

An Introduction To Natural Language Processing

Chapter 2: Linear Text Classification

Jacob Eisenstein

Roadmap for this chapter

- ▶ Basic representations of text data for classification
- ▶ Four linear classifiers
 - ▶ Naïve Bayes
 - ▶ Perceptron
 - ▶ Large-margin (support vector machine)
 - ▶ Logistic regression

The text classification problem

Given a text $\mathbf{w} = (w_1, w_2, \dots, w_T) \in \mathcal{V}^*$
choose a label $y \in \mathcal{Y}$.

The text classification problem

Given a text $\mathbf{w} = (w_1, w_2, \dots, w_T) \in \mathcal{V}^*$
choose a label $y \in \mathcal{Y}$.

- Some direct applications:

Sentiment analysis	$\mathcal{Y} = \{\text{positive, negative, neutral}\}$
Spam filtering	$\mathcal{Y} = \{\text{spam, not-spam}\}$
Language identification	$\mathcal{Y} = \{\text{Mandarin, English, Spanish, } \dots\}$

- Indirectly, methods from text classification apply to a huge range of settings in natural language processing, and will appear again and again throughout the course.

The bag-of-words

- ▶ One challenge is that the sequential representation (w_1, w_2, \dots, w_T) may have a different length T for every document.
- ▶ The bag-of-words is a fixed-length representation, which consists of a vector of word counts:

$$\begin{aligned} \mathbf{w} &= \text{It was the best of times, it was the worst of times} \\ \mathbf{x} &= [\underbrace{\text{aardvark}}_0, \dots, \underbrace{\text{best}}_1, \dots, \underbrace{\text{it}}_2, \dots, \underbrace{\text{of}}_2, \dots, \underbrace{\text{zyther}}_0] \end{aligned}$$

The length of \mathbf{x} is equal to the size of the vocabulary, V .

- ▶ For each \mathbf{x} , there may be many possible \mathbf{w} , depending on word order.

Linear classification on the bag-of-words

Let $\psi(\mathbf{x}, y)$ score the compatibility of bag-of-words \mathbf{x} and label y .
Then,

$$\hat{y} = \operatorname{argmax}_y \psi(\mathbf{x}, y).$$

In a **linear classifier**, this scoring function has a simple form,

$$\psi(\mathbf{x}, y) = \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y) = \sum_{j=1} \theta_j \times f_j(\mathbf{x}, y),$$

where $\boldsymbol{\theta}$ is a vector of weights, and \mathbf{f} is a **feature function**.

Linear classification in python

```
def compute_score(x,y,weights):  
    total = 0  
    for feature,count in feature_function(x,y).items():  
        total += weights[feature] * count  
    return total
```

Feature functions

In classification, the feature function is usually a simple combination of \mathbf{x} and y , such as,

$$f_j(\mathbf{x}, y) = \begin{cases} x_{\text{whale}}, & \text{if } y = \text{FICTION}, \\ 0, & \text{otherwise.} \end{cases}$$

Feature functions

In classification, the feature function is usually a simple combination of \mathbf{x} and y , such as,

$$f_j(\mathbf{x}, y) = \begin{cases} x_{\text{whale}}, & \text{if } y = \text{FICTION}, \\ 0, & \text{otherwise.} \end{cases}$$

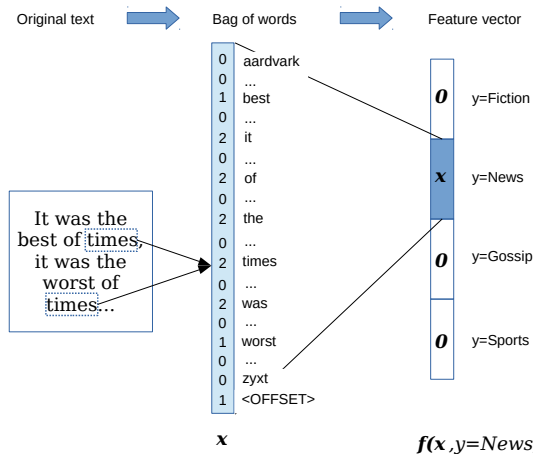
This corresponds to the following formalization as a column vector,

$$\mathbf{f}(\mathbf{x}, y = 1) = [\mathbf{x}; \underbrace{0; 0; \dots; 0}_{(K-1) \times V}]$$

$$\mathbf{f}(\mathbf{x}, y = 2) = [\underbrace{0; 0; \dots; 0}_V; \mathbf{x}; \underbrace{0; 0; \dots; 0}_{(K-2) \times V}]$$

$$\mathbf{f}(\mathbf{x}, y = K) = [\underbrace{0; 0; \dots; 0}_{(K-1) \times V}; \mathbf{x}],$$

Another view of the bag-of-words and feature function



Summary and next steps

- ▶ To summarize, our classification function is,

$$\hat{y} = \operatorname{argmax}_y \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y),$$

where \mathbf{f} is a feature function and \mathbf{x} is the bag-of-words representation.

- ▶ The **learning** problem is to find the right weights $\boldsymbol{\theta}$. The remainder of this chapter describes four learning algorithms:
 - ▶ Naïve Bayes
 - ▶ Perceptron
 - ▶ Large-margin (support vector machine)
 - ▶ Logistic regression

All these methods assume a labeled dataset $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$.

Probabilistic classification

Naïve Bayes is a probabilistic classifier. It takes the following strategy:

- ▶ Define a **probability model**, $p(\mathbf{x}, y)$.
- ▶ Estimate the parameters of the probability model by **maximum likelihood** — that is, by maximizing the likelihood of the dataset.
- ▶ Set the scoring function equal to the log-probability,

$$\psi(\mathbf{x}, y) = \log p(\mathbf{x}, y) = \log p(y \mid \mathbf{x}) + C,$$

where C is constant in y . This ensures that,

$$\hat{y} = \operatorname{argmax}_y p(y \mid \mathbf{x}).$$

A probability model for text classification

- ▶ First, assume each instance is independent of the others,

$$p(\mathbf{x}^{(1:N)}, \mathbf{y}^{(1:N)}) = \prod_{i=1}^N p(\mathbf{x}^{(i)}, y^{(i)}).$$

- ▶ Apply the chain rule of probability,

$$p(\mathbf{x}, y) = p(\mathbf{x} \mid y) \times p(y).$$

- ▶ Define the parametric form of each probability,

$$p(y) = \text{Categorical}(\boldsymbol{\mu}) \quad p(\mathbf{x} \mid y) = \text{Multinomial}(\boldsymbol{\phi}, T).$$

- ▶ The **multinomial** is a distribution over vectors of counts.
- ▶ The **parameters** $\boldsymbol{\mu}$ and $\boldsymbol{\phi}$ are vectors of probabilities.

The multinomial distribution

- ▶ Suppose the word *whale* has probability ϕ_j .
What is the probability that this word appears 3 times?
- ▶ Each word's probability is exponentiated by its count,

$$\text{Multinomial}(\mathbf{x}; \phi, T) = \prod_{j=1}^V \phi_j^{x_j}.$$

The multinomial distribution

- ▶ Suppose the word *whale* has probability ϕ_j .
What is the probability that this word appears 3 times?
- ▶ Each word's probability is exponentiated by its count,

$$\text{Multinomial}(\mathbf{x}; \phi, T) = \frac{\left(\sum_{j=1}^V x_j\right)!}{\prod_{j=1}^V (x_j!)} \prod_{j=1}^V \phi_j^{x_j}.$$

- ▶ The coefficient is the count of the number of possible orderings of \mathbf{x} . Crucially, it does not depend on the frequency parameter ϕ .

Naïve Bayes text classification

Naïve Bayes can be reformulated within our linear classification framework by setting θ equal to the log parameters,

$$\begin{aligned}\theta &= [\theta^{(1)}; \theta^{(2)}; \dots; \theta^{(K)}] \\ \theta^{(y)} &= [\log \phi_{y,1}; \log \phi_{y,2}; \dots; \log \phi_{y,v}; \log \mu_y] \\ \psi(\mathbf{x}, y) &= \theta \cdot \mathbf{f}(\mathbf{x}, y) = \log p(\mathbf{x} \mid y) + \log p(y),\end{aligned}$$

where $\mathbf{f}(\mathbf{x}, y)$ is extended to include an “offset” 1 for each possible label after the word counts.

Estimating Naïve Bayes

In relative frequency estimation, the parameters are set to empirical frequencies:

$$\hat{\phi}_{y,j} = \frac{\text{count}(y,j)}{\sum_{j'=1}^V \text{count}(y,j')} = \frac{\sum_{i:y^{(i)}=y} x_j^{(i)}}{\sum_{j'=1}^V \sum_{i:y^{(i)}=y} x_{j'}^{(i)}}$$

$$\hat{\mu}_y = \frac{\text{count}(y)}{\sum_{y'} \text{count}(y')}.$$

As shown in the book, this turns out to be identical to the **maximum likelihood estimate**,

$$\hat{\phi}, \hat{\mu} = \underset{\phi, \mu}{\operatorname{argmax}} \prod_{i=1}^N p(\mathbf{x}^{(i)}, y^{(i)}) = \underset{\phi, \mu}{\operatorname{argmax}} \sum_{i=1}^N \log p(\mathbf{x}^{(i)}, y^{(i)}). \quad (1)$$

Smoothing, bias, variance

To deal with low counts, it can be helpful to smooth probabilities,

$$\hat{\phi}_{y,j} = \frac{\alpha + \text{count}(y,j)}{V\alpha + \sum_{j'=1}^V \text{count}(y,j')}.$$

- ▶ Smoothing introduces **bias**, moving the parameters away from their maximum-likelihood estimates.
- ▶ But it corrects **variance**, which is the extent to which the parameters depend on idiosyncrasies of a finite dataset.
- ▶ The smoothing term α is a **hyperparameter**, which must be tuned on a **development set**.

Too naïve?

Naïve Bayes is so called because:

- ▶ Bayes rule is used to convert the observation probability $p(\mathbf{x} \mid y)$ into the label probability $p(y \mid \mathbf{x})$.
- ▶ The multinomial distribution naïvely ignores dependencies between words, and treats every word as equally informative.
- ▶ Suppose *naïve* and *Bayes* always occur together. Should we really count them both for classification?
- ▶ **Discriminative** classifiers avoid this problem by not attempting to model the “generative” probability $p(\mathbf{x})$.

The perceptron classifier

A simple learning rule:

- ▶ Run the current classifier on an instance in the training data, obtaining $\hat{y} = \operatorname{argmax}_y \psi(\mathbf{x}^{(i)}, y)$.
- ▶ If the prediction is incorrect:
 - ▶ Increase the weights for the features of the true label $y^{(i)}$.
 - ▶ Decrease the weights for the features of the predicted label \hat{y} .

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y}).$$

- ▶ Repeat until all training instances are correctly classified, or you run out of time.

The perceptron classifier

A simple learning rule:

- ▶ Run the current classifier on an instance in the training data, obtaining $\hat{y} = \operatorname{argmax}_y \psi(\mathbf{x}^{(i)}, y)$.
- ▶ If the prediction is incorrect:
 - ▶ Increase the weights for the features of the true label $y^{(i)}$.
 - ▶ Decrease the weights for the features of the predicted label \hat{y} .

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y}).$$

- ▶ Repeat until all training instances are correctly classified, or you run out of time.

If the dataset is **linearly separable** — if there is some $\boldsymbol{\theta}$ that correctly labels all the training instances — then this method is guaranteed to find it.

Loss functions

Many classifiers can be viewed as minimizing a **loss function** on the weights. Such a function should have two properties:

- ▶ It should be a good proxy for the accuracy of the classifier.
- ▶ It should be easy to optimize.

Do you see why $1 - \text{accuracy}$ is not a good loss function?

Perceptron as gradient descent

The perceptron can be viewed as optimizing the loss function:

$$\ell_{\text{perceptron}}(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) = -\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + \max_{y' \neq y^{(i)}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y')$$

Perceptron as gradient descent

The perceptron can be viewed as optimizing the loss function:

$$\ell_{\text{perceptron}}(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) = -\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + \max_{y' \neq y^{(i)}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y')$$

The gradient of the perceptron loss is part of the perceptron update:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\theta}} \ell_{\text{perceptron}} &= -\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + \mathbf{f}(\mathbf{x}^{(i)}, \hat{y}) \\ \boldsymbol{\theta}^{(t+1)} &\leftarrow \boldsymbol{\theta}^{(t)} - \frac{\partial}{\partial \boldsymbol{\theta}} \ell_{\text{perceptron}} \quad \text{(gradient descent!)} \\ &= \boldsymbol{\theta}^{(t)} + \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y}). \end{aligned}$$

Large margin learning

For better generalization, the correct label should outscore all other labels by a large **margin**:

$$\gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) = \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \max_{y' \neq y^{(i)}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y') \quad (2)$$

Large margin learning

For better generalization, the correct label should outscore all other labels by a large **margin**:

$$\gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) = \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \max_{y' \neq y^{(i)}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y') \quad (2)$$

The margin can be incorporated into a **margin loss**:

$$\begin{aligned} \ell_{\text{margin}} &= \max \left(0, 1 - \gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) \right) \\ &= -\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + \max_{\hat{y} \in \mathcal{Y}} c(y^{(i)}, \hat{y}) + \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, \hat{y}). \end{aligned}$$

Minimizing the margin loss yields an update that is very similar to the perceptron.

Large margin learning

For better generalization, the correct label should outscore all other labels by a large **margin**:

$$\gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) = \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \max_{y' \neq y^{(i)}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y') \quad (2)$$

The margin can be incorporated into a **margin loss**:

$$\begin{aligned} \ell_{\text{margin}} &= \max \left(0, 1 - \gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) \right) \\ &= -\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + \max_{\hat{y} \in \mathcal{Y}} c(y^{(i)}, \hat{y}) + \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, \hat{y}). \end{aligned}$$

Minimizing the margin loss yields an update that is very similar to the perceptron.

Large margin learning

For better generalization, the correct label should outscore all other labels by a large **margin**:

$$\gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) = \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \max_{y' \neq y^{(i)}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y') \quad (2)$$

The margin can be incorporated into a **margin loss**:

$$\begin{aligned} \ell_{\text{margin}} &= \max \left(0, 1 - \gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) \right) \\ &= -\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + \max_{\hat{y} \in \mathcal{Y}} c(y^{(i)}, \hat{y}) + \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, \hat{y}). \end{aligned}$$

Minimizing the margin loss yields an update that is very similar to the perceptron.

Logistic regression

- ▶ Perceptron and large margin classification are **discriminative**: they learn to discriminate correct and incorrect labels.
- ▶ Naïve Bayes is **probabilistic**: it assigns calibrated confidence scores to its predictions.
- ▶ Logistic regression is both discriminative and probabilistic. It directly computes the conditional probability of the label:

$$p(y \mid \mathbf{x}; \boldsymbol{\theta}) = \frac{\exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y))}{\sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y'))}$$

- ▶ Exponentiation ensures that the probabilities are non-negative.
- ▶ Normalization ensures that the probabilities sum to one.

Logistic regression

- ▶ Perceptron and large margin classification are **discriminative**: they learn to discriminate correct and incorrect labels.
- ▶ Naïve Bayes is **probabilistic**: it assigns calibrated confidence scores to its predictions.
- ▶ Logistic regression is both discriminative and probabilistic. It directly computes the conditional probability of the label:

$$p(y \mid \mathbf{x}; \boldsymbol{\theta}) = \frac{\exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y))}{\sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y'))}$$

- ▶ Exponentiation ensures that the probabilities are non-negative.
- ▶ Normalization ensures that the probabilities sum to one.

Logistic regression

- ▶ Perceptron and large margin classification are **discriminative**: they learn to discriminate correct and incorrect labels.
- ▶ Naïve Bayes is **probabilistic**: it assigns calibrated confidence scores to its predictions.
- ▶ Logistic regression is both discriminative and probabilistic. It directly computes the conditional probability of the label:

$$p(y \mid \mathbf{x}; \boldsymbol{\theta}) = \frac{\exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y))}{\sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y'))}$$

- ▶ Exponentiation ensures that the probabilities are non-negative.
- ▶ **Normalization ensures that the probabilities sum to one.**

Learning logistic regression

Two equivalent views of logistic regression learning:

- **Maximization** of the conditional log-likelihood,

$$\log p(\mathbf{y}^{(1:N)} \mid \mathbf{x}^{(1:N)}; \boldsymbol{\theta}) = \sum_{i=1}^N \log p(y^{(i)} \mid \mathbf{x}^{(i)}; \boldsymbol{\theta})$$

- **Minimization** of the logistic loss,

$$\ell_{\text{LOGREG}}(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) = -\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + \log \sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y')).$$

(Compare with the perceptron loss!)

High-dimensional classification

Possible problems:

- ▶ What if the number of features in $\mathbf{f}(\mathbf{x}, y)$ is larger than the number of training instances?
- ▶ What happens in logistic regression if a feature appears only with one label? What will its weight be?

These problems relate to the **variance** of a classifier — its sensitivity to idiosyncratic features of the training data.

Regularization

Learning can often be made more robust by **regularization**: penalizing large weights.

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell_{\text{LOGREG}}(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) + \lambda \|\boldsymbol{\theta}\|_2^2, \quad (3)$$

where

- ▶ $\|\boldsymbol{\theta}\|_2^2 = \sum_j \theta_j^2$.
- ▶ The scalar λ controls the strength of regularization.
- ▶ The **support vector machine** classifier combines regularization with the margin loss.

Regularization

Learning can often be made more robust by **regularization**: penalizing large weights.

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell_{\text{LOGREG}}(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) + \lambda \|\boldsymbol{\theta}\|_2^2, \quad (3)$$

where

- ▶ $\|\boldsymbol{\theta}\|_2^2 = \sum_j \theta_j^2$.
- ▶ The scalar λ controls the strength of regularization.
- ▶ The **support vector machine** classifier combines regularization with the margin loss.

Gradient descent

Logistic regression, perceptron, and large-margin classification all learn by minimizing a **loss function**. A general strategy for minimization is **gradient descent**:

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \eta \frac{\partial}{\partial \boldsymbol{\theta}} \sum_{i=1}^N \ell(\boldsymbol{\theta}^{(t)}; \mathbf{x}^{(i)}, y^{(i)}), \quad (4)$$

where $\eta \in \mathbb{R}_+$ is the **learning rate**.

Stochastic gradient descent

- ▶ Computing the gradient over all instances is expensive.
- ▶ **Stochastic** gradient descent approximates the gradient by its value on a single instance:

$$\frac{\partial}{\partial \theta} \sum_{i=1}^N \ell(\theta^{(t)}; \mathbf{x}^{(i)}, y^{(i)}) \approx C \times \frac{\partial}{\partial \theta} \ell(\theta^{(t)}; \mathbf{x}^{(i)}, y^{(i)}), \quad (5)$$

where $(\mathbf{x}^{(i)}, y^{(i)})$ is sampled at random from the training set.

- ▶ **Minibatch** gradient descent approximates the gradient by its value on small number of instances. This is suited to GPU architectures, and is commonly used in deep learning.

Summary of linear classification

	Pros	Cons
Naïve Bayes	simple, probabilistic, fast	not very accurate
Perceptron	simple, accurate	not probabilistic, may overfit
Large margin	error-driven learning, can be regularized	not probabilistic
Logistic regression	error-driven learning, regularized	more difficult to implement