# Dimension reduction and autoencoders

Hung-Hsuan Chen

1/11/21

1

1

# Dimension reduction

- Let $X \in R^{n \times d}$
- We would like to find a new representation $Z \in R^{n \times k}$, where $k < d$
- We want $Z$ can still well represent the original $X$

1/11/21

2

2

# Why dimension reduction?

- Compress data and preserve useful information
- Data visualization

1/11/21

3

3

# Toy example 1

- Consider the following 3d points
  - (1,2,3), (2,4,6), (3,6,9), (4,8,12), (5,10,15), (6,12,18)
  - If each integer requires 1 byte, we need 1*3*6=18 bytes
- However, we may also store the first point (1,2,3) as the base, and store the multiplier of each point
  - One point (3 bytes) + multipliers (6 bytes)
- Reduced 50% of the storage

1/11/21

4

4

## Toy example (con't)

- Consider the following 3d points
  - (1,2,3), (2,4,6), (3,6,8), (4,8,12), (5,10,15), (6,12,19)
  - If each integer requires 1 byte, we need 1*3*6=18 bytes
- If we store the first point (1,2,3) as the base, and store the multiplier of each point
  - One point (3 bytes) + multipliers (6 bytes)
- Reduced 50% of the storage
- However, we have some small loss

1/11/21                                                    5
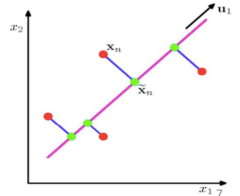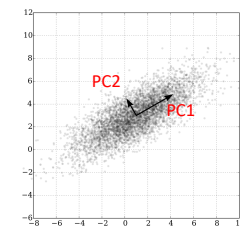
5

# Principal component analysis (PCA)

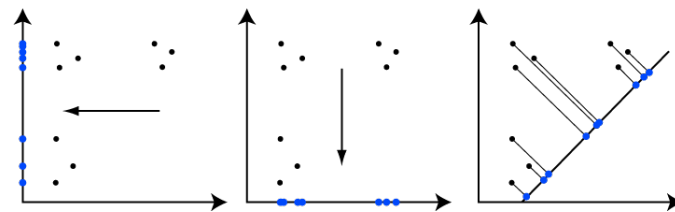1/11/21                                                    6

6

## PCA

- Goal: find base direction(s) that preserves "important" aspects of data
- PCA
  - Allow only linear transforms from the original data point to the new data point
  - Define the goodness by
    - Maximizing the variance of projected data (purple line)
    - Minimize mean squared distance between data points and projections (blue lines)



1/11/21                                                    7

7



1/11/21                                                    8

8

2

## Preprocessing steps

1. Let $\mu = \frac{1}{n}\sum x_i$

2. Replace each $x_i$ with $x_i - \mu$

3. Let $\sigma_j^2 = \frac{1}{n}\sum_i x_{ij}^2$

4. Replace each $x_{ij}$ with $x_{ij}/\sigma_j$

- Step 1 & 2 zero out the mean

- Step 3 & 4 rescale each coordinate to have unit variance

1/11/21                                                                 9

9

## Finding $u_1$

- Given an unit vector $u_1$ and a point $x$, the length of the projection of $x$ onto $u_1$ is given by $x \cdot u_1 = x^T u_1$

- Our task becomes to select a unit-length $u_1$ to maximize

$$\frac{1}{n}\sum(x_i^T u_1)^2 = \frac{1}{n}\sum(u_1^T x_i x_i^T u_1) = u_1^T\left(\frac{1}{n}\sum x_i x_i^T\right)u_1$$
$$= u_1^T \Sigma u_1$$

Subject to $\|u_1\|_2 = 1$

➢ $\Sigma = \frac{1}{n}\sum x_i x_i^T$ is the covariance matrix

    ➢Remember that we replaced each $x_i$ with $x_i - \mu$

1/11/21                                                                 10

10

## Finding $u_1$ (cont')

- By Lagrange multiplier, we have
$$\mathcal{L} = u_1^T \Sigma u_1 + \lambda(1 - u_1^T u_1)$$

- Take derivative and set to 0
$$\frac{\partial \mathcal{L}}{\partial u_1} = 2\Sigma u_1 - 2\lambda u_1 = 0 \Rightarrow \Sigma u_1 = \lambda u_1$$

  So, $u_1$ is an eigenvector of $\Sigma$ with eigenvalue $\lambda$

- Since we want to maximize $u_1^T \Sigma u_1$, $u_1$ must be the eigenvector with maximum eigenvalue of $\Sigma$

1/11/21                                                                 11

11

## Finding $u_2$

- Select a unit-length $u_2$ to maximize
$$u_2^T \Sigma u_2$$

  Subject to $\|u_2\|_2 = 1$ and $u_2^T u_1 = 0$

- Lagrange form
$$\mathcal{L} = u_2^T \Sigma u_2 + \lambda_1(1 - u_2^T u_2) + \lambda_2 u_2^T u_1$$

- Taking derivative and set to 0
$$\frac{\partial \mathcal{L}}{\partial u_2} = 2\Sigma u_2 - 2\lambda_1 u_2 + \lambda_2 u_1 = 0$$
$$\Rightarrow 2u_1^T \Sigma u_2 - 2\lambda_1 u_1^T u_2 + \lambda_2 u_1^T u_1 = 0$$
$$\Rightarrow 0 - 0 + \lambda_2 = 0$$
$$\Rightarrow \lambda_2 = 0$$

1/11/21                                                                 12

12

3

## Finding $u_2$ (cont')

- Taking derivative and set to 0
$$\frac{\partial \mathcal{L}}{\partial u_2} = 2\Sigma u_2 - 2\lambda_1 u_2 = 0$$
$$\Rightarrow \Sigma u_2 = \lambda_1 u_2$$
  So, $u_2$ is an eigenvector of $\Sigma$ with eigenvalue $\lambda_1$
- Since we want to maximize $u_2{}^T \Sigma u_2$ and $u_2 \neq u_1$ ($\because u_2^T u_1 = 0$), $u_2$ must be the eigenvector with second largest eigenvalue of $\Sigma$

1/11/21     13

13

## PCA in general

- Perform PCA by computing the eigenvectors of the $k$ largest eigenvalues of the covariance matrix
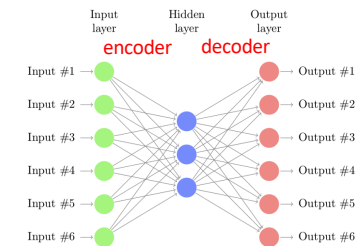
1/11/21     14

14

## Autoencoder

1/11/21     15

15

## Autoencoder

- An autoencoder is a neural network whose outputs are its own inputs
  – Unsupervised (or self-supervised) learning
- Objective: minimize reconstruction error



1/11/21     16

16

## Autoencoder

- Define
$$a_i = g(Wx_i),$$
$$\widetilde{x}_i = Va_i = Vg(Wx_i)$$
- Target: minimize $\frac{1}{n}\sum_{i=1}^{n}(x_i - \widetilde{x}_i)^2$
- If $g$ is linear, then:
$$\widetilde{x}_i = VWx_i$$
- Target: minimize
$$\frac{1}{n}\sum_{i=1}^{n}(x_i - VWx_i)^2 = \frac{1}{n}\sum_{i=1}^{n}(x_i - Ux_i)^2$$
  - Optimal solution is PCA!
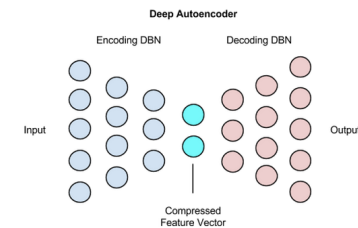
17

## Autoencoder as non-linear PCA

- If $g$ is non-linear, then we have non-linear dimension reduction
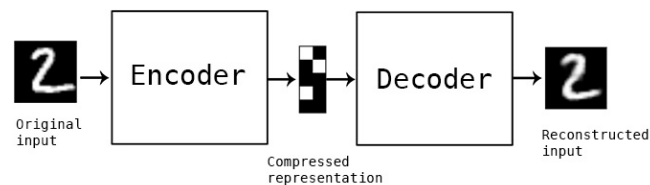- We may further use deep autoencoder to perform non-linear dimension reduction



**Deep Autoencoder**

18

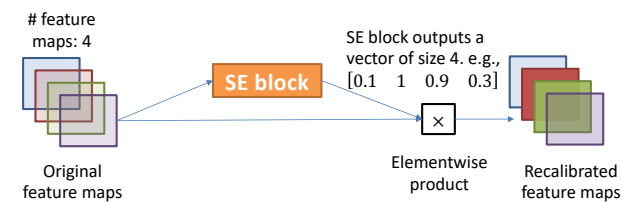## Autoencoder example

19

## Remember the SE network?



- SE block outputs a vector $v = [v_1, \dots, v_d]$
  - $d = $ # feature maps
- Recalibrate the feature map $M^{(i)}_{\text{original}}$ by
  - $M^{(i)}_{\text{recalibrated}} = v_i \times M^{(i)}_{\text{original}}$
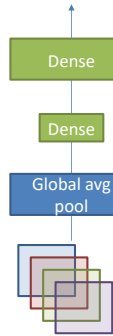
20

## SE block architecture

- Global avg pool: mean for each feature map
  - If $d$ feature maps, get $d$ avg pools
    - E.g., 256 maps, avg pools $\boldsymbol{a} = [a_1, \ldots, a_{256}]^T$
- Consider SE block as an autoencoder
  - Learn $\boldsymbol{W}_1$ and $\boldsymbol{W}_2$ such that $\boldsymbol{W}_2 \mathrm{Relu}(\boldsymbol{W}_1 \boldsymbol{a}) \approx \boldsymbol{a}$
    - $\boldsymbol{W}_1 \in R^{c \times d}, \boldsymbol{W}_2 \in R^{d \times c}, c < d$
      - E.g., $\boldsymbol{W}_1 \in R^{16 \times 256}, \boldsymbol{W}_2 \in R^{256 \times 16}$
  - The small vector, $\boldsymbol{m} = \mathrm{Relu}(\boldsymbol{W}_1 \boldsymbol{a})$, is a compact representation of global avg pool $\boldsymbol{a}$
    - So relationship among $a_i$s are captured by $\boldsymbol{m}$
- The output of SE block is $\sigma\big(\boldsymbol{W}_2 \mathrm{Relu}(\boldsymbol{W}_1 \boldsymbol{a})\big)$

Dense

Dense

Global avg pool

1/11/21     21

21

## Stacked autoencoder



Encoder₁   Encoder₂   Decoder₂   Decoder₁

input layer   bottleneck hidden layer   output layer
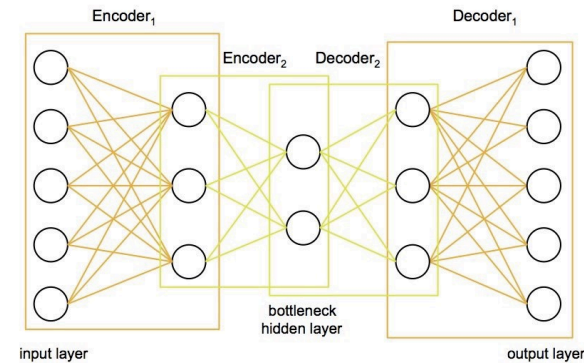
1/11/21     22

22

## Convolutional autoencoder

- Two main structure of a CNN:
  - Convolution
  - Pooling
- Convolutional autoencoder
  - Encoder
    - Consists of convolutional layers and pooling layers
    - Downscale spatial dimensionality (i.e., height and weight) but increase depth (i.e., # feature maps)
  - Decoder:
    - Upscale spatial size and reduce depth back to original dimensions
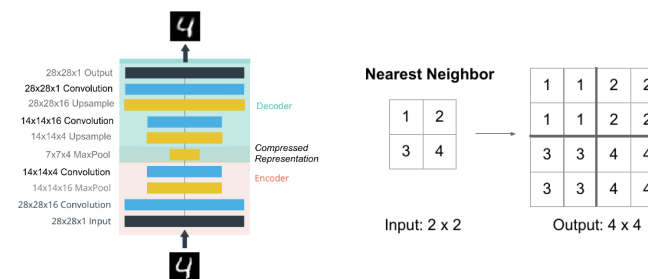
1/11/21     23

23

## Convolutional autoencoder

- Up-sampling



28x28x1 Output
28x28x1 Convolution
28x28x16 Upsample
14x14x16 Convolution
14x14x4 Upsample
7x7x4 MaxPool
14x14x4 Convolution
14x14x16 MaxPool
28x28x16 Convolution
28x28x1 Input

Decoder

Compressed Representation

Encoder

**Nearest Neighbor**

| 1 | 2 |
| 3 | 4 |

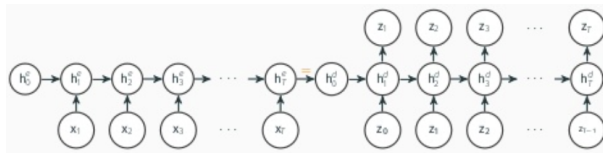| 1 | 1 | 2 | 2 |
| 1 | 1 | 2 | 2 |
| 3 | 3 | 4 | 4 |
| 3 | 3 | 4 | 4 |

Input: 2 x 2     Output: 4 x 4

1/11/21     24

24

## Recurrent autoencoder

- Similar to the encoder-decoder model
- Encoder is a sequence-to-vector RNN
- Decoder is a vector-to-sequence RNN
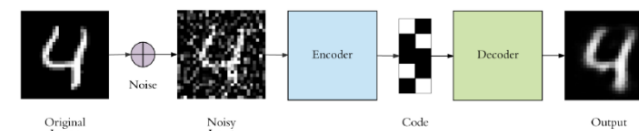


1/11/21                                                          25

25

## Denoising autoencoder

- Add noise to input
- Network tries to recover the original (noise-free) input



1/11/21                                                          26

26

## Variational autoencoder
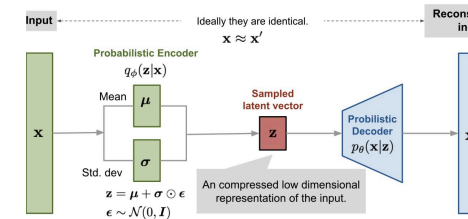
- Variational autoencoder is unique because
  - It is a generative model
  - The output is probabilistic

1/11/21                                                          27

27

## Variational autoencoder



- Encoder generates mean coding $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$
- Actual coding $\boldsymbol{z} \sim N(\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$
- The decoder works as normal
- VAE is a generative model because we can "sample" new $\boldsymbol{z}$s to generate new $\boldsymbol{x}'$

1/11/21                                                          28

28

# Summary

- PCA linearly projects data points into low dimension
  - Good interpretability
  - Can project new data points
- Autoencoder projects the data points non-linearly
  - Interpretability?
  - AE vs Word2Vec
- VAE can generate new instances

1/11/21                                                                                          29

29