Logistic Regression

Sargur N. Srihari University at Buffalo, State University of New York USA

Topics in Linear Classification using Probabilistic Discriminative Models

- Generative vs Discriminative
- 1. Fixed basis functions
- 2. Logistic Regression (two-class)
- 3. Iterative Reweighted Least Squares (IRLS)
- 4. Multiclass Logistic Regression
- 5. Probit Regression
- 6. Canonical Link Functions

Topics in Logistic Regression

- Logistic Sigmoid and Logit Functions
- Parameters in discriminative approach
- Determining logistic regression parameters
 - Error function
 - Gradient of error function
 - Simple sequential algorithm
 - An example
- Generative vs Discriminative Training

– Naiive Bayes vs Logistic Regression

Logistic Sigmoid and Logit Functions

 $f(x) = \log \frac{x}{1-x}$

• In two-class case, *posterior* of class C_1 can be written as as a logistic sigmoid of feature vector $\mathbf{\phi} = [\phi_1, .., \phi_M]^T$

 $p(C_1|\mathbf{\phi}) = y(\mathbf{\phi}) = \sigma(\mathbf{w}^{\mathrm{T}}\mathbf{\phi})$

with $p(C_2|\phi) = 1 - p(C_1|\phi)$ Here σ (.) is the logistic sigmoid function

- Known as logistic regression in statistics
 - Although a model for classification rather than for regression
- Logit function:
 - It is the log of the odds ratio
 - It links the probability to the predictor variables



Srihari

Fewer Parameters in Linear Discriminative Model

- Discriminative approach (Logistic Regression)
 - For *M* -dim feature space ϕ :
 - M adjustable parameters
- Generative based on Gaussians (Bayes/NB)
 - 2*M* parameters for mean
 - *M*(*M*+1)/2 parameters for shared covariance matrix
 - Two class priors
 - Total of M(M+5)/2 + 1 parameters
 - Grows quadratically with M
 - If features assumed independent (naïve Bayes) still 5 needs *M*+3 parameters

Determining Logistic Regression parameters

- Maximum Likelihood Approach for Two classes
- For a data set (ϕ_n, t_n) where $t_n \in \{0,1\}$ and $\phi_n = \phi(x_n), n = 1,..,N$
- Likelihood function can be written as

$$p(t \mid \boldsymbol{w}) = \prod_{n=1}^{N} y_n^{t_n} \left\{ 1 - y_n \right\}^{1 - t_n}$$

where t =(t_1,...,t_N)^T and y_n = p(C_1 \mid \phi_n)

 y_n is the probability that $t_n = 1$

Machine Learning

Error Fn for Logistic Regression

Likelihood function is

$$p(\mathbf{t} \mid \boldsymbol{w}) = \prod_{n=1}^{N} y_n^{t_n} \left\{ 1 - y_n \right\}^{1-t_n}$$

• By taking negative logarithm we get the Cross-entropy Error Function

$$E(\boldsymbol{w}) = -\ln p(t \mid \boldsymbol{w}) = -\sum_{n=1}^{N} \left\{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right\}$$

where $y_n = \sigma(a_n)$ and $a_n = w^T \phi_n$

• We need to minimize E(w)At its minimum, derivative of E(w) is zero So we need to solve for w in the equation

$$\nabla E(\boldsymbol{w}) = 0$$

7

Srihari

What is Cross-entropy?

• Entropy of p(x) is defined as

$$H(p) = -\sum_{x} p(x) \log p(x)$$

- If p(x=1|t)=t and p(x=0|t)=1-t then we can write $p(x)=t^{x}(1-t)^{1-x}$

- Then Entropy of p(x) is $H(p)=t \log t+(1-t)\log(1-t)$
- Cross entropy of p(x) and q(x) is defined as

$$H(p,q) = -\sum p(x)\log q(x)$$

- If q(x=1|y)=y then $H(p,q)=t \log y + (1-t)\log(1-y)$

- In general $H(p,q)=H(p)+D_{KL}(p||q)$
 - where

$$D_{KL}(p,q) = -\sum_{x} p(x) \log \frac{p(x)}{q(x)}$$

Gradient of Error Function

Error function

$$E(\boldsymbol{w}) = -\ln p(t \mid \boldsymbol{w}) = -\sum_{n=1}^{N} \left\{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right\}$$

where $y_n = \sigma(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}_n)$



Using Derivative of logistic sigmoid $d\sigma/da=\sigma(1-\sigma)$ Gradient of the error function

$$\nabla E(\boldsymbol{w}) = \sum_{n=1}^{N} (y_n - t_n) \phi_n$$

Error x Feature Vector

Contribution to gradient by data point *n* is error between target t_n and prediction $y_n = \sigma (w^T \varphi_n)$ times basis φ_n

Proof of gradient expression
Let
$$z = z_1 + z_2$$

where $z_1 = t \ln \sigma(w\phi)$ and $z_2 = (1 - t) \ln[1 - \sigma(w\phi)]$
 $\frac{dz_1}{dw} = \frac{t\sigma(w\phi)[1 - \sigma(w\phi)]\phi}{\sigma(w\phi)}$ $\frac{d\sigma}{da} = \sigma(1 - \sigma)$
and $Using$ $\frac{d}{dx}(\ln ax) = \frac{a}{x}$
 $\frac{dz_2}{dw} = \frac{(1 - t)\sigma(w\phi)[1 - \sigma(w\phi)](-\phi)}{[1 - \sigma(w\phi)]}$ 9
 n
Therefore $\frac{dz}{dw} = (\sigma(w\phi) - t)\phi$

Simple Sequential Algorithm

- Given Gradient of error function
- $\nabla E(\boldsymbol{w}) = \sum_{n=1}^{N} (y_n t_n) \phi_n \quad \text{where } y_n = \sigma(\boldsymbol{w}^{\mathrm{T}} \boldsymbol{\phi}_n)$ Solve using an iterative approach

$$\boldsymbol{w}^{\tau+1} = \boldsymbol{w}^{\tau} - \eta \nabla E_n$$



where

$$\nabla E_n = (y_n - t_n)\phi_n$$

Error x Feature Vector

Takes precisely same form as Gradient of Sum-of-squares error for linear regression

Samples are presented one at a time in which each each of the weight vectors is updated 10



Logistic Regression Code in Python

use sci-kit learn to create a data set.



Prediction: From the code above, you find p. It will be between 0 a

Srihari

 $\sigma(a)$

ML solution can over-fit

- Severe over-fitting for linearly separable data
 - Because ML solution occurs at $\sigma = 0.5$
 - With σ >0.5 and σ < 0.5 for the two classes
 - Solution equivalent to $a = w^T \phi = 0$
 - Logistic sigmoid becomes infinitely steep
 - A Heavyside step function ||*w*||goes to ∞
 - Solution

Machine Learning

- Penalizing wts
- Recall in linear regression





5



An Example of 2-class Logistic Regression

Input Data

C1 =		C2 =	
3.7000	3.0000	1.1000	-1.3000
3.2000	3.5000	0.2000	0.5000
3.5000	5.0000	1.5000	0.3000
3.6000	4.5000	0.2000	0.8000
3.5000	4.0000	1.0000	0.7000



 $\phi_0(\mathbf{x})=1$, dummy feature

Initial Weight Vector, Gradient and Hessian (2-class)

Weight vector

Gradient



Hessian

H =		
3.5000	5.3750	5.2500
5.3750	17.4825	17.4950
5.2500	17.4950	22.4150

Final Weight Vector, Gradient and Hessian (2-class)



Number of iterations : 10

Error (Initial and Final): 15.0642, 1.0000e-009

Generative vs Discriminative Training

Variables $\mathbf{x} = \{x_1, .., x_M\}$ and classifier target y

1. Generative: estimate parameters of variables independently



For classification: Determine joint: $p(y, x) = p(y) \prod_{i=1}^{M} p(x_i | y)$

From joint get required conditional p(y|x)

Simple estimation

independently estimate M sets of parameters But independence is usually false We can estimate M(M+1)/2 covariance matrix

2. Discriminative: estimate joint parameters w_i

Potential Functions (log-linear) $\phi_i(x_i,y) = \exp\{w_i x_i \ I\{y=1\}\}, \phi_0(y) = \exp\{w_0 \ I\{y=1\}\}\}$ Naïve Markov y I has value 1 when y=1, else 0 x_1 x_2 x_M



Logistic Regression is a special architecture of a neural network



Srihari



@TessFerrandez 19