



CSC411/2515 Fall 2018

Homework 4

Soon Chee Loong
999295793

Last Name: Soon
First Name: Chee Loong
cheeloong.soon@mail.utoronto.ca cdf markus: soon chee

1. AlexNet

a) Convolutional Net Architecture

5 Convolutional Layer

3 Fully Connected Layer

↳ includes 1 output layer

Count # of units, # of weights, # of connections in each layer.

Can ignore

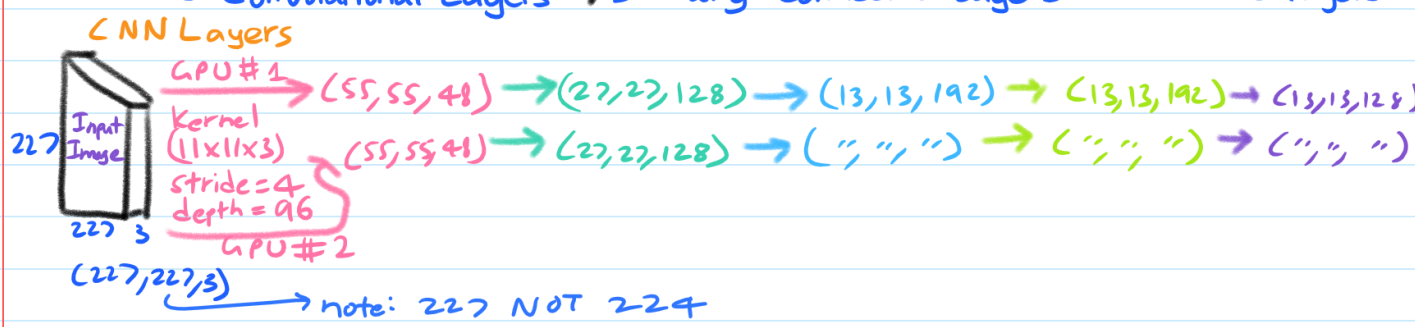
- ↳ units in pooling layers
- ↳ connections between convolution & pooling layers.
- ↳ biases.

When counting the number of connections, adopt convention that when the input to a convolution layer is zero-padded, the connections to the dummy zero values count towards the total

⇒ number of incoming connections is the same for each unit in a given layer.

AlexNet Architecture

5 Convolutional Layers → 3 Fully Connected Layers ⇒ 5 + 3 = 8 layers



CNN Layer 1: $\frac{227 - 11}{4} + 1 = 55$

" " 2: $\frac{55 + 2(1) - 5}{2} + 1 = 27$ (3x3x96)
S=2, d=256, pad=1

" " 3: $\frac{27 - 3}{2} + 1 = 13$ (3x3x256)
S=2, d=384

" " 4: $\frac{13 + 2(1) - 3}{2} + 1 = 13$ (3x3x384)
S=2, d=384, P=1

" " 5: $\frac{13 + 2(1) - 3}{2} + 1 = 13$ (3x3x384)
S=2, d=256, P=1

Fully Connected Layers



FCNN 5: 4096 neurons

" 6: " "

" 7: 1000 "

weights for each convolution operation replicated 96 times

p. =>

Fully Connected Layers

(13, 13, 128) \rightarrow (2048) \rightarrow (2048) \rightarrow 1000 neurons \Rightarrow 1000 way softmax

FCNN S: 4096 neurons

weights for each convolution operation
replicated 96 times

" 6: " "

" 7: 1000 "

Definitions

Input Related

I_h = input height

I_w = input width

I_d = input depth = # of input channels = f_d = filter depth

f_h = filter height

f_w = filter width

N = # of filters = O_d = output depth

O_h = output height

O_w = output width

For AlexNet, size = height = width $\Rightarrow I_h = I_w, f_h = f_w, O_h = O_w$

$$\text{Output Size} = \left\lfloor \frac{\text{Input Size} + 2(\text{Padding Size}) - \text{Filter Size}}{\text{Stride}} \right\rfloor + 1$$

Number of Neurons \Rightarrow number of non-linear activation \Rightarrow number of output units

Fully Connected

$$\# \text{ of output} = O_h \times O_w \times O_d$$

Convolutional

$$O_h \times O_w \times O_d$$

Number of Weights \Rightarrow number of parameters excluding biases
 \Rightarrow space complexity of model

Fully Connected

$$(\# \text{ input}) \times (\# \text{ output})$$

$$= (I_h \times I_w \times I_d) \times (O_h \times O_w \times O_d)$$

Convolutional

$$(f_h \times f_w \times I_d) \times N$$

Number of Connections \Rightarrow number of computations (multiply operations)

\Rightarrow forward pass

Fully Connected $\Rightarrow \forall$ parameter multiplies once

$$(I_h \times I_w \times I_d) \times (O_h \times O_w \times O_d)$$

Fully Connected \Rightarrow \forall parameter multiplies once
 $(I_h \times I_w \times I_d) \times (O_h \times O_w \times O_d)$

Convolutional \Rightarrow Parameters are shared across slides
 $(\# \text{ of parameters}) \times (\# \text{ of slides of convolution})$
 $= (C f_h \times f_w \times I_d \times N) \times (O_h \times O_w)$

Layer	# Units	# Weights	# Connections
Convolution Layer 1	$55 \times 55 \times 96$ $= 290400$	$11 \times 11 \times 3 \times 96$ $= 34848$	$(11 \times 11 \times 3 \times 96) \times (55 \times 55)$ $= 105415200$
Convolution Layer 2	$27 \times 27 \times 256$ $= 186624$	$5 \times 5 \times 48 \times 256$ $= 307200$	$(5 \times 5 \times 48 \times 256) \times (27 \times 27)$ $= 223948800$
Convolution Layer 3	$13 \times 13 \times 384$ $= 64896$	$3 \times 3 \times 256 \times 384$ $= 884736$	$(3 \times 3 \times 256 \times 384) \times (13 \times 13)$ $= 149520384$
Convolution Layer 4	$13 \times 13 \times 384$ $= 64896$	$3 \times 3 \times 192 \times 384$ $= 663552$	$(3 \times 3 \times 192 \times 384) \times (13 \times 13)$ $= 112140288$
Convolution Layer 5	$13 \times 13 \times 256$ $= 43264$	$3 \times 3 \times 192 \times 256$ $= 442368$	$(3 \times 3 \times 192 \times 256) \times (13 \times 13)$ $= 74760192$
Fully Connected Layer 1	4096	9216×4096 $= 37748736$	9216×4096 $= 37748736$
Fully Connected Layer 2	4096	4096×4096 $= 16777216$	4096×4096 $= 16777216$
Output Layer	1000	4096×1000 $= 4096000$	4096×1000 $= 4096000$

b) Suggest a change to the architecture which will help achieve the desired objective.
 \Rightarrow Modify one or more layers.

i) Want to reduce memory usage at test time so network can run on a cellphone.

This requires reducing number of parameters for the network.

Trained parameters need to be stored in memory.

Fully Connected layers account for majority of the # of parameters.

\therefore reduce the size of fully connected layer to reduce # of parameters

\Rightarrow reduce memory usage.

ii) Network needs to make rapid predictions at test time.

Want to reduce the number of connections, since there is approximately one add-multiply operation per connection.

of connections \Rightarrow # of forward pass computation \Rightarrow time complexity for predictions at test time.

Convolutional Layers account for majority of the # of connections.

\therefore reduce number filters for the convolutional layer.



2. Gaussian Naive Bayes

Derive the max. likelihood estimates for Gaussian Naive Bayes, where the features are continuous, and the conditional distribution of each feature given the class is univariate Gaussian rather than Bernoulli. Start with a generative model for a discrete class label $y \in \{1, 2, \dots, k\}$ and a real valued vector of d features $\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$

$$P(y=k) = \alpha_k$$

$$P(\vec{x} | y=k, \vec{\mu}, \vec{\sigma}) = \left(\prod_{i=1}^D \frac{1}{\sqrt{2\pi\sigma_i^2}} \right)^{-\frac{1}{2}} e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ki})^2}$$

each feature is conditionally independent given

α_k = prior on class k the discrete class label.

σ_i^2 = variances for each feature, shared between all classes

μ_{ki} = mean of feature i conditioned on class k

$$\vec{\alpha} = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_k \end{bmatrix}, \quad \vec{\sigma} = \begin{bmatrix} \sigma_1 \\ \sigma_2 \\ \vdots \\ \sigma_d \end{bmatrix}, \quad \vec{\mu} = \begin{bmatrix} \mu_{11} & \mu_{12} & \dots & \mu_{1d} \\ \mu_{21} & \dots & \dots & \mu_{2d} \\ \vdots & & & \vdots \\ \mu_{k1} & \dots & \dots & \mu_{kd} \end{bmatrix}$$

$D = d$ in this question

a) Use Bayes rule to derive an expression for $P(y=k | \vec{x}, \vec{\mu}, \vec{\sigma})$.

Hint: Use law of total probability to derive an expression for $P(\vec{x} | \vec{\mu}, \vec{\sigma})$

$$P(\vec{x}, y=k | \vec{\mu}, \vec{\sigma}) = P(y=k | \vec{x}, \vec{\mu}, \vec{\sigma}) P(\vec{x} | \vec{\mu}, \vec{\sigma})$$

$$= P(\vec{x} | y=k, \vec{\mu}, \vec{\sigma}) P(y=k)$$

$$\Rightarrow P(y=k | \vec{x}, \vec{\mu}, \vec{\sigma}) = \frac{P(\vec{x} | y=k, \vec{\mu}, \vec{\sigma}) P(y=k)}{P(\vec{x} | \vec{\mu}, \vec{\sigma})}$$

$$P(\vec{x} | \vec{\mu}, \vec{\sigma}) = \sum_{j=1}^k P(\vec{x}, y=j | \vec{\mu}, \vec{\sigma}) = \sum_{j=1}^k P(\vec{x} | y=j, \vec{\mu}, \vec{\sigma}) P(y=j), \text{ marginalize out } y$$

$$\Rightarrow P(y=k | \vec{x}, \vec{\mu}, \vec{\sigma}) = \frac{P(\vec{x} | y=k, \vec{\mu}, \vec{\sigma}) P(y=k)}{\sum_{j=1}^k P(\vec{x} | y=j, \vec{\mu}, \vec{\sigma}) P(y=j)}$$

$$= \frac{\left(\prod_{i=1}^D \frac{1}{\sqrt{2\pi\sigma_i^2}} \right)^{-\frac{1}{2}} e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ki})^2} (\alpha_k)}{\sum_{j=1}^k \left(\left(\prod_{i=1}^D \frac{1}{\sqrt{2\pi\sigma_i^2}} \right)^{-\frac{1}{2}} e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ji})^2} (\alpha_j) \right)}$$

$$= \frac{\left(\prod_{i=1}^D \frac{1}{\sqrt{2\pi\sigma_i^2}} \right)^{-\frac{1}{2}} e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ki})^2} (\alpha_k)}{\left(\prod_{i=1}^D \frac{1}{\sqrt{2\pi\sigma_i^2}} \right)^{-\frac{1}{2}} \sum_{j=1}^k \left(e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ji})^2} (\alpha_j) \right)}$$

independent of j

$$= \frac{\left(\prod_{i=1}^D \frac{1}{\sqrt{2\pi\sigma_i^2}} \right)^{-\frac{1}{2}} e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ki})^2} (\alpha_k)}{\sum_{j=1}^k e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ji})^2} (\alpha_j)}$$

a) Use Bayes rule to derive an expression for $p(y=k|\vec{x}, \vec{\mu}, \vec{\sigma})$.

Hint: Use law of total probability to derive an expression for $p(\vec{x}|\vec{\mu}, \vec{\sigma})$

$$P(\vec{x}, y=k|\vec{\mu}, \vec{\sigma}) = P(y=k|\vec{x}, \vec{\mu}, \vec{\sigma}) P(\vec{x}|\vec{\mu}, \vec{\sigma})$$

$$= P(\vec{x}|y=k, \vec{\mu}, \vec{\sigma}) P(y=k)$$

$$\Rightarrow P(y=k|\vec{x}, \vec{\mu}, \vec{\sigma}) = \frac{P(\vec{x}|y=k, \vec{\mu}, \vec{\sigma}) P(y=k)}{P(\vec{x}|\vec{\mu}, \vec{\sigma})}$$

$$P(\vec{x}|\vec{\mu}, \vec{\sigma}) = \sum_{j=1}^K P(\vec{x}, y=j|\vec{\mu}, \vec{\sigma}) = \sum_{j=1}^K P(\vec{x}|y=j, \vec{\mu}, \vec{\sigma}) P(y=j), \text{ marginalize out } y$$

$$\Rightarrow P(y=k|\vec{x}, \vec{\mu}, \vec{\sigma}) = \frac{P(\vec{x}|y=k, \vec{\mu}, \vec{\sigma}) P(y=k)}{\sum_{j=1}^K P(\vec{x}|y=j, \vec{\mu}, \vec{\sigma}) P(y=j)}$$

$$= \frac{\left(\prod_{i=1}^D 2\pi\sigma_i^2\right)^{-\frac{1}{2}} e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ki})^2} (\alpha_k)}{\sum_{j=1}^K \left(\left(\prod_{i=1}^D 2\pi\sigma_i^2\right)^{-\frac{1}{2}} e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ji})^2} (\alpha_j)\right)}$$

$$\sum_{j=1}^K \left(\left(\prod_{i=1}^D 2\pi\sigma_i^2\right)^{-\frac{1}{2}} e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ji})^2} (\alpha_j)\right)$$

$$= \frac{\left(\prod_{i=1}^D 2\pi\sigma_i^2\right)^{-\frac{1}{2}} e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ki})^2} (\alpha_k)}{\left(\prod_{i=1}^D 2\pi\sigma_i^2\right)^{-\frac{1}{2}} \sum_{j=1}^K \left(e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ji})^2} (\alpha_j)\right)}$$

$$\left(\prod_{i=1}^D 2\pi\sigma_i^2\right)^{-\frac{1}{2}} \sum_{j=1}^K \left(e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ji})^2} (\alpha_j)\right)$$

independent of j

$$= \frac{e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ki})^2} (\alpha_k)}{\sum_{j=1}^K \left(e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ji})^2} (\alpha_j)\right)}$$

$$\sum_{j=1}^K \left(e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ji})^2} (\alpha_j)\right)$$

\therefore

$$P(y=k|\vec{x}, \vec{\mu}, \vec{\sigma}) = \frac{(\alpha_k) e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ki})^2}}{\sum_{j=1}^K \left((\alpha_j) e^{-\sum_{i=1}^D \frac{1}{2\sigma_i^2} (x_i - \mu_{ji})^2}\right)}$$

a specific label $= k \neq K = \# \text{ of labels}$

b) Write down an expression for the negative likelihood function

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

Dataset: $D = \{(y^{(1)}, \vec{x}^{(1)}), (y^{(2)}, \vec{x}^{(2)}), \dots, (y^{(N)}, \vec{x}^{(N)})\}$, i.i.d.

Parameters $= \vec{\theta} = \{\vec{\mu}, \vec{\sigma}\}$

$$-\log(P(y^{(1)}, \vec{x}^{(1)}, y^{(2)}, \vec{x}^{(2)}, \dots, y^{(N)}, \vec{x}^{(N)} | \vec{\theta}))$$

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

b) Write down an expression for the negative likelihood function

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

Dataset: $D = \{(y^{(1)}, \tilde{x}^{(1)}), (y^{(2)}, \tilde{x}^{(2)}), \dots, (y^{(N)}, \tilde{x}^{(N)})\}$, i.i.d.

Parameters: $\vec{\theta} = \{\vec{\alpha}, \vec{\mu}, \vec{\sigma}^2\}$

$$-\log(P(y^{(1)}, \tilde{x}^{(1)}, y^{(2)}, \tilde{x}^{(2)}, \dots, y^{(N)}, \tilde{x}^{(N)} | \vec{\theta}))$$

$$= -\log\left(\prod_{i=1}^N P(y^{(i)}, \tilde{x}^{(i)} | \vec{\theta})\right) \quad \text{, since data are i.i.d.}$$

$$= -\sum_{i=1}^N \log(P(y^{(i)}, \tilde{x}^{(i)} | \vec{\theta})) \quad \text{, since } \log(\prod \dots) = \sum (\log(\dots))$$

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

$$= -\sum_{i=1}^N (\log(P(\tilde{x}^{(i)} | y^{(i)} | \vec{\theta})) + \log(P(y^{(i)} | \vec{\theta}))) \quad P(y^{(i)} | \vec{\theta}) = P(y^{(i)}) = \alpha_{y^{(i)}}$$

$$= -\sum_{i=1}^N \left(\log\left(\left(\prod_{j=1}^D 2\pi\sigma_j^2\right)^{-\frac{1}{2}} e^{-\frac{1}{2\sigma_j^2} (x_j^{(i)} - u_{y^{(i)}j})^2}\right) + \log(\alpha_{y^{(i)}}) \right)$$

$$= -\sum_{i=1}^N \left(-\frac{1}{2} \sum_{j=1}^D \log(2\pi\sigma_j^2) + \left(-\frac{1}{2\sigma_j^2} (x_j^{(i)} - u_{y^{(i)}j})^2\right) + \log(\alpha_{y^{(i)}}) \right)$$

\therefore

Negative Log Likelihood

$$l(\theta, D) = -\sum_{i=1}^N \left(-\frac{1}{2} \sum_{j=1}^D \log(2\pi\sigma_j^2) - \sum_{j=1}^D \frac{1}{2\sigma_j^2} (x_j^{(i)} - u_{y^{(i)}j})^2 + \log(\alpha_{y^{(i)}}) \right)$$

c) Take partial derivatives of the likelihood w.r.t. each of the parameters u_{ki} and w.r.t. the shared variances σ_i^2 . Based on this, find the max. likelihood estimates for \vec{u} and $\vec{\sigma}^2$.

Assume each class appears at least once in the dataset.

Partial derivatives w.r.t. u_{ki}

$$\frac{\partial l(\theta, D)}{\partial u_{ki}}$$

$$= \frac{\partial}{\partial u_{ki}} \left(-\sum_{m=1}^N \left(-\frac{1}{2} \sum_{j=1}^D \log(2\pi\sigma_j^2) - \sum_{j=1}^D \frac{1}{2\sigma_j^2} (x_j^{(m)} - u_{y^{(m)}j})^2 + \log(\alpha_{y^{(m)}}) \right) \right)$$

\downarrow since constant w.r.t. u_{ki}
 \downarrow since const. w.r.t. u_{ki}

$$= \frac{\partial}{\partial u_{ki}} \left(\sum_{m=1}^N \sum_{j=1}^D \frac{1}{2\sigma_j^2} (x_j^{(m)} - u_{y^{(m)}j})^2 \right)$$

$$= \sum_{m=1}^N \mathbb{1}[y^{(m)} = k] \frac{2}{2\sigma_i^2} (x_i^{(m)} - u_{ki}) (-1)$$

$$\therefore \frac{\partial l(\theta, D)}{\partial u_{ki}} = \sum_{m=1}^N \mathbb{1}[y^{(m)} = k] \frac{-1}{\sigma_i^2} (x_i^{(m)} - u_{ki})$$

Maximum Likelihood Estimate for \vec{u}

\Rightarrow Set $\frac{\partial l(\theta, D)}{\partial u_{ki}} = 0$ and solve for u_{ki}

$$\frac{\partial l(\theta, D)}{\partial u_{ki}}$$

∂u_{ki} Maximum Likelihood Estimate for \vec{u}

\Rightarrow Set $\frac{\partial \ell(\theta, D)}{\partial u_{ki}} = 0$ and solve for u_{ki}

$$\Rightarrow \sum_{m=1}^N \mathbb{1}[y^{(m)} = k] \frac{-1}{\sigma_i^2} (x_i^{(m)} - u_{ki}) = 0$$

$$\Rightarrow u_{ki \text{ M.L.E.}} = \frac{\sum_{m=1}^N \mathbb{1}[y^{(m)} = k] x_i^{(m)}}{\sum_{m=1}^N \mathbb{1}[y^{(m)} = k]}$$

$$\Rightarrow \vec{u} = \begin{bmatrix} u_{11} & u_{12} & \dots & u_{1d} \\ u_{21} & \dots & \dots & u_{2d} \\ \vdots & & & \vdots \\ u_{k1} & \dots & \dots & u_{kd} \end{bmatrix}_{(k \times d)}, \quad \vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}_{(d \times 1)}, \quad \text{diag}(\vec{x}) = \begin{bmatrix} x_1 & & & \\ & x_2 & & \\ & & \ddots & \\ & & & x_d \end{bmatrix}_{(d \times d)}$$

$$= \begin{bmatrix} \vec{u}_1^T \\ \vec{u}_2^T \\ \vdots \\ \vec{u}_k^T \end{bmatrix}_{(k \times d)}, \quad \vec{u}_k \in (d \times 1)$$

$$\therefore \vec{u}_{\text{M.L.E.}} = \begin{bmatrix} \vec{u}_{1 \text{ M.L.E.}}^T \\ \vec{u}_{2 \text{ M.L.E.}}^T \\ \vdots \\ \vec{u}_{k \text{ M.L.E.}}^T \end{bmatrix}$$

where

$$\vec{u}_{k \text{ M.L.E.}} \in (d \times 1) = \frac{\sum_{m=1}^N \mathbb{1}[y^{(m)} = k] \vec{x}^{(m)}}{\sum_{m=1}^N \mathbb{1}[y^{(m)} = k]}, \quad k \in \{1, 2, \dots, K\}$$

Partial derivatives w.r.t. σ_i^2

$$\frac{\partial \ell(\theta, D)}{\partial \sigma_i^2}$$

$$= \frac{\partial}{\partial \sigma_i^2} \left(-\sum_{m=1}^N \left(-\frac{1}{2} \sum_{j=1}^D \log(2\pi\sigma_j^2) - \frac{1}{2\sigma_j^2} (x_j^{(m)} - u_{y^{(m)}j})^2 + \log(\sigma_{y^{(m)}}) \right) \right)$$

0 since constant w.r.t. σ_i^2

$$= -\sum_{m=1}^N \left(-\frac{1}{2} \sum_{j=1}^D \frac{\partial (\log(2\pi\sigma_j^2))}{\partial \sigma_i^2} - \frac{1}{2\sigma_j^2} (x_j^{(m)} - u_{y^{(m)}j})^2 \frac{\partial (\frac{1}{2\sigma_j^2})}{\partial \sigma_i^2} \right)$$

$$= \sum_{m=1}^N \left(\left(\frac{1}{2} \right) \frac{1}{2\pi\sigma_i^2} 2\pi + (x_i^{(m)} - u_{y^{(m)}i})^2 \left(\frac{-1}{2\sigma_i^4} \right) \right)$$

Partial derivatives w.r.t. σ_i^2

$$\frac{\partial \ell(\theta, D)}{\partial \sigma_i^2}$$

$$= \frac{\partial}{\partial \sigma_i^2} \left(-\sum_{m=1}^N \left(-\frac{1}{2} \sum_{j=1}^D \log(2\pi \sigma_j^2) - \frac{1}{2\sigma_j^2} (x_j^{(m)} - u_{y^{(m)}j})^2 + \log(\sigma_{y^{(m)}j}) \right) \right)$$

0 since constant w.r.t. σ_i^2

$$= -\sum_{m=1}^N \left(-\frac{1}{2} \sum_{j=1}^D \frac{\partial \log(2\pi \sigma_j^2)}{\partial \sigma_i^2} - \frac{1}{2\sigma_j^2} (x_j^{(m)} - u_{y^{(m)}j})^2 \frac{\partial \left(\frac{1}{2\sigma_j^2} \right)}{\partial \sigma_i^2} \right)$$

$$= \sum_{m=1}^N \left(\left(\frac{1}{2} \right) \frac{1}{2\pi \sigma_i^2} 2\pi + (x_i^{(m)} - u_{y^{(m)}i})^2 \left(\frac{-1}{2\sigma_i^4} \right) \right)$$

$$= \sum_{m=1}^N \frac{1}{2\sigma_i^2} - \frac{(x_i^{(m)} - u_{y^{(m)}i})^2}{2\sigma_i^4}$$

$$\therefore \frac{\partial \ell(\theta, D)}{\partial \sigma_i^2} = \sum_{m=1}^N \frac{1}{2\sigma_i^2} - \frac{(x_i^{(m)} - u_{y^{(m)}i})^2}{2\sigma_i^4}$$

Maximum Likelihood Estimate for σ^2

4 try

\Rightarrow Set $\frac{\partial \ell(\theta, D)}{\partial \sigma_i^2} = 0$ and solve for σ_i

$$\Rightarrow \sum_{m=1}^N \frac{1}{2\sigma_i^2} - \frac{(x_i^{(m)} - u_{y^{(m)}i})^2}{2\sigma_i^4} = 0$$

Multiply both sides by $2\sigma_i^4$

$$\Rightarrow \sum_{m=1}^N \sigma_i^2 - (x_i^{(m)} - u_{y^{(m)}i})^2 = 0$$

$$\Rightarrow \sigma_i^2 = \frac{\sum_{m=1}^N (x_i^{(m)} - u_{y^{(m)}i})^2}{N}$$

$$\Rightarrow \sigma_{i, \text{M.L.F.}} = \sqrt{\frac{\sum_{m=1}^N (x_i^{(m)} - u_{y^{(m)}i})^2}{N}}$$

$$\vec{\sigma}_{\text{M.L.F.}} = \begin{bmatrix} \sigma_{1, \text{M.L.F.}} \\ \sigma_{2, \text{M.L.F.}} \\ \vdots \\ \sigma_{D, \text{M.L.F.}} \end{bmatrix}$$

$$\therefore \vec{\sigma}_{\text{M.L.F.}} = \sqrt{\frac{\sum_{m=1}^N (\vec{x}^{(m)} - \vec{u}_{y^{(m)}})^2}{N}}$$

d) Show that the M.L.E. for α_k is given by

$$\alpha_k = \frac{1}{N} \sum_{i=1}^N \mathbb{1}[y^{(i)} = k]$$

Assume each class appears at least once.

Hint: Use Lagrange Multipliers.

Since α_k is the prior probability, $P(y=k)$, will need to satisfy constraint $\sum_{j=1}^K \alpha_j = 1$ to be a valid probability distribution.

Lagrange Multiplier

Use Lagrange Multiplier to turn this constrained optimization problem into an unconstrained optimization problem.

Minimize $f(z)$

s.t. $g(z) = 0$

$$L(z, \lambda) = f(z) + \lambda g(z)$$

$$\Rightarrow f(z) = \ell(\theta, D)$$

$$g(z) = \sum_{j=1}^K \alpha_j - 1 = 0$$

$$\Rightarrow L(z, \lambda) = \ell(\theta, D) + \lambda \left(\sum_{j=1}^K \alpha_j - 1 \right), \text{ s.t. } \lambda = 1$$

$$\frac{\partial L(z, \lambda)}{\partial \alpha_k}$$

$$= \frac{\partial}{\partial \alpha_k} \left(- \sum_{i=1}^N \left(-\frac{1}{2} \sum_{j=1}^D \log(2\pi\sigma_j^2) - \frac{1}{2\sigma_j^2} (x_j^{(i)} - \mu_{y^{(i)}})^2 + \log(\alpha_{y^{(i)}}) \right) + \sum_{i=1}^N \left(\sum_{j=1}^K \alpha_j - 1 \right) \right)$$

\downarrow 0 since constant w.r.t. α_k \downarrow 0 since constant w.r.t. α_k

$$= \sum_{i=1}^N \left(\mathbb{1}[y^{(i)} = k] \left(-\frac{1}{\alpha_k} \right) + 1 \right)$$

$$\frac{\partial L(z, \lambda)}{\partial \alpha_k} = \sum_{i=1}^N \left(\mathbb{1}[y^{(i)} = k] \left(-\frac{1}{\alpha_k} \right) + 1 \right)$$

Set $\frac{\partial L(z, \lambda)}{\partial \alpha_k} = 0$ and solve for α

$$\Rightarrow \sum_{i=1}^N \left(\mathbb{1}[y^{(i)} = k] \left(-\frac{1}{\alpha_k} \right) + 1 \right) = 0$$

$$\Rightarrow \sum_{i=1}^N \mathbb{1}[y^{(i)} = k] \frac{1}{\alpha_k} = \sum_{i=1}^N 1$$

$$\Rightarrow \frac{1}{\alpha_k} = \frac{N}{\sum_{i=1}^N \mathbb{1}[y^{(i)} = k]}$$

$$\Rightarrow \alpha_k = \frac{1}{N} \sum_{i=1}^N \mathbb{1}[y^{(i)} = k], \therefore \text{proven}$$