

# Variational 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
  - 5.1 Bayesian Logistic Regression
  - 5.2 Variational Bayesian Logistic Regression

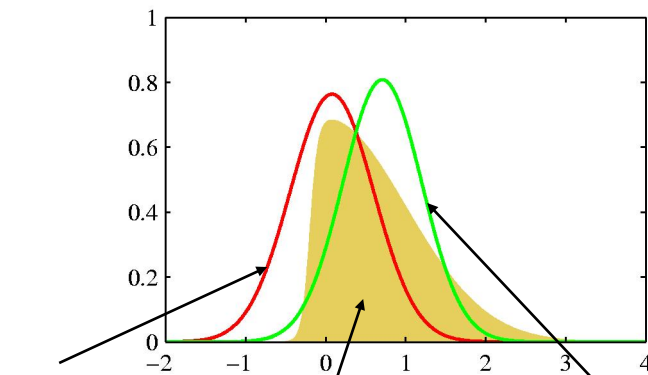
# Topics in Variational Bayesian Logistic Regression

- Bayesian Logistic Regression Posterior
- K-L Divergence
- Variational Inference
- Stochastic Variational Inference

# 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

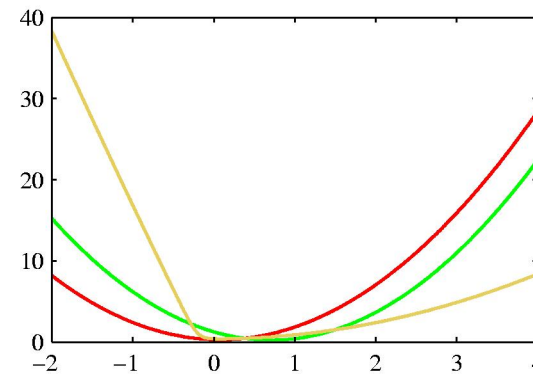
# Example of Variational Method



Laplace

Original  
distribution

Variational distribution is Gaussian  
Optimized with respect to mean and variance



Negative  
Logarithms

# Bayesian Logistic Regression Posterior

- **Data:**  $D = \{(\mathbf{x}^{(1)}, t^1), \dots, (\mathbf{x}^{(N)}, t^N)\}$   $N$  i.i.d. samples  $t \in \{0, 1\}$
- **Parameters:**  $\mathbf{w} = \{w_1, \dots, w_M\}$

– Probabilistic model specifies the joint distribution

$$p(D, \mathbf{w}) = p(D | \mathbf{w}) p(\mathbf{w})$$

- Which is a product of sigmoid likelihood and a Gaussian:

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

$$y_n = p(C_1 | \mathbf{x}^{(n)}) = \sigma(\mathbf{w}^T \mathbf{x}^{(n)})$$

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

– Goal is to approximate posterior distribution  $p(\mathbf{w} | D)$

- which is also a normalized product of sigmoid likelihood and Gaussian

$$p(\mathbf{w} | D) = \frac{p(D, \mathbf{w})}{p(D)} = \frac{p(D | \mathbf{w}) p(\mathbf{w})}{p(D)}$$

– with a Gaussian  $q(\mathbf{w})$ , so we can use it in prediction

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

# Decomposition of Log Marginal Probability

$$\ln p(D) = L(q) + KL(q || p)$$

where

$$L(q) = \int q(\mathbf{w}) \ln \left\{ \frac{p(D, \mathbf{w})}{q(\mathbf{w})} \right\} d\mathbf{w}$$

A functional

and

$$KL\{q || p\} = - \int q(\mathbf{w}) \ln \left\{ \frac{p(\mathbf{w} | D)}{q(\mathbf{w})} \right\} d\mathbf{w}$$

We wish to maximize  $\ln p(D)$  called evidence, by suitable choice of  $q$

We want to minimize K-L Divergence over a family for  $q(\mathbf{w})$

Which finds distribution  $q(\mathbf{w})$  that best approximates  $p(\mathbf{w}|D)$

## Some Observations:

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

Maximizing the lower bound  $L(q)$  wrt distribution  $q(\mathbf{w})$  is equivalent to minimizing KL Divergence

When KL divergence vanishes  $q(\mathbf{w})$  equals the posterior  $p(\mathbf{w}|D)$

## Plan:

We seek that distribution  $q(\mathbf{w})$  for which  $L(q)$  is largest

We consider restricted family for  $q(\mathbf{w})$

Seek member of this family for which KL divergence is minimized

# Variational Inference

- We want to form our Gaussian approximation  $q(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = N(\boldsymbol{w}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$  to the posterior  $p(\boldsymbol{w} | D)$  by minimizing the KL divergence

$$D_{KL} \left( q(\boldsymbol{w}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \parallel p(\boldsymbol{w} | D) \right) = E_{N(\boldsymbol{w}; \boldsymbol{\mu}, \boldsymbol{\Sigma})} \left[ \log \frac{N(\boldsymbol{w}; \boldsymbol{\mu}, \boldsymbol{\Sigma})}{p(\boldsymbol{w} | D)} \right]$$

- We fit the variational parameters  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , not the original parameters  $\boldsymbol{w}$ 
  - Although there is an interpretation that  $\boldsymbol{\mu}$  is an estimate of  $\boldsymbol{w}$ , while  $\boldsymbol{\Sigma}$  indicates a credible range of where the parameters could plausibly be around this estimate.



# Criterion to Minimize

- As we can only evaluate the posterior up to a constant, we write

$$D_{KL}(q || p) = E_{N(\mathbf{w}; \mu, \Sigma)} \left[ \underbrace{\log \frac{N(\mathbf{w}; \mu, \Sigma)}{p(\mathbf{w})p(D|\mathbf{w})}}_J + \log p(D) \right]$$

where we have used  
 $p(\mathbf{w}|D) = p(\mathbf{w}, D) / p(D)$

- And minimize  $J$
- Because  $D_{KL} \geq 0$  we obtain a lower bound on the model likelihood  $\log p(D) \geq -J$ ,
  - $-J$  is called evidence lower bound (ELBO)
- For logistic regression model and prior

$$J = E_{N(\mathbf{w}; \mu, \Sigma)} \left[ \log N(\mathbf{w}; \mu, \Sigma) - \log p(\mathbf{w}) - \sum_{n=1}^N \sigma(\mathbf{w}^T \mathbf{x}^{(n)}) \right]$$

which we wish to minimize

# Method of Optimization

$$J = E_{N(\mathbf{w}; \mu, \Sigma)} \left[ \log N(\mathbf{w}; \mu, \Sigma) - \log p(\mathbf{w}) - \sum_{n=1}^N \log \sigma(\mathbf{w}^T \mathbf{x}^{(n)}) \right]$$

- We could evaluate this expression numerically
  - The first two expectations can be computed analytically, and the remaining  $N$  terms in the sum can be reduced to a one-dimensional integral
- We could similarly evaluate the derivatives wrt  $\mu$  and  $\Sigma$ , and fit the variational parameters with a gradient-based optimizer
  - However the inner loop each function evaluation would require  $N$  numerical integrations, or further approximation (since it expectation of a summation)<sup>10</sup>

# Stochastic Variational Inference

- We can avoid numerical integration by a simple Monte Carlo estimate of the expectation

$$J \approx \frac{1}{S} \sum_s \left[ \log N(\mathbf{w}^{(s)}; \mu, \Sigma) - \log p(\mathbf{w}^{(s)}) - \sum_{n=1}^N \log \sigma(\mathbf{w}^{(s)T} \mathbf{x}^{(n)}) \right], \mathbf{w}^{(s)} \sim N(\mu, \Sigma)$$

- Really cheap estimate is to use one sample
  - Further form an unbiased estimate of the sum over data points by randomly picking one term

$$J \approx \left[ \log N(\mathbf{w}^{(s)}; \mu, \Sigma) - \log p(\mathbf{w}^{(s)}) - N \log \sigma(\mathbf{w}^{(s)T} \mathbf{x}^{(n)}) \right], n \sim \text{Uniform}[1, ..N], \mathbf{w}^{(s)} \sim N(\mu, \Sigma)$$

# Removing dependence on $\mu$ and $\Sigma$

- The goal is to move the variational parameters  $\mu$  and  $\Sigma$ , which express which weights are plausible, so that cost  $J$  gets smaller (on average)
  - If the variational parameters change, the samples we would draw change, which complicates our reasoning.
- We can remove the variational parameters from the random draws by writing down how a Gaussian random generator works:  $\mathbf{w}^{(s)} = (\mu + \Sigma^{1/2}\mathbf{v})$ ,
  - where  $\mathbf{v} \sim N(0,1)$ ,  $\Sigma^{1/2}$  is a matrix square root, such as Cholesky decomposition. Then our cost function becomes
 

$$J \approx \log N((\mu + \Sigma^{1/2}\nu); \mu, \Sigma) - \log p(\mathbf{w} = \mu + \Sigma^{1/2}\nu) - N \log \sigma(\mu + \Sigma^{1/2}\nu)^T \mathbf{x}^{(n)},$$

$$n \sim \text{Uniform}[1, ..N], \nu \sim N(0, I)$$
  - Where all dependence on variational parameters  $\mu$  and  $\Sigma$  are now in the estimator itself

# Minimizing cost function $J$

- We can minimize the cost function  $J$  by stochastic gradient descent. We draw a datapoint  $n$  at random, some Gaussian white noise  $v$ , and evaluate  $J$  and its derivatives wrt  $\mu$  and  $\Sigma$
- We then make a small change to  $\mu$  and  $\Sigma$  to improve the current estimate of  $J$  and repeat.

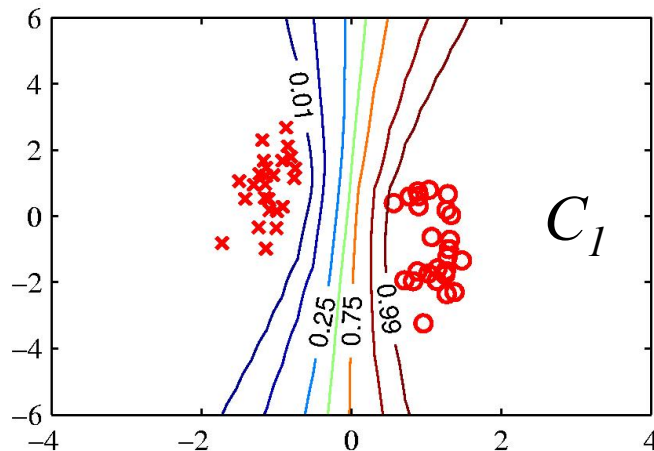
# Applicability to neural networks

- While some of the fine details are slightly complicated, none of them depend on the fact we were considering logistic regression.
- We could easily replace the logistic function  $\sigma(\cdot)$  with another model likelihood, such as from a neural network. As long as we can differentiate the log-likelihood, we can apply stochastic variational inference.

# Example of Variational Bayes Logistic Regression

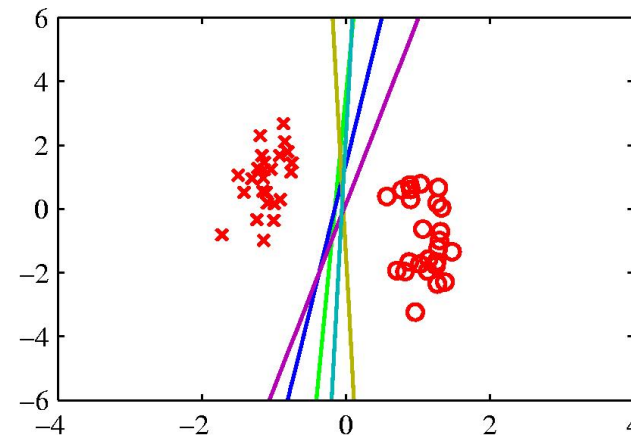
Linearly separable data set

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