

HW1

1. $\frac{\partial E(v, h | W, b, c)}{\partial b_i}$

$$E(v, h | W, b, c) = -b^T v - c^T h - b^T W c$$

$$\begin{aligned} \therefore c^T h &= \sum_{i=1}^n c_i h_i &= -\sum b_i v_i - \sum c_i h_i - [b_1, b_2, \dots, b_n] \begin{bmatrix} w_{11} & \dots & w_{1n} \\ w_{21} & & \\ \vdots & & \\ & & w_{nn} \end{bmatrix} \begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix} \\ & &= -\sum b_i v_i - \sum c_i h_i - \left[\sum_1 b_1 w_{11}, \sum_1 b_2 w_{12}, \dots \right] \begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix} \end{aligned}$$

$$= -\sum b_i v_i - \sum c_i h_i -$$

$$\begin{aligned} b^T W c &= \cancel{\left[\sum_{i=1}^n b_i w_{ni} \right]} \\ &= \left[\underbrace{\sum_{i=1}^n b_i w_{i1}}_{\text{scalar value}}, \underbrace{\sum_{i=1}^n b_i w_{i2}}_{\text{row vector}}, \dots, \underbrace{\sum_{i=1}^n b_i w_{in}}_{\text{row vector}} \right] \times \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix} \end{aligned}$$

$1 \times n \times n \times 1 = 1 \rightarrow \text{scalar}$

$$\begin{aligned} \text{Scalar} \rightarrow b^T W c &= c_1 \cdot \sum_{i=1}^n b_i w_{i1} + c_2 \cdot \sum_{i=1}^n b_i w_{i2} + \dots + c_n \cdot \sum_{i=1}^n b_i w_{in} \\ &= \sum_{i=1}^n \sum_{j=1}^n c_j b_i w_{ij} \end{aligned}$$

$$A) \frac{\partial E(v, h, w, b, c)}{\partial b_i} = \frac{\partial (-\sum b_i v_i - \sum c_i h_i - \sum \sum c_j b_i w_{ij})}{\partial b_i}$$

$$= -v_i - 0 - c_j w_{ij}$$

$$B) \frac{\partial E(v, h, w, b, c)}{\partial c_i} = \frac{\partial (-\sum b_i v_i - \sum c_i h_i - \sum \sum c_j b_i w_{ij})}{\partial c_i}$$

$$= -0 - h_i - b_i w_{ij}$$

$$C) \frac{\partial E(v, h, w, b, c)}{\partial w_{ij}} = \frac{\partial (-\sum b_i v_i - \sum c_i h_i - \sum \sum c_j b_i w_{ij})}{\partial w_{ij}}$$

$$= -0 - 0 - c_j b_i$$

D) The dimension of

$$1) \frac{\partial E(v, h, w, b, c)}{\partial w} \Rightarrow \frac{\partial (-b^T v - c^T h - b^T W c)}{\partial w}$$

$$\therefore \frac{\partial a^T x b}{\partial x} \dots \Rightarrow -0 - 0 - b c^T$$

$$\Rightarrow -b c^T$$

$b \rightarrow$ column vector $= n \times 1$ / $c \Rightarrow$ Transpose of column vector
 \downarrow
 $\Rightarrow 1 \times n$

$$\therefore \text{Dimension of } \frac{\partial E}{\partial w} = n \times 1 \times 1 \times n = \underline{\underline{n \times n}}$$

$$2) \frac{\partial E(v, h; W, b, c)}{\partial b} = \frac{\partial (-b^T v - c^T h - b^T W c)}{\partial b}$$

$$\Rightarrow -v - 0 - Wc$$



$n \times 1$

$$\hookrightarrow n \times n \times n \times 1$$

$$\Rightarrow n \times 1$$

$$\therefore \text{Dimension of } \frac{\partial E}{\partial b} \Rightarrow \underline{n \times 1}$$

$$3) \frac{\partial E(v, h; W, b, c)}{\partial c} = \frac{\partial (-b^T v - c^T h - b^T W c)}{\partial c}$$

$$= -0 - h - b^T W$$



$n \times 1$



$$1 \times n \times n \times n = 1 \times n$$

To add 2 vectors of different dimensions, you can take the transpose of one, which will get you same dimension vector. if they have same no. of elements i.e. 'n'

$$\therefore \text{The dimension of } \frac{\partial E}{\partial c} = \underline{n \times 1} \text{ or } \underline{1 \times n}$$

2. random variable X follows geometric distribution (parameter ' p ')

$$\text{probability of } X = k = (1-p)^{k-1} p$$

\downarrow
no. of trials needed before getting first success.

$$\text{Dataset } D = \{x_1, x_2, \dots, x_n\}$$

Given above details, the Likelihood function can be written as;

$$\begin{aligned} L(D|p) &= \prod_{i=1}^n (1-p)^{x_i-1} \times p \\ &= (1-p)^{x_1-1} \times p \times (1-p)^{x_2-1} \times p \times \dots \times (1-p)^{x_n-1} \times p \end{aligned}$$

$$L(D|p) = p^n (1-p)^{\sum_{i=1}^n (x_i - n)} \quad \dots (1)$$

As its easier with logarithm :- product \rightarrow sum.

\therefore Taking log on both sides;

$$\ln L(D|p) = \ln \left((1-p)^{\sum_{i=1}^n (x_i - n)} \times p^n \right)$$

$$\ln(L(D|p)) = \sum_{i=1}^n (x_i - n) \ln(1-p) + n \ln p \quad \dots (2)$$

Further to get the Maximum likelihood estimation on parameter ' p ' we will differentiate the above equation (2) and equate it to '0'.

This is because maximum of a function occurs when its derivative is equals to zero.

$$\frac{d[\ln L(D|p)]}{dp} = \frac{d \sum_{i=1}^n (x_i - n) \ln(1-p)}{dp} + \frac{d n \ln p}{dp}$$

$$\frac{d[\ln L(D|p)]}{dp} = - \sum_{i=1}^n (x_i - n) \frac{x_i}{(1-p)} + n \frac{1}{p} = 0$$

$$\therefore \frac{n}{p} = \frac{\sum_{i=1}^n (x_i - n)}{1-p}$$

$$n(1-p) = p \left(\sum_{i=1}^n (x_i - n) \right)$$

$$n - np = p \sum_{i=1}^n x_i - np$$

$$n = p \sum_{i=1}^n x_i$$

$$\therefore p = \frac{n}{\left(\sum_{i=1}^n x_i \right)}$$

$$\text{The MLE of } p \Rightarrow P = \frac{n}{\sum_{i=1}^n x_i} \approx \frac{1}{\bar{x}}$$

that is 1 success for each set of trials \bar{x}

3.

where $X \rightarrow$ follows geometric distribution with 'p' parameter

Dataset $D = \{x_1, x_2, \dots, x_n\}$

We need to find a prior distribution along with the likelihood function.

As per Bayes Theorem;

$$\Pr(p|D) = \frac{\Pr(p) \Pr(D|p)}{\Pr(D)}$$

↓ posterior
↓ prior
→ likelihood
→ Marginal.

To find the posterior distribution on parameter ' p ', we can choose any prior distribution, which might lead to expensive numerical computation. Thus, to avoid the massive computation for Bayesian inference it is preferred to consider a conjugate prior.

Conjugate prior are probability distribution that belongs to the same family of likelihood distribution or it simply mimics the form of likelihood.

Beta distribution, is known as the conjugate prior of the geometric distribution that is followed by random variable X .

Set of Trials \Rightarrow geometric distribution $\Rightarrow X(k) = (1-p)^{k-1} p$

Probability of success \Rightarrow Beta distribution $\rightarrow f(p) = \frac{1}{B(\alpha, \beta)} (1-p)^{\beta-1} p^{\alpha-1}$

I would choose the Beta distribution as the prior, as it is a conjugate prior of geometric distribution and thus greatly simplifies the Bayesian Inference calculations.

The posterior distribution of 'p' for dataset 'D' can be written as;

$$\Pr(p|D) = \underbrace{\left(\prod (1-p)^{x_i-1} \times p \right)}_{\text{likelihood}} \times \underbrace{\text{Beta}(\alpha, \beta)}_{\text{prior}}$$

↓
Posterior

$\Pr(D) \rightarrow \text{Marginal.}$

This can be further simplified as; $x_i \approx k$ & marginal likelihood is discarded \rightarrow normalization constant $= 1$

$$\Pr(p|D) = \left((1-p)^{k-1} \times \underbrace{p}_{\text{Beta}(\alpha, \beta)} \times p^{\alpha-1} (1-p)^{\beta-1} \right) \times \text{constant}$$

$p \times p^{\alpha-1} = p^{\alpha-1+1} = p^\alpha$

Taking log on both sides;

$$\log(\Pr(p|D)) = (k-1) \log(1-p) + \alpha \log p + (\beta-1) \log(1-p) + \text{constant}$$

To find maximum of this function, take derivative w.r.t 'p' and equate it to '0'.

$$\frac{d \log(\Pr(p|D))}{dp} = \frac{-1 \times (k-1)}{(1-p)} + \frac{\alpha}{p} - 1 \frac{(\beta-1)}{(1-p)} = 0$$

\therefore to find 'p' ;

$$\frac{\alpha}{p} = \frac{\beta-1}{(1-p)} + \frac{k-1}{(1-p)}$$

$$\frac{\beta-1+k-1}{1-p} = \frac{\beta+k-2}{1-p}$$

$$\alpha(1-p) = p\beta + pk - 2p$$

$$\alpha - \alpha p = p\beta + pk - 2p$$

$$\alpha = \alpha p + \beta p + kp - 2p$$

$$\alpha = p(\alpha + \beta + k - 2)$$

This is the BPE for \rightarrow
the geometric distribution
given the prior as beta distribution
with parameters α & β

$$p = \frac{\alpha}{\alpha + \beta + k - 2}$$