

BLG 561 E FALL 2021

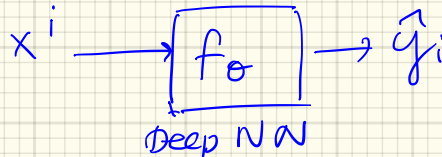
Deep Learning

30.11.2021

Görde ÜNAL

Supervised Learning: We've done this so far

$\{x^i, y^i\}_{i=1}^m$   
Input Data      Labels



We require lots of

generate → **LABELLED DATA**

Classification, Regression, Obj. Detection, Segmentation, Tracking.

Learn to map the Input → Label prediction

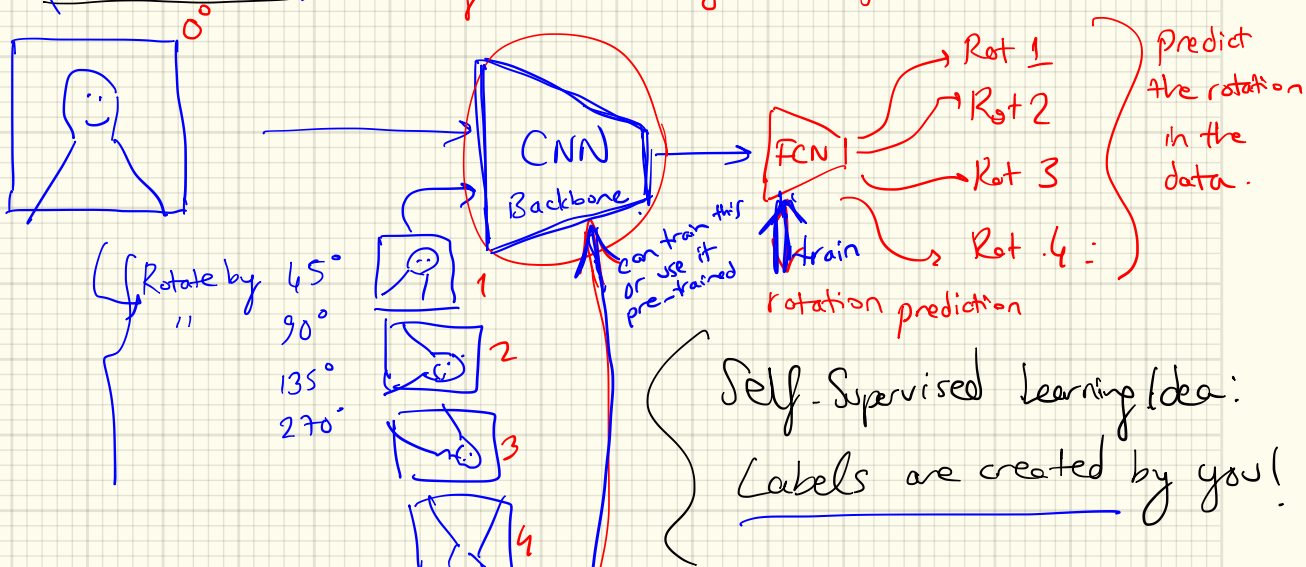
Self-Supervised Learning:

Let the data generate labels

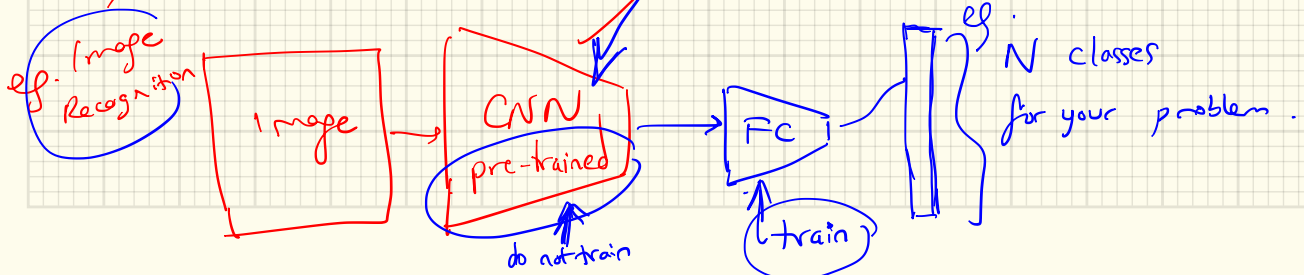
1) Pretext Task ←

2) Downstream Task ← final goal :

# 1) Pretext Task: Predicting Rotation of the Image

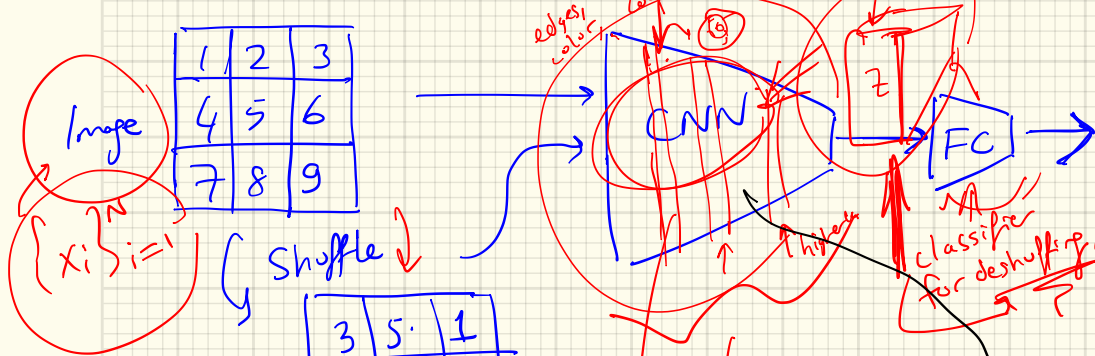
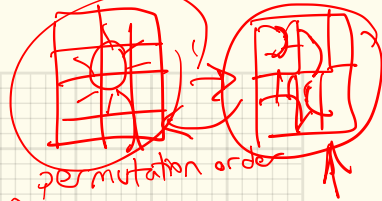


# 2) Downstream Task : → your real goal



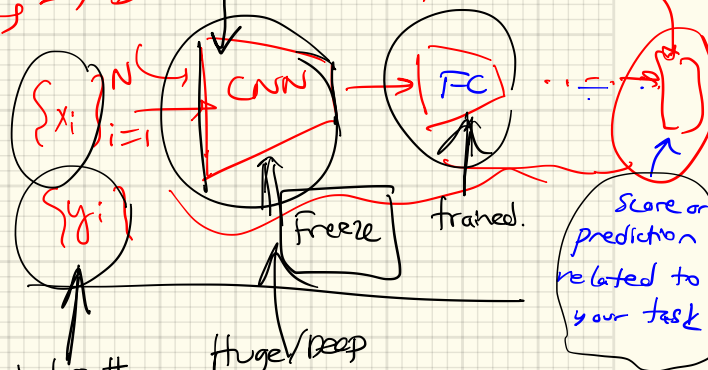
# 1) Another Pretext Task: Shuffling the data

auxiliary task used in self-supervised learning



learns to deshuffle

## 2) Downstream task



Idea: Now, this CNN captures the input data characteristics better w/ self-supervision.

→ For a certain downstream task, I still need some labeled data, but not as much of labeled data as before.

Huge/Deep CNN is now pre-trained

NLP (GPT) : Sentence : "I am going to school".

Self-supervision pretext task : Shuffle the words in the sentence

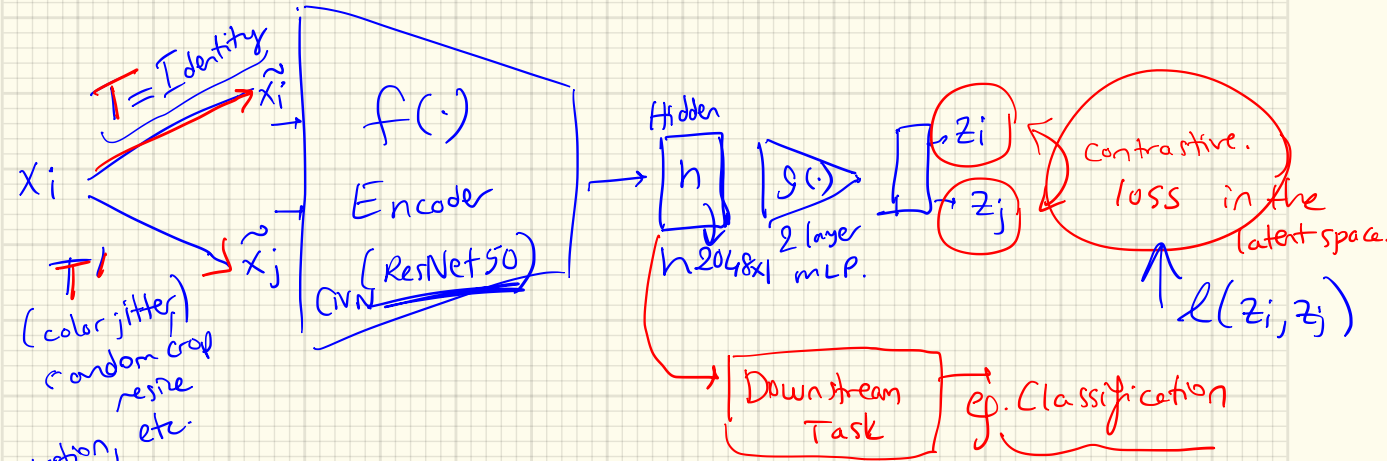
pretext task : 1.1) NV predicts the correct order of the sentence.

1.2) Delete a word from the sentence & predict the missing word(s).

Open Directions : Designing appropriate pretext tasks for your own problem.

- Incorporating self-supervision idea into your research work in a novel way.

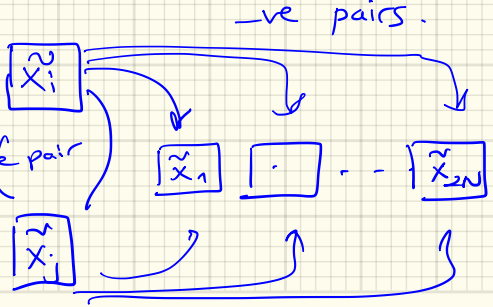
→ Contrastive Learning (SimCLR paper Hinton et al., 2020.)



$\{\tilde{X}_i, \tilde{X}_j\}$  : +ve pair

$N$  Batches :  $2N-1$  -ves  
 $2N$  pairs

$\{\tilde{X}_i, \tilde{X}_k\}$  +ve pair



$$\min_{\text{anchor}} \text{loss}(z_i, z_j) = -\text{sim}\left(\frac{z_i}{\tau}, \frac{z_j}{\tau}\right) + \log \sum_{\substack{k=1 \\ k \neq i}}^N \exp\left(\frac{\text{sim}(z_i, z_k)}{\tau}\right)$$

↑
↑
↑
↑
↑

anchor
true pair
temperature

$$\log \left( \exp\left(\frac{\text{sim}(z_i, z_1)}{\tau}\right) + \exp\left(\frac{\text{sim}(z_i, z_2)}{\tau}\right) + \dots + \exp\left(\frac{\text{sim}(z_i, z_N)}{\tau}\right) \right)$$

(Unsupervised Learning) : No Labels !!!

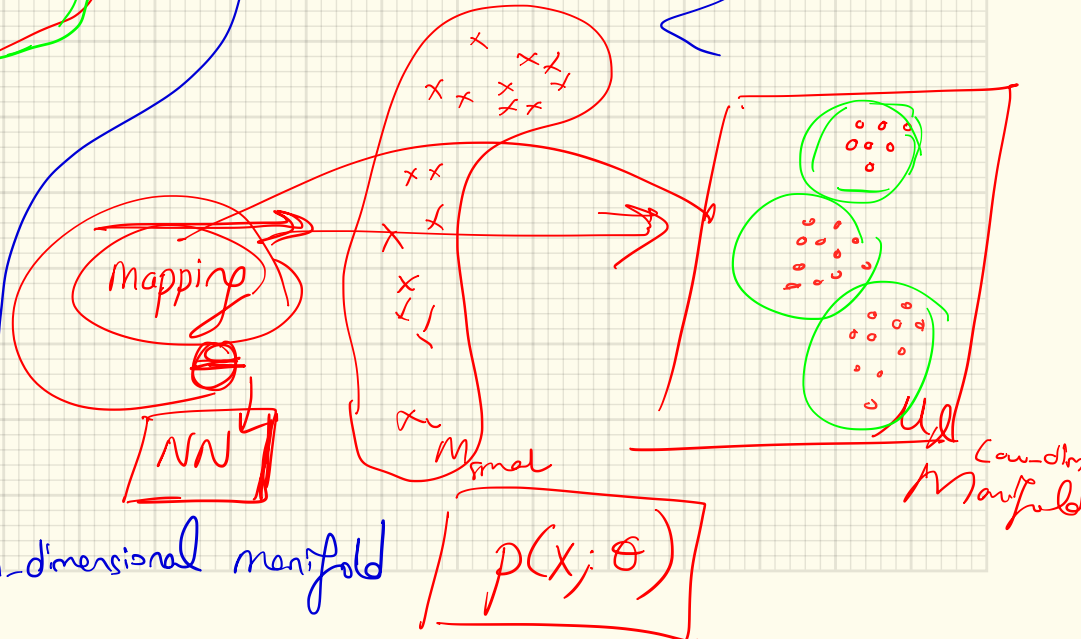
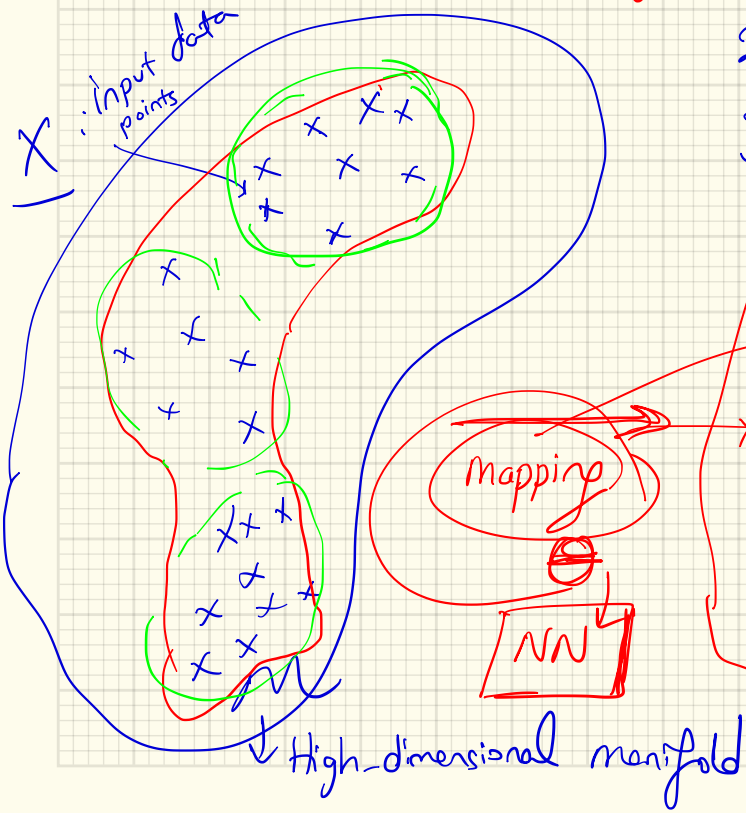
Unsupervised  
Segmentation.

Representation Learning

1) Clustering

2) Recommendation

3) Dimensionality Reduction





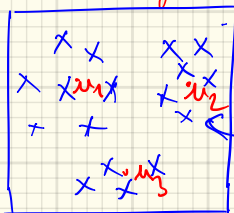
Goal : Learn the distrib. of the input data  $p(x; \theta)$   
from training data  $\left\{ x_i \right\}_{i=1}^m$  w/ NO LABELS.

## Unsupervised Learning Tasks:

1) Clustering; 2) Dimensionality Reduction.

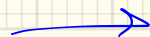
[1] K-means Clustering : Hypothesis function.

Reformulate K-means in an ML framework



$\theta$  : cluster means:  $\mu_1, \dots, \mu_K$   
 $\left\{ x_i \right\}_{i=1}^m$  : Dataset w/ no Labels

① Hypothesis Class:  $\underset{j}{\text{arg min}} \|x^i - \mu_j\|_2^2 = h_\theta(x^i)$



② Loss function:  $\mathcal{L}(x^i, h_\theta(x^i)) =$

$$\mathcal{L} = \sum_{i=1}^m \left\| x^i - \mu_{\left(\arg \min_j \|x^i - \mu_j\|_2\right)} \right\|_2^2$$

③ Optimization: Q. Is the loss  $\mathcal{L}$  a convex loss fn? No!

$\arg \min_{\theta} \mathcal{L}(x, h_\theta(x)) \rightarrow$  non-convex optim. problem.

$$\theta = \{\mu_1, \dots, \mu_k\}$$

We solve it approximately by an iterative approach.

Alternate between (i) & (ii):

(i) Find optimal  $\mu_i$  for each <sup>data</sup> instance  $i = 1, \dots, k$ .

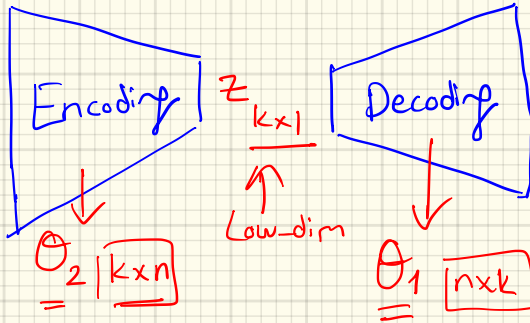
(ii) Then we update  $\mu_i$  to be average of all samples assigned to  $\mu_i$  until "Convergence".

## 2 PCA (Dimensionality Reduction)

Given  $\{x_i\}_{i=1}^m$ ,  $x_i \in \mathbb{R}^n$ ,  $n \uparrow \uparrow$ .

Input Data

$\underline{X}$   
 $n \times 1$   
 high-dim



Reconstructed  
 data

Idea: Reconstruct the data so that "most of the information in  $\underline{x}$  is preserved."

① Hypothesis Class:  $\hat{x} = h_{\Theta}(x) = \underline{\Theta}_1 \underline{\Theta}_2 x$  ✓

② (Reconstruction) Loss:  $d_r(\hat{x}, x) = \|x - h_{\Theta}(x)\|_2^2$

\* most widely used Unsupervised Learning loss

③ Minimize the Reconstruction Loss:  $\underline{\Theta} = [\underline{\Theta}_1 \ \underline{\Theta}_2]$

$\min_{\underline{\Theta}_1, \underline{\Theta}_2} \sum_{i=1}^m \|x^i - \underline{\Theta}_1 \underline{\Theta}_2 x^i\|_2^2$  ;  $\Theta$ : Is this a convex problem?  $\rightarrow$

$\hat{x}^i$ : reconstructed data

→ Not a convex problem: but  $\exists$  a solution through EVD  
(Eigen Value Decomposition)

→ Given  $\{x_i\}_{i=1}^m$  :  $\underline{X} = \begin{bmatrix} x^1 & x^2 & \dots & x^n \\ \vdots & \vdots & & \vdots \end{bmatrix}_{n \times m}$

Q. How do you perform EVD on  $\underline{X}$ ?  
analysis

Calculate COVARIANCE matrix of the input (training) data matrix.

$\underline{C} = E \left[ \underbrace{\tilde{X}}_{n \times m} \underbrace{\tilde{X}^T}_{m \times n} \right]$  ← of course, there's centering of the data.

perform EVD on

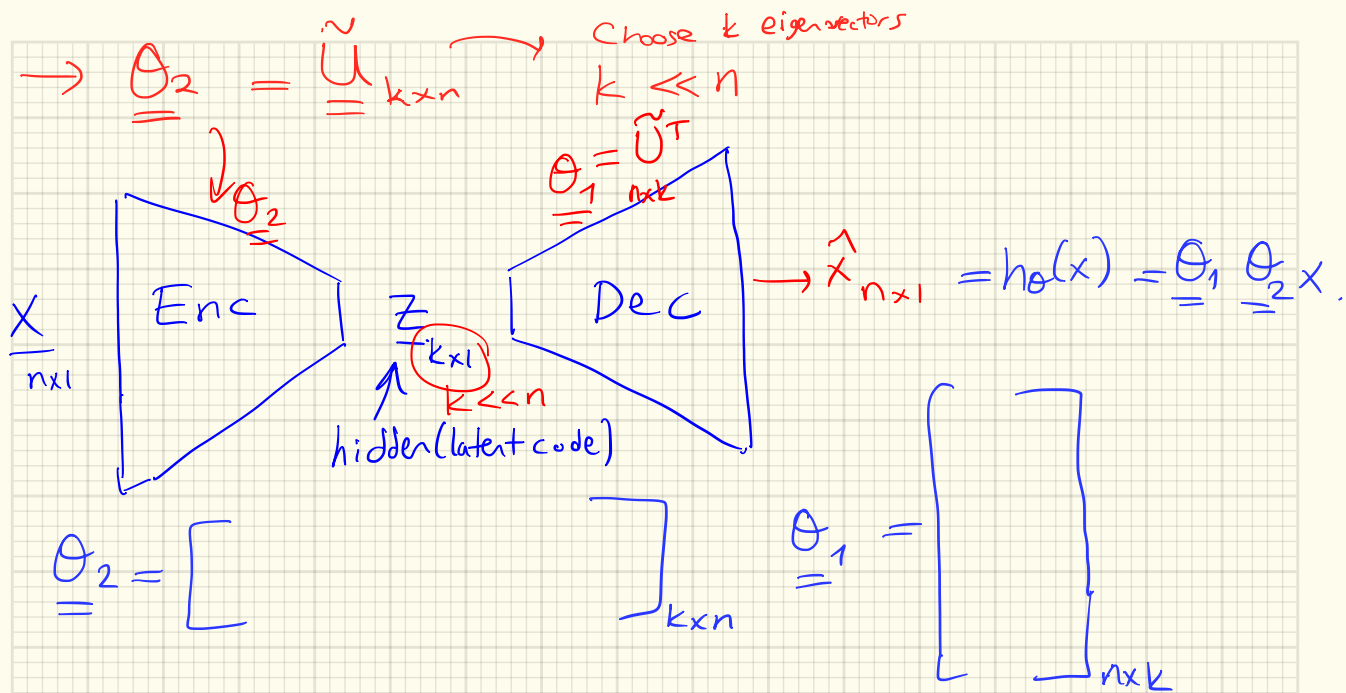
$\underline{C} = \underline{U} \underline{\Sigma} \underline{U}^T$   
matrix of vectors: orthogonal matrix

$\begin{bmatrix} d_1 & & & 0 \\ & d_2 & & \\ & & \dots & \\ 0 & & & d_n & & \\ & & & & & 0 \end{bmatrix}_{n \times n}$   
e. value  
min

Set  $\underline{O}_2 = \underline{U}$

→  $\underline{\tilde{U}}$  → eliminating "small" dimensions w.r.t small e. values  
 $(u_1 \dots u_k)_{k \times n}$

$\underline{\Sigma} = \begin{bmatrix} d_1 & & & 0 \\ & d_2 & & \\ & & \dots & \\ & & & d_k & & \\ & & & & & 0 \end{bmatrix}$

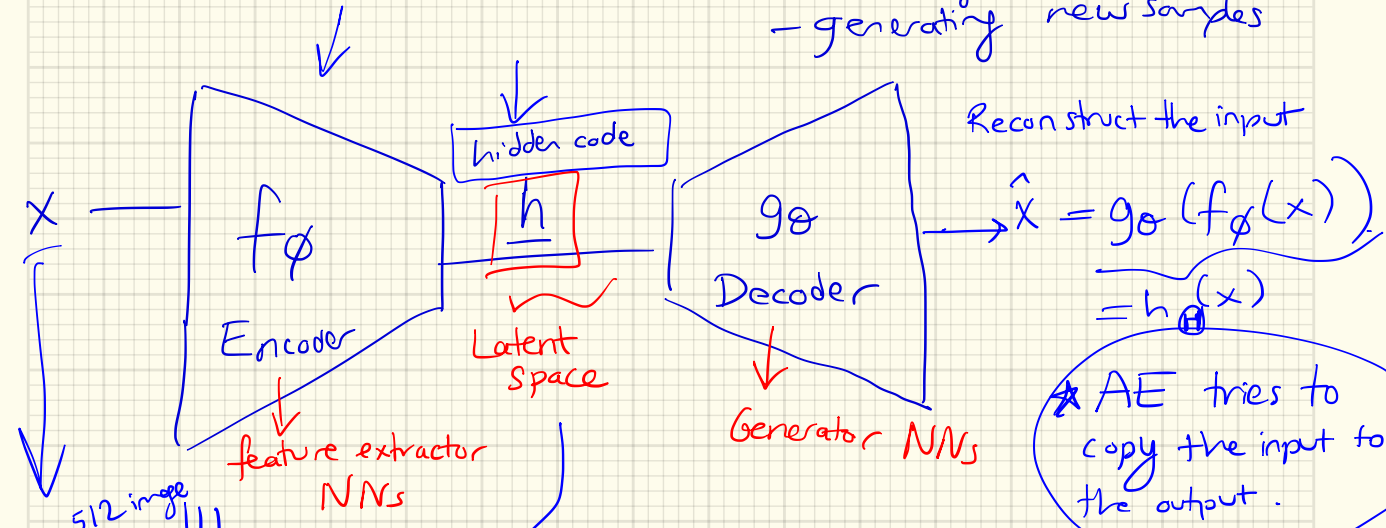


☆ PCA:  $\approx$  a NN w/o nonlinear activation functions  
 without

$\downarrow$  AutoEncoders are like PCA but w/ NONLINEAR NN models.

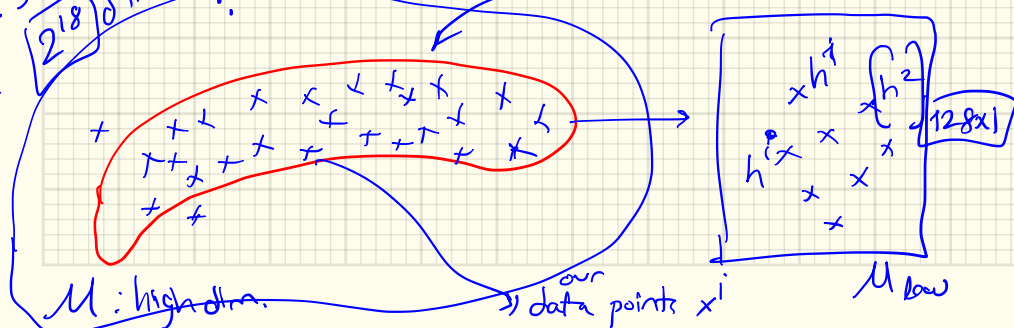
# AUTOENCODERS (AEs) → feature extraction

- dimensionality reduction
- clustering
- generating new samples

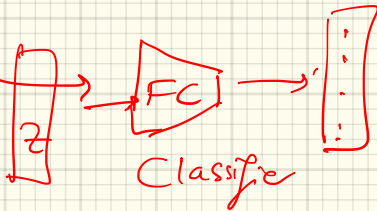
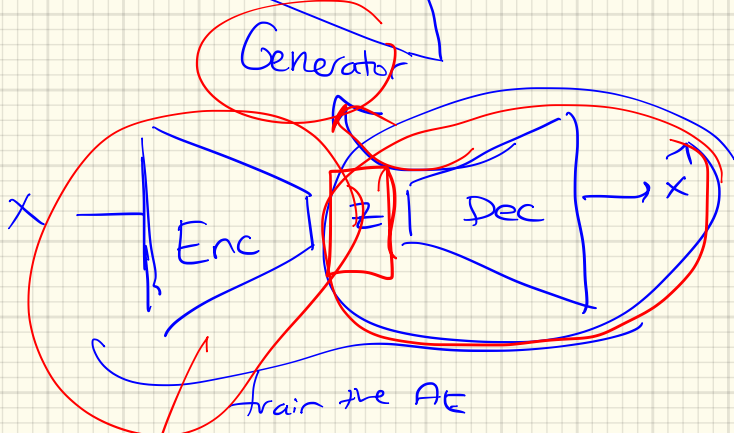
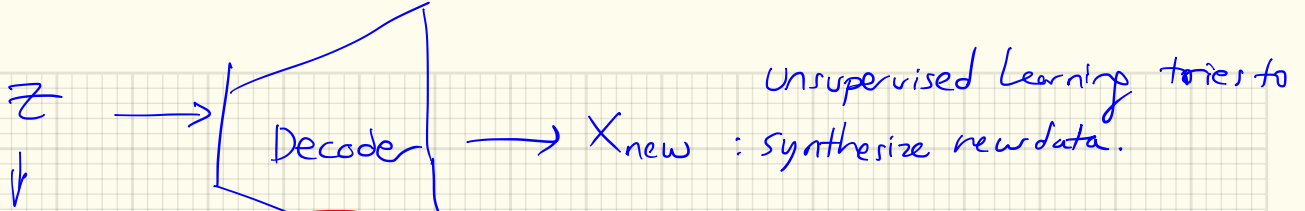


★ AE tries to copy the input to the output.

eg.  $512 \times 512$  image  
 $M: \begin{bmatrix} 218 \\ \dots \end{bmatrix}$  dim. !!!



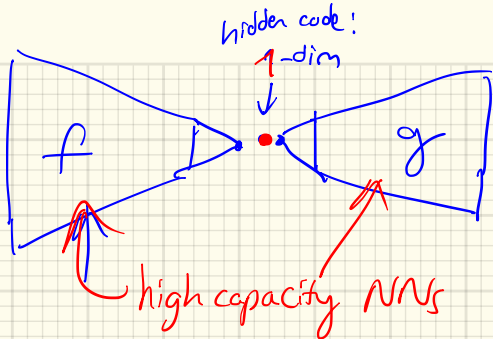
- ★ Input Data:
- Images (Face)
  - Places
  - Text + Document
  - music (Audio Libraries)
  - Time Series Data



Extreme ex: AE:

$$x^i_{n \times 1}$$

$\{x^i\}_{i=1}^m$ : millions of data



$$\hat{x}^i \approx x^i \leftarrow \text{normally}$$

$$\hat{x}^i = x^i$$

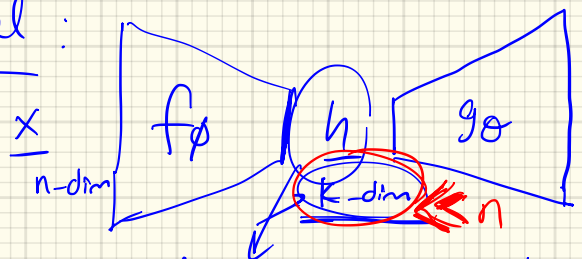
exact reconstruct.

This AE does not learn anything useful about the data

!!  
 $f$ : maps each input in the training set to its INDEX in the dataset.  
 $g$ : maps the index back to  $x^i$  again.

\* Think about the generalization problem of this extreme case.

Hourglass model:



reconstruct.

$$\hat{x} = h_{\theta}(x) \approx x$$

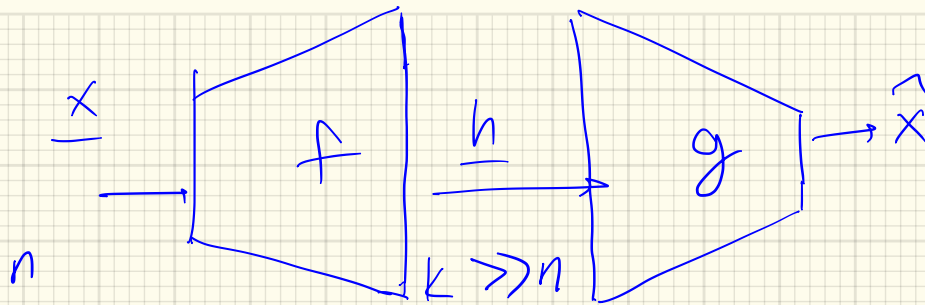
n-dim.

x goes thru INFORMATION bottleneck to get reconstructed.

\* UNDERCOMPLETE Representation:



Q.



over-complete  
representation

? not preferred?

Typically  $h$  spans a subspace of the  $X$ .

Undercomplete representation:  $\rightarrow$  tends to capture the

$\rightarrow$  "most salient features"  
of the input data distribution

$\rightarrow$  How?

By  $\rightarrow$  ① Limiting the capacity of the AE (hourglass model)

$\rightarrow$  ② Regularizing the AE  $\equiv$  adding some constraints

↳ Unconstrained Optimization Objective:

$$L(\theta) = \underbrace{L_{\text{data}}(\theta)} + \lambda \cdot L_{\text{regularize}}(\theta)$$

Data Term of the Loss (Objective) function:

$$L_{\text{data}}(\theta) = L_{\text{reconstruction}}(x^i, h_{\theta}(x^i)) = \sum_{i=1}^m \|x^i - \underbrace{h_{\theta}(x^i)}_{\text{typically}}\|_2^2$$

1) Hypothesis fn for the AE,

$$g. \underbrace{h_{\theta}(x)}_{L\text{-layer FCN (MLP) or CNN}} = f^L(\underline{W}^L * (\underbrace{f^{L-1}(\underline{W}^{L-1} * (\dots f^2(\underline{W}^2 * f^1(\underline{W}^1 * \underbrace{x}_{a^0}})) \dots)) \dots))$$

$\xleftarrow{\text{upconv}} \quad \xleftarrow{\text{conv-layers}}$

# Regularizers for AEs:

- i) Sparse representation : Sparse AE
- ii) Contracting AE : CAE ←
- iii) Denoising AE : → robustness to noise

✓ push the AE to have other properties in addition to its ability to copy the input to the input.

i) Sparse AE: adds a "sparsity" penalty in  $\underline{h}$  code:

$$\underset{\phi, \theta}{\text{arg min}} \mathcal{L} = \mathcal{L}_{\text{rec}}(x, \underbrace{g_{\theta}(f_{\phi}(x))}_{\hat{x}}) + \boxed{d} \text{Reg}(\underline{h})$$

$$\rightarrow \text{SAE} : \text{Reg}(\underline{h}) = \|\underline{h}\|_1$$

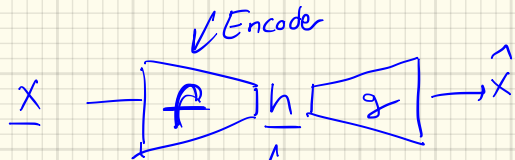
put sparseness constraint.  
→ L1 norm.

d: regularizer weight → important hyperparameter in the optimization.  
to be tuned.

→  $L_{reg} = \|\underline{h}\|_1 = \sum_i |h_i|$  : L1 regularization on the latent space.

(Recall: L1 norm is induced by a Laplacian prior on the latent code  $h$  from probabilistic MC interpretation before)

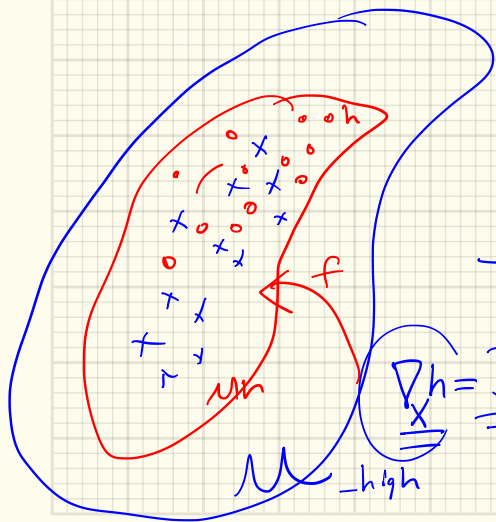
(ii) Contractive AE: (CAE)



Want  $\frac{\partial f}{\partial x}$  small!

Jacobian:

tells us contracting space  
expanding  
identity



$$h = f(x) \Rightarrow \frac{\partial f}{\partial x} = \underline{\underline{J}} = \text{Jacobian}$$

$$\nabla_x h = \underline{\underline{J}} = \begin{bmatrix} \frac{\partial f(x_1)}{\partial x_1} & \frac{\partial f(x_2)}{\partial x_1} & \dots & \frac{\partial f(x_n)}{\partial x_1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(x_1)}{\partial x_n} & \dots & \dots & \frac{\partial f(x_n)}{\partial x_n} \end{bmatrix} \quad n \times n$$

$$\mathcal{L} = \mathcal{L}_{\text{data}}(x, g(f(x))) + \alpha \cdot \underbrace{\text{Reg}(h)}_{\text{CAA}} \sum_i \|\nabla_x h_i\|^2$$

Intuition:  $f(x + \Delta x) \approx f(x) + \Delta x \underbrace{\left( \frac{\partial f}{\partial x} \right)}_{\underline{J}}$

— A linear approximation to the nonlinear  $f(\cdot)$  encoder.

— For a linear operator, it is contractive if  $\|Jx\| \leq 1 \quad \forall \|x\|=1$

$f \approx \text{RLTS} \rightarrow f$  is  $\approx$  contractive.

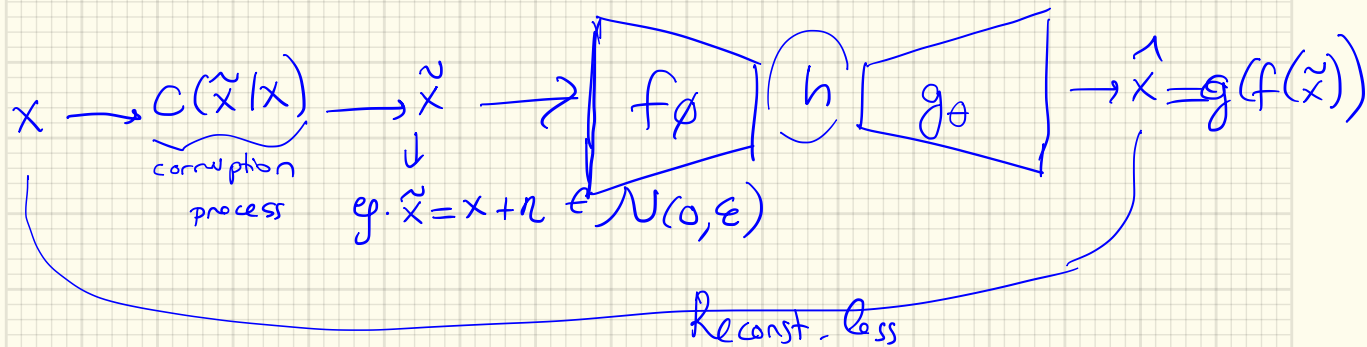
→ CAE introduces the explicit regularizer on the hidden code  $h$ :

$$\underbrace{\mathcal{L}_{\text{reg}}(h)}_{\underline{J}} = \left\| \underbrace{\frac{\partial f(x)}{\partial x}}_{\underline{J}} \right\|_F^2$$

Squared Frobenius-norm of the Jacobian matrix.

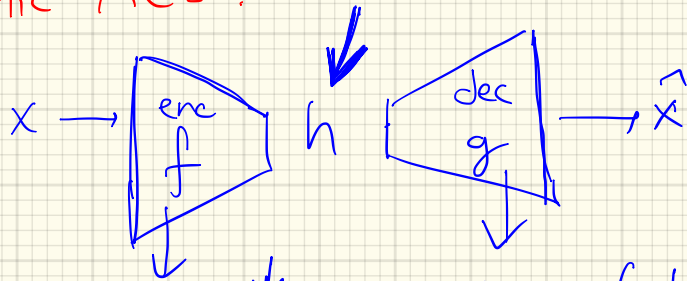
make derivatives of the encoder  $f$  as small as possible.

(iii) **Denoising AEs (DAE)**: Assumes the input data is noisy,  
 Tries to make the AE less sensitive to noise perturbations in the input.



- DAE:
- 1) Samples a training example  $x^{(i)}$  from training data
  - 2) Sample a corrupted version of  $\tilde{x}^{(i)}$  ←
  - 3) Train AE w/  $(x, \tilde{x})$  w/  $\text{rec. Loss}(x, g_\theta(f_\theta(\tilde{x})))$   
 $\underbrace{\hspace{10em}}_{L(x, \hat{x})}$   
 $\uparrow$

# Stochastic AEs:

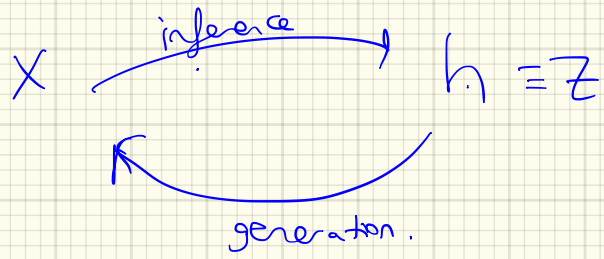


$P_{\text{encoder}}(h|x)$

= posterior distrib./density  
(inference)

$P_{\text{decoder}}(x|h)$   
= generator density

ideally  $\approx$  inverse of  
the encoder  
distribution.



Next time:  
VAEs