

1 Mixture of Experts

1.1

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

$$= \prod_{n=1}^N \sum_{k=1}^K \pi_{nk} \cdot \text{Exponential}(y_n | \lambda = \exp(\theta_k^\top \mathbf{x}_n)) \quad (2)$$

$$\ln p(y | \mathbf{X}, \Theta, \Phi) = \sum_{n=1}^N \ln \sum_{k=1}^K \pi_{nk} \cdot \text{Exponential}(y_n | \lambda = \exp(\theta_k^\top \mathbf{x}_n)) \quad (3)$$

1.2

Responsibility of expert i for datapoint n :

$$r_{ni} = p(z_n = i | \mathbf{x}_n, y_n) \quad (4)$$

posterior = likelihood \times prior / evidence, so if we interpret the given routing mechanism as a prior over z_n , then

$$r_{ni} = \frac{p(y_n | \mathbf{x}_n, \theta_i, z_n = i) p(z_n = i | \mathbf{x}_n, \Phi)}{\sum_{k=1}^K p(y_n | \mathbf{x}_n, \theta_k, z_n = k) p(z_n = k | \mathbf{x}_n, \Phi)} \quad (5)$$

$$= \frac{\pi_{ni} \cdot \text{Exponential}(y_n | \lambda = \exp(\theta_i^\top \mathbf{x}_n))}{\sum_{k=1}^K \pi_{nk} \cdot \text{Exponential}(y_n | \lambda = \exp(\theta_k^\top \mathbf{x}_n))} \quad (6)$$

1.3

$$\frac{\partial}{\partial \theta_i} \ln p(y | \mathbf{X}, \Theta, \Phi) = \frac{\partial}{\partial \theta_i} \sum_{n=1}^N \ln \sum_{k=1}^K p(z_n = k | \mathbf{x}_n, \Phi) p(y_n | \mathbf{x}_n, \theta_k, z_n = k) \quad (7)$$

$$= \sum_{n=1}^N \frac{\partial}{\partial \theta_i} \ln \sum_{k=1}^K p(z_n = k | \mathbf{x}_n, \Phi) p(y_n | \mathbf{x}_n, \theta_k, z_n = k) \quad (8)$$

$$= \sum_{n=1}^N \frac{\frac{\partial}{\partial \theta_i} \sum_{k=1}^K p(z_n = k | \mathbf{x}_n, \Phi) p(y_n | \mathbf{x}_n, \theta_k, z_n = k)}{\sum_{k=1}^K p(z_n = k | \mathbf{x}_n, \Phi) p(y_n | \mathbf{x}_n, \theta_k, z_n = k)} \quad (9)$$

$$= \sum_{n=1}^N \frac{\frac{\partial}{\partial \theta_i} p(z_n = i | \mathbf{x}_n, \Phi) p(y_n | \mathbf{x}_n, \theta_i, z_n = i)}{\sum_{k=1}^K p(z_n = k | \mathbf{x}_n, \Phi) p(y_n | \mathbf{x}_n, \theta_k, z_n = k)} \quad (10)$$

$$= \sum_{n=1}^N \frac{p(z_n = i | \mathbf{x}_n, \Phi) \frac{\partial}{\partial \theta_i} p(y_n | \mathbf{x}_n, \theta_i, z_n = i)}{\sum_{k=1}^K p(z_n = k | \mathbf{x}_n, \Phi) p(y_n | \mathbf{x}_n, \theta_k, z_n = k)} \quad (11)$$

$$= \sum_{n=1}^N \frac{p(z_n = i | \mathbf{x}_n, \Phi) p(y_n | \mathbf{x}_n, \theta_i, z_n = i)}{\sum_{k=1}^K p(z_n = k | \mathbf{x}_n, \Phi) p(y_n | \mathbf{x}_n, \theta_k, z_n = k)} \frac{\frac{\partial}{\partial \theta_i} p(y_n | \mathbf{x}_n, \theta_i, z_n = i)}{p(y_n | \mathbf{x}_n, \theta_i, z_n = i)} \quad (12)$$

$$= \sum_{n=1}^N r_{ni} \frac{\partial}{\partial \theta_i} \ln p(y_n | \mathbf{x}_n, \theta_i, z_n = i) \quad (13)$$

$$\frac{\partial}{\partial \phi_i} \ln p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\Theta}, \boldsymbol{\Phi}) = \frac{\partial}{\partial \phi_i} \sum_{n=1}^N \ln \sum_{k=1}^K p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\Phi}) p(y_n \mid \mathbf{x}_n, \boldsymbol{\theta}_k, z_n = k) \quad (14)$$

$$= \sum_{n=1}^N \frac{\partial}{\partial \phi_i} \ln \sum_{k=1}^K p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\Phi}) p(y_n \mid \mathbf{x}_n, \boldsymbol{\theta}_k, z_n = k) \quad (15)$$

$$= \sum_{n=1}^N \frac{\frac{\partial}{\partial \phi_i} \sum_{k=1}^K p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\Phi}) p(y_n \mid \mathbf{x}_n, \boldsymbol{\theta}_k, z_n = k)}{\sum_{k=1}^K p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\Phi}) p(y_n \mid \mathbf{x}_n, \boldsymbol{\theta}_k, z_n = k)} \quad (16)$$

$$= \sum_{n=1}^N \frac{\sum_{k=1}^K p(y_n \mid \mathbf{x}_n, \boldsymbol{\theta}_k, z_n = k) \frac{\partial}{\partial \phi_i} p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\Phi})}{\sum_{k=1}^K p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\Phi}) p(y_n \mid \mathbf{x}_n, \boldsymbol{\theta}_k, z_n = k)} \quad (17)$$

$$= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \frac{\frac{\partial}{\partial \phi_i} p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\Phi})}{p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\Phi})} \quad (18)$$

$$= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \frac{\partial}{\partial \phi_i} \ln p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\Phi}) \quad (19)$$

1.4

$$\frac{\partial}{\partial \theta_i} \ln p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\Theta}, \boldsymbol{\Phi}) = \sum_{n=1}^N r_{ni} \frac{\partial}{\partial \theta_i} \ln p(y_n \mid \mathbf{x}_n, \boldsymbol{\theta}_i, z_n = i) \quad (20)$$

$$= \sum_{n=1}^N r_{ni} \frac{\partial}{\partial \theta_i} \ln \left[\exp(\boldsymbol{\theta}_i^\top \mathbf{x}_n) \exp(-\exp(\boldsymbol{\theta}_i^\top \mathbf{x}_n) y_n) \right] \quad (21)$$

$$= \sum_{n=1}^N r_{ni} \frac{\partial}{\partial \theta_i} \left[\boldsymbol{\theta}_i^\top \mathbf{x}_n - \exp(\boldsymbol{\theta}_i^\top \mathbf{x}_n) y_n \right] \quad (22)$$

$$= \sum_{n=1}^N r_{ni} \left[\mathbf{x}_n^\top - y_n \frac{\partial}{\partial \theta_i} \exp(\boldsymbol{\theta}_i^\top \mathbf{x}_n) \right] \quad (23)$$

$$= \sum_{n=1}^N r_{ni} \left[\mathbf{x}_n^\top - y_n \exp(\boldsymbol{\theta}_i^\top \mathbf{x}_n) \mathbf{x}_n^\top \right] \quad (24)$$

$$= \sum_{n=1}^N r_{ni} \left[1 - y_n \exp(\boldsymbol{\theta}_i^\top \mathbf{x}_n) \right] \mathbf{x}_n^\top \quad (25)$$

$\frac{\partial}{\partial \phi_i} \ln p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\Theta}, \boldsymbol{\Phi})$: next page.

$$\frac{\partial}{\partial \phi_i} \ln p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\Theta}, \boldsymbol{\Phi}) = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \frac{\partial}{\partial \phi_i} \ln p(z_n = k \mid \mathbf{x}_n, \boldsymbol{\Phi}) \quad (26)$$

$$= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \frac{\partial}{\partial \phi_i} \ln \frac{\exp(\phi_k^\top \mathbf{x}_n)}{\sum_j \exp(\phi_j^\top \mathbf{x}_n)} \quad (27)$$

$$= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \frac{\partial}{\partial \phi_i} \left[\phi_k^\top \mathbf{x}_n - \ln \sum_j \exp(\phi_j^\top \mathbf{x}_n) \right] \quad (28)$$

$$= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \left[\mathbf{x}^\top \delta_{ik} - \frac{\frac{\partial}{\partial \phi_i} \sum_j \exp(\phi_j^\top \mathbf{x}_n)}{\sum_j \exp(\phi_j^\top \mathbf{x}_n)} \right] \quad (29)$$

$$= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \left[\mathbf{x}^\top \delta_{ik} - \frac{\exp(\phi_i^\top \mathbf{x}_n) \mathbf{x}_n^\top}{\sum_j \exp(\phi_j^\top \mathbf{x}_n)} \right] \quad (30)$$

$$= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \left[\mathbf{x}^\top \delta_{ik} - \pi_{ni} \mathbf{x}_n^\top \right] \quad (31)$$

$$= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \left[\delta_{ik} - \pi_{ni} \right] \mathbf{x}^\top \quad (32)$$

$$= \sum_{n=1}^N \left[r_{ni} - \pi_{ni} \right] \mathbf{x}^\top \quad (33)$$

2 Quadratic Discriminant Analysis (QDA)

2.1

$$p(\mathbf{x}_n, \mathcal{C}_k) = p(\mathcal{C}_k) p(\mathbf{x}_n \mid \mathcal{C}_k) \quad (34)$$

$$= \pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (35)$$

2.2

Here we use $\boldsymbol{\theta}$ to denote $\pi_1, \dots, \pi_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K$.

$$p(\mathbf{T}, \mathbf{X} \mid \boldsymbol{\theta}) = \prod_{n=1}^N p(\mathbf{t}_n, \mathbf{x}_n \mid \boldsymbol{\theta}) \quad (36)$$

$$= \prod_{n=1}^N p(\mathbf{t}_n \mid \boldsymbol{\theta}) p(\mathbf{x}_n \mid \mathbf{t}_n, \boldsymbol{\theta}) \quad (37)$$

$$= \prod_{n=1}^N \prod_{k=1}^K \left[p(t_{nk} \mid \boldsymbol{\theta}) p(\mathbf{x}_n \mid t_{nk}, \boldsymbol{\theta}) \right]^{t_{nk}} \quad (38)$$

$$= \prod_{n=1}^N \prod_{k=1}^K \left[\pi_k p(\mathbf{x}_n \mid \mathcal{C}_k) \right]^{t_{nk}} \quad (39)$$

$$= \prod_{n=1}^N \prod_{k=1}^K \left[\pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right]^{t_{nk}} \quad (40)$$

$$\ln p(\mathbf{T}, \mathbf{X} \mid \boldsymbol{\theta}) = \sum_{n=1}^N \sum_{k=1}^K \ln \left[\pi_k p(\mathbf{x}_n \mid \mathcal{C}_k) \right]^{t_{nk}} \quad (41)$$

$$= \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln \left[\pi_k p(\mathbf{x}_n \mid \mathcal{C}_k) \right] \quad (42)$$

$$= \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln \left[\pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right] \quad (43)$$

2.3

Lagrangian:

$$L = \ln p(\mathbf{T}, \mathbf{X} \mid \boldsymbol{\theta}) + \lambda(1 - \sum_{k=1}^K \pi_k) \quad (44)$$

$$= \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln \left[\pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right] + \lambda(1 - \sum_{k=1}^K \pi_k) \quad (45)$$

2.4

$$\frac{\partial L}{\partial \pi_k} = \sum_{n=1}^N \frac{\partial}{\partial \pi_k} \sum_{k'=1}^K t_{nk'} \ln \left[\pi_{k'} \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'}) \right] + \lambda \frac{\partial}{\partial \pi_k} (1 - \sum_{k'}^K \pi_{k'}) \quad (46)$$

$$= \sum_{n=1}^N t_{nk} \frac{\partial}{\partial \pi_k} \ln \left[\pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right] - \lambda \frac{\partial}{\partial \pi_k} \sum_{k'}^K \pi_{k'} \quad (47)$$

$$= \sum_{n=1}^N t_{nk} \frac{1}{\pi_k} - \lambda \quad (48)$$

$$= \frac{N_k}{\pi_k} - \lambda, \quad N_k = \sum_{n=1}^N t_{nk} \quad (49)$$

$$\stackrel{!}{=} 0 \quad (50)$$

$$\pi_k = \frac{N_k}{\lambda} \quad (51)$$

$$(52)$$

Summing over all k gives us λ :

$$\sum_{k=1}^K \pi_k = 1 = \frac{\sum_{k=1}^K N_k}{\lambda} \quad (53)$$

$$\lambda = \sum_{k=1}^K N_k = N \quad (54)$$

So for π_{ML} we get:

$$\pi_k = \frac{N_k}{N} \quad (55)$$

2.5

In the following we use D to denote the size of vector \mathbf{x}_n .

$$\frac{\partial L}{\partial \boldsymbol{\mu}_k} = \sum_{n=1}^N \frac{\partial}{\partial \boldsymbol{\mu}_k} \sum_{k'=1}^K t_{nk'} \ln \left[\pi_{k'} \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'}) \right] \quad (56)$$

$$= \sum_{n=1}^N t_{nk} \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln \left[\pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right] \quad (57)$$

$$= \sum_{n=1}^N t_{nk} \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (58)$$

$$= \sum_{n=1}^N t_{nk} \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln \left(\frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}_k|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right] \right) \quad (59)$$

$$= \sum_{n=1}^N t_{nk} \frac{\partial}{\partial \boldsymbol{\mu}_k} \ln \exp \left[-\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right] \quad (60)$$

$$= \sum_{n=1}^N -\frac{t_{nk}}{2} \frac{\partial}{\partial \boldsymbol{\mu}_k} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (61)$$

$$= \sum_{n=1}^N -\frac{t_{nk}}{2} \left[-2 (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k \right] \quad (62)$$

$$= \sum_{n=1}^N t_{nk} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k \quad (63)$$

$$= \sum_{n=1}^N t_{nk} \mathbf{x}_n^\top \boldsymbol{\Sigma}_k - \sum_{n=1}^N t_{nk} \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}_k \quad (64)$$

$$= \sum_{n \in \mathcal{C}_k}^N \mathbf{x}_n^\top \boldsymbol{\Sigma}_k - N_k \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}_k \quad (65)$$

$$\stackrel{!}{=} 0 \quad (66)$$

$$N_k \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}_k = \sum_{n \in \mathcal{C}_k}^N \mathbf{x}_n^\top \boldsymbol{\Sigma}_k \quad (67)$$

We can take the transpose on each side, and eliminate $\boldsymbol{\Sigma}_k$ from the equation. So for $\boldsymbol{\mu}_{\text{ML}}$ we get:

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k}^N \mathbf{x}_n \quad (68)$$

2.6

Using Eq. 2.122 in Bishop (p.94) we find for $\boldsymbol{\Sigma}_{\text{ML}}$:

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k}^N (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top, \quad (69)$$

where $\boldsymbol{\mu}_k$ is the ML-estimator for the k 'th class mean. Since each class has its own covariance matrix, we made Eq. 2.122 class-specific.

2.7

So we have:

$$\pi_k = \frac{N_k}{N} \quad (70)$$

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k}^N \mathbf{x}_n \quad (71)$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n \in \mathcal{C}_k}^N (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top \quad (72)$$

π_k tells us what the probability is of class k , and the ML-estimator (Eq. 70) tells us that that probability can be estimated by taking the number of data points that belong to class k , and dividing it by the total number of data points N . The second one is also very intuitive. $\boldsymbol{\mu}_k$ tells us what the mean is of the data points belonging to class k , and The ML-estimator (Eq. 71) shows that that mean can be estimated by summing up all data points belonging to class k , and dividing that by the number of data points belonging to class k . Lastly, $\boldsymbol{\Sigma}_k$ tells us how the points in class k are correlated, and the ML-estimator (Eq. 72) shows that these correlations can be estimated by comparing each data point in that class with the mean of that class (and dividing the sum by the number of points in that class, N_k).

3 Principal Component Analysis

3.1

$$z_{ni} = \mathbf{u}_i^\top \mathbf{x}_n \quad (73)$$

3.2

Since our data set has zero mean for each dimension, the empirical mean of the projection z_i across all points \mathbf{x}_n is zero.

$$\bar{z}_i = \frac{1}{N} \sum_{n=1}^N z_{ni} = \mathbf{u}_i^\top \bar{\mathbf{x}} = 0 \quad (74)$$

3.3

$$\mathbb{V}(z_i) = \frac{1}{N} \sum_{n=1}^N (\mathbf{u}_i^\top \mathbf{x}_n - \mathbf{u}_i^\top \bar{\mathbf{x}})^2 = \mathbf{u}_i^\top \mathbf{S} \mathbf{u}_i \quad (75)$$

3.4

$$\mathbb{V}(z_i) = \mathbf{u}_i^\top \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top \mathbf{u}_i \quad (76)$$

Since \mathbf{U} is an orthonormal matrix, we get $\mathbf{u}_i^\top \mathbf{U} = \mathbf{e}_i^\top$, where \mathbf{e}_i is a vector of zeros, with the i 'th element being equal to 1. Also, because $\boldsymbol{\Lambda}$ is a diagonal matrix consisting of eigenvalues λ_i of \mathbf{S} , we get

$$\mathbb{V}(z_i) = \mathbf{e}_i^\top \boldsymbol{\Lambda} \mathbf{e}_i \quad (77)$$

$$= \lambda_i \quad (78)$$

3.5

K should be chosen such that the following holds:

$$\frac{\sum_{k=1}^K \lambda_k}{\sum_{d=1}^D \lambda_d} \geq 0.99 \quad (79)$$

This condition is most easily met if the eigenvalues are sorted in descending order.