CSC 411 Fall 2018
Machine Learning and Data Mining

Homework 4

Family name: Bhatia          Given name: Pooja

# Solution 1

(a)

Input:

$W \times H \times C$      Kernal size K

                 Output Maps M

Units = pixels × output units

$\qquad = WH \times M = WHM$

Weights = size of each piece × No of filters
                     of input

$\qquad\qquad = K^2 C \times M$

$\qquad\qquad = K^2 C M$

Connections = values processed × output units
                   by a kernal

$\qquad\qquad = K^2 C \times WHM$

$\qquad\qquad = K^2 C\, WHM.$

## Layer 1

$M = 96$ Kernals $\rightarrow$ $\overset{K}{11} \times 11 \times 3$ each

stride $= 4$

Input $= \overset{W}{224} \times \overset{H}{224} \times \overset{c}{3}$

New Iput: $W = H = 55$ $\quad c = 3$

units $= 55 \times 55 \times 96 = \boxed{290,400}$

weight $= K^2 CM = 11^2 \times 3 \times 96 = \boxed{34,848}$

Connections $= K^2 CMWH = 11^2 \times 3 \times 96 \times 55 \times 55$

$\qquad\qquad\qquad\qquad = \boxed{105,415,200}$

## Layer 2

$M = 256$ $\quad$ size $\rightarrow 5 \times 5 \times 48 = K^2 C$

$W = H = 55/2 = 27$ (max pooling)

units $= WHM = 27^2 \times 256 = \boxed{186,624}$

weight $= K^2 CM = 5^2 \times 48 \times 256 = \boxed{307,200}$

Connections $= K^2 C WHM = 5^2 \times 48 \times 27^2 \times 256$

$\qquad\qquad\qquad = \boxed{223,948,800}$

## Layer 3

$M = 384$ $\quad K^2C = 3 \times 3 \times 256$

$W = H = 27/2 = 13$ (max pooling)

Units $= WHM = 13^2 \times 384 = \boxed{64,896}$

Weights $= K^2CM = 3 \times 3 \times 256 \times 384$

$\qquad = \boxed{884,736}$

Connections $= K^2C\,WHM = 3^2 \times 256 \times 13^2 \times 384$

$\qquad = \boxed{149,520,384}$

## Layer 4

$M = 384$ $\quad K^2C = 3 \times 3 \times 192$

units $= 13^2 \times 384 = \boxed{64,896}$

weight $= 3^2 \times 192 \times 384 = \boxed{663,552}$

Connection $= 3^2 \times 192 \times 13^2 \times 384$

$\qquad = \boxed{112,140,288}$

## layer 5

$$M = 256 \qquad K^2C = 3 \times 3 \times 192$$

units $= 13^2 \times 256 = \boxed{43,264}$

weight $= 3^2 \times 192 \times 256 = \boxed{442,368}$

connection
$$= 3^2 \times 192 \times 13^2 \times 256$$
$$= \boxed{74,760,192}$$

## layer 6

units $= 4096$

weight $= 6^2 \times 256 \times 4096 = \boxed{37,748,736}$

connection $= 6^2 \times 256 \times 4096 = \boxed{37,748,736}$

## layer 7

unit $= 4096$

weight $= 4096 \times 4096 = \boxed{16,777,216}$

connection $= 4096^2 = \boxed{16,777,216}$

## layer 8

Units = $\boxed{1000}$

weight = $4096 \times 1000 = \boxed{4,096,000}$

Connection = $4096 \times 1000 = \boxed{4,096 \times 000}$

| layer | #Units | #weights | #Connections |
|---|---|---|---|
| convolution Layer1 | 290,400 | 34,848 | 105,415,200 |
| Convolution layer2 | 186,624 | 307,200 | 223,948,800 |
| Convolution layer3 | 64,896 | 884,736 | 149,520,384 |
| Convolution layer4 | 64,896 | 663,552 | 112,140,288 |
| Convolution layer5 | 43,264 | 442,368 | 74,760,192 |
| Fully connected layer 1 | 4096 | 37,748,736 | 37,748,736 |
| Fully connected layer2 | 4096 | 16,777,216 | 16,777,216 |
| Output Layer | 1000 | 4,096,000 | 4,096,000 |

Solution 1.

(b)

(i)

Fully connected layers contribute to the most number of paramets i.e. the weights.

$$weights = K^2 C M$$

(hence reducing K and M would help)
Also, Alexnet works with same accuracy even if its parameters are reduced certain times.

(b)
(ii) Convolution Layers contributes to most number of connections

$$connection = K^2 C M W H$$

hence, suggestion: Reduce the No. of Kernals; or reduce size of kernals. would both help.

**Q2.** Gaussian Naive Bayes

soln

(a) $X$ is $(x_1 \ldots x_d)$ i.e $d$ feature $X$.

$X = (x_1, x_2 \ldots x_d)$, where each $X_i$ is a continuous random variable.

$Y \in (1, 2 \ldots k)$ but to start and for simplicity let's first Assume that $Y$ is Boolean and has parameter say $\pi = P(Y=1)$

Given, that for each $X_i$, $P(X_i | Y = \alpha_k)$ is a Gaussian Distribution of the form $N(\mu_{ik}, \sigma_i)$

For all $i$ and $j \neq i$, $X_i$ and $X_j$ are conditionally independent given $Y$.

$\sigma_i$ varies from attribute to attribute but doesn't depend on $Y$.

NOW,

$$P(Y=1 | X) = \frac{P(Y=1) P(X | Y=1)}{P(Y=1) P(X | Y=1) + P(Y=0) P(X | Y=0)}$$

Divide both numerator and denominator by numerator.

$$p(Y=1|X) = \frac{1}{1 + \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}}$$

or equivalently,

$$p(Y=1|X) = \frac{1}{1 + \exp\left(\ln\left(\frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}\right)\right)}$$

As taken conditional independence assumption, we can say

$$p(Y=1|X) = \frac{1}{1 + \exp\left(\ln\frac{P(Y=0)}{P(Y=1)} + \sum_i \ln\frac{P(X_i|Y=0)}{P(X_i|Y=1)}\right)}$$

$$= \frac{1}{1 + \exp\left(\ln\frac{1-\pi}{\pi} + \sum_i \ln\frac{P(X_i|Y=0)}{P(X_i|Y=1)}\right)}$$

$\hookrightarrow \; ①$

Since it is given that $P(X_i | Y = \alpha_K)$
Is Gaussian, hence

$$\sum_i \ln \frac{P(X_i | Y = 0)}{P(X_i | Y = 1)} = \sum_i \ln \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(X_i - \mu_{i0})^2}{2\sigma_i^2}\right)$$

$$= \sum_i \ln \exp\left(\frac{(X_i - \mu_{i1})^2 - (X_i - \mu_{i0})^2}{2\sigma_i^2}\right)$$

$$= \sum_i \left(\frac{(X_i - \mu_{i1})^2 - (X_i - \mu_{i0})^2}{2\sigma_i^2}\right)$$

$$= \sum_i \left(\frac{(X_i^2 - 2X_i\mu_{i1} + \mu_{i1}^2) - (X_i^2 - 2X_i\mu_{i0} + \mu_{i0}^2)}{2\sigma_i^2}\right)$$

$$= \sum_i \left(\frac{2X_i(\mu_{i0} - \mu_{i1}) + \mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}\right)$$

$$= \sum_i \left(\frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} X_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}\right)$$

$\hookrightarrow$ (2)

Now, this expression is a linear weighted sum of $X_i$'s. Substituting expression (2) in (1), we get

$$P(Y=1|X) = \frac{1}{1 + \exp\left(\ln\frac{1-\pi}{\pi} + \sum_i \left(\frac{\mu_{i0}-\mu_{i1}}{\sigma_i^2}X_i + \frac{\mu_{i1}^2-\mu_{i0}^2}{2\sigma_i^2}\right)\right)}$$

$$= \quad P(Y=1|X) = \frac{1}{1 + \exp\left(w_0 + \sum_{i=1}^{d} w_i X_i\right)}$$

$w_i \ldots \ldots w_d$ are given by

$$w_i = \frac{\mu_{i0}-\mu_{i1}}{\sigma_i^2}$$

and

$$w_0 = \ln\frac{1-\pi}{\pi} + \sum_i \frac{\mu_{i1}^2-\mu_{i0}^2}{2\sigma_i^2}$$

Now for
$$P(Y=0|X) = 1 - P(Y=1|X)$$

$$P(Y=0|X) = 1 - P(Y=1|X) = \frac{\exp\left(w_0 + \sum_{i=1}^{d} w_i X_i\right)}{1 + \exp\left(w_0 + \sum_{i=1}^{d} w_i X_i\right)}$$

we considered only cases where $Y$ was a boolean variable, but Now If $Y$ can take any of discrete class laybel values i.e $y \in (1 \ldots k)$ that is $Y$ can take values $(Y_1 - - - Y_k)$

Then

$P(Y=\alpha_k|X)$ for $Y=\alpha_1$ $Y=\alpha_2$ . . . .

$Y=\alpha_{k-1}$ is :

$$P(Y=\alpha_k|X) = \frac{\exp\left(w_{k0} + \sum_{i=1}^{d} w_{ki} X_i\right)}{1 + \sum_{j=1}^{k-1} \exp\left(w_{j0} + \sum_{i=1}^{n} w_{ji} X_i\right)}$$

$$P(Y=k|X, \mu, \sigma) = \frac{\exp\left(w_{k0} + \sum_{i=1}^{d} w_{ki} X_i\right)}{1 + \sum_{j=1}^{k-1} \exp\left(w_{j0} + \sum_{i=1}^{n} w_{ji} X_i\right)}$$

Ans

Ans (2)

$$l(\theta; D) = -\log P(y^{(1)}, x^{(1)}, y^{(2)}, x^{(2)}$$
$$\cdots \cdots \cdots y^{(N)}, x^{(N)} | \theta)$$

Assuming data are i.i.d

$$l(\theta; D) = -\log\left[ \prod_{i=1}^{N} P(y^{(i)}, x^{(i)} | \theta) \right)$$

we can take $\Sigma$ as i.i.d data.

$$= -\sum_{i=1}^{N} \log P(y^{(i)}, x^{(i)} | \theta)$$

$$= -\sum_{N} \log P(x^{(i)}, y^{(i)} | \theta) \rightarrow \textcircled{1}$$

Since we can write

$$P(x, y | \theta) = \log P(y | \theta)$$
$$+ \log P(x | y, \theta)$$

putting it in equation ①

$$\ell(\theta; D) = -\sum_N \log P(x^{(i)}, y^{(i)} | \theta)$$

$$\ell(\theta; D) = -\sum_N \left( \log P(y^{(i)} | \theta) + \log P(x^{(i)} | y^{(i)}, \theta) \right)$$

Ans

Q2
C. Using the Multivarate Gaussian Distribution Equation since in the question it menhons shared variance $(\sigma_b^2)$.

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}$$

$$\hookrightarrow \text{①}$$

where

$x$ is $D$-dimensional vector

$\mu$ is $D$ dimensional mean vector

$\Sigma$ is $D \times D$ covariance matrix with determinant $|\Sigma|$

$\Sigma$ (covariance matrix) is a matrix whose $(i,j)$ entry is covariance i.e

$$\Sigma_{ij} = cov(x_i, x_j)$$

$$= E[(x_i - \mu_i)(x_j - \mu_j)]$$

$$= E[(x_i x_j)] - \mu_i \mu_j$$

Diagonal entries are variance of

each elements, hence using equation ① in the question.

$$\theta = [\mu, \Sigma, \alpha]$$

$$Z = \sqrt{(2\pi)^D |\Sigma|} = (2\pi)^{D/2} |\Sigma|^{1/2}$$

$$p(x|y) = \frac{1}{Z} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

$$\log L(\theta) = \log P(x, y|\theta) = \log P(y|\theta) + \log p(x|y, \theta)$$

$$= \sum_{i=1}^{N} \log \alpha_{y^{(i)}} - \log Z - \frac{1}{2}\left(X^{(i)} - \mu_{y^{(i)}}\right)^T \Sigma^{-1}_{y^{(i)}}\left(X^{(i)} - \mu_{y^{(i)}}\right) \quad \longrightarrow ①$$

since we are aiming at maximum likelihood that is $\arg\max_\theta \log L(\theta)$ this can happen when

$$\sum_K \alpha_K = 1$$

Taking derivative w.r.t $\mu$

$$\frac{\partial \log L}{\partial \mu_K} = -\sum_{i=0}^{N} \mathbb{1}(y^{(i)} = K) \Sigma^{-1}(x^{(i)} - \mu_K) = 0$$

$$\therefore \boxed{\mu_K = \frac{\sum_{i=1}^{N} \mathbb{1}(y^{(i)} = K) x^{(i)}}{\sum_{i=1}^{N} \mathbb{1}(y^{(i)} = K)}} \quad \underline{Ans}$$

Now take derivative of equation ① w.r.t $\Sigma^{-1}$

$$\frac{\partial \log L}{\partial \Sigma_K^{-1}} = -\sum_{i=0}^{N} \mathbb{1}(y^{(i)} = K)\left[-\frac{\partial \log Z_K}{\partial \Sigma_K^{-1}} - \frac{1}{2}(x^{(i)} - \mu_K)(x^{(i)} - \mu_K)^T\right]$$

$$= 0$$
$$\hookrightarrow ②$$

$$Z = (2\pi)^{D/2} |\Sigma_K|^{1/2}$$

$$\frac{\partial \log Z_K}{\partial \Sigma_K^{-1}} = \frac{1}{Z_K}\frac{\partial Z_K}{\partial \Sigma_K^{-1}}$$

$$= (2\pi)^{-D/2} \, |\Sigma_K|^{-1/2} \, (2\pi)^{D/2} \, \frac{\partial \left( | \Sigma_K^{-1} | \right)^{-1/2}}{\partial \Sigma_K^{-1}}$$

$$= | \Sigma_K^{-1} |^{1/2} \left( -\tfrac{1}{2} \right) | \Sigma_K^{-1} |^{3/2} | \Sigma_K^{-1} | \Sigma_K^{T} = -\frac{1}{2} \Sigma_K$$

$$\left( \text{using property that} \; \Sigma = \Sigma^{T} \right)$$

②

Now substituting this back in equation ②
we get

$$\frac{\partial \log L}{\partial \Sigma_K^{-1}} = - \sum_{i=0}^{N} \mathbb{1} \left( y^{(i)} = K \right) \left[ \tfrac{1}{2} \Sigma_K - \tfrac{1}{2} (x^{(i)} - \mu_K)(x^{(i)} - \mu_K)^{T} \right] = \bar{0}$$

$$\boxed{\Sigma_K = \frac{\sum_{i=1}^{N} \mathbb{1} \left( y^{(i)} = K \right) (x^{(i)} - \mu_K)(x^{(i)} - \mu_K)^{T}}{\sum_{i=1}^{N} \mathbb{1} \left( y^{(i)} = K \right)}}$$

Ans.

Q(2)
(d)

As suggested in the question using Lagrange multiplier

$$\frac{\partial L(\theta)}{\partial \alpha_K} + \lambda \frac{\partial \sum_K \alpha_K}{\partial \alpha_K} = 0$$

Note: Since we need to maximize we are adding the $\lambda$ term.

$$\lambda = -\sum_{i=1}^{N} \mathbb{1}(y^{(i)} = K) \frac{1}{\alpha_K}$$

$$\alpha_K = \frac{-\sum_{i=1}^{N} \mathbb{1}(y^{(i)} = K)}{\lambda} \longrightarrow ①$$

Since again we are wanting to maximize the likelihood as that is possible when $\sum_K \alpha_K = 1 \Rightarrow \lambda = -N$

Substituting value of $\lambda = -N$ in
equation ①

$$\alpha_K = \frac{-\sum\limits_{i=1}^{N} \mathbb{1}\left(y^{(i)} = K\right)}{-N}$$

$$\boxed{\alpha_K = \frac{1}{N} \sum\limits_{i=1}^{N} \mathbb{1}\left(y^{(i)} = K\right)}$$ Ans.