# CS181 Practice Questions: Linear Regression, Continued

1. **Posterior Weight Distribution Using Bayes' Rule for Linear Gaussian Systems**

   **Some background:** In section (2.3.3), Bishop derives the following facts about linear Gaussian systems: assuming we have a marginal distribution on $x$ and a conditional distribution on $y$ given by

   $$p(x) = \mathcal{N}(x \mid \mu, \Lambda^{-1}) \tag{2.113}$$
   $$p(y \mid x) = \mathcal{N}(y \mid Ax + b, L^{-1}) \tag{2.114}$$

   then

   $$p(x \mid y) = \mathcal{N}(x \mid \Sigma(A^\mathsf{T} L(y - b) + \Lambda\mu), \Sigma) \tag{2.116}$$

   where

   $$\Sigma = (\Lambda + A^\mathsf{T} LA)^{-1}. \tag{2.117}$$

   Now, we know from (3.10) in Bishop that the regression likelihood can be written as

   $$p(t \mid w) = \prod_{n=1}^{N} \mathcal{N}(t_n \mid w^\mathsf{T}\phi(x_n), \beta^{-1}),$$

   where $\beta = \frac{1}{\sigma^2}$. If the prior distribution on $w$ is given by $p(w) = \mathcal{N}(w \mid m_0, S_0)$, derive that the posterior distribution $p(w \mid t)$ is given by

   $$p(w \mid t) = \mathcal{N}(w \mid m_N, S_N)$$

   where

   $$m_N = S_N(S_0^{-1}m_0 + \beta\Phi^\mathsf{T} t)$$
   $$S_N^{-1} = S_0^{-1} + \beta\Phi^\mathsf{T}\Phi$$

   in the following way:

   (a) Write down the likelihood in the form of a multivariate Gaussian.

   (b) Explain what we need to substitute for $x, y, \mu, \Lambda, A, b, L$ (respectively) in equations (2.113)-(2.117) to derive the posterior.

2. **Posterior Weight Distribution By Completing the Square (Bishop 3.7)**

We know from (3.10) in Bishop that the regression likelihood can be written as

$$p(t \mid w) = \prod_{n=1}^{N} \mathcal{N}(t_n \mid w^\mathsf{T} \phi(x_n), \beta^{-1})$$
$$\propto \exp\left( -\frac{\beta}{2}(t - \Phi w)^\mathsf{T}(t - \Phi w) \right),$$

where $\beta = \frac{1}{\sigma^2}$ and in the second line above we have ignored the Gaussian normalizing constants. By completing the square, show that with a prior distribution on $w$ given by $p(w) = \mathcal{N}(w \mid m_0, S_0)$, the posterior distribution $p(w \mid t)$ is given by

$$p(w \mid t) = \mathcal{N}(w \mid m_N, S_N)$$

where

$$m_N = S_N(S_0^{-1} m_0 + \beta \Phi^\mathsf{T} t)$$
$$S_N^{-1} = S_0^{-1} + \beta \Phi^\mathsf{T} \Phi$$

3. **Predictive Distribution**

Bishop notes in section (3.2.2) that if our prior distribution on $w$ is

$$p(w) = \mathcal{N}(w \mid 0, \alpha^{-1}I),$$

and if we assume again that our likelihood involves an inverse variance parameter $\beta$, then the predictive distribution of $t$ for a new datapoint $x$ is given by

$$p(t \mid t, \alpha, \beta) = \int p(t \mid w, \beta) p(w \mid t, \alpha, \beta) \, dw \qquad (3.57)$$

Does the equation in (3.57) make any independence assumptions about the variables involved? If so, which?

4. **Deriving Lasso Regularization with Lagrange Multipliers**

Show that minimization of the unregularized sum-of-squares error function given by

$$E_D(\boldsymbol{w}) = \frac{1}{2} \sum_{n=1}^{N} (t_n - \boldsymbol{w}^\mathsf{T} \boldsymbol{\phi}(\boldsymbol{x}_n))^2,$$

subject to the constraint

$$\sum_{j=1}^{M} |w_j| \leq \eta,$$

is equivalent to minimizing the regularized error function

$$\frac{1}{2} \sum_{n=1}^{N} (t_n - \boldsymbol{w}^\mathsf{T} \boldsymbol{\phi}(\boldsymbol{x}_n))^2 + \frac{\lambda}{2} \sum_{j=1}^{M} |w_j|$$

5. **Connection between Priors and Regularization**

Consider the Bayesian linear regression model given in Bishop 3.3.1. The prior is given by

$$p(w|\alpha) = \mathcal{N}(w|0, \alpha^{-1}I),$$

where $\alpha$ is the precision parameter that controls the variance of the Gaussian prior. The likelihood can be written as

$$p(t\,|\,w) = \prod_{n=1}^{N} \mathcal{N}(t_n\,|\,w^{\mathsf{T}}\phi(x_n), \beta^{-1}),$$

Using the fact that the posterior is the product of the prior and the likelihood, show that maximizing the log posterior (i.e. $\ln p(w\,|\,t) = \ln p(w|\alpha) + \ln p(t\,|\,w)$) is equivalent to minimizing the regularized error term given by $E_D(w) + \lambda E_W(w)$ with

$$E_D(w) = \frac{1}{2}\sum_{n=1}^{N}(t_n - w^{\mathsf{T}}\phi(x_n))^2$$

$$E_W(w) = \frac{\lambda}{2}w^{\mathsf{T}}w$$

Do this by writing $\ln p(w\,|\,t)$ as a function of $E_D(w)$ and $E_W(w)$, dropping constant terms if necessary.

Conclude that maximizing this posterior is equivalent to minimizing the regularized error term given by $E_D(w) + \lambda E_W(w)$.

(Hint: take $\lambda = \alpha/\beta$)

6. **Bayesian Updates in Linear Regression, Bishop 3.8**

Suppose we have the standard Bayesian linear regression model and we have already observed $N$ data points, so the posterior distribution is the same as the one derived in problem 2,

$$p(w \mid t) = \mathcal{N}(w \mid m_N, S_N)$$

where

$$m_N = S_N(S_0^{-1}m_0 + \beta\Phi^\mathsf{T}t)$$
$$S_N^{-1} = S_0^{-1} + \beta\Phi^\mathsf{T}\Phi$$

Suppose we observe a new data point $(x_{N+1}, t_{N+1})$. Show that the resulting posterior distribution is of the same form with $m_{N+1}$ and $S_{N+1}$.