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

Srihari

Example of Variational Method



Bayesian Logistic Regression Posterior

- **Data:** $D = \{(x^{(1)}, t^{(1)}), ..., (x^{(N)}, t^{(N)})\}$ N i.i.d.samples $t \in \{0, 1\}$
- **Parameters:** $w = \{w_1, ..., w_M\}$
 - Probabilistic model specifies the joint distribution
 p(D,w) =p (D/w)p(w)
 - Which is a product of sigmoid likelihood and a Gaussian:

 $p(D \mid \boldsymbol{w}) = \prod_{n=1}^{N} y_n^{t_n} \{1 - y_n\}^{1 - t_n} \quad y_n = p(C_1 \mid \boldsymbol{x}^{(n)}) = \sigma(\boldsymbol{w}^T \boldsymbol{x}^{(n)}) \quad p(\boldsymbol{w}) = N(\boldsymbol{w} \mid \mathbf{m}_0, \mathbf{S}_0)$

- Goal is to approximate posterior distribution p(w|D)
 - which is also a normalized product of sigmoid likelihood and Gaussian $p(w | D) = \frac{p(D,w)}{p(D)} = \frac{p(D/w)p(w)}{p(D)}$

- with a Gaussian q(w), so we can use it in prediction

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



Srihari

Empirical distribution and 2 approximations

Observations: Worm teeth





 $X = \{x_1, x_2, ... x_{10}\}$ $P(X) = \{p(x_i)\}$ $\sum_{i=1}^{10} p(x_i) = 1$

Uniform distribution: assign prob 1/11 to 11 values This approximation wipes out nuances in the data



Which one is the better approximation?

Binomial distribution: E(x)=np; with n=10, mean=5.7 implies that p=0.57Approximation has more subtleity, but is not perfect



https://www.countbayesie.com/blog/2017/5/9/kullback-leibler-divergence-explained

7

K-L Divergence between distributions

- Entropy of our distribution:
 - $H[p] = -\sum_{i=1}^{N} p(x_i) \cdot \log p(x_i)$
 - If we take log to base 2 it is the minimum no of bits needed to encode the information
 - E.g., one bit needed for $p(x_1) = p(x_2) = 1/2$
- K-L divergence: a slight modification of entropy $D_{KL}(p||q) = \sum_{i=1}^{N} p(x_i) \cdot (\log p(x_i) - \log q(x_i))$ $= E [\log p(x) - \log q(x)]$
 - It is the no of bits we expect to lose in encoding p by q

KL for the two approximations

- $D_{KL}(Observed||Uniform)=0.338$ $D_{KL}(Observed||Binomial)=0.477$
- Information lost from binomial approximation is more than with uniform
 - So we are better off using a Uniform approximation
- KL divergence is not a distance measure since it is not symmetric
 - $-E.g., D_{KL}(Binomial||Observed)=0.330$

Optimizing using KL Divergence



KL for our empirical distribution and Binomial distribution (n,p)

Optimal value of p marked by dot is the best approximation Although we didn't try to minimize KL

An ad-hoc example

Distribution is defined as follows: Interval [6,10]: $p(x_i)=p/5$ Interval [0,5] $p(x_i)=(1-p)/6$



The minimum KL divergence of 0.338 occurs when p=0.47.

Since this distribution is close to the uniform distribution we don't save anything Key point is that KL divergence can be used as an objective function to find the optimal value for any approximating distribution While this example is optimizing a single parameter, we can extend it to a model with many parameters

Function versus Functional

- Function takes a value of a variable as input and returns the function value as output
- Functional takes a function as input and returns the functional value as output
- An example is entropy

$$H[p(x)] = \int p(x) \ln p(x) \, dx$$

Calculus of Variations

- *Function*: input *x* is a value, returns a value y(x)
- *Functional*: input is a function y(x), returns a value F[y(x)]
- Standard Calculus concerns Derivatives of functions
 - How output changes in response to small changes in input
- Variational Calculus: concerns Functional derivative
 - How output changes in response to infinitesimal changes in input function



Functional Derivative

- Defined by considering:
 - How the value of a functional F[y] changes when the function y(x) is changed to $y(x) + \epsilon \eta(x)$
 - where $\eta(x)$ is an arbitrary function of x



Maximizing a Functional

- Common problem in conventional calculus:
 - Find a value x that maximizes a function y(x)
- In calculus of variations find a function y(x) that maximizes a functional F[y]

– E.g., Find a function that maximizes entropy

Can be used to show that:

1.Shortest path between two points is a straight line

2.Distribution that maximizes entropy $H[p(x)] = \int p(x) \ln p(x) dx \quad \text{is Gaussian}$

Variational methods

- Nothing intrinsically approximate
- But naturally lend themselves to approximation
 - By restricting range of functions, e.g.,
 - Quadratic
 - Linear combination of fixed basis functions
 - In which only the coefficients of the linear combination vary
 - Factorization assumptions
 - Method used for probabilistic inference

Application to Probabilistic Inference

- Consider a fully Bayesian model
- All parameters are given prior distributions
- Model may have latent variables as well as parameters
 - Denote all parameters and latent variables by \boldsymbol{Z}
 - Denote set of observed variables by X
 - Given *N* iid samples for which $X = \{x_1, .., x_N\}$ and $Z = \{z_1, .., z_N\}$, probabilistic model specifies posterior distribution p(Z|X) as well as model evidence p(X)17

Machine Learning Decomposition of Log Marginal Probability

$$\ln p(D) = L(q) + KL(q || p)$$
where
$$L(q) = \int q(w) \ln \left\{ \frac{p(D, w)}{q(w)} \right\} dw$$
A functional
and
$$KL\{q || p\} = -\int q(w) \ln \left\{ \frac{p(w | D)}{q(w)} \right\} dw$$
We want to minimize K-L
Divergence over a family
for $q(w)$
Which finds distribution
 $q(w)$ that best

Some Observations:

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

Maximizing the lower bound L(q) wrt distribution q(w)

is equivalent to *minimizing* KL Divergence

When KL divergence vanishes q(w) equals the posterior p(w|D)Plan:

We seek that distribution q(w) for which L(q) is largestWe consider restricted family for q(w)Seek member of this family for which KL divergence is minimized

approximates p(w|D)

Variational Inference

 We want to form our Gaussian approximation q(μ, Σ)=N(w; μ,Σ) to the posterior p(w|D) by minimizing the KL divergence

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

- We fit the variational parameters μ and Σ , not the original parameters w
 - Although there is an interpretation that μ is an estimate of *w*, while Σ 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

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

$$J$$

$$J$$

$$J$$

$$J$$

$$J$$

where we have used p(w|D)=p(w,D)/p(D)

• Because $D_{\text{KL}} \ge 0$ we obtain a lower bound on the model likelihood $\log p(D) \ge -J$,

-J is called evidence lower bound (ELBO)

• For logistic regression model and prior

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

which we wish to minimize

Method of Optimization

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

- 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 μ and Σ , 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)²¹

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

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

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

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

Machine Learning Removing dependence on μ and Σ

- The goal is to move the variational parameters μ and Σ, 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: $w^{(s)} = (\mu + \Sigma^{\frac{1}{2}}v)$,
 - where $v \sim N(0,1), \Sigma^{\frac{1}{2}}$ is a matrix square root, such as Cholesky decomposition. Then our cost function becomes

$$J \approx \log N((\mu + \Sigma^{1/2}\nu); \mu, \Sigma) - \log p(\boldsymbol{w} = \mu + \Sigma^{1/2}\nu) - N \log \sigma(\mu + \Sigma^{1/2}\nu)^T \boldsymbol{x}^{(n)}),$$
$$n \sim \text{Uniform}[1, ..N], \nu \sim N(0, I)$$

- Where all dependence on variational parameters μ and Σ 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 ν, and evaluate J and its derivatives wrt μ and Σ
- We then make a small change to μ and Σ 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 σ(.) 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 ₂₆ result to the Bayesian solution