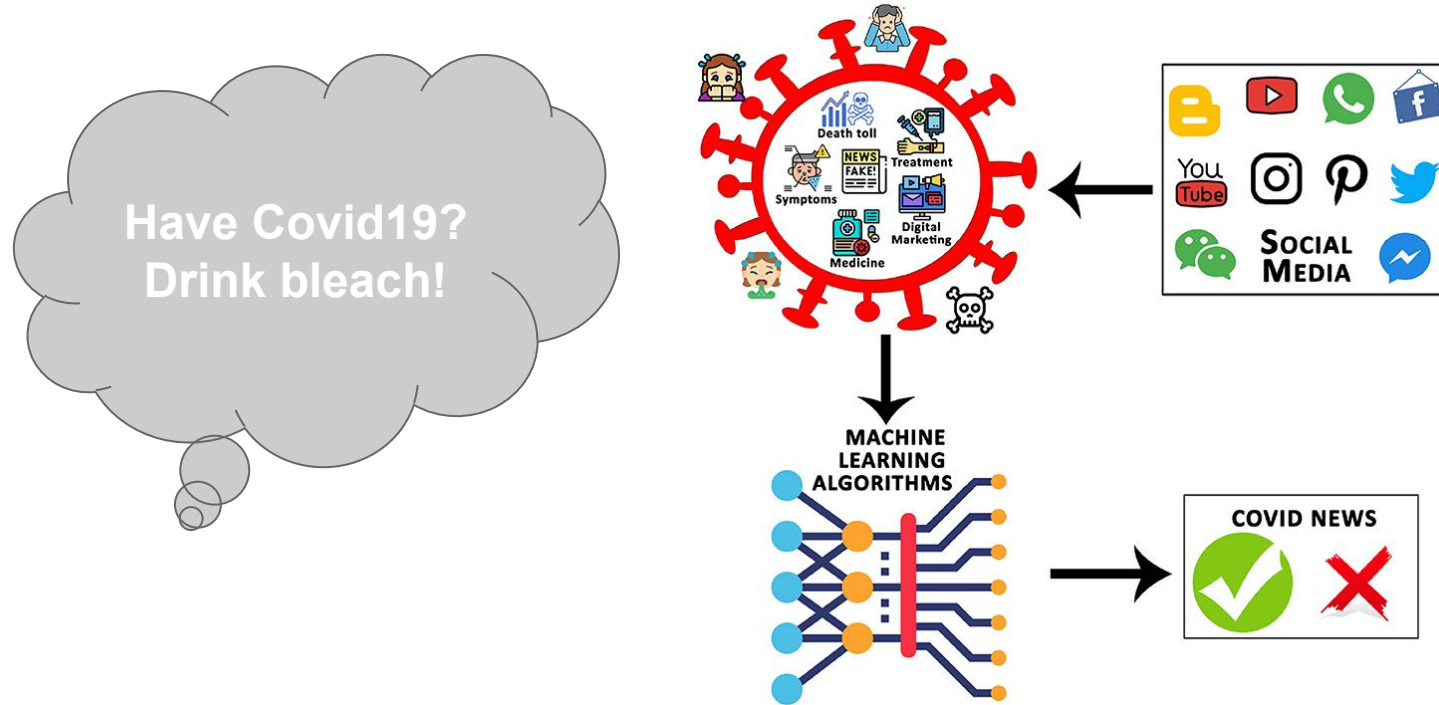
- Antagonists and Inhibitors
- Blood Supply
- Chemistry
- Drug Therapy
- Embryology
- Epidemiology
- ...

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 Cochinchina; 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...

Fact verification: trustworthy or fake?



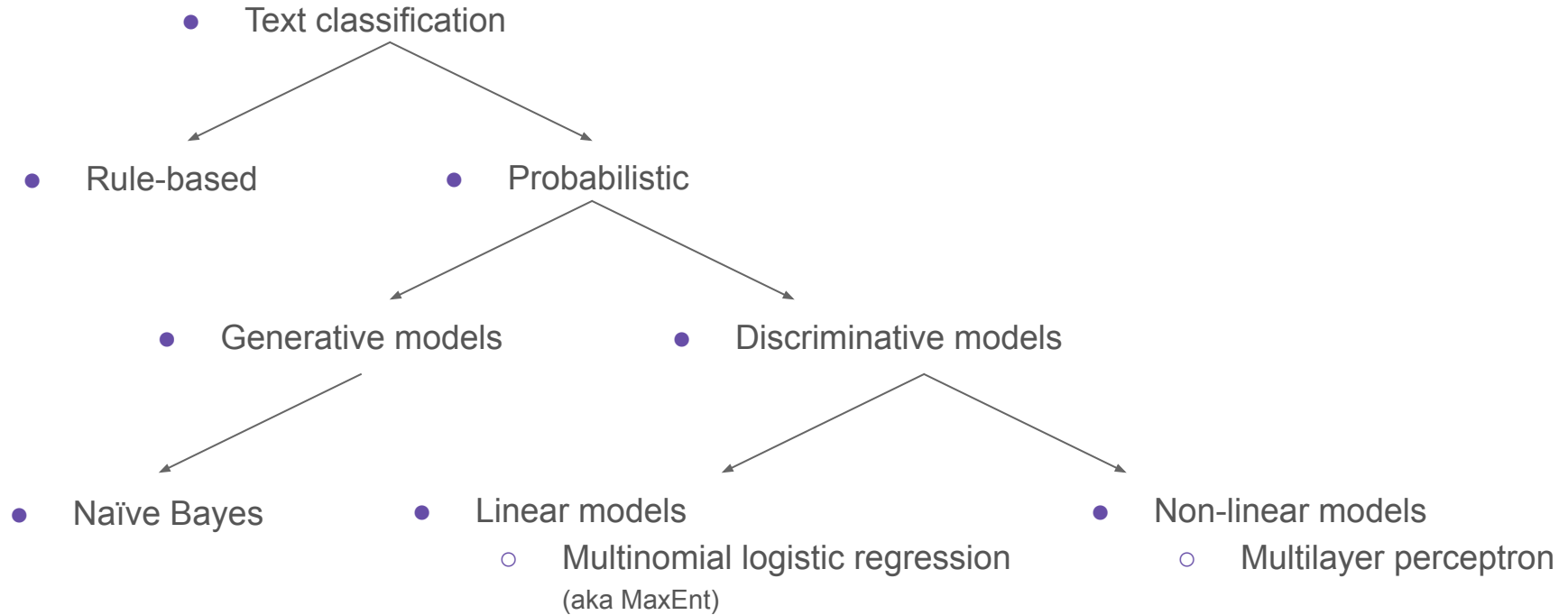
Detecting COVID-19-Related Fake News Using Feature Extraction

Suleman Khan, Saqib Hakak, N. Deepa, B. Prabadevi, Kapal Dev and Silvia Trelova

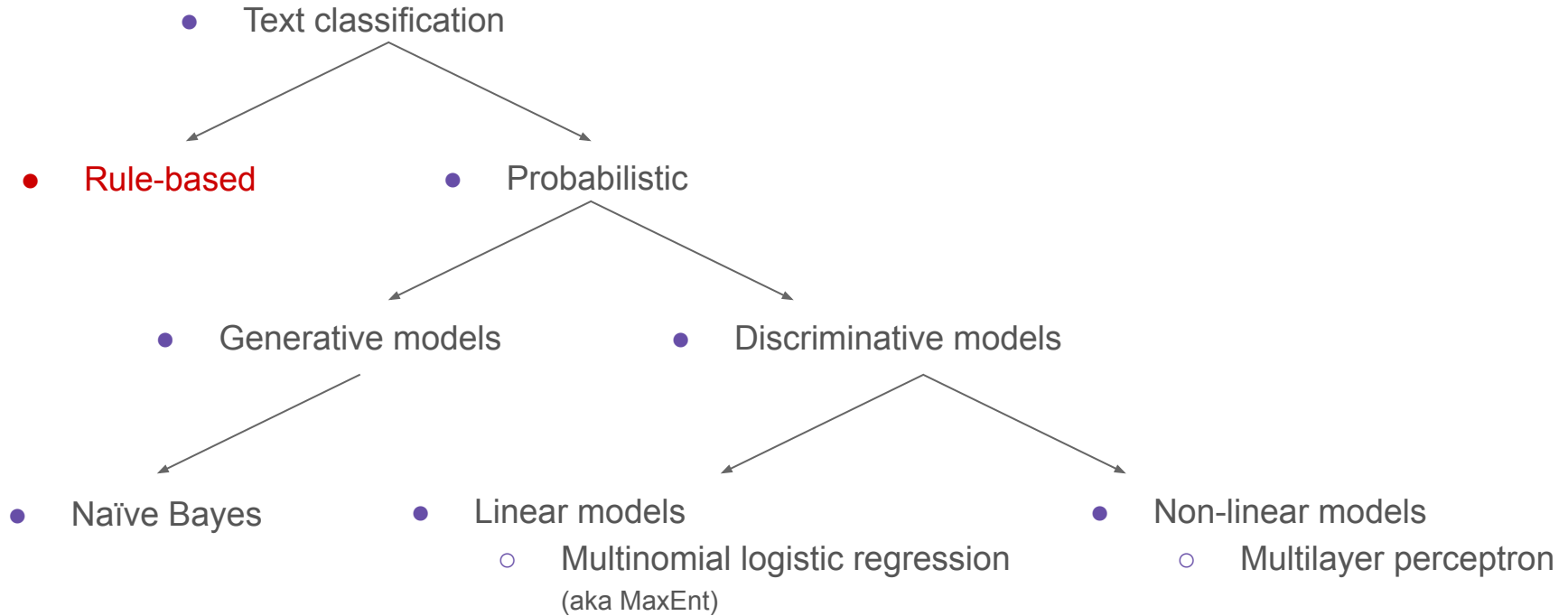
Text classification

- 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)
 - ...

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.

→ word order 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.

→ word order matters, but hard to encode in rules!

Language ID: All falter, stricken in kind.

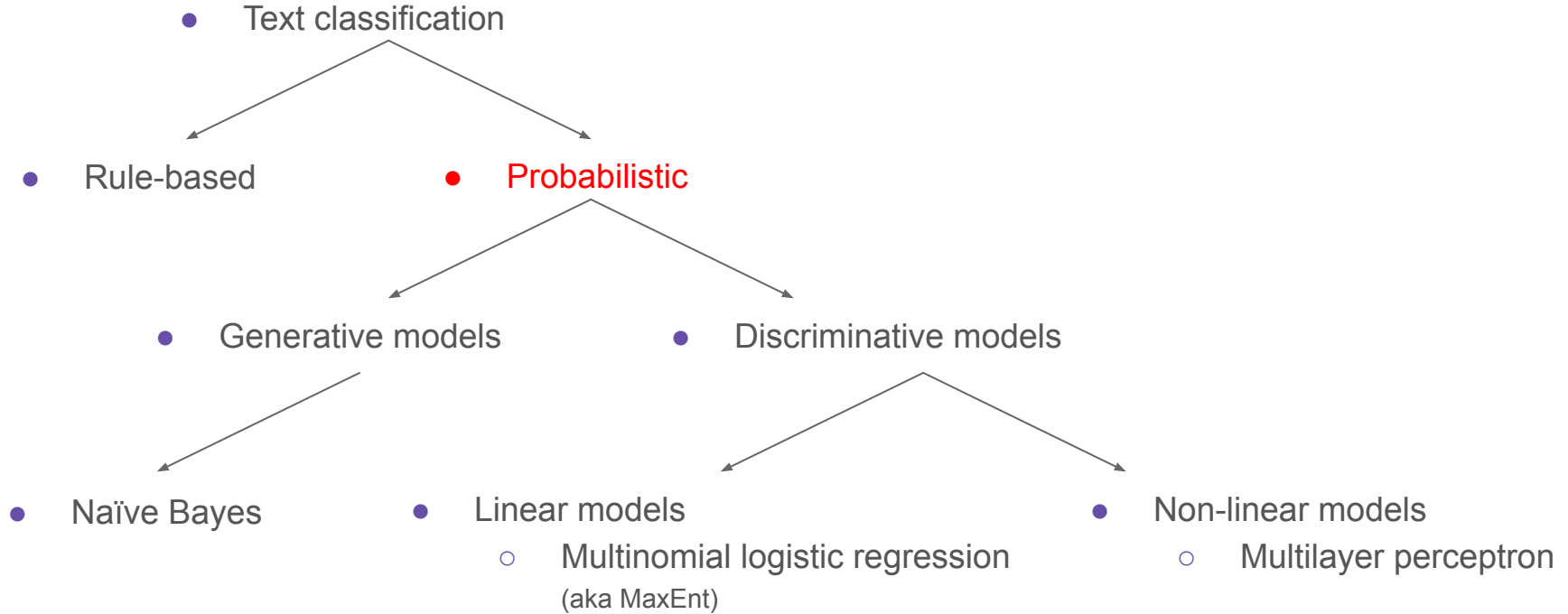
“LINGERIE SALE”

→ simple features can be misleading!

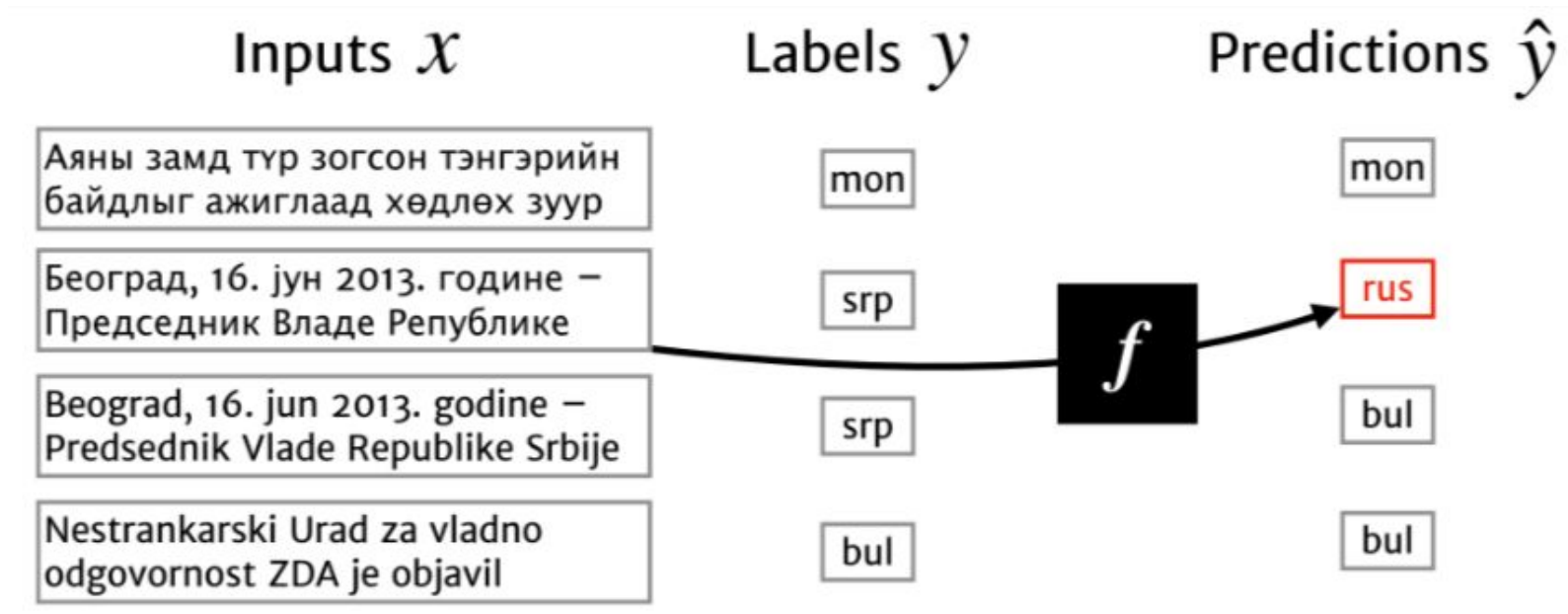
Rule-based classification

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



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

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

Government
proceedings

Product
reviews

Noisy
labels

Domain
names

Link text

Expert
annotation

Treebanks

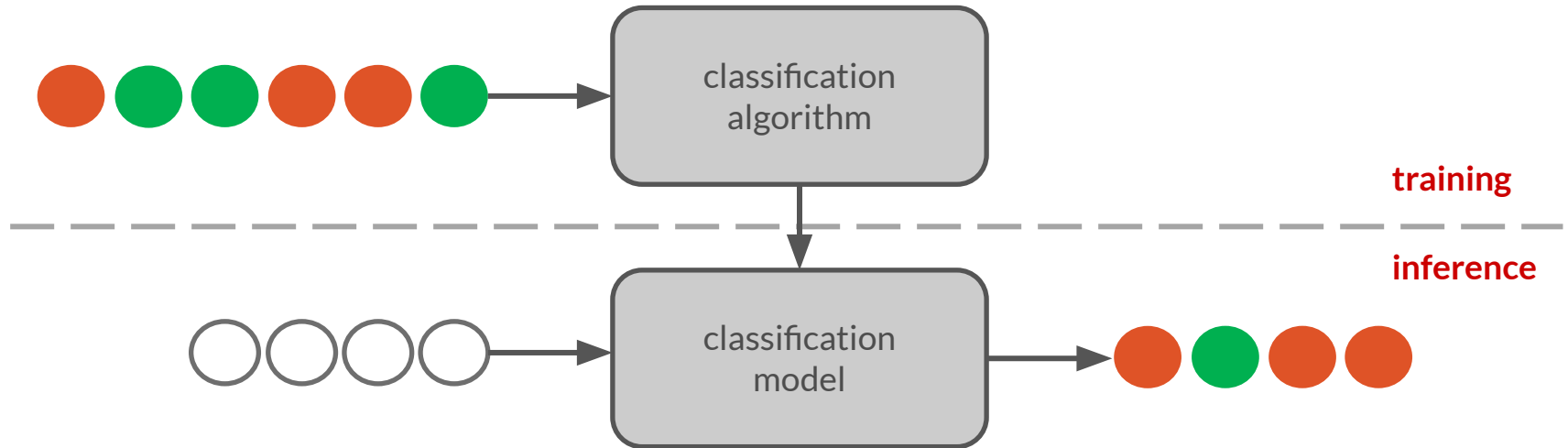
Biomedical
corpora

Crowd
workers

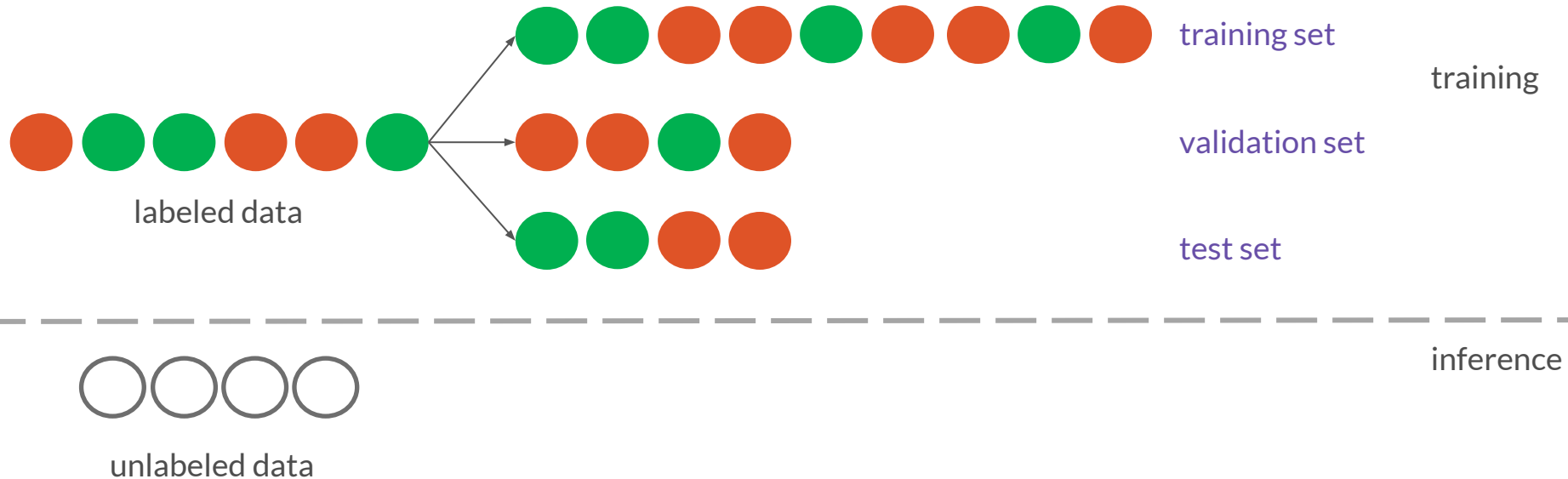
Question
answering

Image
captions

Supervised classification



Training, validation, and test sets



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

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**

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.

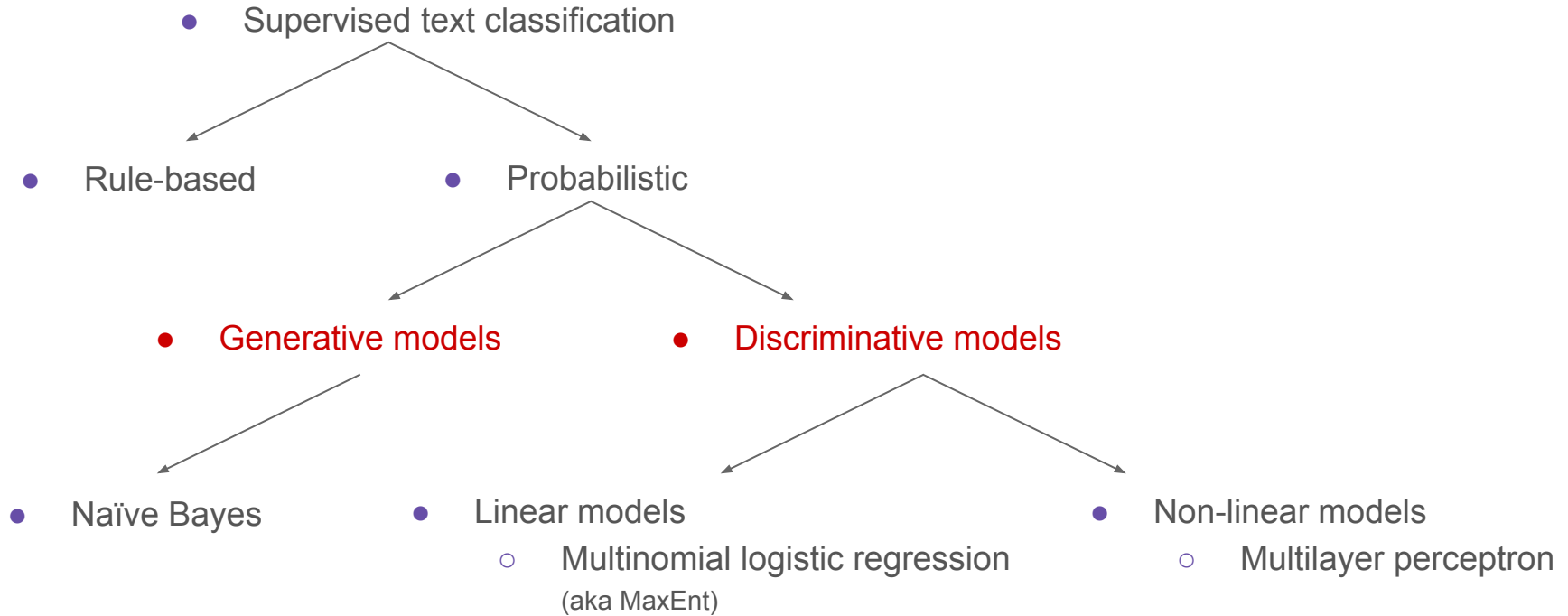
(almost) the entire lexicon

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
- ...

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

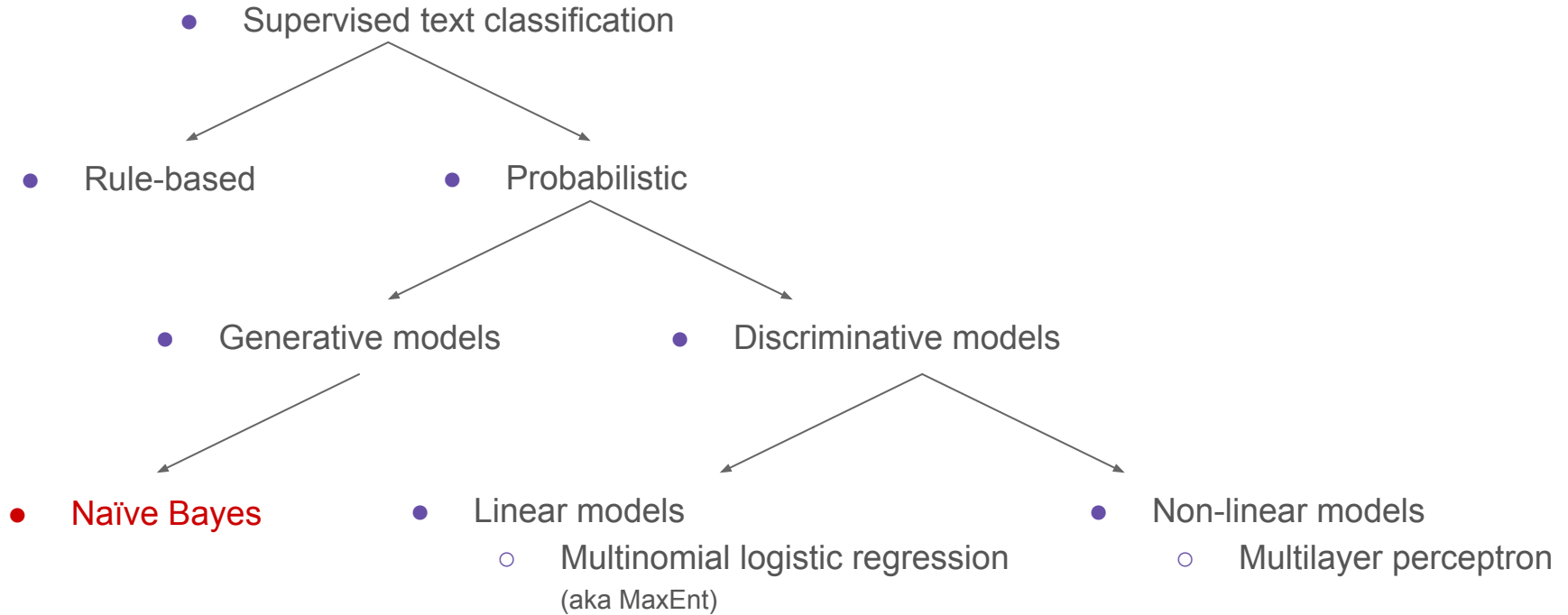
- Generative text classification: Learn a model of the joint $P(\mathbf{X}, \mathbf{y})$, and find

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

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

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

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 a 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)

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

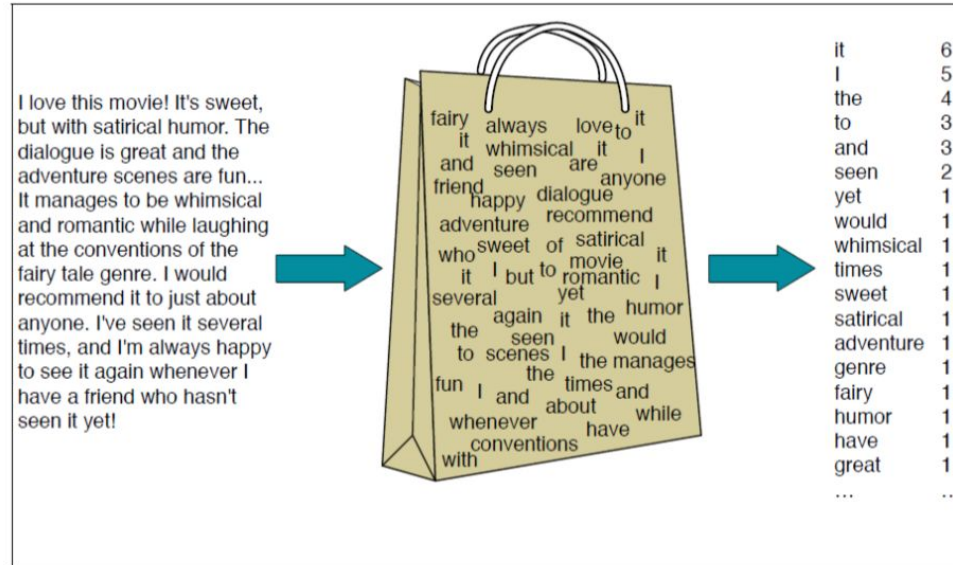


Figure 7.1 Intuition of the multinomial naive Bayes classifier applied to a movie review. The position of the words is ignored (the *bag of words* assumption) and we make use of the frequency of each word.

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

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

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



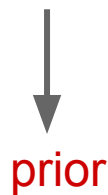
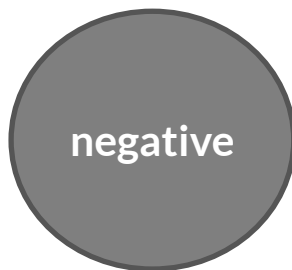
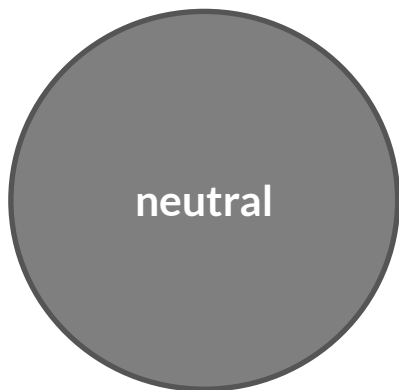
prior

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

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

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

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

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

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

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

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

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

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**?

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

- 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

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

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_j
 - 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

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

	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	?

Priors:

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

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

Conditional Probabilities:

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

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

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

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

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

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

Choosing a class:

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

$$\approx 0.0003$$

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

$$\approx 0.0001$$

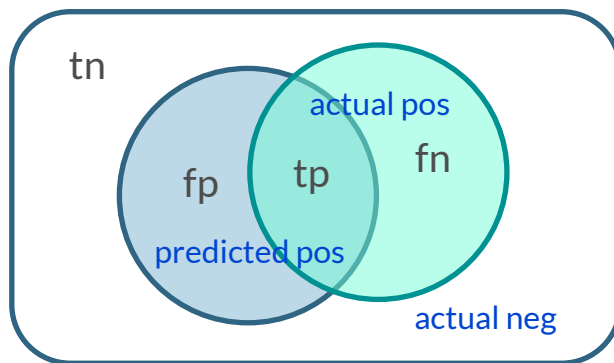
Summary: naïve Bayes is not so naïve

- 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

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)



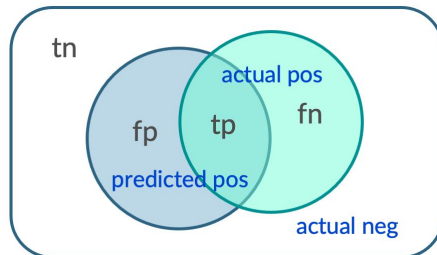
Classification evaluation

- 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

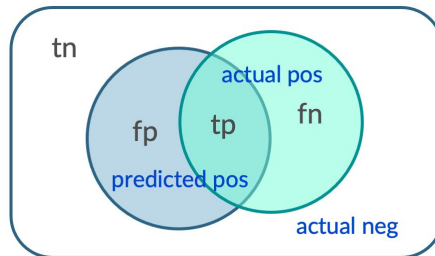
$$\text{precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}} = \frac{tp}{tp + fp}$$



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

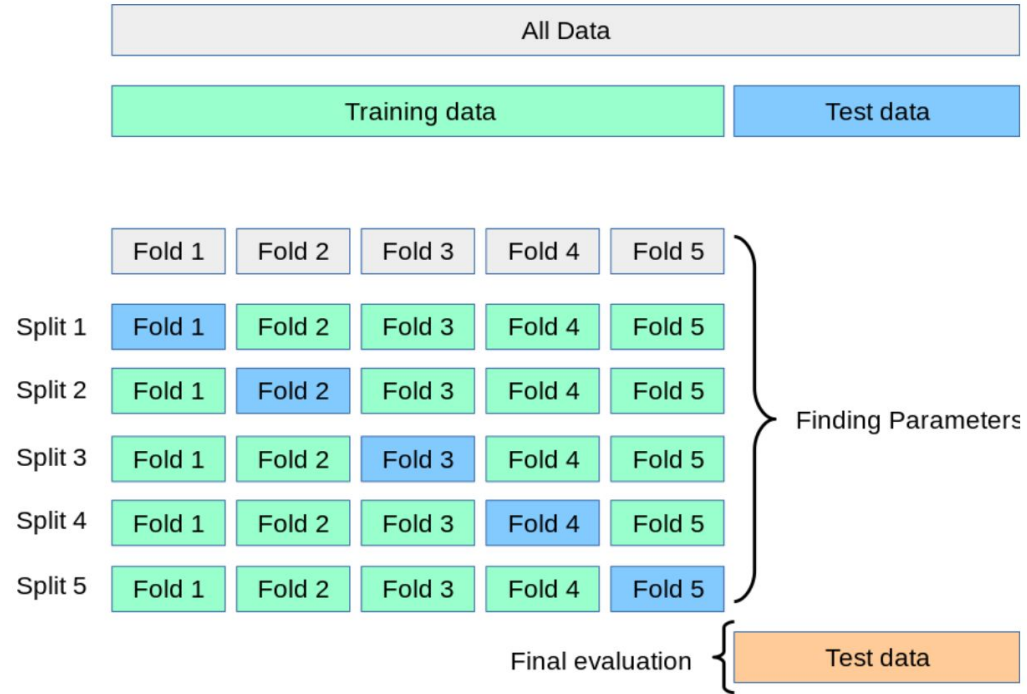
If we have more than one class, how do we combine multiple performance measures into one quantity?

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

Classification common practices

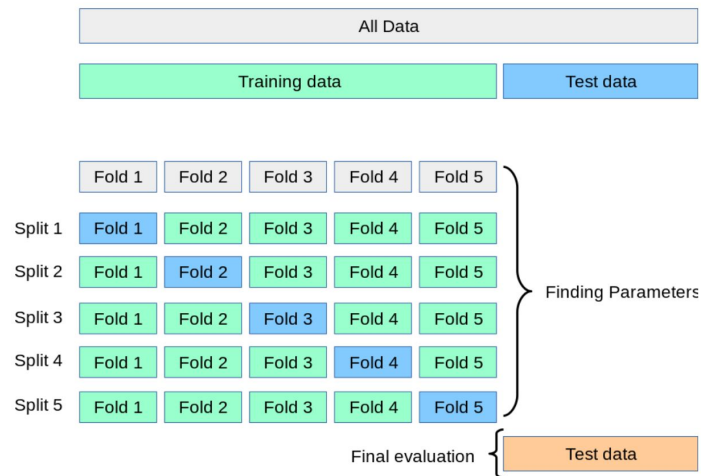
- Divide the training data into k folds (e.g., $k=10$)
- Repeat k times: train on $k-1$ folds and test on the holdout fold, cyclically
- Average over the k folds' results

K-fold cross-validation

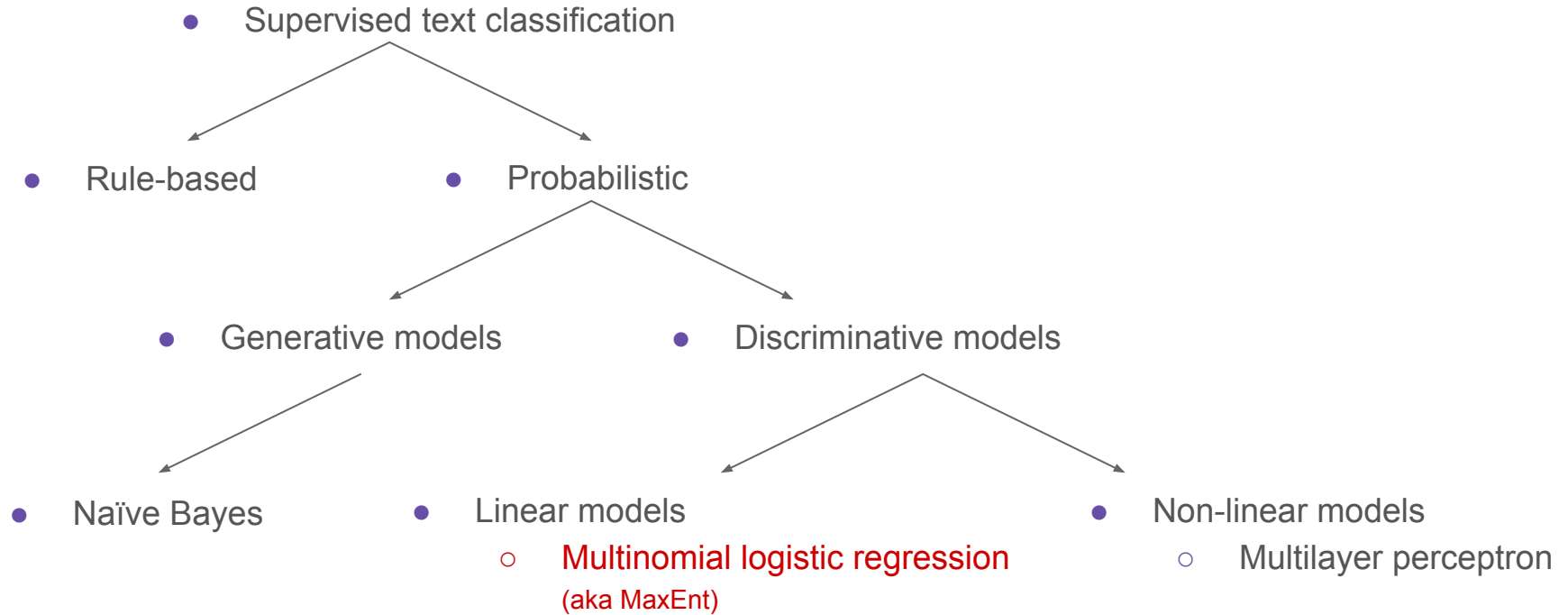


K-fold cross-validation

- 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.