# 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

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

   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.

(a) The likelihood can be written $p(t \mid w) = \mathcal{N}(t \mid \Phi w, \beta^{-1} I)$.

(b) Substituting $x = w, y = t, \mu = m_0, \Lambda = S_0^{-1}, A = \Phi, b = 0, L = \beta I$, then using (2.117) we first get that $\Sigma = S_N = (S_0^{-1} + \beta \Phi^\mathsf{T} \Phi)^{-1}$, and plugging into (2.116) we get that $m_N = S_N(\beta \Phi^\mathsf{T} I(t - 0) + S_0^{-1} m_0) = S_N(S_0^{-1} m_0 + \beta \Phi^\mathsf{T} t)$.

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

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

Multiplying the likelihood and prior, we get

$$p(w \mid t) \propto p(t \mid w)p(w)$$

$$\propto \exp\left(-\frac{\beta}{2}(t - \Phi w)^\mathsf{T}(t - \Phi w)\right)\exp\left(-\frac{1}{2}(w - m_0)^\mathsf{T}S_0^{-1}(w - m_0)\right) \quad (1)$$

$$= \exp\left(-\frac{1}{2}w^\mathsf{T}(S_0^{-1} + \beta\Phi^\mathsf{T}\Phi)w - \beta t^\mathsf{T}\Phi w - \beta w^\mathsf{T}\Phi^\mathsf{T}t + \beta t^\mathsf{T}t - m_0^\mathsf{T}S_0^{-1}w\right.$$

$$\tag{2}$$

$$\left. - w^\mathsf{T}S_0^{-1}m_0 + m_0^\mathsf{T}S_0^{-1}m_0\right)$$

$$= \exp\left(-\frac{1}{2}w^\mathsf{T}(S_0^{-1} + \beta\Phi^\mathsf{T}\Phi)w - (S_0^{-1}m_0 + \beta\Phi^\mathsf{T}t)^\mathsf{T}w - w^\mathsf{T}(S_0^{-1}m_0 + \beta\Phi^\mathsf{T}t)\right.$$

$$\tag{3}$$

$$\left. + \beta t^\mathsf{T}t + m_0^\mathsf{T}S_0^{-1}m_0\right)$$

$$= \exp\left(-\frac{1}{2}(w - m_N)^\mathsf{T}S_N^{-1}(w - m_N)\right)\exp\left(-\frac{1}{2}(\beta t^\mathsf{T}t + m_0^\mathsf{T}S_0^{-1}m_0 - m_N^\mathsf{T}S_N^{-1}m_N)\right),$$

$$\tag{4}$$

where the first exponential in (4) has the form of the desired (unnormalized) Gaussian, and the second exponential in (4) doesn't involve $w$, and so can be absorbed into the Gaussian's normalization constant.

To derive (4) from (3), note that plugging in the definitions of $m_N$ and $S_N^{-1}$ into the first exponential in (4) and multiplying through gives us all the terms in (3) involving $w$, except with an additional $m_N^\mathsf{T} S_N^{-1} m_N$ term. We subtract this off in the second exponential, and also collect there the other terms not involving $w$, which are absorbed into the normalization.

3. **Predictive Distribution**

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

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

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

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

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

Yes, it assumes that $t$ is independent of $t$ (and $\alpha$) given $w$. This can be seen, for instance, by writing out $p(t, w \,|\, t, \alpha, \beta)$ and using the chain rule of probability.

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

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

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

subject to the constraint

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

is equivalent to minimizing the regularized error function

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

Rewrite the constraint as

$$\sum_{j=1}^{M} |w_j| - \eta \leq 0$$

We get the Lagrangian function

$$L(w, \lambda) = \frac{1}{2} \sum_{n=1}^{N} (t_n - w^\mathsf{T} \phi(x_n))^2 + \frac{\lambda}{2} \sum_{j=1}^{M} (|w_j| - \eta)$$

where we introduce the factor of $1/2$ in front of the second term for convenience. We see immediately that the above function is equal to the regularized error function plus the terms of $\eta$ which do not depend on $w$. Therefore, minimizing the Lagrangian with respect to $w$ will give the same $w^*$ as minimizing the regularized error function.

5. **Connection between Priors and Regularization**

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

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

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

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

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

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

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

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

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

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

Expanding the posterior gives

$$\begin{aligned}
p(\boldsymbol{w}\,|\,\boldsymbol{t}) &= p(\boldsymbol{w}|\alpha) + p(\boldsymbol{t}\,|\,\boldsymbol{w}) \implies \\
\ln p(\boldsymbol{w}\,|\,\boldsymbol{t}) &= \ln p(\boldsymbol{w}|\alpha) + \ln p(\boldsymbol{t}\,|\,\boldsymbol{w}) \\
&= \sum_{n=1}^{N} \ln \mathcal{N}(t_n\,|\,\boldsymbol{w}^\mathsf{T}\boldsymbol{\phi}(\boldsymbol{x}_n), \beta^{-1}) + \ln \mathcal{N}(\boldsymbol{w}\,|\,\boldsymbol{0}, \alpha^{-1}\boldsymbol{I}) \\
&= \sum_{n=1}^{N} -\frac{\beta}{2}(t_n - \boldsymbol{w}^\mathsf{T}\boldsymbol{\phi}(\boldsymbol{x}_n))^2 + const - \frac{1}{2}\alpha\boldsymbol{w}^\mathsf{T}\boldsymbol{w} + const \\
&= -\beta E_D(\boldsymbol{w}) - \alpha E_W(\boldsymbol{w})
\end{aligned}$$

Taking $\lambda = \alpha/\beta$, note that minimizing this function is equivalent to maximizing its negative, and that multiplying by a constant positive factor $(1/\beta)$ has no effect on the extrema, and we have the result.

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 $S_{N+1}$.

We proceed by completing the square, as in problem 2.

The prior is of the form

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

and the likelihood is given by

$$p(t_{N+1} \mid x_{N+1}, w) = (\beta/2\pi)^{1/2} \exp\left(-\beta/2(t_{N+1} - w^\mathsf{T}\phi_{N+1})^2\right)$$

where $\phi_{N+1} = \phi(x_{N+1})$. Then the posterior is

$$p(w \mid t) \propto \exp\left(-\frac{1}{2}(w - m_N)^\mathsf{T} S_N^{-1}(w - m_N) - \frac{1}{2}\beta(t_{N+1} - w^\mathsf{T}\phi_{N+1})^2\right)$$

Completing the square in the exponential just as we did in problem 2, we get

$$(w - m_N)^\mathsf{T} S_N^{-1}(w - m_N) + \beta(t_{N+1} - w^\mathsf{T}\phi_{N+1})^2$$
$$= w^\mathsf{T} S_N^{-1} w - 2w^\mathsf{T} S_N^{-1} m_N + \beta w^\mathsf{T} \phi_{N+1}\phi_{N+1}^\mathsf{T} w - 2\beta w^\mathsf{T}\phi_{N+1}t_{N+1} + \text{const}$$
$$= w^\mathsf{T}(S_N^{-1} + \beta\phi_{N+1}\phi_{N+1}^\mathsf{T})w - 2w^\mathsf{T}(S_N^{-1} m_N + \beta\phi_{N+1}t_{N+1}) + \text{const}$$

Plugging this back in to the expression for the posterior we get

$$p(w \mid t) = \mathcal{N}(w \mid m_{N+1}, S_{N+1})$$

where

$$m_{N+1} = S_{N+1}(S_N^{-1} m_N + \beta\phi_{N+1}t_{N+1})$$
$$S_{N+1}^{-1} = S_N^{-1} + \beta\phi_{N+1}\phi_{N+1}^\mathsf{T}$$