# Representation learning : Learn useful representation

1) Dimensionality Reduction
2) Clustering / Compression
3) Generative learning

Unsupervised Learning: Extract useful information  ★ Data is about variance

---

## Dimensionality Reduction - Big Idea

$x \in \mathbb{R}^d \rightarrow$ $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$, $k < d$ $\rightarrow z = f(x) \in \mathbb{R}^k$

1) Generally speaking, we need $f$ to preserve as much information about $x$ as possible

$$\|z^{(i)} - z^{(j)}\|^2 \approx \|x^{(i)} - x^{(j)}\|^2$$

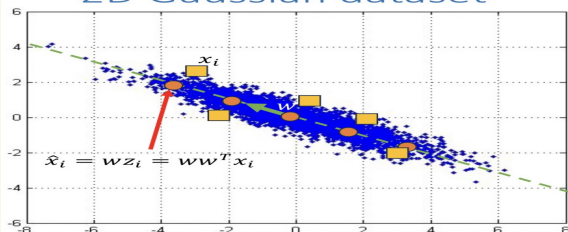$$f(x) = w^T x$$

$$z = w^T x^{(i)} \quad s.t \quad \|w\| = 1$$

① Find $w$ s.t maximizes the variance of the projected data

② Find $w$ s.t minimizes the reconstruction error



### 2D Gaussian dataset

$x_i$

$w$

$\hat{x}_i = w z_i = w w^T x_i$

① = ②

What's $w^T x$? : Reduced data

---

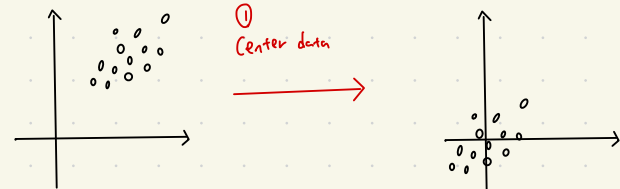## Principal Component Analysis (PCA)

Given data : X

① Center data

$B = X - \bar{X}$

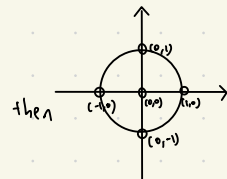Centering data: Subtract mean of each variable from the dataset



---

② Covariance Matrix $\Sigma$ : $\Sigma = B^T B$

$$\Sigma = \begin{pmatrix} Var(X) & Cov(X,Y) \\ Cov(X,Y) & Var(Y) \end{pmatrix}$$

Ex)

Let $\Sigma = \begin{pmatrix} 9 & 4 \\ 4 & 3 \end{pmatrix}$, then


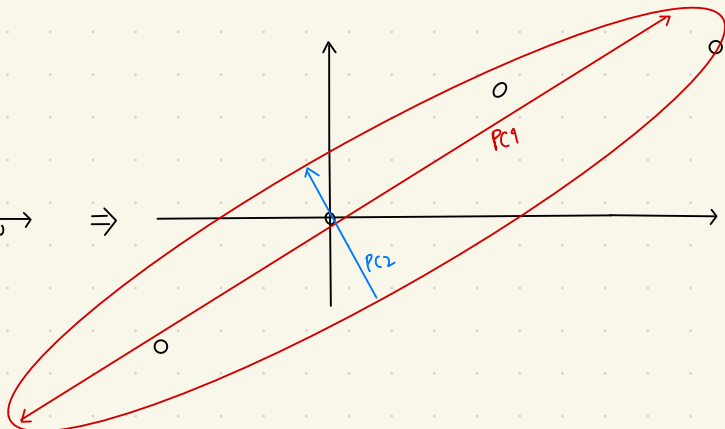
$\Rightarrow$



PC1

PC2

$(X,Y) \rightarrow (9x+4y, 4x+3y)$

$(0,0) \rightarrow (0,0)$
$(1,0) \rightarrow (9, 4)$
$(0,1) \rightarrow (4, 3)$
$(-1,0) \rightarrow (-9, -4)$
$(0,-1) \rightarrow (-4, -3)$

③ Compute eigenvectors & eigenvalues $Sv = \lambda v$

④ Sort eigenvectors w.r.t eigenvalues in descending order (The eigenvalue indicate the variance of the projected data)

⑤ Take top $k$ eigenvectors

How these steps are related to finding $w$ that maximizes variance of projected data?

1) The reduced data is $\{z^{(i)} = w^T x^{(i)}\}_{i=1}^N$

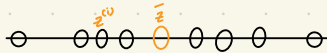2) Calculate the mean of the reduced data

$$\bar{z} = \frac{1}{N}\sum_{i=1}^{N} z^{(i)} = \frac{1}{N}\sum_{i=1}^{N} w^T x^{(i)} = w^T\left(\frac{1}{N}\sum_{i=1}^{N} x^{(i)}\right) = w^T \bar{x}$$

3) Calculate the variance of the reduced data

$$Var = \frac{1}{N}\sum_{i=1}^{N}(z^{(i)} - \bar{z})^2 = \frac{1}{N}\sum_{i=1}^{N}(w^T x^{(i)} - w^T \bar{x})^2 = \frac{1}{N}\sum_{i=1}^{N}(w^T(x^{(i)} - \bar{x}))^2$$

$$= \frac{1}{N}\sum_{i=1}^{N} w^T(x^{(i)} - \bar{x})(x^{(i)} - \bar{x})w$$

$$= w^T\left(\underbrace{\frac{1}{N}\sum_{i=1}^{N}(x^{(i)} - \bar{x})(x^{(i)} - \bar{x})}_{S}\right)w = w^T S w$$

Hence, maximizing the variance can be written as

$$\max_{w} w^T S w \quad \Rightarrow \quad \max_{w} w^T S w$$
$$\text{s.t } \|w\| = 1 \qquad\qquad \text{s.t } w^T w = 1$$

Use Lagrangian function

$$\mathcal{L}(w, a) = w^T S w + a(1 - w^T w) = w^T S w + a - a w^T w$$

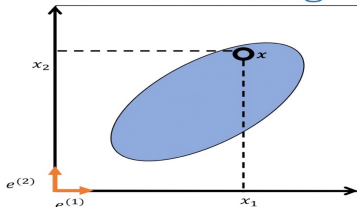$$\frac{\partial \mathcal{L}}{\partial w} = 2Sw - 2aw = 0$$

$$= Sw - aw = 0$$

$$= Sw = aw : \text{At } \frac{\partial \mathcal{L}}{\partial w} = 0, \text{ we derive to equation for getting eigen vectors \& eigen values}$$
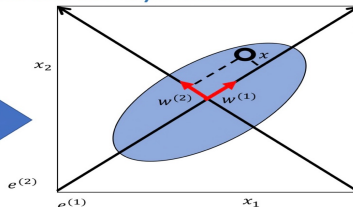
★ $S = \frac{1}{N}\sum_{i=1}^{N}(x^{(i)} - \bar{x})^2$
     Centering data



PCA as a change of coordinate system

From Vanderbilt University Lecture note

$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$  ➡  $x = x_1 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + x_2 \begin{bmatrix} 0 \\ 1 \end{bmatrix}$
                                    $e^{(1)}$     $e^{(2)}$

$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$  ➡  $x = \bar{x} + (x^T w^{(1)})w^{(1)} + (x^T w^{(2)})w^{(2)}$

10/6/23     Foundations of Machine Learning - Week 4     13

Reconstruction : $\bar{X} + ww^T x^{(i)} - ww^T \bar{x}$
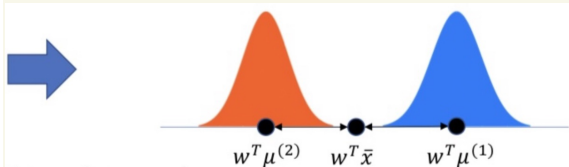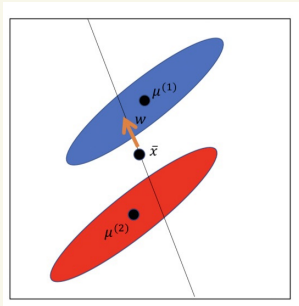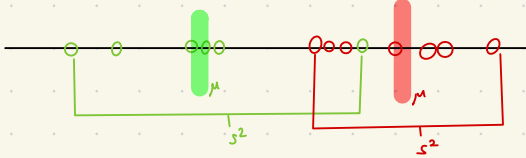
$= \bar{X} + w z^{(i)} - w \bar{z}$

# Linear Discriminant Analysis Derivations (LDA) : Supervised dimensionality reduction

: LDA focuses on maximizing the separability among <u>known categories</u>

How LDA creates a new axis?

1) Maximize the distance between means

2) Minimize the variation (which LDA calls "scatter", and is represented $s^2$) within each categories.





① Distance between class means    * K: # of classes

$$: \frac{1}{K}\sum_K \left(w^T\mu^{(K)} - w^T\bar{x}\right)^2 = \frac{1}{K}\sum_K \left(w^T\left(\mu^{(K)} - \bar{x}\right)\right)^2$$

$$= w^T\left(\frac{1}{K}\sum_K \left(\mu^{(K)} - \bar{x}\right)\left(\mu^{(K)} - \bar{x}\right)^T\right)w = w^T S_b w$$

$$\underbrace{\phantom{\frac{1}{K}\sum_K \left(\mu^{(K)} - \bar{x}\right)\left(\mu^{(K)} - \bar{x}\right)^T}}_{S_b}$$

② Within class variance

$$: \frac{1}{N}\sum_K \sum_{i \in C_K} \left(w^T x^{(i)} - w^T\mu^{(K)}\right)^2$$

$$= w^T\left(\frac{1}{N}\sum_K \sum_{i \in C_K} \left(x^{(i)} - \mu^{(K)}\right)\left(x^{(i)} - \mu^{(K)}\right)^T\right)w = w^T S_w w$$

$$\underbrace{\phantom{\frac{1}{N}\sum_K \sum_{i \in C_K} \left(x^{(i)} - \mu^{(K)}\right)\left(x^{(i)} - \mu^{(K)}\right)^T}}_{S_w}$$

LDA Objective : $\underset{w}{\operatorname{argmax}} \dfrac{w^T S_b w}{w^T S_w w}$   is equivalent to   $\underset{w}{\operatorname{argmax}} \; w^T S_b w$

$\qquad\qquad\qquad\quad$ s.t. $\|w\|=1$ $\qquad\qquad\qquad\qquad\qquad$ s.t. $w^T S_w w = 1$

Lagrangian function:

$$\mathcal{L}(w,\lambda) = w^T S_b w - \lambda(w^T S_w w - 1)$$

$$\nabla_w \mathcal{L} = 2 S_b w - 2\lambda S_w w = 0$$

$\qquad \Rightarrow S_b w - \lambda S_w w = 0$

$\qquad \Rightarrow \left(S_w^{-1} S_b\right)w = \lambda w$    Hence, finding eigenvector & eigenvalues of $S_w^{-1} S_b$

# Auto-Encoder (AE).

## PCA   vs   AE


PCA

$X$ | $W \in \mathbb{R}^{d \times k}$ | $z$ | $W^T$ | $\hat{X}$

$z \in \mathbb{R}^k$

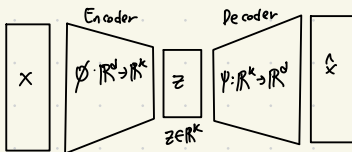PCA: *Linear dimensionality Reduction*

- Forward transform: $z = w^T x$

- Inverse transform: $\hat{X} = Wz$

Objective: Maximizing Variance or Minimizing Reconstruction Error (they're equivalent)

$$\min_{W} \mathbb{E}_x \left[ \| X - \hat{X} \|^2 \right] = \mathbb{E}_x \left[ \| X - WW^T X \|^2 \right]$$
$$\text{s.t. } W^T W = I_{k \times k}$$

## Auto Encoder



Encoder          Decoder

$X$ | $\phi: \mathbb{R}^d \to \mathbb{R}^k$ | $z$ | $\psi: \mathbb{R}^k \to \mathbb{R}^d$ | $\hat{X}$
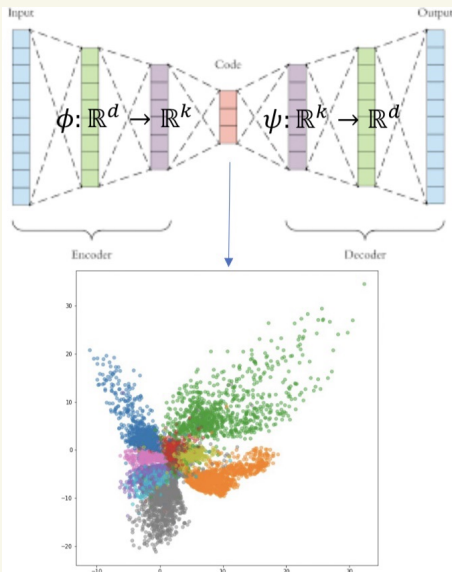
$z \in \mathbb{R}^k$

AE: *Non linear dimensionality Reduction*

- Forward transform: $z = \phi(x)$

- Inverse transform: $\hat{X} = \psi(x)$

Objective: Minimizing Reconstruction Error.

$$\min_{\phi, \psi} \mathbb{E}_x \left[ \| X - \hat{X} \|^2 \right] = \mathbb{E}_x \left[ \| X - \psi(\phi(x)) \|^2 \right]$$



From Vanderbilt University
   Lecture note

- Reconstruction Loss: $\mathbb{E}_x [\| x - \psi(\phi(x)) \|^2]$

- $\hat{y}$ = Likelihood = Softmax $(w^T \phi(x))$

- Negative Log-likelihood = $-\sum_i y^{(i)} \log(y^{(i)})$, which is the Cross entropy.