

Logistic Regression

Il-Chul Moon
Dept. of Industrial and Systems Engineering
KAIST

icmoon@kaist.ac.kr

Weekly Objectives

- Learn the logistic regression classifier
 - Understand why the logistic regression is better suited than the linear regression for classification tasks
 - Understand the logistic function
 - Understand the logistic regression classifier
 - Understand the approximation approach for the open form solutions
- Learn the gradient descent algorithm
 - Know the Taylor expansion
 - Understand the gradient descent/ascent algorithm
- Learn the difference between the naïve Bayes and the logistic regression
 - Understand the similarity of the two classifiers
 - Understand the differences of the two classifiers
 - Understand the performance differences

NAÏVE BAYES VS. LOGISTIC REGRESSION

Gaussian Naïve Bayes

- We want to compare the performance of the two classifiers
 - Logistic regression handles the continuous features
 - Why not naïve Bayes?
- Naïve Bayes Classifier Function
 - $f_{NB}(x) = \operatorname{argmax}_{Y=y} P(Y = y) \prod_{1 \leq i \leq d} P(X_i = x_i | Y = y)$
- What-if the feature is a continuous random variable?
 - We can assume that the variable follows the Gaussian distribution with the mean of μ and the variance of σ^2
 - $P(X_i | Y, \mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(X_i - \mu)^2}{2\sigma^2}}$
 - In addition, let's use more shortened terms
 - $P(Y = y) = \pi_1$
 - $P(Y) \prod_{1 \leq i \leq d} P(X_i | Y) = \pi_k \prod_{1 \leq i \leq d} \frac{1}{\sigma_k^i C} \exp\left(-\frac{1}{2} \left(\frac{X_i - \mu_k^i}{\sigma_k^i}\right)^2\right)$

Derivation to Logistic Regression (1)

- Derivation from the naïve Bayes to the logistic regression

- $P(Y) \prod_{1 \leq i \leq d} P(X_i|Y) = \pi_k \prod_{1 \leq i \leq d} \frac{1}{\sigma_k^i} \exp\left(-\frac{1}{2} \left(\frac{X_i - \mu_k^i}{\sigma_k^i}\right)^2\right)$

- With naïve Bayes assumption

- $$P(Y = y|X) = \frac{P(X|Y = y)P(Y=y)}{P(X)} = \frac{P(X|Y = y)P(Y=y)}{P(X|Y = y)P(Y=y) + P(X|Y = n)P(Y=n)}$$
$$= \frac{P(Y = y) \prod_{1 \leq i \leq d} P(X_i|Y = y)}{P(Y = y) \prod_{1 \leq i \leq d} P(X_i|Y = y) + P(Y = n) \prod_{1 \leq i \leq d} P(X_i|Y = n)}$$

Derivation to Logistic Regression (2)

- With naïve Bayes assumption

$$\begin{aligned}
 \bullet \quad P(Y = y|X) &= \frac{P(X|Y = y)P(Y=y)}{P(X)} = \frac{P(X|Y = y)P(Y=y)}{P(X|Y = y)P(Y=y)+P(X|Y = n)P(Y=n)} \\
 &= \frac{P(Y = y) \prod_{1 \leq i \leq d} P(X_i|Y = y)}{P(Y = y) \prod_{1 \leq i \leq d} P(X_i|Y = y) + P(Y = n) \prod_{1 \leq i \leq d} P(X_i|Y = n)}
 \end{aligned}$$

$$\begin{aligned}
 \bullet \quad P(Y = y|X) &= \frac{\pi_1 \prod_{1 \leq i \leq d} \frac{1}{\sigma_1^i C} \exp(-\frac{1}{2} \left(\frac{X_i - \mu_1^i}{\sigma_1^i} \right)^2)}{\pi_1 \prod_{1 \leq i \leq d} \frac{1}{\sigma_1^i C} \exp(-\frac{1}{2} \left(\frac{X_i - \mu_1^i}{\sigma_1^i} \right)^2) + \pi_2 \prod_{1 \leq i \leq d} \frac{1}{\sigma_2^i C} \exp(-\frac{1}{2} \left(\frac{X_i - \mu_2^i}{\sigma_2^i} \right)^2)} \\
 &= \frac{1}{1 + \frac{\pi_2 \prod_{1 \leq i \leq d} \frac{1}{\sigma_2^i C} \exp(-\frac{1}{2} \left(\frac{X_i - \mu_2^i}{\sigma_2^i} \right)^2)}{\pi_1 \prod_{1 \leq i \leq d} \frac{1}{\sigma_1^i C} \exp(-\frac{1}{2} \left(\frac{X_i - \mu_1^i}{\sigma_1^i} \right)^2)}}
 \end{aligned}$$

Derivation to Logistic Regression (3)

- Assuming the same variable of the two classes, $\sigma_2^i = \sigma_1^i$

$$\begin{aligned}
 P(Y = y|X) &= \frac{1}{1 + \frac{\pi_2 \prod_{1 \leq i \leq d} \frac{1}{\sigma_2^i} \exp(-\frac{1}{2} \left(\frac{X_i - \mu_2^i}{\sigma_2^i} \right)^2)}{\pi_1 \prod_{1 \leq i \leq d} \frac{1}{\sigma_1^i} \exp(-\frac{1}{2} \left(\frac{X_i - \mu_1^i}{\sigma_1^i} \right)^2)}} = \frac{1}{1 + \frac{\pi_2 \prod_{1 \leq i \leq d} \exp(-\frac{1}{2} \left(\frac{X_i - \mu_2^i}{\sigma_2^i} \right)^2)}{\pi_1 \prod_{1 \leq i \leq d} \exp(-\frac{1}{2} \left(\frac{X_i - \mu_1^i}{\sigma_1^i} \right)^2)}} \\
 &= \frac{1}{1 + \frac{\pi_2 \exp(-\sum_{1 \leq i \leq d} \{ \frac{1}{2} \left(\frac{X_i - \mu_2^i}{\sigma_2^i} \right)^2 \})}{\pi_1 \exp(-\sum_{1 \leq i \leq d} \{ \frac{1}{2} \left(\frac{X_i - \mu_1^i}{\sigma_1^i} \right)^2 \})}} \\
 &= \frac{1}{1 + \frac{\exp(-\sum_{1 \leq i \leq d} \left\{ \frac{1}{2} \left(\frac{X_i - \mu_2^i}{\sigma_2^i} \right)^2 \right\} + \log \pi_2)}{\exp(-\sum_{1 \leq i \leq d} \left\{ \frac{1}{2} \left(\frac{X_i - \mu_1^i}{\sigma_1^i} \right)^2 \right\} + \log \pi_1)}}
 \end{aligned}$$

Derivation to Logistic Regression (4)

- Assuming the same variable of the two classes, $\sigma_2^i = \sigma_1^i$

$$P(Y = y|X) = \frac{1}{1 + \exp(-\sum_{1 \leq i \leq d} \left\{ \frac{1}{2} \left(\frac{X_i - \mu_2^i}{\sigma_2^i} \right)^2 \right\} + \log \pi_2 + \sum_{1 \leq i \leq d} \left\{ \frac{1}{2} \left(\frac{X_i - \mu_1^i}{\sigma_1^i} \right)^2 \right\} - \log \pi_1)}$$

$$= \frac{1}{1 + \exp(-\frac{1}{2(\sigma_1^i)^2} \sum_{1 \leq i \leq d} \{ (X_i - \mu_1^i)^2 - (X_i - \mu_2^i)^2 \} + \log \pi_2 - \log \pi_1)}$$

$$= \frac{1}{1 + \exp(-\frac{1}{2(\sigma_1^i)^2} \sum_{1 \leq i \leq d} \{ 2(\mu_2^i - \mu_1^i)X_i + \mu_2^{i^2} - \mu_1^{i^2} \} + \log \pi_2 - \log \pi_1)}$$

Naïve Bayes vs. Logistic Regression

- Naïve Bayes classifier

- $$P(Y|X) = \frac{1}{1 + \exp\left(-\frac{1}{2(\sigma_1^i)^2} \sum_{1 \leq i \leq d} \{2(\mu_2^i - \mu_1^i)X_i + \mu_2^{i^2} - \mu_1^{i^2}\} + \log \pi_2 - \log \pi_1\right)}$$

- Assumption to get this formula

- Naïve Bayes assumption, Same variance assumption between classes
 - Gaussian distribution for $P(X|Y)$
 - Bernoulli distribution for $P(Y)$

Together, modeling joint prob.

- # of parameters to estimate = $2 \times 2 \times d + 1 = 4d + 1$

- With the different variances between classes

- Logistic Regression

- $$P(Y|X) = \frac{1}{1 + e^{-\theta^T x}}$$

- Assumption to get this formula

- Fitting to the logistic function

- # of parameters to estimate = $d + 1$

- Who is the winner?

- Really??? There is no winner... Why?

Generative-Discriminative Pair

- Generative model, $P(Y|X)=P(X,Y)/P(X)=P(X|Y)P(Y)/P(X)$
 - Full probabilistic model of all variables
 - Estimate the parameters of $P(X|Y)$, $P(Y)$ from the data
 - Characteristics: Bayesian, Prior, Modeling the joint probability
 - Naïve Bayes Classifier
- Discriminative model, $P(Y|X)$
 - Do not need to model the distribution of the observed variables
 - Estimate the parameters of $P(Y|X)$ from the data
 - Characteristics: Modeling the conditional probability
 - Logistic Regression
- Pros and Cons [Ng & Jordan, 2002]
 - Logistic regression is less biased
 - Probably approximately correct learning: Naïve Bayes learns faster