# CS 6770 Natural Language Processing

## Text Classification (I): Logistic Regression

Yangfeng Ji

Information and Language Processing Lab
Department of Computer Science
University of Virginia

UNIVERSITY of VIRGINIA | ENGINEERING

# Overview

# Problem Definition

# Case I: Sentiment Analysis



[Pang et al., 2002]

Example topics

- ▶ Business
- ▶ Arts
- ▶ Technology
- ▶ Sports
- ▶ . . .

A text classifier demo on Hugging Face webpage.



Link

# Classification

- **Input**: a text $x$
  - Example: a product review on Amazon
- **Output**: $y \in \mathcal{Y}$, where $\mathcal{Y}$ is the predefined category set (sample space)
  - Example: $\mathcal{Y} = \{\textsc{Positive}, \textsc{Negative}\}$

---

[1] In this course, we use $x$ for both text and its representation with no distinction

- **Input**: a text $x$
  - Example: a product review on Amazon
- **Output**: $y \in \mathcal{Y}$, where $\mathcal{Y}$ is the predefined category set (sample space)
  - Example: $\mathcal{Y} = \{\text{Positive}, \text{Negative}\}$

The pipeline of text classification:[1]

Text $\longrightarrow$ | Numeric Vector $x$ | $\longrightarrow$ | Classifier | $\longrightarrow$ Category $y$

---

[1] In this course, we use $x$ for both text and its representation with no distinction

## Probabilistic Formulation

With the conditional probability $P(Y \mid X)$, the prediction on $Y$ for a given text $X = x$ is

$$\hat{y} = \underset{y \in \mathcal{Y}}{\operatorname{argmax}}\, P(Y = y \mid X = x) \tag{1}$$

# Probabilistic Formulation

With the conditional probability $P(Y \mid X)$, the prediction on $Y$ for a given text $X = x$ is

$$\hat{y} = \operatorname*{argmax}_{y \in \mathcal{Y}} P(Y = y \mid X = x) \tag{1}$$

Or, for simplicity

$$\hat{y} = \operatorname*{argmax}_{y \in \mathcal{Y}} P(y \mid x) \tag{2}$$

Recall

- The formulation defined in the previous slide

$$\hat{y} = \underset{y \in \mathcal{Y}}{\arg\max}\, P(Y = y \mid X = x) \qquad (3)$$

- The pipeline of text classification

$$\text{Text} \longrightarrow \boxed{\text{Numeric Vector } x} \longrightarrow \boxed{\text{Classifier}} \longrightarrow \text{Category } y$$

# Key Questions

Recall

- ▶ The formulation defined in the previous slide

$$\hat{y} = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} P(Y = y \mid X = x) \qquad (3)$$

- ▶ The pipeline of text classification

Text $\longrightarrow$ Numeric Vector $x$ $\longrightarrow$ Classifier $\longrightarrow$ Category $y$

Building a text classifier is about answering the following two questions

1. How to represent a text as $x$?

2. How to estimate $P(y \mid x)$?

Recall

- ▶ The formulation defined in the previous slide

$$\hat{y} = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} P(Y = y \mid X = x) \tag{3}$$

- ▶ The pipeline of text classification

Text $\longrightarrow$ ⟦Numeric Vector $x$⟧ $\longrightarrow$ ⟦Classifier⟧ $\longrightarrow$ Category $y$

Building a text classifier is about answering the following two questions

1. How to represent a text as $x$?
   - ▶ Bag-of-words representation
2. How to estimate $P(y \mid x)$?

Recall

▶ The formulation defined in the previous slide

$$\hat{y} = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} P(Y = y \mid X = x) \tag{3}$$

▶ The pipeline of text classification

Text $\longrightarrow$ | Numeric Vector $x$ | $\longrightarrow$ | Classifier | $\longrightarrow$ Category $y$

Building a text classifier is about answering the following two questions

1. How to represent a text as $x$?
   ▶ Bag-of-words representation
2. How to estimate $P(y \mid x)$?
   ▶ Logistic regression models

# Key Questions

Recall

▶ The formulation defined in the previous slide

$$\hat{y} = \underset{y \in \mathcal{Y}}{\arg\max}\, P(Y = y \mid X = x) \qquad (3)$$

▶ The pipeline of text classification

Text ⟶ Numeric Vector $x$ ⟶ Classifier ⟶ Category $y$

Building a text classifier is about answering the following two questions

1. How to represent a text as $x$?
   ▶ Bag-of-words representation
2. How to estimate $P(y \mid x)$?
   ▶ Logistic regression models
   ▶ Neural network classifiers

# Key Questions

Recall

▶ The formulation defined in the previous slide

$$\hat{y} = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \, P(Y = y \mid X = x) \qquad (3)$$

▶ The pipeline of text classification

Text $\longrightarrow$ | Numeric Vector $x$ | $\longrightarrow$ | Classifier | $\longrightarrow$ Category $y$

Building a text classifier is about answering the following two questions

1. How to represent a text as $x$?
   ▶ Bag-of-words representation
2. How to estimate $P(y \mid x)$?
   ▶ Logistic regression models
   ▶ Neural network classifiers
   ▶ Other classifiers: Naive Bayes classifier, support vector classifier, random forest, etc.

# Bag-of-Words Representation

# Bag-of-Words Representation

## Example Texts

Text 1: I love coffee.

Text 2: I don't like tea.

# Bag-of-Words Representation

## Example Texts

Text 1: I love coffee.

Text 2: I don't like tea.

**Step I**: convert a text into a collection of tokens (e.g., tokenization)

## Tokenized Texts

Tokenized text 1: `I love coffee`

Tokenized text 2: `I don t like tea`

# Bag-of-Words Representation

## Example Texts

Text 1: I love coffee.

Text 2: I don't like tea.

**Step I**: convert a text into a collection of tokens (e.g., tokenization)

## Tokenized Texts

Tokenized text 1: I love coffee

Tokenized text 2: I don t like tea

**Step II**: build a dictionary/vocabulary

## Vocabulary

{I love coffee don t like tea}

# Bag-of-Words Representations

**Step III**: based on the vocab, convert each text into a numeric representation as

## Bag-of-Words Representations

|  | I | love | coffee | don | t | like | tea |
|---|---|---|---|---|---|---|---|
| $x^{(1)} =$ | [1 | 1 | 1 | 0 | 0 | 0 | 0]$^{\mathsf{T}}$ |
| $x^{(2)} =$ | [1 | 0 | 0 | 1 | 1 | 1 | 1]$^{\mathsf{T}}$ |

# Bag-of-Words Representations

**Step III**: based on the vocab, convert each text into a numeric representation as

## Bag-of-Words Representations

|  | I | love | coffee | don | t | like | tea |
|---|---|---|---|---|---|---|---|
| $x^{(1)} =$ | $[1$ | 1 | 1 | 0 | 0 | 0 | $0]^{\mathsf{T}}$ |
| $x^{(2)} =$ | $[1$ | 0 | 0 | 1 | 1 | 1 | $1]^{\mathsf{T}}$ |

The pipeline of text classification:

Text $\longrightarrow$ | Numeric Vector $x$ | $\longrightarrow$ | Classifier | $\longrightarrow$ Category $y$

Bag-of-words

Representation

# Preprocessing for Building Vocab

1. Convert all characters to lowercase

$$UVa, UVA \rightarrow uva$$

# Preprocessing for Building Vocab

1. Convert all characters to lowercase

$$\texttt{UVa, UVA} \rightarrow \texttt{uva}$$

Shall we always convert all words to lowercase?

$$\texttt{Apple vs. apple}$$

1. Convert all characters to lowercase

$$\text{UVa}, \text{UVA} \rightarrow \text{uva}$$

Shall we always convert all words to lowercase?

$$\text{Apple vs. apple}$$

2. Map low frequency words to a special token ⟨unk⟩



Zipf's law: $\text{freq}(w_t) \propto 1/r_t$

where $\text{freq}(w_t)$ is the frequency of word $w_t$ and $r_t$ is the rank of this word

It is critical to keep in mind about what information is preserved in bag-of-words representations:

- ▶ Keep:
  - ▶ words in texts

It is critical to keep in mind about what information is preserved in bag-of-words representations:

- ▶ Keep:
  - ▶ words in texts
- ▶ Lose:
  - ▶ word order

```
I love coffee don t like tea
```

# Information Embedded in BoW Representations

It is critical to keep in mind about what information is preserved in bag-of-words representations:

- ▶ Keep:
  - ▶ words in texts
- ▶ Lose:
  - ▶ word order

    ```
    I love coffee don t like tea
    ```

  - ▶ sentence boundary
  - ▶ sentence order
  - ▶ . . .

Read between the lines ...

From Regina Barzilay's lecture note

# Why Sentence Order Matters?

Read between the lines ...



Pool For Members Only. Use The Toilets, Not The Pool.

# Why Sentence Order Matters?

Read between the lines ...



Use The Toilets, Not The Pool. Pool For Members Only.

# Case Study: Sentiment Analysis

# A Dummy Predictor

Consider the following toy example (adding one more example to make it more interesting)

## Tokenized Texts

|                  | Text $X$         | Label $Y$ |
| ---------------- | ---------------- | --------- |
| Tokenized text 1 | I love coffee    | POSITIVE  |
| Tokenized text 2 | I don t like tea | NEGATIVE  |
| Tokenized text 3 | I like coffee    | POSITIVE  |

---

[2]The evaluation of classifiers will be discussed in one of the future lectures.

[3]It can be a competitive baseline in practice, particularly for unbalanced datasets

# A Dummy Predictor

Consider the following toy example (adding one more example to make it more interesting)

## Tokenized Texts

|  | Text $X$ | Label $Y$ |
|---|---|---|
| Tokenized text 1 | I love coffee | POSITIVE |
| Tokenized text 2 | I don t like tea | NEGATIVE |
| Tokenized text 3 | I like coffee | POSITIVE |

What is the simplest classifier that we can constructed based on this small dataset?

---

[2]The evaluation of classifiers will be discussed in one of the future lectures.

[3]It can be a competitive baseline in practice, particularly for unbalanced datasets

# A Dummy Predictor

Consider the following toy example (adding one more example to make it more interesting)

## Tokenized Texts

|  | Text $X$ | Label $Y$ |
|---|---|---|
| Tokenized text 1 | I love coffee | POSITIVE |
| Tokenized text 2 | I don t like tea | NEGATIVE |
| Tokenized text 3 | I like coffee | POSITIVE |

What is the simplest classifier that we can constructed based on this small dataset?

▶ Predict every text as POSITIVE

_____

[2] The evaluation of classifiers will be discussed in one of the future lectures.
[3] It can be a competitive baseline in practice, particularly for unbalanced datasets

# A Dummy Predictor

Consider the following toy example (adding one more example to make it more interesting)

## Tokenized Texts

|  | Text $X$ | Label $Y$ |
|---|---|---|
| Tokenized text 1 | I love coffee | Positive |
| Tokenized text 2 | I don t like tea | Negative |
| Tokenized text 3 | I like coffee | Positive |

What is the simplest classifier that we can constructed based on this small dataset?

▶ Predict every text as Positive
▶ 66.7% prediction accuracy on this dataset[2]

---

[2]The evaluation of classifiers will be discussed in one of the future lectures.
[3]It can be a competitive baseline in practice, particularly for unbalanced datasets

# A Dummy Predictor

Consider the following toy example (adding one more example to make it more interesting)

## Tokenized Texts

|  | Text $X$ | Label $Y$ |
|---|---|---|
| Tokenized text 1 | I love coffee | Positive |
| Tokenized text 2 | I don t like tea | Negative |
| Tokenized text 3 | I like coffee | Positive |

What is the simplest classifier that we can constructed based on this small dataset?

- ▶ Predict every text as Positive
- ▶ 66.7% prediction accuracy on this dataset[2]
- ▶ It has a name: majority baseline[3]

[2]The evaluation of classifiers will be discussed in one of the future lectures.
[3]It can be a competitive baseline in practice, particularly for unbalanced datasets

# A Simple Predictor

Consider the following toy example, again

## Tokenized Texts

Tokenized text 1: `I love coffee`
Tokenized text 2: `I don t like tea`
Tokenized text 3: `I like coffee`

Consider the following toy example, again

## Tokenized Texts

Tokenized text 1: I love coffee
Tokenized text 2: I don t like tea
Tokenized text 3: I like coffee

|  | I | love | coffee | don | t | like | tea |
|---|---|---|---|---|---|---|---|
| $x^{(1)}$ | [1 | 1 | 1 | 0 | 0 | 0 | 0 ]$^\mathsf{T}$ |
| $w_{\mathrm{Pos}}$ | [0 | 1 | 0 | 0 | 0 | 1 | 0 ]$^\mathsf{T}$ |
| $w_{\mathrm{Neg}}$ | [0 | 0 | 0 | 1 | 0 | 0 | 0 ]$^\mathsf{T}$ |

Consider the following toy example, again

## Tokenized Texts

Tokenized text 1: I love coffee
Tokenized text 2: I don t like tea
Tokenized text 3: I like coffee

|  | I | love | coffee | don | t | like | tea |
|---|---|---|---|---|---|---|---|
| $x^{(1)}$ | [1 | 1 | 1 | 0 | 0 | 0 | 0 ]$^\mathsf{T}$ |
| $w_{\text{Pos}}$ | [0 | 1 | 0 | 0 | 0 | 1 | 0 ]$^\mathsf{T}$ |
| $w_{\text{Neg}}$ | [0 | 0 | 0 | 1 | 0 | 0 | 0 ]$^\mathsf{T}$ |

The prediction of sentiment polarity can be formulated as the following

$$w_{\text{Pos}}^{\mathsf{T}} x = 1 > w_{\text{Neg}}^{\mathsf{T}} x = 0 \qquad (4)$$

# A Simple Predictor

Consider the following toy example, again

## Tokenized Texts

Tokenized text 1: `I love coffee`

Tokenized text 2: `I don t like tea`

Tokenized text 3: `I like coffee`

|  | I | love | coffee | don | t | like | tea |
|---|---|---|---|---|---|---|---|
| $x^{(1)}$ | [1 | 1 | 1 | 0 | 0 | 0 | 0 ]$^\mathsf{T}$ |
| $w_{\text{Pos}}$ | [0 | 1 | 0 | 0 | 0 | 1 | 0 ]$^\mathsf{T}$ |
| $w_{\text{Neg}}$ | [0 | 0 | 0 | 1 | 0 | 0 | 0 ]$^\mathsf{T}$ |

The prediction of sentiment polarity can be formulated as the following

$$w_{\text{Pos}}^{\mathsf{T}} x = 1 > w_{\text{Neg}}^{\mathsf{T}} x = 0 \qquad (4)$$

Essentially, it is equivalent to counting the positive and negative words.

17

# Another Example

The limitation of word counting

|  | I | love | coffee | don | t | like | tea |
|---|---|---|---|---|---|---|---|
| $x^{(2)}$ | [1 | 0 | 0 | 1 | 1 | 1 | 1 ]$^\mathsf{T}$ |
| $w_{\text{Pos}}$ | [0 | 1 | 0 | 0 | 0 | 1 | 0 ]$^\mathsf{T}$ |
| $w_{\text{Neg}}$ | [0 | 0 | 0 | 1 | 0 | 0 | 0 ]$^\mathsf{T}$ |

The limitation of word counting

|  | I | love | coffee | don | t | like | tea |
|---|---|---|---|---|---|---|---|
| $x^{(2)}$ | [1 | 0 | 0 | 1 | 1 | 1 | 1 ]$^\mathsf{T}$ |
| $w_{\text{Pos}}$ | [0 | 1 | 0 | 0 | 0 | 1 | 0 ]$^\mathsf{T}$ |
| $w_{\text{Neg}}$ | [0 | 0 | 0 | 1 | 0 | 0 | 0 ]$^\mathsf{T}$ |

▶ Different words should contribute differently. e.g., not vs. dislike

## Another Example

The limitation of word counting

|  | I | love | coffee | don | t | like | tea |
|---|---|---|---|---|---|---|---|
| $x^{(2)}$ | [1 | 0 | 0 | 1 | 1 | 1 | 1 ]$^\mathsf{T}$ |
| $w_{\text{Pos}}$ | [0 | 1 | 0 | 0 | 0 | 1 | 0 ]$^\mathsf{T}$ |
| $w_{\text{Neg}}$ | [0 | 0 | 0 | 1 | 0 | 0 | 0 ]$^\mathsf{T}$ |

▶ Different words should contribute differently. e.g., not vs. dislike
▶ Sentiment word lists are definitely incomplete

### A Positive Review of Coffee without Sentiment Words

*Its aroma was of earth and smoke. The first sip was an abrupt, bitter jolt that commanded my full attention. Any trace of morning fatigue vanished. I finished the entire cup without pause and immediately brewed another. This is the coffee I will be drinking from now on.*

# Logistic Regression

# Linear Models

Directly modeling a linear classifier as

$$h_y(x) = w_y^\mathsf{T} x + b_y \tag{5}$$

with

- $x \in \mathbb{N}^V$: vector, bag-of-words representation
- $w_y \in \mathbb{R}^V$: vector, classification weights associated with label $y$
- $b_y \in \mathbb{R}$: scalar, label bias in the training set $y$

# Linear Models

Directly modeling a linear classifier as

$$h_y(x) = w_y^\mathsf{T} x + b_y \tag{5}$$

with

- $x \in \mathbb{N}^V$: vector, bag-of-words representation
- $w_y \in \mathbb{R}^V$: vector, classification weights associated with label $y$
- $b_y \in \mathbb{R}$: scalar, label bias in the training set $y$

## About Label Bias

Consider a case with highly-imbalanced examples, where we have 90 positive examples and 10 negative examples in the training set. With

$$b_{\text{Pos}} > b_{\text{Neg}},$$

a classifier can get 90% predictions correct without even resorting the texts.

# Logistic Regression

Rewrite the linear decision function in the log probabilitic form

$$\log P(y \mid x) \propto \underbrace{w_y^\mathsf{T} x + b_y}_{h_y(x)} \tag{6}$$

# Logistic Regression

Rewrite the linear decision function in the log probabilitic form

$$\log P(y \mid x) \propto \underbrace{w_y^\mathsf{T} x + b_y}_{h_y(x)} \tag{6}$$

or, the probabilistic form is

$$P(y \mid x) \propto \exp(w_y^\mathsf{T} x + b_y) \tag{7}$$

## Logistic Regression

Rewrite the linear decision function in the log probabilitic form

$$\log P(y \mid x) \propto \underbrace{w_y^\mathsf{T} x + b_y}_{h_y(x)} \tag{6}$$

or, the probabilistic form is

$$P(y \mid x) \propto \exp(w_y^\mathsf{T} x + b_y) \tag{7}$$

To make sure $P(y \mid x)$ is a valid definition of probability, we need to make sure $\sum_y P(y \mid x) = 1$,

$$P(y \mid x) = \frac{\exp(w_y^\mathsf{T} x + b_y)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\mathsf{T} x + b_{y'})} \tag{8}$$

## Alternative Form

Rewriting $x$ and $w$ as

- $x^\mathsf{T} = [x_1, x_2, \cdots, x_V, 1]$
- $w_y^\mathsf{T} = [w_1, w_2, \cdots, w_V, b_y]$

allows us to have a more concise form

$$P(y \mid x) = \frac{\exp(w_y^\mathsf{T} x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\mathsf{T} x)} \tag{9}$$

# Alternative Form

Rewriting $x$ and $w$ as

- $x^\mathsf{T} = [x_1, x_2, \cdots, x_V, 1]$
- $w_y^\mathsf{T} = [w_1, w_2, \cdots, w_V, b_y]$

allows us to have a more concise form

$$P(y \mid x) = \frac{\exp(w_y^\mathsf{T} x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\mathsf{T} x)} \tag{9}$$

Comments:

- $\frac{\exp(a)}{\sum_{a'} \exp(a')}$ is the softmax function
- This form works with any size of $\mathcal{Y}$ — it does not have to be a binary classification problem.

# Binary Classifier

Assume $\mathcal{Y} = \{\text{NEG}, \text{POS}\}$, then the corresponding logistic regression classifier with $Y = \text{Pos}$ is

$$P(Y = \text{Pos} \mid x) = \frac{1}{1 + \exp(-w^\mathsf{T} x)} \tag{10}$$

where $w$ is the only parameter.

# Binary Classifier

Assume $\mathcal{Y} = \{\text{NEG}, \text{POS}\}$, then the corresponding logistic regression classifier with $Y = \text{Pos}$ is

$$P(Y = \text{Pos} \mid x) = \frac{1}{1 + \exp(-w^\mathsf{T} x)} \tag{10}$$

where $w$ is the only parameter.

- $P(Y = \text{NEG} \mid x) = 1 - P(Y = \text{POS} \mid x)$

# Binary Classifier

Assume $\mathcal{Y} = \{\text{NEG}, \text{POS}\}$, then the corresponding logistic regression classifier with $Y = \text{Pos}$ is

$$P(Y = \text{Pos} \mid x) = \frac{1}{1 + \exp(-w^\mathsf{T} x)} \tag{10}$$

where $w$ is the only parameter.

▶ $P(Y = \text{NEG} \mid x) = 1 - P(Y = \text{POS} \mid x)$
▶ $\frac{1}{1+\exp(-z)}$ is the Sigmoid function

# Demo

Link to the demo

... of building a logistic regression classifier

$$P(y \mid x) = \frac{\exp(w_y^\mathsf{T} x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\mathsf{T} x)} \tag{11}$$

▶ How to learn the parameters $W = \{w_y\}_{y \in \mathcal{Y}}$?

... of building a logistic regression classifier

$$P(y \mid x) = \frac{\exp(w_y^\mathsf{T} x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\mathsf{T} x)} \tag{11}$$

▶ How to learn the parameters $W = \{w_y\}_{y \in \mathcal{Y}}$?
▶ Can $x$ be better than the bag-of-words representations?

With a collection of training examples $\{(x^{(i)}, y^{(i)})\}_{i=1}^{m}$, the likelihood function of $\{w_y\}_{y \in \mathcal{Y}}$ is

$$L(W) = \prod_{i=1}^{m} P(y^{(i)} \mid x^{(i)}) \tag{12}$$

and the log-likelihood function is

$$\ell(\{w_y\}) = \sum_{i=1}^{m} \log P(y^{(i)} \mid x^{(i)}) \tag{13}$$

With the definition of a LR model

$$P(y \mid x) = \frac{\exp(w_y^\mathsf{T} x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\mathsf{T} x)} \tag{14}$$

the log-likelihood function is

$$\ell(W) = \sum_{i=1}^{m} \log P(y^{(i)} \mid x^{(i)}) \tag{15}$$

$$= \sum_{i=1}^{m} \left\{ w_{y^{(i)}}^\mathsf{T} x^{(i)} - \log \sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\mathsf{T} x^{(i)}) \right\} \tag{16}$$

Given the training examples $\{(x^{(i)}, y^{(i)})\}_{i=1}^{m}$, $\ell(W)$ is a function of $W = \{w_y\}$.

MLE is equivalent to minimize the Negative Log-Likelihood (NLL) as

$$\begin{aligned}
\text{NLL}(W) &= -\ell(W) \\
&= \sum_{i=1}^{m} \Big\{ -w_{y^{(i)}}^{\mathsf{T}} x^{(i)} + \log \sum_{y' \in \mathcal{Y}} \exp(w_{y'}^{\mathsf{T}} x) \Big\}
\end{aligned}$$

then, the parameter $w_y$ associated with label $y$ can be updated as

$$w_y \leftarrow w_y - \eta \cdot \frac{\partial \text{NLL}(\{w_y\})}{\partial w_y}, \quad \forall y \in \mathcal{Y} \tag{17}$$

where $\eta$ is called **learning rate**.

# Optimization with Gradient (II)

Two questions answered by the update equation

(1) which direction?
(2) how far it should go?

[Jurafsky and Martin, 2019]

Two questions answered by the update equation

(1) which direction?
(2) how far it should go?

$$w_y \leftarrow w_y - \underbrace{\eta}_{(2)} \cdot \underbrace{\frac{\partial \mathrm{NLL}(\{w_y\})}{\partial w_y}}_{(1)} \tag{18}$$

[Jurafsky and Martin, 2019]

Two questions answered by the update equation

(1) which direction?
(2) how far it should go?

$$w_y \leftarrow w_y - \underbrace{\eta}_{(2)} \cdot \underbrace{\frac{\partial \text{NLL}(\{w_y\})}{\partial w_y}}_{(1)} \tag{18}$$



Loss

one step
of gradient
descent

slope of loss at $w^1$
is negative

$w^1$       $w^{min}$       $w$

*0*       *(goal)*

[Jurafsky and Martin, 2019]

Steps for parameter estimation, given the current parameter $\{\boldsymbol{w}_y\}$

1. Compute the derivative

$$\frac{\partial \text{NLL}(\{\boldsymbol{w}_y\})}{\partial \boldsymbol{w}_y}, \quad \forall y \in \mathcal{Y}$$

2. Update parameters with

$$\boldsymbol{w}_y \leftarrow \boldsymbol{w}_y - \eta \cdot \frac{\partial \text{NLL}(\{\boldsymbol{w}_y\})}{\partial \boldsymbol{w}_y}, \quad \forall y \in \mathcal{Y}$$

3. If not done, retrun to step 1

Review: the pipeline of text classification:

Text $\longrightarrow$ | Numeric Vector $x$ | $\longrightarrow$ | Classifier | $\longrightarrow$ Category $y$

Bag-of-words          Logistic

Representation       Regression

# $L_2$ Regularization

# $L_2$ Regularization

The commonly used regularization trick is the $L_2$ regularization. For that, we need to redefine the objective function of LR by adding an additional item

$$\text{Loss}(\boldsymbol{W}) = \underbrace{\sum_{i=1}^{m} \big\{ - \boldsymbol{w}_{y^{(i)}}^{\mathsf{T}} \boldsymbol{x}^{(i)} + \log \sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{w}_{y'}^{\mathsf{T}} \boldsymbol{x}^{(i)}) \big\}}_{\text{NLL}}$$

<div align="right">(19)</div>

# $L_2$ Regularization

The commonly used regularization trick is the $L_2$ regularization. For that, we need to redefine the objective function of LR by adding an additional item

$$\text{Loss}(W) = \underbrace{\sum_{i=1}^{m} \big\{ -w_{y^{(i)}}^{\mathsf{T}} x^{(i)} + \log \sum_{y' \in \mathcal{Y}} \exp(w_{y'}^{\mathsf{T}} x^{(i)}) \big\}}_{\text{NLL}} + \underbrace{\frac{\lambda}{2} \cdot \sum_{y \in \mathcal{Y}} \|w_y\|_2^2}_{L_2 \text{ reg}}$$

$$(19)$$

- $\lambda$ is the regularization parameter

# $L_2$ Regularization in Gradient Descent

- ▶ The gradient of the loss function

$$\frac{\partial \text{Loss}(\boldsymbol{W})}{\partial \boldsymbol{w}_y} = \frac{\partial \text{NLL}(\boldsymbol{W})}{\partial \boldsymbol{w}_y} + \lambda \boldsymbol{w}_y \qquad (20)$$

# $L_2$ Regularization in Gradient Descent

▶ The gradient of the loss function

$$\frac{\partial \text{Loss}(\boldsymbol{W})}{\partial \boldsymbol{w}_y} = \frac{\partial \text{NLL}(\boldsymbol{W})}{\partial \boldsymbol{w}_y} + \lambda \boldsymbol{w}_y \tag{20}$$

▶ To minimize the loss, we need update the parameter as

$$\boldsymbol{w}_y \quad \leftarrow \quad \boldsymbol{w}_y - \eta \Big( \frac{\partial \text{NLL}(\boldsymbol{W})}{\partial \boldsymbol{w}_y} + \lambda \boldsymbol{w}_y \Big) \tag{21}$$

# $L_2$ Regularization in Gradient Descent

▶ The gradient of the loss function

$$\frac{\partial \text{Loss}(W)}{\partial w_y} = \frac{\partial \text{NLL}(W)}{\partial w_y} + \lambda w_y \tag{20}$$

▶ To minimize the loss, we need update the parameter as

$$w_y \quad \leftarrow \quad w_y - \eta \Big( \frac{\partial \text{NLL}(W)}{\partial w_y} + \lambda w_y \Big) \tag{21}$$

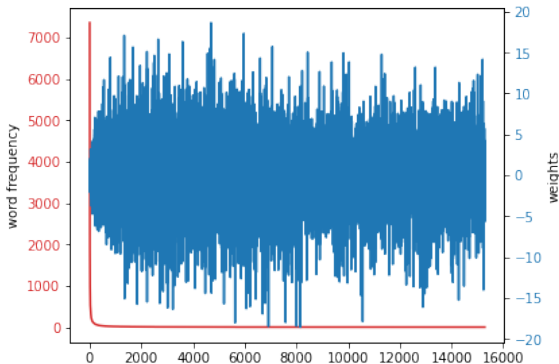$$= \quad (1 - \eta \lambda) \cdot w_y - \eta \frac{\partial \text{NLL}(W)}{\partial w_y}$$

▶ Depending on the strength (value) of $\lambda$, the regularization term tries to keep the parameter values close to 0, which to some extent can help avoid overfitting

In the demo code, we chose $\lambda = \frac{1}{C} = 0.001$ to approximate the case without regularization.

▶ Training accuracy: 99.89%
▶ Val accuracy: 52.21%

## Classification Weights without Regularization

Here are some word features and their classification weights from the previous model without regularization. Positive weights indicate the word feature contribute to positive sentiment classification and negative weights indicate the opposite contribution

|             | interesting | pleasure | boring | zoe   | write | workings |
|-------------|-------------|----------|--------|-------|-------|----------|
| Without Reg | 0.011       | -5.63    | 1.80   | -5.68 | -8.20 | 14.16    |

# Classification Weights without Regularization

Here are some word features and their classification weights from the previous model without regularization. Positive weights indicate the word feature contribute to positive sentiment classification and negative weights indicate the opposite contribution

|             | interesting | pleasure | boring | zoe   | write  | workings |
| ----------- | ----------- | -------- | ------ | ----- | ------ | -------- |
| Without Reg | 0.011       | -5.63    | 1.80   | -5.68 | -8.20  | 14.16    |

▶ NEGATIVE: woody allen can `write` and deliver a one liner as
  well as anybody .

# Classification Weights without Regularization

Here are some word features and their classification weights from the previous model without regularization. Positive weights indicate the word feature contribute to positive sentiment classification and negative weights indicate the opposite contribution

|             | interesting | pleasure | boring | zoe   | write | workings |
|-------------|-------------|----------|--------|-------|-------|----------|
| Without Reg | 0.011       | -5.63    | 1.80   | -5.68 | -8.20 | 14.16    |

- ▶ NEGATIVE: woody allen can `write` and deliver a one liner as well as anybody .
- ▶ POSITIVE: soderbergh , like kubrick before him , may not touch the planet 's skin , but understands the `workings` of its spirit .

# Learning with Regularization

We chose $\lambda = \frac{1}{C} = 10^2$

- Training accuracy: 62.54%
- Val accuracy: 63.17%

# Classification Weights with Regularization

With regularization, the classification weights make more sense to us

|              | interesting | pleasure | boring | zoe     | write   | workings |
|--------------|-------------|----------|--------|---------|---------|----------|
| Without Reg  | 0.011       | -5.63    | 1.80   | -5.68   | -8.20   | 14.16    |
| With Reg     | 0.16        | 0.36     | -0.21  | -0.057  | -0.066  | 0.040    |

# Classification Weights with Regularization

With regularization, the classification weights make more sense to us

|             | interesting | pleasure | boring | zoe    | write  | workings |
|-------------|-------------|----------|--------|--------|--------|----------|
| Without Reg | 0.011       | -5.63    | 1.80   | -5.68  | -8.20  | 14.16    |
| With Reg    | 0.16        | 0.36     | -0.21  | -0.057 | -0.066 | 0.040    |

## Regularization for Avoiding Overfitting

Reduce the correlation between class label and some noisy features.

# Demo Code

What we are going to review from this demo code

- ▶ NLP
  - ▶ Bag-of-words representations
  - ▶ Text classifiers
- ▶ Machine Learning
  - ▶ Overfitting
  - ▶ $L_2$ regularization

# Reference

Jurafsky, D. and Martin, J. (2019).
Speech and language processing.

Pang, B., Lee, L., and Vaithyanathan, S. (2002).
Thumbs up?: sentiment classification using machine learning techniques.
In *Proceedings of the ACL-02 conference on Empirical methods in natural language processing-Volume 10*, pages 79–86. Association for Computational Linguistics.