

# Variational Autoencoder

11-785 Introduction to Deep Learning  
– lecture 21 –

TAVE Research DL001  
Heeji Won

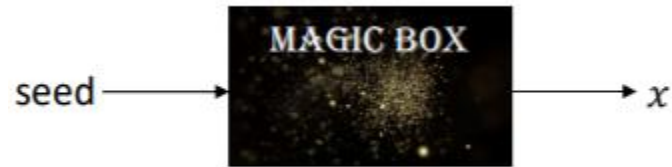
# Contents

1. Generative model
2. How to deal with incomplete data
3. Expectation Maximization
4. PCA
5. VAE

# Contents

1. Generative model
2. How to deal with incomplete data
3. Expectation Maximization
4. PCA
5. VAE

# 01. Generative model




a model that can generate data with a distribution similar to the given data  $x$

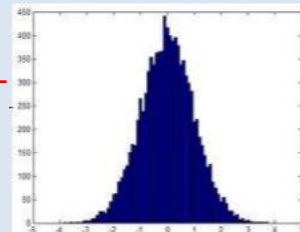
## Learning a generative model

"Estimate the  $\theta$  such that  $P(x; \theta)$  best 'fits' the observations  $X = \{x\}$ "

$$\operatorname{argmax}_{\theta} P(X; \theta) = \operatorname{argmax}_{\theta} \log(P(X; \theta))$$
$$\operatorname{argmax}_{\{p_1, p_2, p_3, p_4, p_5, p_6\}} \sum_i n_i \log(p_i) \quad \leftarrow$$
$$p_i = \frac{n_i}{N} \text{ (N is the total number of observations)}$$



$$\operatorname{argmax}_{\mu, \sigma^2} \sum_{x \in X} \log \text{Gaussian}(x; \mu, \sigma^2) \quad \leftarrow$$
$$\mu = \frac{1}{N} \sum_{x \in X} x; \quad \sigma^2 = \frac{1}{N} \sum_{x \in X} (x - \mu)^2$$

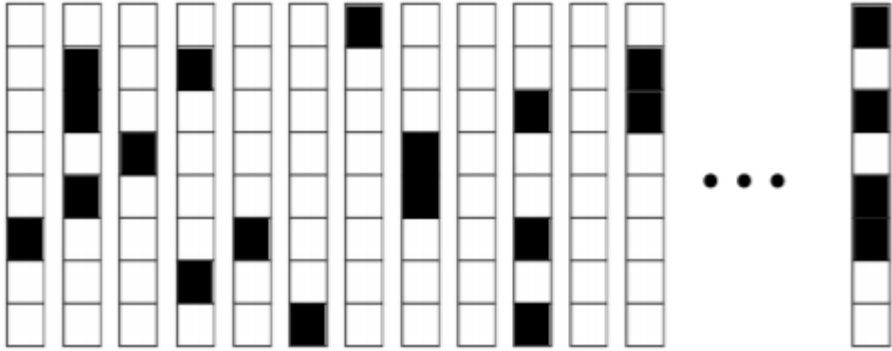


# Contents

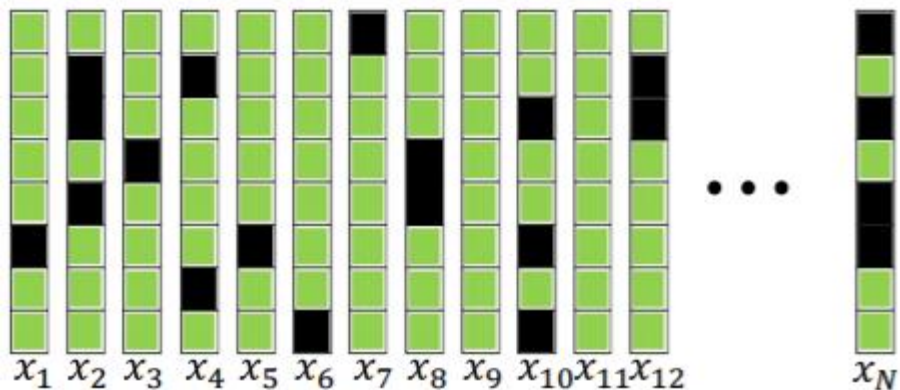
1. Generative model
2. How to deal with incomplete data
3. Expectation Maximization
4. PCA
5. VAE

## 02. How to deal with incomplete data

if the data have missing components



Blacked-out components are missing from data



- Complete data includes the observed & missing components

$$X = \{x_1, \dots, x_N\}, \quad x_i = (o_i, m_i)$$

- Original problem :

$$\operatorname{argmax}_{\mu, \sigma^2} \sum_{x \in X} \log P(x)$$

where  $X$  is the entire data!

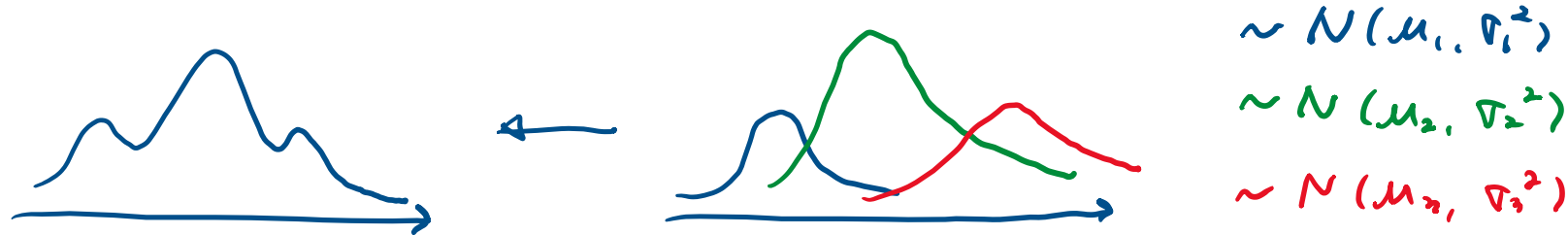
But, there are missing values

$$\operatorname{argmax}_{\mu, \sigma^2} \log(P(O)) = \operatorname{argmax}_{\mu, \sigma^2} \sum_{o \in O} \log P(o)$$

$$= \operatorname{argmax}_{\mu, \sigma^2} \sum_{o \in O} \log \int_{-\infty}^{\infty} P(o, m) dm$$

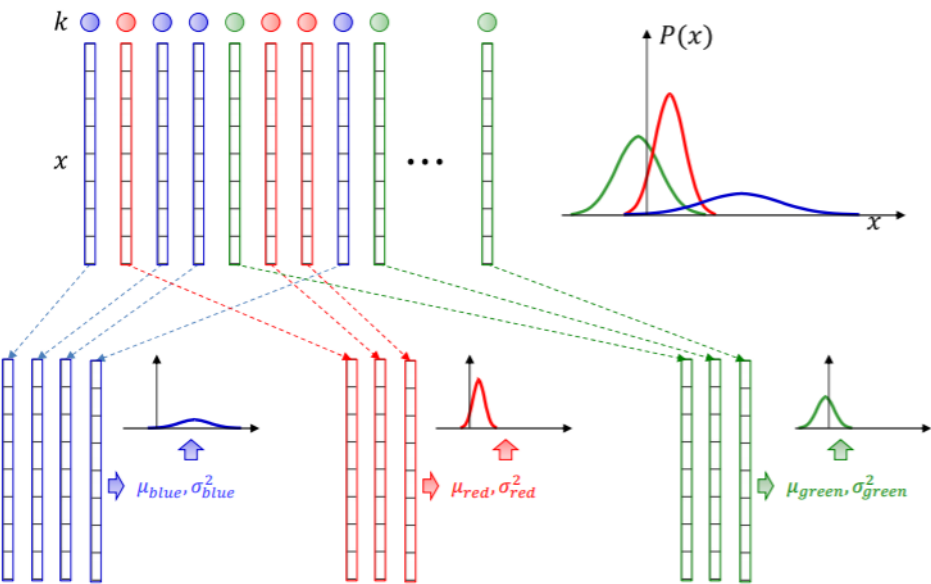
# 02. How to deal with incomplete data

cf) The Gaussian Mixture model

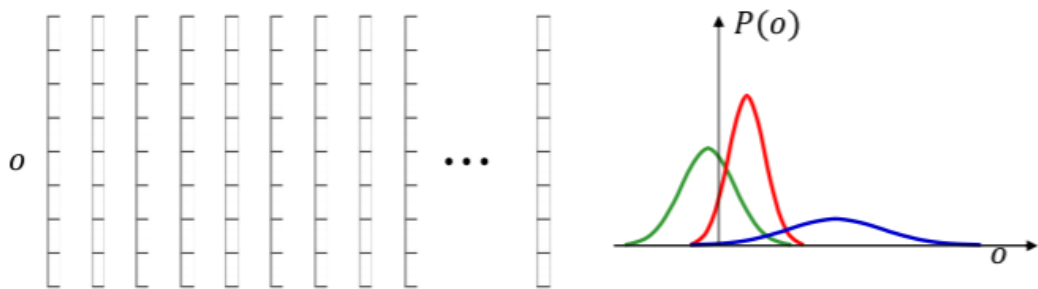


## The structure of the network

– if learning a GMM with 'complete' data



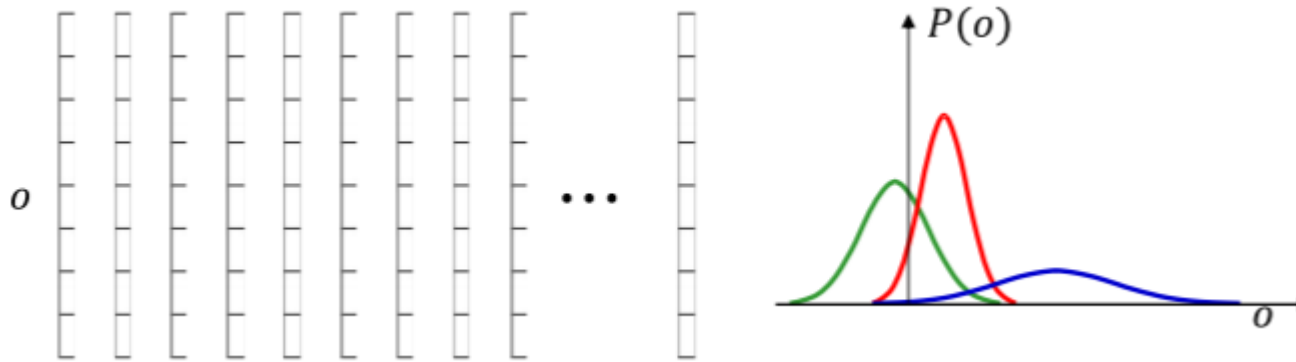
✓ But, we are not given the actual Gaussian for each  $\mathbf{o}_i$



- What we want :  $(o_1, k_1), (o_2, k_2), \dots$
- What we have :  $\mathbf{o}_1, \mathbf{o}_2, \dots$

## 02. How to deal with incomplete data

### The structure of the network



we are not given the actual  
Gaussian for each  $o_i$

– MLE with only observed data

$$\operatorname{argmax}_{\{(\mu_k, \sigma_k^2), \forall k\}} \log(P(O)) = \operatorname{argmax}_{\{(\mu_k, \sigma_k^2), \forall k\}} \sum_{o \in O} \log P(o), \quad P(o) = \sum_k P(k) N(o; \mu_k, \sigma_k^2)$$

$$= \operatorname{argmax}_{\{(\mu_k, \sigma_k^2), \forall k\}} \sum_{o \in O} \log \sum_k P(k) N(o; \mu_k, \sigma_k^2)$$

challenging!  $\Rightarrow$  EM algorithm

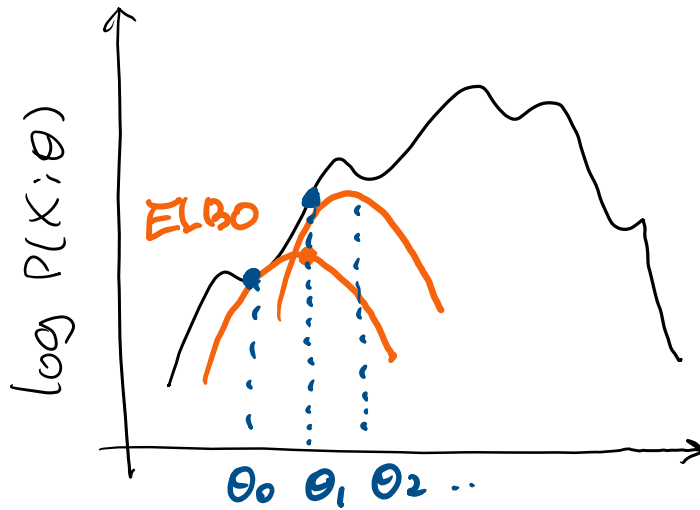
- ✓ no closed form solutions
- ✓ need efficient iterative algorithms



# Contents

1. Generative model
2. How to deal with incomplete data
- 3. Expectation Maximization**
4. PCA
5. VAE

# 03. Expectation Maximization



- initialize  $\theta^0$
  - $k = 0$
  - iterate (over  $k$ ) until  $\sum_{o \in O} P(o; \theta)$  converges:
    - Expectation Step:  
Compute  $P(h|o; \theta)$  for all  $o \in O$  for all  $k$
    - Maximization Step:  

$$\theta^{k+1} \leftarrow \underset{\theta}{\operatorname{argmax}} \underbrace{\sum_{o \in O} P(h|o; \theta^k) \log P(h, o; \theta)}_{\text{ELBO}}$$
- \*  $h$  : missing components

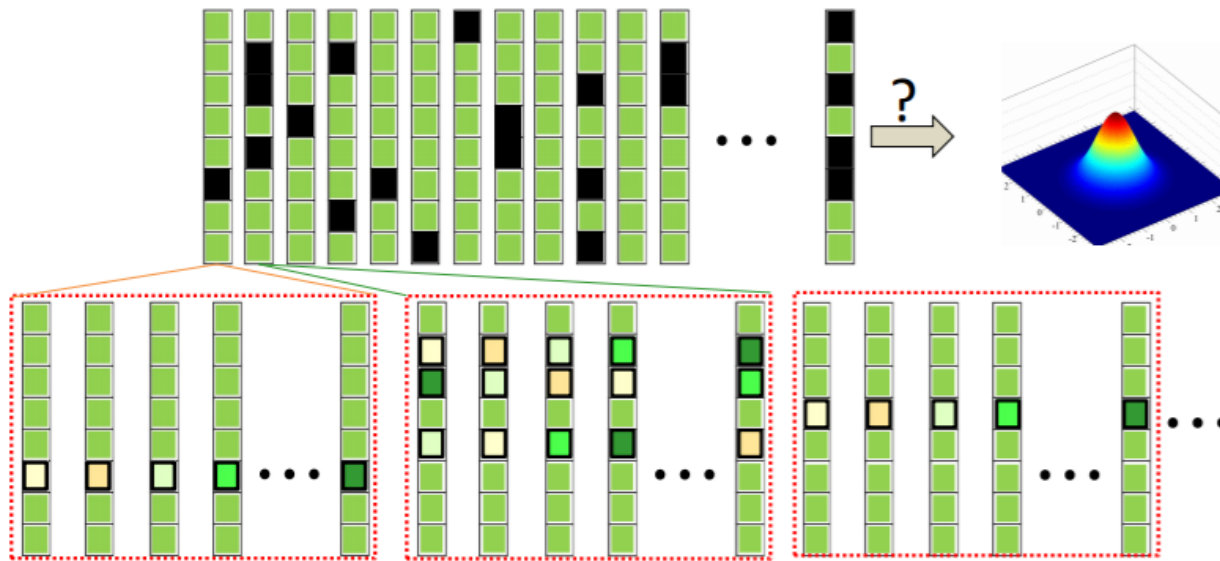
<Construct an ELBO(empirical lower bound function)  $J(\theta, \theta^k)$ >

$$J(\theta, \theta^k) = \underbrace{\sum_{o \in O} \sum_h P(h|o; \theta^k) \log P(h, o; \theta)}_{\text{not a function of } \theta} - \sum_{o \in O} \sum_h P(h|o; \theta^k) \log P(h|o; \theta^k)$$

# 03. Expectation Maximization

if the data have missing components

- Completing incomplete vector



- Let  $x_i(m)$  be the 'completed' version of the observation  $o_i$

$$x_i(m) = (m, o_i)$$

- Estimate from the expanded data

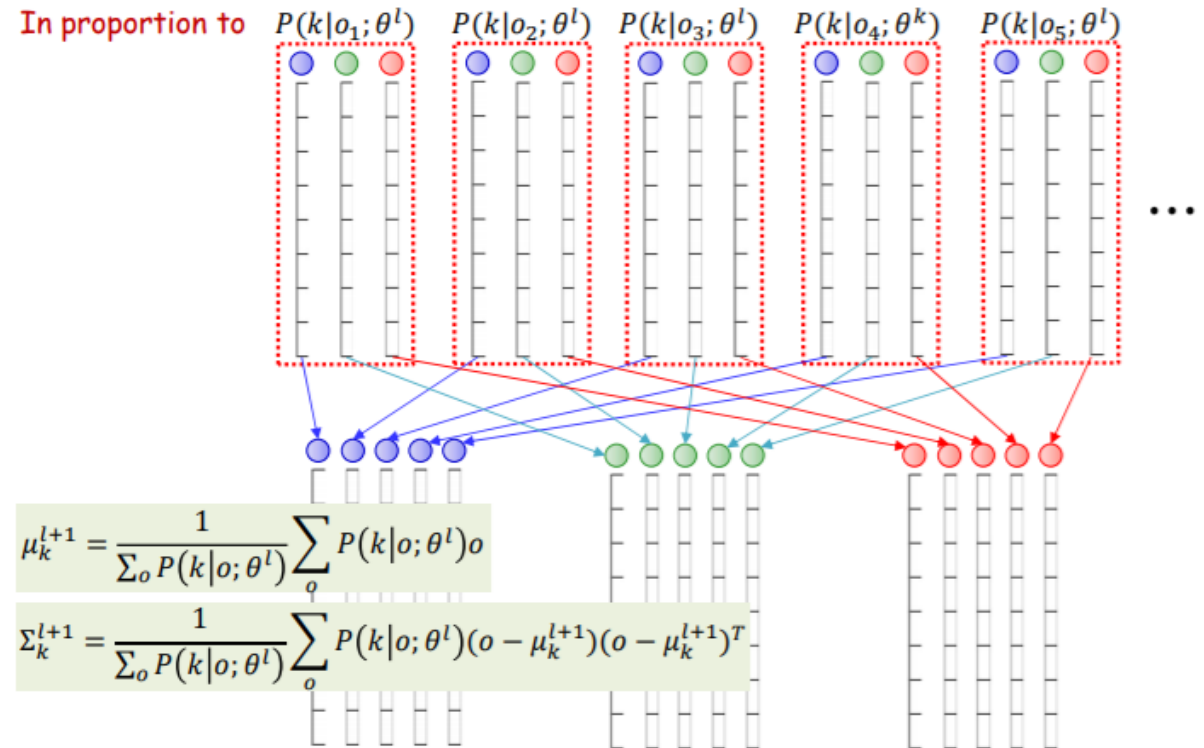
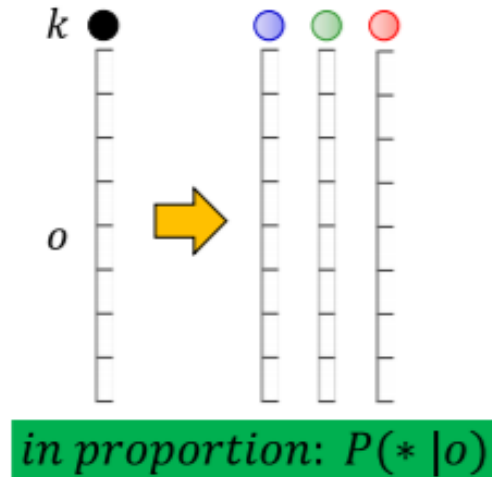
$$\mu^{k+1} = \frac{1}{N} \sum_{o \in O} \int_{-\infty}^{\infty} P(m|o; \theta^k) x_i(m) dm$$

$$\Sigma^{k+1} = \frac{1}{N} \sum_{o \in O} \int_{-\infty}^{\infty} P(m|o; \theta^k) (x_i(m) - \mu^{k+1})(x_i(m) - \mu^{k+1})^T dm$$

- Expand every incomplete vector out into all possibilities
- in proportion:  $P(m|o)$  from a previous estimate of the model

# 03. Expectation Maximization

## The structure of network



Proportion to  $P(k|o)$  which can be computed if we know  $P(k)$  and  $P(o|k)$

from previous estimate of model

iterate!

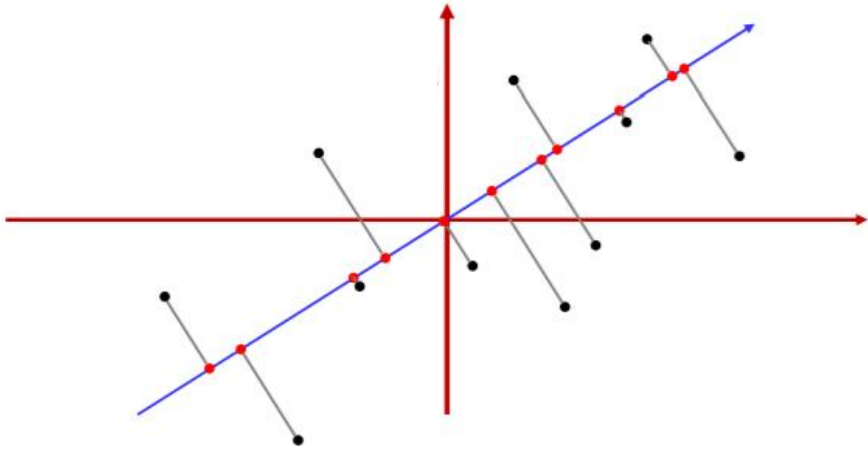
< EM principle >

- ✓ 'Complete' the data by considering every possible value for missing data in proportion to posterior prob.
- ✓ Re-estimate parameters

# Contents

1. Generative model
2. How to deal with incomplete data
3. Expectation Maximization
- 4. PCA**
5. VAE

# 04. PCA

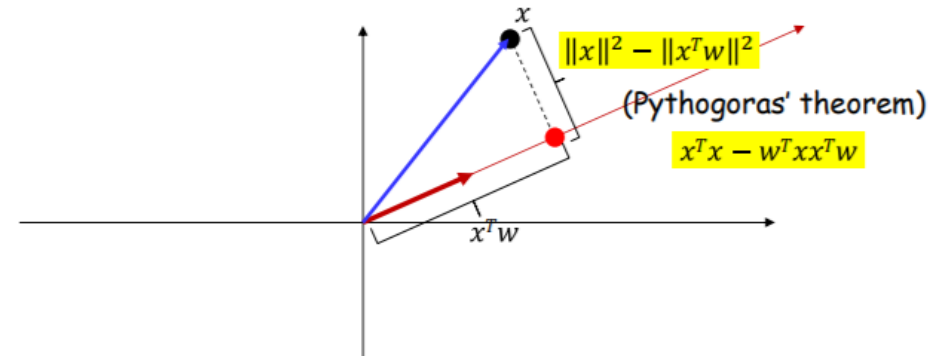


- Find the principal subspace that can be projected
- Minimize the sum of the squared lengths
- There are several method to find

## 1) Search method

"search through all subspaces with minimum projection error"

## 2) Close form



minimizing  $L_2$  error :

$$L = \sum_x x^T x - w^T x x^T w$$

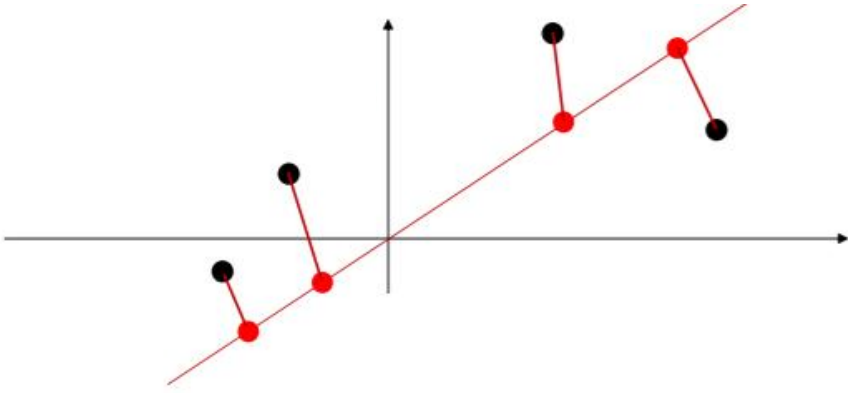
minimizing L

$$\left( \sum_x x^T x \right) w = \lambda w$$

eigenvalue equation 14

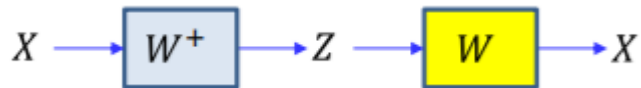
# 04. PCA

## 3) The iterative algorithm



"Let  $W$  rotate and stretch/shrink, keeping the arrangement of  $Z$  location fixed"

## Drawing this differently

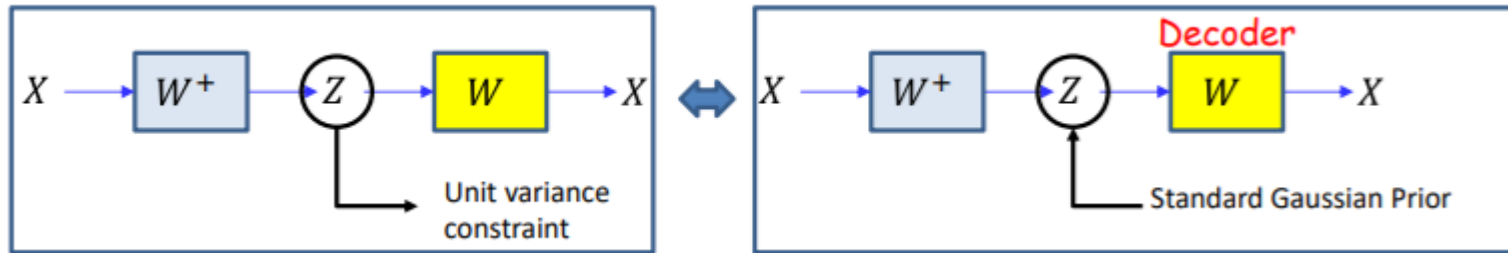


- Autoencoder with linear activation!
- But, the solution is not unique!
  - ✓ Scale invariance
  - ✓ Rotation invariance

# 04. PCA

Find a unique  $W$

1. Orthogonal & unit eigen vector : standard eigen vector
2. Constrain the variance of  $Z$  to be unity



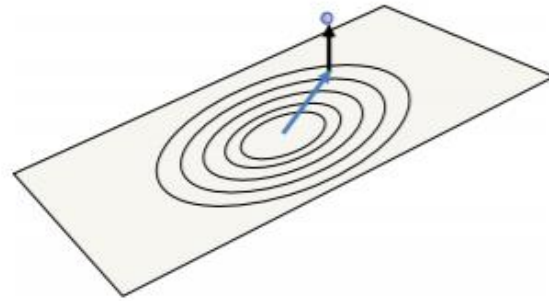
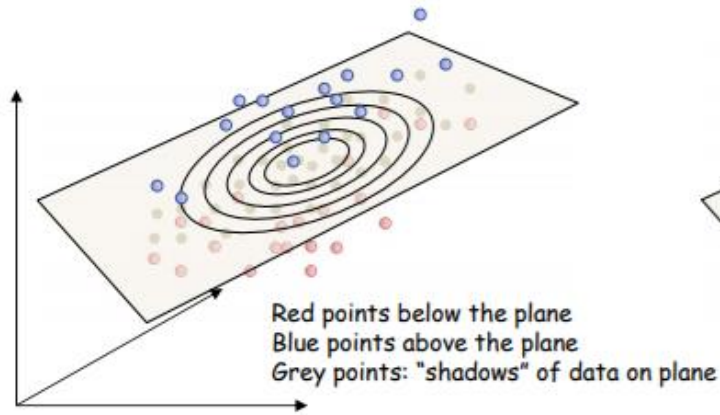
- Imposing the constraint that  $z$  must have unit variance is the same as assuming that is **drawn from a standard Gaussian**
- The decoder of AE with the unit-variance constraint on  $z$  is in fact a **Generative model**



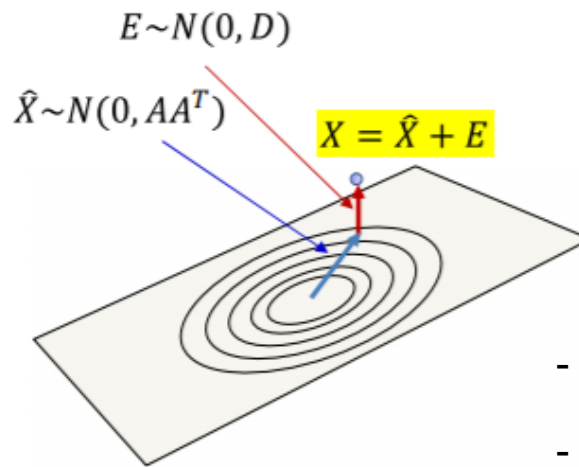
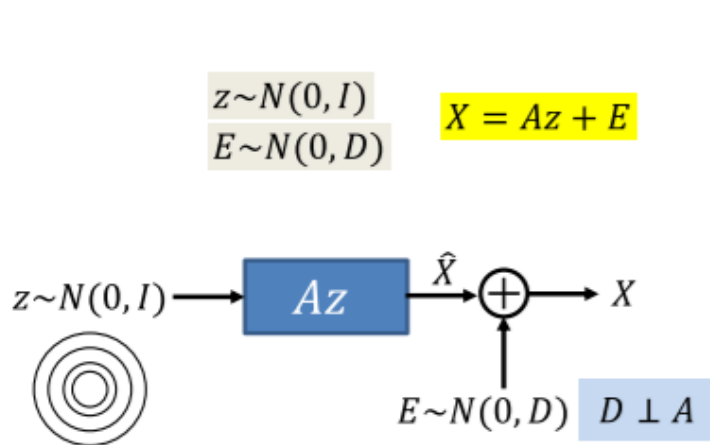
# Contents

1. Generative model
2. How to deal with incomplete data
3. Expectation Maximization
4. PCA
5. VAE

## 05. VAE



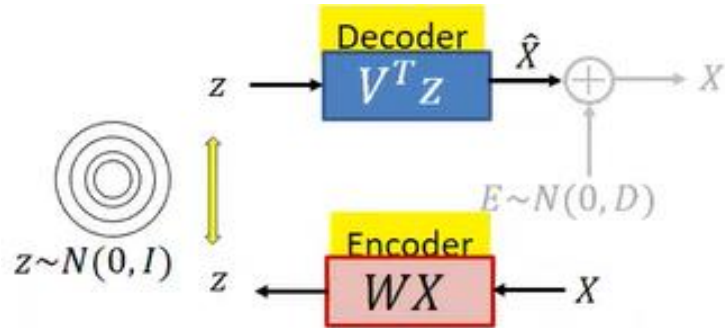
- take a Gaussian step on the principal plane
- take a orthogonal Gaussian step where we land to generate a point



- $z$  is drawn from  $K$ -dim isotropic Gaussian
- $A$  is a basis matrix
- $E$  is a Gaussian noise that is orthogonal

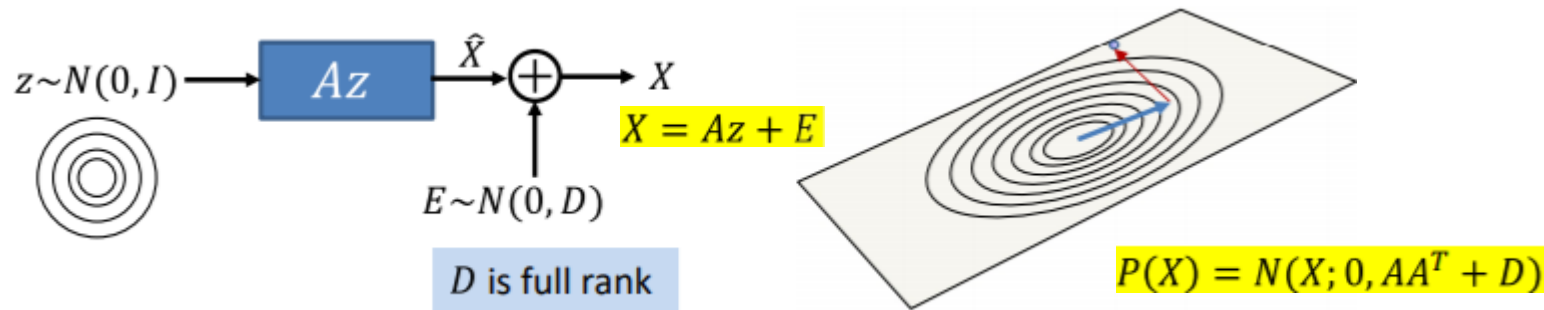
\* covariance is low-rank

## 05. VAE



- The decoder weights are just the PCA basis matrix
- Encoder: transforms input  $X$  into Gaussian  $z$
- Decoder: transforms Gaussian  $z$  into principal subspace reconstruction  $\hat{X}$

## The Linear Gaussian Model



also a generative model!

Also known as **Factor Analysis**

- Update the model : The noise can lie in any direction
  - Noise is drawn from full-rank uncorrelated Gaussian distribution
- ⇒ The way to produce any data instance is no longer unique!



- $A$  is the loading matrix
- $z$  are the factors
- $D$  is diagonal

Thank you