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

*Naïve Bayes*

# Outline

- Probabilistic classification

- Logistic regression

- Generative vs. Discriminative models

# Outline

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

K-nearest neighbor

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

## Recap

K-nearest neighbor

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

Perceptron

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

# Probabilistic Models

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

## Probabilistic Models

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

# Probabilistic Classification

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

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

**Goal**: $h : X \to Y$

*[handwritten: $P(y|x)$]*

For each class $c_j$, estimate

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

Assign to $\boldsymbol{x}$ the class with the highest probability

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

*[handwritten: Perceptron]*

*[handwritten: $\hat{y} = \text{sign}(w \cdot x + b)$]*

# 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 $\beta_i$
- Observations $X_i$
- "Bias" $\beta_0$ (like intercept in linear regression)

$P(y|x)$

$P(Y=0|x) + P(Y=(1|x) = 1$

$$\sum_{j=1}^{c} P(y=c|x) = 1$$

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

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

# Logistic Regression: Definition

- Weight vector $\beta_i$
- Observations $X_i$
- "Bias" $\beta_0$ (like intercept in linear regression)

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

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

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

# Logistic Regression: Definition

- Weight vector $\beta_i$
- Observations $X_i$
- "Bias" $\beta_0$ (like intercept in linear regression)

$$\hat{y} = \arg\max_{c} P(Y = c | x)$$

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

$$= \frac{\exp(-\beta_0 - \sum_i \beta_i X_i)}{1 + \exp(-\beta_0 - \sum_i \beta_i X_i)}$$

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

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

$$P(Y = 1 | X) = P(Y = 0 | X) = \frac{1}{2}$$

What is the decision boundary?

$$\beta_0 + \sum_i \beta_i X_i = 0$$

$$y = \begin{cases} 1 & \beta_0 + \sum_i \beta_i X_i > 0 \\ 0 & \beta_0 + \sum_i \beta_i X_i < 0 \end{cases}$$

# Logistic Regression: Definition

- Weight vector $\beta_i$

- Observations $X_i$

- For shorthand, we'll say that

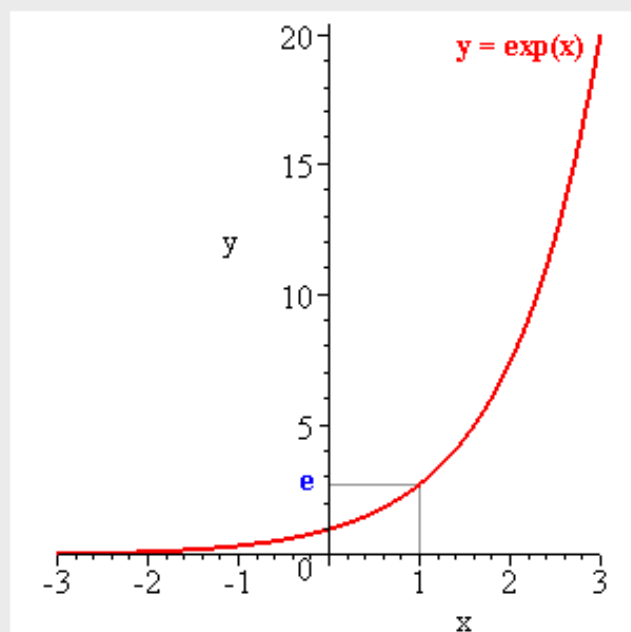$$P(Y = 1|X) = \sigma((\beta_0 + \sum_i \beta_i X_i)) = \frac{1}{1 + exp(-z)} \quad (3)$$

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

- Where $\sigma(z) = \frac{1}{1 + exp[-z]}$
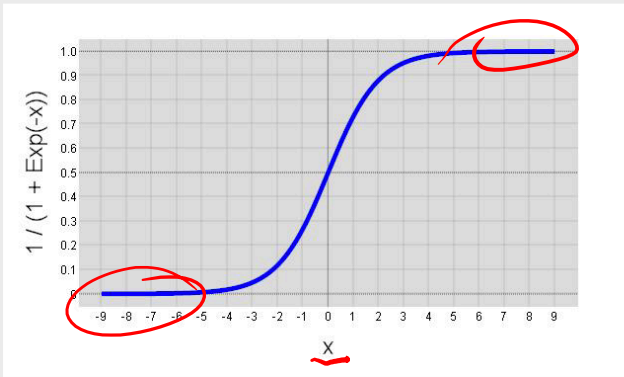
# What's this "exp" doing?

### 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 $\sigma(z) = \dfrac{1}{1+e^{-z}}$
- 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 | $\beta_0$ | 0.1 |
| "viagra" | $\beta_1$ | 2.0 |
| "mother" | $\beta_2$ | $-1.0$ |
| "work" | $\beta_3$ | $-0.5$ |
| "nigeria" | $\beta_4$ | 3.0 |

- What does $Y = 1$ mean?

Spam

Example 1: Empty Document?

$$X = \{\}$$

$$P(Y=1 \mid \{\}) = \frac{1}{1 + \exp(-0.1)} > 0.5$$

# Logistic Regression Example

| feature | coefficient | weight |
|---------|-------------|--------|
| bias | $\beta_0$ | 0.1 |
| "viagra" | $\beta_1$ | 2.0 |
| "mother" | $\beta_2$ | -1.0 |
| "work" | $\beta_3$ | -0.5 |
| "nigeria" | $\beta_4$ | 3.0 |

- $Y = 1$: spam

## Example 1: Empty Document?

$X = \{\}$

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

- $P(Y = 1) = \dfrac{\exp[0.1]}{1 + \exp[0.1]} = \dfrac{1}{\exp(-21) + 1}$

# Logistic Regression Example

| feature | coefficient | weight |
|---------|-------------|--------|
| bias | $\beta_0$ | 0.1 |
| "viagra" | $\beta_1$ | 2.0 |
| "mother" | $\beta_2$ | −1.0 |
| "work" | $\beta_3$ | −0.5 |
| "nigeria" | $\beta_4$ | 3.0 |

- $Y = 1$: spam

$X = \{\}$

- $P(Y = 0) = \frac{1}{1+\exp[0.1]} = 0.48$

- $P(Y = 1) = \frac{\exp[0.1]}{1+\exp[0.1]} = 0.52$

- Bias $\beta_0$ encodes the prior probability of a class

# Logistic Regression Example

| feature | coefficient | weight |
| --- | --- | --- |
| bias | $\beta_0$ | 0.1 |
| "viagra" | $\beta_1$ | 2.0 |
| "mother" | $\beta_2$ | -1.0 |
| "work" | $\beta_3$ | -0.5 |
| "nigeria" | $\beta_4$ | 3.0 |

- $Y = 1$: spam

**Example 2**

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

$$-1 + 3 + 0.1 = 2.1$$

$$P(y=1|x) = \frac{1}{1 + \exp(-2.1)} > 0.5$$

# Logistic Regression Example

| feature | coefficient | weight |
|---------|-------------|--------|
| bias | $\beta_0$ | 0.1 |
| "viagra" | $\beta_1$ | 2.0 |
| "mother" | $\beta_2$ | $-1.0$ |
| "work" | $\beta_3$ | $-0.5$ |
| "nigeria" | $\beta_4$ | 3.0 |

- $Y = 1$: spam

## Example 2

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

- $P(Y = 0) = \dfrac{1}{1 + \exp\left[0.1 - 1.0 + 3.0\right]} =$

- $P(Y = 1) = \dfrac{\exp\left[0.1 - 1.0 + 3.0\right]}{1 + \exp\left[0.1 - 1.0 + 3.0\right]} =$

- Include bias, and sum the other weights

# Logistic Regression Example

| feature | coefficient | weight |
|---------|-------------|--------|
| bias | $\beta_0$ | 0.1 |
| "viagra" | $\beta_1$ | 2.0 |
| "mother" | $\beta_2$ | $-1.0$ |
| "work" | $\beta_3$ | $-0.5$ |
| "nigeria" | $\beta_4$ | 3.0 |

- $Y = 1$: spam

## Example 2

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

- $P(Y = 0) = \dfrac{1}{1 + \exp\left[0.1 - 1.0 + 3.0\right]} = 0.11$

- $P(Y = 1) = \dfrac{\exp\left[0.1 - 1.0 + 3.0\right]}{1 + \exp\left[0.1 - 1.0 + 3.0\right]} = 0.89$

- Include bias, and sum the other weights

# Logistic Regression Example

| feature | coefficient | weight |
|---------|-------------|--------|
| bias | $\beta_0$ | 0.1 |
| "viagra" | $\beta_1$ | 2.0 |
| "mother" | $\beta_2$ | $-1.0$ |
| "work" | $\beta_3$ | $-0.5$ |
| "nigeria" | $\beta_4$ | 3.0 |

- $Y = 1$: spam

Example 3

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

$-1 + -0.5 + 2 + -1 + 0.1 = -0.4$

## Logistic Regression Example

| feature | coefficient | weight |
|---------|-------------|--------|
| bias | $\beta_0$ | 0.1 |
| "viagra" | $\beta_1$ | 2.0 |
| "mother" | $\beta_2$ | $-1.0$ |
| "work" | $\beta_3$ | $-0.5$ |
| "nigeria" | $\beta_4$ | 3.0 |

- $Y = 1$: spam

### Example 3

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

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

# Logistic Regression Example

| feature | coefficient | weight |
|---------|-------------|--------|
| bias | $\beta_0$ | 0.1 |
| "viagra" | $\beta_1$ | 2.0 |
| "mother" | $\beta_2$ | $-1.0$ |
| "work" | $\beta_3$ | $-0.5$ |
| "nigeria" | $\beta_4$ | 3.0 |

- $Y = 1$: spam

- $P(Y = 0) = \frac{1}{1+\exp\left[0.1-1.0-0.5+2.0-1.0\right]} = 0.60$

- $P(Y = 1) = \frac{\exp\left[0.1-1.0-0.5+2.0-1.0\right]}{1+\exp\left[0.1-1.0-0.5+2.0-1.0\right]} = 0.40$

- 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 $x$.

$$\sum p(x,y) \quad p_x$$

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:

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

$$= P(c) \, P(d|c)$$

the class is happening

happening vs the class

$$\prod_{1 \le i \le n_d} P(w_i|c)$$

$$P(d|c) = P(w_1 \cdots, w_{n_d}|c)$$

$$= P(w_1|c) \cdots P(w_{n_d}|c)$$

$$= \prod_{1 \le i \le n_d} P(w_i|c)$$

## 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 $P(d)$.

# 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* $c_{\mathrm{MAP}}$:

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

- 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)

Using chain rule,

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

Using chain rule,

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

$$= \frac{P(green, round, size = 2 \mid apple)P(apple)}{\sum_{fruits} P(green, round, size = 2 \mid fruit\,j)P(fruit\,j)}$$

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

$$\times P(size = 2 \mid apple)P(apple)$$

$$P(c, d) = P(c)\,P(d|c)$$

$$\frac{}{\sum_{c \in fruits} P(c, d)}$$

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

## Naive Bayes Classifier: More examples

Idea: assume conditional independence for all features given class,

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

Idea: assume conditional independence for all features given class,

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

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

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

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

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



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

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

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

Logistic

$$C = \arg\max_{c} \hat{P}(c|d) \overset{\text{argmax}}{=_c} C \cdot \hat{P}(c=1|d) + (1-c) \hat{P}(c=0|d)$$

$$= (\beta_0 + \widehat{\sum_j \beta_i x_i}) > 0$$

$$C_{map} = \arg\max_{c_j \in C} \log \hat{P}(c_j) + \left(\sum_{i=1}^{d} \log \hat{P}(w_i|c_j)\right)$$

$$\beta_0 = \log \hat{P}(c_j)$$

$$\max_{c_j} \log P(c_j)$$

$$\beta_0 = \log \hat{P}_{(1)} - \log \hat{P}_{(0)} = \log \frac{\hat{P}_{(1)}}{\hat{P}_{(0)}}$$

## Contrasting Naïve Bayes and Logistic Regression

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

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

# Contrasting Naïve Bayes and Logistic Regression

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

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

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

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

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

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