

CS 6501 Natural Language Processing

Logistic Regression

Yangfeng Ji

September 12, 2019

Department of Computer Science
University of Virginia



ENGINEERING

Overview

1. Naive Bayes Classifier, revisited
2. Logistic Regression
3. Classification Evaluation
4. Example
5. Building Better LR Models
6. Discriminative vs. Generative Classifiers

Naive Bayes Classifier, revisited

An Alternative View

- ▶ Prior probability: $p(y; \theta) = \theta_y$
- ▶ Likelihood: $p(x | y; \theta) \propto \prod_{i=1}^V \theta_{i,y}^{x_i}$

$$p(x, y; \theta) \propto \theta_y \cdot \prod_{i=1}^V \theta_{i,y}^{x_i} \quad (1)$$

Or

$$\log p(x, y; \theta) \propto \log \theta_y + \sum_{i=1}^V (x_i \cdot \log \theta_{i,y}) \quad (2)$$

An Alternative View (II)

$$\log p(x, y; \theta) \propto \log \theta_y + \sum_{i=1}^V (x_i \cdot \log \theta_{i,y}) \quad (3)$$

$$= b_y + \mathbf{w}_y^\top \mathbf{x} \quad (4)$$

where

$$b_y = \log \theta_y$$

$$\mathbf{w}_y^\top = [\log \theta_{1,y}, \dots, \log \theta_{V,y}]$$

$$\mathbf{x}^\top = [x_1, \dots, x_V]$$

It's a linear classifier (with respect to \mathbf{x})

Logistic Regression

Log-linear Models

Directly modeling a linear classifier as

$$\log p(y \mid \mathbf{x}) \propto \mathbf{w}_y^\top \mathbf{x} + b_y \quad (5)$$

with

- ▶ $\mathbf{x} \in \mathbb{N}^V$: bag-of-words representation
- ▶ $\mathbf{w}_y \in \mathbb{R}^V$: classification weights associated with label y
- ▶ $b_y \in \mathbb{R}$: classification bias associated with label y

Probabilistic Form

Given

$$\log P(y \mid \mathbf{x}) \propto \mathbf{w}_y^\top \mathbf{x} + b_y, \quad (6)$$

the probabilistic form is

$$P(y \mid \mathbf{x}) \propto \exp(\mathbf{w}_y^\top \mathbf{x} + b_y) \quad (7)$$

or

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{\exp(\mathbf{w}_{\mathbf{y}}^\top \mathbf{x} + b_{\mathbf{y}})}{\sum_{y' \in \mathcal{Y}} \exp(\mathbf{w}_{y'}^\top \mathbf{x} + b_{y'})} \quad (8)$$

Alternative Form

Rewrite \mathbf{x} and \mathbf{w} as

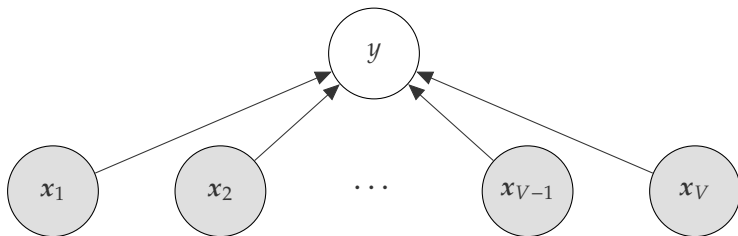
- ▶ $\mathbf{x}^\top = [x_1, x_2, \dots, x_V, \mathbf{1}]$
- ▶ $\mathbf{w}^\top = [w_1, w_2, \dots, w_V, \mathbf{b}_y]$

then,

$$P(y \mid \mathbf{x}) = \frac{\exp(\mathbf{w}_y^\top \mathbf{x})}{\sum_{y' \in \mathcal{Y}} \exp(\mathbf{w}_{y'}^\top \mathbf{x})} \quad (9)$$

Probabilistic Graphical Models

$$P(y \mid \mathbf{x}) = \frac{\exp(\mathbf{w}_y^\top \mathbf{x})}{\sum_{y' \in \mathcal{Y}} \exp(\mathbf{w}_{y'}^\top \mathbf{x})} \quad (10)$$



Binary Classifier

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

$$P(\text{POS} \mid \mathbf{x}) = \frac{\exp(\mathbf{w}_{\text{POS}}^\top \mathbf{x})}{\exp(\mathbf{w}_{\text{POS}}^\top \mathbf{x}) + \exp(\mathbf{w}_{\text{NEG}}^\top \mathbf{x})} \quad (11)$$

$$= \frac{1}{1 + \frac{\exp(\mathbf{w}_{\text{NEG}}^\top \mathbf{x})}{\exp(\mathbf{w}_{\text{POS}}^\top \mathbf{x})}} \quad (12)$$

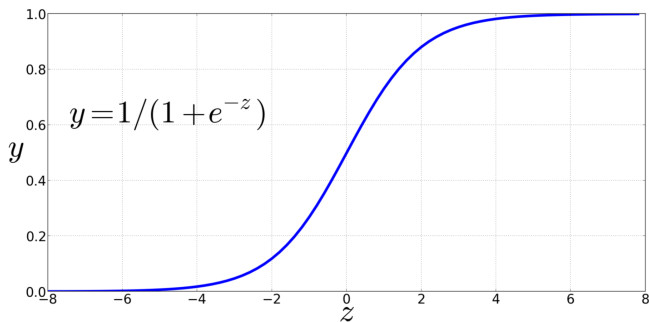
$$= \frac{1}{1 + \exp(-(\mathbf{w}_{\text{POS}} - \mathbf{w}_{\text{NEG}})^\top \mathbf{x})} \quad (13)$$

$$= \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \quad (14)$$

where $\mathbf{w} = \mathbf{w}_{\text{POS}} - \mathbf{w}_{\text{NEG}}$

Sigmoid Function

$$y = \frac{1}{1 + \exp(-z)} \quad (15)$$



$$z \in (-\infty, \infty), y \in (-1, 1)$$

Two Questions

... of using a logistic regression classifier

$$P(y \mid x) = \frac{\exp(w_y^\top x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top x)} \quad (16)$$

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

(Log)-likelihood Function

With a collection of training examples $\{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$, the likelihood function of $\{\mathbf{w}_y\}_{y \in \mathcal{Y}}$ is

$$L(\{\mathbf{w}_y\}) = \prod_{n=1}^N P(y^{(n)} \mid \mathbf{x}^{(n)}) \quad (17)$$

and the **log**-likelihood function is

$$\ell(\{\mathbf{w}_y\}) = \sum_{n=1}^N \log P(y^{(n)} \mid \mathbf{x}^{(n)}) \quad (18)$$

Log-likelihood Function (II)

With

$$P(y \mid \mathbf{x}) = \frac{\exp(\mathbf{w}_y^\top \mathbf{x})}{\sum_{y' \in \mathcal{Y}} \exp(\mathbf{w}_{y'}^\top \mathbf{x})} \quad (19)$$

the log-likelihood function is

$$\ell(\{\mathbf{w}_y\}) = \sum_{n=1}^N \log P(y^{(n)} \mid \mathbf{x}^{(n)}) \quad (20)$$

$$= \sum_{n=1}^N \left\{ \mathbf{w}_{y^{(n)}}^\top \mathbf{x}^{(n)} - \log \sum_{y' \in \mathcal{Y}} \exp(\mathbf{w}_{y'}^\top \mathbf{x}^{(n)}) \right\} \quad (21)$$

Given the training examples $\{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$, $\ell(\{\mathbf{w}_y\})$ is a function of $\{\mathbf{w}_y\}$.

Optimization with Gradient

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

$$\begin{aligned}\text{NLL}(\{w_y\}) &= -\ell(\{w_y\}) \\ &= \sum_{n=1}^N \left\{ -w_{y^{(n)}}^\top x^{(n)} + \log \sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top x) \right\}\end{aligned}$$

then

$$w_y \leftarrow w_y - \eta \cdot \frac{\partial \text{NLL}(\{w_y\})}{\partial w_y}, \quad \forall y \in \mathcal{Y} \quad (22)$$

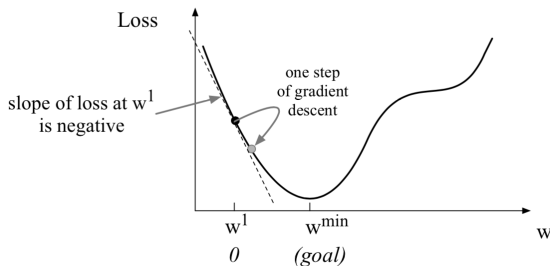
where η is called **learning rate**.

Optimization with Gradient (II)

Two questions answered by the update equation

- ▶ which direction?
- ▶ how far it should go?

$$w_y \leftarrow w_y - \eta \cdot \frac{\partial \text{NLL}(\{w_y\})}{\partial w_y} \quad (23)$$



Training Procedure

Steps for parameter estimation, given the current parameter $\{w_y\}$

1. Compute the derivative

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

2. Update parameters with

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

3. If not **done**, rerun to step 1

Classification Evaluation

A Development Set

- ▶ Training set $\mathcal{T} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{i=1}^N$
- ▶ Development set $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^M$
- ▶ Test set $\mathcal{U} = \{(\mathbf{x}^{(j)}, y^{(j)})\}_{j=1}^L$

Cross-validation

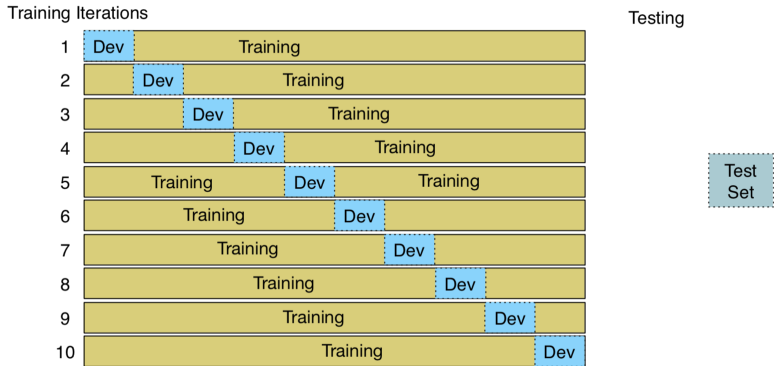


Figure: 10-fold cross validation.

Evaluation Measurements

- ▶ Accuracy
- ▶ Precision, recall, and F-measure

[Eisenstein, 2018, Sec 4.4]

Accuracy

Given M examples, with $y^{(i)}$ is the ground-truth label of the i -th example, and $\hat{y}^{(i)}$ is the predicted label

$$\text{ACC}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{M} \sum_{i=1}^M \delta(y^{(i)}, \hat{y}^{(i)}) \quad (24)$$

δ function is defined as

$$\delta(y, \hat{y}) = \begin{cases} 1 & y = \hat{y} \\ 0 & y \neq \hat{y} \end{cases} \quad (25)$$

Confusion Matrix

		Ground truth	
		POSITIVE	NEGATIVE
Prediction	POSITIVE	True Positive (TP)	False Positive (FP)
	NEGATIVE	False Negative (FN)	True Negative (TN)

Accuracy:

$$\text{acc}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (26)$$

Recall, Precision and F Measure

		Ground truth	
		POSITIVE	NEGATIVE
Prediction	POSITIVE	True Positive (TP)	False Positive (FP)
	NEGATIVE	False Negative (FN)	True Negative (TN)

Recall:

$$r(\mathbf{y}, \hat{\mathbf{y}}) = \frac{TP}{TP + FN} \quad (27)$$

Precision:

$$p(\mathbf{y}, \hat{\mathbf{y}}) = \frac{TP}{TP + FP} \quad (28)$$

F measure:

$$F(\mathbf{y}, \hat{\mathbf{y}}) = \frac{2 \cdot p \cdot r}{p + r} \quad (29)$$

Example: Balanced case

		Ground truth	
		POSITIVE	NEGATIVE
Prediction	POSITIVE	480	30
	NEGATIVE	20	470


- ▶ Accuracy: $\text{acc} = \frac{480+470}{480+20+30+470} = 0.95$
- ▶ Precision: $p = \frac{480}{480+30} \approx 0.94$
- ▶ Recall: $r = \frac{480}{480+20} \approx 0.96$
- ▶ F-measure: $F = \frac{2 \times 0.94 \times 0.96}{0.94 + 0.96} \approx 0.95$

Example: Unbalanced case

		Ground truth	
		POSITIVE	NEGATIVE
Prediction	POSITIVE	80	30
	NEGATIVE	20	870

- ▶ Accuracy: $\text{acc} = \frac{80+870}{80+20+30+870} = 0.95$
- ▶ Precision: $p = \frac{80}{80+30} \approx 0.73$
- ▶ Recall: $r = \frac{80}{80+20} = 0.80$
- ▶ F-measure: $F = \frac{2 \times 0.94 \times 0.96}{0.94 + 0.96} \approx 0.76$

Example

A horizontal line spanning the width of the slide, divided into two equal halves. The left half is a solid blue line, and the right half is a solid white line.

Example Dataset

A subset of the Yelp Dataset

<https://www.yelp.com/dataset/challenge>

	Training	Development	Test
Documents	40K	5K	5K
Words	4.7M	0.5M	0.6M

- ▶ 5 classes (user rating from 1 to 5)
- ▶ Code available on <https://github.com/jiyfeng/textclassification>

Building BoW Representations

Sklearn function

- ▶ `sklearn.feature_extraction.text.CountVectorizer`
 - ▶ `sklearn.linear_model.LogisticRegression`
-
- ▶ Given a collection of texts, it will build a vocab and also convert all texts into numeric vectors
 - ▶ Classification accuracy on the development data is 61.4%

How Far We Can Go with BoW?

Tricks to reduce vocab size

- ▶ remove punctuation (default)
- ▶ lowercase (↗)
- ▶ remove low-frequency words (↗)
- ▶ remove high-frequency words (↘)
- ▶ replace numbers with a special token

Comments

- ▶ Not always helpful (these are empirical tricks)
- ▶ Not always the case (it depends on the data/domain)

Interpretability

Weights learned from training data

Vocab	$w_{\text{rating}=1}$	$w_{\text{rating}=5}$
$\left(\begin{array}{c} \text{SUPER} \\ \dots \\ \text{QUICK} \\ \text{FOOD} \\ \text{FRIENDLY} \\ \text{EAT} \\ \dots \\ \text{DELICIOUS} \end{array} \right)$	$\left[\begin{array}{c} 0.33 \\ \dots \\ -1.26 \\ 0.08 \\ -2.57 \\ -0.47 \\ \dots \\ -3.60 \end{array} \right]$	$\left[\begin{array}{c} -0.09 \\ \dots \\ -0.01 \\ -0.09 \\ 0.16 \\ 0.00 \\ \dots \\ 0.64 \end{array} \right]$

Interpretability (II)

Top features

rating = 5	rating = 1
exceptional	worst
incredible	joke
phenomenal	disgusted
body	unprofessional
<i>regret</i>	garbage
worried	disgusting
skeptical	<i>luck</i>
hesitate	pathetic
happier	apologies
mike	horrible

Building Better LR Models

Ways to Improve LR Models

- ▶ A better text representation x
- ▶ Regularization on w

A Generic LR Model

Replace x in the following equation

$$P(y \mid x) = \frac{\exp(w_y^\top x)}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top x)} \quad (30)$$

with a **generic** feature function $f(x)$

$$P(y \mid x) = \frac{\exp(w_y^\top f(x))}{\sum_{y' \in \mathcal{Y}} \exp(w_{y'}^\top f(x))} \quad (31)$$

Rich Features: bi-grams

Combine words together to form high-order features

- ▶ Uni-grams:

$\{\text{FAST, SLOW, SUPER}\}$

- ▶ Bi-grams:

$\{\text{SUPER FAST, SUPER SLOW}\}$

Extended feature set

$\{\text{FAST, SLOW, SUPER, SUPER FAST, SUPER SLOW}\}$

Other Features

For example: other rich features from $f(x)$ for sentiment analysis

- ▶ the length of the text
- ▶ the number of sentences
- ▶ the number of positive words
- ▶ the number of negative words
- ▶ the number of pronouns

Combining these features with the uni- and bi-gram features

Sklearn function

```
sklearn.feature_extraction.text.CountVectorizer  
with ngram_range=(1, 2)
```

- ▶ Even larger vocab size: from 60K to 1M
- ▶ Performance change: from 61.4% to 62.4% on the dev data

Problems with Large Feature Set

- ▶ Overlap among features

Example

{FAST, SLOW, SUPER, SUPER FAST, SUPER SLOW}

- ▶ Learning in high dimensional space — overfitting

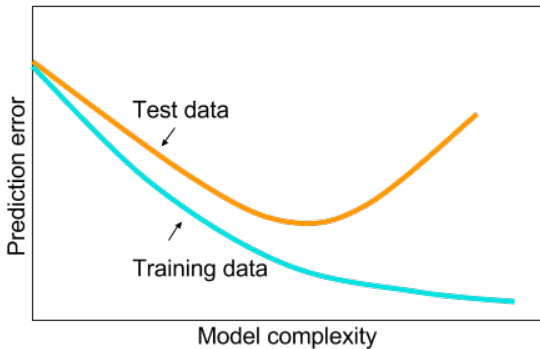
Example

62.4% on dev data

vs.

~ 100% on training data

Overfitting



[Abu-Mostafa et al., 2012]

Ways to Improve LR Models

- ✓ A better text representation x
- ▶ Regularization on w

L_2 Regularization

Add an additional constraint on $\{\mathbf{w}_y\}$

$$\ell_{L_2}(\boldsymbol{\theta}) = - \sum_{n=1}^N \log P(y^{(n)} | \mathbf{x}^{(n)}; \mathbf{w}) + \frac{\lambda}{2} \sum_y \|\mathbf{w}_y\|_2^2 \quad (32)$$

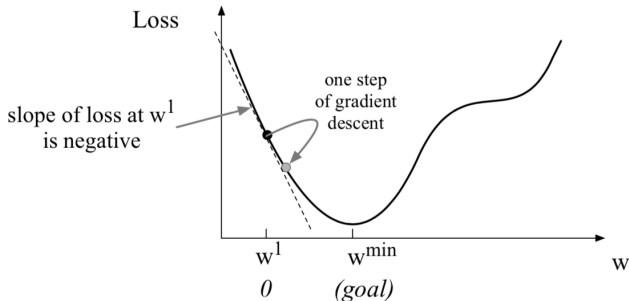
The gradient of the new loss function

$$\frac{\partial \ell_{L_2}(\{\mathbf{w}_y\})}{\partial \mathbf{w}_y} = - \sum_{n=1}^N \frac{\partial}{\partial \mathbf{w}_y} \log P(y^{(n)} | \mathbf{x}^{(n)}; \mathbf{w}) + \lambda \mathbf{w}_y \quad (33)$$

[Eisenstein, 2018, Sec. 2.4.1]

L_2 Regularization (Cont.)

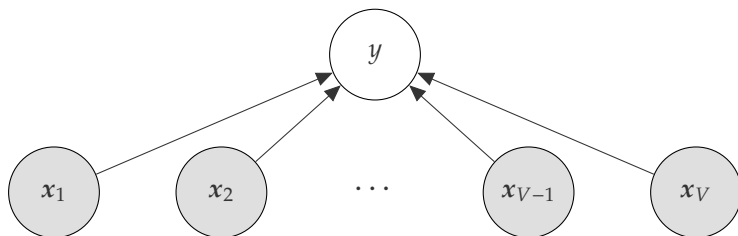
L_2 Regularization introduces a trade-off between the likelihood function and the norm of $\{w_y\}$



Discriminative vs. Generative Classifiers

Discriminative Models

Logistic models

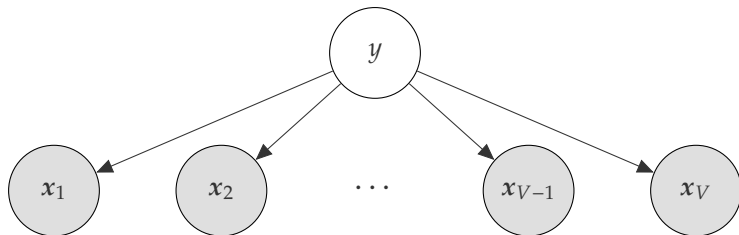


- ▶ No conditional independence assumption
 - ▶ Rich feature set
- ▶ Performance depends the choice of optimization methods

[Jurafsky and Martin, 2019]

Generative Models

Naive Bayes classifiers



- ▶ Conditional independence
- ▶ Easy to implement
- ▶ Works better on *small* training sets or *short* texts

[Jurafsky and Martin, 2019]

Reference



Abu-Mostafa, Y. S., Magdon-Ismail, M., and Lin, H.-T. (2012).
Learning from data, volume 4.
AMLBook New York, NY, USA:.



Eisenstein, J. (2018).
Natural Language Processing.
MIT Press.



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