

# Bayesian Logistic Regression

Sargur N. Srihari

University at Buffalo, State University of New York  
USA

# Topics

- 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  $\phi$  , 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  $\phi$  is a  $M$ -dimensional feature vector  
 $\sigma(\cdot)$  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 for Two classes

Data set  $(\phi_n, t_n)$

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

Since  $t_n$  is binary we can use Bernoulli

Let  $y_n$  be the probability that  $t_n = 1$

- Likelihood function associated with  $N$  observations

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

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

# Simple sequential solution

- Error function is the negative of the log-likelihood

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

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

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:

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

$\Phi$  is  $N \times M$  design matrix whose  $n$ th row is  $\phi_n^T$

$\mathbf{R}$  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:

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

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

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

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

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

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 discriminative probabilistic linear classification  $p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \phi)$
  - Exact Bayesian inference for Logistic Regression is intractable
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)  $p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$
    - Apply Laplace approximation to get Gaussian  $q(\mathbf{w})$
  2. Evaluation of predictive distribution  $p(C_1 | \phi, \mathbf{t}) \approx \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w}$ 
    - Convolution of sigmoid and Gaussian
    - Approximate Sigmoid by Probit

# Laplace Approximation (summary)

- Need mode  $w_0$  of posterior distribution  $p(w|t)$ 
  - Done by a numerical optimization algorithm
- Fit a Gaussian centered at the mode

$$q(w) = \frac{1}{W} f(w) = \frac{|A|^{1/2}}{(2\pi)^{M/2}} \exp\left\{-\frac{1}{2}(w - w_0)^T A(w - w_0)\right\}$$

$$= N(w|w_0, A^{-1})$$

- Needs second derivatives of log posterior  $A = -\nabla\nabla \ln f(w)|_{w=w_0}$ 
  - Equivalent to finding Hessian matrix

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

# Evaluation of Posterior Distribution

- Gaussian prior

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

– where  $\mathbf{m}_0$  and  $\mathbf{S}_0$  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  $p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$

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

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

# Gaussian Approximation of Posterior

- Maximize posterior  $p(\mathbf{w}|\mathbf{t})$  to give
  - MAP solution  $\mathbf{w}_{\text{map}}$ 
    - Done by numerical optimization
  - Defines mean of the Gaussian
- Covariance given by
  - Inverse of matrix of 2nd derivatives of negative log-likelihood
 
$$\mathcal{S}_N = -\nabla\nabla \ln p(\mathbf{w}|\mathbf{t}) = \mathcal{S}_0^{-1} + \sum_{i=1}^n y_n(1 - y_n)\phi_n\phi_n^T$$
- Gaussian approximation to posterior
 
$$q(\mathbf{w}) = N(\mathbf{w} | \mathbf{w}_{\text{map}}, \mathcal{S}_N)$$
- 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  $p(\mathbf{w}|\mathbf{t})$

$$\begin{aligned}
 p(C_1 | \phi, \mathbf{t}) &= \int p(C_1, \mathbf{w} | \phi, \mathbf{t}) d\mathbf{w} \quad \text{Sum rule} \\
 &= \int p(C_1 | \phi, \mathbf{t}, \mathbf{w}) p(\mathbf{w} | \mathbf{t}) d\mathbf{w} \quad \text{Product rule} \\
 &= \int p(C_1 | \phi, \mathbf{w}) p(\mathbf{w} | \mathbf{t}) d\mathbf{w} \quad \text{Given } \phi \text{ and } \mathbf{w}, C_1 \text{ is indep of } \mathbf{t} \\
 &\approx \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w} \quad \text{Approximate } p(\mathbf{w} | \mathbf{t}) \text{ by Gaussian } q(\mathbf{w})
 \end{aligned}$$

corresponding probability for class  $C_2$

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

# Predictive distrib. is a Convolution

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

- Function  $\sigma(\mathbf{w}^T \phi)$  depends on  $\mathbf{w}$  only through its projection onto  $\phi$
- Denoting  $a = \mathbf{w}^T \phi$  we have  $\sigma(\mathbf{w}^T \phi) \approx \int \delta(a - \mathbf{w}^T \phi) \sigma(a) da$ 
  - where  $\delta$  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

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

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

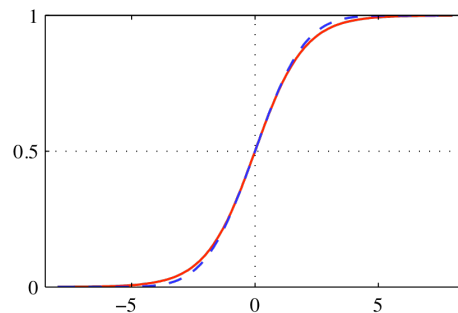
$$= \int q(\mathbf{w}) \{(\mathbf{w}^T \phi)^2 - (\mathbf{m}_N^T \phi)^2\} d\mathbf{w} = \phi^T S_N \phi$$

# Variational Approximation

- Predictive distribution is

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

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



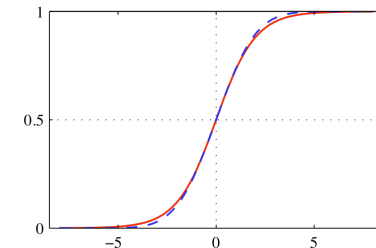
# Approximation using Probit

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

- 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  $\lambda$  such that two functions have same slope at origin

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

Find suitable value of  $\lambda$  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

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

where

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

# Probit Classification

Applying it to

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

We have

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

where

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

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

$$\mu_a = 0$$

This is the same solution as

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

Thus marginalization has no effect!

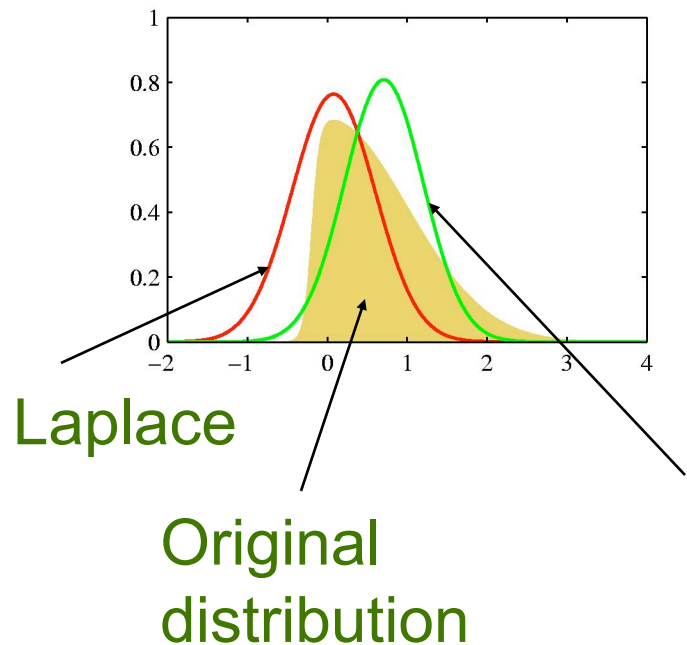
When minimizing misclassification rate with equal prior probabilities

For more complex decision criteria it plays important role<sub>17</sub>

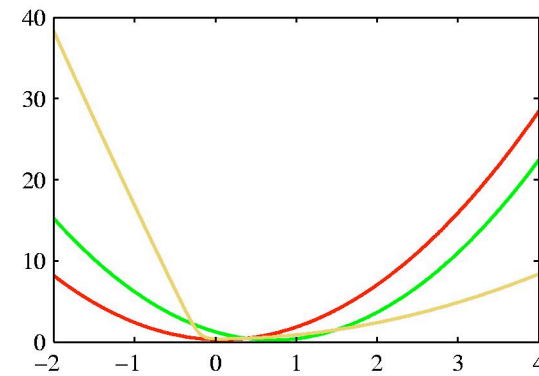
# Variational Logistic Regression

- Variational: Based on calculus of variations
  - How a derivative of a functional changes
    - Functional takes function as input and returns a value
- Approximation is more flexible than Laplace with additional variational parameters

# Essence of Variational Method



Variational distribution is Gaussian  
Optimized with respect to mean and variance



Negative Logarithms

# Inference Problem

- Observed Variables

$$X = \{x_1, \dots, x_N\} \text{ N i.i.d. data}$$

- Latent Variables and Parameters

$$Z = \{z_1, \dots, z_N\}$$

- Joint distribution  $p(X, Z)$  is specified

– E.g., full set of tables or pdfs given

- Goal is to find approximation for posterior distribution  $p(Z|X)$

# Decomposition of Log Marginal Probability

$$\ln p(X) = L(q) + KL(q \parallel p)$$

where

$$L(q) = \int q(Z) \ln \left\{ \frac{p(X, Z)}{q(Z)} \right\} dZ \quad \text{A functional}$$

and

$$KL\{q \parallel p\} = - \int q(Z) \ln \left\{ \frac{p(Z | X)}{q(Z)} \right\} dZ$$

← Also applicable to discrete distributions  
By replacing integrations with summations

← Kullback-Leibler Divergence

## Some Observations:

Lower bound on  $\ln p(X)$  is  $L(q)$

Maximizing the lower bound  $L(q)$  wrt distribution  $q(Z)$  is equivalent to minimizing KL Divergence

When KL divergence vanishes  $q(Z)$  equals the posterior  $p(Z|X)$

## Plan:

We seek that distribution  $q(Z)$  for which  $L(q)$  is largest

Since true posterior is intractable we consider restricted family for  $q(Z)_1$

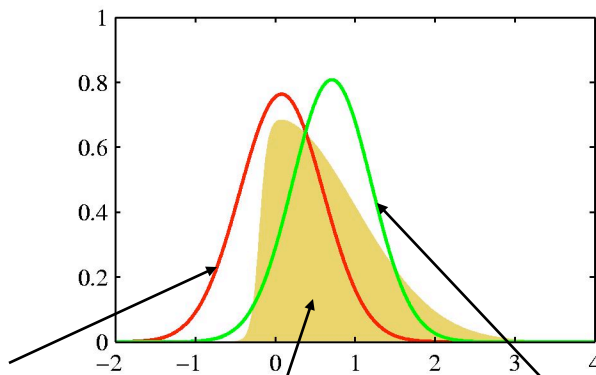
Seek member of this family for which KL divergence is minimized

# Variational Approximation Example

Use parametric distribution  $q(Z|w)$

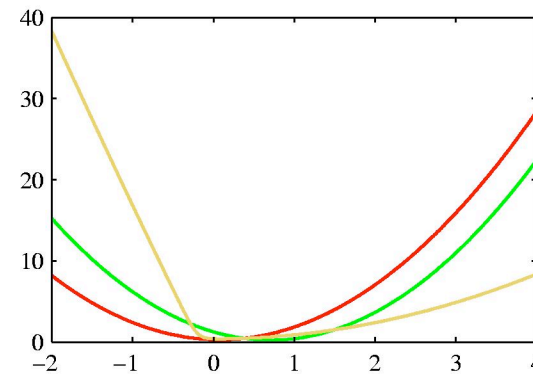
Lower bound  $L(q)$  is a function of  $w$  and

can use standard nonlinear optimization to determine optimal  $w$



Laplace

Original  
distribution



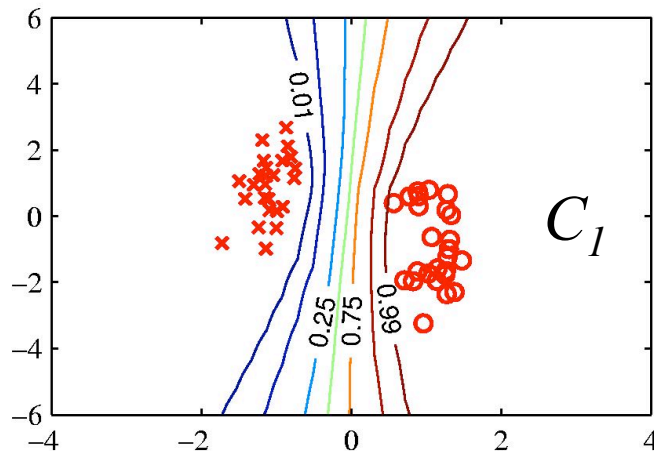
Negative  
Logarithms

Variational distribution is Gaussian  
Optimized with respect to mean and variance

# Example of Variational Bayes Logistic Regression

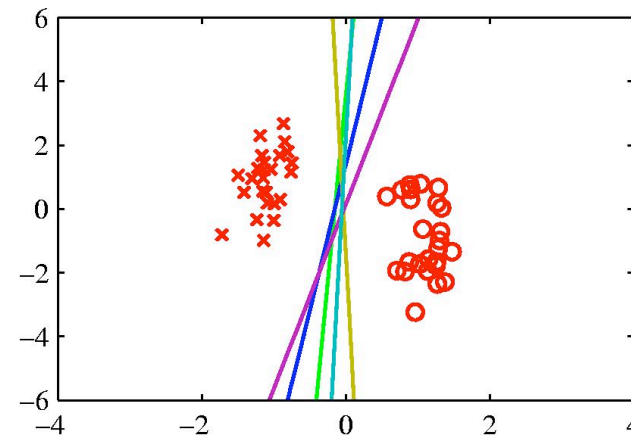
Linearly separable data set

Predictive distribution



$p(C_1|\phi)$  is plotted  
from 0.01 to 0.99

Decision boundaries



Corresponding to five samples of  
parameter vector  $w$

Large margin solution of SVM has qualitatively similar  
result to the Bayesian solution 23

# 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  $p(C_1 | \phi, \mathbf{t}) \approx \int \sigma(\mathbf{w}^T \phi) q(\mathbf{w}) d\mathbf{w}$ 
  - Probit yields convolution as probit
- Variational approach is more flexible