

EE5904. lecture 4.

GD Method

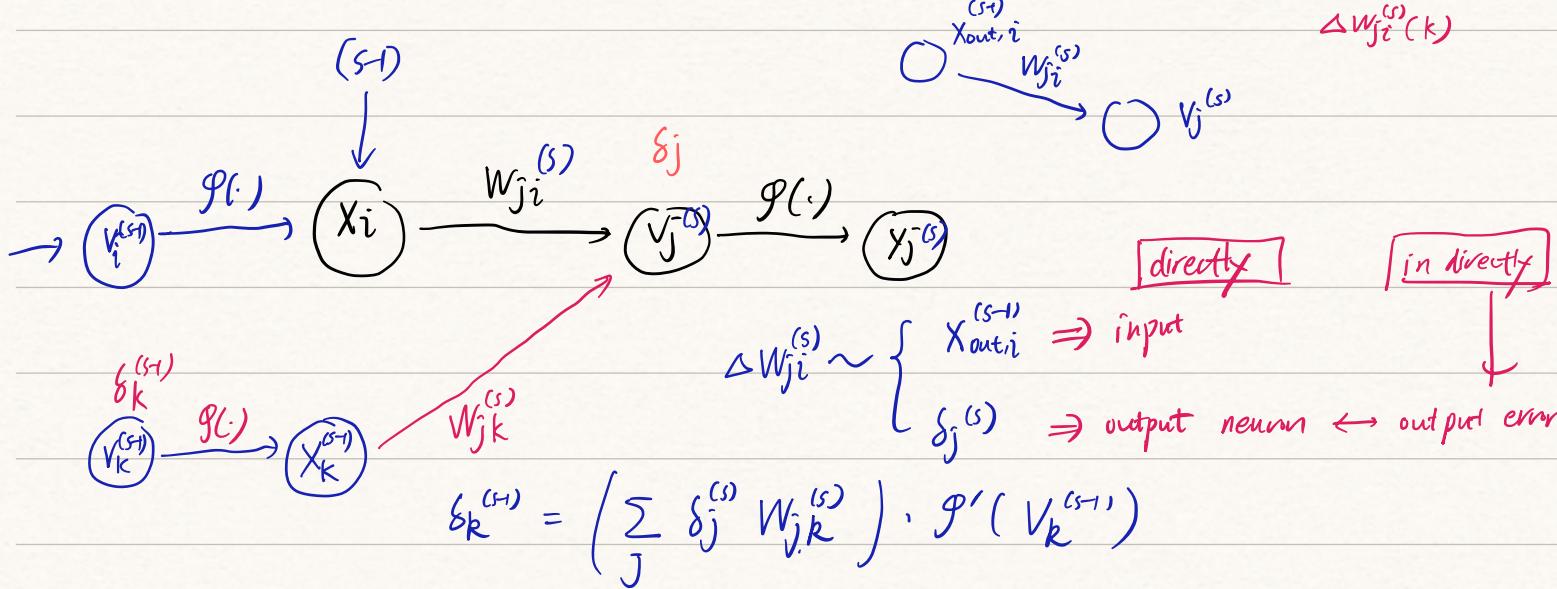
$$w(k+1) = w(k) - \eta g(k)$$

Single layer perceptron: $w(k+1) = w(k) + \eta \underline{e(k)} \overbrace{x(k)}^{\text{input signal}}$

output error input signal

$$\text{MLP: BP} \rightarrow w_{ji}^{(s)}(k+1) = w_{ji}^{(s)}(k) + \eta \underline{s_j^{(s)}} \overbrace{x_{\text{out},i}^{(s-1)}}^{\text{input signal}}$$

$$\Delta w_{ji}^{(s)}(k)$$



Algorithm Framework

① Initialization $w_{ij}^{(s)}(0)$

② use $w_{ij}^{(s)}(k)$ $k=0, 1, 2, \dots$

to compute $\boxed{x_{\text{out},i}^{(s)}} \xrightarrow{w_{ki}^{(s+1)}} \boxed{v_k^{(s+1)}} \xrightarrow{g(\cdot)} \boxed{x_{\text{out},k}^{(s+1)}} \rightarrow \dots \dots$

forward computation

③ In step 2, we attain $\begin{bmatrix} \underline{\text{output}}(k) \\ \underline{d}(k) \end{bmatrix}$

\Rightarrow error $e(k) = \underline{d}(k) - \underline{\text{output}}(k)$

||

$$\text{optimization objective } E(k) = \frac{1}{2} \|e(k)\|_2^2$$



update

process

$$W(k+1) = W(k) - \eta \nabla_w E(k) \quad \Big|_{w=W(k)}$$

$$W_{ji}^{(s)}(k+1) = W_{ji}^{(s)}(k) - \eta^{(s)} \frac{\partial E(k)}{\partial W_{ji}^{(s)}} \quad \Big|_{W_{ji}^{(s)} = W_{ji}^{(s)}(k)}$$

Step 1

the output layer

$$= W_{ji}^{(s)}(k) + \eta^{(s)} \delta_j^{(s)}(k) x_{out,i}^{(s-1)}(k)$$

$$\delta_j^{(e)}(k) = e_j(k) \cdot g'(V_j^{(e)}(k))$$

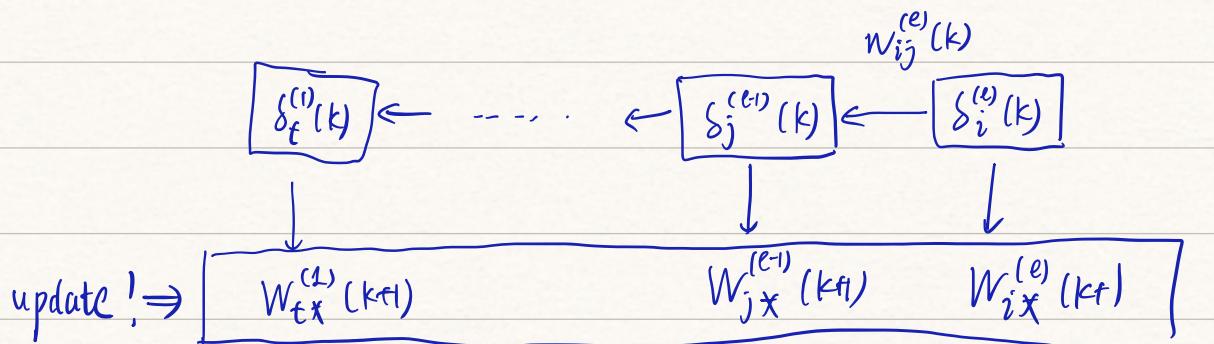
$$\text{Step 2: } \delta_i^{(e-1)}(k) = \left(\sum_j \delta_j^{(e)}(k) W_{ji}^{(e)}(k) \right) g'(V_i^{(e-1)}(k))$$

:

$$\delta_i^{(1)}(k) = \dots$$

Step 3: update with

$$W_{ji}^{(s)}(k+1) = W_{ji}^{(s)}(k) + \eta^{(s)} \delta_j^{(s)}(k) x_{out,i}^{(s-1)}(k)$$



Backwards \Rightarrow
$$\begin{bmatrix} \delta_i^{(e)}(k) \end{bmatrix}$$

$(e-1)$

Some limits:

⇒ the cost surface is typically non-quadratic & non-convex

& high-dimensional

→ { convergence slow
not converge
converge to a bad sol[±]

→ How to construct a NN to avoid so many questions?

① How to use data?

{ Sequential → online
Batch

1. sequential

$$\Delta w_{ij}(1) = -\eta \frac{\partial E(1)}{\partial w_{ij}} \quad | \quad w_{ij} = w_{ij}(1)$$

$$w_{ij}(2) = w_{ij}(1) + \Delta w_{ij}(1)$$

$x(1), d(1) \Rightarrow x(2), d(2) \Rightarrow \dots \Rightarrow x(N), d(N)$

$$\Delta w_{ij}(2) = -\eta \frac{\partial E(2)}{\partial w_{ij}(2)}$$

$$w_{ij}(3)$$

one epoch

↓
change the sequence

事实上这不是 GD Method

↓
objective always change

2. Batch

$(x(1), d(1)), \dots, (x(N), d(N))$

↓

↓

E_{batch}

$E(i)$

$E(N)$

$$E_{av} = \frac{1}{N} \sum_{i=1}^N E(i) \rightarrow \text{one epoch}$$



Steepest descent

$$\Delta w_{ij} = -\eta \frac{\partial E_{av}}{\partial w_{ij}} \quad \left| \begin{array}{l} w_{ij} = w_{ij}(s) \end{array} \right.$$

standard form!

$$= -\eta \frac{1}{N} \sum_{k=1}^N \frac{\partial E(k)}{\partial w_{ij}} \quad \left| \begin{array}{l} w_{ij} = w_{ij}(s) \end{array} \right.$$

{ Batch : One epoch \rightarrow one update

Sequential : one epoch \rightarrow N updates

Which is BETTER? (one output case)

Sequential mode : $E(k) = \frac{1}{2} e(k)^2$

$$\Delta w_{ij}(k) = -\eta \frac{\partial E(k)}{\partial w_{ij}} \rightarrow \text{true estimate of the truth gradient}$$

NO!
a good estimation?

Batch mode : $E_{av} = \frac{1}{N} \sum_{k=1}^N E(k)$

$$\Delta w_{ij} = -\eta \frac{1}{N} \sum_{k=1}^N \frac{\partial E(k)}{\partial w_{ij}}$$

\rightarrow theoretically better

However, sequential learning generally provides better result.

Why?



local minimum problem!

\hookrightarrow sequential learning can release this!

Since in sequential

learning, the estimate of
truth gradient is very

noise!

\downarrow
jump out of local
minima

noisy!

give us more advantage!

Amazing!

[chance] to have better

ER $\hat{z} \Rightarrow$ 优 \hat{z}

results!

just chance

批处理

② Batch method advantage \Rightarrow guarantee to converge

Sequential learning \Rightarrow no guarantee to converge



- { ① let $\eta^{(s)} \rightarrow \eta^{(s)}(k) \rightarrow 0$ as $k \rightarrow \infty$
- { ② { mini-batches
 - { firstly sequential
 - { then batch

③ Another advantages of Batch Method \Rightarrow the biggest advantage

use second-order information to increase convergence

- { small # of weights \rightarrow Newton Method, LM alg.
- { moderate \rightarrow quasi-Newton
- { large \rightarrow CG method

② normalize inputs

reasonable

if $x_i \in [0,1] \rightarrow \text{normalization}$

$$\Delta w_{ji}^{(l)}(k) = \eta s_j^{(l)}(k) x_i(k)$$

focus on $j=1$

\downarrow
same for fixed j

$$x_i \circ \rightarrow v_i^{(l)}$$

(Heuristic)

→ if $x_i(k) \geq 0$ for all i ,

then it output ≥ 0 , all change $\Delta w_{ji}^{(l)}(k) \geq 0$

→ not good

not efficient

$[1,1]$ likely

we should normalize the inputs so their average ≈ 0

Moreover, the bigger the input, the larger the change!

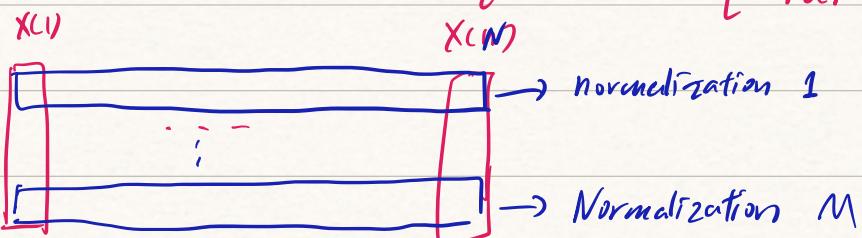
Advantage : Normalization \leadsto faster convergence

Normalization Method :-

$$\bar{x}_i = \frac{\sum_{n=1}^N x_i(n)}{N}$$

$$s_i = \sqrt{\frac{\sum (x_i(n) - \bar{x}_i)^2}{N}}$$

$$\Rightarrow x_i'(n) = \frac{x_i(n) - \bar{x}_i}{s_i} \sim \begin{cases} \text{mean 0} \\ \text{var 1} \end{cases}$$



③ The choice of activation function

① hidden neuron

$\begin{cases} \text{logistic} & \text{logsig: } \mathbb{R} \rightarrow [0, 1] \\ \text{hyperbolic tangent} & \mathbb{R} \rightarrow [-1, 1] \end{cases} \rightarrow \text{tansig}(x) = \frac{2}{1+e^{-2x}} - 1$

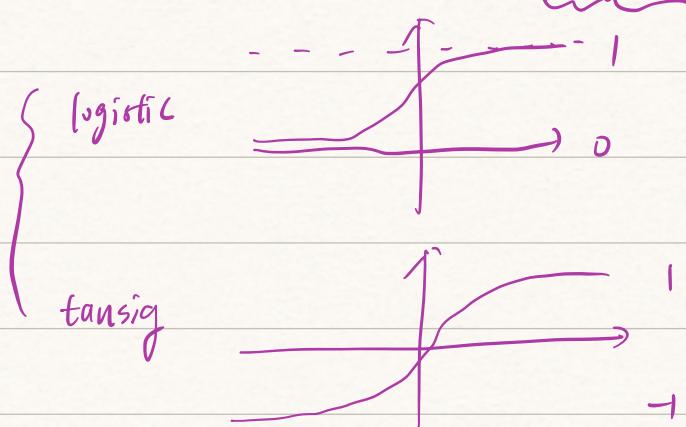
$\log \text{sig}(x) = \frac{1}{1+e^{-x}}$

[preferred] → [similar reason as normalization]

② output neuron

$\begin{cases} \text{reg} \rightarrow \text{identity mapping} & (\text{we can still use logsig}) \\ \text{PR} \rightarrow \text{logistic} & \text{or tansig to normalize} \end{cases}$

④ For PR, how to design target value?



it set output

$\begin{cases} 0 \\ 1 \end{cases}$ or $\begin{cases} -1 \\ 1 \end{cases}$

(Disadvantage)

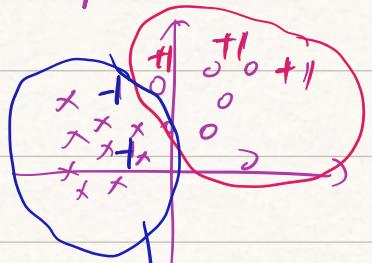
① then the weights are forced to be AS BIG AS POSSIBLE

② and if Induced local field $\rightarrow \infty$

\Rightarrow gradient $\rightarrow 0$ (update slow!)

③ Large weights will make all the output value

close to $\{0\}$ or $\{1\}$



\Rightarrow Decision Boundary are meaningless!

not what we want

\rightarrow As a consequence, we don't set $\{0\}$ or $\{1\}$

instead we set the target value to the
maximum second derivative on the sigmoid function!

最陡峭的地方

$$\left\{ \begin{array}{l} 0.2 \\ 0.8 \end{array} \right.$$

$$\left\{ \begin{array}{l} -0.6 \\ 0.6 \end{array} \right.$$

Recommended!

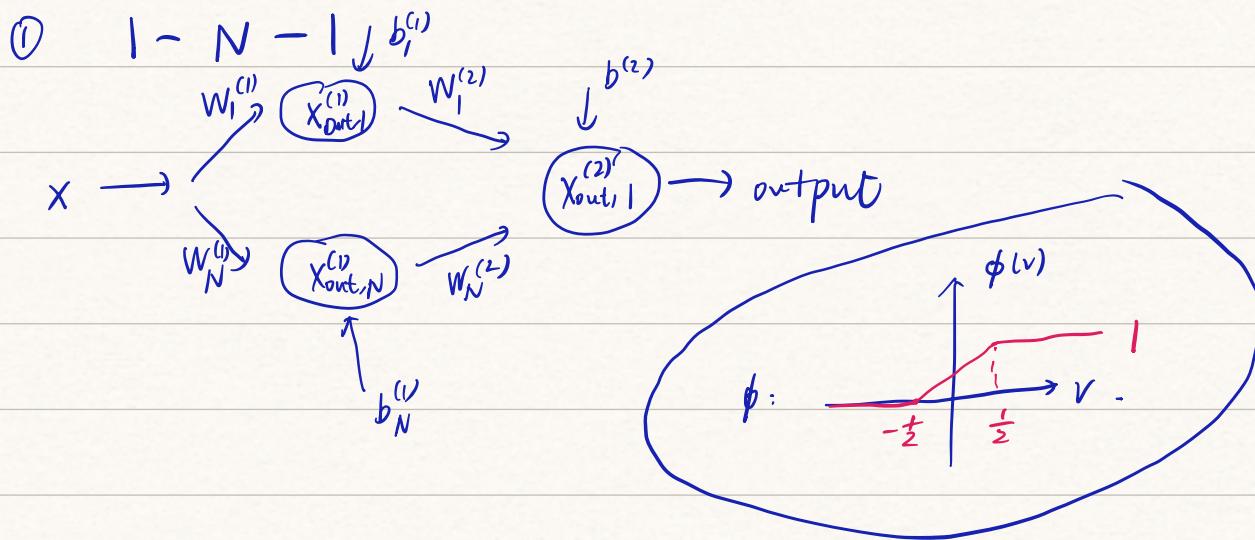
⑤

How many hidden layers
neurons?

Section 4.14

⇒ Geometrical meanings of { # of neurons ?
weights ?
bias ? }

① Hidden Neuron (MLP)



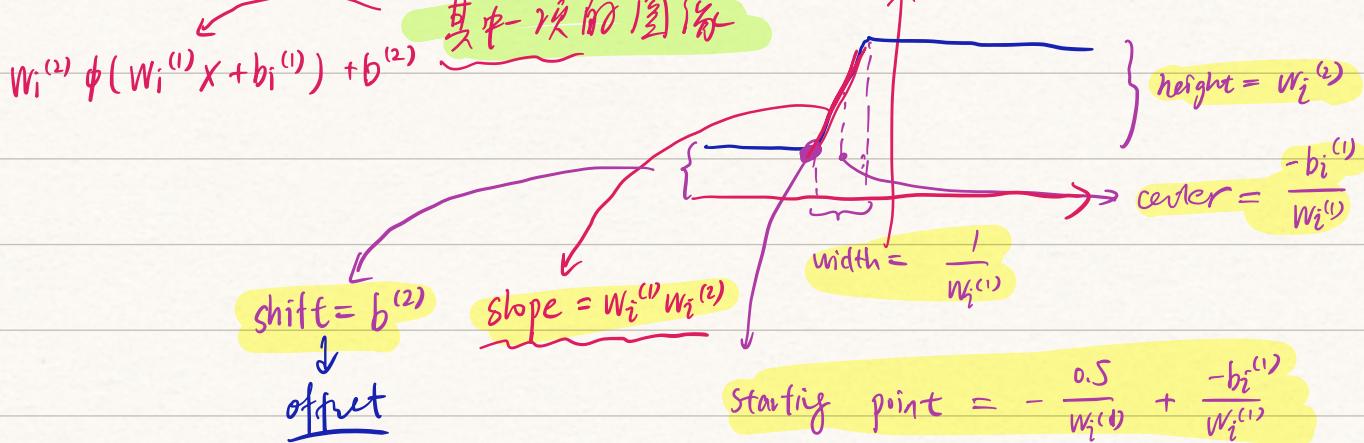
$$y(x) = \sum_{i=1}^N w_i^{(2)} \phi(w_i^{(1)} x + b_i^{(1)}) + b^{(2)}$$

平移 ↓

拉伸, 翻转, 平移



一共有 N 项

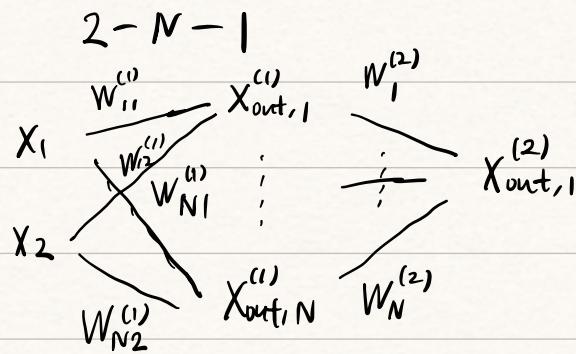


$$\text{Starting point} = -\frac{0.5}{w_i^{(1)}} + \frac{-b_i^{(1)}}{w_i^{(1)}}$$

Guideline! → guideline! → I-N-1 MLP

of neurons (hidden) \approx # of line segments !

Consider for 2-dimensional case.



$$\Rightarrow y(\underline{x}) = \sum_{i=1}^N w_i^{(2)} \phi(w_{i1}^{(1)}x_1 + w_{i2}^{(1)}x_2 + b_i^{(1)}) + b^{(2)}$$

↓
imagine what this like!

Gaussian Hill $\rightarrow \phi(\cdot) \rightarrow \text{logsig} \Rightarrow 2-3-1$ can give good approximation

$\phi(\cdot) \rightarrow \text{linear} \Rightarrow 2-20-1$

② More Hidden Layer!

$$\left\{ \begin{array}{l} 1-3-1 \rightarrow \text{params: } 18 + 10 \times 1 = 28 \\ 1-3-3-1 \rightarrow \text{params: } 2 \times 3 + 4 \times 3 + 4 \times 1 = 22 \end{array} \right.$$

\downarrow decrease!

but structure more complicated

more likely to trap into local minima!

Reason:

cost surface is more complicated

Motivation

① If a task can be decomposed into several sub-tasks



simple t- solve

Multi-layers ≥ 2

② else:

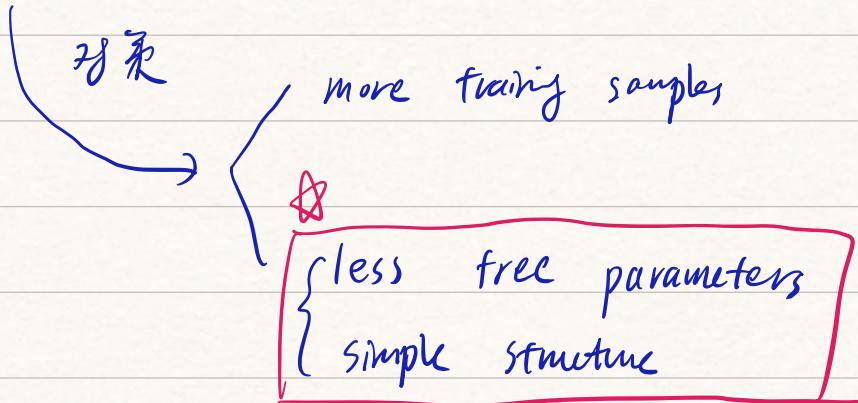
1 hidden layers

⑥ Generalization { under-fitting
over-fitting

{ # of samples (training) small ⇒ over-fitting
of free parameters
Training process

Guideline:

$$N \approx O\left(\frac{w}{\epsilon}\right)$$
 { w: # of free parameters
ε, error



GOAL: choose parameters as less as possible!

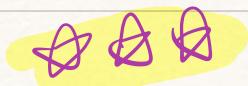
Method:  use line segments to estimate

Metric

the least # of parameters

& structure!

①



the value of weights is more important than
the size of network!

more important

can be explained by the geometric meaning of MLP

Strategy : Penalty / Regularization

previously $\lambda = 0$!

no penalty

$$F = E_D + \lambda E_W$$

Model Modification (objective)

$$\left\{ \begin{array}{l} ① E_W = \sum_{i=1}^N [w_i^{(1)}{}^2 + w_i^{(2)}{}^2 + b_i^{(1)}{}^2 + b_i^{(2)}{}^2] \\ ② E_W = \sum_{i=1}^N [w_i^{(2)} w_i^{(1)}]^2 \end{array} \right.$$

slope
more reasonable!

Occam's Razor



We are interested in minimal structure of NN.

↓
difficult to estimate when dimension is
high!

⇒ Our strategy to overcome over-fitting ⇒

- ① minimal structure of NN
- ② Regularization!

SVD!

$$H_{N \times N} = U_{N \times N} \Sigma_{N \times N} V_{N \times N}^T$$

$$= \sum_{i=1}^k b_i u_i v_i^T$$

if rank(H) = k

low rank approximation

↓
numerical rank

numerical rank
actual

$$A = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix} \quad B = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$$

numerical rank 3

3

actual rank

3

2!

$$SVD = \begin{pmatrix} 2.247 & & \\ & 0.8 & \\ & & 0.55 \end{pmatrix}$$

effective rank = 3

$$SVD = \begin{pmatrix} 2.1358 & & \\ & 0.6622 & \\ & & 0.0071 \end{pmatrix}$$

effective rank = 2

⇒ Goal: Define Effective rank of hidden Neurons

↓
find the redundant hidden neuron

→ N training samples $\in \mathbb{R}^M$ n hidden neurons

inputs $\left\{ \begin{array}{l} x_1(1) \\ \vdots \\ x_M(1) \end{array} \right.$ $\xrightarrow{\begin{array}{c} w_{11} \\ \vdots \\ w_{1n} \end{array}} v_i = \sum_{j=1}^M w_{ij} x_j(k) + b_i \rightarrow f_i(v_i(k)) = h_{ki}$

$$h_{ki} = f_i \left(\sum_{j=1}^M w_{ij} x_j(k) + b_i \right) \in H_{N \times n} \rightarrow \text{hidden neuron matrix}$$

$$\downarrow \quad \left(\begin{array}{c} h(1)^T \\ \vdots \\ h(N)^T \end{array} \right)_{N \times n} \quad h(i)^T = [h_{1(i)} \dots h_{n(i)}]$$

$$H_{N \times n} = V_{N \times N} \Sigma_{N \times n} V_{N \times n}^T$$

↓ cutoff
define a threshold!

Framework

singular value $b_1 \geq b_2 \geq \dots \geq b_n \geq 0$

normalization

$$\tilde{b}_1 \geq \tilde{b}_2 \geq \dots \geq \tilde{b}_n \geq 0$$

$$\textcircled{1} \quad \tilde{b}_i = \frac{b_i}{\sum b_i}$$

$$\textcircled{2} \quad \tilde{b}_i = \frac{b_i^2}{\sum b_i}$$

↓

threshold γ
 $(95\%, 99\%)$

$$\min_k \sum_{i=1}^k b_i \geq \gamma$$

\rightarrow cutoff

① start with n (# of hidden neurons) big enough

↓ cut

② cut off

↓ cut off

:

:

terminate!

for large-dimension

Conclusion: use SVD on $H_{N \times n} := \begin{pmatrix} h^{(1)} \\ \vdots \\ h^{(N)} \end{pmatrix}_{N \times n}$

to shrink hidden neurons !

ΣΙΣ

: Determine # of hidden neurons :

- ① dimension low: minimal # of line segments
- ② high dimension: start with large number

↓

use SVD to cut !



2:53:00 Η το ΣΙΣ

Deal with over-fitting

- { ① minimal structure
- ② Regularization

main questions:

- ① # of hidden layers?
- ② # of hidden neurons?
- ③ choice of activation $f = \begin{cases} \text{hidden} \rightarrow \text{tansig} \\ \text{output} \end{cases}$
- ④ normalization \rightarrow preprocess
- ⑤ { Batch
Sequential }
- ⑥ over-fitting { Minimal structure
Regularization }