

# Text Classification

He He

New York University

September 15, 2021

# Text classification

- ▶ Input: text (sentence, paragraph, document)
- ▶ Predict the **category or property** of the input text
  - ▶ Sentiment classification: Is the review positive or negative?
  - ▶ Spam detection: Is the email/message spam or not?
  - ▶ Hate speech detection: Is the tweet/post toxic or not?
  - ▶ Stance classification: Is the opinion liberal or conservative?

# Text classification

- ▶ Input: text (sentence, paragraph, document)
- ▶ Predict the **category or property** of the input text
  - ▶ Sentiment classification: Is the review positive or negative?
  - ▶ Spam detection: Is the email/message spam or not?
  - ▶ Hate speech detection: Is the tweet/post toxic or not?
  - ▶ Stance classification: Is the opinion liberal or conservative?
- ▶ Predict the **relation** of two pieces of text
  - ▶ Textual entailment (HW1): does the premise entail the hypothesis?  
Premise: The dogs are running in the park.  
Hypothesis: There are dogs in the park.
  - ▶ Paraphrase detection: are the two sentences paraphrases?  
Sentence 1: The dogs are in the park.  
Sentence 2: There are dogs in the park.

# Table of Contents

Generative models: naive Bayes

Discriminative models: logistic regression

Regularization, model selection, evaluation

# Intuition

**Example:** sentiment classification for movie reviews

*Action. Comedy. Suspense. This movie has it all. The Plot goes that 4 would be professional thieves are invited to take part in a heist in a small town in Montana. every type of crime movie archetype character is here. Frank, the master mind. Carlos, the weapons expert. Max, the explosives expert. Nick, the safe cracker and Ray, the car man. Our 4 characters meet up at the train station and from the beginning none of them like or trust one another. Added to the mix is the fact that Frank is gone and they are not sure why they have called together. Now Frank is being taken back to New Jersey by the 2 detectives but soon escapes on foot and tries to make his way back to the guys who are having all sorts of problems of their own. Truly a great film loaded with laughs and great acting. Just an overall good movie for anyone looking for a laugh or something a little different*

# Intuition

**Example:** sentiment classification for movie reviews

*Action. Comedy. Suspense. This movie has it all. The Plot goes that 4 would be professional thieves are invited to take part in a heist in a small town in Montana. every type of crime movie archetype character is here. Frank, the master mind. Carlos, the weapons expert. Max, the explosives expert. Nick, the safe cracker and Ray, the car man. Our 4 characters meet up at the train station and from the beginning none of them like or trust one another. Added to the mix is the fact that Frank is gone and they are not sure why they have called together. Now Frank is being taken back to New Jersey by the 2 detectives but soon escapes on foot and tries to make his way back to the guys who are having all sorts of problems of their own. Truly a great film loaded with laughs and great acting. Just an overall good movie for anyone looking for a laugh or something a little different*

**Idea:** count the number of positive/negative words

- ▶ What is a “word”?
- ▶ How do we know which are positive/negative?

## Preprocessing: tokenization

**Goal:** Splitting a string of text  $s$  to a sequence of **tokens**  $[x_1, \dots, x_n]$ .

### Language-specific solutions

- ▶ Regular expression: “I didn’t watch the movie”.  $\rightarrow$  [“I”, “did”, “n’t”, “watch”, “the”, “movie”, “.”]
  - ▶ Special cases: U.S., Ph.D. etc.
- ▶ Dictionary / sequence labeler: “我没有去看电影。”  $\rightarrow$  [“我”, “没有”, “去”, “看”, “电影”, “。”]

## Preprocessing: tokenization

**Goal:** Splitting a string of text  $s$  to a sequence of **tokens**  $[x_1, \dots, x_n]$ .

### Language-specific solutions

- ▶ Regular expression: “I didn’t watch the movie”.  $\rightarrow$  [“I”, “did”, “n’t”, “watch”, “the”, “movie”, “.”]
  - ▶ Special cases: U.S., Ph.D. etc.
- ▶ Dictionary / sequence labeler: “我没有去看电影。”  $\rightarrow$  [“我”, “没有”, “去”, “看”, “电影”, “。”]

### General solutions: don’t split by words

- ▶ Characters: [“u”, “n”, “a”, “f”, “f”, “a”, “b”, “l”, “e”] (pros and cons?)



## Preprocessing: tokenization

**Goal:** Splitting a string of text  $s$  to a sequence of **tokens**  $[x_1, \dots, x_n]$ .

### Language-specific solutions

- ▶ Regular expression: “I didn’t watch the movie”.  $\rightarrow$  [“I”, “did”, “n’t”, “watch”, “the”, “movie”, “.”]
  - ▶ Special cases: U.S., Ph.D. etc.
- ▶ Dictionary / sequence labeler: “我没有去看电影。”  $\rightarrow$  [“我”, “没有”, “去”, “看”, “电影”, “。”]

### General solutions: don’t split by words

- ▶ Characters: [“u”, “n”, “a”, “f”, “f”, “a”, “b”, “l”, “e”] (pros and cons?)
- ▶ Subword (e.g., byte pair encoding): [“un”, “aff”, “able#”] [board]

## Classification: problem formulation

- ▶ **Input:** a sequence of tokens  $x = (x_1, \dots, x_n)$  where  $x_i \in \mathcal{V}$ .
- ▶ **Output:** binary label  $y \in \{0, 1\}$ .
- ▶ **Probabilistic model:**

$$f(x) = \begin{cases} 1 & \text{if } p_\theta(y = 1 | x) > 0.5 \\ 0 & \text{otherwise} \end{cases},$$

where  $p_\theta$  is a distribution parametrized by  $\theta \in \Theta$ .

- ▶ **Question:** how to choose  $p_\theta$ ?

## Model $p(y | x)$

How to write a review:

1. Decide the sentiment by flipping a coin:  $p(y)$
2. Generate word sequentially conditioned on the sentiment  $p(x | y)$

## Model $p(y | x)$

How to write a review:

1. Decide the sentiment by flipping a coin:  $p(y)$
2. Generate word sequentially conditioned on the sentiment  $p(x | y)$

$$p(y) = \quad (1)$$

$$p(x | y) = \quad (2)$$

$$= \quad (3)$$

## Model $p(y | x)$

How to write a review:

1. Decide the sentiment by flipping a coin:  $p(y)$
2. Generate word sequentially conditioned on the sentiment  $p(x | y)$

$$p(y) = \text{Bernoulli}(\alpha) \tag{1}$$

$$p(x | y) = \tag{2}$$

$$= \tag{3}$$

## Model $p(y | x)$

How to write a review:

1. Decide the sentiment by flipping a coin:  $p(y)$
2. Generate word sequentially conditioned on the sentiment  $p(x | y)$

$$p(y) = \text{Bernoulli}(\alpha) \quad (1)$$

$$p(x | y) = \prod_{i=1}^n p(x_i | y) \quad (\text{independent assumption}) \quad (2)$$

$$= \quad (3)$$

## Model $p(y | x)$

How to write a review:

1. Decide the sentiment by flipping a coin:  $p(y)$
2. Generate word sequentially conditioned on the sentiment  $p(x | y)$

$$p(y) = \text{Bernoulli}(\alpha) \quad (1)$$

$$p(x | y) = \prod_{i=1}^n p(x_i | y) \quad (\text{independent assumption}) \quad (2)$$

$$= \prod_{i=1}^n \text{Categorical}(\underbrace{\theta_{1,y}, \dots, \theta_{|\mathcal{V}|,y}}_{\text{sum to 1}}) \quad (3)$$

## Model $p(y | x)$

How to write a review:

1. Decide the sentiment by flipping a coin:  $p(y)$
2. Generate word sequentially conditioned on the sentiment  $p(x | y)$

$$p(y) = \text{Bernoulli}(\alpha) \quad (1)$$

$$p(x | y) = \prod_{i=1}^n p(x_i | y) \quad (\text{independent assumption}) \quad (2)$$

$$= \prod_{i=1}^n \text{Categorical}(\underbrace{\theta_{1,y}, \dots, \theta_{|\mathcal{V}|,y}}_{\text{sum to 1}}) \quad (3)$$

## Bayes rule

$$p(y | x) = \frac{p(x | y)p(y)}{p(x)} = \frac{p(x | y)p(y)}{\sum_{y \in \mathcal{Y}} p(x | y)p(y)}$$



# Naive Bayes models

## Naive Bayes assumption

The input features are **conditionally independent** given the label:

$$p(x | y) = \prod_{i=1}^n p(x_i | y) .$$

- ▶ A strong assumption, but works surprisingly well in practice.
- ▶  $p(x_i | y)$  doesn't have to be a categorical distribution (e.g., Gaussian distribution)

**Inference:**  $y = \arg \max_{y \in \mathcal{Y}} p_{\theta}(y | x)$  [board]

## Maximum likelihood estimation

**Task:** estimate parameters  $\theta$  of a distribution  $p(y; \theta)$  given i.i.d. samples  $D = (y_1, \dots, y_N)$  from the distribution.

**Goal:** find the parameters that make the observed data most probable.

## Maximum likelihood estimation

**Task:** estimate parameters  $\theta$  of a distribution  $p(y; \theta)$  given i.i.d. samples  $D = (y_1, \dots, y_N)$  from the distribution.

**Goal:** find the parameters that make the observed data most probable.

**Likelihood function** of  $\theta$  given  $D$ :

$$L(\theta; D) \stackrel{\text{def}}{=} p(D; \theta) = \prod_{i=1}^N p(y_i; \theta) .$$

**Maximum likelihood estimator:**

$$\hat{\theta} = \arg \max_{\theta \in \Theta} L(\theta; D) = \arg \max_{\theta \in \Theta} \sum_{i=1}^N \log p(y_i; \theta) \quad (4)$$

## MLE and ERM

ERM:

$$\min \sum_{i=1}^N \ell(x^{(i)}, y^{(i)}, \theta)$$

## MLE and ERM

ERM:

$$\min \sum_{i=1}^N \ell(x^{(i)}, y^{(i)}, \theta)$$

MLE:

$$\max \sum_{i=1}^N \log p(y^{(i)} | x^{(i)}; \theta)$$

## MLE and ERM

ERM:

$$\min \sum_{i=1}^N \ell(x^{(i)}, y^{(i)}, \theta)$$

MLE:

$$\max \sum_{i=1}^N \log p(y^{(i)} | x^{(i)}; \theta)$$

What's the connection between MLE and ERM?

MLE is equivalent to ERM with the **negative log-likelihood** (NLL) loss function:

$$\ell_{\text{NLL}}(x^{(i)}, y^{(i)}, \theta) \stackrel{\text{def}}{=} -\log p(y^{(i)} | x^{(i)}; \theta)$$

## MLE for our Naive Bayes model

[board]

## MLE solution for our Naive Bayes model

$\text{count}(w, y) \stackrel{\text{def}}{=} \text{frequency of } w \text{ in documents with label } y$

$$p_{\text{MLE}}(w | y) = \frac{\text{count}(w, y)}{\sum_{w \in \mathcal{V}} \text{count}(w, y)}$$

= how often the word occur in positive/negative documents

$$p_{\text{MLE}}(y = k) = \frac{\sum_{i=1}^N \mathbb{I}(y^{(i)} = k)}{N}$$

= fraction of positive/negative documents



## MLE solution for our Naive Bayes model

$\text{count}(w, y) \stackrel{\text{def}}{=} \text{frequency of } w \text{ in documents with label } y$

$$p_{\text{MLE}}(w | y) = \frac{\text{count}(w, y)}{\sum_{w \in \mathcal{V}} \text{count}(w, y)}$$

= how often the word occur in positive/negative documents

$$p_{\text{MLE}}(y = k) = \frac{\sum_{i=1}^N \mathbb{I}(y^{(i)} = k)}{N}$$

= fraction of positive/negative documents

**Smoothing:** reserve probability mass for unseen words

$$p(w | y) = \frac{\alpha + \text{count}(w, y)}{\sum_{w \in \mathcal{V}} \text{count}(w, y) + \alpha |\mathcal{V}|}$$

Laplace smoothing:  $\alpha = 1$

## Feature design

Naive Bayes doesn't have to use single words as features

- ▶ Lexicons, e.g., LIWC.
- ▶ Task-specific features, e.g., is the email subject all caps.
- ▶ Bytes and characters, e.g., used in language ID detection.

## Summary of Naive Bayes models

- ▶ Modeling: the conditional independence assumption simplifies the problem
- ▶ Learning: MLE (or ERM with negative log-likelihood loss)
- ▶ Inference: very fast (adding up scores of each word)

# Table of Contents

Generative models: naive Bayes

Discriminative models: logistic regression

Regularization, model selection, evaluation

## Discriminative models

Idea: directly model the conditional distribution  $p(y | x)$

## Discriminative models

Idea: directly model the conditional distribution  $p(y | x)$

	generative models	discriminative models
modeling	joint: $p(x, y)$	conditional: $p(y   x)$
assumption on $y$	yes	yes
assumption on $x$	yes	no
development	generative story	feature extractor

## Model $p(y | x)$

How to model  $p(y | x)$ ?

$y$  is a Bernoulli variable:

$$p(y | x) = \alpha^y (1 - \alpha)^{(1-y)}$$

## Model $p(y | x)$

How to model  $p(y | x)$ ?

$y$  is a Bernoulli variable:

$$p(y | x) = \alpha^y (1 - \alpha)^{(1-y)}$$

Bring in  $x$ :

$$p(y | x) = h(x)^y (1 - h(x))^{(1-y)} \quad h(x) \in [0, 1]$$



## Model $p(y | x)$

How to model  $p(y | x)$ ?

$y$  is a Bernoulli variable:

$$p(y | x) = \alpha^y (1 - \alpha)^{(1-y)}$$

Bring in  $x$ :

$$p(y | x) = h(x)^y (1 - h(x))^{(1-y)} \quad h(x) \in [0, 1]$$

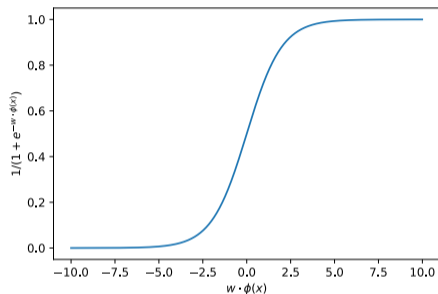
Parametrize  $h(x)$  using a linear function:

$$h(x) = w \cdot \phi(x) + b \quad \phi: \mathcal{X} \rightarrow \mathbb{R}^d$$

Problem:  $h(x) \in \mathbb{R}$  (score)

## Logistic regression

Map  $w \cdot \phi(x) \in \mathbb{R}$  to a probability by the **logistic function**



$$p(y = 1 \mid x; w) = \frac{1}{1 + e^{-w \cdot \phi(x)}} \quad (y \in \{0, 1\})$$

$$p(y = k \mid x; w) = \frac{e^{w_k \cdot \phi(x)}}{\sum_{i \in \mathcal{Y}} e^{w_i \cdot \phi(x)}} \quad (y \in \{1, \dots, K\})$$

“softmax”

# Inference

$$\hat{y} = \arg \max_{k \in \mathcal{Y}} p(y = k \mid x; w) \quad (5)$$

$$= \arg \max_{k \in \mathcal{Y}} \frac{e^{w_k \cdot \phi(x)}}{\sum_{i \in \mathcal{Y}} e^{w_i \cdot \phi(x)}} \quad (6)$$

$$= \arg \max_{k \in \mathcal{Y}} e^{w_k \cdot \phi(x)} \quad (7)$$

$$= \arg \max_{k \in \mathcal{Y}} \underbrace{w_k \cdot \phi(x)}_{\text{score for class } k} \quad (8)$$

# MLE for logistic regression

[board]

# MLE for logistic regression

[board]

- ▶ Likelihood function is concave
- ▶ No closed-form solution
- ▶ Use gradient ascent

## BoW representation

**Feature extractor:**  $\phi: \mathcal{V}^n \rightarrow \mathbb{R}^d$ .

**Idea:** a sentence is the “sum” of words.

Example:

$\mathcal{V} = \{the, a, an, in, for, penny, pound\}$

sentence = *in for a penny, in for a pound*

$x = (in, for, a, penny, in, for, a, pound)$

$$\phi_{\text{BoW}}(x) = \sum_{i=1}^n \phi_{\text{one-hot}}(x_i)$$

[board]

## Compare with naive Bayes

- ▶ Our naive Bayes model ( $x_i \in \{1, \dots, |\mathcal{V}|\}$ ):

$$X_i \mid Y = y \sim \text{Categorical}(\theta_{1,y}, \dots, \theta_{|\mathcal{V}|,y}) .$$

- ▶ The naive Bayes generative story produces a BoW vector following a multinomial distribution:

$$\phi_{\text{BoW}}(X) \mid Y = y \sim \text{Multinomial}(\theta_{1,y}, \dots, \theta_{|\mathcal{V}|,y}, n) .$$

## Compare with naive Bayes

- ▶ Our naive Bayes model ( $x_i \in \{1, \dots, |\mathcal{V}|\}$ ):

$$X_i \mid Y = y \sim \text{Categorical}(\theta_{1,y}, \dots, \theta_{|\mathcal{V}|,y}) .$$

- ▶ The naive Bayes generative story produces a BoW vector following a multinomial distribution:

$$\phi_{\text{BoW}}(X) \mid Y = y \sim \text{Multinomial}(\theta_{1,y}, \dots, \theta_{|\mathcal{V}|,y}, n) .$$

- ▶ Both multinomial naive Bayes and logistic regression learn a linear separator  $w \cdot \phi_{\text{BoW}}(x) + b = 0$ .



## Compare with naive Bayes

- ▶ Our naive Bayes model ( $x_i \in \{1, \dots, |\mathcal{V}|\}$ ):

$$X_i \mid Y = y \sim \text{Categorical}(\theta_{1,y}, \dots, \theta_{|\mathcal{V}|,y}) .$$

- ▶ The naive Bayes generative story produces a BoW vector following a multinomial distribution:

$$\phi_{\text{BoW}}(X) \mid Y = y \sim \text{Multinomial}(\theta_{1,y}, \dots, \theta_{|\mathcal{V}|,y}, n) .$$

- ▶ Both multinomial naive Bayes and logistic regression learn a linear separator  $w \cdot \phi_{\text{BoW}}(x) + b = 0$ .

Question: what's the advantage of using logistic regression?

## Feature extractor

Logistic regression allows for richer features.

Define each feature as a function  $\phi_i: \mathcal{X} \rightarrow \mathbb{R}$ .

$$\phi_1(x) = \begin{cases} 1 & x \text{ contains "happy"} \\ 0 & \text{otherwise} \end{cases},$$

$$\phi_2(x) = \begin{cases} 1 & x \text{ contains words with suffix "yyy"} \\ 0 & \text{otherwise} \end{cases}.$$

In practice, use a dictionary

```
feature_vector["prefix=un+suffix=ing"] = 1
```

# Feature vectors for multiclass classification

Multinomial logistic regression

$$p(y = k \mid x; w) = \frac{e^{w_k \cdot \phi(x)}}{\sum_{i \in \mathcal{Y}} e^{w_i \cdot \phi(x)}} \quad (y \in \{1, \dots, K\})$$



$$p(y = k \mid x; w) = \frac{e^{w \cdot \Psi(x, k)}}{\sum_{i \in \mathcal{Y}} e^{w \cdot \Psi(x, i)}} \quad (y \in \{1, \dots, K\})$$

scores of each class  $\rightarrow$  compatibility of an input and a label

# Feature vectors for multiclass classification

Multinomial logistic regression

$$p(y = k | x; w) = \frac{e^{w_k \cdot \phi(x)}}{\sum_{i \in \mathcal{Y}} e^{w_i \cdot \phi(x)}} \quad (y \in \{1, \dots, K\})$$



$$p(y = k | x; w) = \frac{e^{w \cdot \Psi(x, k)}}{\sum_{i \in \mathcal{Y}} e^{w \cdot \Psi(x, i)}} \quad (y \in \{1, \dots, K\})$$

scores of each class  $\rightarrow$  compatibility of an input and a label

Multivector construction of  $\Psi(x, y)$ :

$$\Psi(x, 1) \stackrel{\text{def}}{=} \left[ \underbrace{0}_{y=0}, \underbrace{\phi(x)}_{y=1}, \dots, \underbrace{0}_{y=K} \right]$$

## N-gram features

Potential problems with the the BoW representation?

## N-gram features

Potential problems with the the BoW representation?

**N-gram** features:

*in for a penny , in for a pound*

- ▶ Unigram: in, for, a, ...
- ▶ Bigram: in/for, for/a, a/penny, ...
- ▶ Trigram: in/for/a, for/a/penny, ...

What's the pros/cons of using higher order n-grams?

# Table of Contents

Generative models: naive Bayes

Discriminative models: logistic regression

Regularization, model selection, evaluation

## Error decomposition

(ignoring optimization error)

$$\text{risk}(h) - \text{risk}(h^*) = \text{approximation error} + \text{estimation error}$$

[board]



## Error decomposition

(ignoring optimization error)

$$\text{risk}(h) - \text{risk}(h^*) = \text{approximation error} + \text{estimation error}$$

[board]

Larger hypothesis class: approximation error  $\downarrow$ , estimation error  $\uparrow$

Smaller hypothesis class: approximation error  $\uparrow$ , estimation error  $\downarrow$

How to control the size of the hypothesis class?

## Reduce the dimensionality

Linear predictors:  $\mathcal{H} = \{w : w \in \mathbb{R}^d\}$

Reduce the number of features:  
[discussion]

## Reduce the dimensionality

Linear predictors:  $\mathcal{H} = \{w : w \in \mathbb{R}^d\}$

Reduce the number of features:  
[discussion]

Other predictors:

- ▶ Depth of decision trees
- ▶ Degree of polynomials
- ▶ Number of decision stumps in boosting

# Regularization

Reduce the “size” of  $w$ :

$$\min_w \underbrace{\frac{1}{N} \sum_{i=1}^N L(x^{(i)}, y^{(i)}, w)}_{\text{average loss}} + \underbrace{\frac{\lambda}{2} \|w\|_2^2}_{\ell_2 \text{ norm}}$$

Why is small norm good? Small change in the input doesn't cause large change in the output.

[board]

## Gradient descent with $\ell_2$ regularization

Run SGD on

$$\min_w \underbrace{\frac{1}{N} \sum_{i=1}^N L(x^{(i)}, y^{(i)}, w)}_{\text{average loss}} + \underbrace{\frac{\lambda}{2} \|w\|_2^2}_{\ell_2 \text{ norm}}$$

Also called **weight decay** in the deep learning literature:

$$w \leftarrow w - \eta(\nabla_w L(x, y, w) + \lambda w)$$

Shrink  $w$  in each update.

## Hyperparameter tuning

**Hyperparameters:** parameters of the learning algorithm (not the model)

Example: use MLE to learn a logistic regression model using BoW features

## Hyperparameter tuning

**Hyperparameters:** parameters of the learning algorithm (not the model)

Example: use MLE to learn a logistic regression model using BoW features

# Hyperparameter tuning

**Hyperparameters:** parameters of the learning algorithm (not the model)

Example: use MLE to learn a logistic regression model using BoW features

How do we select hyperparameters?

- Pick those minimizing the training error?

- Pick those minimizing the test error?



# Validation

**Validation set:** a subset of the training data reserved for tuning the learning algorithm (also called the **development set**).

***K*-fold cross validation**

[board]

# Validation

**Validation set:** a subset of the training data reserved for tuning the learning algorithm (also called the **development set**).

***K*-fold cross validation**

[board]

It's important to look at the data and errors during development, but **not the test set**.