

Kaggle Competetion

1. (1%) 請附上你在 kaggle 競賽上表現最好的降維以及分群方式，並條列**五種**不同降維維度的設定對應到的表現(public / private accuracy)

註1: *auto-encoder* 和 *PCA* 只要任一維度不一樣即可算是一種組合。

註2: 不限於以上方法，同學也可以使用任何其他 *embedding algorithm* 實現降維。

我使用的autoencoder的架構如下，有三層結構對稱的convolution layer以及convolution transpose layer，夾在中間的是兩層linear的layer，負責產生latent vector。產生的latent vector 會放進pca演算法再降至更低的維度，之後再用K-means 做clustering。

```
Net(
  (encoder): Sequential(
    (0): Conv2d(3, 8, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1))
    (1): ReLU()
    (2): Conv2d(8, 16, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1))
    (3): ReLU()
    (4): Conv2d(16, 32, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1))
    (5): ReLU()
  )
  (fc1): Linear(in_features=512, out_features=32, bias=True)
  (fc2): Linear(in_features=32, out_features=512, bias=True)
  (decoder): Sequential(
    (0): ConvTranspose2d(32, 16, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), output_padding=(1, 1))
    (1): ReLU()
    (2): ConvTranspose2d(16, 8, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), output_padding=(1, 1))
    (3): ReLU()
    (4): ConvTranspose2d(8, 3, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), output_padding=(1, 1))
    (5): ReLU()
  )
)
```


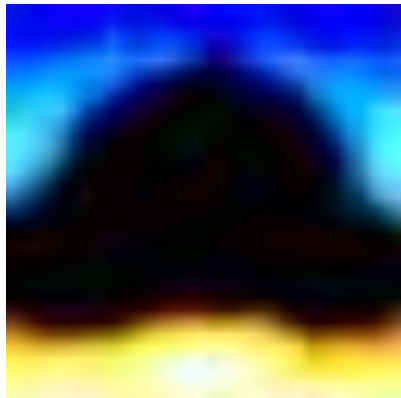
下面我試著用不同的維度產生出降維的資料，我發現怎麼做調整好像都不會比原本的好，維度太高的話可能會造成下面一個分類器分類困難，維度太低的話可能又會損失一些重要的資訊。

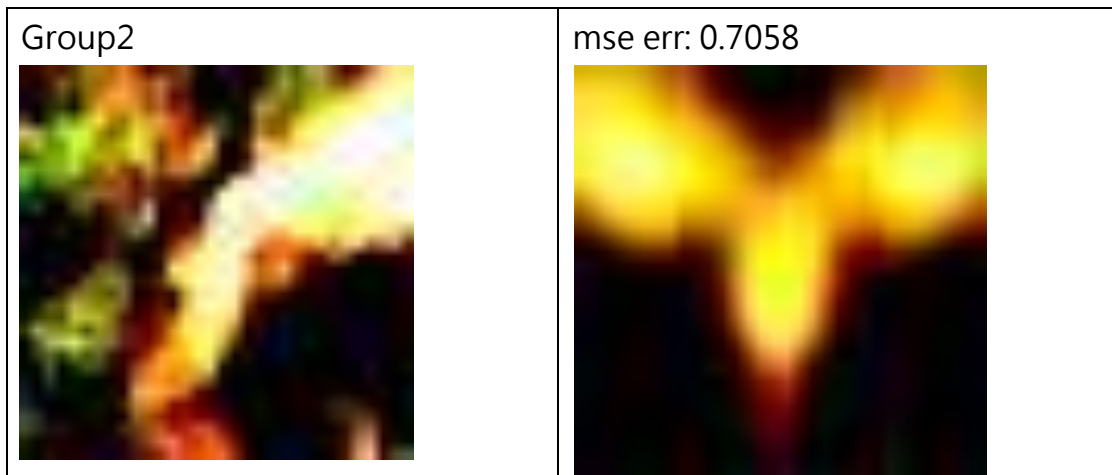
Autoencoder 維度	PCA 維度	Private	Public
32	8	0.81333	0.80666
32	12	0.78889	0.78067

32	4	0.77888	0.76933
64	8	0.78223	0.77712
16	8	0.77022	0.764

2. (1%) 從 trainX.npy 選出不同類別的 2 張圖，貼上原圖以及你的 autoencoder reconstruct 的圖片。用 Mean Square Error 計算這兩張圖的 reconstruction error, 並說明該 error 與 kaggle score 的關係。

下面的第一組圖片主要是風景照的類別，產生的圖片跟原圖比較相似，error 也比較小。第二組圖片是屬於另一種類別，reconstruction 的結果跟原圖不太一樣，error 也比上一個 model 還高。在 kaggle score 的表現上，我的觀察是不一定 rmse 越小 kaggle score 就會越好。rmse 越小代表 latent vector 越能代表一張圖片。然而題目是要做分類器，如果兩種圖片都學到類似的特徵，在做分類的話就會比較困難。舉例來說，我在 model 加入 batchnorm 之後 rmse 有下降，但是 kaggle score 變差了，而且兩種圖片的 rmse 也較為接近。我認為最好的方式應該是第一組的 rmse 很低，然後跟第二組的差很多。

Original	Reconstructed image
Group1 	mse err: 0.5353 



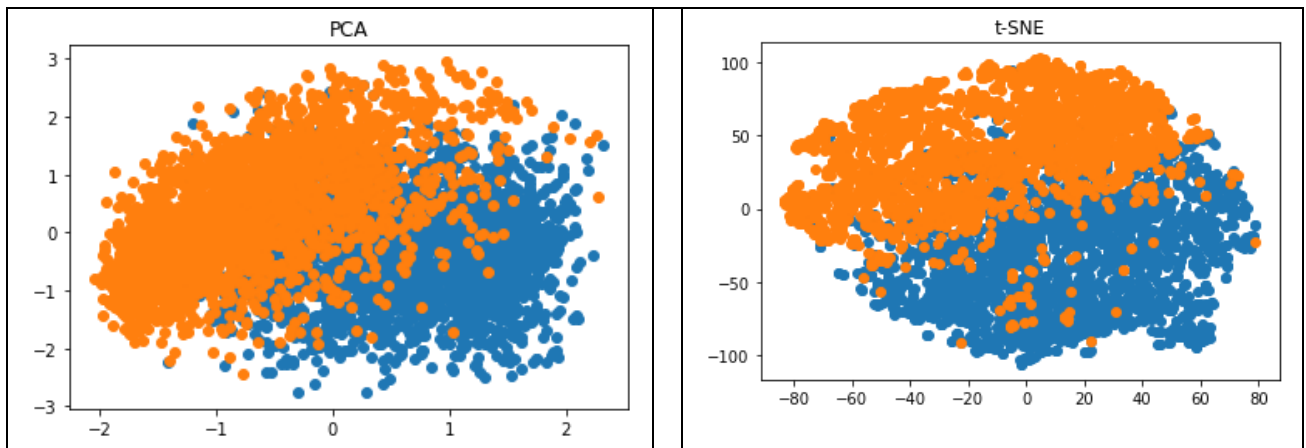
3. (2%) 請使用 `pca` 以及 `tsne` 兩種方法, 將 `visualization.npy` 的圖片經過 `autoencoder` 降維後得到之 `latent vector`, 進一步降維至二維平面並作圖。並說明兩張圖之差異。

註1: `visualization.npy` 前 2500 張 `label` 為 0 ; 後 2500 張 `label` 為 1

註2: 一共要貼上2張圖片。

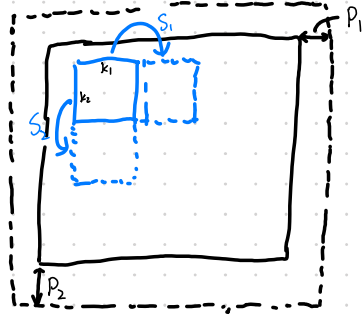
註3: 範例圖片如下 (顏色、分佈不用完全一樣)

下圖為使用不同的方法將`latent vector`做分類所得到的結果，橘色是風景照的類別，可以發現不管用哪種方法都分佈的比較密集，可見他們具有比較共同的特徵，而這種特徵是可以被學習出來的。藍色的類別就比較分散，可見不是風景照的類別比較沒有共同的特徵。至於兩種降維方法的比較，我發現`t-SNE`對於風景照的分佈更為密集，而且看起來也分得比較散。推測可能的原因是`t-SNE`降維在低維度的時候會採用`t`分佈，比起高斯模型更注重長尾分佈，可以有效的將靠近中間的資料在映射後有較大的距離，可以稍微解決在降維後的擁擠問題。然而需要訓練的時間比起`PCA`差非常多，對於使用`colab`的我非常不友善，因此沒辦法在`kaggle`上有很好的實踐。



4. (4%) Refer to math problem :

1. convolution



(B, W, H, input_channels)

1. padding increase the both boundary
2. kernel size decrease on side boundary
3. stride would down sampling the image

$$H_{out} = \left\lceil \frac{H + 2 \times P_1 - (k_1 - 1)}{S_1} + 1 \right\rceil$$

$$W_{out} = \left\lceil \frac{W + 2 \times P_2 - (k_2 - 1)}{S_2} + 1 \right\rceil$$

$$2. \frac{\partial \mathcal{L}}{\partial \hat{x}_i} = \frac{\partial y_i}{\partial \hat{x}_i} \frac{\partial \mathcal{L}}{\partial y_i} = \frac{\partial r \hat{x}_i + \beta}{\partial \hat{x}_i} = r \frac{\partial \mathcal{L}}{\partial y_i}$$

$$\frac{\partial \mathcal{L}}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial \hat{x}_i}{\partial \sigma_B^2} \frac{\partial \mathcal{L}}{\partial \hat{x}_i} = \sum_{i=1}^m \frac{\partial \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}}{\partial \sigma_B^2} r \frac{\partial \mathcal{L}}{\partial y_i} = \sum_{i=1}^m -\frac{r}{2} (x_i - \mu_B) (\sigma_B^2 + \epsilon)^{-\frac{3}{2}} \frac{\partial \mathcal{L}}{\partial y_i}$$

$$\frac{\partial \mathcal{L}}{\partial \mu_B} = \sum_{i=1}^m \left(\frac{\partial \mathcal{L}}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial \mu_B} + \frac{\partial \mathcal{L}}{\partial \sigma_B^2} \frac{\partial \sigma_B^2}{\partial \mu_B} \right) = \sum_{i=1}^m \left(r \frac{\partial \mathcal{L}}{\partial y_i} \cdot \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial \mathcal{L}}{\partial \sigma_B^2} \frac{-2(x_i - \mu_B)}{m} \right)$$

$$\frac{\partial \mathcal{L}}{\partial x_i} = \frac{\partial \mathcal{L}}{\partial \hat{x}_i} \frac{\partial \hat{x}_i}{\partial x_i} + \frac{\partial \mathcal{L}}{\partial \mu_B} \frac{\partial \mu_B}{\partial x_i} + \frac{\partial \mathcal{L}}{\partial \sigma_B^2} \frac{\partial \sigma_B^2}{\partial x_i}$$

$$= \frac{\partial \mathcal{L}}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial \mathcal{L}}{\partial \mu_B} \cdot \frac{1}{m} + \frac{\partial \mathcal{L}}{\partial \sigma_B^2} \cdot \frac{2(x_i - \mu_B)}{m}$$

$$\frac{\partial \mathcal{L}}{\partial \beta} = \sum_{i=1}^m \frac{\partial y_i}{\partial \beta} \frac{\partial \mathcal{L}}{\partial y_i} = \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial y_i}$$

$$\frac{\partial \mathcal{L}}{\partial \gamma} = \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial y_i} \frac{\partial y_i}{\partial \gamma} = \sum_{i=1}^m \hat{\eta}_i \cdot \frac{\partial \mathcal{L}}{\partial y_i}$$

3. softmax

$$1. s_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} \Rightarrow \frac{\partial s_i}{\partial z_j} = \frac{\partial \frac{e^{z_i}}{\sum_{k=1}^N e^{z_k}}}{\partial z_j}$$

$$\text{if } i \neq j \quad e^{z_i} = \text{const} \quad \frac{\partial s_i}{\partial z_j} = \frac{-e^{z_i} e^{z_j}}{(\sum_{k=1}^N e^{z_k})^2} = -s_i s_j$$

$$\text{if } i = j \quad e^{z_i} = e^{z_j} \quad \frac{\partial s_i}{\partial z_j} = s_i - \frac{e^{z_i} e^{z_j}}{(\sum_{k=1}^N e^{z_k})^2} = (1 - s_j) s_i$$

$$\text{import } \delta_{ij} \Rightarrow \frac{\partial s_i}{\partial z_j} = (\delta_{ij} - s_j) s_i \quad \#$$

$$2. L(y, \hat{y}) = - \sum_i y_i \log(\hat{y}_i)$$

$$\frac{\partial L}{\partial z_i} = - \sum y_i \frac{\partial \log(s_i)}{\partial z_i} = - \sum y_i \frac{(1-s_i)s_i}{s_i} = - \sum y_i - s_i$$

$$= - \sum (y_i - \hat{y}_i) \quad \#$$

4. 1. let $M=0$

$$\text{then } \frac{1}{n} \sum_{i=1}^n (x_i)(x_i)^T = \bar{\Sigma} = U \Lambda U^T$$

$\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, and Λ is a symmetric matrix

$$\Rightarrow \Lambda^T = \Lambda$$

$$\bar{\Sigma} = U \Lambda^{\frac{1}{2}} \Lambda^{\frac{1}{2}} U^T = (U \Lambda^{\frac{1}{2}})(\Lambda^{\frac{1}{2}} U^T) \text{ where } \Lambda^{\frac{1}{2}} = \text{diag}(\pm \sqrt{\lambda_1}, \pm \sqrt{\lambda_2}, \dots, \pm \sqrt{\lambda_n}),$$

$$x_i = \sqrt{n} U \Lambda^{\frac{1}{2}} \quad x_i = \sqrt{n} U \Lambda^{\frac{1}{2}} \quad \text{and } \lambda_1 \leq \lambda_2 \leq \lambda_3 \leq \dots \leq \lambda_m$$

$$2. \text{Tr}(\Phi \bar{\Sigma} \Phi) = \frac{1}{n} \|\Phi^T U \Lambda^{\frac{1}{2}}\|_F^2$$

$$\text{let } \Phi = [u_1, u_2, u_3, \dots, u_k] \in \mathbb{R}^{m \times k}$$

$$\text{Tr}(\Phi \bar{\Sigma} \Phi) = \left\| \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_k \end{bmatrix}^{\frac{k \times m}{} \quad m \times m} \Lambda^{\frac{1}{2}} \right\|_F^2$$

$$= \left\| \begin{bmatrix} I_k & 0 \end{bmatrix}^{\frac{k \times k}{} \quad k \times m-k} \begin{bmatrix} \pm \sqrt{\lambda_1} & & \\ & \pm \sqrt{\lambda_2} & \\ & & \ddots \\ & & & \pm \sqrt{\lambda_m} \end{bmatrix} \right\|_F^2$$

$$= \sum_{i=1}^k \lambda_i, \text{ where } \lambda_1 < \lambda_2 < \lambda_3 < \dots < \lambda_m$$

this result give the smallest set of eigen value, therefore, it's a minimum solution