

學號：R10922111 系級：資工所碩一 姓名：黃冠瑋

model checkpoint 連結：<https://drive.google.com/file/d/1Tmwdx0TSONjj5HXOaKYnXWYfb8trzWGo/view?usp=sharing>

1. (1%) 請以**block diagram**或是文字的方式說明這次表現最好的**model**使用哪些**layer module**(如 **Conv/Linear** 和各類 **normalization layer**) 及連接方式(如一般**forward** 或是使用 **skip/residual connection**)，並概念性逐項說明選用該 **layer module** 的理由。

conv --> normalize --> leakyReLU --> maxPool

因為圖片有局部的特性，也就是假設要在圖中尋找鳥的嘴巴，不一定要觀察整張圖，只看一小部分也能辨認出那是一個鳥嘴，因此選擇conv達到這個目的。

normalize 則是要標準化不同pattern，已加快訓練

leakyReLU 相較於linear可以訓練得更快，而選擇leaky是因為它可以防止某些神經元一直沒有使用到

maxPool 則是因為大部分的情況下，將一張圖縮小並不會影響到我們辨識圖中物體。

2. (1%) 嘗試使用 **augmentation/early-stopping/ensemble** 三種訓練 **trick** 中的兩種，說明實作細節並比較有無該 **trick** 對結果表現的影響(**validation** 或是 **testing** 擇一即可)。

early-stopping: 記錄下valid資料精準度的最大值，當連續五次都沒有超過它就停止繼續訓練。

我的訓練次數設定為60，有使用early-stopping的模型最後的valid acc為0.6456，而沒有使用的最後則為0.5849

augmentation: 載入訓練資料時加上transforms，我選擇的有水平翻轉、左右旋轉15度
使用augmentation的模型在train data的精準度只能達到0.84多，但在valid data的精準度能達到0.6456。而沒有使用augmentation的模型在train data的精準度能夠高達0.97，但在 valid data最高只能到0.58。

3. (1%) 畫出 **confusion matrix** 分析哪些類別的圖片容易使 **model** 搞混，並簡單說明。

(ref: https://en.wikipedia.org/wiki/Confusion_matrix)

	0	1	2	3	4	5	6
0	0.54	0.13	0.12	0.03	0.124	0.023	0.099
1	0.031	0.826	0.02	0.003	0.014	0.003	0.001

2	0.122	0.043	0.48	0.025	0.166	0.078	0.096
3	0.045	0	0.046	0.832	0.034	0.049	0.072
4	0.136	0	0.157	0.024	0.497	0.011	0.141
5	0.03	0	0.123	0.016	0.02	0.82	0.037
6	0.09	0	0.048	0.07	0.145	0.014	0.554

最常把恐懼、難過、中立這三者搞錯，而厭惡、高興、驚訝搞錯的比例則非常低，應該是因為像是高興很明顯會有嘴巴笑的樣子、驚訝則會有眼睛睜大，而恐懼、難過的特徵較相似。

4. (1%) 請統計訓練資料中不同類別的數量比例，並說明：

對 **testing** 或是 **validation** 來說，不針對特定類別，直接選擇機率最大的類別會是最好的結果嗎？針對上述內容，是否存在更好的方式來提升表現？例如設置不同條件來選擇預測結果/變更訓練資料抽樣的方式，或是直接回答「否」（但需要給出支持你論點的論述）

我覺得當機器訓練出來時，當機率彼此都很接近時在直接選擇機率最大的類別會比較好的結果，舉例來說，當某張圖對於恐懼、難過、中立的機率都很高且相近時，就從這三類中挑圖片總數量最多的那個做為結果。

5. (3%) Refer to math problem

https://hackmd.io/@GfOkB4kgS66YhhM7j6TJew/SJy_akYUK

1. Each image's size from (W, H) to $(W+2p_1, H+2p_2)$

$$\begin{cases} k_1 + (x-1)s_1 \leq W+2p_1 \\ k_2 + (y-1)s_2 \leq H+2p_2 \end{cases} \Rightarrow \begin{aligned} x &\leq \frac{W+2p_1-k_1}{s_1} + 1 \\ y &\leq \frac{H+2p_2-k_2}{s_2} + 1 \end{aligned}$$

So, after convolution layer, it's shape will become $(B, \frac{W+2p_1-k_1}{s_1} + 1, \frac{H+2p_2-k_2}{s_2} + 1)$
(output channel)

2. To do the optimization process of loss, let loss function $l(y_i, \hat{y}_i)$, and $y_i = x_i \hat{\sigma}_i + \beta$

$$\frac{\partial l}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \cdot \frac{\partial y_i}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \cdot \gamma$$

$$\frac{\partial l}{\partial \hat{\sigma}_i} = \sum_{j=1}^m \frac{\partial l}{\partial \hat{x}_j} \cdot \frac{\partial \hat{x}_j}{\partial \hat{\sigma}_i}, \quad \text{and} \quad \frac{\partial \hat{x}_i}{\partial \hat{\sigma}_i} = \frac{\partial (x_i - \mu_B)(\hat{\sigma}_i^2 + \epsilon)^{-\frac{1}{2}}}{\partial \hat{\sigma}_i^2} = -\frac{1}{2} (x_i - \mu_B)(\hat{\sigma}_i^2 + \epsilon)^{-\frac{3}{2}}$$

$$(b) \quad \therefore \frac{\partial l}{\partial \hat{\sigma}_i} = \sum_{j=1}^m \frac{\partial l}{\partial y_j} \cdot \gamma \cdot -\frac{1}{2} (x_j - \mu_B)(\hat{\sigma}_i^2 + \epsilon)^{-\frac{3}{2}}$$

$$\frac{\partial l}{\partial \mu_B} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial \mu_B} + \frac{\partial l}{\partial \hat{\sigma}_i} \cdot \frac{\partial \hat{\sigma}_i}{\partial \mu_B}, \quad \text{and} \quad \frac{\partial \hat{x}_i}{\partial \mu_B} = \frac{-1}{\sqrt{\hat{\sigma}_i^2 + \epsilon}}, \quad \frac{\partial \hat{\sigma}_i^2}{\partial \mu_B} = \frac{1}{m} \sum_{j=1}^m 2(x_j - \mu_B)(\hat{\sigma}_i^2 + \epsilon)^{-\frac{1}{2}}$$

$$\text{so, } \frac{\partial l}{\partial \mu_B} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \cdot \gamma \cdot \frac{-1}{\sqrt{\hat{\sigma}_i^2 + \epsilon}} + \frac{\partial l}{\partial \hat{\sigma}_i} \cdot \frac{1}{m} \sum_{j=1}^m 2(x_j - \mu_B)(\hat{\sigma}_i^2 + \epsilon)^{-\frac{1}{2}}$$

$$\text{so } \frac{\partial l}{\partial x_i} = \frac{\partial l}{\partial \hat{x}_i} \cdot \frac{\partial \hat{x}_i}{\partial x_i} + \frac{\partial l}{\partial \mu_B} \cdot \frac{\partial \mu_B}{\partial x_i} + \frac{\partial l}{\partial \hat{\sigma}_i} \cdot \frac{\partial \hat{\sigma}_i}{\partial x_i}$$

$$\therefore \frac{\partial \hat{x}_i}{\partial x_i} = \frac{1}{\sqrt{\hat{\sigma}_i^2 + \epsilon}}, \quad \frac{\partial \mu_B}{\partial x_i} = \frac{1}{m}, \quad \frac{\partial \hat{\sigma}_i^2}{\partial x_i} = \frac{2(x_i - \mu_B)}{m}$$

$$\text{so, } \frac{\partial l}{\partial x_i} = \frac{1}{\sqrt{\hat{\sigma}_i^2 + \epsilon}} \cdot \frac{\partial l}{\partial \hat{x}_i} + \frac{1}{m} \frac{\partial l}{\partial \mu_B} + \frac{2(x_i - \mu_B)}{m} \frac{\partial l}{\partial \hat{\sigma}_i}$$

$$\frac{\partial l}{\partial y} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \cdot \frac{\partial y_i}{\partial y} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \cdot \hat{x}_i, \quad \frac{\partial l}{\partial \beta} = \sum_{i=1}^m \frac{\partial l}{\partial y_i}$$

3.

$$\frac{\partial L}{\partial z_t} = - \frac{\partial}{\partial z_t} \sum y_i \log \hat{y}_i = - \sum y_i \frac{\partial}{\partial z_t} \log \hat{y}_i$$

$$\text{Since } \frac{\partial \log \hat{y}_i}{\partial z_t} = \frac{\partial (\log e^{z_i} - \log \sum_j e^{z_j})}{\partial z_t} = \frac{\partial z_i}{\partial z_t} - \frac{e^{z_t}}{\sum_j e^{z_j}} = \frac{\partial z_i}{\partial z_t} - \hat{y}_t$$

$$\begin{aligned} \text{So, } \frac{\partial L}{\partial z_t} &= - \sum y_i \left(\frac{\partial z_i}{\partial z_t} - \hat{y}_t \right) = - \sum y_i \frac{\partial z_i}{\partial z_t} + \sum y_i \hat{y}_t \\ &= -y_t + \hat{y}_t \quad \checkmark \quad \because \sum y_i = 1 \end{aligned}$$

4.

$$(a) m^t = \beta_1 m^{t-1} + (1-\beta_1) g^t$$

$$= \beta_1^2 m^{t-2} + \beta_1(1-\beta_1) g^{t-1} + (1-\beta_1) g^t = \dots = \beta_1^t m^0 + \beta_1^{t-1}(1-\beta_1) g^1 + \dots + (1-\beta_1) g^t$$

$$\therefore m^t = (1-\beta_1) \sum_{i=1}^t \beta_1^{t-i} g^i$$

$$v^t = \beta_2 v^{t-1} + (1-\beta_2) (g^t)^2$$

$$= \beta_2 (\beta_2 v^{t-2} + (1-\beta_2) (g^{t-1})^2) + (1-\beta_2) (g^t)^2$$

$$= \beta_2^2 v^{t-2} + \beta_2(1-\beta_2) (g^{t-1})^2 + (1-\beta_2) (g^t)^2$$

$$= \dots$$

$$= (\beta_2)^t v^0 + (\beta_2)^{t-1}(1-\beta_2) (g^1)^2 + \dots + (1-\beta_2) (g^t)^2$$

$$= (1-\beta_2) \sum_{i=1}^t \beta_2^{t-i} (g^i)^2$$

$$= \beta_1 (\beta_1 w^{t-1} + (1-\beta_1) g^t)$$

$$\beta_1^3 w^{t-3} + \beta_1^2 (1-\beta_1) g^t +$$

(b) Adam:

$$w^t = w^{t-1} - \frac{\eta_0 \cdot \frac{1}{J_t}}{\sqrt{\frac{v^t}{1-\beta_2}}} \cdot \frac{w^t}{1-\beta_1}$$

$$= w^{t-1} - \frac{\frac{\eta_0}{J_t}}{\sqrt{\sum (g^i)^2}} g^t$$

by (a)

↓

$$= w^{t-1} - \frac{\frac{\eta_0}{J_t}}{\sqrt{\sum \beta_1^{t-i} (g^i)^2}} \cdot \sum_{i=1}^t \beta_1^{t-i} g^i$$

$$0^0 = 1$$

↓

Adagrad:

$$w^t = w^{t-1} - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$