

1 Mixture of Experts

a) Write down the data likelihood of $p(y|X, \Theta, \Phi)$ and its log form

Since we have N data points in the training dataset (i.i.d.), data likelihood can be written out as product of each data point:

$$p(y|X, \Theta, \Phi) = \prod_{n=1}^N p(y_n|x_n, \Theta, \Phi) \quad (1)$$

According to product rule and sum rule of probability, Equation 1 can be written as:

$$p(y|X, \Theta, \Phi) = \prod_{n=1}^N p(y_n|x_n, \Theta, \Phi) = \prod_{n=1}^N \sum_{k=1}^K p(y_n|x_n, \theta_k, z_n = k) p(z_n = k|x_n, \Phi) \quad (2)$$

And the log-likelihood is given by

$$\begin{aligned} \log p(y|X, \Theta, \Phi) &= \log \left(\prod_{n=1}^N \sum_{k=1}^K p(y_n|x_n, \theta_k, z_n = k) p(z_n = k|x_n, \Phi) \right) \\ &= \sum_{n=1}^N \log \left(\sum_{k=1}^K p(y_n|x_n, \theta_k, z_n = k) p(z_n = k|x_n, \Phi) \right) \end{aligned} \quad (3)$$

b) Write down the posterior probability r_{ni} of expert i producing the label y for datapoint n

Posterior probability r_{ni} of expert i producing the label y for datapoint n , is also referred as the responsibility of expert i for datapoint n , thus:

$$r_{ni} = p(z_n = i|y_n) = \frac{p(y_n|z_n = i)p(z_n = i)}{p(y_n)} = \frac{p(y_n|z_n = i)p(z_n = i)}{\sum_{j=1}^K p(y_n|z_n = j)p(z_n = j)} \quad (4)$$

Now we just need to plug the conditional variables into Equation 4, therefore:

$$r_{ni} = \frac{p(y_n|x_n, \theta_i, z_n = i)p(z_n = i|x_n, \Phi)}{\sum_{j=1}^K p(y_n|x_n, \theta_j, z_n = j)p(z_n = j|x_n, \Phi)} \quad (5)$$

c) Take the derivative of the log-likelihood w.r.t. the parameters of each expert θ_i and the parameters of the routing mechanism for each expert ϕ_i

We first build some preliminaries:

$$\frac{\partial f(x)}{\partial x} = f(x) \frac{\partial \log f(x)}{\partial x} \rightarrow \frac{\partial \log f(x)}{\partial x} = \frac{1}{f(x)} \frac{\partial f(x)}{\partial x} \quad (6)$$

We let

$$f(\theta) = \sum_{k=1}^K p(y_n|x_n, \theta_k, z_n = k) p(z_n = k|x_n, \Phi) \quad (7)$$

$$f(\phi) = \sum_{k=1}^K p(y_n|x_n, \theta_k, z_n = k) p(z_n = k|x_n, \Phi) \quad (8)$$

$$\mathcal{Z} = \log p(y|X, \Theta, \Phi) = \sum_{n=1}^N \log f(\theta) \quad (9)$$

Based on the above preliminaries, we can solve for the derivative of the log-likelihood w.r.t. the parameters of each expert θ_i as below:

$$\begin{aligned}
\frac{\partial \mathcal{Z}}{\partial \theta_i} &= \sum_{n=1}^N \frac{1}{f(\theta)} \frac{\partial f(\theta)}{\partial \theta_i} \\
&= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K p(y_n | x_n, \theta_k, z_n = k) p(z_n = k | x_n, \Phi)} \frac{p(z_n = i | x_n, \Phi) \partial p(y_n | x_n, \theta_i, z_n = i)}{\partial \theta_i} \\
&= \sum_{n=1}^N \frac{p(z_n = i | x_n, \Phi) p(y_n | x_n, \theta_i, z_n = i)}{\sum_{k=1}^K p(y_n | x_n, \theta_k, z_n = k) p(z_n = k | x_n, \Phi)} \frac{\partial \log p(y_n | x_n, \theta_i, z_n = i)}{\partial \theta_i} \\
&= \sum_{n=1}^N r_{ni} \frac{\partial \log p(y_n | x_n, \theta_i, z_n = i)}{\partial \theta_i}
\end{aligned} \tag{10}$$

Based on the above preliminaries, we can solve for the derivative of the log-likelihood w.r.t. the parameters of the routing mechanism for each expert ϕ_i as below:

$$\begin{aligned}
\frac{\partial \mathcal{Z}}{\partial \phi_i} &= \sum_{n=1}^N \frac{1}{f(\phi)} \frac{\partial f(\phi)}{\partial \phi_i} \\
&= \sum_{n=1}^N \frac{1}{\sum_{k=1}^K p(y_n | x_n, \theta_k, z_n = k) p(z_n = k | x_n, \Phi)} \frac{\sum_{k=1}^K p(y_n | x_n, \theta_i, z_n = k) \partial p(z_n = k | x_n, \Phi)}{\partial \phi_i} \\
&= \sum_{n=1}^N \sum_{k=1}^K \frac{p(z_n = k | x_n, \Phi) p(y_n | x_n, \theta_i, z_n = k)}{\sum_{k=1}^K p(y_n | x_n, \theta_k, z_n = k) p(z_n = k | x_n, \Phi)} \frac{\partial \log p(z_n = k | x_n, \Phi)}{\partial \phi_i} \\
&= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \frac{\partial \log p(z_n = k | x_n, \Phi)}{\partial \phi_i}
\end{aligned} \tag{11}$$

d) Replace the expressions for each of the respective probability distributions and compute the final derivatives for θ_i, ϕ_i

We first look at the final derivatives for θ_i :

$$\log p(y_n | x_n, \theta_i, z_n = i) = \log \exp(y_n | \lambda = \exp(\theta_i^T x_n)) = \log \lambda \exp(-\lambda y) = \log \lambda - \lambda y_n = \theta_i^T x_n - y_n \exp(\theta_i^T x_n) \tag{12}$$

$$\frac{\partial \log p(y_n | x_n, \theta_i, z_n = i)}{\partial \theta_i} = x_n^T - y_n \exp(\theta_i^T x_n) x_n^T \tag{13}$$

$$\frac{\partial \mathcal{Z}}{\partial \theta_i} = \sum_{n=1}^N r_{ni} (1 - y_n \exp(\theta_i^T x_n)) x_n^T \tag{14}$$

We then look at the final derivatives for ϕ_i :

$$\log p(z_n = k | x_n, \Phi) = \log \pi_{nk} = \log \frac{\exp(\phi_k^T x_n)}{\sum_{j=1}^{j=K} \exp(\phi_j^T x_n)} = \phi_k^T x_n - \log \sum_{j=1}^{j=K} \exp(\phi_j^T x_n) \tag{15}$$

$$\frac{\partial \log p(z_n = k | x_n, \Phi)}{\partial \phi_i} = x_n^T - \frac{\exp(\phi_i^T x_n) x_n^T}{\sum_{j=1}^{j=K} \exp(\phi_j^T x_n)} = (I[k = i] - \pi_{ni}) x_n^T \tag{16}$$

$$\frac{\partial \mathcal{Z}}{\partial \phi_i} = \sum_{n=1}^N \sum_{k=1}^K r_{nk} (I[k = i] - \pi_{ni}) x_n^T = \sum_{n=1}^N (r_{nk} - \pi_{ni}) x_n^T \tag{17}$$

2 Quadratic Discriminant Analysis

a) Write down the joint probability $p(\mathbf{x}_n, C_k)$ for a single datapoint using a Gaussian class-conditional density and prior

According to the product rule and Gaussian distribution for class observations:

$$p(\mathbf{x}_n, C_k) = p(\mathbf{x}_n | C_k) p(C_k) = \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k \quad (18)$$

b) Write down the likelihood function and its log form

Using the result obtained in a) for a single datapoint and i.i.d. assumption:

$$p(\mathbf{t}_n, \mathbf{x}_n | \pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K) = \prod_{n=1}^N \prod_{k=1}^K (\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k)^{t_{nk}} \quad (19)$$

Using the property of log function $\log(xy) = \log x + \log y$, we have:

$$\log p(\mathbf{t}_n, \mathbf{x}_n | \pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K) = \sum_{n=1}^N \sum_{k=1}^K \log \left((\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k)^{t_{nk}} \right) \quad (20)$$

$$= \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log(\mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \pi_k) \quad (21)$$

$$= \sum_{n=1}^N \sum_{k=1}^K t_{nk} (\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) \quad (22)$$

c) Write down the Lagrangian function using the log likelihood, a Lagrange multiplier and the equality constraint $\sum_{k=1}^K \pi_k = 1$

The equality constrain can be written as:

$$\sum_{k=1}^K \pi_k - 1 = 0 \quad (23)$$

Therefore the Lagrange function can be written as:

$$L(\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \lambda) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} (\log \pi_k + \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \quad (24)$$

where λ is the Lagrange multiplier.

d) Write down π_{ML}

We only have to look at the terms related to π_k in the Lagrange function, and take the derivative w.r.t. to it to zero:

$$\frac{\partial L}{\partial \pi_k} = \frac{\partial \left(\sum_{n=1}^N \sum_{k=1}^K t_{nk} \log \pi_k + \lambda \sum_{k=1}^K \pi_k \right)}{\partial \pi_k} = \lambda + \sum_{n=1}^N \frac{t_{nk}}{\pi_k} = \lambda + \frac{N_k}{\pi_k} = 0 \rightarrow \pi_k = \frac{-N_k}{\lambda} \quad (25)$$

where N_k represents number of samples in class k.

Now if we plug in π_k into Equation 23, we can solve for λ :

$$\sum_{k=1}^K \frac{-N_k}{\lambda} - 1 = 0 \rightarrow \lambda = -N \quad (26)$$

where N represents total number of samples in input space.

Therefore, we plug in λ back, we can obtain the final form of π_k :

$$\pi_{k,ML} = \frac{-N_k}{\lambda} = \frac{N_k}{N} \quad (27)$$

e) Write down μ_{ML}

We only have to look at the terms related to μ_k in the Lagrange function, and take the derivative w.r.t. to it to zero:

$$\frac{\partial L}{\partial \mu_k} = \frac{\partial}{\partial \mu_k} \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log N(\mathbf{x}_n | \mu_k, \Sigma_k) \quad (28)$$

As we only care about the quadratic form of Gaussian when setting the derivative to zero, we have:

$$\begin{aligned} \frac{\partial L}{\partial \mu_k} &= \frac{\partial}{\partial \mu_k} \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log N(\mathbf{x}_n | \mu_k, \Sigma_k) \\ &= \frac{\partial}{\partial \mu_k} \left\{ \sum_{n=1}^N \sum_{k=1}^K \frac{-t_{nk}}{2} (\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) + \text{const} \right\} \\ &= \sum_{n=1}^N \sum_{k=1}^K t_{nk} (\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} = 0 \rightarrow \sum_{n=1}^N (t_{nk} \mathbf{x}_n^T - t_{nk} \mu_k^T) = 0 \rightarrow \sum_{n=1}^N t_{nk} \mathbf{x}_n = \sum_{n=1}^N t_{nk} \mu_k = N_k \mu_k \end{aligned} \quad (29)$$

Therefore, we have the final form:

$$\mu_{k,ML} = \frac{1}{N_k} \sum_{n=1}^N t_{nk} \mathbf{x}_n \quad (30)$$

f) Write down Σ_{ML}

According to Equation 2.122 in Bishop book, we know the MLE solution Σ_{ML} for a multivariate Gaussian distribution is given by

$$\Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu_{ML})(\mathbf{x}_n - \mu_{ML})^T \quad (31)$$

In the case of class conditional Gaussian distribution, we only need to replace the mean and sample counts by the mean and sample counts of a class, therefore, we have:

$$\begin{aligned} \frac{\partial L}{\partial \Sigma_k} &= \frac{\partial}{\partial \Sigma_k} \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log N(\mathbf{x}_n | \mu_k, \Sigma_k) \\ &= \frac{\partial}{\partial \Sigma_k} \left\{ \sum_{n=1}^N \sum_{k=1}^K \frac{-t_{nk}}{2} (\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) + \text{const} \right\} = 0 \\ &\rightarrow \Sigma_{k,ML} = \frac{N_k}{N} \sum_{n=1}^N (\mathbf{x}_n - \mu_{k,ML})(\mathbf{x}_n - \mu_{k,ML})^T \end{aligned} \quad (32)$$

g) Write a single-sentence interpretation for each of the solutions: π_{ML} , μ_{ML} and Σ_{ML}

- for π_{ML} , it can be interpreted as the frequency of class k in the entire input data set.
- for μ_{ML} , it is the mean of all the input vector assigned to class K.
- for Σ_{ML} , as we can see from $\Sigma_{k,ML}$, the final covariance is a weighted average of the covariance for each class k.

3 Principal Component Analysis

a) Projection z_{ni} of given data point \mathbf{x}_n onto eigen vector \mathbf{u}_i

It is easy to see that:

$$z_{ni} = \mathbf{u}_i^T \mathbf{x}_n \quad (33)$$

where $\mathbf{u}_i^T = (u_1, u_2, \dots, u_D)$ and $\mathbf{x}_n^T = (x_1, x_2, \dots, x_D)$.

b) Empirical mean of the projection z_i across all points \mathbf{x}_n

Based on the assumption that all input dimensions have zero mean, we know that $\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = 0$, therefore, we have

$$\mathbb{E}(z_i) = \frac{1}{N} \sum_{n=1}^N z_{ni} = \frac{1}{N} \sum_{n=1}^N \mathbf{u}_i^T \mathbf{x}_n = \frac{\mathbf{u}_i^T}{N} \sum_{n=1}^N \mathbf{x}_n = 0 \quad (34)$$

c) Empirical variance of the projection z_i across all points \mathbf{x}_n

According to the definition of variance $\text{VAR}(X) = \mathbb{E}[(X - \mu)^2]$, we have

$$\text{VAR}(z_i) = \frac{1}{N} \sum_{n=1}^N (z_{ni} - \mathbb{E}(z_i))^2 = \frac{1}{N} \sum_{n=1}^N (\mathbf{u}_i^T \mathbf{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N \mathbf{u}_i^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{u}_i = \mathbf{u}_i^T \left(\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{u}_i = \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i \quad (35)$$

where $\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$

e) Replace covariance matrix \mathbf{S} with its eigen decomposition and simply the variance above

We just need to plug in the eigen decomposition:

$$\text{VAR}(z_i) = \mathbf{u}_i^T \mathbf{S} \mathbf{u}_i = \mathbf{u}_i^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T \mathbf{u}_i = (\mathbf{U}^T \mathbf{u}_i)^T \mathbf{\Lambda} (\mathbf{U}^T \mathbf{u}_i) = \mathbf{e}_i^T \mathbf{\Lambda} \mathbf{e}_i = \lambda_i \quad (36)$$

where $\mathbf{e}_i^T = (0, \dots, \mathbf{u}_i^T \mathbf{u}_i, \dots, 0) = (0, \dots, 1, \dots, 0)$

d) Select proper K ($K < D$) such that 99% of variance is captured

We can do the following:

- Do eigen decomposition of input data \mathbf{X} to get eigen values of $\mathbf{\Lambda} = [\lambda_1, \lambda_2, \dots, \lambda_D]$
- Sort $\mathbf{\Lambda}$ by descending order
- Select eigen values from sorted $\mathbf{\Lambda}$ until $\sum_{i=1}^K \lambda_i \geq 0.99$
- K has been found