

Exercise 1: Approximation interpretation of PCA

Suppose we want to find an orthogonal set of M linear basis vectors $\mathbf{u}_j \in \mathbb{R}^D$, and the corresponding scores $\mathbf{z}_i \in \mathbb{R}^M$, such that we minimize the average reconstruction error (Equation 1)

$$J = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2 = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{U}\mathbf{z}_i\|^2 \quad (1)$$

where \mathbf{U} is an orthonormal matrix with \mathbf{u}_j as its j -th column. Show that the optimal solution is obtained by setting $\mathbf{U}^* = \mathbf{V}_M$, where \mathbf{V}_M contains the M eigenvectors with largest eigenvalues of the empirical covariance matrix $\mathbf{\Sigma} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top$. (We assume the \mathbf{x}_i have zero mean, for notational simplicity.) Furthermore, the optimal low-dimensional encoding of the data is given by $\mathbf{z}_i = \mathbf{U}^\top \mathbf{x}_i$, which is an orthogonal projection of the data onto the column space spanned by the eigenvectors. (Hint: start with the case of $M = 1, 2$ and then proof by induction; Use Lagrange multipliers to include the orthonormality constraints $\mathbf{u}_i^\top \mathbf{u}_j = \delta_{ij}$.)

Solution: We use $\mathbf{u}_j \in \mathbb{R}^D$ to denote the j -th principal direction, $\mathbf{x}_i \in \mathbb{R}^D$ to denote the i -th high-dimensional observation, $\mathbf{z}_i \in \mathbb{R}^L$ to denote the i -th low-dimensional representation (i.e., the projection), and $\tilde{\mathbf{z}}_j \in \mathbb{R}^N$ to denote the $[z_{1j}, \dots, z_{Nj}]$, which is the j -th component of all the low-dimensional vectors.

Let us start by estimating the best 1d solution, $\mathbf{u}_1 \in \mathbb{R}^D$, and the corresponding projected points $\tilde{\mathbf{z}}_1 \in \mathbb{R}^N$. The reconstruction error is given by:

$$\begin{aligned} J(\mathbf{u}_1, \mathbf{z}_1) &= \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - z_{i1} \mathbf{u}_1\|^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - z_{i1} \mathbf{u}_1)^\top (\mathbf{x}_i - z_{i1} \mathbf{u}_1) \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i^\top \mathbf{x}_i - 2z_{i1} \mathbf{u}_1^\top \mathbf{x}_i + z_{i1}^2 \mathbf{u}_1^\top \mathbf{u}_1) \\ &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i^\top \mathbf{x}_i - 2z_{i1} \mathbf{u}_1^\top \mathbf{x}_i + z_{i1}^2) \end{aligned}$$

where we have used the orthonormality $\mathbf{u}_1^\top \mathbf{u}_1 = 1$.

Taking derivatives w.r.t. z_{i1} and equating to zero gives:

$$\frac{\partial}{\partial z_{i1}} J(\mathbf{u}_1, \mathbf{z}_1) = \frac{1}{N} [-2\mathbf{u}_1^\top \mathbf{x}_i + 2z_{i1}] \stackrel{!}{=} 0 \quad \Rightarrow \quad z_{i1} = \mathbf{u}_1^\top \mathbf{x}_i$$

i.e., the optimal reconstruction weights are obtained by orthogonally projecting the data onto \mathbf{u}_1 .

Plugging the optimal \mathbf{z}_1 into the expression of J , we obtain:

$$J(\mathbf{u}_1) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i^\top \mathbf{x}_i - z_{i1}^2) = \text{const} - \frac{1}{N} \sum_{i=1}^N z_{i1}^2$$

Also, we see that the variance of the projected data is given by:

$$\text{Var}[\tilde{z}_1] = \mathbb{E}[\tilde{z}_1^2] - (\mathbb{E}[\tilde{z}_1])^2 = \frac{1}{N} \sum_{i=1}^N z_{i1}^2 - 0 = \frac{1}{N} \sum_{i=1}^N z_{i1}^2$$

since $\mathbb{E}[z_{i1}] = \mathbb{E}[\mathbf{x}_i^\top \mathbf{u}_1] = \mathbb{E}[\mathbf{x}_i]^\top \mathbf{u}_1 = 0$

Therefore, minimizing the reconstruction error $J(\mathbf{u}_1)$ is equivalent to maximizing the variance of the projected data, i.e.,

$$\underset{\mathbf{u}_1}{\text{argmin}} J(\mathbf{u}_1) = \underset{\mathbf{u}_1}{\text{argmax}} \text{Var}[\tilde{z}_1]$$

The variance of the projected data can also be written as:

$$\frac{1}{N} \sum_{i=1}^N z_{i1}^2 = \frac{1}{N} \sum_{i=1}^N \mathbf{u}_1^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u}_1 = \mathbf{u}_1^\top \hat{\Sigma} \mathbf{u}_1$$

which is exactly the objective in conventional PCA, i.e.,

$$\mathbf{u}_1 = \underset{\mathbf{u}_1}{\text{argmax}} \mathbf{u}_1^\top \hat{\Sigma} \mathbf{u}_1 \quad \text{s.t.} \quad \mathbf{u}_1^\top \mathbf{u}_1 = 1$$

Solving this by using Lagrange multipliers, we see that \mathbf{u}_1 is the eigenvector of the covariance matrix $\hat{\Sigma}$ with the largest associated eigenvalue.

Assume that it holds that $\forall j \leq M-1$ that $z_{ij} = \mathbf{u}_j^\top \mathbf{x}_i$, and \mathbf{u}_j is the eigenvector of the covariance matrix $\hat{\Sigma}$ with the j -th largest associated eigenvalue (And the orthogonormality holds $\mathbf{u}_i \mathbf{u}_j = \delta_{ij}$ for $i, j \leq M-1$).

Now prove the case for $j = M$.

$$J = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - z_{i1}\mathbf{u}_1 - z_{i2}\mathbf{u}_2 - \dots - z_{iM}\mathbf{u}_M\|^2 \quad (2)$$

Optimizing w.r.t. z_{iM} (setting the derivative of J w.r.t. z_{iM} equals to zero) gives:

$$\frac{\partial J}{\partial z_{iM}} = \frac{1}{N} [-2\mathbf{u}_M^\top \mathbf{x}_i + 2z_{iM}] \stackrel{!}{=} 0 \quad \Rightarrow \quad z_{iM} = \mathbf{u}_M^\top \mathbf{x}_i$$

Substituting the solutions for all z_{i1}, \dots, z_{iM} and $\mathbf{u}_1, \dots, \mathbf{u}_{M-1}$, we have:

$$\begin{aligned} J(\mathbf{u}_M) &= \frac{1}{N} \sum_{i=1}^N [\mathbf{x}_i^\top \mathbf{x}_i - \mathbf{u}_1^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u}_1 - \dots - \mathbf{u}_{M-1}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u}_{M-1} - \mathbf{u}_M^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u}_M] \\ &= \frac{1}{N} \sum_{i=1}^N (\text{const} - \mathbf{u}_M^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u}_M) \\ &= \text{const} - \mathbf{u}_M^\top \hat{\Sigma} \mathbf{u}_M \end{aligned}$$

Incorporating the orthogonormality constraints $\mathbf{u}_i \mathbf{u}_j = \delta_{ij}$ via Langrange multipliers:

$$\tilde{J}(\mathbf{u}_M) = -\mathbf{u}_M^\top \hat{\Sigma} \mathbf{u}_M + \lambda_M (\mathbf{u}_M^\top \mathbf{u}_M - 1) + \sum_{j=1}^{M-1} \lambda_{jM} (\mathbf{u}_M^\top \mathbf{u}_j - 0)$$

The stationary points occur when

$$0 = 2\hat{\Sigma} \mathbf{u}_M - 2\lambda_M \mathbf{u}_M + \sum_{j=1}^{M-1} \lambda_{jM} \mathbf{u}_j.$$

Left multiplying with \mathbf{u}_j^\top , and using the orthogonality constraints, we see that $\lambda_{jM} = 0$ for $j = 1, \dots, M-1$.

We therefore obtain

$$\hat{\Sigma} \mathbf{u}_M = \lambda_M \mathbf{u}_M$$

and so \mathbf{u}_M must be an eigenvector of $\hat{\Sigma}$ with eigenvalue λ_M . So the reconstruction error $\text{const} - \mathbf{u}_M^\top \hat{\Sigma} \mathbf{u}_M$ is minimized by choosing \mathbf{u}_M to be the eigenvector having the largest eigenvalue amongst those not previously selected.

Exercise 2: PCA and Kernel PCA

Show that the conventional linear PCA algorithm is recovered as a special case of kernel PCA if we choose the linear kernel function given by $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$.

Solution: W.l.o.g. assuming that the data is centered. For kernel PCA, the eigenvectors \mathbf{a} of the kernel matrix K (associated with largest eigenvalues) are computed:

$$K \mathbf{a} = \lambda \mathbf{a}$$

For linear kernel function $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$, the kernel matrix is equivalent to $K = \mathbf{X} \mathbf{X}^\top$.

Hence,

$$\begin{aligned} K \mathbf{a} &= \lambda \mathbf{a} \\ \iff \mathbf{X} \mathbf{X}^\top \mathbf{a} &= \lambda \mathbf{a} \\ \iff \mathbf{X}^\top \mathbf{X} \mathbf{X}^\top \mathbf{a} &= \lambda \mathbf{X}^\top \mathbf{a} \\ \iff \mathbf{C}' \mathbf{X}^\top \mathbf{a} &= \lambda \mathbf{X}^\top \mathbf{a} \\ \iff \mathbf{C}' \mathbf{u} &= \lambda \mathbf{u} \end{aligned}$$

where $\mathbf{C}' = \mathbf{C}$ corresponds to the (scaled) covariance matrix, $\mathbf{u} = \mathbf{X}^\top \mathbf{a}$ corresponds to the eigenvector of \mathbf{C}' .

The projection of data matrix in conventional PCA recovers that of linear kernel PCA:

$$\mathbf{X} \mathbf{U} = \mathbf{X} (\mathbf{X}^\top \mathbf{A}) = (\mathbf{X} \mathbf{X}^\top) \mathbf{A} = K \mathbf{A}$$

Exercise 3: Probabilistic PCA

Probabilistic PCA is a simple example of the linear-Gaussian framework. First, a latent variable \mathbf{z} is introduced (which corresponds to the principal-component subspace). Next, we define a Gaussian prior distribution $p(\mathbf{z})$ (Equation 3) over the latent variable, together with a Gaussian conditional distribution $p(\mathbf{x}|\mathbf{z})$ (Equation 4) for the observed variable \mathbf{x} conditioned on the value of the latent variable. Specifically, the mean of \mathbf{x} is a general linear function of \mathbf{z} governed by the $D \times M$ matrix \mathbf{W} and the D -dimensional vector $\boldsymbol{\mu}$.

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | 0, \mathbf{I}) \tag{3}$$

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x} | \mathbf{W} \mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \tag{4}$$

- i) Compute the marginal distribution $p(\mathbf{x})$ and posterior distribution $p(\mathbf{z}|\mathbf{x})$.

- ii) Given a data set $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ of observed data points, write down the corresponding log likelihood function.
- iii) Verify that the maximum likelihood solution for the parameter $\boldsymbol{\mu}$ is given by $\boldsymbol{\mu}_{\text{ML}} = \bar{\mathbf{x}}$, where $\bar{\mathbf{x}}$ is the mean of the data vectors.
- iv) The maximum likelihood solution for the parameter \mathbf{W} has a form of

$$\mathbf{W}_{\text{ML}} = \mathbf{U}_M(\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R} \quad (5)$$

where \mathbf{U}_M is a $D \times M$ matrix whose columns are given by the M eigenvectors of the data covariance matrix $\boldsymbol{\Sigma}$ with the largest eigenvalues, the $M \times M$ diagonal matrix \mathbf{L}_M has elements given by the corresponding eigenvalues λ_i , and \mathbf{R} is an arbitrary $M \times M$ orthogonal matrix.

Show that in the limit $\sigma^2 \rightarrow 0$, the posterior mean $\mathbb{E}[\mathbf{z}|\mathbf{x}]$ for the probabilistic PCA model becomes an orthogonal projection onto the principal subspace, as in conventional PCA.

Solution:

- i) Both the marginal and posterior will also be Gaussian as a result of the linear Gaussian model:

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}$$

And we have

$$\begin{aligned} \mathbb{E}[\mathbf{x}] &= \mathbb{E}[\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\varepsilon}] = \boldsymbol{\mu} \\ \text{Cov}[\mathbf{x}] &= \mathbb{E}[(\mathbf{W}\mathbf{z} + \boldsymbol{\varepsilon})(\mathbf{W}\mathbf{z} + \boldsymbol{\varepsilon})^\top] \\ &= \mathbb{E}[\mathbf{W}\mathbf{z}\mathbf{z}^\top \mathbf{W}^\top] + \mathbb{E}[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top] = \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I} \end{aligned}$$

Hence,

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C}) \quad \text{where} \quad \mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I} \quad (6)$$

Making use of matrix inversion identity (Equation (C.7) in [1])

$$(\mathbf{A} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}$$

to compute the \mathbf{C}^{-1} , we have

$$\mathbf{C}^{-1} = \sigma^{-1} \mathbf{I} - \sigma^{-2} \mathbf{W}\mathbf{M}^{-1}\mathbf{W}^\top$$

where $\mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}$ is an $M \times M$ matrix.

Then, making use of Equation (2.116) from [1] (See Figure 1), we have:

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mathbf{M}^{-1}\mathbf{W}^\top(\mathbf{x} - \boldsymbol{\mu}), \sigma^{-2}\mathbf{M}^{-1}) \quad (7)$$

- ii) The log likelihood function is given by:

$$\begin{aligned} \log p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \sum_{n=1}^N \log p(\mathbf{x}_n|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) \\ &= -\frac{ND}{2} \log(2\pi) - \frac{N}{2} \log |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^\top \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \quad (8) \end{aligned}$$

- iii) The log likelihood is a quadratic function of $\boldsymbol{\mu}$, thus taking the derivative and letting it equal to zero leads to the unique maximum solution $\boldsymbol{\mu}_{\text{ML}} = \bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$.

iv) Taking the limit $\sigma^2 \rightarrow 0$, we have $\mathbf{M} = \mathbf{W}^\top \mathbf{W}$

From the posterior in Equation 7, we have that:

$$\mathbb{E}[z|\mathbf{x}] = \mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x} - \boldsymbol{\mu}) = \mathbf{M}^{-1} \mathbf{W}_{\text{ML}}^\top (\mathbf{x} - \bar{\mathbf{x}}) = (\mathbf{W}_{\text{ML}}^\top \mathbf{W}_{\text{ML}})^{-1} \mathbf{W}_{\text{ML}}^\top (\mathbf{x} - \bar{\mathbf{x}})$$

Substituting for \mathbf{W}_{ML} using Equation 4, in which we take $\mathbf{R} = \mathbf{I}$ for compatibility with conventional PCA. Using the orthogonality property $\mathbf{U}_M^\top \mathbf{U}_M = \mathbf{I}$ and setting $\sigma^2 = 0$, we have

$$\begin{aligned} \mathbf{W}_{\text{ML}} &= \mathbf{U}_M \mathbf{L}_M^{1/2} \\ \mathbb{E}[z|\mathbf{x}] &= \mathbf{L}_M^{-1/2} \mathbf{U}_M^\top (\mathbf{x} - \bar{\mathbf{x}}) \end{aligned}$$

Note that this corresponds to the whitening operation:

- | | |
|-----------------|---|
| 1. centering : | $\mathbf{x} - \bar{\mathbf{x}}$ |
| 2. projection : | $\mathbf{U}_M^\top (\mathbf{x} - \bar{\mathbf{x}})$ |
| 2. rescaling : | $\mathbf{L}_M^{-1/2} \mathbf{U}_M^\top (\mathbf{x} - \bar{\mathbf{x}})$ |

Exercise 4: Neural Network: Properties of activation functions

- Show that the derivative of the *sigmoid* (i.e., $\sigma(z) = \frac{1}{1+e^{-z}}$) and the *tanh* (i.e., $\sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$) activation function can be expressed in terms of the function value itself.
- Show that the derivative of the binary cross-entropy error function (Equation 10) with respect to the activation a_k for an output unit having a logistic sigmoid activation function satisfies Equation 9.
- Show that the derivative of the multiclass cross-entropy error function (Equation 11) with respect to the activation a_k for output units having a softmax activation function satisfies Equation 9.

$$\frac{\partial E}{\partial a_k} = y_k - t_k \tag{9}$$

Solution:

- The derivative of *sigmoid*:

$$\begin{aligned} \sigma'(z) &= \frac{d}{dz} \left(\frac{1}{1+e^{-z}} \right) \\ &= - \left(\frac{1}{1+e^{-z}} \right)^2 \cdot \frac{d}{dz} (1+e^{-z}) \\ &= - \left(\frac{1}{1+e^{-z}} \right)^2 (-e^{-z}) \\ &= \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{1+e^{-z}} \\ &= \sigma(z)(1 - \sigma(z)) \end{aligned}$$

The derivative of *tanh*:

$$\begin{aligned}
tanh'(z) &= \frac{d}{dz} \left(\frac{e^z - e^{-z}}{e^z + e^{-z}} \right) \\
&= \frac{(e^z - e^{-z})'(e^z + e^{-z}) - (e^z - e^{-z})(e^z + e^{-z})'}{(e^z + e^{-z})^2} \\
&= \frac{(e^z + e^{-z})(e^z + e^{-z}) - (e^z - e^{-z})(e^z - e^{-z})}{(e^z + e^{-z})^2} \\
&= \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2} \\
&= 1 - \frac{(e^z - e^{-z})^2}{(e^z + e^{-z})^2} \\
&= 1 - tanh^2(z)
\end{aligned}$$

Both derivatives can be expressed by the function value itself.

- ii) We know that $y_k = \sigma(a_k)$ where σ is the logistic sigmoid function. As shown above, we have the derivative $\sigma' = \sigma(1 - \sigma)$. Thus, differentiating Equation 10 w.r.t. the activation a_k corresponding to a particular data point k , we obtain

$$\begin{aligned}
\frac{\partial E(\mathbf{w})}{\partial a_k} &= \frac{\partial}{\partial a_k} [- (t_k \cdot \log y_k + (1 - t_k) \cdot \log (1 - y_k))] \\
&= - \left(t_k \frac{1}{y_k} y_k (1 - y_k) + (1 - t_k) \frac{1}{1 - y_k} (-y_k (1 - y_k)) \right) \\
&= -t_k (1 - y_k) + (1 - t_k) y_k \\
&= -t_k + y_k
\end{aligned}$$

- iii) Similar to ii), we first denote $y_{kn} = y_k(\mathbf{x}_n, \mathbf{w})$ the k -th entry of the output vector on data point n . We know that $y_{kn} = \frac{\exp(a_{kn})}{\sum_{j=1}^K \exp(a_{jn})}$ and the derivative of softmax activation function is:

$$\frac{\partial y_{kn}}{\partial a_{jn}} = y_{kn} (\delta_{kj} - y_{jn})$$

Therefore,

$$\begin{aligned}
\frac{\partial E(\mathbf{w})}{\partial a_{jn}} &= - \sum_{k=1}^K t_{kn} \frac{1}{y_{kn}} (y_{kn} (\delta_{kj} - y_{jn})) \\
&= - \sum_{k=1}^K t_{kn} (\delta_{kj} - y_{jn}) \\
&= - \sum_{k=1}^K t_{kn} \delta_{kj} + \sum_{k=1}^K t_{kn} y_{jn} \\
&= -t_{jn} + y_{jn}
\end{aligned}$$

where we have used the fact that $\sum_{k=1}^K t_{kn} = 1$.

Exercise 5: Neural Network: Probabilistic Interpretation of Classification Models

- i) Consider a binary classification problem in which the target values are $t \in \{0, 1\}$, with a network output $y(\mathbf{x}, \mathbf{w})$ that represents $p(t = 1|\mathbf{x})$, and suppose that there is a probability ε

that the class label on a training data point has been incorrectly set. Assuming independent and identically distributed data, write down the error function corresponding to the negative log likelihood. Verify that when $\varepsilon = 0$, the error function is reduced to the usual cross-entropy error function:

$$E(\mathbf{w}) = - \sum_{n=1}^N \{t_n \cdot \log y_n + (1 - t_n) \cdot \log(1 - y_n)\} \quad (10)$$

Note that this error function (that consider mislabelling) makes the model robust to incorrectly labelled data, in contrast to the usual cross-entropy error function.

- ii) Show that maximizing likelihood for a multiclass neural network model in which the network outputs have the interpretation $y_k(\mathbf{x}, \mathbf{w}) = p(t_k = 1 | \mathbf{x})$ is equivalent to the minimization of the cross-entropy error function:

$$E(\mathbf{w}) = - \sum_{n=1}^N \sum_{k=1}^K \{t_{kn} \cdot \log y_k(\mathbf{x}_n, \mathbf{w})\} \quad (11)$$

Solution:

- i) First, we use t to denote the observed target label, and t_r to denote the real label. Then, we have that:

$$\begin{aligned} p(t = 1 | \mathbf{w}, \mathbf{x}) &= (1 - \varepsilon) \cdot p(t_r = 1 | \mathbf{w}, \mathbf{x}) + \varepsilon \cdot p(t_r = 0 | \mathbf{w}, \mathbf{x}) \\ p(t = 0 | \mathbf{w}, \mathbf{x}) &= (1 - \varepsilon) \cdot p(t_r = 0 | \mathbf{w}, \mathbf{x}) + \varepsilon \cdot p(t_r = 1 | \mathbf{w}, \mathbf{x}) \end{aligned}$$

Note that the network is aimed to predict the real label t_r instead of the noisy one t , i.e., we model $p(t_r = 1 | \mathbf{w}, \mathbf{x}) = y(\mathbf{w}, \mathbf{x})$.

Hence,

$$\begin{aligned} p(t = 1 | \mathbf{w}, \mathbf{x}) &= (1 - \varepsilon) \cdot y(\mathbf{w}, \mathbf{x}) + \varepsilon \cdot (1 - y(\mathbf{w}, \mathbf{x})) \\ p(t = 0 | \mathbf{w}, \mathbf{x}) &= (1 - \varepsilon) \cdot (1 - y(\mathbf{w}, \mathbf{x})) + \varepsilon \cdot y(\mathbf{w}, \mathbf{x}) \end{aligned}$$

Combining the two cases ($t = 0/1$), we have:

$$p(t | \mathbf{w}) = (1 - \varepsilon) \cdot y^t (1 - y)^{1-t} + \varepsilon (1 - y)^t y^{1-t}$$

Given a data set with N points, the error function corresponding to the negative log likelihood function is then:

$$E(\mathbf{w}) = - \sum_{n=1}^N \log \left((1 - \varepsilon) \cdot y_n^{t_n} (1 - y_n)^{1-t_n} + \varepsilon (1 - y_n)^{t_n} y_n^{1-t_n} \right) \quad (12)$$

When $\varepsilon = 0$, it is obvious that the equation above will reduce to Equation 10.

- ii) For the given interpretation of $y_k(\mathbf{x}, \mathbf{w})$, the probability that the observed sample has target vector \mathbf{t} is given by:

$$p(\mathbf{t} | \mathbf{w}) = \prod_{k=1}^K y_k^{t_k}$$

Then, for a data set of N points, the log likelihood function will be:

$$l(\mathbf{w}) = \sum_{n=1}^N \log p(\mathbf{t}_n | \mathbf{w}) = \sum_{n=1}^N \log \left(\prod_{k=1}^K y_k^{t_{kn}} \right) = \sum_{n=1}^N \sum_{k=1}^K t_{kn} \log y_{kn}$$

Maximizing the log likelihood function w.r.t. \mathbf{w} is exactly minimizing the cross-entropy in Equation 11.

Appendix

Marginal and Conditional Gaussians

Given a marginal Gaussian distribution for \mathbf{x} and a conditional Gaussian distribution for \mathbf{y} given \mathbf{x} in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \quad (2.113)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{Ax} + \mathbf{b}, \mathbf{L}^{-1}) \quad (2.114)$$

the marginal distribution of \mathbf{y} and the conditional distribution of \mathbf{x} given \mathbf{y} are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T) \quad (2.115)$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma}) \quad (2.116)$$

where

$$\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1}. \quad (2.117)$$

Figure 1: Commonly used results for linear Gaussian models.

References

- [1] C. M. Bishop. *Pattern recognition and machine learning*. springer, 2006.
- [2] J. Friedman, T. Hastie, R. Tibshirani, et al. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.
- [3] K. P. Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.