

# SM2001 Study Problems

1. <sup>(1)</sup> Def full SVD; <sup>(2)</sup> Expression of economic SVD, <sup>(3)</sup> labeling size, name

Given a Dataset  $\tilde{X} \in \mathbb{R}^{n \times m}$

$$\tilde{X} = [\tilde{x}_1 \ \tilde{x}_2 \ \dots \ \tilde{x}_m] \quad \text{snapshots, } \tilde{x}_i \in \mathbb{R}^{n \times 1}$$

- Definition of SVD: A unique matrix decomposition exists for every complex-valued Matrix:  $\tilde{X} \in \mathbb{C}^{n \times m}$

- Full SVD:

$$U \Sigma V^* = \tilde{X}$$

- $U \in \mathbb{C}^{n \times n}$ ,  $V \in \mathbb{C}^{m \times m}$ ,  $\begin{cases} U \cdot U^* = I \\ V \cdot V^* = I \end{cases}$   
unitary matrix with orthonormal columns
- $\Sigma \in \mathbb{R}^{n \times m}$   
columns of  $U$ : left singular vectors,  $\Sigma$ : singular vectors

Real, non-negative Entries on diagonal and zeros off diag

Ordered from Largest  $\rightarrow$  Smallest  
Diagonal Value: Singular Value

- Economic SVD: when  $n \gg m$ ,  $\Sigma$  has most  $m$  non-zero elements on diagonal  $\hat{\Sigma}$

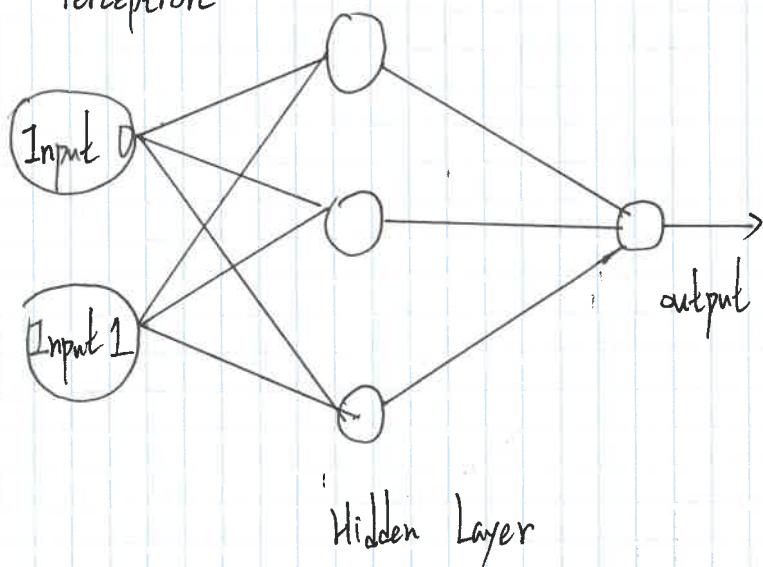
$$\text{Therefore: } \tilde{X} = U \hat{\Sigma} V^* = [\hat{U} \ \hat{0}^\perp] \begin{bmatrix} \hat{\Sigma} \\ 0 \end{bmatrix} V^*$$

$$\tilde{X} = \hat{U} \cdot \hat{\Sigma} \cdot \hat{V}^* \quad \hat{U} \text{ and } \hat{V}^\perp \text{ are complementary}$$

2. Schematic Representation of MLP, LSTM  
Show Inputs & Outputs, strengths & weakness.

### 2-1 MLP: Multi-Layer Perceptron

• Schematic :



• Strengths. 1. Good performance of linear regression / classification  
Gives a good prediction of point data (no time information)

2. Easier to interpret and understand Algorithm

• Weakness:

1. Can not treat data as sequence. Can not predict data in time series

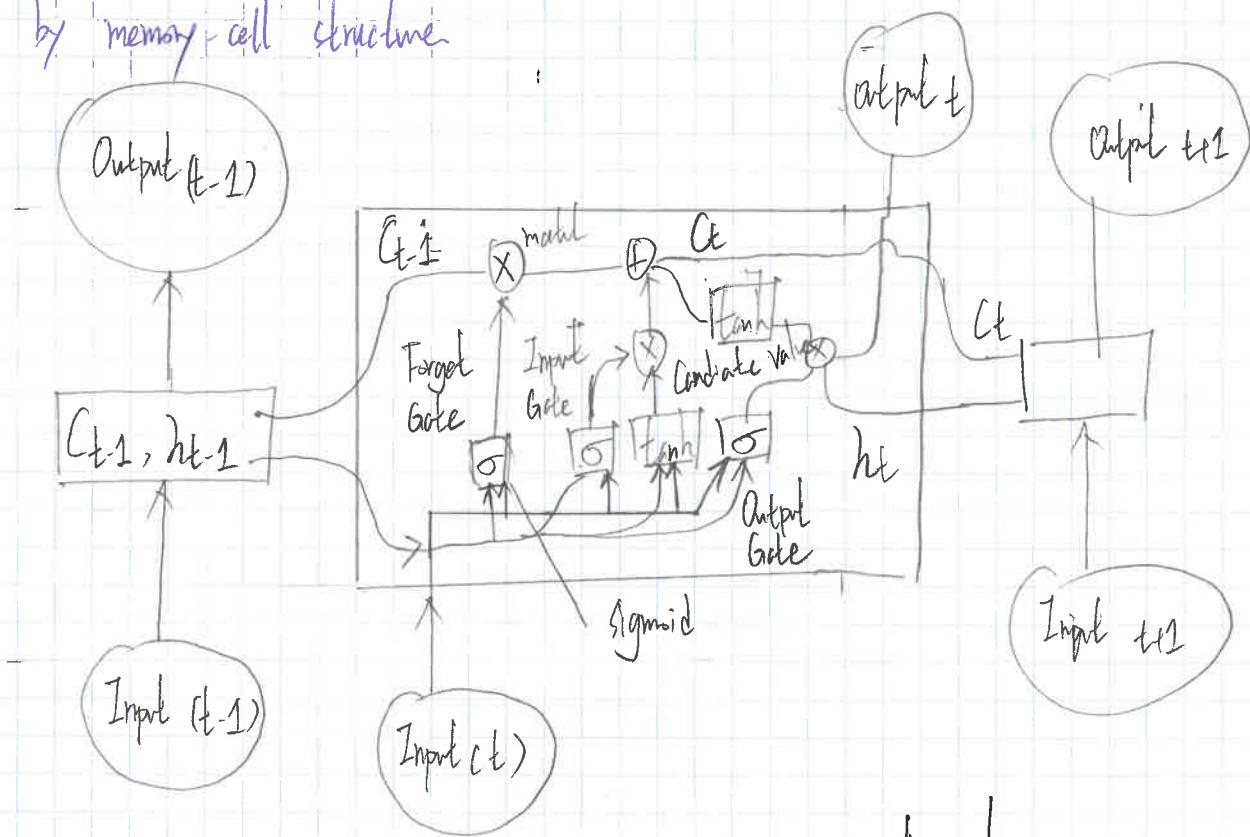
2. Not able to get spatial information between data point, e.g. Image Data, thus it's difficult to use MLP deal with high-dimensional data.

## 2.2 LSTM: Long short term Memory

LSTM is a kind of Recurrent Neural Network (RNN)

RNN ~~can~~ can deal with sequence by iterating through sequence elements and maintaining a state that contains information about what it has "seen" so far. Map a sequence of data points to another sequence

Compared with simple RNN LSTM solved Gradient Vanish Problem  
by memory cell structure



**Strength:**

- Treat data as sequential, thus avoiding temporal information loss
- it will be useful when predicting time-series data.
- Also, since it treats data as sequential, bidirectional LSTM
- performs well in NLP

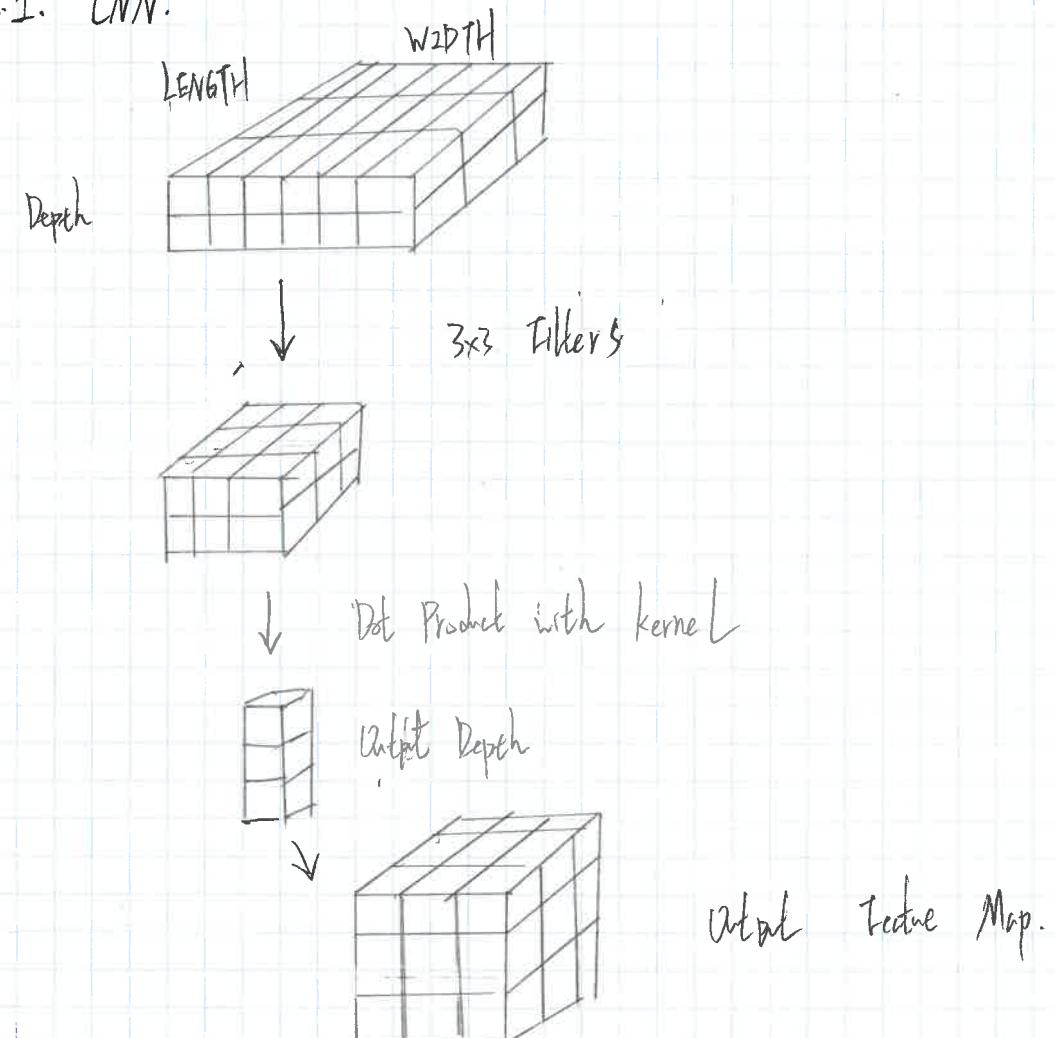
**Weakness:**

1. Easier to Overfitting. Need dropout which built-in in tensorflow)
2. At tensorflow, LSTM compute via CPU instead of GPU.

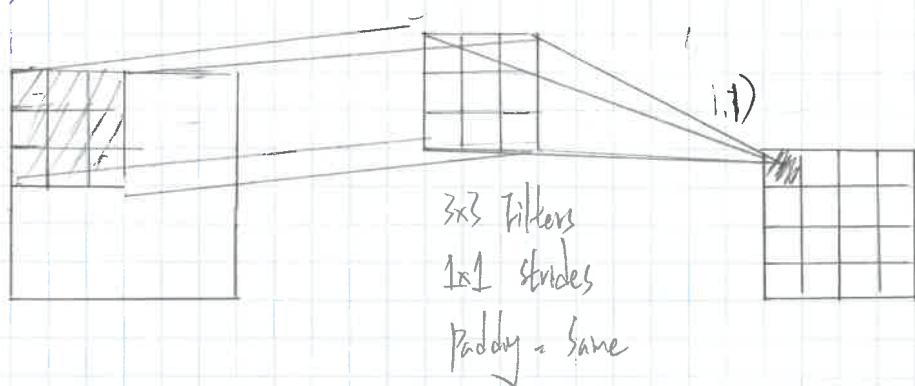


### 3. Schematic a LSTM, a CNN, anAE, list Application

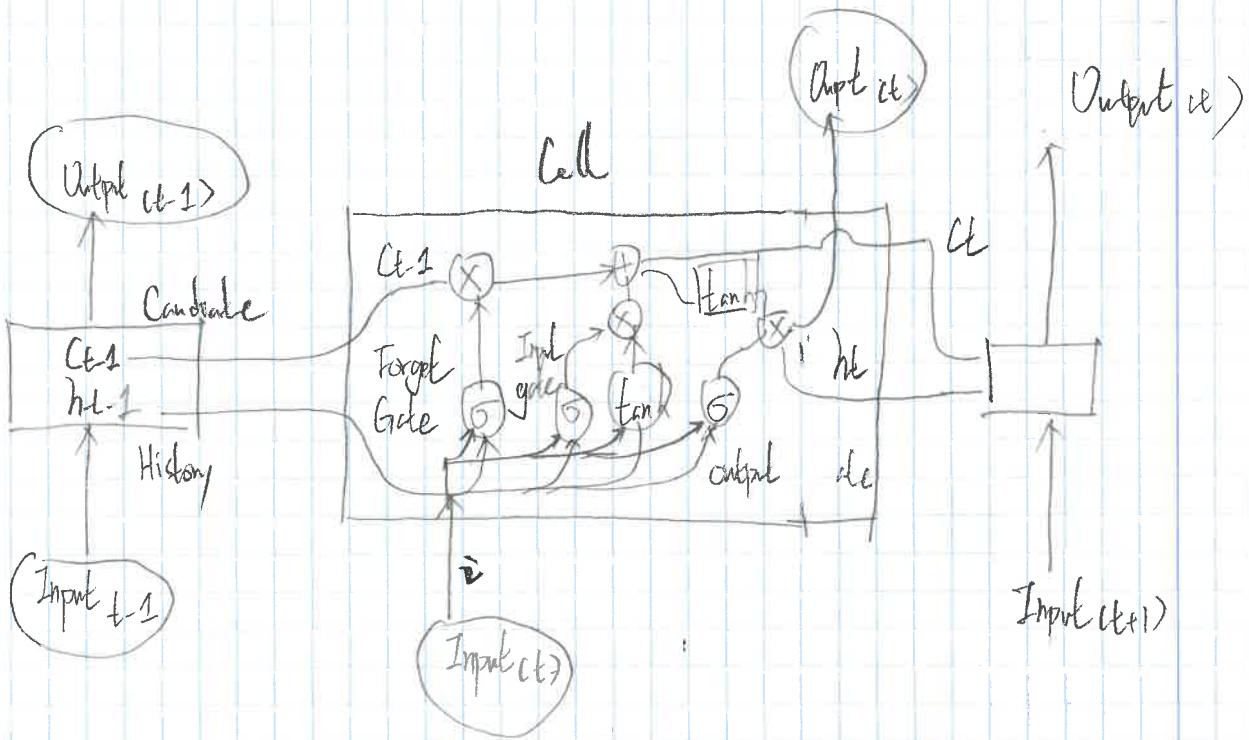
#### 3.1. CNN.



In 2D:

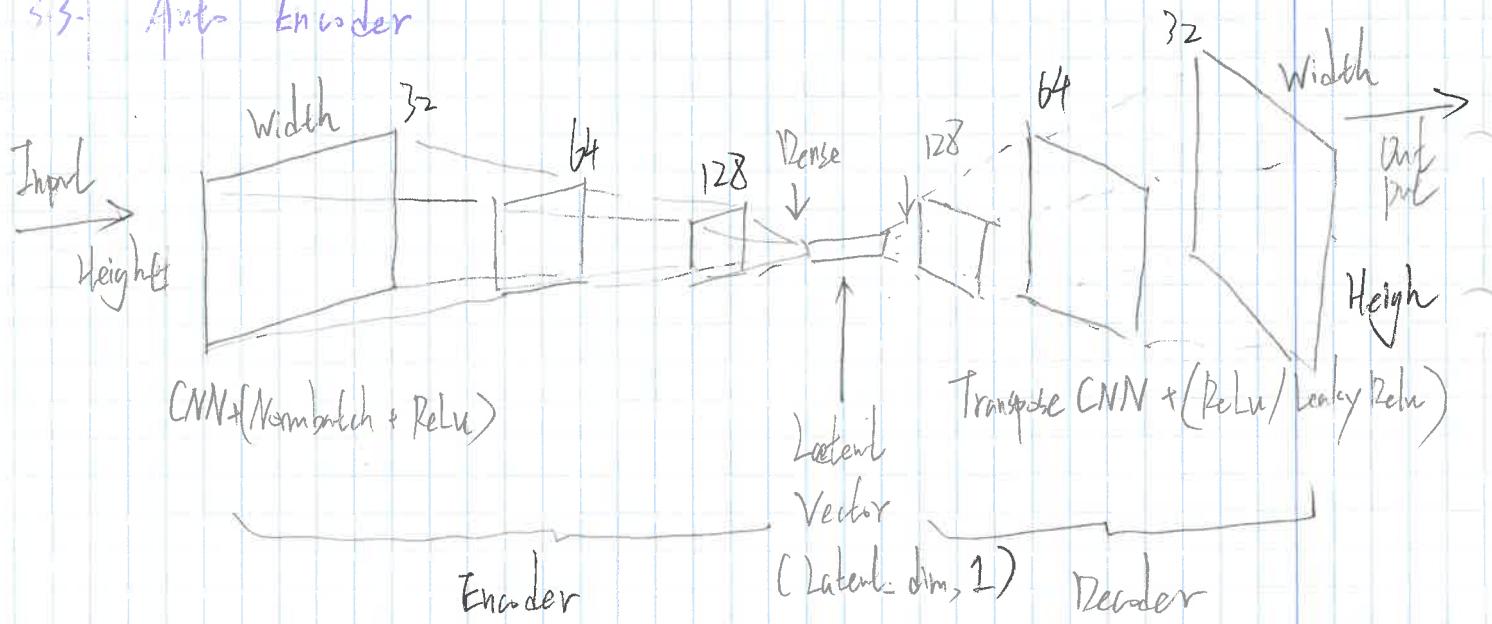


### 3.2. LSTM



Time series prediction    Weather Prediction

### 3.3. Auto Encoder

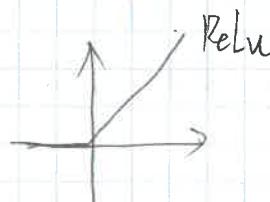
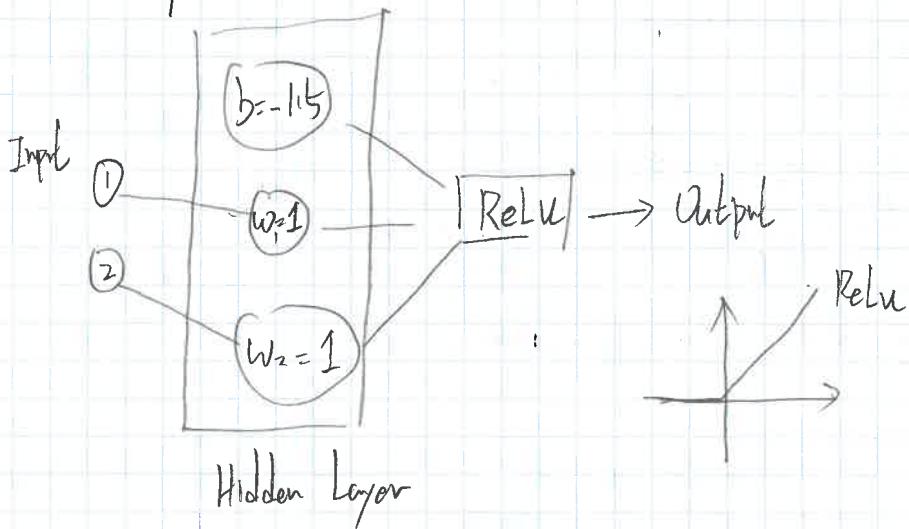


Mode decomposition and reconstruction

4. A neuron, 2 inputs, 1 <sup>output</sup><sub>Input</sub>, Activation = "ReLU"

weights  $w_1=1$ ,  $w_2=1$ , Bias  $b=-1.5$

If Input =  $(0,0)$ ,  $(1,0)$ ,  $(0,1)$ , ~~or~~  $(1,1)$   
Output?



- (a) Input  $(0,0)$        $\text{ReLU}(0 \cdot 1 + 0 \cdot 1 - 1.5) = 0$
- (b)  $(1,0)$        $\text{ReLU}(1 \cdot 1 + 0 \cdot 1 - 1.5) = 0$
- (c)  $(0,1)$        $\text{ReLU}(0 \cdot 1 + 1 \cdot 1 - 1.5) = 0$
- d       $(1,1)$        $\text{ReLU}(1 + 1 - 1.5) = 0.5$

5. Consider loss  $E(\underline{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x, \underline{w}) - t_n\}^2$

$y$ : prediction,  $\underline{w}$ : weights,  $t_n$ : targets

For Prediction  $y(x, \underline{w}) = w_0 + w_1 x + w_2 x^2 \dots = \sum_{j=0}^M w_j x^j$

Show that:  $\underline{w} = \{w_i\}$  minimize loss function  $E(\underline{w})$  are given by the solutions to the following linear functions

$$\boxed{\sum_{j=0}^M A_{ij} w_j = T_i} \quad (1)$$

where  $A_{ij} = \sum_{n=1}^N x_n^{ij}$ ,  $T_i = \sum_{n=1}^N x_n^i \cdot t_n$  (2)

$i, j$  denote index of component

At iteration  $i$  snapshot of  $x^i$  is  $x_n^i = [x_1^{i,0}, x_1^{i,1}, \dots, x_1^{i,M}]$

$$x^i = \begin{bmatrix} -x_1^{i,0} \\ -x_1^{i,1} \\ \vdots \\ -x_1^{i,M} \end{bmatrix} \quad N \quad x^i \in \mathbb{R}^{N \times M}$$

For

thus in (1), LHS =  $\sum_{j=0}^M \sum_{n=1}^N x_n^{ij} \cdot w_j$

For  $A_{ij}$ , it means at certain iter, a certain  $x_1^i$  summed through  $N$

$$x^i = \begin{bmatrix} x_1^{i,0} & x_1^{i,1} & \dots & x_1^{i,M} \\ x_2^{i,0} & x_2^{i,1} & \dots & x_2^{i,M} \\ \vdots & \vdots & \ddots & \vdots \\ x_N^{i,0} & x_N^{i,1} & \dots & x_N^{i,M} \end{bmatrix} \quad N$$

$$\underline{w}^i = \begin{bmatrix} w_0^i \\ w_1^i \\ \vdots \\ w_M^i \end{bmatrix} \quad M$$

$$\text{LHS} = [A_{i,0} \dots A_{i,M}] \cdot \begin{bmatrix} w_0^i \\ w_1^i \\ \vdots \\ w_M^i \end{bmatrix} \quad (7)$$

Solution for #5:

Substitute (2) into (1)

$$\sum_{j=0}^M \sum_{n=1}^N x_n^{i+j} w_j = \sum_{n=1}^N x_n^i t_n \quad (3)$$

$$\text{For (3). LHS} = \sum_{j=0}^M (x_1^{i+j} + x_2^{i+j} + \dots + x_N^{i+j}) w_j$$

$$= x_1^i \underbrace{\sum_{j=0}^M x_1^j w_j}_N + x_2^i \underbrace{\sum_{j=0}^M x_2^j w_j}_N + \dots + x_N^i \underbrace{\sum_{j=0}^M x_N^j w_j}_N$$

$$\Rightarrow \left\{ \sum_{j=0}^M x_n^j w_j = y_n \right\}$$

$$= x_1^i y_1 + \dots + x_N^i y_N = \sum_{n=1}^N x_n^i y_n$$

$$\text{LHS} = \text{RHS}$$

$$\sum_{n=1}^N x_n^i y_n = \sum_{n=1}^N x_n^i t_n$$

$$\sum_{n=1}^N y_n = \sum_{n=1}^N t_n \quad (4)$$

(4) It shows that for optimal  $\underline{w}$ , prediction is identical to target @ n

In another word, the  $\underline{w}$  minimize the loss function are given by

Solution (1).

b. Show that derivative of loss function

$$E(w) = - \sum_{n=1}^N \{ t_n \ln(y_n) + (1-t_n) \ln(1-y_n) \}$$

with respect to activation value  $a_k$  (which is argument of activation function)

for an output unit having a logistic sigmoid activation function activation function satisfies

$$\frac{\partial E}{\partial a_k} = y_k - t_k$$

Note that sigmoid is  $\sigma(a) = \{1 + \exp(-a)\}^{-1}$

$$\frac{\partial E}{\partial a_k} = - \left[ \frac{t_k}{y_k} \times \frac{\partial \sigma(a_k)}{\partial a_k} - (1-t_k) \times \frac{1}{1-y_k} \times \frac{\partial \sigma(a_k)}{\partial a_k} \right] \quad (1)$$

For sigmoid:

$$\begin{aligned} \frac{\partial \sigma(a)}{\partial a} &= -1 [1 + \exp(-a)]^{-2} \times [-\exp(-a)] \\ &= \frac{e^{-a}}{(1+e^{-a})^2} = \frac{e^{-a}}{1+e^{-a}} \times \frac{1}{1+e^{-a}} = \frac{1}{1+e^{-a}} (1 - \frac{1}{1+e^{-a}}) \end{aligned}$$

$$\boxed{\frac{\partial \sigma(a)}{\partial a} = \sigma(a) \cdot [1 - \sigma(a)]}$$

Substitute into (1)

$$\frac{\partial E}{\partial a_k} = - \left[ \frac{t_k}{\sigma(a_k)} \cdot \sigma(a_k) \cdot [1 - \sigma(a_k)] - \frac{1-t_k}{1-\sigma(a_k)} \cdot \sigma(a_k) \cdot [1 - \sigma(a_k)] \right]$$

$$\frac{\partial E}{\partial a_k} = -1 \times [t_k - t_k \cdot y_k - (1-t_k) \cdot y_k] \quad \checkmark$$

$$\boxed{\frac{\partial E}{\partial a_k} = y_k - t_k} \quad \checkmark$$

7. Show that sigmoid  $\sigma(a)$ ,  $\frac{d\sigma}{da}$  can be expressed as function of  $\sigma$ ,  $\sigma(a) = \{1 + \exp(a)\}^{-1}$

$$\begin{aligned}\frac{d\sigma}{da} &= -1 \cdot \{1 + \exp(-a)\}^{-2} \times -1 \times \exp(-a) \\ &= \frac{e^{-a}}{1 + e^{-a}} \times \frac{1}{1 + e^{-a}} \\ &= \left(1 - \frac{1}{1 + e^{-a}}\right) \left(\frac{1}{1 + e^{-a}}\right) \\ \text{where } \frac{1}{1 + e^{-a}} &= \sigma(a)\end{aligned}$$

thus  $\frac{d\sigma}{da} = (1 - \sigma)\sigma$  ✓

8. Show that for activation function  $\tanh(a) = g(a)$ .

$\frac{dg}{da}$  can be a function of  $g$ .

$$\tanh(a) = \left\{ \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)} \right\}$$

$$\begin{aligned} \frac{dy}{da} &= \frac{[\exp(a) - \exp(-a)]^2 [\exp(a) + \exp(-a)] - [-] [-]^2}{[\exp(a) + \exp(-a)]^2} \\ &= \frac{(e^a + e^{-a})^2 - (e^a - e^{-a})^2}{[\exp(a) + \exp(-a)]^2} \\ &= 1 - \frac{[e^a - e^{-a}]^2}{[e^a + e^{-a}]^2} \end{aligned}$$

$$\boxed{\frac{dy}{da} = 1 - g^2} \quad \checkmark$$

9. A matrix  $\underline{\underline{X}}$ , following eigenvector equation

$$\underline{\underline{X}} \cdot \underline{\underline{v}_i} = \lambda_i \underline{\underline{v}_i}$$

$\lambda_i$ : Eigenvalue  
 $\underline{\underline{v}_i}$ : Eigenvectors

Eigenvectors : Orthonormal set

- therefore, an arbitrary vector  $\underline{\underline{v}}$  we have:

$$\underline{\underline{v}}^T \underline{\underline{X}} \cdot \underline{\underline{v}} = \sum_i c_i^2 \lambda_i$$

- $c_i$  are coefficients of  $\underline{\underline{v}}$  in basis of eigenvectors
- By setting  $\underline{\underline{v}}$  equal to each  $\underline{\underline{v}_i}$ , show that  $\underline{\underline{X}}$  is positive definite, if and only if all eigenvalues are positive  
(ie. if & only if  $\underline{\underline{v}}^T \underline{\underline{X}} \cdot \underline{\underline{v}} > 0$  for all  $\underline{\underline{v}}$ )

Since  $\underline{\underline{v}_i}$  are orthonormal, which means they are orthogonal & all of unit length mutually.

thus  $\underline{\underline{v}_i}^T \underline{\underline{X}} \cdot \underline{\underline{v}_i} = \lambda_i$  for all  $\underline{\underline{v}_i}$  &  $\lambda_i$

If  $\underline{\underline{X}}$  is positive-definite, if and only if all  $\lambda_i > 0$

$$\underline{\underline{v}_i}^T \underline{\underline{X}} \cdot \underline{\underline{v}_i} \text{ since } \underline{\underline{X}} \cdot \underline{\underline{v}_i} = \lambda_i \underline{\underline{v}_i}$$

$\Rightarrow \underline{\underline{v}_i}^T \cdot \lambda_i \underline{\underline{v}_i}$  the eigenvectors are orthonormal,

$$\underline{\underline{v}_i}^T \cdot \underline{\underline{v}_i} = 1$$

thus  $\underline{\underline{v}_i}^T \underline{\underline{X}} \cdot \underline{\underline{v}_i} = \lambda_i$

If  $\underline{\underline{X}}$  is positive definite, if and only if all eigenvectors are positive

10. Consider a quadratic loss function defined by

$$E(\underline{w}) = E(\underline{w}^*) + \frac{1}{2} (\underline{w} - \underline{w}^*)^T \underline{H} (\underline{w} - \underline{w}^*)$$

where  $\underline{H}$  is an known operator evaluated at the set of weights  $\underline{w}^*$   
 if  $\underline{H}$  has a eigenvalue equation  $\underline{H} \underline{v}_i = \lambda_i \underline{v}_i$

Show that the contour of constant error  $C$  are ellipses.  
 the axes are aligned with eigenvectors  $\underline{v}_i$

(Hint: you can start by expressing  $(\underline{w} - \underline{w}^*)$  as a linear combination  
 of eigenvectors)

Express the linear combination of  $(\underline{w} - \underline{w}^*)$  with eigenvectors  $\underline{v}_i$

$$(\underline{w} - \underline{w}^*) = \sum_{i=0}^N k_i \underline{v}_i \quad [ \underline{w}_0 - \underline{w}^*_0 = \underline{v}_{k_0} - \underline{v}_{k_0} \cdot \underline{v}_{k_0} ]$$

Thus the  $\frac{1}{2} (\underline{w} - \underline{w}^*)^T \underline{H} (\underline{w} - \underline{w}^*)$  can expressed as

$$C = \frac{1}{2} \sum_{i=0}^N k_i \underline{v}_i \cdot \sum_{j=0}^N \lambda_j k_j \underline{v}_j \quad \text{since } \underline{H} \underline{v}_i = \lambda_i \underline{v}_i$$

thus  ~~$C = \frac{1}{2} \sum_{i=0}^N k_i^2$~~

Since  $\underline{v}_i$  and  $\underline{v}_j$  are orthogonal, and all the  
 eigenvectors are normalized.

thus  $\underline{v}_i \cdot \underline{v}_j = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{else } i \neq j \end{cases}$

$$C = \frac{1}{2} \sum_{i=0}^N k_i^2$$

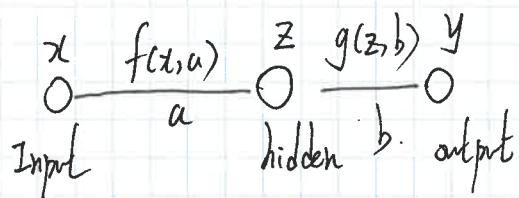
which is ellipses contain in hyperplane, and aligned with  $\underline{v}_i$

12. Consider One-node, one-hidden-layer network

Assume linear activation function  $f(\xi, \alpha) = g(\xi, \alpha) = \alpha\xi$

For loss function  $E = \frac{1}{2} (y_0 - y)^2$

Calculate expressions for  $\frac{\partial E}{\partial \alpha}$  and  $\frac{\partial E}{\partial b}$



$y_0$ : target

$a, b$ : weight

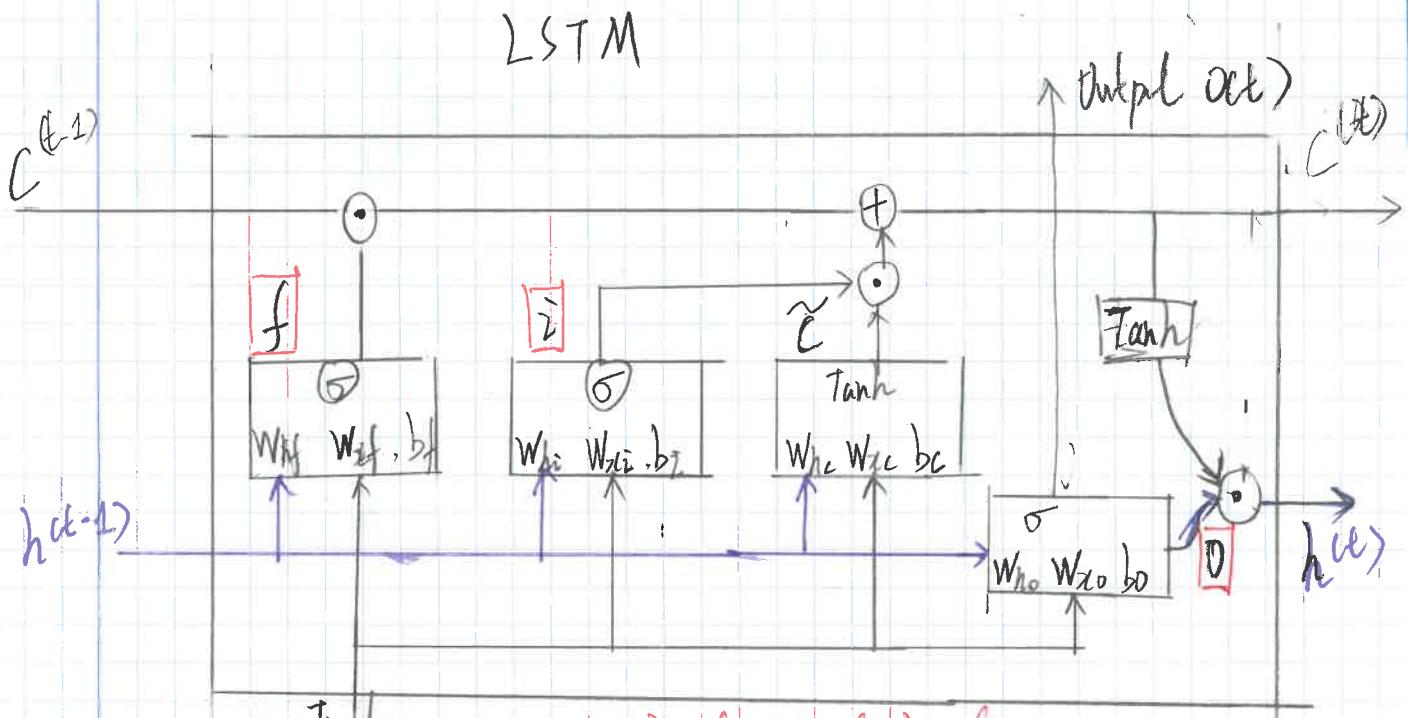
$x$ : input,  $z$ : hidden layer output,  $y$  output

$$\begin{aligned}\frac{dE}{d\alpha} &= \frac{dE}{dy} \frac{dy}{dz} \frac{dz}{d\alpha} = \frac{dE}{dy} \cdot \frac{dy(z, b)}{dz} \cdot \frac{d(f(x, a))}{d\alpha} \\ &= -(y_0 - y) \cdot b \cdot x = (y - y_0) b \cdot x\end{aligned}$$

$$\frac{dE}{db} = \frac{dE}{dy} \frac{dy}{db} = \frac{dE}{dy} \cdot \frac{dg(z, b)}{db}$$

$$\begin{aligned}&= -(y_0 - y) \cdot z \\ &= (y - y_0) \cdot a \cdot x\end{aligned}$$

12. Describe different gates in LSTM, write algorithm for output



Input  
 $\underline{x}^{(t)}$

$$\star \left\{ \begin{array}{l} C^{(t)} = (\underline{C}^{(t-1)} \odot f_t) \oplus (\underline{i_t} \odot \tilde{c}_t) \\ h^{(t)} = \text{tanh}(C^{(t)}) \odot o_t \end{array} \right.$$

Operator

$\odot$ : Element-wise Product.

$\oplus$ : element-wise Summary

Aktiveringsfunktion

$\sigma$ : sigmoid.

[Tanh]: tanh,

Variable State

$C$ : Cell state.

$h$ : history

1.  $f$ : forget gate: Decide which information is allowed to go through reset cell state not going infinitely.

$$\star f_t = \sigma(W_{hf} \cdot h^{(t-1)} + W_{xf} x^{(t)} + b_f)$$

2. Input Gate

$i_t$ : Candidate Value

updating cell state

$$\star i_t = \sigma(W_{hi} h^{(t-1)} + W_{xi} x^{(t)} + b_i)$$

$$\tilde{C} = \tanh(W_{hc} h^{(t-1)} + W_{xc} x^{(t)} + b_c)$$

3.

Output Gate

$$\star o_t = \sigma(W_{ho} h^{(t-1)} + W_{xo} x^{(t)} + b_o)$$

157

13. Consider a CNN where we want Global output  
Do we need any-fully-connected layer?  
Motivate your answer.

- It depends on the problem CNN is going to solve.
- In theory, once we use Global Max Pooling layer or Global Average Pooling layer, we do not need Fully-connected layer anymore, if we just need a global output.  
it is suitable for AutoEncoder or GAN
- However, if suppose we are dealing with image classification  
we need linear output, say with Sigmoid activation function  
thus loss function Binary Crossentropy / Softmax can be implemented.  
In this case, we do need fully connected layer producing a linear output.

14. Consider on CNN with 5 input channels of  $256^2$  grids point each and a kernel size  $(5 \times 5)$ .

If we apply total 64 filters in first layers

How many parameters will each of filters have in that layer?

Parameters in CNN Layer totally will be

$$(\text{Num. filter}) \times [(\text{filter Height}) \times (\text{filter width}) \times \text{Num. channels} + 1] \text{ incl. } \uparrow \\ \text{Bias}$$

For each filter

$$(\text{Height} \times \text{Width} \times \text{Num. channels} + 1) \text{ s}$$

$$= 5 \times 5 \times 5 + 1 = \boxed{126 \text{ paras}}$$

15. In a snapshot Matrix  $\underline{X} \in \mathbb{R}^{n \times m}$   $n > m$ .

If we want to apply method of snapshot to compute SVD  
which correlation matrix is more convenient to use?

- The  $\underline{X}^* \cdot \underline{X}$  is more convenient since  $\underline{X}^* \cdot \underline{X} \in \mathbb{R}^{m \times m}$   
is very low-dimensional matrix.  $m = \text{no. of snapshots}$
- Then by eigenvalue decomposition of  $\underline{X}^* \cdot \underline{X}$   
we can get right singular vectors  $\underline{v}_j$   
&  
singular values  $\sigma_j$
- Then by definition of SVD:

$$\underline{X} = \underline{U} \cdot \Sigma \cdot \underline{V}^*$$

$$\underline{X} \underline{v}_j = \underline{u}_j \cdot \sigma_j, \underline{u}_j \text{ is left singular vector}$$

$$\underline{u}_j = \underline{X} \underline{v}_j / \sigma_j$$

get first  $m$  left singular vector,  
 $m$  corresponds to number of snapshot

alternatives get, maintain  $\underline{U}^*$  &  $\Sigma$  which only has non-zero value

$$\underline{U}^* = \underline{X} \underline{V}^* \Sigma^{-1}$$

1b. Compute Eigenvalues & Eigenvectors of Matrices

$$\underline{A} = \begin{bmatrix} 5 & 7 \\ -2 & -4 \end{bmatrix}, \underline{B} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \underline{C} = \begin{bmatrix} 1 & 2 & 1 \\ 6 & -1 & 0 \\ -1 & 2 & 1 \end{bmatrix}$$

①

$$\underline{A}: \det(\underline{A} - \lambda \cdot \underline{I}) = 0$$

$$\det\left(\begin{bmatrix} 5 & 7 \\ -2 & -4 \end{bmatrix} - \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}\right) = 0$$

$$\begin{bmatrix} (5-\lambda_1) & 7 \\ -2 & -4-\lambda_2 \end{bmatrix} = 0$$

$$(5-\lambda_1)(-4+\lambda_2) + 14 = 0$$

$$-20 - 5\lambda_1 + 4\lambda_2 + \lambda_1^2 + 14 = 0$$

$$\lambda_1^2 + \lambda_1 - 6 = 0 \Rightarrow \begin{cases} \lambda_1 = 3 \\ \lambda_2 = -2 \end{cases}$$

$$\begin{bmatrix} 5 & 7 \\ -2 & -4 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 3 \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\begin{cases} 5v_1 + 7v_2 = 3v_1 \\ -2v_1 - 4v_2 = 3v_2 \end{cases} \begin{cases} -2v_1 + 7v_2 = 0 \\ -7v_2 = 2v_1 \end{cases}$$

$$v_1 = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$$

$$\begin{cases} 5v_1 + 7v_2 = -2v_1 \\ -2v_1 - 4v_2 = -2v_2 \end{cases} \Rightarrow \begin{cases} v_1 = -v_2 \\ v_1 = -v_2 \end{cases}$$

$$v_2 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$\underline{Q} = \begin{bmatrix} -1 & 1 \\ \frac{2}{7} & -1 \end{bmatrix}$$

$$\underline{A} = \begin{bmatrix} 3 & 0 \\ 0 & -2 \end{bmatrix}$$

$$\underline{A} \underline{Q} = \underline{Q} \underline{A} \quad \checkmark$$

17. Compute full SVD & Pseudoinverse of  $\underline{A}$

$$\underline{A} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

① Full SVD

$$\underline{A} = \underline{U} \cdot \underline{\Sigma} \cdot \underline{V}^*$$

$$\underline{A}^* = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 0 \end{bmatrix} \text{ GR } 3 \times 2$$

let  $\underline{X} = \underline{A}^* \cdot \underline{A}$  GR  $3 \times 3$  to find  $\underline{V}$  &  $\underline{\Sigma}$

$$\underline{X} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad \underline{X} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

eigenvalue of  $\underline{X}$  {  $\lambda_1 = 1$   
 $\lambda_2 = 1$   
 $\lambda_3 = 0$  }

eigenvectors of  $\underline{X}$   $\begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

thus  $\underline{V} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$  GR  $3 \times 3$

$$\underline{\Sigma} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \text{ GR } 2 \times 3$$

Then let  $\underline{X}_2 = \underline{A} \cdot \underline{A}^* = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

eigenvectors of  $\underline{X}_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

$$\underline{U} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \text{ GR } 2 \times 2$$

$$\underline{A} = \underline{U} \cdot \underline{\Sigma} \cdot \underline{V}^*$$

checked!

② Pseudo Inverse

$$A^+ = \underline{V} \underline{\Sigma}^{-1} \underline{U}^*$$

where  $\underline{V}$ ,  $\underline{\Sigma}$  &  $\underline{U}$  are truncated

$$\underline{\Sigma}^{-1} = \begin{bmatrix} 4 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

$$U^* = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$A^+ = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 1 \\ -1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\boxed{A^+ = \begin{bmatrix} 0 & 1 \\ -1 & 0 \\ 0 & 0 \end{bmatrix}}$$

check.

$$A^+ A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} = I^{2 \times 2}$$

checked ✓

2o. Mention strengths & weakness of POD & DMD, motivate your answer.

### Proper-Orthogonal Decomposition POD

Pros:

- Fast convergence
- "Optimal Orthogonality"  
 $\Rightarrow$  Good reconstruction with relative low noise

Cons:

- Not always interpretable
- physical meaning
- Temporal & Spatial info. are always mixed together

### Dynamic Mode Decomposition DMD

- Modes have coherent behavior in time
- Modes are orthogonal in time
- Each mode can be interpreted

• Result are less optimal

- Oscillation appears @ certain frequency
- the growth will decay with time

2). Describe STLS algorithm in SINDy. Does this algorithm lead to a parsimonious model? Motivate your answer.

For SINDy, it solve coefficients for ODE equation  $\dot{x} = \theta(x)$

STLS is sequential-threshold least square method.

Once the matrix of ODEs coefficient was obtained through SINDy.

( $\Phi$ )

Assume  $\Phi \in \mathbb{R}^{n \times m}$ , n means the number of terms in library  $\theta(x)$   
 m is the number of ODE we want

thus for  $i = 1$  to  $k$ : self defined iteration

for:  $j = 1 : m$

$\text{small} = \text{where } |\Phi_{i,j}| \leq \lambda$  (threshold, self defined)

$\Phi_{i,\text{small},j} = 0$

$\text{big} = \sim \text{small}$

$\Phi_{i,\text{big},j} = \text{Least\_square\_regression}(\theta(x)|_{j:\text{big}}, \dot{x}|_{j:j})$

end

end

By iteration and setting threshold, it can ensure most of coefficient will be filtered without affect the result.

In this way, the  $\Phi$  matrix is obtained as sparse as possible

which means the ODE will has less coefficients

it lead to a parsimonious model

22. Describe LASSO, elastic-net & Ridge regression regularization  
Write corresponding equations and Norms.

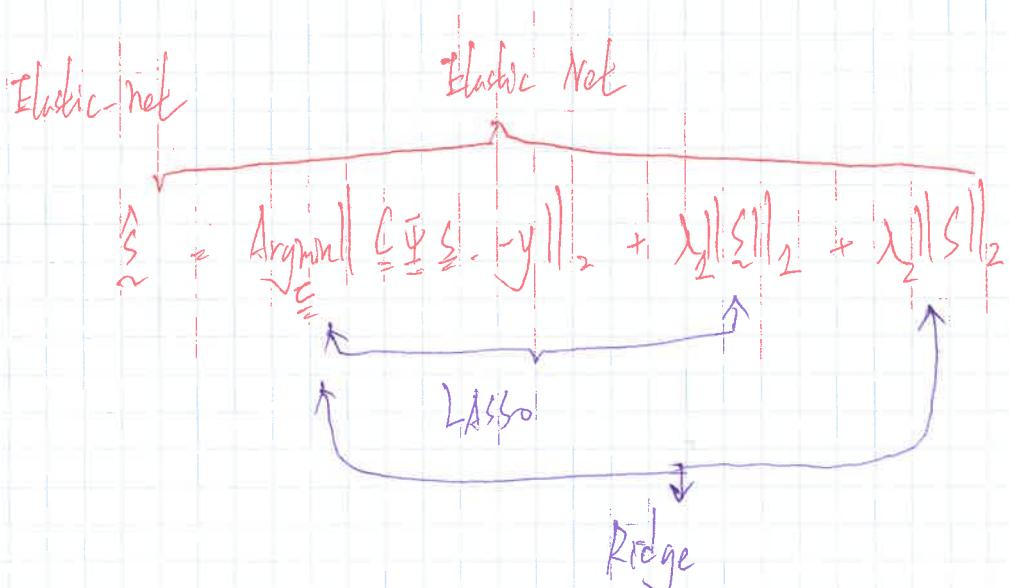
When finding the sparse vector in compress sensing problem, in order to improve the accuracy & reduce noise as much as possible, the  $\ell_2$ -norm regularization may not sufficient. thus we consider adding hyper parameter  $\lambda$  to strengthen the penalty from regularization of  $\underline{\Sigma}$ .  $\ell_1$  norm

then we introduce LASSO Least absolute shrinkage and selection operator

$$\hat{\underline{\Sigma}} = \underset{\underline{\Sigma}}{\operatorname{Argmin}} \| \underline{\Sigma} - \underline{y} \|_2 + \lambda \| \underline{\Sigma} \|_1$$

where  $\lambda$  is the hyperparameter.

It can effectively regularizes the model and avoid overfitting and increase sparsity.



### 23. Describe the DMD algorithm

DMD is able to discover dynamic systems from high-dimensional data.  
 The data can be decomposed into a simple representation based on temporal-spatial features

$$u(x, t) = \sum_{m=1}^{\infty} a_m \cdot v_m(x) \exp[(\lambda_m + i\omega_m) \cdot t]$$

Amplitude ↑      mode ↑      growth ↑      Frequency ↑

Algorithm: We define Matrix  $\underline{X} = \left[ \begin{array}{c|c|c|c} & & & \\ \hline x_0 & \dots & x_{n-1} & \end{array} \right] \in \mathbb{R}^{n \times m}$   
 $\underline{X}^T = \left[ \begin{array}{c|c|c|c} & & & \\ \hline x_1 & \dots & x_n & \end{array} \right] \in \mathbb{R}^{m \times n}$

For exact DMD, a representation of  $\underline{A}$  is

$$\underline{A} \underline{X} = \underline{X}^T \Leftrightarrow \underline{A} = \underline{X}^T \underline{X}^{-1}, \quad A \in \mathbb{R}^{m \times m}$$

① We take truncated SVD on  $\underline{X}$  with  $r \leq m$

$$\underline{X} = \underline{U}_r \cdot \Sigma_r \cdot \underline{V}_r^*$$

According to definition,  $\underline{X}^T = \underline{V}_r \cdot \Sigma_r^{-1} \cdot \underline{U}_r^*$

$$\textcircled{2} \quad \underline{A} = \underline{X}^T \cdot \underline{V}_r \cdot \Sigma_r^{-1} \cdot \underline{U}_r^*, \quad A \in \mathbb{R}^{m \times m}$$

However,  $\underline{A}$  could be truncated as  $r \times r$  via

$$\tilde{\underline{A}} = \underline{V}_r^* \underline{A} \underline{U}_r$$

$$= \underline{V}_r^* \cdot \underline{X}^T \cdot \underline{V}_r \cdot \Sigma_r^{-1} \cdot \underline{V}_r^* \cdot \underline{V}_r, \quad \& \quad \underline{V}_r^* \cdot \underline{V}_r = I_r$$

$$\text{thus } \tilde{\underline{A}} = \underline{V}_r^* \underline{X}^T \cdot \underline{V}_r \cdot \Sigma_r^{-1}$$

(29)

PoD-reduced state

③ Solve eigenvalue decomposition on  $\tilde{A}$

$$\tilde{A} \cdot \underline{W} = \underline{\Lambda} \cdot \underline{W} \quad \left\{ \begin{array}{l} \underline{W} : \text{eigen vectors} \\ \underline{\Lambda} : \text{Eigenvalues} \end{array} \right.$$

④ Compute DMD modes

Modes are orthogonal in temporal space not in spatial  
That's why not use  $\underline{V}_r^*$  in projection

Eigenvectors  $\underline{W}$  have been calculated in different basis  
thus we need to project  $\underline{W}$  onto original basis

DMD modes:  $\underline{\Phi} = \underline{X}^T \underline{V} \underline{\Sigma}^{-1} \underline{W}$ ;  $\underline{b} = \underline{\Phi}^+ \underline{x}_0$

⑤ DMD Expansion

No input:

$$K(x, t) = \underline{\Phi}(x) \exp(-\underline{\Lambda} \cdot t) \cdot \underline{b}$$



$$U(x, t) = \sum_{k=1}^r b_k \phi_k(x) \exp(\lambda_k t)$$

24. Describe snapshot method to compute SVD.

When we are dealing with Data  $\underline{X} \in \mathbb{R}^{n \times m}$   $n > m$

$$\underline{X} = \begin{bmatrix} | & | & | \\ x_1 & x_2 & \dots & x_m \\ | & | & \dots & | \\ n & & & n \end{bmatrix}$$

each  $x_i$  for  $i \in m$  is called a snapshot

We find  $\underline{X}^* \in \mathbb{R}^{m \times n}$

- thus  $\underline{X}^* \cdot \underline{X} \in \mathbb{R}^{m \times m}$ , which is a very low-dimensional matrix

it means for finding right singular vectors & singular value

• cost much less computational power and time

- thus we do eigenvalue decomposition on  $\underline{X}^* \cdot \underline{X}$  first

to obtain  $\underline{V} \in \mathbb{R}^{m \times m}$  &  $\Sigma \in \mathbb{R}^{n \times m}$

Then for each snapshot  $(j)$ , the left singular vectors can be yield by

$$u_j = \underline{X} \cdot v_j / \sigma_j \quad \text{since for each } j,$$

$$\underline{X} = u_j \cdot \sigma_j \cdot v_j^*$$

$$\underline{X} \cdot v_j = v_j \cdot \sigma_j \cdot 1$$

$$v_j = \underline{X} \cdot v_j / \sigma_j$$

25. Describle calculation of PoD in case of non-uniform sampling in space & time. Consider that  $\underline{M}$  &  $\underline{N}$  are mass matrix and matrix containing temporal weights respectively.

The POD (proper-Orthogonal Decomposition) is used for get deterministic & Orthogonal basis function  $\underline{\Phi}$  and their time collections  $\underline{\Sigma}$  to decompose snapshot matrix  $\underline{X} \in \mathbb{R}^{n \times m}$

① According to the definition

$$\underline{X} = \underline{\Phi} \cdot \underline{\Sigma} = \underline{\Phi} \cdot \underline{\Sigma} \cdot \underline{V}^* \quad (*)$$

② We defined that inner product for  $x_j$  snapshot

$$\langle x_j, x_j \rangle_{\underline{M}} = \underline{x}_j^T \cdot \underline{M} \cdot \underline{x}_j \quad \text{where } M \text{ is matrix of mass} \\ \text{for non-uniform meshes, also valid for matrix } \underline{N} \stackrel{\text{diagonal}}{\in} \mathbb{R}^{n \times m}$$

③ In (\*)  $\|\underline{\Phi}\|_M = \underline{I}^{n \times n}$

$$\|\underline{V}\|_N = \underline{I}^{m \times m}$$

④ Denote  $M^{\frac{1}{2}}$ ,  $N^{\frac{1}{2}}$  as normalized  $\underline{M}$  &  $\underline{N}$  in spatial & temporal

$\underline{\Phi}^*$  &  $\underline{V}^*$  as normalized  $\underline{\Phi}$ ,  $\underline{V}$

$$\underline{\Phi}^* \cdot \underline{\Sigma} \cdot \underline{V}^* = \underline{I}^{n \times m}$$

$$\underline{\Phi}^* \cdot \underline{\Sigma} \cdot \underline{V}^* = \underline{I}^{m \times m}$$

then (\*) is

$$SVD \left( \underline{M}^{\frac{1}{2}} \cdot \underline{X} \cdot \underline{N}^{\frac{1}{2}} \right) = \underline{\Phi}^* \cdot \underline{\Sigma} \cdot \underline{V}^*$$

2b. Derive the Galerkin Projection after apply POD on a partial differential equation.

For a dynamic equation  $\dot{\underline{x}} = f(\underline{x}) \quad \underline{x} \in \mathbb{R}^n$  (1)

Suppose we have original data of  $\underline{x}$ .

$\underline{X} \in \mathbb{R}^{n \times m}$  contain  $m$  snapshots.

① First we apply economic / truncated SVD on  $\underline{X}$

$$\underline{X} = \underline{U}_r \Sigma_r \underline{V}_r^* \quad \text{which } r < n$$

② then we can obtain low-dimensional approximation of  $\underline{X}$

$$\underline{U}_r^* \underline{X} = \underline{Z} \Leftrightarrow \underline{X} = \underline{U}_r \cdot \underline{Z}$$

Substitute (2) into (1)

$$\frac{d \underline{X}}{dt} = f(\underline{X})$$

$$\frac{\partial}{\partial t} \frac{d}{dt} (\underline{U}_r \underline{Z}) = f(\underline{U}_r \underline{Z})$$

since  $\underline{U}_r$  is the spatial modes, thus it is not related to time.

$$\underline{U}_r \frac{d}{dt} (\underline{Z}) = f(\underline{U}_r \underline{Z})$$

$$\underline{U}_r^* \underline{U}_r \frac{d}{dt} (\underline{Z}) = \underline{U}_r^* f(\underline{U}_r \underline{Z})$$

$$\dot{\underline{Z}} = \underline{U}_r^* f(\underline{U}_r \underline{Z}) = g(\underline{Z}) \quad (3)$$

(3) is called Galerkin Projection onto a POD subspace.

26. Derive the Galerkin Projection after apply POD on a partial differential equation.

For a dynamic equation  $\dot{\underline{x}} = f(\underline{x}) \quad \underline{x} \in \mathbb{R}^n$  (1)

Suppose we have original data of  $\underline{x}$ .

$\underline{X} \in \mathbb{R}^{n \times m}$  contain  $m$  snapshots.

① First we apply economic / truncated SVD on  $\underline{X}$

$$\underline{X} = \underline{U}_r \underline{\Sigma}_r \underline{V}_r^* \quad \text{which } r < n$$

② then we can obtain low-dimensional approximation of  $\underline{X}$

$$\underline{U}_r^* \underline{X} = \underline{Z} \Leftrightarrow \underline{X} = \underline{U}_r \cdot \underline{Z}$$

Substitute (2) into (1)

$$\frac{d \underline{X}}{dt} = f(\underline{X})$$

$$\frac{\partial}{\partial t} \frac{d}{dt} (\underline{U}_r \underline{Z}) = f(\underline{U}_r \underline{Z})$$

since  $\underline{U}_r$  is the spatial modes, thus it is not related to time.

$$\underline{U}_r \frac{d}{dt} (\underline{Z}) = f(\underline{U}_r \underline{Z})$$

$$\underline{U}_r^* \underline{U}_r \frac{d}{dt} (\underline{Z}) = \underline{U}_r^* f(\underline{U}_r \underline{Z})$$

$$\dot{\underline{Z}} = \underline{U}_r^* f(\underline{U}_r \underline{Z}) = g(\underline{Z}) \quad (3)$$

(3) is called Galerkin Projection onto a POD subspace.

27. Discuss the difference between variational autoencoder, denoising autoencoder and a sparse autoencoder.

- Variational AutoEncoder, instead of turn data into a latent vector, it assumes data generated by normal distribution and compress input as its mean and variance parameters in latent space.  
Then the mean & variance are generated statistically with randomness into a vector, finally the decoder will reconstruct data by this vector
- Denoising encoder holds the same architecture as AutoEncoder, the difference is, the inputs of denoising encoder is added with some random noise, while the target is original data. In this way, a trained denoising autoencoder can eliminate noise of input data
- Sparse autoencoder : Holds the sparse information when it is compressed into latent vector. the difference mainly at the loss function.  
Sparse autoencoder use  $L_1$  norm regularization / penalty when computing reconstruction loss.  
In this way, some coefficient can be shrink to ZERO, but  $L_2$  norm can only let coefficient close to ZERO, which ensure the sparsity.



28. Derive ROM for Navier - Stokes Equation by POD  
and Galerkin Projection

Ans

For compressible Navier - Stokes Equation

$$\text{Mass: } \nabla \cdot \underline{\underline{U}} = 0$$

$$\text{Momentum: } \frac{\partial \underline{\underline{U}}}{\partial t} + (\underline{\underline{U}} \cdot \nabla) \underline{\underline{U}} = -\nabla \cdot \underline{\underline{P}} + \frac{1}{Re} \nabla^2 \underline{\underline{U}}$$

① First we do POD on  $\underline{\underline{U}}(x, t)$ , which can  
be expressed as:

$$\begin{aligned} \underline{\underline{U}}(x, t) &= \sum_{k=0}^R a_k(t) \underline{\underline{U}_k(x)} \\ &= \sum_{k=1}^R a_k(t) \underline{\underline{U}_k(x)} + \underline{\underline{U}_0(x)} \end{aligned}$$

where  $R \ll M$

② Apply Galerkin Projection by the form of inner product

for each  $\underline{\underline{U}_j}$  in  $\underline{\underline{V}}^*$   $\underline{\underline{V}}^* \in \mathbb{R}^{R \times N}$

$$\dot{a}(x, t) = \sum_{k=0}^R b_{ik} \dot{a}_k(t) + \sum_{k=1}^R \sum_{j=0}^R l_{ik} a_k(t) \cdot a_j(t)$$

where  $b_{ik} = \langle \nabla \underline{\underline{U}_i}(x), \nabla \underline{\underline{U}_k}(x) \rangle$

$$l_{ik} = \langle \underline{\underline{U}_i}(x), [\underline{\underline{U}_i}(x) \nabla] \underline{\underline{U}_k}(x) \rangle$$

29. Describe different kind of learning covered in this course, provide examples

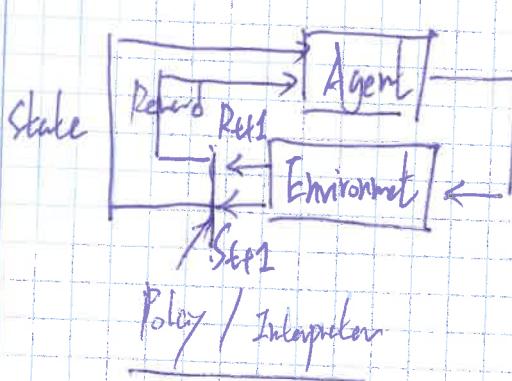
Mainly talking about supervised & unsupervised learning in this course.

- Supervised learning is the learning algorithms are trained through labeled target. Such as Linear Regression, logistic Regression, SVM belongs to supervised learning.
- Unsupervised learning are learning algorithms are trained without target or labeled data. The k-Means, Decision trees are belongs to unsupervised learning.
- In the course, since autoencoder is introduced as well, I suppose it could be categorized as semi-supervised learning to some extent.

3o Describe Reinforcement Learning and provide a schematic representation with the various interacting part.

- Reinforcement Learning is a kind of machine learning that enables an agent to learn an interactive environment by trial & error using feed back from its own actions & experiences
- Unlike Supervise Learning, it does not correct the actions it uses reward & punishment for pos/neg behavior

Unlike Unsupervised learning, it aims at maximize total cumulative reward of agent



1. Environm.: Physical Word where agent operates
2. State: current state of agent
3. Reward: feed back from environment
4. Policy: map state & Action
5. Value: Future reward that an agent would receive by taking action