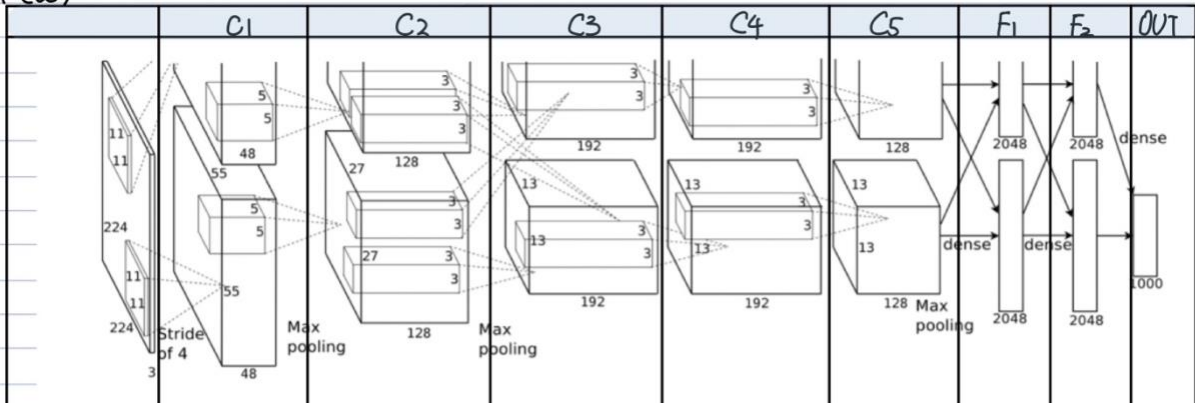


## CSC411 Assignment 4

1. (a)

	C1	C2	C3	C4	C5	F1	F2	OUT
								
# Units	$(55^2 \times 48) \times 2$ = 290,400	$(27^2 \times 128) \times 2$ = 186,624	$(13^2 \times 192) \times 2$ = 64,896	$(13^2 \times 192) \times 2$ = 64,896	$(13^2 \times 128) \times 2$ = 43,264	$2048 \times 2$ = 4096	$2048 \times 2$ = 4096	1000
# Weights	$(11^2 \times 3) \times (48 \times 2)$ = 34,848	$(5^2 \times 48) \times (128 \times 2)$ = 307,200	$(3^2 \times 128 \times 2) \times (192 \times 2)$ = 884,736	$(3^2 \times 192 \times 2) \times (192 \times 2)$ = 663,552	$(3^2 \times 192 \times 2) \times (128 \times 2)$ = 442,368	$(13 \times 128 \times 2) \times (2048 \times 2)$ = 177,209,344	$(2048 \times 2) \times (2048 \times 2)$ = 16,777,216	$(2048 \times 2) \times 1000$ = 4,096,000
# Connections	$(11^3 \times 3) \times (55^2 \times 48 \times 2)$ = 105,415,200	$(5^3 \times 48) \times (27^2 \times 128 \times 2)$ = 223,948,800	$(3^3 \times 128 \times 2) \times (13^2 \times 192 \times 2)$ = 149,520,384	$(3^3 \times 192 \times 2) \times (13^2 \times 192 \times 2)$ = 112,140,288	$(3^3 \times 192 \times 2) \times (13^2 \times 128 \times 2)$ = 74,760,192	$(13^3 \times 128 \times 2) \times (2048 \times 2)$ = 177,209,344	$(2048^3 \times 2) \times (2048 \times 2)$ = 16,777,216	$(2048^3 \times 2) \times 1000$ = 4,096,000

	# Units	# Weights	# Connections
Convolution Layer 1	290,400	34,848	105,415,200
Convolution Layer 2	186,624	307,200	223,948,800
Convolution Layer 3	64,896	884,736	149,520,384
Convolution Layer 4	64,896	663,552	112,140,288
Convolution Layer 5	43,264	442,368	74,760,192
Fully Connected Layer 1	4096	177,209,344	177,209,344
Fully Connected Layer 2	4096	16,777,216	16,777,216
Output Layer	1000	4,096,000	4,096,000

(b)

- (i) To reduce the number of parameters for the network  $\rightarrow$  reduce the number of weights
- Consider reducing the number of units generated by increasing the pooling size in the convolution layers.
  - Consider increasing the stride to reduce the output units.
- (ii) To reduce the number of connections:
- Reduce the kernel size in the convolution layers.
  - Add "bottle-neck" layer between fully-connected layers.

$$\begin{aligned}
 2.(a) \quad p(y=k | \vec{x}, \vec{\mu}, \vec{\sigma}) &= \frac{P(\vec{x} | y=k, \vec{\mu}, \vec{\sigma}) P(y=k | \vec{\mu}, \vec{\sigma})}{P(\vec{x} | \vec{\mu}, \vec{\sigma})} \quad \# \text{ By Bayes' Rule: } P(A|B) = \frac{P(B|A) P(A)}{P(B)} \\
 &= \frac{P(\vec{x} | y=k, \vec{\mu}, \vec{\sigma}) P(y=k | \vec{\mu}, \vec{\sigma})}{\sum_{j=1}^K P(\vec{x} | y=j, \vec{\mu}, \vec{\sigma}) P(y=j | \vec{\mu}, \vec{\sigma})} \\
 &= \frac{\left(\prod_{d=1}^D \frac{1}{\sqrt{2\pi\sigma_d^2}} \exp\left\{-\sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d - \mu_{kd})^2\right\} \cdot \alpha_k}{\sum_{j=1}^K \left[\left(\prod_{d=1}^D \frac{1}{\sqrt{2\pi\sigma_d^2}} \exp\left\{-\sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d - \mu_{jd})^2\right\} \cdot \alpha_j\right]}\right)} \quad \# \text{ By (1), (2).}
 \end{aligned}$$

$$\begin{aligned}
 (b) \quad l(\vec{\theta}; D) &= -\log p(y^{(1)}, \vec{x}^{(1)}, y^{(2)}, \vec{x}^{(2)}, \dots, y^{(N)}, \vec{x}^{(N)} | \vec{\theta}) \\
 &= -\log [p(y^{(1)}, \vec{x}^{(1)} | \vec{\theta}) \cdot p(y^{(2)}, \vec{x}^{(2)} | \vec{\theta}) \cdot \dots \cdot p(y^{(N)}, \vec{x}^{(N)} | \vec{\theta})] \quad \# \text{ By data are i.i.d.} \\
 &= -[\log p(y^{(1)}, \vec{x}^{(1)} | \vec{\theta}) + \log p(y^{(2)}, \vec{x}^{(2)} | \vec{\theta}) + \dots + \log p(y^{(N)}, \vec{x}^{(N)} | \vec{\theta})] \\
 &= -\sum_{i=1}^N \log p(y^{(i)}, \vec{x}^{(i)} | \vec{\theta}) \\
 &= -\sum_{i=1}^N [\log(p(\vec{x}^{(i)} | y^{(i)}, \theta) p(y^{(i)} | \vec{\theta}))] \quad \# \text{ By Chain Rule.} \\
 &= -\sum_{i=1}^N [\log\left(\left(\prod_{d=1}^D \frac{1}{\sqrt{2\pi\sigma_d^2}} \exp\left\{-\sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{y^{(i)}d})^2\right\} \cdot \alpha_{y^{(i)}}\right)\right)] \\
 &= -\sum_{i=1}^N \left[-\frac{1}{2} \log\left(\prod_{d=1}^D \frac{1}{2\pi\sigma_d^2}\right) - \sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{y^{(i)}d})^2 + \log \alpha_{y^{(i)}}\right] \\
 &= \sum_{i=1}^N \left[\frac{1}{2} \sum_{d=1}^D \log(2\pi\sigma_d^2) + \sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{y^{(i)}d})^2 - \log \alpha_{y^{(i)}}\right]
 \end{aligned}$$

$$\begin{aligned}
 (c) \quad l(\vec{\theta}; D) &= \sum_{i=1}^N \left[\frac{1}{2} \sum_{d=1}^D \log(2\pi\sigma_d^2) + \sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{y^{(i)}d})^2 - \log \alpha_{y^{(i)}}\right] \\
 &= \frac{1}{2} \sum_{i=1}^N \sum_{d=1}^D \log(2\pi\sigma_d^2) + \sum_{i=1}^N \sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{y^{(i)}d})^2 - \sum_{i=1}^N \log \alpha_{y^{(i)}}
 \end{aligned}$$

$$\begin{aligned}
 \bullet \quad \frac{\partial l(\vec{\theta}; D)}{\partial \mu_{y^{(i)}d}} &= 0 + \sum_{i=1}^N \sum_{d=1}^D \frac{1}{\sigma_d^2} (x_d^{(i)} - \mu_{y^{(i)}d}) \cdot (-1) = 0 \\
 &= \sum_{i=1}^N \sum_{d=1}^D \frac{1}{\sigma_d^2} (\mu_{y^{(i)}d} - x_d^{(i)})
 \end{aligned}$$

$$\Rightarrow \frac{\partial l(\vec{\theta}; D)}{\partial \mu_{kd}} = \sum_{i=1}^N \frac{1}{\sigma_d^2} (\mu_{kd} - x_d^{(i)}) \mathbb{I}[y^{(i)} = k].$$

$$\begin{aligned}
 \bullet \quad \frac{\partial l(\vec{\theta}; D)}{\partial \sigma_d^2} &= \frac{1}{2} \sum_{i=1}^N \frac{1}{\sigma_d^2} + \sum_{i=1}^N \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{y^{(i)}d})^2 - 0 \\
 &= \frac{N}{2} \cdot \frac{1}{\sigma_d^2} - \sum_{i=1}^N \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{y^{(i)}d})^2
 \end{aligned}$$

$$\begin{aligned}
 \bullet \quad \text{Set } \frac{\partial l(\vec{\theta}; D)}{\partial \mu_{kd}} = 0 &\Rightarrow \sum_{i=1}^N \frac{1}{\sigma_d^2} (\mu_{kd} - x_d^{(i)}) \mathbb{I}[y^{(i)} = k] = 0 \\
 &\Rightarrow \sum_{i=1}^N \mu_{kd} \mathbb{I}[y^{(i)} = k] = \sum_{i=1}^N x_d^{(i)} \mathbb{I}[y^{(i)} = k] \\
 &\Rightarrow \mu_{kd} = \frac{\sum_{i=1}^N x_d^{(i)} \mathbb{I}[y^{(i)} = k]}{\sum_{i=1}^N \mathbb{I}[y^{(i)} = k]} \quad \# \text{ Let } N_k = \sum_{i=1}^N \mathbb{I}[y^{(i)} = k] \\
 &\Rightarrow \mu_{kd} = \frac{\sum_{i=1}^N x_d^{(i)} \mathbb{I}[y^{(i)} = k]}{N_k} \quad \# \text{ maximum likelihood estimate for } \mu_{kd} \text{ where } 1 \leq d \leq D
 \end{aligned}$$

$$\begin{aligned}
 \bullet \quad \text{Set } \frac{\partial l(\vec{\theta}; D)}{\partial \sigma_d^2} = 0 &\Rightarrow \frac{N}{2} \cdot \frac{1}{\sigma_d^2} - \sum_{i=1}^N \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{y^{(i)}d})^2 = 0 \\
 &\Rightarrow \frac{N}{2} \cdot \frac{1}{\sigma_d^2} = \sum_{i=1}^N \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{y^{(i)}d})^2 \\
 &\Rightarrow \frac{6_d}{2} = \frac{1}{N} \sum_{i=1}^N (x_d^{(i)} - \mu_{y^{(i)}d})^2 \quad \# \text{ maximum likelihood estimate for } \sigma_d^2 \\
 &\Rightarrow \sigma_d^2 = \frac{1}{N} \sum_{i=1}^N (x_d^{(i)} - \mu_{y^{(i)}d})^2 \quad \# \text{ where } 1 \leq d \leq D
 \end{aligned}$$

$$(d) \quad l(\vec{\theta}; D) = \frac{1}{2} \sum_{i=1}^N \sum_{d=1}^D \log(2\pi\sigma_d^2) + \sum_{i=1}^N \sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{y^{(i)}d})^2 - \sum_{i=1}^N \log \alpha_{y^{(i)}}$$

$$\begin{aligned}
 \Rightarrow \frac{\partial l(\vec{\theta}; D)}{\partial \alpha_k} &= 0 + 0 - \sum_{i=1}^N \frac{1}{\alpha_k} \mathbb{I}[y^{(i)} = k] \\
 &= -\frac{1}{\alpha_k} \sum_{i=1}^N \mathbb{I}[y^{(i)} = k] = -\frac{N_k}{\alpha_k} \quad \# \text{ Let } N_k = \sum_{i=1}^N \mathbb{I}[y^{(i)} = k]
 \end{aligned}$$

$$\begin{aligned}
 \bullet \quad \text{Let } g(d_1, d_2, \dots, d_k) &= \sum_{i=1}^k d_i = 1 \Rightarrow \frac{\partial g}{\partial d_i} = 1 \quad (*) \\
 \rightarrow \text{By Lagrange Multiplier: } (f' = \lambda g' \text{ with constraint } g(d_1, \dots, d_k) &= \sum_{i=1}^k d_i = 1) \\
 \left(-\frac{N_1}{d_1}, -\frac{N_2}{d_2}, \dots, -\frac{N_k}{d_k}\right) &= \lambda (1, 1, \dots, 1) \quad \# \text{ By } (*)
 \end{aligned}$$

$$\Rightarrow d_i = -\frac{N_i}{\lambda} \text{ where } 1 \leq i \leq k$$

$$\begin{aligned}
 \rightarrow \text{Use constraint } \sum_{i=1}^k d_i &= 1, \text{ we have:} \\
 -\frac{1}{\lambda} \sum_{i=1}^k N_i &= 1 \Rightarrow -\frac{N}{\lambda} = 1 \Rightarrow \lambda = -N
 \end{aligned}$$

$$\begin{aligned}
 \bullet \quad \text{Therefore, } d_i &= -\frac{N_i}{\lambda} = -\frac{N_i}{-N} = \frac{N_i}{N} \\
 \Rightarrow d_k &= \frac{N_k}{N} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[y^{(i)} = k]
 \end{aligned}$$