d.2.    Poisson Distribution:- $y \sim \text{Poisson}(\lambda)$

(a)

$$\Rightarrow P(y = k_\lambda) = \frac{\lambda^k e^{-\lambda}}{k!}$$

$$\Rightarrow p(y) = \frac{\lambda^y e^{-\lambda}}{y!} \qquad -①$$

Now, if $y \sim \text{exp-family}(\eta)$

$$\Rightarrow P(y;\eta) = b(y) \, e^{(\eta y - a(\eta))} \qquad -②$$

In ①, we can write

$$\log[P(y;\lambda)] = y \log \lambda - \lambda - \log y! \qquad -③$$

Taking log in ② we get

$$\log P(y;\eta) = \log b(y) + \eta y - a(\eta) \qquad -④$$

Equating ③ & ④ we get

$$\eta = \log \lambda \qquad - (a)$$

$$a(\eta) = \lambda = e^\eta \qquad -(b)$$

$$\log b(y) = -\log y! \Rightarrow b(y) = \frac{1}{y!} \qquad -(c)$$

$$\Rightarrow y \sim \text{Poisson}(\lambda) \text{ belongs to exponential}$$

U

feunly

$$y = \theta^T x \Rightarrow \log \lambda = \theta^T x \Rightarrow \lambda = e^{\theta^T x}$$

(b)

$$\log P(y; \eta) = \log b(y) + \eta y - a(\eta)$$

substituting values for $b(y)$ $\&$ $\eta$ we get:-

$$\sum_{i=1}^{m} \log P(y^{(i)}; \theta) = \sum_{i=1}^{m} \Big[ -\log(y^{(i)}!) + \theta^T x^{(i)} y^{(i)} - e^{\theta^T x^{(i)}} \Big]$$

$$= LL(\theta)$$

$$\Rightarrow \nabla_\theta LL(\theta) = \sum_{i=1}^{m} \Big[ 0 + x^{(i)} y^{(i)} - e^{\theta^T x^{(i)}} \cdot x^{(i)} \Big]$$

$$= \sum_{i=1}^{m} \Big[ y^{(i)} - e^{\theta^T x^{(i)}} \Big] x^{(i)}$$

(c)  $$\nabla_\theta LL(\theta) = \sum_{i=1}^{m} \Big( y^{(i)} - e^{\theta^T x^{(i)}} \Big) x^{(i)}$$

$$\Rightarrow \nabla_\theta^2 LL(\theta) = \sum_{i=1}^{m} \Big[ - e^{\theta^T x^{(i)}} \Big] x^{(i)} x^{(i)T}$$

$$\Big[ \because \frac{d LL(\theta)}{d \theta_i} = \sum_{i=1}^{m} \big( y^{(i)} - e^{\theta^T x^{(i)}} \big) x_j^{(i)}$$

$$\Rightarrow \frac{d^2 LL(\theta)}{d\theta_j \, d\theta_k} = \sum_{i=1}^{m} \left(-e^{\theta^T x^{(i)}}\right) x_j^{(i)} x_k^{(i)}$$

$\Rightarrow$ Hessian matrix: $H = \sum_{i=1}^{m} \left(-e^{\theta^T x^{(i)}}\right) x^{(i)} x^{(i)T}$

$Z \in R^n$

$$\Rightarrow Z^T H Z = \sum_{i=1}^{m} \underbrace{Z^T x^{(i)} x^{(i)T} Z}_{\left[-e^{\theta^T x^{(i)}}\right]}$$

$$= \sum_{i=1}^{m} \underbrace{\left[Z^T x^{(i)}\right]^2}_{\geq 0} \underbrace{\left[e^{\theta^T x^{(i)}}\right](-1)}_{\Downarrow \atop > 0}$$

$$\Rightarrow Z^T H Z \leq 0 \Rightarrow H \text{ is } -\text{ve} \atop \text{semi-definite}$$

$\Rightarrow LL(\theta)$ is concave function of $\theta$.

Hence, proved

0.2. (i) The model is overfitting. The improvement advice would be to try a simpler model first. Alternatively, if additional data is available they can train original model with additional data to see if generalization accuracy improves.

(ii) The model seems to be doing well in both training as well as generalization

Student can try to improve [...] [...]
by trying a more sophisticated model
for possible improvement in [...] val
accuracy. This might require training
with additional data to prevent overfitting

(III) Both training & validation accuracy are
low. Model is underfitting & student
should increase the complexity of
their model to better represent
patterns in underlying data.

(b)

$$\theta_{y=k} = \frac{\sum_{i=1}^{m} 1\{y^{(i)}=k\} 1\{x_j^{(i)}=1\} + C \cdot 1}{\sum_{i=1}^{m} 1\{y^{(i)}=k\} + C \cdot L}$$

c controls overfitting because if does not
allow the model to overfit the
training data in extreme noisy scenarios. For example
(by having a prior)
if frequency count of some word $w_c$ is
0 is in training for some class
$y=k$, the introduction of c term
makes sure that the probability
of $(y=k) \neq 0$ even if $w_c$ does not
appear in the document. Similarly, if frequency
[...] [...] [...] to

count of we is very very low compared t... other words. For some class $j=k$, the introduction of c term (smoothing) ensures that the probability is increased by adding a prior in form of c term in numerator (& denominator). Larger the value of c, larger the mean of this term is. Overfitting is prevented by increasing the effective count for each attribute (word) by c, for each class $j=k$, by effectively adding a document in which each word we appears c times. By carefully choosing this mechanism, a uniform prior over words (attributes) is introduced reducing overfitting.

Large value of c:- Model is too heavily biased by prior & it will tend to underfit

Small value of c:- prior is very weak & model might start overfitting.

(a)

$$\hat{\Sigma}_1 = \begin{bmatrix} \sigma_{11}^2 & & \bigcirc \\ & \ddots & \\ \bigcirc & & \sigma_{1n}^2 \end{bmatrix}$$

variances of individual components

Assume (w.l.o.g)

$x^{(i)} \in R^n$

cross diagonal entries are zero, since attribute are independent (given the class) & hence un-correlated with each other. (i.e covariance

$$\overline{\text{...terms will be zero}}.$$

Similarly
$$\hat{\Sigma}_2 = \begin{bmatrix} \sigma_{21}^2 & & & 0 \\ & \ddots & & \\ 0 & & & \sigma_{2n}^2 \end{bmatrix}$$

$$\log\left[\prod_{i=1}^{m} P(y^{(i)}, x^{(i)}; \theta)\right]$$

$$= \sum_{i=1}^{m} \log\left[P(y^{(i)}; \phi) \left[p(x^{(i)}/y^{(i)}; \theta)\right]\right]$$

$$= \sum_{i=1}^{m} 1\{y^{(i)} = k\}[\phi_k] \quad e^{-\frac{[x^{(i)} - \mu_{y^{(i)}}]^T \hat{\Sigma}_{y^{(i)}}^{-1}(x^{(i)} - \mu_{y^{(i)}})}{2}}$$
$$+ \log \frac{1}{(2\pi)^{n/2}|\hat{\Sigma}_{y^{(i)}}|^{1/2}}$$

Since $\hat{\Sigma}_{y^{(i)}}$ is diagonal $\forall (i) \{y^{(i)} \varepsilon \{1 - - n\})$

we can decompose above expression as: -
$$\left(\mu_k \equiv \mu_{k1} - - \mu_{kn}\right)$$

$$= \sum_{i=1}^{m} \sum_{k=1}^{r} 1\{y^{(i)} = k\} \log \phi_k$$
$$+ \sum_{i=1}^{m} \left[\log \prod_{j=1}^{n}\left(\frac{1}{(2\pi) \sigma_{y^{(i)}j}^2}\right)\right] + \log e^{\sum_{j=1}^{n} - \left[\frac{(x_j^{(i)} - \mu_{y^{(i)}j})^2}{2\sigma_{y^{(i)}j}^2}\right]}$$

$$= \sum_{i=1}^{m} \sum_{k=1}^{r} 1\{y^{(i)} = k\} \log \phi_k + \sum_{i=1}^{m}\left[\sum_{j=1}^{n} \log \frac{1}{\sqrt{2\pi} \sigma_{y^{(i)}j}^2}\right.$$

$$+ \sum_{j=1}^{n} (-1/2)\left[\left(x_j^{(i)} - \mu_{y^{(i)}j}\right)\right]$$

Part (b)

$$\cdots \sum_{j=1}^{?} \left(\cdots\right) \left[ \frac{\cdots}{\sigma_{y h_j}{}^2} \right] \cdots$$

(i) Differentiating wrt $\mu_{kj}$

$$\frac{\partial LL(\theta)}{\partial \mu_{kj}} = 0 + 0 + \left\{ (1/2) \right\} \sum_{i=1}^{m} \sum_{k=1}^{?} 1\{y^{h)} = k\} \frac{d}{d\mu_{kj}} \left[ \frac{(x^{h)}{}_j - \mu_{kj})^2}{(\sigma_{kj})^2} \right]$$

$$= \sum_{i=1}^{m} (1/2) \sum_{k=1}^{?} 1\{y^{h)} = k\} \; 2 \; \frac{(x^{h)}{}_j - \mu_{kj})}{(\sigma_{kj})^2}$$

Equating to zero, we get

$$\frac{\partial LL(\theta)}{\partial \mu_{kj}} = \sum_{i=1}^{m} \left[ \sum_{k=1}^{?} 1\{y^{h)} = k\} (x^{h)}{}_j - \mu_{kj}) \right] = 0$$

$$\Rightarrow \mu_{kj} = \frac{\sum_{i=1}^{m} 1\{y^{h)} = k\} x^{hi}_j}{\sum_{i=1}^{m} 1\{y^{h)} = k\}}$$

$$\Rightarrow \mu_k = \frac{\sum_{i=1}^{m} 1\{y^{h)} = k\} x^{h)}}{\sum_{i=1}^{m} 1\{y^{h)} = k\}}$$

$\Rightarrow$ which is same as derived
in class.

(ii) Differentiating wrt $\sigma_{kj}$, we get

$$\frac{\partial LL(\theta)}{\cdots} = \frac{d}{\cdots} \left[ \sum_{i=1}^{m} 1\{y^{h)} \in \{1, 3 \cdots \gamma-1\}\} \cdots \right]$$

$$\frac{\partial}{\partial \sigma_{1\delta}} \quad \frac{\partial \|\theta\|}{\partial \theta}\Big|_{\theta=L_{m}} \sum_{\delta=1}^{m} \log \frac{1}{\sqrt{(2\pi \sigma_{1\delta})}} + \sum_{\delta=1}^{2} - \frac{(-1)}{2\sigma_{1\delta}^{2}}$$

$$= \frac{\partial}{\partial \sigma_{1\delta}} \sum_{r=1}^{m} \Bigg[ 1\{y^{(r)} \text{ is odd}\}$$

$$\cdot \sum_{\delta=1}^{n} \log \frac{1}{\sqrt{2\pi}} - \log(\sigma_{1\delta})$$

$$+ \sum_{\delta=1}^{n} - \frac{(x_{\delta}^{(r)} - \mu_{y^{(r)}\delta})^{2}}{(2\sigma_{1\delta}^{2})} \Bigg]$$

$$= \sum_{r=1}^{m} 1\{y^{(r)} \text{ is odd}\}\Bigg\{ \frac{(-1)}{(\sigma_{1\delta})}$$

$$+ \frac{(x_{\delta}^{(r)} - \mu_{y^{(r)}\delta})^{2}}{2(\sigma_{1\delta})^{3}} \cdot 2 \Bigg\}$$

Equating to zero we get:-

$$\sigma_{1\delta}^{2} = \frac{\sum\limits_{r=1}^{m} 1\{y^{(r)} \text{ is odd}\} (x_{\delta}^{(r)} - \mu_{y^{(r)}\delta})^{2}}{\sum\limits_{i=1}^{m} 1\{y^{(r)} \text{ is odd}\}}$$

→ Empirical variance computed across classes which chose the $\Sigma$ parameter

→ same as expression derived in the

for the general case

$$\hat{\Sigma}_j = \frac{\sum_{i=1}^{m} \mathbb{1}\{y^{(i)} \text{ is odd}\} (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T}{\sum_{i=1}^{m} \mathbb{1}\{y^{(i)} \text{ is odd}\}}$$

Since cross diagonal entries are zero
(due to naive Bayes assumption)

we get

$$\sigma_{1\delta}^2 = \frac{\sum_{i=1}^{m} \mathbb{1}\{y^{(i)} \text{ is odd}\} (x_\delta^{(i)} - \mu_{y^{(i)}\delta})^2}{\sum_{i=1}^{m} \mathbb{1}\{y^{(i)} \text{ is odd}\}}$$

Hence the two expressions are identical

Similarly for $\Sigma_2$ i.e.

$$\sigma_{\delta 2}^t = \frac{\sum_{i=1}^{m} \mathbb{1}\{y^{(i)} \text{ is even}\} (x_\delta^{(i)} - \mu_{y^{(i)}\delta})^2}{\sum_{i=1}^{m} \mathbb{1}\{y^{(i)} \text{ is even}\}}$$

## Q4.
### (a)

$\overline{\theta^T A \theta}$ :-

Following derivation works for both when A is symmetric or not symmetric.

$$\frac{\partial \theta^T A \theta}{\partial \theta_j} = \frac{\partial \sum_{k,k'} \theta_k A_{kk'} \theta_{k'}}{\partial \theta_j}$$

where $k = j$ or $k' = j$ will remain.

$$\frac{\partial}{\partial \theta_j} \left[ \sum_{k=j, k' \neq j} \theta_j A_{jk'} \theta_{k'} + \sum_{k \neq j, k'=j} \theta_k A_{kj} \theta_j \right.$$
$$\left. + \sum_{k=j, k'=j} \theta_j A_{jj} \theta_j \right]$$

$$= \sum_{k' \neq j} A_{jk'} \theta_{k'} + \sum_{k \neq j} \theta_k A_{kj}$$

$$+ 2 \theta_j A_{jj}$$

$\overline{A^T := \text{transpose of } A}$

$$= \sum_{k'} A_{jk'} \theta_{k'} + \sum_{k} \theta_k A_{kj} = \sum_{k'} A_{jk'} \theta_{k'}$$
$$+ \sum_{k} A^T_{jk} \theta_k$$

$$= (A\theta + A^T\theta)_j$$

$$\Rightarrow \nabla_\theta \theta^T A \theta = (A + A^T) \theta$$

$$\Rightarrow$$

$$\frac{\partial^2 \theta^T A \theta}{\partial \theta_j \partial \theta_\ell} = \frac{\partial}{\partial \theta_\ell} \left( \sum_{k'} A_{jk'} \theta_{k'} + \sum_{k} A^T_{jk} \theta_k \right)$$

$$= (A_{j\ell} + A^T_{j\ell})$$

$$\Rightarrow \nabla^2_\theta \theta^T A \theta = \overline{(A + A^T)}$$

$$f(\theta) = \theta^T A \theta + a^T \theta + b$$

$$\Rightarrow \text{Newton's update:-}$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - H^{-1} \nabla_\theta f(\theta)$$

$$= \theta^{(t)} - (A + A^T)^{-1} \left[ (A + A^T)\theta^{(t)} + a \right]$$

$$= \cancel{\theta^{(t)}} - \cancel{I\theta^{(t)}} - (A + A^T)^{-1} a$$

$$= -(A + A^T)^{-1} a \quad \text{(constant wrt } \theta)$$

$\Rightarrow$ Newton's method converges in single iteration

**(b)**
(i) GD:-

$$\theta^{(t+1)} \leftarrow \theta^{(t)} + \eta \cdot \frac{1}{m} \nabla_\theta LL(\theta)$$

$$\left[ \theta^{(t+1)} - \theta^{(t)} \right] \leftarrow \eta \cdot \frac{1}{m} \nabla_\theta LL(\theta)$$

SGD:-

$$\theta^{(t+1)} \leftarrow \theta^{(t)} + \eta \cdot \frac{1}{\gamma} \nabla_\theta LL_b(\theta)$$

$$\Rightarrow \theta^{(t+1)} - \theta^{(t)} \leftarrow \eta \cdot \frac{1}{\gamma} \nabla_\theta LL_b(\theta)$$

$$\mathbb{E}_{\{x^{(i)}, y^{(i)}\}_{i=1}^\gamma} \left[ \theta^{(t+1)} - \theta^{(t)} \right] \leftarrow \eta \cdot \frac{1}{\gamma} \mathbb{E}_{\{x^{(i)}, y^{(i)}\}_{i=1}^\gamma}$$

$\forall i :\sim (x^{(i)}, y^{(i)}) \sim D$    $\forall i. (x^{(i)}, y^{(i)}) \sim D$

D:- Training data distribution    $\left[ \nabla_\theta LL_b(\theta) \right]$

RHS

$$= \eta \cdot \frac{1}{\gamma} \mathbb{E}_{\{x^{(i)}, y^{(i)}\}_{i=1}^m} \nabla_\theta \left[ \sum_{i=1}^{\gamma} LL_i(\theta) \right]$$

$\forall i : (x^{(i)}, y^{(i)}) \sim D$

where $LL_i(\theta)$ denotes the log-likelihood computed over the $i$th example in mini-batch

$$= \eta \cdot \frac{1}{\delta} \quad \nabla_\theta \; E_{\{x^{(i)}, y^{(i)}\}_{i=1}^{\delta}} \; \sum_{i=1}^{\delta} LL_i(\theta)$$
$$\forall i : (x^{(i)}, y^{(i)}) \sim D$$

$$= \eta \cdot \frac{1}{\delta} \quad \nabla_\theta \sum_{i=1}^{\delta} E_{(x^{(i)}, y^{(i)}) \sim D} \left[ LL_i(\theta) \right]$$

$$= \eta \cdot \frac{1}{\delta} \quad \nabla_\theta \sum_{i=1}^{\delta} \left[ \frac{1}{m} \sum_{\ell=1}^{m} LL_\ell(\theta) \right]$$

$$= \eta \cdot \frac{1}{\delta} \nabla_\theta \; \delta \; \frac{1}{m} \sum_{\ell=1}^{m} LL_\ell(\theta)$$

$$\underbrace{\qquad\qquad}_{LL(\theta)}$$

$$= \eta \cdot \nabla_\theta \; LL(\theta)$$
$$\Rightarrow \text{Same as } GD \quad \underline{\textbf{update}}$$

(ii)

In SGD:-

$$\theta^{(t+1)} \leftarrow \theta^{(t)} + \eta \cdot \frac{1}{\delta} \; \nabla_\theta \; LL_s(\theta)$$

$$\theta^{(t+1)} - \theta^{(t)} \leftarrow \eta \cdot \frac{1}{\delta} \; \nabla_\theta LL_s(\theta)$$

$$\text{Var}\left( \theta^{(t+1)} - \theta^{(t)} \right) \leftarrow \eta^2 \frac{1}{\delta^2} \; \text{Var}\left[ \nabla_\theta \; LL_s(\theta) \right]$$

$$RHS = \eta^2 \frac{1}{\delta^2} \; \text{Var}\left[ \nabla_\theta LL_s(\theta) \right]$$

$$= \eta^2 \frac{1}{\delta^2} \; \text{Var}\left[ \nabla_\theta \sum_{i=1}^{\delta} LL_i(\theta) \right]$$

$\dfrac{LL_i(\theta)}{\Downarrow}$
Log-likelihood
for $i$th
example

where the variance is computed wrt
distribution $\{x^{(i)}, y^{(i)}\}_{r=1}^{\delta}$, $(x^{(i)}, y^{(i)}) \sim D$
(training data dist)

$$\eta^2 \frac{1}{\delta^2} \; \text{Var}\left[ \sum_{i=1}^{\delta} \nabla_\theta \; LL_i(\theta) \right]$$

$$\underbrace{\phantom{xxxx}}_{\substack{\gamma \text{ independent} \\ \text{random variables}}}$$

$$\Rightarrow \quad \gamma^2 \frac{1}{\gamma^2} \sum_{i=1}^{\gamma} \mathrm{Var}\left[\nabla_\theta LL_i(\theta)\right]$$

But $\quad \mathrm{Var}\left[\nabla_\theta LL_i(\theta)\right] = \mathrm{Var}\left[\nabla_\theta LL_1(\theta)\right]$

$\qquad$ By definition (each of $\nabla_\theta LL_i(\theta)$ are i.i.d.)

$$\Rightarrow \quad \mathrm{Var}\left(\theta^{(t+1)} - \theta^{(t)}\right) = \gamma^2 \frac{1}{\gamma^2} \times \gamma \; \mathrm{Var}\left[\nabla_\theta LL_1(\theta)\right]$$

$$= \gamma^2 \frac{1}{\gamma} \; \mathrm{Var}\left[\nabla_\theta LL_1(\theta)\right]$$