

Logistic Regression

A brief tutorial



Jay

January 28, 2015

Outline

- 1 Definition
 - Basics
 - Parameters
 - Basic Model
 - Training
 - Testing
- 2 Motivation
 - Maximum Likelihood
 - Maximum Entropy
- 3 Derivatives
 - Maximum Entropy Model
 - Naive Bayes
 - CRF
- 4 Code Review



Outline

- 1 **Definition**
 - Basics
 - Parameters
 - Basic Model
 - Training
 - Testing
- 2 **Motivation**
 - Maximum Likelihood
 - Maximum Entropy
- 3 **Derivatives**
 - Maximum Entropy Model
 - Naive Bayes
 - CRF
- 4 **Code Review**



Basics

- Counterpart of Linear Regression
- Basic model for Classification
- Logistic: Uses the logit function
- Regression: Predicts the probability of an outcome

Outline

- 1 **Definition**
 - Basics
 - **Parameters**
 - Basic Model
 - Training
 - Testing
- 2 **Motivation**
 - Maximum Likelihood
 - Maximum Entropy
- 3 **Derivatives**
 - Maximum Entropy Model
 - Naive Bayes
 - CRF
- 4 **Code Review**



Parameters

\mathbf{x}^l	l instance에 대한 \mathbb{R}^n 크기의 feature vector
y^l	l instance에 대한 \mathbb{R} 크기의 label
w^l	Model의 \mathbb{R}^n 크기의 weight vector
b	Model의 \mathbb{R} 크기의 bias value

Outline

- 1 **Definition**
 - Basics
 - Parameters
 - **Basic Model**
 - Training
 - Testing
- 2 **Motivation**
 - Maximum Likelihood
 - Maximum Entropy
- 3 **Derivatives**
 - Maximum Entropy Model
 - Naive Bayes
 - CRF
- 4 **Code Review**

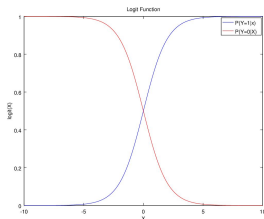


Logistic Regression Model

Modeling of log of odds ratio

$$\log \frac{P(y^l = 1|x)}{1 - P(y^l = 1|x)} = w^T x + b$$

- $P(y^l = 1|x^l) = \frac{e^{w^T x + b}}{1 + e^{w^T x + b}}$
- $P(y^l = 0|x^l) = \frac{1}{1 + e^{w^T x + b}}$



Outline

- 1 **Definition**
 - Basics
 - Parameters
 - Basic Model
 - **Training**
 - Testing
- 2 **Motivation**
 - Maximum Likelihood
 - Maximum Entropy
- 3 **Derivatives**
 - Maximum Entropy Model
 - Naive Bayes
 - CRF
- 4 **Code Review**



Training

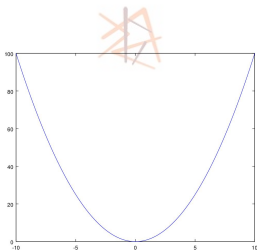
Cost function

$$J = - \sum_l y^l \ln(P(y^l = 1 | \mathbf{x}^l)) + (1 - y^l) \ln(1 - P(y^l = 0 | \mathbf{x}^l))$$

- If $y^l = 1$
 $\ln(P(y^l = 1 | \mathbf{x}^l))$ 을 최대화하는 w
- If $y^l = 0$
 $\ln(1 - P(y^l = 0 | \mathbf{x}^l))$ 을 최대화하는 w

Gradient Descent (because not strictly convex)

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial J}{\partial \mathbf{w}}$$



Outline

- 1 Definition
 - Basics
 - Parameters
 - Basic Model
 - Training
 - **Testing**
- 2 Motivation
 - Maximum Likelihood
 - Maximum Entropy
- 3 Derivatives
 - Maximum Entropy Model
 - Naive Bayes
 - CRF
- 4 Code Review



Testing

- $P(y^l = 1 | \mathbf{x}^l) = \frac{e^{\mathbf{w}^T \mathbf{x} + b}}{1 + e^{\mathbf{w}^T \mathbf{x} + b}}$
- $P(y^l = 0 | \mathbf{x}^l) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x} + b}}$

Example:

- $x_i = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$
- $w = \begin{pmatrix} 3 \\ 7 \end{pmatrix}$
- $b = 1$
- $P(y^l = 1 | \mathbf{x}^l) = \frac{e^{-3}}{1 + e^{-3}} = 0.0474$
- $P(y^l = 0 | \mathbf{x}^l) = \frac{1}{1 + e^{-3}} = 0.953$

Thus, x_i is classified as $y^l = 0$



Outline

- 1 Definition
 - Basics
 - Parameters
 - Basic Model
 - Training
 - Testing
- 2 Motivation
 - **Maximum Likelihood**
 - Maximum Entropy
- 3 Derivatives
 - Maximum Entropy Model
 - Naive Bayes
 - CRF
- 4 Code Review



Maximum Likelihood

Bayes' formula

$$\begin{aligned}
 \text{Posterior} &= \frac{\text{likelihood} \times \text{prior}}{\text{evidence}} \\
 P(x|y) &= \frac{P(y|x) \times P(x)}{P(y)}
 \end{aligned} \tag{1}$$

Assuming i.i.d

$$\begin{aligned}
 \text{likelihood}(\mathbf{w}) &= \prod_l P(Y^l = y_k | \mathbf{x}^l, \mathbf{w}) \\
 l(\mathbf{w}) &= -\ln \prod_l P(y^l = y_k | \mathbf{x}^l, \mathbf{w}) \\
 &= -\sum_l \ln P(y^l = y_k | \mathbf{x}^l)
 \end{aligned} \tag{2}$$

Outline

- 1 Definition
 - Basics
 - Parameters
 - Basic Model
 - Training
 - Testing
- 2 **Motivation**
 - Maximum Likelihood
 - **Maximum Entropy**
- 3 Derivatives
 - Maximum Entropy Model
 - Naive Bayes
 - CRF
- 4 Code Review



Maximum Entropy



$$P(y = \textit{scissors}) = 33\%$$

$$P(y = \textit{rock}) = 33\%$$

$$P(y = \textit{paper}) = 33\%$$



$$P(y = \textit{scissors}) = 98\%$$

$$P(y = \textit{rock}) = 1\%$$

$$P(y = \textit{paper}) = 1\%$$

Entropy : $H = - \sum_i p_i \log_2 p_i$

Conditional Entropy : $H(y|x) = - \sum_{x,y} \tilde{p}(x)p(y|x) \log p(y|x)$

Outline

- 1 Definition
 - Basics
 - Parameters
 - Basic Model
 - Training
 - Testing
- 2 Motivation
 - Maximum Likelihood
 - Maximum Entropy
- 3 **Derivatives**
 - **Maximum Entropy Model**
 - Naive Bayes
 - CRF
- 4 Code Review



Maximum Entropy Model

Logistic Regression is equivalent to the basic Maximum Entropy model

$$\begin{aligned}
 \xi(p, \lambda, \gamma) = & - \left(\sum_{x,y} \tilde{p}(x)p(y|x) \log p(y|x) \right) \\
 & + \left(\sum_i \lambda_i \sum_{x,y} \tilde{p}(x)p(y|x) f_i(x, y) - \tilde{p}(x, y) f_i(x, y) \right) \\
 & + \gamma \left(1 - \sum_y p(y|x) \right)
 \end{aligned}
 \tag{3}$$

Maximum Entropy Model

$$\frac{\partial \xi}{\partial p(y|x)} = 0$$

Primal Function

$$p(y|x) = \exp\left(\sum_i \lambda_i f_i(x, y)\right) \exp\left(-\frac{\gamma}{\tilde{p}(x)} - 1\right)$$

Outline

- 1 Definition
 - Basics
 - Parameters
 - Basic Model
 - Training
 - Testing
- 2 Motivation
 - Maximum Likelihood
 - Maximum Entropy
- 3 **Derivatives**
 - Maximum Entropy Model
 - **Naive Bayes**
 - CRF
- 4 Code Review



Naive Bayes

Naive Bayes (Generative) vs Logistic Regression (Discriminative)

- Discriminative:

$$\operatorname{argmax}_y p(y|x)$$

- Generative:

$$\operatorname{argmax}_y p(x|y)p(y)$$

$$P(x_1, x_2, \dots, x_n|y) = \prod_i P(x_i|y)$$

$$y \leftarrow \operatorname{argmax}_k \left(P(Y = y_k) \prod_i P(x_i|Y = y_k) \right)$$

Naive Bayes

Assuming a Gaussian Naive Bayes with parameters $\mu_{ik}, \sigma_{ik}^2, \pi_k$

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

$$P(Y = 1|X) = \frac{1}{1 + \exp\left(\ln\left(\frac{1-\pi}{\pi}\right) + \sum_i \left(x_i \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}\right)\right)}$$

Equivalent model assuming

- Infinite i.i.d data
- Naive bayes assumption

Outline

- 1 Definition
 - Basics
 - Parameters
 - Basic Model
 - Training
 - Testing
- 2 Motivation
 - Maximum Likelihood
 - Maximum Entropy
- 3 **Derivatives**
 - Maximum Entropy Model
 - Naive Bayes
 - **CRF**
- 4 Code Review



CRF

Basic equation

$$p(\mathbf{x}, \mathbf{y}) = \frac{1}{Z} \prod_A \psi_A(\mathbf{x}_A, \mathbf{y}_A)$$

Definition of psi

$$\psi_A(\mathbf{x}_A, \mathbf{y}_A) = \exp \left(\sum_k \theta_{Ak} f_{Ak}(\mathbf{x}_A, \mathbf{y}_A) \right)$$

HMM implementation of CRF

$$P(y, \mathbf{x}) = \exp \left\{ \sum_t \left(\sum_{i, j \in S} 1_{y_{t-1}=j} 1_{y_t=i} \ln p(y_t | y_{t-1}) + \sum_{i \in S} \sum_{o \in O} 1_{x_t=o} 1_{y_t=i} \ln p(x_t | y_t) \right) \right\}$$

Outline

- 1 Definition
 - Basics
 - Parameters
 - Basic Model
 - Training
 - Testing
- 2 Motivation
 - Maximum Likelihood
 - Maximum Entropy
- 3 Derivatives
 - Maximum Entropy Model
 - Naive Bayes
 - CRF
- 4 Code Review



Code Review

- 1 Preprocessing
- 2 Running Algorithm
- 3 Evaluation

