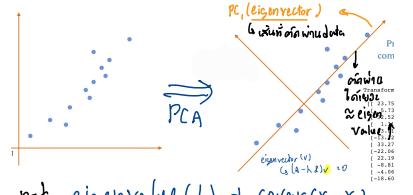


## Reduced Dim

(1) PCA : မျှသုတေသန၏ Vars (ပုံမှန်) (Vars များ၏ ၂၁)



note eigenvalue ( $\lambda$ ) & covar( $x_i, x_j$ )

PCA step

(1) Standardize dataset (ပို့ပို့ပြုပါ)

$$x_{\text{new}} = \frac{x - \mu}{\sigma} \quad (\text{optional})$$

(2) Calculate the Covariance matrix

↳ Covar တဲ့ အကျဉ်းချုပ်များ၏ trend

ရိုးစွဲကြီး (Relationship)

For Sample

$$\text{Cov}(x,y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{(N-1)}$$

① Std Dataset. (optional)

| f1      | f2       | f3       | f4       |
|---------|----------|----------|----------|
| -1      | -0.63246 | 0        | 0.26062  |
| 0.33333 | 1.26491  | 1.73205  | 1.56374  |
| -1      | 0.63246  | -0.57735 | -0.17375 |
| 0.33333 | 0        | -0.57735 | -1.04249 |
| 1.33333 | -1.26491 | -0.57735 | -0.60812 |

② Covar.

| f1 | f2         | f3         | f4         |            |
|----|------------|------------|------------|------------|
| f1 | var(f1)    | cov(f1,f2) | cov(f1,f3) | cov(f1,f4) |
| f2 | cov(f2,f1) | var(f2)    | cov(f2,f3) | cov(f2,f4) |
| f3 | cov(f3,f1) | cov(f3,f2) | var(f3)    | cov(f3,f4) |
| f4 | cov(f4,f1) | cov(f4,f2) | cov(f4,f3) | var(f4)    |

| f1       | f2       | f3      | f4       |
|----------|----------|---------|----------|
| 0.8      | -0.25298 | 0.03849 | -0.14479 |
| -0.25298 | 0.8      | 0.51121 | 0.4945   |
| 0.03849  | 0.51121  | 0.8     | 0.75236  |
| -0.14479 | 0.4945   | 0.75236 | 0.8      |

(3) Eigenvalue ( $\lambda$ ) & eigenvector ( $v$ )

① eigenvalue ( $\lambda$ )

$$\det(A - \lambda I) = 0 ; A = \text{Matrix}, I = \text{Identity} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

② eigenvector ( $v$ ) (pc)

$$(A - \lambda I)v = 0 \quad \text{if } \lambda = 2.516$$

$$\begin{pmatrix} 0.800000 - \lambda & -0.252982 & 0.038490 & -0.144791 \\ -0.252982 & 0.800000 - \lambda & 0.511208 & 0.494498 \\ 0.038490 & 0.511208 & 0.800000 - \lambda & 0.752355 \\ -0.144791 & 0.494498 & 0.752355 & 0.800000 - \lambda \end{pmatrix} \times \begin{pmatrix} v_1 \\ v_2 \\ v_3 \\ v_4 \end{pmatrix} = 0$$

Result

$$\begin{array}{cccc} \lambda_1 = 2.516 & e_1 = & e_2 = & e_3 = \\ 0.161960 & -0.917059 & -0.307071 & 0.196162 \\ 0.524048 & 0.206922 & -0.817319 & 0.120610 \\ 0.585896 & -0.320539 & 0.188250 & -0.720099 \\ 0.596547 & -0.115935 & 0.449733 & 0.654547 \end{array}$$

$$\lambda = 2.5179324, 1.0652885, 0.39388704, 0.02503121$$

(4) ဒေါက်တန် အားလုံး  $\rightarrow$   $\lambda$  &  $V$  ဖော်ဆိုပါ။

(5) မြော်  $V$  အကျဉ်းချုပ်များ၏ n-component

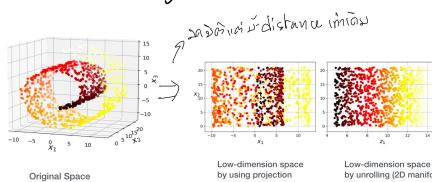
(a) Transform to  $X_{PC}$

$$X_{PC} = \begin{pmatrix} f_1 & f_2 & f_3 & f_4 \end{pmatrix} \begin{pmatrix} e_1 & e_2 & e_3 & e_4 \end{pmatrix}^T = \begin{pmatrix} 0.161960 & 0.917059 & -0.307071 & 0.196162 \\ 0.524048 & 0.206922 & -0.817319 & 0.120610 \\ 0.585896 & -0.320539 & 0.188250 & -0.720099 \\ 0.596547 & -0.115935 & 0.449733 & 0.654547 \end{pmatrix}$$

Proportion of total variance explained by first r PCs:

$$f(r) = \frac{\lambda_1 + \lambda_2 + \dots + \lambda_r}{\lambda_1 + \lambda_2 + \dots + \lambda_n} \quad \text{total}$$

## Manifold learning



↳ မြတ်စွာ ပို့ပို့ပြုတဲ့ distance [သွယ်တဲ့]  
2 type

↳ Linear manifold learning uses direct pairwise distances in original space.

$$\text{number of pair} = \frac{n(n-1)}{2}$$

↳ Non-linear manifold learning uses geodesic distance

or gives more weights to closer neighbors in original space.

## Advance

limitation of KMean → အောက်ပါတဲ့ cluster တွေများ၏ အားလုံး

(1) GMM → ကိုယ် KMean Tavarageများ၏ အားလုံး

↳ cluster ⇒ ဂေါ်ဆုန် dist (bell shape)

(2) Univariate single feature (1 D-Dim)

1: Input:  $n$  data points  $\{x_1, x_2, \dots, x_n\}$  and the number of clusters  $K$

2: Output:  $K$  model parameters  $\{(\mu_k, \sigma_k)\}, 1 \leq k \leq K$

3:

4: Initialization

5: Set  $\hat{\mu}_k$  to be randomly chosen from dataset.

$$\hat{\mu}_k^2 \leftarrow \frac{1}{n} \sum_i (x_i - \bar{x})^2$$

$$\hat{\sigma}_k \leftarrow 1/K$$

$$w_k \leftarrow \frac{1}{K}$$

8: repeat

10: Expectation Step

$$\text{Compute for all data points } x_i, w_{ki} \leftarrow \frac{w_k f(x_i | \theta_k)}{\sum_{j=1}^K w_j f(x_i | \theta_j)}, \forall k$$

$$w_k \leftarrow (1/n) \sum_i w_{ki}, \forall k$$

14: Maximization Step

$$\hat{\mu}_k \leftarrow \frac{1}{n} \sum_i \frac{w_{ki}}{w_k} x_i$$

$$\hat{\sigma}_k^2 \leftarrow \frac{1}{n} \sum_i \frac{w_{ki}}{w_k} (x_i - \hat{\mu}_k)^2$$

17: until  $\hat{\mu}_k$  converges

18: return

$\hat{\mu}_k, \hat{\sigma}_k$

$$L(\Theta, w | X) = \prod_{i=1}^n \sum_{k=1}^K w_k f(x_i; \theta_k)$$

parameter in each k  
↳  $M_k = n$

↳ covar =  $n + \frac{K-1}{2}$   
weight in each k  
 $w_k = K$   
K cluster

num params =  $K(M_k+6)$

↳ Multivariate: (n feature, n Dim)

1: Input:  $n$  data points  $X = \{x_1, x_2, \dots, x_n\}$ ,  $x_i \in \mathbb{R}^d$ ,  $X \in \mathbb{R}^{n \times d}$ , and the number of clusters  $K$

2: Output:  $K$  model parameters  $\{(\mu_k, \Sigma_k)\}, 1 \leq k \leq K$

3:

4: Initialization

5: Set  $\hat{\mu}_k$  to be randomly chosen from dataset.

$$\hat{\Sigma}_k \leftarrow \frac{1}{n} (X - \bar{X})(X - \bar{X})^T$$

$$w_k \leftarrow 1/K$$

$$\theta_k \leftarrow (\hat{\mu}_k, \hat{\Sigma}_k)$$

9:

10: repeat

11: Expectation Step

12: Compute for all data points  $x_i$ ,  $w_{ki} \leftarrow \frac{w_k f(x_i | \theta_k)}{\sum_{j=1}^K w_j f(x_i | \theta_j)}$

$$w_k \leftarrow (1/n) \sum_i w_{ki}$$

14:

15: Maximization Step

$$\hat{\mu}_k \leftarrow \frac{1}{n} \sum_i \frac{w_{ki}}{w_k} x_i$$

$$\hat{\Sigma}_k \leftarrow \frac{1}{n} \sum_i \frac{w_{ki}}{w_k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$

17: until  $\hat{\mu}_k$  converges

18: return

$\hat{\mu}_k, \hat{\Sigma}_k$

↳ Hierarchical Clustering

↳ MIN (Single linkage)

↳ MAX (Complete linkage)

↳ Group Average (Average linkage)

↳ Ward's Method

① find distance Pair & ② find min

P1 P2 P3 P4 P5

| P1      | P2      | P3      | P4      | P5 |
|---------|---------|---------|---------|----|
| 0       | 1.04139 | 0       |         |    |
| 0.59304 | 0.77369 | 0       |         |    |
| 0.46098 | 0.61612 | 0.30232 | 0       |    |
| 0.81841 | 0.32388 | 0.45222 | 0.35847 | 0  |

③ merge & recompute distance

P1 P2 P3,P4 P5

| P1      | P2                      | P3,P4                    | P5                      |
|---------|-------------------------|--------------------------|-------------------------|
| 0       | min C(P3,P1), DCP4, P1) |                          |                         |
| 1.04139 | 0                       | min C(DP3,P2), DCP4, P2) |                         |
| 0.46098 | 0.61612                 | 0                        | min(DCP3,P2), DCP4, P3) |
| 0.81841 | 0.32388                 | 0.35847                  | 0                       |

④ repeat above step

Dendrogram with Single linkage

repeat?

↳ Ward's Method (Centroid linkage)

Let be the SS error of every points from the centroid of cluster .

Determine all combinations of clusters to join such that is smallest.

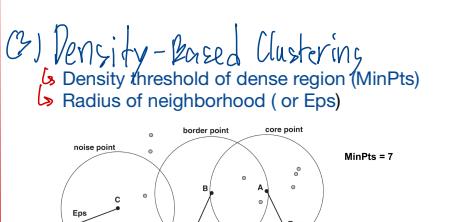
Repeat until only one cluster is left.

For  $K$  clusters, the total SS is  $ESS_K = \sum_{i=1}^K ESS_i$

↳ Density-based Clustering

↳ Density threshold of dense region (MinPts)

↳ Radius of neighborhood (or Eps)



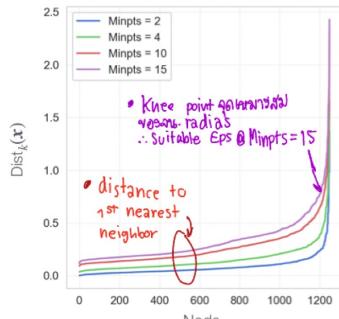
## choice of Minpts & Eps

1. เลือกค่า Minpts  $\rightarrow k = \text{Minpts}$

2. หาจุดที่ distance  $k^*$  ห่างกว่ามากที่สุด

- small values of dist<sub>k</sub>(x) กลุ่ม cluster 8 ห่างกันมาก!

- knee point of dist<sub>k</sub>(x) ห่างกันมากที่สุด Eps ทำให้ Minpts=k



## Factor Analysis.

↳ advance PCA ; ประมาณ 1st PC

Step 1 Std dataset

Step 2 factor loading(L) = eigenvector's PCA  
(n-component)

### Mathematic Behind.

|                   | $X_1$ | $X_2$ | $X_p$ |
|-------------------|-------|-------|-------|
| wheel-base        | 88.6  | 168.8 | 64.1  |
| length            | 88.6  | 168.8 | 48.8  |
| width             | 94.5  | 171.2 | 55.5  |
| height            | 99.8  | 176.6 | 66.2  |
| curb-weight       | 2548  | 2823  | 54.3  |
| engine-size       | 130   | 109   | 2337  |
| compression-ratio | 9.0   | 10.0  | 9.0   |
| city-mpg          | 21    | 24    | 19    |
| highway-mpg       | 27    | 30    | 22    |
| ...               | ...   | ...   | ...   |
| 200               | 109.1 | 188.8 | 68.9  |
| 201               | 109.1 | 188.8 | 68.8  |
| 202               | 109.1 | 188.8 | 68.9  |
| 203               | 109.1 | 188.8 | 68.9  |
| 204               | 109.1 | 188.8 | 68.9  |

$$\hat{\rho}_{12} \approx \frac{1}{n} \mathbf{x}_1 \cdot \mathbf{x}_2 \\ = \frac{1}{208} (88.6 \cdot 168.8 + 88.6 \cdot 168.8 + 94.5 \cdot 171.2 + \dots + 109.1 \cdot 188.8)$$

### (1) Factor Analysis : eigenvector von 1st PC

$$\begin{aligned} X_1 - \mu_1 &= X_1^{(s)} = \ell_{11} F_1 + \ell_{12} F_2 + \dots + \ell_{1m} F_m + \varepsilon_1 \\ X_2 - \mu_2 &= X_2^{(s)} = \ell_{21} F_1 + \ell_{22} F_2 + \dots + \ell_{2m} F_m + \varepsilon_2 \\ &\vdots &&\vdots \\ X_p - \mu_p &= X_p^{(s)} = \ell_{p1} F_1 + \ell_{p2} F_2 + \dots + \ell_{pm} F_m + \varepsilon_p \end{aligned}$$

$$\sum X = \sum L \cdot F + \sum \varepsilon$$

เจตนา relationship ของ feature กับ cluster merge  
(Correlation Matrix)  $\text{Corr}\{X\} \triangleq R = LL^T + \Psi$

$$LL^T + \Psi = \left( \begin{array}{cc} .969 & -.231 \\ .519 & .807 \\ .785 & -.587 \\ .971 & -.210 \\ .704 & .667 \end{array} \right) \left( \begin{array}{ccccc} .969 & .519 & .785 & .971 & .704 \\ .519 & .807 & -.587 & -.210 & .667 \\ -.231 & -.587 & .785 & .971 & .704 \\ .785 & .971 & .971 & .704 & .667 \\ .971 & .704 & .704 & .667 & .704 \end{array} \right) \left( \begin{array}{ccccc} .007 & 0 & 0 & 0 & 0 \\ 0 & .079 & 0 & 0 & 0 \\ 0 & 0 & .040 & 0 & 0 \\ 0 & 0 & 0 & .013 & 0 \\ 0 & 0 & 0 & 0 & .060 \end{array} \right) \left( \begin{array}{ccccc} X_1 = .969 F_1 - .231 F_2 \\ X_2 = .519 F_1 + .807 F_2 \\ X_3 = .784 F_1 - .587 F_2 \\ X_4 = .971 F_1 - .210 F_2 \\ X_5 = .704 F_1 + .667 F_2 \end{array} \right)$$

Specific variances  
Kind Intelligent Happy Likable Just

[1, 5, 5, 1, 1], [8, 9, 7, 9, 8], [9, 8, 9, 9, 8], [9, 9, 9, 9, 9], [1, 9, 1, 1, 9], [9, 7, 7, 9, 9], [9, 7, 9, 9, 7]

dataset = [1.000 .317 .896 .990 .528, .317 1.000 -.066 .335 .904, -.896 -.066 1.000 .885 .161, .990 .335 .885 1.000 .543, .528 .904 .161 .543 1.000]

## Result factor Analysis

| Observed variables | Loadings F1 | Loadings F2 | Communalities | Specific var. |
|--------------------|-------------|-------------|---------------|---------------|
| Kind               | 0.9695      | -0.2311     | 0.9933        | 0.0067        |
| Intelligent        | 0.5194      | 0.8069      | 0.9209        | 0.0791        |
| Happy              | 0.7845      | -0.5872     | 0.9603        | 0.0397        |
| Likable            | 0.9709      | -0.2099     | 0.9867        | 0.0133        |
| Just               | 0.7040      | 0.6669      | 0.9404        | 0.0596        |

$$\text{Variance } F_i = \sum \lambda_i^2 \quad \text{Prop. of Var} = \frac{\sum \lambda_i^2}{p}$$

ผลรวมของค่าถักจะแสดงถึงทั้งหมดมีค่าเท่ากับ  $n$  (จำนวน Observed Variables)  $\Rightarrow P$

### Choosing the number of factor

① เลือก factor = หาค่าถักที่ใหญ่ที่สุด = (eigenvalue) มีค่ามากกว่า 1

② Q scree Plot ของ Eigenvalue

③ ล็อกถักที่ใหญ่ที่สุด = Cumulative var proportion

# ถ้าถักแรกไปจนถึง 50% ถือว่าดี ปัจจัยไม่สามารถอธิบายได้ในรูปแบบเดียว Factor Interpretation = คำอธิบายโครงสร้างของ factor

Factor Rotation มากช่วยให้ความหมายในตัว loading ไม่เข้าใจ

Group 1: King, Happy, Likeable

Group 2: Intelligent, Just

Group 1 : ค. เก่ง น่ารัก / Amiability

G. 2 = Logical ภารกิจ

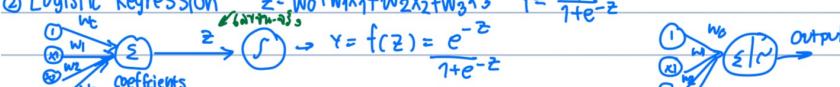
## Lec 11 Neural Networks & Deep Learning

### Neural networks basics

- Activation : ① Linear Regression Model

$$Y = W_0 + W_1 X_1 + W_2 X_2 + W_3 X_3$$

② Logistic Regression  $Z = W_0 + W_1 X_1 + W_2 X_2 + W_3 X_3 \quad Y = \frac{1}{1+e^{-Z}}$



③ Artificial Neuron = a single-layer neural network

$$Z = W_0 + W_1 X_1 + W_2 X_2 + W_3 X_3 \quad Y = f(Z) = h(X) \quad \text{sigmoid logistic } f(z) = \frac{1}{1+e^{-z}}$$

Sigmoid 用于 transform ให้เป็น output [0, 1]

- Loss/Cost: MSE ของ weight ที่ fit train data ให้ได้มากที่สุด  $h = f(\xi; w_i x_i)$

$$\text{in best-fit weight } L(w_1, w_2) = \frac{1}{2N} \sum (w_1 x_1^{(i)} + w_2 x_2^{(i)} - y^{(i)})^2 \quad \text{Error}$$

• Loss Function landscape = Error : ของ lost มากสุด

• Gradient Descent Algorithm ของ weight ที่ fit ให้ได้มากที่สุด gradience มากสุด stop criteria  $\| \nabla L(w_1, w_2) \| < \epsilon$  tolerance  $L$  ค่าที่ต้องการ loss ที่ต้องการ

Sum of squared elements

ใน gradience ของ loss

- ปรับ weight :  $w_{ji} \leftarrow w_{ji} - \alpha \cdot \frac{\partial \text{Loss}}{\partial w_{ji}}$

• ③ ก่อ Overfit ④ L1 L2 Regularization

• Dropout Layer (ลด weight ของ node บางส่วน)

• Feature Importance using Permutation สำหรับใน col  $\rightarrow$  predict  $\rightarrow$  Accuracy?

ผลลัพธ์ของ feature ที่ acc ลดลงมาก = สำคัญ

## Classification Loss/Cost : Cross-Entropy : អាជីវការពិនិត្យរបាយការណ៍

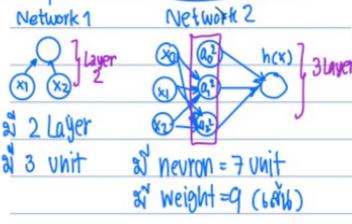
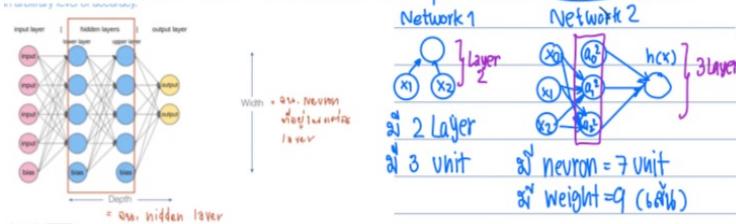
$$\text{Error}_i \text{ or } L_i(\mathbf{w}) = -(y_i \log_2(h(x_i; \mathbf{w})) + (1 - y_i) \log_2(1 - h(x_i; \mathbf{w})))$$

True value Predict prob.

| Inputs | $h$ | $y$ | Squared error | Cross-entropy |
|--------|-----|-----|---------------|---------------|
|        | 0.8 | 1   | 0.04          | 0.10          |
|        | 0.1 | 1   | 0.81          | 1.00          |
|        | 0.1 | 0   | 0.01          | 0.05          |

វាគ្នាំនេះ : សោរផ្តល់ទៅ  $\rightarrow$  សៀវភៅនៃការការពិនិត្យរបាយការណ៍

- Neural Network (NN) perception រាយការណ៍ សម្រាប់ hidden layers



(Q17) Want a single-layer perceptron network with 3 input with no bias unit 1 sigmoid-activation output. With  $w_1 = -1.5, w_2 = 1, w_3 = 1$ . For a dataset  $\{(x_1, x_2, x_3, y)\}$  with a data points  $\{(1, 0, 0, 0.3), (1, 1, 0, 0.5), (1, 0, 1, 0.6)\}$  What's the total squared errors of the network?

$$\text{Soln } z = (x_1 \times w_1) + (x_2 \times w_2) + (x_3 \times w_3) \quad \text{Weight sum data point}$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \text{Sigmoid activation function}$$

$$(g(z) - y)^2 \quad \text{Square error of each data point}$$

Data Point 1  $(1, 0, 0, 0.3)$

$$z = (1 \times 1.5) + (0 \times 1) + (0 \times 1) = -1.5$$

$$g(z) = \frac{1}{1 + e^{(-1.5)}} \approx 0.182$$

$$\text{square error} = (0.182 - 0.3)^2 \approx 0.01386$$

Data Point 2  $(1, 1, 0, 0.5)$

$$z = (1 \times 1.5) + (1 \times 1) + (0 \times 1) = -0.5$$

$$g(z) = \frac{1}{1 + e^{(-0.5)}} \approx 0.372$$

$$\text{square error} = (0.372 - 0.5)^2 \approx 0.0151$$

Data Point 3  $(1, 0, 1, 0.6)$

$$z = (1 \times 1.5) + (0 \times 1) + (1 \times 1) = -0.5$$

$$g(z) = \frac{1}{1 + e^{(-0.5)}} \approx 0.377$$

$$\text{square error} = (0.377 - 0.6)^2 \approx 0.04968$$

$$\therefore \text{Total Squared Error} = 0.01386 + 0.0151 + 0.04968 \approx 0.07864 \approx 0.08 \# Ans$$

(Q17) Want a feed-forward neural network with one hidden layer with 4 nodes + bias unit. At output layer 3 nodes vs input layer 2 nodes

+ bias unit. How many weight parameters are there?

Soln Input Layers : 2 nodes + 1 bias = 3 nodes (Input  $\rightarrow$  Hidden  $\rightarrow$  Output)

Output : 3 nodes

hidden layers : 4 nodes + 1 bias = 5 nodes no bias, vs hidden.

[1] Input Layer to Hidden Layer =  $3 \times 4 = 12$  weights bias

[2] Hidden Layer to Output Layer =  $3 \times 5 = 15$  weights include bias

$\hookrightarrow$  Total Weight Parameter =  $12 + 15 = 27$  total weights # Ans

■ Deep Neural Network(DNN) : Neural networks with 2 នៃជាមុន ការពិនិត្យរបាយការណ៍

◦ Shallow neural network Python  $\rightarrow$  # 1 sklearn, # 2 tensorflow, # 3 keras

■ Neural network Regression with keras  $\rightarrow$  overfit/underfit

Recommendation Tuning

$\rightarrow$  min lost > training lost  
= overfit!

Hyper-parameter Tuning Recommendations

→ 1 - 5 hidden layers (with same # hidden states/units in layer)

→ Cone-like topology : (2 layers)  $2 \times 2 \times 2 \times 2 \times 2 \approx 320 \rightarrow$

→ Batch size  $\leq 32$

→ Make it deeper tends to be more helpful than wider.  $\hookrightarrow$  deep (wider)

→ Learning rate -- Start large and reduce till converge

Fixing Unstable Gradients 1) unstable

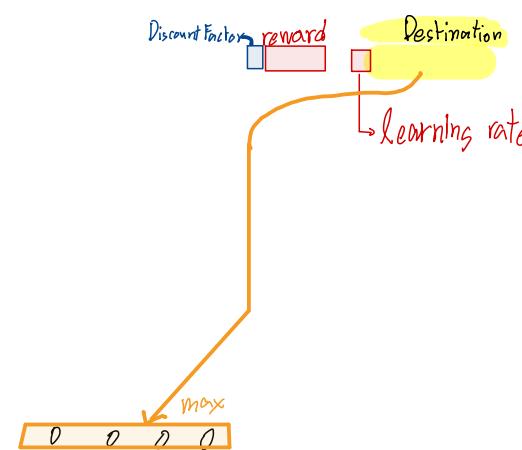
2) instability : transference Learning

① Weight initialization

② Non-saturating Activation

③ Batch Normalization : Scale  $\rightarrow$  bias output normalization layer

避免 unstable gradient



# Quiz 1) PCA

QUIZ 1.1  $\text{var}(x_1, x_2)$  នៃអាជីវកម្ម

$$\begin{bmatrix} x_1 & x_2 \\ 5 & 4 \\ 1 & -2 \\ -1 & 1 \\ 3 & 3 \\ 3 & 0 \end{bmatrix} \quad \frac{\text{var}_2 =}{12}$$

$$\bar{x}_1 = \frac{5+1-1+3}{4} = 2, \bar{x}_2 = \frac{4-2+1+3+0}{5} = 1$$

$$\rightarrow \begin{bmatrix} 3 & 3 \\ 3 & 0 \end{bmatrix} = \begin{bmatrix} 9 & 9 \\ 9 & 0 \end{bmatrix} \quad \text{np.cov}(X, \text{rowvar}=\text{False})$$

$$\therefore \text{var}(x_1, x_2) = \frac{9+3+0+0}{5-1} = \frac{12}{4} = 3$$

2. ដំណោះស្រាយ Total sample variance នៃក្នុងព័ត៌មាន

$$\begin{bmatrix} 6.67 & 4 \\ 4 & 6 \end{bmatrix} \quad \text{Var}(x_1) = (3^2 + (-1)^2 + (-3)^2 + (1)^2) / 4 = 6.67$$

$$\text{Var}(x_2) = \frac{(3^2 + (-1)^2 + (0)^2 + (2)^2)}{4} = 6$$

$$\text{total var} = 6.67 + 6 = 12.67$$

## Clustering

- From a one-dimensional data set  $\{2, 4, 10, 12, 3, 20\}$ , if we apply K-means clustering algorithm using  $K=2$ , and that we randomly pick the initial centroid as  $(\mu_1, \mu_2) = (1, 6)$ . What is the centroid in the next iteration?

| <u>centroid</u> | <u>cluster</u>      |
|-----------------|---------------------|
| $\mu_1 = 1$     | $\{2, 3\}$          |
| $\mu_2 = 6$     | $\{4, 10, 12, 20\}$ |

### update

$$\begin{aligned} \mu_{1\text{new}} &\rightarrow \mu_1 = \frac{2+3}{2} = 2.5 \\ \mu_{2\text{new}} &\rightarrow \mu_2 = \frac{4+10+12+20}{4} = 11.5 \end{aligned}$$

- ត្រូវរាយៈ cluster ទាំងអស់មួលបើជាបិទ  $\{(5, 4), (3, 1), (1, 1)\}$  គោលនយោបាយ SSE ទាំងអស់ cluster មិនចាត់ទៅទេ \*

$$\begin{aligned} \text{centroid } \mu_1 &= \frac{5+3+1}{3} = \frac{9}{3} = 3 \\ \text{centroid } \mu_2 &= \frac{4+1+1}{3} = \frac{6}{3} = 2 \\ \text{SSE} &= \sqrt{\left(\sqrt{(5-3)^2 + (4-2)^2}\right)^2 + \left(\sqrt{(3-3)^2 + (1-2)^2}\right)^2 + \left(\sqrt{(1-3)^2 + (1-2)^2}\right)^2} \\ &= \sqrt{8+1+5} = \sqrt{14} \end{aligned}$$

Consider the likelihood function of a mixture of three bivariate Gaussian distributions below. How many model parameters to estimate in this mixture model?

$$L(\Theta, w | \mathcal{X}) = \prod_{i=1}^n \sum_{k=1}^3 w_k f(x_i; \theta_k)$$

$$f(x; \theta_k) = \frac{\exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))}{\sqrt{(2\pi)^2 |\Sigma|}}, \quad \theta_k = (\mu_k, \Sigma_k)$$

→ Variable of each Gaussian dist.

$$\begin{cases} \mu = \mu_1, \mu_2 \\ \text{cov} = \Sigma_1, \Sigma_2, \text{cov}(X_1, X_2) \end{cases}$$

→  $K = 3$  cluster

→ Weight  $= w_1, w_2, w_3 = 3$  នាក់  
numbers of param  $= 3(5) + 3 = 18$  នាក់

From the current proximity matrix, which of the two clusters/points are to be merged next by using a single linkage?

Current proximity matrix

|   | 1 | 2     | 3     | 4     | 5     | 6     | 7     | 8     |
|---|---|-------|-------|-------|-------|-------|-------|-------|
| 1 | 0 | 1.414 | 4.472 | 5.657 | 6.708 | 5.099 | 7.280 |       |
| 2 |   | 0     | 3.162 | 4.243 | 5.000 | 3.162 | 5.385 |       |
| 3 |   |       | 0     | 4.472 | 5.000 | 3.162 | 5.385 |       |
| 4 |   |       |       | 0     | 4.123 | 4.243 | 5.385 |       |
| 5 |   |       |       |       | 0     | 2.236 | 3.162 | 5.606 |
| 6 |   |       |       |       |       | 0     | 2.236 | 1.414 |
| 7 |   |       |       |       |       |       | 0     | 2.236 |
| 8 |   |       |       |       |       |       |       | 0     |

Pairwise distance matrix

|   | 1 | 2     | 3     | 4     | 5     | 6     | 7     | 8     |
|---|---|-------|-------|-------|-------|-------|-------|-------|
| 1 | 0 | 1.414 | 2.000 | 4.472 | 5.657 | 6.708 | 5.099 | 7.280 |
| 2 |   | 0     | 1.414 | 3.162 | 4.243 | 5.385 | 4.000 | 6.083 |
| 3 |   |       | 0     | 4.000 | 4.472 | 5.000 | 3.162 | 5.385 |
| 4 |   |       |       | 0     | 2.000 | 4.123 | 4.243 | 5.385 |
| 5 |   |       |       |       | 0     | 2.236 | 3.162 | 5.606 |
| 6 |   |       |       |       |       | 0     | 2.236 | 1.414 |
| 7 |   |       |       |       |       |       | 0     | 2.236 |
| 8 |   |       |       |       |       |       |       | 0     |

$$\text{single linkage} = 1 \text{ និង } \{2, 3\} \text{ នឹង } 6, 8$$

From each of the possible merges in the previous question, what is the distance of the closest point to the merged cluster?

$$\min(D(5, 6), D(5, 8)) = \min(2.23, 3.606)$$

$$= 2.23$$

$$\min(D(7, 6), D(7, 8)) = \min(2.23, 2.23)$$

$$= 2.23$$

Distance  $= 2.23$

## Factor Analysis

| Observed variables | Loadings F1 | Loadings F2 | Communalities | Specific var. |
|--------------------|-------------|-------------|---------------|---------------|
| Kind               | 0.9695      | -0.2311     | 0.9933        | 0.0067        |
| Intelligent        | 0.5194      | 0.8069      | 0.9209        | 0.0791        |
| Happy              | 0.7845      | -0.5872     | 0.9603        | 0.0397        |
| Likable            | 0.9709      | -0.2099     | 0.9867        | 0.0133        |
| Just               | 0.7040      | 0.6669      | 0.9404        | 0.0596        |

What is the amount of variance in the data explained by factor F1? (Choose the nearest value)

$$\text{Var at F1} = (0.96)^2 + (0.52)^2 + (0.78)^2 + (0.97)^2 + (0.7)^2$$

$$= 3.261$$

What is the proportion of variance in the data explained by factor F1? (Choose the nearest value)

$$\text{prop} = \frac{\text{Var F1}}{P} = \frac{3.261}{5} = 65\%$$

What is the proportion of variance in the data explained by factor F2? (Choose the nearest value)

$$\text{prop} = \frac{(-0.23)^2 + (0.8)^2 + (-0.58)^2 + (-0.21)^2 + (0.67)^2}{5}$$

$$= 30\%$$