

## PML Theory Questions:

Q 1: (a) If we consider the product of the likelihood & prior:

$$\begin{aligned} p(x|C) &= \frac{p(C|x)p(x)}{p(C)} \propto p(C|x)p(x) \\ &= \left( \frac{N!}{\prod_{k=1}^K n_k!} \prod_{k=1}^K x_k^{n_k} \right) \cdot \left( \frac{\Gamma(\sum \alpha_k)}{\prod \Gamma(\alpha_k)} \prod_k x_k^{\alpha_k-1} \right) \\ &\propto \prod_k x_k^{n_k + \alpha_k - 1} \end{aligned}$$

which we recognize as the functional form of the Dirichlet( $x|\alpha+n$ )

whence Dirichlet is conjugate prior of the multinomial (in)  
(i.e.  $x|C \sim \text{Dir}(-, \alpha+n)$  with  $n := (n_1, \dots, n_K)$ )

(b) WLOG let  $i=1, j=2$

we aggregate by marginalizing out  $j=2$  & summing  $x_1, x_2$

So let  $y := x_1 + x_2$ . Then the pdf is given by:

$$\begin{aligned} p(x_1, x_2, x_3, \dots, x_K) &= \int_0^y \frac{\Gamma(\sum \alpha_k)}{\prod \Gamma(\alpha_k)} z^{\alpha_1-1} (y-z)^{\alpha_2-1} \left( \prod_{j=3}^K x_j^{\alpha_j-1} \right) dz \\ &\propto \left( \int_0^y z^{\alpha_1-1} (y-z)^{\alpha_2-1} dz \right) \prod_{k=3}^K x_k^{\alpha_k-1} \end{aligned}$$

the <sup>above</sup> integral is evaluated using the substitution  $Q := \frac{z}{y}$

$$\Rightarrow \mathcal{I} = \int_0^1 y^{\alpha_1 + \alpha_2 - 1} Q^{\alpha_1-1} (1-Q)^{\alpha_2-1} dQ$$

which we recognize as the beta integral  $\frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)}$

$$\Rightarrow p(x_1, x_2, x_3, \dots, x_K) \propto (x_1 + x_2)^{\alpha_1 + \alpha_2 - 1} x_3^{\alpha_3-1} \dots x_K^{\alpha_K-1}$$

$$\Rightarrow (x_1 + x_2, x_3, \dots, x_K) \sim \text{Dir}(\alpha_1 + \alpha_2, \alpha_3, \dots, \alpha_K)$$

(in)

• The cross entropy loss is given by: (part (c))

$$CE = \underset{\underline{x}}{\operatorname{argmax}} \log p(C|\underline{x}) = \underset{\underline{x}}{\operatorname{argmax}} \log \text{likelihood}$$

we have

$$\log p(C|\underline{x}) = \log N! + \sum_{k=1}^K (\log(x_k)) n_k - \sum_{k=1}^K \log n_k$$

in order to find the  $\operatorname{argmax}$ , one must employ Lagrange multipliers and optimize  $\log p(C|\underline{x}) + \lambda \left(1 - \sum_{k=1}^K x_k\right) =: \mathcal{F}(\underline{x}, \lambda)$

$$\text{then } \partial_{x_k} \mathcal{F} = \frac{n_k}{x_k} - \lambda \quad \forall k \in \{1, \dots, K\}$$

$$\text{optimizing } (\partial_{x_k} \mathcal{F} = 0) \Rightarrow x_k = \frac{n_k}{\lambda} \quad \forall k$$

$$\Rightarrow \sum x_k = \frac{\sum n_k}{\lambda} \quad \forall k \Rightarrow \lambda = \frac{\sum n_k}{\sum x_k} = \sum n_k = N$$

$$\text{thus } \underline{x}_{\text{MLE}} = \frac{1}{N} (n_1, \dots, n_K)$$

• Our posterior is  $\mathcal{Dir}(-, \underline{\alpha} + \underline{n})$

$$\text{so } p(\underline{x}|C) = \frac{\Gamma(\sum \alpha_k + n_k)}{\prod \Gamma(\alpha_k + n_k)} \prod x_k^{\alpha_k + n_k - 1}$$

we use the same (Lagrange multiplier) trick to solve the max

$$\log p(\underline{x}|C) = \log(\sum \alpha_k + n_k) - \sum \log(\alpha_k + n_k) + \sum (\alpha_k + n_k - 1) \log x_k$$

$$\mathcal{G}(\underline{x}, \lambda) := \log p(\underline{x}|C) + \lambda (1 - \sum x_k)$$

$$\partial_{x_k} \mathcal{G} = \frac{\alpha_k + n_k - 1}{x_k} - \lambda \quad \text{Setting equal to zero gives}$$

$$x_k = \frac{\alpha_k + n_k - 1}{\lambda} \quad \text{By the same logic as above; } \rightarrow$$

$$\lambda = \sum \alpha_k + \sum n_k - K \quad (\text{summing } k \text{ expressions})$$

where the max is given by

$$\hat{x}_k = \frac{\alpha_k + n_k - 1}{\sum \alpha_k + n_k - K} \quad \forall k \in \{1, \dots, K\}$$

The two forms are very similar but the posterior takes into account the prior information (parametrized by  $\underline{\alpha}$ )

(d) (i) we have  $\text{MLE}(\underline{x} | C_1) = (0, 0, 0, 0, 1)$   
 $\text{MLE}(\underline{x} | C_2) = \left\{ \frac{1}{410} (1, 0, 12, 43, 354) \right\}$

(ii)  $p(\underline{x} | C_1, \underline{\alpha}) = D(\underline{x} | (1, 1, 1, 1, 4))$

$$p(\underline{x} | C_2, \underline{\alpha}) = D\left(\underline{x} \mid \frac{2}{415}, \frac{1}{415}, \frac{13}{415}, \frac{44}{415}, \frac{355}{415}\right)$$

we have  $p(C_{n+1, \vec{z}} = (0, 0, 0, 0, 1) | \underline{x}, C, \underline{\alpha}) = x_5 \quad \forall i \in \{1, 4\}$   
 (likelihood of single event is just probability of that event)

$$\Rightarrow p((0, 0, 0, 0, 1) | C, \underline{\alpha}, \underline{x}) = \int_{[0,1]^5} p(C_{n+1} | \underline{x}) p(\underline{x} | C, \underline{\alpha}, \underline{x}) d\underline{x}$$

$$= \int_{[0,1]^5} x_5 \mathcal{D}(\underline{x} | \underline{\alpha} + \underline{n}) d\underline{x} = \mathbb{E}_{\underline{x} \sim D(\underline{n} + \underline{\alpha})}(x_5)$$

if we apply the aggregation property to indices 1-4 of  $\underline{x} = (x_1, \dots, x_5)$

we get  $(x_1 + x_2 + x_3 + x_4, x_5) \sim \text{Dir}(\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4, \alpha_5) \xrightarrow{\text{Beta}} \text{Beta}\left(\sum_{j=1}^4 \alpha_j, \alpha_5\right)$



our probability then reduces to  $\frac{\alpha_5}{\sum \alpha_j}$  (Beta integral has ~~known~~ solution)  
(and for 4-stars,  $\frac{\alpha_4}{\sum \alpha_j}$ )

Thus for  $C_1$ , we have:

$$p(C_{n+1} = 5 \star s \mid \text{data, prior}) = \frac{4}{8} = \frac{1}{2}$$

$$p(C_{n+1} \in \{4 \star s, 5 \star s\} \mid \text{info}) = \frac{1}{2} + \frac{1}{8} = \frac{5}{8}$$

and for  $C_2$ :

$$p(C_{n+1} = 5 \star s \mid \text{info}) = \frac{355}{415}$$

$$p(C_n \in \{4 \star s, 5 \star s\} \mid \text{info}) = \frac{355}{415} + \frac{44}{415} = \frac{399}{415}$$

So  $C_2$  seems a better choice with a reasonable prior!