

1. Uniform distribution of a continuous variable  $x$  is

$$p(x; a, b) = \frac{1}{b-a} \quad a \leq x \leq b$$

$$\text{now } \int_a^b \frac{1}{b-a} dx = \left| \frac{x}{b-a} \right|_a^b = \frac{b-a}{b-a} = 1$$

Therefore, the distribution is normalized.

$$\text{Mean} = E(x) = \int_a^b \frac{x}{b-a} dx = \left| \frac{x^2}{2(b-a)} \right|_a^b = \frac{b^2 - a^2}{2(b-a)} = \frac{a+b}{2}$$

$$\text{Var}(x) = E(x^2) - E(x)^2 \quad \text{Now, } E(x^2) = \int_a^b \frac{x^2}{b-a} dx = \frac{b^3 - a^3}{3(b-a)}$$

$$\begin{aligned} \text{So, } \text{Var}(x) &= \frac{b^2 + ab + a^2}{3} - \frac{(a+b)^2}{4} \\ &= \frac{b^2 + ab + a^2}{3} - \frac{a^2 + 2ab + b^2}{4} \\ &= \frac{1}{12} (4b^2 + 4ab + 4a^2 - 3a^2 - 6ab - 3b^2) \\ &= \frac{1}{12} (a^2 - 2ab + b^2) \\ &= \frac{(a-b)^2}{12} \end{aligned}$$

2. PMF of Poisson random variable,

$$p(x; \lambda) = \frac{\lambda^x e^{-\lambda}}{x!}$$

$$L(\lambda) = \prod_{i=1}^n \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} \quad \left[ \text{likelihood of i.i.d samples} \right]$$

taking log likelihood,

$$\begin{aligned} LL(\lambda) &= \sum_{i=1}^n (\log \lambda^{x_i} + \log e^{-\lambda} - \log(x_i!)) \\ &= \log \lambda \sum_{i=1}^n x_i - \sum_{i=1}^n \lambda - \sum_{i=1}^n \log(x_i!) \end{aligned}$$

now for MLE,

$$\frac{\partial LL(\lambda)}{\partial \lambda} = \frac{1}{\lambda} (x_1 + x_2 + \dots + x_n) - n = 0$$

now

$$\lambda = \frac{x_1 + x_2 + \dots + x_n}{n}$$

(1)

3.  $\text{randn}(d,1)$  generates multivariate normal  $x \in \mathbb{R}^d$  with 0 mean and  $I$  covariance.

therefore  $x \sim \mathcal{N}(0, I)$

Applying affine transformation on  $x$ , ~~for~~ in the form  $Ax + b$ .

From property of normal distribution,

$$Ax + b \sim \mathcal{N}(b, AA^T) \quad \text{--- (1)}$$

Our goal is to generate random variable from  $\mathcal{N}(\mu, \Sigma)$  --- (2)

Comparing (1) and (2)

$$\mu = b.$$

$$\Sigma = AA^T.$$

Therefore to generate the required distribution, we need to perform Affine transformation on  $x$ ,  $(Ax + \mu)$ .

$\mu \rightarrow$  given as input.

$A \rightarrow$  we can derive from  $\Sigma$  input.

Note,  $\Sigma$  is a covariance matrix, positive-definite matrix.

We can perform Cholesky Factorization to find matrix  $A$ .

4.

$$\underset{\theta}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n r_i (y_i - \phi_i^T \theta)^2$$

This can be thought of as weighted least square.

The ordinary least square,  $r_i \forall i \in \{1, n\} = 1$

Now, we write matrix  $R$ , with  $r_i$  on the diagonal. This could be rewritten,

$$\text{WMSE} = n^{-1} (Y - \Phi^T \theta)^T R (Y - \Phi^T \theta)$$

$$= \frac{1}{n} (Y^T R Y - Y^T R \Phi^T \theta - \theta^T \Phi R Y + \theta^T \Phi R \Phi^T \theta)$$

Differentiating in terms of  $\theta$ ,

$$\nabla_{\theta} \text{WMSE} = \frac{2}{n} (-\Phi R Y + \Phi R \Phi^T \theta)$$

Setting this to zero,

$$\hat{\theta} = (\Phi R \Phi^T)^{-1} \Phi R Y.$$

### 5. a) Naive Bayes,

a document  $D$ , whose class is given by  $C$ . So we choose the class  $C$ , which has the highest posterior Probability.

$$P(C|D) = \frac{P(D|C) P(C)}{P(D)} \propto P(D|C) P(C)$$

Now in bernoulli distribution, to calculate the likelihood of the data, we do not consider the frequency of the word, we only consider the availability of data.

$$P(D|C) = P(w_1|C) \dots P(w_m|C)$$

① We would take the label and data for train to calculate the prior probability.  $P(C)$  gives us the number of documents percentage.

② For each class, we would calculate the occurrence of a word across the documents. If any word ( $w_i$ ) belongs to 100 documents of class ( $C_i$ ) and if that class has total 200 documents,  $P(w_i|C_i) = \frac{1}{2}$ .  
Then the likelihood  $P(D|C) = \prod_{i=1}^m P(w_i|C_i)$

③ Now for ~~training~~ testing, we use the data generated from ① and ② to predict. If a document is present, we would use the probability of appearance and if it's absent we would consider the probability of not occurring.

$$f(D_j) = \underset{C}{\operatorname{argmax}} \prod_{i=1}^m [b_{ji} P(w_i|C) + (1-b_{ji})(1-P(w_i|C))] * P(C)$$

$b_{ji} \rightarrow$  whether  $i$ th word is present in document  $j$

$$\text{Accuracy} = 62.33\%$$

b) Similar approach to bernoulli, however while calculating likelihood and prior, we are ~~not~~ going to consider the frequencies as well.

① Prior Probability  $P(C=k) = \frac{N_k}{N}$  ; if there are  $N$  documents in the training set and  $N_k$  belongs to class  $k$

② Estimate the likelihoods 
$$P(w_j | C=k) = \frac{\sum_{i=1}^N x_{ij} Z_{ik}}{\sum_{j=1}^V \sum_{i=1}^N x_{ij} Z_{ik}}$$

$V \rightarrow$  n.o of words in histogram

$Z_{ik} = 1$  , when  $D_i$  has class  $C=k$ .

③ Test data,

$$f(D_j) = \operatorname{argmax}_c \left( \sum_{i=1}^m (x_{ji} \log P(w_i | c)) + \log P(c) \right)$$

$m \rightarrow$  n.o of words in histogram.

$x_{1j}, \dots, x_{mj} \rightarrow$  count of occurrences of  $x_1, \dots, x_m$  in  $j$ th document.  
of test data.  $\rightarrow$  accuracy 77.76%

c)  $P(C|D) \propto P(D|C) \cdot P(C)$   
 $\downarrow$   
 $D$  is a document vector ( $x$ )  $\left| \begin{array}{l} P(D|C) \text{ having normal distribution.} \\ \text{Assuming univariate normal distribution} \end{array} \right.$

$$P(D|C) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (x - \mu_c)^T \Sigma^{-1} (x - \mu_c) \right)$$

so,

$d \rightarrow$  length of histogram

$$P(C|D) = \frac{1}{N_c} \exp \left( -\frac{1}{2} (x - \mu_c)^T \Sigma^{-1} (x - \mu_c) \right) \cdot P(C)$$

now,  $f(x) = \operatorname{argmax}_c \exp \left( -\frac{1}{2} (x - \mu_c)^T \Sigma^{-1} (x - \mu_c) \right) \cdot P(C)$

taking log,

$$= \operatorname{argmin}_c (x - \mu_c)^T \Sigma^{-1} (x - \mu_c) + \operatorname{argmax}_c P(C)$$

$$= \operatorname{argmin}_c x^T \Sigma^{-1} x + \mu_c^T \Sigma^{-1} \mu_c - 2\mu_c^T \Sigma^{-1} x + b_c$$

$$= \operatorname{argmax}_c (\mu_c^T \Sigma^{-1} x + b_c)$$

$\mu, \Sigma \rightarrow$  we will get from TF-IDF, we will use the parameters to estimate

$f(m)$  using