# CS181 Practice Questions: Probability and Linear Regression Basics

1. **Mean of Gaussian (Bishop, 1.8, part 1)**

   By using a change of variables, verify that the univariate Gaussian distribution satisfies

   $$\mathbb{E}[x] = \int (2\pi\sigma^2)^{-1/2} \exp\{-\frac{1}{2\sigma^2}(x-\mu)^2\} x\, dx$$
   $$= \mu.$$

   Changing variables $y = x - \mu$, gives us:

   $$\mathbb{E}[x] = \int (2\pi\sigma^2)^{-1/2} \exp\{-\frac{1}{2\sigma^2}y^2\}(y+\mu)dy.$$

2. **Mode of Gaussian (Bishop, 1.9)**

Show that the mode (i.e. the maximum) of the Gaussian distribution

$$\mathcal{N}(x|\mu,\sigma^2) = (2\pi\sigma^2)^{-1/2}\exp\{-(2\sigma^2)^{-1}(x-\mu)^2\}$$

is given by $\mu$.

We differentiate with respect to $x$ to obtain:

$$\frac{d}{dx}\mathcal{N}(x|\mu,\sigma^2) = -\mathcal{N}(x|\mu,\sigma^2)\frac{x-\mu}{\sigma^2}.$$

Setting this to zero, we obtain $x = \mu$.

3. **Gaussian MLE**

Suppose we have $N$ iid values $x_n \sim \mathcal{N}(\mu, \sigma^2)$, where $n = 1, \ldots, N$.

(a) Write down the likelihood function.

(b) Write down the log-likelihood function.

(c) Find the maximum likelihood estimator for $\mu_{\text{ML}}$.

(d) Find the maximum likelihood estimator for $\sigma^2_{\text{ML}}$.

(e) Show that the $\mu_{ML}$ is unbiased.

(f) Show that the $\sigma^2_{ML}$ is biased.

(g) Give an unbiased estimator for the variance parameter.

(a) The likelihood function is

$$p(\boldsymbol{x}|\mu, \sigma^2) = \prod_{n=1}^{N} \mathcal{N}(x_n|\mu, \sigma^2).$$

(b) The log-likelihood function is

$$
\begin{aligned}
\log p(\boldsymbol{x}|\mu, \sigma^2) &= \log(\prod_{n=1}^{N} \mathcal{N}(x_n|\mu, \sigma^2)) \\
&= \sum_{n=1}^{N} \log \mathcal{N}(x_n|\mu, \sigma^2) \\
&= -\frac{1}{2\sigma^2} \sum_{n=1}^{N} (x_n - \mu)^2 - \frac{N}{2} \log \sigma^2 - \frac{N}{2} \log(2\pi).
\end{aligned}
$$

(c) To find the MLE, we take the derivative with respect to $\mu$ and set it equal to zero. Solving for $\sigma^2$, we get that the MLE of the mean estimator is the sample mean:

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^{N} x_n.$$

(d) To find the MLE, we take the derivative with respect to $\sigma^2$ and set it equal to zero. Solving for $\sigma^2$, we get that the MLE of the variance estimator is the sample variance:

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^{N} (x_n - \hat{\mu})^2.$$

(e) The bias for the mean estimator is:

$$\text{bias}(\mu_{\text{ML}}) = \mathbb{E}[\mu_{\text{ML}}] - \mu$$
$$= \frac{1}{N} \sum_{n=1}^{N} \mathbb{E}[x_n] - \mu$$
$$= \mu - \mu = 0.$$

(f) The bias for the variance estimator is:

$$\text{bias}(\sigma_{\text{ML}}^2) = \mathbb{E}[\sigma_{\text{ML}}^2] - \sigma^2$$
$$= \left(\frac{N-1}{N}\right)\sigma^2 - \sigma^2$$
$$= \frac{N-1}{N} - 1.$$

(g) To get an unbiased estimator for the variance, we multiply by the bias:

$$\frac{N}{N-1}\sigma_{\text{ML}}^2.$$

4. **MLE Estimate of the Bias Term (Bishop (3.19))**

Let $\mathbf{\Phi}$ be our $N \times J$ design matrix, $\mathbf{t}$ our vector of $N$ target values, $\mathbf{w}$ our vector of $J$ parameters, and $w_0$ our bias parameter. As Bishop notes in (3.18), the sum-of-squares error function of $\mathbf{w}$ and $w_0$ can be written as follows

$$E(\mathbf{w}, w_0) = \frac{1}{2} \sum_{n=1}^{N} \left( t_n - w_0 - \sum_{j=1}^{J-1} w_j \cdot \phi_j(x_n) \right)^2.$$

Show that the value of $w_0$ that minimizes $E$ is

$$w_{0_{MLE}} = \frac{1}{N} \sum_{n=1}^{N} t_n - \sum_{j=1}^{J-1} w_j \cdot \left( \frac{1}{N} \sum_{n=1}^{N} \phi_j(x_n) \right)$$

$$= \bar{t} - \sum_{j=1}^{J-1} w_j \cdot \overline{\phi_j(x)} \qquad \text{[compare Bishop (3.19)]}$$

We have that $\frac{\partial E}{\partial w_0} = -\sum_{n=1}^{N} (t_n - w_0 - \sum_{j=1}^{J-1} w_j \cdot \phi_j(x_n))$.

Thus, we set $\sum_{n=1}^{N} t_n - N w_0 - \sum_{n=1}^{N} \sum_{j=1}^{J-1} w_j \cdot \phi_j(x_n) = 0$, and solving for $w_0$ gives the result.

5. **Simple Linear Regression (Bishop, 1.1)**

Consider the sum-of-squares error function given by:

$$E(\boldsymbol{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \boldsymbol{w}) - t_n\}^2,$$

in which the function $y(x, \boldsymbol{w})$ is given by the polynomial

$$y(x, \boldsymbol{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^{M} w_j x^j.$$

Show that the coefficients $\boldsymbol{w} = \{w_i\}$ that minimize this error function are given by the solution to the following set of linear equations

$$\sum_{j=0}^{M} A_{ij} w_j = T_i$$

where

$$A_{ij} = \sum_{n=1}^{N} (x_n)^{i+j},$$

$$T_i = \sum_{n=1}^{N} (x_n)^i t_n.$$

Here a suffix $i$ or $j$ denotes the index of a component, where as $(x)^i$ denotes $x$ reaised to the power of $i$.

Substituting the second equation into the first equation and then differentiating with respect to $w_i$, we obtain

$$\sum_{n=1}^{N} \left( \sum_{j=0}^{M} w_j x_n^j - t_n \right) x_n^i = 0.$$

6. **Multivariate Regression (Adapted from Stanford CS229)**

Suppose we have $\mathbf{\Phi} \in \mathbb{R}^{N \times J}$ as our design matrix, but that instead of predicting scalar values $t_n$, we'd like to use least squares regression to predict vector-valued targets $\mathbf{t}_n \in \mathbb{R}^M$ for each row $\phi(\mathbf{x}_n)$ in $\mathbf{\Phi}$. To do this, we can introduce a parameter matrix $\mathbf{W} \in \mathbb{R}^{N \times M}$ and attempt to minimize the following sum-of-squares error function:

$$E(\mathbf{W}) = \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{M} \left( (\mathbf{W}^\mathsf{T} \phi(\mathbf{x}_n))_m - t_{nm} \right)^2$$

$$= \frac{1}{2} \sum_{n=1}^{N} (\mathbf{W}^\mathsf{T} \phi(\mathbf{x}_n) - \mathbf{t}_n)^\mathsf{T} (\mathbf{W}^\mathsf{T} \phi(\mathbf{x}_n) - \mathbf{t}_n)$$

$$= \frac{1}{2} \mathrm{tr}((\mathbf{\Phi} \mathbf{W} - \mathbf{T})^\mathsf{T} (\mathbf{\Phi} \mathbf{W} - \mathbf{T})),$$

where $\mathbf{T} = \begin{bmatrix} - & \mathbf{t}_1^\mathsf{T} & - \\ & \vdots & \\ - & \mathbf{t}_N^\mathsf{T} & - \end{bmatrix}$.

(a) Rewrite $E(\mathbf{W})$ in terms of the traces of 3 matrices.

(b) Derive the gradient of your answer to part (a) with respect to $\mathbf{W}$.

(c) Derive $\mathbf{W}_{MLE}$ by setting the gradient of your answer to part (b) to 0.

Hint: in answering the above questions, it may be useful to keep in mind the following facts:

1. $\mathrm{tr}(\mathbf{A} + \mathbf{B}) = \mathrm{tr}(\mathbf{A}) + \mathrm{tr}(\mathbf{B})$
2. $\mathrm{tr}(\mathbf{A}) = \mathrm{tr}(\mathbf{A}^\mathsf{T})$
3. $\mathrm{tr}(\mathbf{A}\mathbf{B}) = \mathrm{tr}(\mathbf{B}\mathbf{A})$
4. $\frac{\partial}{\partial \mathbf{A}} \mathrm{tr}(\mathbf{A}\mathbf{B}) = \mathbf{B}^\mathsf{T}$
5. $\frac{\partial}{\partial \mathbf{A}} \mathrm{tr}(\mathbf{A}^\mathsf{T} \mathbf{B}) = \mathbf{B}$

(a)

$$\frac{1}{2} \mathrm{tr}((\mathbf{\Phi}\mathbf{W} - \mathbf{T})^\mathsf{T}(\mathbf{\Phi}\mathbf{W} - \mathbf{T})) = \frac{1}{2} \mathrm{tr}(\mathbf{W}^\mathsf{T}\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi}\mathbf{W} - \mathbf{W}^\mathsf{T}\mathbf{\Phi}^\mathsf{T}\mathbf{T} - \mathbf{T}^\mathsf{T}\mathbf{\Phi}\mathbf{W} + \mathbf{T}^\mathsf{T}\mathbf{T})$$

$$= \frac{1}{2} \left[ \mathrm{tr}(\mathbf{W}^\mathsf{T}\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi}\mathbf{W}) - \mathrm{tr}(\mathbf{W}^\mathsf{T}\mathbf{\Phi}^\mathsf{T}\mathbf{T}) - \mathrm{tr}(\mathbf{T}^\mathsf{T}\mathbf{\Phi}\mathbf{W}) + \mathrm{tr}(\mathbf{T}^\mathsf{T}\mathbf{T}) \right]$$

$$= \frac{1}{2} \left[ \mathrm{tr}(\mathbf{W}^\mathsf{T}\mathbf{\Phi}^\mathsf{T}\mathbf{\Phi}\mathbf{W}) - 2\mathrm{tr}(\mathbf{T}^\mathsf{T}\mathbf{\Phi}\mathbf{W}) + \mathrm{tr}(\mathbf{T}^\mathsf{T}\mathbf{T}) \right],$$

where in the last line we used the fact that $\mathrm{tr}(\mathbf{A}) = \mathrm{tr}(\mathbf{A}^\mathsf{T})$.

(b) We have $\nabla_W E = \frac{1}{2} \left[ \mathbf{\Phi}^\mathsf{T} \mathbf{\Phi} W + \mathbf{\Phi}^\mathsf{T} \mathbf{\Phi} W - 2\mathbf{\Phi}^\mathsf{T} T \right] = \mathbf{\Phi}^\mathsf{T} \mathbf{\Phi} W - \mathbf{\Phi}^\mathsf{T} T$, where the first equality is obtained by using the product rule on the first trace, and the 4th identity on the second trace.

(c) Setting our answer to 0, we get $W_{MLE} = (\mathbf{\Phi}^\mathsf{T} \mathbf{\Phi})^{-1} \mathbf{\Phi}^\mathsf{T} T$, which is almost exactly the same as in the scalar prediction case.