

Mixture of Experts

1

The likelihood for the entire dataset is given by:

$$\begin{aligned} p(y_n | \mathbf{x}_n, \Theta, \Phi) &= \prod_{n=1}^N \sum_{k=1}^k p(y_n | \mathbf{x}_n, \theta_k = \Theta \mathbf{z}_n) p(z_n = k | \mathbf{x}_n, \Phi) \\ &= \prod_{n=1}^N \sum_{k=1}^k \pi_{nk} \text{Exponential}(y_n | \lambda = \exp(\theta_k^T \mathbf{x}_n)) \end{aligned}$$

Since $\text{Exponential}(y_n | \lambda = \lambda \exp(-\lambda y))$:

$$= \prod_{n=1}^N \sum_{k=1}^k \pi_{nk} \exp(\theta_k^T \mathbf{x}_n) \exp(-\exp(\theta_k^T \mathbf{x}_n) y_n)$$

Taking the log on both sides:

$$\begin{aligned} \ln p(y_n | \mathbf{x}_n, \Theta, \Phi) &= \ln \prod_{n=1}^N \sum_{k=1}^k \pi_{nk} \exp(\theta_k^T \mathbf{x}_n) \exp(-\exp(\theta_k^T \mathbf{x}_n) y_n) \\ &= \sum_{n=1}^N \ln \sum_{k=1}^k \pi_{nk} e^{(\theta_k^T \mathbf{x}_n)} \exp(-e^{(\theta_k^T \mathbf{x}_n) y_n}) \end{aligned}$$

2

The formula for posterior in general is given by:

$$P(w | \mathcal{D}) = \frac{P(\mathcal{D} | w) P(w)}{P(\mathcal{D})}$$

The posterior probability r_{ni} of expert i producing label y for datapoint n is given by:

$$\begin{aligned} r_{ni} &= \frac{p(z_n = i | \mathbf{x}_n, \Phi) p(y_n | \mathbf{x}_n, \theta_i = \Theta \mathbf{z}_n)}{p(y_n | \mathbf{x}_n, \Theta, \Phi)} \\ &= \frac{\pi_{ni} \exp(-\exp(\theta_i^T \mathbf{x}_n) y_n)}{\sum_{i=1} \pi_{ni} \exp(-\exp(\theta_i^T \mathbf{x}_n) y_n)} \end{aligned}$$

3

Rewriting the log-likelihood equation in terms of its probability functions we get:

$$\begin{aligned} \ln p(y | \mathbf{x}, \Theta, \Phi) &= \sum_{n=1}^N \ln \sum_{k=1}^k \pi_{nk} e^{(\theta_k^T \mathbf{x}_n)} \exp(-e^{(\theta_k^T \mathbf{x}_n) y_n}) \\ &= \sum_{n=1}^N \ln \sum_{k=1}^k p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = k | \mathbf{x}_n, \Phi) \end{aligned}$$

First, we take the derivative with respect to θ_i :

$$\begin{aligned}\frac{\partial L}{\partial \theta_i} &= \frac{\partial}{\partial \theta_i} \left[\sum_{n=1}^N \ln \sum_{k=1}^K p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = k | \mathbf{x}_n, \Phi) \right] \\ &= \sum_{n=1}^N \left[\frac{\frac{\partial}{\partial \theta_i} \sum_{k=1}^K p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = k | \mathbf{x}_n, \Phi)}{\sum_{k=1}^K p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = k | \mathbf{x}_n, \Phi)} \right] \\ &= \sum_{n=1}^N \left[\frac{p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = i | \mathbf{x}_n, \Phi)}{\sum_{k=1}^K p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = k | \mathbf{x}_n, \Phi)} \right]\end{aligned}$$

From 1.2, we know that: $r_{ni} = \frac{p(z_n=i|\mathbf{x}_n, \Phi)p(y_n|\mathbf{x}_n, z_n, \Theta)}{p(y_n|\mathbf{x}_n, \Theta, \Phi)}$. Moreover, using the fact that: $\frac{\partial f(x)}{\partial x} = f(x) \frac{\partial \ln f(x)}{\partial x}$, we can re-write the equation as:

$$\begin{aligned}&= \sum_{n=1}^N \left[\frac{p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = i | \mathbf{x}_n, \Phi)}{\sum_{k=1}^K p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = k | \mathbf{x}_n, \Phi)} \right] \\ &= \sum_{n=1}^N r_{ni} \frac{\partial}{\partial \theta_i} [\ln p(y_n | \mathbf{x}_n, z_n, \Theta)]\end{aligned}$$

Next, we take the derivative with respect to ϕ_i :

$$\begin{aligned}\frac{\partial L}{\partial \phi_i} &= \frac{\partial}{\partial \phi_i} \left[\sum_{n=1}^N \ln \sum_{k=1}^K p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = k | \mathbf{x}_n, \Phi) \right] \\ &= \sum_{n=1}^N \frac{\partial}{\partial \phi_i} \ln \left[\sum_{k=1}^K p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = k | \mathbf{x}_n, \Phi) \right] \\ &= \sum_{n=1}^N \left[\frac{\sum_{k=1}^K p(y_n | \mathbf{x}_n, z_n, \Theta) \frac{\partial}{\partial \phi_i} p(z_n = k | \mathbf{x}_n, \Phi)}{\sum_{k=1}^K p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = k | \mathbf{x}_n, \Phi)} \right] \\ &= \sum_{n=1}^N \sum_{k=1}^K \left[\frac{p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = i | \mathbf{x}_n, \Phi)}{\sum_{k=1}^K p(y_n | \mathbf{x}_n, z_n, \Theta) p(z_n = k | \mathbf{x}_n, \Phi)} \right] \\ &= \sum_{n=1}^N \sum_{k=1}^K r_{ni} \frac{\partial}{\partial \phi_i} [\ln p(z_n = i | \mathbf{x}_n, \Phi)]\end{aligned}$$

4

For the derivative with respect to theta $\frac{\partial L}{\partial \theta_i}$:

$$\begin{aligned}\frac{\partial L}{\partial \theta_i} &= \sum_{n=1}^N r_{ni} \frac{\partial}{\partial \theta_i} [\ln p(y_n | \mathbf{x}_n, z_n, \Theta)] \\ &= \sum_{n=1}^N r_{ni} \frac{\partial}{\partial \theta_i} \left[\sum_{k=1}^K \ln e^{(\theta_k^T \mathbf{x}_n)} \exp(-e^{(\theta_k^T \mathbf{x}_n)} y_n) \right] \\ &= \sum_{n=1}^N r_{ni} \frac{\partial}{\partial \theta_i} \left[\sum_{k=1}^K \theta_k^T \mathbf{x}_n - e^{(\theta_k^T \mathbf{x}_n)} y_n \right] \\ &= \sum_{n=1}^N r_{ni} \left[\frac{\partial}{\partial \theta_i} \sum_{k=1}^K \theta_k^T \mathbf{x}_n \right] - \left[\frac{\partial}{\partial \theta_i} \sum_{k=1}^K e^{(\theta_k^T \mathbf{x}_n)} y_n \right] \\ &= \sum_{n=1}^N r_{ni} \left[\mathbf{x}_n^T - \mathbf{x}_n^T e^{(\theta_k^T \mathbf{x}_n)} y_n \right] \\ &= \sum_{n=1}^N r_{ni} \mathbf{x}_n^T \left[1 - e^{(\theta_k^T \mathbf{x}_n)} y_n \right]\end{aligned}$$

For the derivative with respect to $\boldsymbol{\phi}_i$:

$$\begin{aligned}
\frac{\partial L}{\partial \boldsymbol{\phi}_i} &= \sum_{n=1}^N \sum_{k=1}^K r_{ni} \frac{\partial}{\partial \boldsymbol{\phi}_i} [\ln p(z_n = i | \mathbf{x}_n, \boldsymbol{\Phi})] \\
&= \sum_{n=1}^N \sum_{k=1}^K r_{ni} \frac{\partial}{\partial \boldsymbol{\phi}_i} \left[\ln \frac{e^{\boldsymbol{\Phi}_i^T \mathbf{x}_n}}{\sum_j e^{\boldsymbol{\Phi}_j^T \mathbf{x}_n}} \right] \\
&= \sum_{n=1}^N \sum_{k=1}^K r_{ni} \frac{\partial}{\partial \boldsymbol{\phi}_i} \left[\boldsymbol{\Phi}_i^T \mathbf{x}_n - \ln \sum_j e^{\boldsymbol{\Phi}_j^T \mathbf{x}_n} \right] \\
&= \sum_{n=1}^N \sum_{k=1}^K r_{ni} \frac{\partial}{\partial \boldsymbol{\phi}_i} [\boldsymbol{\Phi}_i^T \mathbf{x}_n] - \frac{\partial}{\partial \boldsymbol{\phi}_i} \left[\ln \sum_j e^{\boldsymbol{\Phi}_j^T \mathbf{x}_n} \right] \\
&= \sum_{n=1}^N \sum_{k=1}^K r_{ni} [x_n^T] - \frac{\partial}{\partial \boldsymbol{\phi}_i} \left[\ln \sum_j e^{\boldsymbol{\Phi}_j^T \mathbf{x}_n} \right] \\
&= \sum_{n=1}^N \sum_{k=1}^K r_{ni} [x_n^T] - \left[\frac{e^{\boldsymbol{\Phi}_i^T \mathbf{x}_n}}{\sum_j e^{\boldsymbol{\Phi}_j^T \mathbf{x}_n}} \right] [x_n^T] \\
&= \sum_{n=1}^N \sum_{k=1}^K r_{ni} [x_n^T] \left[1 - \frac{e^{\boldsymbol{\Phi}_i^T \mathbf{x}_n}}{\sum_j e^{\boldsymbol{\Phi}_j^T \mathbf{x}_n}} \right] \text{ or} \\
&= \sum_{n=1}^N \sum_{k=1}^K r_{ni} [x_n^T] [1 - p(z_n = i | \mathbf{x}_n, \boldsymbol{\Phi})]
\end{aligned}$$

Quadratic Discriminant Analysis

1

The probability density function for a multivariate gaussian is given by:

$$f(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|_k}} \exp \left\{ \frac{-(\mathbf{x} - \mu_k)^T \Sigma_k^{-1} (\mathbf{x} - \mu_k)}{2} \right\}$$

The prior $p(\mathcal{C}) = \pi_k$ and thus the joint probability is altogether given by:

$$p(\mathbf{x}_n | \mathcal{C}_k) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|_k}} \exp \left\{ \frac{-(\mathbf{x} - \mu_k)^T \Sigma_k^{-1} (\mathbf{x} - \mu_k)}{2} \right\} \pi_k$$

2

The likelihood is given by:

$$p(\mathbf{T}, \mathbf{X} | \boldsymbol{\pi}_1, \dots, \boldsymbol{\pi}_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K) = \prod_{n=1}^N \prod_{k=1}^K \left(\frac{1}{\sqrt{(2\pi)^d |\Sigma|_k}} \exp \left\{ \frac{-(\mathbf{x} - \mu_k)^T \Sigma_k^{-1} (\mathbf{x} - \mu_k)}{2} \right\} \pi_k \right)^{\mathcal{I}(t_n=k)}$$

The log likelihood is then given by:

$$\begin{aligned}
\ln p(\mathbf{T}, \mathbf{X} | \boldsymbol{\pi}_1, \dots, \boldsymbol{\pi}_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K) &= \ln \prod_{n=1}^N \prod_{k=1}^K \left(\frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp \left\{ \frac{-(x - \mu_k)^T \boldsymbol{\Sigma}_k^{-1} (x - \mu_k)}{2} \right\} \pi_k \right)^{I(t_n=k)} \\
&= \sum_{n=1}^N \sum_{k=1}^K \ln \left(\frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp \left\{ \frac{-(x - \mu_k)^T \boldsymbol{\Sigma}_k^{-1} (x - \mu_k)}{2} \right\} \pi_k \right) \\
&= \sum_{n=1}^N \sum_{k=1}^K \ln \left(\frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \right) - \left(\frac{(x - \mu_k)^T \boldsymbol{\Sigma}_k^{-1} (x - \mu_k)}{2} \right) + \ln(\pi_k)
\end{aligned}$$

3

The equality constraint is given by: $\sum_k^K \pi_k = 1$ The Lagrangian is therefore given by:

$$\begin{aligned}
\mathcal{L} &= \ln p(\mathbf{T}, \mathbf{X} | \boldsymbol{\pi}_1, \dots, \boldsymbol{\pi}_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K) + \lambda [1 - \sum_k^K \pi_k] \\
&= \sum_{n=1}^N \sum_{k=1}^K \ln \left(\frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \right) \left(\frac{-(x - \mu_k)^T \boldsymbol{\Sigma}_k^{-1} (x - \mu_k)}{2} \right) + \ln(\pi_k) + \lambda \left[1 - \sum_k^K \pi_k \right] \\
&= \sum_{n=1}^N \sum_{k=1}^K \ln \left(\frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \right) \left(\frac{-(x - \mu_k)^T \boldsymbol{\Sigma}_k^{-1} (x - \mu_k)}{2} \right) + \ln(\pi_k) + \lambda - \lambda \sum_k^K \pi_k
\end{aligned}$$

4

The partial derivative of the Lagrangian with respect to π_k is:

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \pi_k} &= \frac{\partial}{\partial \pi_k} \sum_{n=1}^N \sum_{k=1}^K \ln \left(\frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \right) \left(\frac{-(x - \mu_k)^T \boldsymbol{\Sigma}_k^{-1} (x - \mu_k)}{2} \right) + \ln(\pi_k) + \lambda - \lambda \sum_k^K \pi_k \\
&= \frac{\partial}{\partial \pi_k} \sum_{n=1}^N \sum_{k=1}^K \ln(\pi_k) - \lambda \sum_k^K \pi_k \\
&= \sum_{n=1}^N \sum_{k=1}^K \left(\frac{1}{\pi_k} \right) - \lambda K \\
0 &= \frac{NK}{\pi_k} - \lambda K \\
\frac{NK}{\pi_k} &= \lambda K \\
\pi_k &= \frac{N}{\lambda}
\end{aligned}$$

5

The partial derivative of the Lagrangian with respect to μ_k is:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mu_k} &= \frac{\partial}{\partial \mu_k} \sum_{n=1}^N \sum_{k=1}^K \left(\frac{-(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)}{2} \right) \\ &= \sum_{n=1}^N \sum_{k=1}^K \frac{\partial}{\partial \mu_k} \left(\frac{-(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)}{2} \right)\end{aligned}$$

Using $\frac{\partial}{\partial x} x^T A x = x^T (A^T + A)$

$$\begin{aligned}0 &= \sum_{n=1}^N \sum_{k=1}^K -\frac{1}{2} (x - \mu_k)^T (\Sigma_k^{-T} + \Sigma_k^{-1}) \\ 0 &= - \sum_{n=1}^N \sum_{k=1}^K (x^T - \mu_k^T) (\Sigma_k^{-T} + \Sigma_k^{-1}) \\ 0 &= - \sum_{n=1}^N \sum_{k=1}^K x^T (\Sigma_k^{-T} + \Sigma_k^{-1}) - \mu_k^T (\Sigma_k^{-T} + \Sigma_k^{-1}) \\ &= - \sum_{n=1}^N \sum_{k=1}^K x^T \Sigma_k^{-T} + x^T \Sigma_k^{-1} - \mu_k^T \Sigma_k^{-T} - \mu_k^T \Sigma_k^{-1} \\ &= - \sum_{n=1}^N \sum_{k=1}^K x^T \Sigma_k^{-T} - \sum_{n=1}^N \sum_{k=1}^K x^T \Sigma_k^{-1} + \sum_{n=1}^N \sum_{k=1}^K \mu_k^T \Sigma_k^{-T} + \sum_{n=1}^N \sum_{k=1}^K \mu_k^T \Sigma_k^{-1} \\ \sum_{n=1}^N \sum_{k=1}^K \mu_k^T \Sigma_k^{-T} + \sum_{n=1}^N \sum_{k=1}^K \mu_k^T \Sigma_k^{-1} &= \sum_{n=1}^N \sum_{k=1}^K x^T \Sigma_k^{-T} + \sum_{n=1}^N \sum_{k=1}^K x^T \Sigma_k^{-1} \\ \sum_{n=1}^N \sum_{k=1}^K \mu_k^T (\Sigma_k^{-T} + \Sigma_k^{-1}) &= \sum_{n=1}^N \sum_{k=1}^K x^T (\Sigma_k^{-T} + \Sigma_k^{-1}) \\ \mu_k &= x\end{aligned}$$

6

The partial derivative of the Lagrangian with respect to Σ_k is:

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \Sigma_k} &= \frac{\partial}{\partial \Sigma_k} \sum_{n=1}^N \sum_{k=1}^K \ln \left(\frac{1}{\sqrt{(2\pi)^d |\Sigma_k|}} \right) \left(\frac{-(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)}{2} \right) \\ &= \sum_{n=1}^N \sum_{k=1}^K \frac{\partial}{\partial \Sigma_k} \ln ((2\pi)^d |\Sigma_k|)^{-\frac{1}{2}} - \frac{\partial}{\partial \Sigma_k} \left(\frac{(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)}{2} \right) \\ &= \sum_{n=1}^N \sum_{k=1}^K -\frac{1}{2} (\ln (2\pi)^d |\Sigma_k|)^{-\frac{1}{2}} - \frac{1}{2} (x - \mu_k)^T (x - \mu_k) \\ &= - \sum_{n=1}^N \sum_{k=1}^K \frac{1}{2} (\ln (2\pi)^d |\Sigma_k|)^{-\frac{1}{2}} - \sum_{n=1}^N \sum_{k=1}^K \frac{1}{2} (x - \mu_k)^T (x - \mu_k) \\ \sum_{n=1}^N \sum_{k=1}^K \frac{1}{2} (\ln (2\pi)^d |\Sigma_k|)^{-\frac{1}{2}} &= - \sum_{n=1}^N \sum_{k=1}^K \frac{1}{2} (x - \mu_k)^T (x - \mu_k)\end{aligned}$$

Using equation 2.122 from Bishop 2.3.4:

$$\Sigma_k = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (x - \mu_k)(x - \mu_k)^T$$

7

π_k equals to the total N divided by the lagrange multiplier. μ_k the mean of the estimate is the observed data point itself. Finally, the Σ_k depends on the symmetric difference between the observed data points and the

means.

Principal Component Analysis

1

The projection z_{ni} is given by:

$$z_{ni} = u_i^T x_n$$

2

The empirical mean of the projection z_i across all points is given by:

$$\bar{z}_i = u_i^T \bar{x}$$

3

The empirical variance of the projection z_i in terms of covariance matrix \mathcal{S} is given by:

$$\begin{aligned} Var(z_i) &= \frac{1}{N} \sum_{n=1}^N \{u_i^T x_n - u_i^T \bar{z}_i\}^2 \\ &= u_i^T \mathcal{S} u_i \end{aligned}$$

4

$$\begin{aligned} Var(z_i) &= u_i^T \mathcal{S} u_i \\ \text{Since } \mathcal{S} &= U \Lambda U^T \\ Var(z_i) &= u_i^T U \Lambda U^T u_i \\ &= \lambda_i \end{aligned}$$

5

Reducing dimensionality from D to K such that 99% of variance is maintained by picking a K such that the proportion of variance explained is the highest:

$$\frac{\sum_{i=1}^K \lambda_i}{Tr(\mathcal{S})} = \frac{\sum_{i=1}^K \lambda_i}{\sum_{i=1}^D \lambda_i} > 0.99$$