



University of Colorado **Boulder**

Department of Computer Science

CSCI 5622: Machine Learning

Chenhao Tan

Lecture 7: Logistic regression

Slides adapted from Chris Ketelsen, Jordan Boyd-Graber,
and Noah Smith

Administrivia

- HW2 is out
- Grading issues for HW1

Learning Objectives

- Understand probabilistic classification
- Understand logistic regression
- Understand generative models vs. discriminative models

Outline

- Probabilistic classification
- Logistic regression
- Generative vs. Discriminative models

Outline

- Probabilistic classification
- Logistic regression
- Generative vs. Discriminative models

Recap

K-nearest neighbor

- Find $\mathcal{N}_K(\mathbf{x})$: the set of K training examples nearest to \mathbf{x}
- Predict \hat{y} to be majority label in $\mathcal{N}_K(\mathbf{x})$
- Admits a probabilistic interpretation of class given data: $p(y = c \mid \mathbf{x})$

Recap

K-nearest neighbor

- Find $\mathcal{N}_K(\mathbf{x})$: the set of K training examples nearest to \mathbf{x}
- Predict \hat{y} to be majority label in $\mathcal{N}_K(\mathbf{x})$
- Admits a probabilistic interpretation of class given data: $p(y = c \mid \mathbf{x})$

Perceptron

- Learn weights \mathbf{w} and b via the perceptron algorithm
- Predict \hat{y} via $\hat{y} = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$
- Has no probabilistic interpretation

Probabilistic Models

- hypothesis function $h : X \rightarrow Y$.

Probabilistic Models

- hypothesis function $h : X \rightarrow Y$.
In this special case, we define h based on estimating a probabilistic model $P(X, Y)$.

Probabilistic Classification

Input: $S_{\text{train}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ training examples

$$y_i \in \{c_1, c_2, \dots, c_J\}$$

Goal: $h : X \rightarrow Y$

For each class c_j , estimate

$$P(y = c_j \mid \mathbf{x}, S_{\text{train}})$$

Assign to \mathbf{x} the class with the highest probability

$$\hat{y} = h(\mathbf{x}) = \arg \max_c P(y = c \mid \mathbf{x}, S_{\text{train}})$$

Outline

- Probabilistic classification
- **Logistic regression**
- Generative vs. Discriminative models

What are we talking about?

- Probabilistic classification: $p(y|x)$
- Classification uses: ad placement, spam detection
- Building block of other machine learning methods

Logistic Regression: Definition

- Weight vector β_i
- Observations X_i
- “Bias” β_0 (like intercept in linear regression)

$$P(Y = 0|X) = \frac{1}{1 + \exp [\beta_0 + \sum_i \beta_i X_i]} \quad (1)$$

$$P(Y = 1|X) = \frac{\exp [\beta_0 + \sum_i \beta_i X_i]}{1 + \exp [\beta_0 + \sum_i \beta_i X_i]} \quad (2)$$

Logistic Regression: Definition

- Weight vector β_i
- Observations X_i
- “Bias” β_0 (like intercept in linear regression)

$$P(Y = 0|X) = \frac{1}{1 + \exp [\beta_0 + \sum_i \beta_i X_i]} \quad (1)$$

$$P(Y = 1|X) = \frac{\exp [\beta_0 + \sum_i \beta_i X_i]}{1 + \exp [\beta_0 + \sum_i \beta_i X_i]} \quad (2)$$

$$\beta_0 + \sum_i \beta_i X_i = \log \frac{P(Y = 1|X)}{P(Y = 0|X)}$$

Logistic Regression: Definition

- Weight vector β_i
- Observations X_i
- “Bias” β_0 (like intercept in linear regression)

$$P(Y = 0|X) = \frac{1}{1 + \exp [\beta_0 + \sum_i \beta_i X_i]} \quad (1)$$

$$P(Y = 1|X) = \frac{\exp [\beta_0 + \sum_i \beta_i X_i]}{1 + \exp [\beta_0 + \sum_i \beta_i X_i]} \quad (2)$$

$$\beta_0 + \sum_i \beta_i X_i = \log \frac{P(Y = 1|X)}{P(Y = 0|X)}$$

What is the decision boundary?

Logistic Regression: Definition

- Weight vector β_i
- Observations X_i
- For shorthand, we'll say that

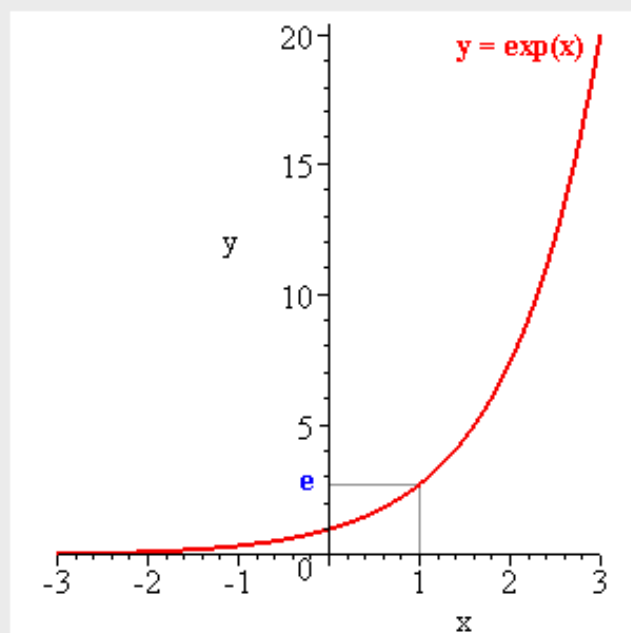
$$P(Y = 1|X) = \sigma((\beta_0 + \sum_i \beta_i X_i)) \quad (3)$$

$$P(Y = 0|X) = 1 - \sigma((\beta_0 + \sum_i \beta_i X_i)) \quad (4)$$

- Where $\sigma(z) = \frac{1}{1+\exp[-z]}$

What's this “exp” doing?

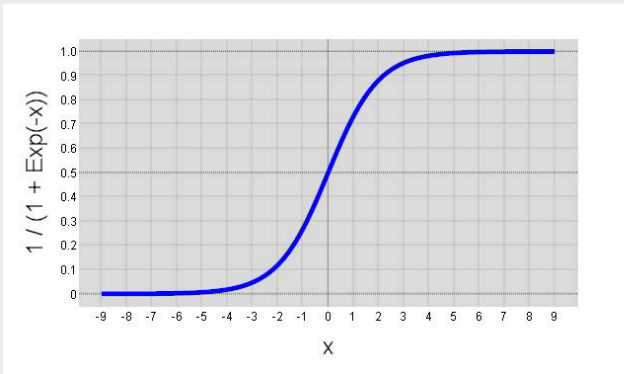
Exponential function



- $\exp[x]$ is shorthand for e^x
- e is a special number, about 2.71828
 - e^x is the limit of compound interest formula as compounds become infinitely small
 - It's the function whose derivative is itself
- The “logistic” function is $\sigma(z) = \frac{1}{1+e^{-z}}$
- Looks like an “S”
- Always between 0 and 1.

What's this “exp” doing?

Logistic function



- $\exp[x]$ is shorthand for e^x
- e is a special number, about 2.71828
 - e^x is the limit of compound interest formula as compounds become infinitely small
 - It's the function whose derivative is itself
- The “logistic” function is $\sigma(z) = \frac{1}{1 + e^{-z}}$
- Looks like an “S”
- Always between 0 and 1.
 - Allows us to model probabilities
 - Different from **linear** regression

Logistic Regression Example

feature	coefficient	weight
bias	β_0	0.1
“viagra”	β_1	2.0
“mother”	β_2	-1.0
“work”	β_3	-0.5
“nigeria”	β_4	3.0

- What does $Y = 1$ mean?

Example 1: Empty Document?

$$X = \{\}$$

Logistic Regression Example

feature	coefficient	weight
bias	β_0	0.1
“viagra”	β_1	2.0
“mother”	β_2	-1.0
“work”	β_3	-0.5
“nigeria”	β_4	3.0

- $Y = 1$: spam

Example 1: Empty Document?

$$X = \{\}$$

- $P(Y = 0) = \frac{1}{1 + \exp[0.1]} =$
- $P(Y = 1) = \frac{\exp[0.1]}{1 + \exp[0.1]} =$

Logistic Regression Example

feature	coefficient	weight
bias	β_0	0.1
“viagra”	β_1	2.0
“mother”	β_2	-1.0
“work”	β_3	-0.5
“nigeria”	β_4	3.0

- $Y = 1$: spam

Example 1: Empty Document?

$$X = \{\}$$

- $P(Y = 0) = \frac{1}{1 + \exp[0.1]} = 0.48$
- $P(Y = 1) = \frac{\exp[0.1]}{1 + \exp[0.1]} = 0.52$
- Bias β_0 encodes the prior probability of a class

Logistic Regression Example

feature	coefficient	weight
bias	β_0	0.1
“viagra”	β_1	2.0
“mother”	β_2	-1.0
“work”	β_3	-0.5
“nigeria”	β_4	3.0

- $Y = 1$: spam

Example 2

$X = \{\text{Mother, Nigeria}\}$

Logistic Regression Example

feature	coefficient	weight
bias	β_0	0.1
“viagra”	β_1	2.0
“mother”	β_2	-1.0
“work”	β_3	-0.5
“nigeria”	β_4	3.0

- $Y = 1$: spam

Example 2

$X = \{\text{Mother, Nigeria}\}$

- $P(Y = 0) = \frac{1}{1 + \exp[0.1 - 1.0 + 3.0]} =$
- $P(Y = 1) = \frac{\exp[0.1 - 1.0 + 3.0]}{1 + \exp[0.1 - 1.0 + 3.0]} =$
- Include bias, and sum the other weights

Logistic Regression Example

feature	coefficient	weight
bias	β_0	0.1
“viagra”	β_1	2.0
“mother”	β_2	-1.0
“work”	β_3	-0.5
“nigeria”	β_4	3.0

- $Y = 1$: spam

Example 2

$X = \{\text{Mother, Nigeria}\}$

- $P(Y = 0) = \frac{1}{1 + \exp[0.1 - 1.0 + 3.0]} = 0.11$
- $P(Y = 1) = \frac{\exp[0.1 - 1.0 + 3.0]}{1 + \exp[0.1 - 1.0 + 3.0]} = 0.89$
- Include bias, and sum the other weights

Logistic Regression Example

feature	coefficient	weight
bias	β_0	0.1
“viagra”	β_1	2.0
“mother”	β_2	-1.0
“work”	β_3	-0.5
“nigeria”	β_4	3.0

- $Y = 1$: spam

Example 3

$X = \{\text{Mother, Work, Viagra, Mother}\}$

Logistic Regression Example

feature	coefficient	weight
bias	β_0	0.1
“viagra”	β_1	2.0
“mother”	β_2	-1.0
“work”	β_3	-0.5
“nigeria”	β_4	3.0

- $Y = 1$: spam

Example 3

$X = \{\text{Mother, Work, Viagra, Mother}\}$

- $P(Y = 0) = \frac{1}{1 + \exp[0.1 - 1.0 - 0.5 + 2.0 - 1.0]} =$
- $P(Y = 1) = \frac{\exp[0.1 - 1.0 - 0.5 + 2.0 - 1.0]}{1 + \exp[0.1 - 1.0 - 0.5 + 2.0 - 1.0]} =$
- Multiply feature presence by weight

Logistic Regression Example

feature	coefficient	weight
bias	β_0	0.1
“viagra”	β_1	2.0
“mother”	β_2	-1.0
“work”	β_3	-0.5
“nigeria”	β_4	3.0

- $Y = 1$: spam

Example 3

$X = \{\text{Mother, Work, Viagra, Mother}\}$

- $P(Y = 0) = \frac{1}{1 + \exp[0.1 - 1.0 - 0.5 + 2.0 - 1.0]} = 0.60$
- $P(Y = 1) = \frac{\exp[0.1 - 1.0 - 0.5 + 2.0 - 1.0]}{1 + \exp[0.1 - 1.0 - 0.5 + 2.0 - 1.0]} = 0.40$
- Multiply feature presence by weight

How is Logistic Regression Used?

- Given a set of weights $\vec{\beta}$, we know how to compute the conditional likelihood $P(y|\beta, x)$
- Find the set of weights $\vec{\beta}$ that maximize the conditional likelihood on training data (next week)
- **Intuition:** higher weights mean that this feature implies that this feature is a good feature for the positive class

Outline

- Probabilistic classification
- Logistic regression
- **Generative vs. Discriminative models**

Generative vs. Discriminative Models

Discriminative

Model only conditional probability $p(y|x)$, excluding the data x .

Logistic regression

- Logistic: A special mathematical function it uses
- Regression: Combines a weight vector with observations to create an answer
- General cookbook for building conditional probability distributions

Generative

Model joint probability $p(x, y)$ including the data x .

Naïve Bayes

- Uses Bayes rule to reverse conditioning $p(x|y) \rightarrow p(y|x)$
- Naïve because it ignores joint probabilities within the data distribution

The Naïve Bayes classifier

- The Naïve Bayes classifier is a probabilistic classifier.
- We compute the probability of a document d being in a class c as follows:

$$P(c|d) \propto P(c, d) = P(c) \prod_{1 \leq i \leq n_d} P(w_i|c)$$

The Naïve Bayes classifier

- The Naïve Bayes classifier is a probabilistic classifier.
- We compute the probability of a document d being in a class c as follows:

$$P(c|d) \propto P(c, d) = P(c) \prod_{1 \leq i \leq n_d} P(w_i|c)$$

The Naïve Bayes classifier

- The Naïve Bayes classifier is a probabilistic classifier.
- We compute the probability of a document d being in a class c as follows:

$$P(c|d) \propto P(c, d) = P(c) \prod_{1 \leq i \leq n_d} P(w_i|c)$$

- n_d is the length of the document. (number of tokens)
- $P(w_i|c)$ is the conditional probability of term w_i occurring in a document of class c
- $P(w_i|c)$ as a measure of how much evidence w_i contributes that c is the correct class.
- $P(c)$ is the prior probability of c .
- If a document's terms do not provide clear evidence for one class vs. another, we choose the c with higher $P(c)$.

Maximum a posteriori class

- Our goal is to find the “best” class.
- The best class in Naïve Bayes classification is the most likely or *maximum a posteriori (MAP) class* c_{MAP} :

$$c_{\text{MAP}} = \arg \max_{c_j \in \mathbb{C}} \hat{P}(c_j|d) = \arg \max_{c_j \in \mathbb{C}} \hat{P}(c_j) \prod_{1 \leq i \leq n_d} \hat{P}(w_i|c_j)$$

- We write \hat{P} for P since these values are *estimates* from the training set.

Naive Bayes Classifier: More examples

This works because the coin flips are independent given the coin parameter. What about this case:

- want to identify the type of fruit given a set of features: color, shape and size
- color: red, green, yellow or orange (discrete)
- shape: round, oval or long+skinny (discrete)
- size: diameter in inches (continuous)



Naive Bayes Classifier: More examples

Using chain rule,

$$P(\textit{apple} \mid \textit{green}, \textit{round}, \textit{size} = 2)$$

Naive Bayes Classifier: More examples

Using chain rule,

$$\begin{aligned} &P(\text{apple} \mid \text{green}, \text{round}, \text{size} = 2) \\ &= \frac{P(\text{green}, \text{round}, \text{size} = 2 \mid \text{apple})P(\text{apple})}{\sum_{\text{fruits}} P(\text{green}, \text{round}, \text{size} = 2 \mid \text{fruit } j)P(\text{fruit } j)} \\ &\propto P(\text{green} \mid \text{round}, \text{size} = 2, \text{apple})P(\text{round} \mid \text{size} = 2, \text{apple}) \\ &\quad \times P(\text{size} = 2 \mid \text{apple})P(\text{apple}) \end{aligned}$$

But computing conditional probabilities is hard! There are many combinations of $(\text{color}, \text{shape}, \text{size})$ for each fruit.

Naive Bayes Classifier: More examples

Idea: assume conditional independence for all features given class,

$$P(\textit{green} \mid \textit{round}, \textit{size} = 2, \textit{apple}) = P(\textit{green} \mid \textit{apple})$$

$$P(\textit{round} \mid \textit{green}, \textit{size} = 2, \textit{apple}) = P(\textit{round} \mid \textit{apple})$$

$$P(\textit{size} = 2 \mid \textit{green}, \textit{round}, \textit{apple}) = P(\textit{size} = 2 \mid \textit{apple})$$

Naive Bayes Classifier: More examples

Idea: assume conditional independence for all features given class,

$$P(\textit{green} \mid \textit{round}, \textit{size} = 2, \textit{apple}) = P(\textit{green} \mid \textit{apple})$$

$$P(\textit{round} \mid \textit{green}, \textit{size} = 2, \textit{apple}) = P(\textit{round} \mid \textit{apple})$$

$$P(\textit{size} = 2 \mid \textit{green}, \textit{round}, \textit{apple}) = P(\textit{size} = 2 \mid \textit{apple})$$

$$P(\textit{apple} \mid \textit{green}, \textit{round}, \textit{size} = 2) \propto P(\textit{apple})P(\textit{green} \mid \textit{apple})P(\textit{round} \mid \textit{apple})P(\textit{size} = 2 \mid \textit{apple})$$

Naive Bayes Classifier: More examples

Conditioned on type of fruit, these features are not necessarily independent:



Given category “apple,” the color “green” has a higher probability given “size < 2”:

$$P(\text{green} \mid \text{size} < 2, \text{apple}) > P(\text{green} \mid \text{apple})$$

Contrasting Naïve Bayes and Logistic Regression

$$c_{\text{MAP}} = \arg \max_{c_j \in \mathbb{C}} \hat{P}(c_j|d) = \arg \max_{c_j \in \mathbb{C}} \hat{P}(c_j) \prod_{1 \leq i \leq n_d} \hat{P}(w_i|c_j)$$

Contrasting Naïve Bayes and Logistic Regression

$$c_{\text{MAP}} = \arg \max_{c_j \in \mathbb{C}} \hat{P}(c_j|d) = \arg \max_{c_j \in \mathbb{C}} \hat{P}(c_j) \prod_{1 \leq i \leq n_d} \hat{P}(w_i|c_j)$$

$$\arg \max_{c_j \in \mathbb{C}} [\ln \hat{P}(c_j) + \sum_{1 \leq i \leq n_d} \ln \hat{P}(w_i|c_j)]$$

Contrasting Naïve Bayes and Logistic Regression

$$c_{\text{MAP}} = \arg \max_{c_j \in \mathbb{C}} \hat{P}(c_j|d) = \arg \max_{c_j \in \mathbb{C}} \hat{P}(c_j) \prod_{1 \leq i \leq n_d} \hat{P}(w_i|c_j)$$

$$\arg \max_{c_j \in \mathbb{C}} [\ln \hat{P}(c_j) + \sum_{1 \leq i \leq n_d} \ln \hat{P}(w_i|c_j)]$$

Naïve Bayes is a special case of logistic regression that uses Bayes rule and conditional probabilities to set these weights

Contrasting Naïve Bayes and Logistic Regression

$$c_{\text{MAP}} = \arg \max_{c_j \in \mathbb{C}} \hat{P}(c_j|d) = \arg \max_{c_j \in \mathbb{C}} \hat{P}(c_j) \prod_{1 \leq i \leq n_d} \hat{P}(w_i|c_j)$$

$$\arg \max_{c_j \in \mathbb{C}} [\ln \hat{P}(c_j) + \sum_{1 \leq i \leq n_d} \ln \hat{P}(w_i|c_j)]$$

Naïve Bayes is a special case of logistic regression that uses Bayes rule and conditional probabilities to set these weights

Contrasting Naïve Bayes and Logistic Regression

- Naïve Bayes is easier for learning
- Naïve Bayes works better on smaller datasets
- Logistic regression works better on medium-sized datasets
- On huge datasets, both algorithms perform about the same (data always win)
- Logistic regression allows for arbitrary features (biggest difference!)

Contrasting Naïve Bayes and Logistic Regression

- Naïve Bayes is easier for learning
- Naïve Bayes works better on smaller datasets
- Logistic regression works better on medium-sized datasets
- On huge datasets, both algorithms perform about the same (data always win)
- Logistic regression allows for arbitrary features (biggest difference!)
- Don't need to memorize (or work through) previous slide—just understand that naïve Bayes is a special case of logistic regression