

Bayesian Logistic Regression

Sargur N. Srihari

University at Buffalo, State University of New York
USA

Topics in Linear Models for Classification

- Overview
 1. Discriminant Functions
 2. Probabilistic Generative Models
 3. Probabilistic Discriminative Models
 4. The Laplace Approximation
 5. Bayesian Logistic Regression

Topics in Bayesian Logistic Regression

- Recap of Logistic Regression
- Roadmap of Bayesian Logistic Regression
- Laplace Approximation
- Evaluation of posterior distribution
 - Gaussian approximation
- Predictive Distribution
 - Convolution of Sigmoid and Gaussian
 - Approximate sigmoid with probit
- Variational Bayesian Logistic Regression

Recap of Logistic Regression

- Feature vector ϕ , two-classes C_1 and C_2
- A *posteriori* probability $p(C_1 | \phi)$ can be written as
$$p(C_1 | \phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$
where ϕ is a M -dimensional feature vector
 $\sigma(.)$ is the logistic sigmoid function
- Goal is to determine the M parameters
- Known as logistic regression in statistics
 - Although a model for classification rather than for regression

Determining Logistic Regression parameters

- Maximum Likelihood Approach to determine \mathbf{w}

Data set $\{\varphi_n, t_n\}$

where $t_n \in \{0,1\}$ and $\varphi_n = \varphi(\mathbf{x}_n)$, $n=1, \dots, N$

Since t_n is binary we can use Bernoulli

Let y_n be the probability that $t_n = 1$, i.e., $y_n = p(C_1 | \varphi_n)$

Denote $\mathbf{t} = (t_1, \dots, t_N)^T$

- Likelihood function associated with N observations

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

Simple sequential solution

- Error function is the negative of the log-likelihood

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

Cross-entropy error function

- No closed-form maximum likelihood solution
for determining \mathbf{w}

- Given Gradient of error function

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

- Solve using an iterative approach

- where
$$\nabla E(\mathbf{w}) = \sum_{n=1}^N \underbrace{(y_n - t_n) \phi_n}_{\text{Error x Feature Vector}}$$

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

Solution has severe over-fitting problems for linearly separable data
So use IRLS algorithm

IRLS for Logistic Regression

- Posterior probability of class C_1 is

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

- Likelihood Function

for data set $\{\phi_n, t_n\}$, $t_n \in \{0,1\}$, $\phi_n = \phi(\mathbf{x}_n)$

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

1. Error Function

Log-likelihood yields Cross-entropy

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

IRLS for Logistic Regression

2. Gradient of Error Function:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = \Phi^T (\mathbf{y} - \mathbf{t})$$

3. Hessian:

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1 - y_n) \phi_n \phi_n^T = \Phi^T \mathbf{R} \Phi$$

Φ is $N \times M$ design matrix whose n^{th} row is ϕ_n^T

\mathbf{R} is $N \times N$ diagonal matrix with elements

$$R_{nn} = y_n (1 - y_n)$$

Hessian is not constant and depends on \mathbf{w} through \mathbf{R}

Since \mathbf{H} is positive-definite (i.e., for arbitrary \mathbf{u} , $\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$)

error function is a concave function of \mathbf{w} and so has a unique minimum

IRLS for Logistic Regression

4. Newton-Raphson update:

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

Substituting $H = \Phi^T R \Phi$ and $\nabla E(\mathbf{w}) = \Phi^T (\mathbf{y} - \mathbf{t})$

$$\begin{aligned} \mathbf{w}^{(\text{new})} &= \mathbf{w}^{(\text{old})} - (\Phi^T R \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) \\ &= (\Phi^T R \Phi)^{-1} \{ \Phi \Phi^T \mathbf{w}^{(\text{old})} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\Phi^T R \Phi)^{-1} \Phi^T R \mathbf{z} \end{aligned}$$

where \mathbf{z} is a N -dimensional vector with elements

$$\mathbf{z} = \Phi \mathbf{w}^{(\text{old})} - R^{-1} (\mathbf{y} - \mathbf{t})$$

Update formula is a set of normal equations

Since Hessian depends on \mathbf{w}

Apply them iteratively each time using the new weight vector

Roadmap of Bayesian Logistic Regression

- Logistic regression is a discriminative probabilistic linear classifier: $p(C_1 | \phi) = \sigma(\mathbf{w}^T \phi)$
- Exact Bayesian inference for Logistic Regression $p(C_1 | \phi) = \int \sigma(\mathbf{w}^T \phi) p(\mathbf{w}) d\mathbf{w}$ is intractable, because:

1. Evaluation of posterior distribution $p(\mathbf{w} | \mathbf{t})$

- Needs normalization of prior $p(\mathbf{w}) = N(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$ times likelihood (a product of sigmoids) $p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$

- Solution: use Laplace approximation to get Gaussian $q(\mathbf{w})$

2. Evaluation of predictive distribution

$$p(C_1 | \phi) \approx \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w}$$

- Convolution of sigmoid and Gaussian
- Solution: Approximate Sigmoid by Probit

Laplace Approximation (summary)

- Need mode \mathbf{w}_0 of posterior distribution $p(\mathbf{w}|\mathbf{t})$
 - Done by a numerical optimization algorithm
- Fit a Gaussian centered at the mode

$$q(\mathbf{w}) = \frac{1}{W} f(\mathbf{w}) = \frac{|A|^{1/2}}{(2\pi)^{M/2}} \exp \left\{ -\frac{1}{2} (\mathbf{w} - \mathbf{w}_0)^T A (\mathbf{w} - \mathbf{w}_0) \right\}$$

$$= N(\mathbf{w} | \mathbf{w}_0, A^{-1})$$

- Needs second derivatives of log posterior $A = -\nabla \nabla \ln f(\mathbf{w})|_{\mathbf{w}=\mathbf{w}_0}$
 - Equivalent to finding Hessian matrix

$$S_N = -\nabla \nabla \ln p(\mathbf{w} | \mathbf{t}) = S_0^{-1} + \sum_{i=1}^n y_n (1 - y_n) \phi_n \phi_n^T$$

Evaluation of Posterior Distribution

- Gaussian prior

$$p(\mathbf{w}) = N(\mathbf{w} | \mathbf{m}_0, S_0)$$

– Where \mathbf{m}_0 and S_0 are hyper-parameters

- Posterior distribution

$$p(\mathbf{w} | \mathbf{t}) \propto p(\mathbf{w}) p(\mathbf{t} | \mathbf{w})$$

where $\mathbf{t} = (t_1, \dots, t_N)^T$

– Substituting

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

$$\ln p(\mathbf{w} | \mathbf{t}) = -\frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^T S_0^{-1} (\mathbf{w} - \mathbf{m}_0) + \sum_{i=1}^n (t_n \ln y_n + (1 - t_n) \ln(1 - y_n)) + \text{const}$$

• where $y_n = \sigma(\mathbf{w}^T \phi_n)$

Gaussian Approximation of Posterior

- Maximize posterior $p(\mathbf{w}|\mathbf{t})$ to give
 - MAP solution \mathbf{w}_{map}
 - Done by numerical optimization
 - Defines mean of the Gaussian
- Covariance given by
 - Inverse of matrix of 2nd derivatives of negative log-likelihood $S_N = -\nabla\nabla \ln p(\mathbf{w}|\mathbf{t}) = S_0^{-1} + \sum_{i=1}^n y_i(1 - y_i)\phi_i\phi_i^T$
- Gaussian approximation to posterior
$$q(\mathbf{w}) = N(\mathbf{w} | \mathbf{w}_{\text{map}}, S_N)$$
- Need to marginalize wrt this distribution to make predictions

Predictive Distribution

- Predictive distribution for class C_1 , given new feature vector $\phi(\mathbf{x})$
 - Obtained by marginalizing wrt posterior $p(\mathbf{w}|\mathbf{t})$

$$p(C_1 | \phi, \mathbf{t}) = \int p(C_1, \mathbf{w} | \phi, \mathbf{t}) d\mathbf{w} \quad \text{Sum rule}$$

$$= \int p(C_1 | \phi, \mathbf{t}, \mathbf{w}) p(\mathbf{w} | \mathbf{t}) d\mathbf{w} \quad \text{Product rule}$$

$$= \int p(C_1 | \phi, \mathbf{w}) p(\mathbf{w} | \mathbf{t}) d\mathbf{w} \quad \text{Given } \phi \text{ and } \mathbf{w}, C_1 \text{ is indep of } \mathbf{t}$$

$$\simeq \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} \quad \text{Approximate } p(\mathbf{w} | \mathbf{t}) \text{ by Gaussian } q(\mathbf{w})$$

corresponding probability for class C_2

$$p(C_2 | \phi, \mathbf{t}) = 1 - p(C_1 | \phi, \mathbf{t})$$

Predictive distrib. is a Convolution

$$p(C_1 | \phi, t) \simeq \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w}$$

- Function $\sigma(\mathbf{w}^T \phi)$ depends on \mathbf{w} only through its projection onto ϕ
- Denoting $a = \mathbf{w}^T \phi$ we have $\sigma(\mathbf{w}^T \phi) \simeq \int \delta(a - \mathbf{w}^T \phi) \sigma(a) da$
 - where δ is the Dirac delta function
- Thus $\int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} = \int \sigma(a) p(a) da$ where $p(a) = \int \delta(a - \mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w}$
 - Can evaluate $p(a)$ because
 - the delta function imposes a linear constraint on \mathbf{w}
 - Since $q(\mathbf{w})$ is Gaussian, its marginal is also Gaussian
 - Evaluate its mean and covariance

$$\mu_a = E[a] = \int p(a) da = \int q(\mathbf{w}) \mathbf{w}^T \phi d\mathbf{w} = \mathbf{w}_{map}^T \phi$$

$$\sigma_a^2 = \text{var}[a] = \int p(a) \{a^2 - E[a]^2\} da$$

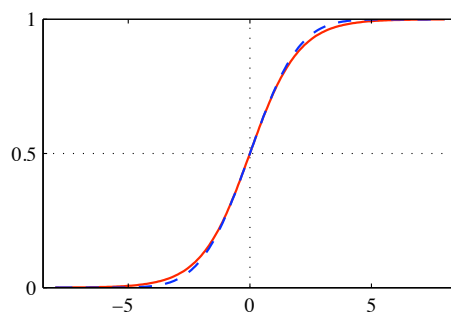
$$= \int q(\mathbf{w}) \{(\mathbf{w}^T \phi)^2 - (m_N^T \phi)^2\} d\mathbf{w} = \phi^T S_N \phi$$

Variational Approximation to Predictive Distribution

- Predictive distribution is

$$\begin{aligned} p(C_1 | \mathbf{t}) &= \int \sigma(a) p(a) da \\ &= \int \sigma(a) N(a | \mu_a, \sigma_a^2) da \end{aligned}$$

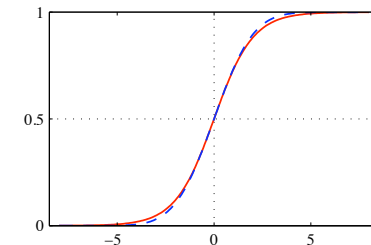
- Convolution of Sigmoid-Gaussian is intractable
- Use probit instead of logistic sigmoid



Approximation using Probit

$$p(C_1 | t) = \int \sigma(a) N(a | \mu_a, \sigma_a^2) da$$

- Use probit which is similar to Logistic sigmoid
 - Defined as $\Phi(a) = \int_{-\infty}^a N(\theta | 0, 1) d\theta$
- Approximate $\sigma(a)$ by $\Phi(\lambda a)$
 - Find λ such that two functions have same slope at origin



Approximate $\sigma(a)$ by $\Phi(\lambda a)$

Find suitable value of λ by requiring that two have same slope at origin, which yields $\lambda^2 = \pi/8$

- Convolution of probit with Gaussian is a probit

$$\int \Phi(\lambda a) N(a | \mu, \sigma^2) da = \Phi\left(\frac{\mu}{(\lambda^{-2} + \sigma^2)^{1/2}}\right)$$

– Thus

$$\begin{aligned} p(C_1 | \phi, t) &= \int \sigma(a) N(a | \mu_a, \sigma_a^2) da \\ &\simeq \sigma(\kappa(\sigma_a^2) \mu_a) \end{aligned}$$

where

$$\kappa(\sigma^2) = (1 + \pi \sigma^2 / 8)^{-1/2}$$

Probit Classification

Applying it to

$$p(C_1 | \mathbf{t}) = \int \sigma(a) N(a | \mu_a, \sigma_a^2) da$$

We have

$$p(C_1 | \phi, \mathbf{t}) = \sigma(\kappa(\sigma_a^2) \mu_a)$$

where

$$\mu_a = \mathbf{w}_{map}^T \phi \quad \sigma_a^2 = \phi^T S_N \phi$$

Decision boundary corresponding to $p(C_1 | \phi, \mathbf{t}) = 0.5$ is given by

$$\mu_a = 0$$

This is the same solution as

$$\mathbf{w}_{map}^T \phi = 0$$

Thus marginalization has no effect!

When minimizing misclassification rate with equal prior probabilities

For more complex decision criteria it plays important role₁₈

Summary

- Logistic regression is a linear probabilistic discriminative model $p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \phi)$
- Bayesian Logistic Regression is intractable
- Using Laplacian the posterior parameter distribution $p(\mathbf{w} | \mathbf{t})$ can be approximated as a Gaussian
- Predictive distribution is convolution of sigmoids and Gaussian $p(C_1 | \phi) \approx \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w}$
 - Probit yields convolution as probit