

ML1 HOME ASSIGNMENT 4

id: 12179078

VOLODYMYR MEDENTSIR

w1

K-experts, N-datapoints $\{(x_i, y_i)\}_{i=1}^N$ with $x_i \in \mathbb{R}^D$

y_i - label of x_i ; z_n - k -of- K vector.

$\Theta \in \mathbb{R}^{D \times K}$ - vector of parameters for each expert.

$$p(y_n | x_n, z_n, \Theta) = \text{Exp}(y_n | \lambda = \exp(\Theta_K^\top x_n))$$

$$p(z_n = k | x_n, \varphi) = \pi_{nk} = \frac{\exp(\varphi_k^\top x_n)}{\sum_j \exp(\varphi_j^\top x_n)}$$

($\varphi \in \mathbb{R}^{D \times K}$ - parameters of the routing mechanism.
 $z_n = k$ means that $z_{nk} = 1$).

I Part

1.1. The likelihood of the entire dataset $p(y | X, \varphi, \Theta)$ and its log under i.i.d assumption.

$$p(y | X, \varphi, \Theta) = [\text{i.i.d. assumption}] = \prod_{n=1}^N p(y_n | x_n, \Theta, \varphi) =$$

$$= \prod_{n=1}^N \sum_{k=1}^K p(z_n = k | x_n, \varphi) \cdot p(y_n | x_n, z_n, \Theta) =$$

$$= \prod_{n=1}^N \sum_{k=1}^K \frac{\exp(\varphi_k^\top x_n)}{\sum_{j=1}^K \exp(\varphi_j^\top x_n)} \cdot \text{Exp}(y_n | \lambda = \exp(\Theta_K^\top x_n)) =$$

$$= \prod_{n=1}^N \sum_{k=1}^K \pi_{nk} \text{Exp}(y_n | \lambda = \exp(\Theta_K^\top x_n))$$

$$\log p(y | X, \varphi, \Theta) = \sum_{n=1}^N \log \sum_{k=1}^K \frac{\exp(\varphi_k^\top x_n)}{\sum_{j=1}^K \exp(\varphi_j^\top x_n)} \text{Exp}(y_n | \lambda = \exp(\Theta_K^\top x_n)) =$$

$$= \sum_{n=1}^N \log \sum_{k=1}^K \pi_{nk} \text{Exp}(y_n | \lambda = \exp(\Theta_K^\top x_n)).$$

1.2. Find the responsibility (r_{ni}) of the expert i for datapoint n .

$$\begin{aligned}
 r_{ni} &= p(z_n = i | y_n, x_n, \varphi, \Theta) = \frac{p(y_n | z_n = i, x_n, \Theta) p(z_n = i | x_n, \varphi)}{p(y_n | x_n, \Theta, \varphi)} = \\
 &= \frac{p(y_n | z_n = i, x_n, \Theta) \cdot p(z_n = i | x_n, \varphi)}{\sum_{j=1}^K p(z_n = j | x_n, \varphi) \cdot p(y_n | z_n = j, x_n, \Theta)} = \\
 &= \frac{\pi_{ni} \text{Exp}(y_n | \lambda = \exp(\Theta_i^\top x_n))}{\sum_{j=1}^K \pi_{nj} \text{Exp}(y_n | \lambda = \exp(\Theta_j^\top x_n))}
 \end{aligned}$$

1.3.

$$\frac{d}{d\Theta_i} \log p(y | X, \varphi, \Theta) = \sum_{n=1}^N \frac{1}{\sum_{k=1}^K p(y_n | x_n, z_n = k, \Theta) p(z_n = k | x_n, \varphi)}.$$

$$\bullet \frac{d}{d\Theta_i} \sum_{k=1}^K p(y_n | x_n, z_n = k, \Theta) \cdot p(z_n = k | x_n, \varphi) =$$

$$= \sum_{n=1}^N \frac{r_{ni}}{p(z_n = i | x_n, \varphi) \cdot p(y_n | x_n, z_n = i, \Theta)} \cdot p(z_n = i | x_n, \varphi).$$

$$\frac{d}{d\Theta_i} p(y_n | x_n, z_n = i, \Theta) = \sum_{n=1}^N \frac{r_{ni}}{p(y_n | x_n, z_n = i, \Theta)} \frac{d}{d\Theta_i} p(y_n | x_n, z_n = i, \Theta) =$$

$$= \sum_{n=1}^N r_{ni} \frac{d}{d\Theta_i} \log(p(y_n | x_n, z_n = i, \Theta))$$

$$\frac{d}{d\varphi_i} \log p(y|X, \varphi, \Theta) = \sum_{n=1}^N \frac{1}{\sum_{k=1}^K p(y_n|x_n, z_n=k, \Theta) p(z_n=k|x_n, \varphi)} =$$

$$\begin{aligned} \cdot \frac{d}{d\varphi_i} \sum_{k=1}^K p(y_n|x_n, z_n=k, \Theta) \cdot p(z_n=k|x_n, \varphi) &= \\ = \sum_{n=1}^N \frac{z_{ni}}{p(z_n=i|x_n, \varphi) \cdot p(y_n|x_n, z_n=i, \Theta)} \cdot \sum_{k=1}^K p(y_n|x_n, z_n=k) \frac{d}{d\varphi_i} p(z_n=k|x_n, \varphi) & \end{aligned}$$

1.4.

$$\begin{aligned} 1) \quad \frac{d}{d\theta_i} \log p(y|X, \varphi, \Theta) &= \sum_{n=1}^N \frac{z_{ni}}{p(y_n|x_n, z_n=i, \Theta)} \frac{d}{d\theta_i} p(y_n|x_n, z_n=i, \Theta) = \\ = \sum_{n=1}^N \frac{z_{ni}}{\exp(\theta_i^\top x_n) \exp(-\exp(\theta_i^\top x_n) y_n)} \frac{d \exp(\theta_i^\top x_n) \exp(-\exp(\theta_i^\top x_n) y_n)}{d\theta_i} &= \\ = \sum_{n=1}^N \frac{z_{ni}}{\exp(\theta_i^\top x_n) \exp(-\exp(\theta_i^\top x_n) y_n)} \cdot (\exp(\theta_i^\top x_n) \exp(-\exp(\theta_i^\top x_n) y_n) \cdot x_n^\top - & \\ - \exp(\theta_i^\top x_n) \cdot \exp(-\exp(\theta_i^\top x_n) y_n) \cdot \exp(\theta_i^\top x_n) \cdot y_n x_n^\top) &= \\ = \sum_{n=1}^N z_{ni} \cdot (1 - \exp(\theta_i^\top x_n) y_n) x_n^\top. & \end{aligned}$$

2)

$$\begin{aligned} \frac{d}{d\varphi_i} \log p(y|X, \varphi, \Theta) &= \sum_{n=1}^N \frac{z_{ni}}{\pi_{ni} \exp(\theta_i^\top x_n) \exp(-\exp(\theta_i^\top x_n) y_n)} \cdot \sum_{k=1}^K (\exp(\theta_k^\top x_n) \cdot & \\ \cdot \exp(-\exp(\theta_k^\top x_n) y_n) \cdot \frac{d}{d\varphi_i} \frac{\exp(\varphi_i^\top x_n)}{\sum_{j=1}^K \exp(\varphi_j^\top x_n)}) & \circledcirc \end{aligned}$$

$$\frac{d}{d\varphi_i} \pi_{nk} \stackrel{k \neq i}{=} \exp(\varphi_n^\top x_n) (-1) \cdot \frac{1}{\left(\sum_{j=1}^k \exp(\varphi_j^\top x_n) \right)^2} \cdot \exp(\varphi_i^\top x_n) \cdot x_n^\top =$$

$$= - \frac{\exp(\varphi_k^\top x_n + \varphi_i^\top x_n)}{\left(\sum_{j=1}^k \exp(\varphi_j^\top x_n) \right)^2} \cdot x_n^\top = - \pi_{nk} \cdot \pi_{ni} x_n^\top$$

$$\frac{d}{d\varphi_i} \pi_{ni} = \frac{\exp(\varphi_i^\top x_n) \cdot x_n^\top}{\sum_{j=1}^k \exp(\varphi_j^\top x_n)} - \frac{\exp(\varphi_i^\top x_n + \varphi_i^\top x_n)}{\left(\sum_{j=1}^k \exp(\varphi_j^\top x_n) \right)^2} \cdot x_n^\top =$$

$$= \frac{\exp(\varphi_i^\top x_n)}{\sum_{j=1}^k \exp(\varphi_j^\top x_n)} \cdot \left(1 - \frac{\exp(\varphi_i^\top x_n)}{\sum_{j=1}^k \exp(\varphi_j^\top x_n)} \right) x_n^\top = \pi_{ni} (1 - \pi_{ni}) x_n^\top$$

$$\textcircled{=} \sum_{n=1}^N \frac{\tau_{ni}}{\pi_{ni} \exp(\Theta_i^\top x_n) \exp(-\exp(\Theta_i^\top x_n) y_n)} \cdot \sum_{k=1}^K (\exp(\Theta_k^\top x_n) \cdot \exp(-\exp(\Theta_k^\top x_n) y_n))$$

$$\cdot \pi_{ni} (\delta_{ik} - \pi_{nk}) x_n^\top =$$

$$= \sum_{n=1}^N \tau_{ni} \cdot \sum_{k=1}^K (\exp((\Theta_k^\top - \Theta_i^\top) x_n) \cdot \exp(-(\exp \Theta_k^\top x_n - \exp \Theta_i^\top x_n) y_n) \cdot (\delta_{ik} - \pi_{nk})) x_n^\top$$

$$\text{with } \delta_{ik} = \begin{cases} 0, & i \neq k \\ 1, & i = k \end{cases}$$

Or we can also rewrite it as:

$$\sum_{n=1}^N \frac{\tau_{ni}}{\pi_{ni} \exp(\Theta_i^\top x_n) \exp(-\exp(\Theta_i^\top x_n) y_n)} \cdot \sum_{k=1}^K (\exp(\Theta_k^\top x_n) \cdot \exp(-\exp(\Theta_k^\top x_n) y_n))$$

$$\cdot \pi_{ni} (\delta_{ik} - \pi_{nk}) x_n^\top =$$

$$= \sum_{n=1}^N \frac{\pi_{ni} \left(-\sum_{k=1}^K \exp(\theta_k^T x_n) \cdot \exp(-\exp(\theta_k^T x_n) y_n) \cdot \pi_{nk} + ((1 - \pi_{ni}) + \pi_{ni}) \exp(\theta_i^T x_n) \right)}{\sum_{j=1}^N \pi_{nj} \exp(\theta_j^T x_n) \exp(-\exp(\theta_j^T x_n) y_n)}$$

$$\cdot \exp(-\exp(\theta_i^T x_n) y_n) \cdot x_n^T = \sum_{n=1}^N \pi_{ni} \left(-1 + \frac{\exp(\theta_i^T x_n) \cdot \exp(-\exp(\theta_i^T x_n))}{\sum_{j=1}^N \pi_j \exp(\theta_j^T x_n) \exp(-\exp(\theta_j^T x_n) y_n)} \right) x_n^T$$

$$= \sum_{n=1}^N \pi_{ni} \left(-1 + \frac{c_{ni}}{\pi_{ni}} \right) x_n^T = \sum_{n=1}^N (c_{ni} - \pi_{ni}) x_n^T.$$

1.5. iterative algorithm:

1. Randomly initialize \varPhi and Θ . Compute π_{nk} and c_{nn} .
2. Repeat until convergence ($\Delta L < \text{epsilon}$).

E-step: update π_{nk} and c_{nn}

M-step: update Θ and \varPhi using gradients obtained in 1.4. (Θ and \varPhi could be updated with gradient descent algorithm).

Compute L . If $\Delta L < \text{epsilon} \Rightarrow$ Algorithm converged.

II Part.

1. Likelihood: $p(\tilde{y}|\tilde{X}, \varphi, \Theta) = \prod_{n=1}^N p(y_n|x_n, \Theta, \varphi) \cdot \prod_{m=1}^M p(y_m|x_m, \varphi, \Theta)$

$$= \prod_{n=1}^N \sum_{k=1}^K \pi_{nk} \text{Exp}(y_n|\lambda = \exp(\Theta_k^\top x_n)) \cdot \prod_{m=1}^M \text{Exp}(y_m|\lambda = \exp((\Theta z_m)^\top x_m))$$

Or we can also write it as:

$$\prod_{n=1}^N \sum_{k=1}^K \pi_{nk} \text{Exp}(y_n|\lambda = \exp(\Theta_k^\top x_n)) \prod_{m=1}^M \sum_{k=1}^K \pi_{mk}^{z_{mk}} (\text{Exp}(y_m|\lambda = \exp(\Theta_k^\top x_m)))^{z_{mk}}$$

log-Likelihood:

$$\begin{aligned} \log p(\tilde{y}|\tilde{X}, \varphi, \Theta) &= \sum_{n=1}^N \log \sum_{k=1}^K \pi_{nk} p(y_n|x_n, \Theta_k = \Theta z_n) + \\ &+ \sum_{m=1}^M \log p(y_m|x_m, \Theta_m = \Theta z_m) = \\ &= \sum_{n=1}^N \log \sum_{k=1}^K \pi_{nk} \text{Exp}(y_n|\lambda = \exp(\Theta_k^\top x_n)) + \sum_{m=1}^M \log \text{Exp}(y_m|\lambda = \exp((\Theta z_m)^\top x_m)) \\ &= \sum_{n=1}^N \log \sum_{k=1}^K \pi_{nk} \text{Exp}(y_n|\lambda = \exp(\Theta_k^\top x_n)) - \\ &- \sum_{m=1}^M (\Theta z_m)^\top x_m - \exp((\Theta z_m)^\top x_m) y_m. \end{aligned}$$

2. $\frac{d}{d\varphi_i} \log p(\tilde{y}|\tilde{X}, \varphi, \Theta) = \frac{d}{d\varphi_i} \sum_{n=1}^N \log \sum_{k=1}^K \pi_{nk} \text{Exp}(y_n|\lambda = \exp(\Theta_k^\top x_n)) =$

$$= \sum_{n=1}^N (\pi_{ni} - \bar{\pi}_{ni}) x_n^\top$$

$$\frac{d}{d\theta_i} \log p(\tilde{y} | \tilde{X}, \varphi, \theta) = \sum_{n=1}^N z_{ni} \cdot (1 - \exp(\theta_i^T x_n)) x_n^T + \\ + \sum_{m=1}^M \mathbb{I}(z_{mi}=1) \cdot (1 - y_m \exp((\theta^T z_m)^T x_m)) \cdot x_m^T.$$

3. The derivatives $\frac{d}{d\varphi_i}$ do not change. The derivative $\frac{d}{d\theta_i}$

is using new data points. Though the individual expert is a linear model — the overall model is non-linear (or soft piece-wise linear).

12

$$\{x_i\}_{i=1}^N, x_i \in \mathbb{R}^d, E x_i = \bar{O}, S = U \Lambda U^T$$

2.1.

$$a) z_{ni} = [z_n = U_K^T \tilde{x}_n = U_K^T x_n] = u_i^T x_n,$$

with u_i — i -th column of U , U_K consists of first K columns of U .

b) Empirical Mean

$$E z_{ni} = E u_i^T x_n = u_i^T E x_n = \bar{O}$$

c) Empirical Variance.

Let's find covariance C_{ij} and $\text{Variance}(z_i) = C_{ii}$.

$$C_{ij} = \frac{1}{N} \sum_{n=1}^N z_n^{(i)} z_n^{(j)} = \frac{1}{N} \sum_{n=1}^N u_i^T x_n x_n^T u_j = \left[\frac{1}{N} \sum_{n=1}^N x_n x_n^T = S \right] =$$

$$= u_i^T S u_j$$

So Variance(z_i) = $u_i^T S u_i$

d) Variance(z_i) = $u_i^T S u_i = u_i^T U \Lambda U^T u_i =$

$$= \left[\begin{array}{l} U^T U = U U^T = I \Rightarrow \\ \Rightarrow \underbrace{U^T}_{K \times K} \underbrace{U \Lambda U^T}_{D \times D} \underbrace{U_k}_{K \times K} = \Lambda_K \in \mathbb{R}^{K \times K} \text{ with diagonal elements of } 1 \end{array} \right] = \lambda_i.$$

e) The proportion of explained variance is equal to:

$$\frac{\sum_{i=1}^K \lambda_i}{\sum_{j=1}^D \lambda_j}.$$

So K should be the smallest possible such that it satisfies:

$$\frac{\sum_{i=1}^K \lambda_i}{\sum_{j=1}^D \lambda_j} \geq 0.99.$$

2.2.

$$C = \frac{1}{N} \sum_{n=1}^N z_n z_n^T = \frac{1}{N} \sum_{n=1}^N U^T x_n \cdot (U x_n)^T = \frac{1}{N} \sum_{n=1}^N U^T x_n x_n^T U =$$

$$= U^T \left(\frac{1}{N} \sum_{n=1}^N x_n x_n^T \right) U U^T S U = U^T U \Lambda U^T U = \left[U \text{-orthogonal} \Rightarrow U^T = U^{-1} \right] = \Lambda$$

$\Rightarrow c_{ij} = 0$ for $i \neq j \Rightarrow$ dimensions are de-correlated

Q. 3. $E z_{ni} = 0$ and Variance $z_{ni} = \lambda_i$

Then for $\tilde{z}_{ni} = \sqrt{\frac{\gamma}{\lambda_i}} z_{ni} + m$:

$$E \tilde{z}_{ni} = \sqrt{\frac{\gamma}{\lambda_i}} E z_{ni} + E m = m$$

$$\text{Var } \tilde{z}_{ni} = \frac{\gamma}{\lambda_i} \text{Var } z_{ni} = \gamma.$$