Recap:

Last lecture ⇒ tricks to improve performance

① Ensemble ⟶ average over multiple models

② Dropout ⤵ can be viewed as stochastic version of ensemble

⇕

cooperate with Neural Networks

③ Batch Normalization ⟶ Decouple

④ Data Augmentation ⟶ change-invariant prior

⤷ more data, less generalization gap

⑤ Learning Rate Decay Scheme

⇓

{ High learning rate at begining stage

Then decay over training

⑥ Adversarial Training          Defn: Robust

⇓

$\hat{f}(x') = \hat{f}(x)$   for   $\forall \|x'-x\| \ll 1$

small perturbation

DSA5204   Lec 8

Continuous topic : ( one last topic)

⟶ Adversarial Training ⟷ adversarial example

① Definition (adversarial example)          ⟶ the worse example

given parameter $\hat{\theta}$ , adversarial example is $x' = \underset{z:\|z-x\|\leq\delta}{\arg\max}\ L(\hat{y}(z,\hat{\theta}), y)$

② Definition (adversarial training)

How can we reduce the effect of "Adversarial example" ?

can be achieved by gradient ascent

1-sample   $\underset{\theta}{\min}\ \underset{z:\|x-z\|\leq\delta}{\max}\ L(\hat{y}(z;\theta), y)$

$$\boxed{\text{multi-sample}} \quad \min_{\theta} \quad \frac{1}{N} \sum_{i=1}^{N} \max_{z_i : \|x_i - z_i\| \leq \delta} L(\hat{y}(z_i ; \theta), y)$$

③ **Algorithm** ( Fast Gradient Sign Method [FGSM] )

$\rightarrow$ For $k = 0, 1, \ldots$

$z_0 = x$

find the "Adversarial Example"

gradient ascent in input space $\begin{cases} \text{For } \hat{j} = 0, 1, \ldots, J-1 \\ \qquad z_{j+1} = z_j + \varepsilon_2 \, \text{sign} \left( \nabla_z L(\hat{y}(z_j ; \theta_k), y) \right) \end{cases}$

$$\theta_{k+1} = \theta_k - \varepsilon_1 \nabla_\theta L(\hat{y}(z_J ; \theta_k), y)$$

in argmax

Here, $\delta = J \times \varepsilon_2$ ( we want $\|z - x\| \leq \delta$ )

**Rmk:** Generally speaking, we do not need to pre-train model first, then find adversarial examples and store them, lastly re-train model.

Instead, we modify the loss function and train model directly.

$$\min_{\theta} \max_{z : \|z - x\| \leq \delta} L(\hat{y}(z ; \theta), y)$$

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{N}$$

Something like Data Augmentation

④ **1-D Example**

Model: $\hat{y}(x ; \theta) = \theta x$

Data $\{x = 1, y = 0\} = \mathcal{D}$

Loss function: $L(\hat{y}, y) = \begin{cases} \frac{1}{2}(\hat{y} - y - \text{sign}(\hat{y} - y))^2 & \text{o/w} \\ 0 & |\hat{y} - y| \leq 1 \end{cases}$

Then $L(\hat{y}, y) = \begin{cases} \frac{1}{2}(\theta - \text{sign}(\theta))^2 & |\theta| > 1 \\ 0 & |\theta| \leq 1 \end{cases}$

$\Downarrow$

Trivial Loss, via GD, $\boxed{\hat{\theta}_{GD}^{\infty} = 1}$ if we

start from $\theta^{(0)} > 1$

$\longrightarrow$ Adversarial Loss $L_{adv}(\theta) = \max_{z : \|1-z\| \leq \delta} L(\hat{y}(z, \theta), y)$

$= \max_{z : \|1-z\| \leq \delta} \frac{1}{2}(\theta z - \text{sign}(\theta z))^2 \mathbb{1}\{|\theta z| \geq 1\}$

$= \frac{1}{2}(\theta(1+\delta) - 1)^2 \mathbb{1}\{|\theta(1+\delta)| \geq 1\}$

$\boxed{y = \theta x}$ Model

$y = 1$ and $x = 0$

$\Rightarrow$ we want $\theta \approx 0$

$\Rightarrow \hat{\theta}_{adv}^{\infty} = \frac{1}{1+\delta}$ is more close to $0$

and $L_{adv}(\hat{\theta}_{GD}^{\infty}) = \frac{\delta^2}{2} > 0$

---

Today's topic : Unsupervised Learning + Semi-supervised Learning

GT: $f^*$

1. Task

a) Supervised Learning $\longrightarrow$ input : $x$     output : $y$

b) **Unsupervised Learning** $\longrightarrow$ no output !

learn some task-agnostic patterns

① Dimensionality Reduction

② Generative Model

③ Clustering

④ Density Estimation

$z \longrightarrow \boxed{f(z)}$

c) **Semi-supervised Learning**

## 2. PCA

① Design Matrix $X \in \mathbb{R}^{N \times d}$

② Covariance Matrix $S = \frac{1}{N} X^T X$

③ Eigenvalue Decomp. $S = U \Sigma U^T$

$\Updownarrow$

$SU = U \Sigma$ $\boxed{U = (u_1, \dots, u_d)}$

m **Principle component direction** $u_i \in \mathbb{R}^d$

$U_m = (u_1, \dots, u_m)$ $m < d$

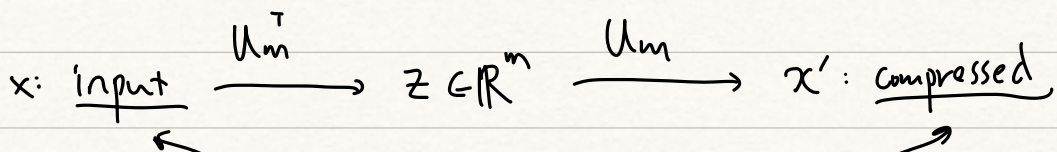$\boxed{Z_m = X U_m \in \mathbb{R}^{N \times m}} \rightarrow$ **Principle Scoring**

Reconstruction : $X \approx X U_m U_m^T = Z_m \cdot U_m^T \in \underline{\mathbb{R}^{N \times d}}$

$$= Z_m \begin{pmatrix} u_1^T \\ \vdots \\ u_m^T \end{pmatrix}$$

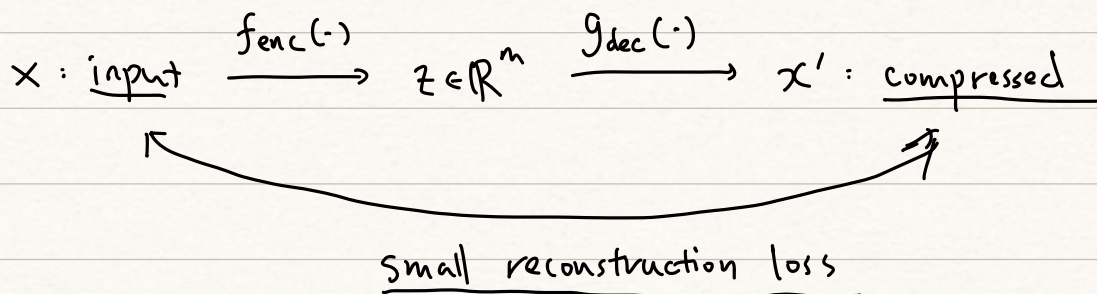explained variance : $\boxed{\sum_{i=1}^{m} \lambda_i}$

# 3. Auto-Encoder (AE)
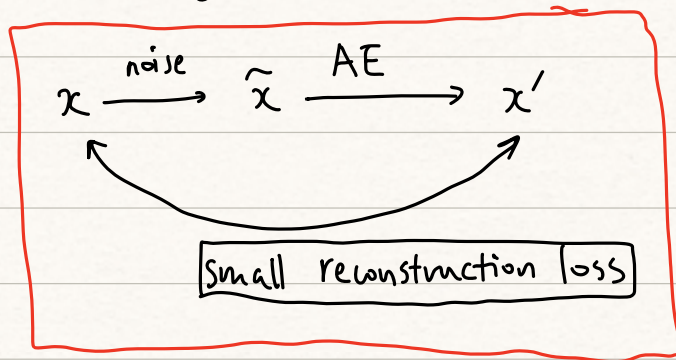
Re-cap : [ PCA framework ]

$$x: \underline{\text{input}} \xrightarrow{U_m^T} z \in \mathbb{R}^m \xrightarrow{U_m} x': \underline{\text{compressed}}$$

$\Rightarrow$ **Generalization** (AE): $\underline{\text{Small reconstruction loss}}$

$$x: \underline{\text{input}} \xrightarrow{f_{enc}(\cdot)} z \in \mathbb{R}^m \xrightarrow{g_{dec}(\cdot)} x': \underline{\text{compressed}}$$

$\underline{\text{Small reconstruction loss}}$

## "Denoising AE"



$$x \xrightarrow{\text{noise}} \tilde{x} \xrightarrow{AE} x'$$

[Small reconstruction loss]

# 4. Semi-supervised Learning

a) $\Big\{$ transductive

inductive $\qquad$ $\underline{\text{task}}$

b) Method to $\underline{\text{label}}$ $\underline{\text{unlabelled data}}$

① Naive approach ( $\underline{\text{Self-learning}}$ )

train model $\rightarrow$ label unlabelled $\rightarrow$ re-train

② Label propagation

Rmk: different from some clustering algorithms like
K-means etc.