

CSE 802
Spring 2017
Logistic Regression

Inci M. Baytas
Computer Science
Michigan State University

March 29, 2017

Introduction

- ▶ Consider two-class classification problem, the posterior probability of class C_1 can be written as:

$$p(C_1|\Phi) = y(\Phi) = \sigma(\mathbf{w}^T\Phi) \quad (1)$$

- ▶ $\sigma(\cdot)$ is the logistic sigmoid function.
- ▶ $p(C_2|\Phi) = 1 - p(C_1|\Phi)$
- ▶ Φ is a feature vector, a non-linear transformation on original observation space \mathbf{x} .
- ▶ The model in Eq.1 is called as *Logistic Regression* in the terminology of statistics.

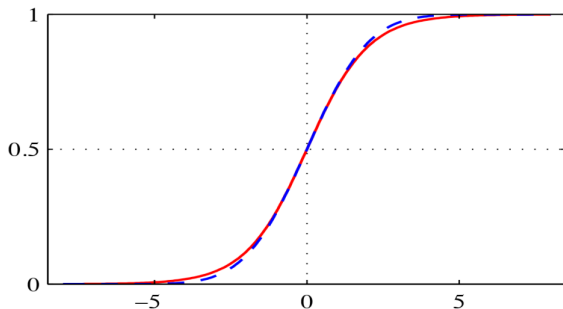
Logistic Regression I

- ▶ A classification model rather than regression
- ▶ A probabilistic discriminative model
 - ▶ We estimate the parameter \mathbf{w} directly.
- ▶ Comparison of logistic regression and generative model in M -dimensional space Φ :
 - ▶ Logistic regression: M adjustable parameters.
 - ▶ Generative models: Assume we fit Gaussian class conditional densities using maximum likelihood; $M(M+5)/2 + 1 =$
Means: $2M$ + Shared covariance: $(M+1)M/2$ + Prior $p(C_1)$: 1
- ▶ **Maximum likelihood** is used to determine the parameters of logistic regression model.

Logistic Regression II

- ▶ Definition and properties of logistic sigmoid function:

$$\begin{aligned}\sigma(a) &= \frac{1}{1 + \exp(-a)} \\ \sigma(-a) &= 1 - \sigma(a) \\ \frac{d\sigma}{da} &= \sigma(1 - \sigma)\end{aligned}\tag{2}$$



Logistic Regression III - How to Estimate \mathbf{w}

- ▶ For a training data set $\{\Phi_n, t_n\}$, where $t_n \in \{0, 1\}$ and $\Phi_n = \Phi(x_n)$, with $n = 1, \dots, N$, the log likelihood can be written as:

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n} \quad (3)$$

where $\mathbf{t} = (t_1, \dots, t_N)^T$ and $y_n = p(C_1|\Phi_n)$

- ▶ The error function is the negative logarithm of the likelihood, known as *Cross-entropy* error function:

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln (1 - y_n)\} \quad (4)$$

where $y_n = \sigma(a_n)$ and $a_n = \mathbf{w}^T \Phi_n$.

Logistic Regression IV - How to Estimate \mathbf{w}

- ▶ There is no analytical (closed-form) solution.
- ▶ The cross entropy loss is a convex function.
 - ▶ There is a global minimum.
 - ▶ Can use an iterative approach.
- ▶ Calculate the gradient with respect to \mathbf{w} :

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \Phi_n \quad (5)$$

- ▶ Use gradient descent (batch or online):

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E(\mathbf{w}^{\tau}) \quad (6)$$

Logistic Regression V - How to Estimate w

- ▶ Newton-Raphson Algorithm

$$\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \mathbf{H}^{-1} \nabla E(\mathbf{w}) \quad (7)$$

- ▶ It uses a local quadratic approximation to the cross-entropy error function to update w iteratively.
- ▶ Newton-Raphson algorithm is also known as **iterative reweighted least squares**.
- ▶ Convexity: \mathbf{H} is positive definite (eigenvalues of \mathbf{H} are non-negative).

Multi-class Logistic Regression

- ▶ Cross-entropy for multi-class classification problem:

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk} \quad (8)$$

where $y_k(\Phi) = p(C_k|\Phi) = \frac{\exp(\mathbf{w}_k^T \Phi)}{\sum_j \exp(\mathbf{w}_j^T \Phi)}$ which is called *softmax function*.

- ▶ Use maximum likelihood to estimate the parameters.
- ▶ Use an iterative approach such as Newton-Rapson.

Over-fitting in Logistic Regression

- ▶ Maximum likelihood can suffer from severe over-fitting.
- ▶ This can be overcome by finding a MAP solution for \mathbf{w} (Bayesian treatment).
- ▶ Another alternative is to use regularization.
- ▶ Add regularizers to the loss function, regularized log-likelihood.
 - ▶ l_2 norm
 - ▶ l_1 norm (Lasso)

References

- ▶ Classification lecture of Dr. Jiayu Zhou.
- ▶ Christopher Bishop, Pattern Recognition and Machine Learning, *Information Science and Statistics*, Springer-Verlag New York, 2006.