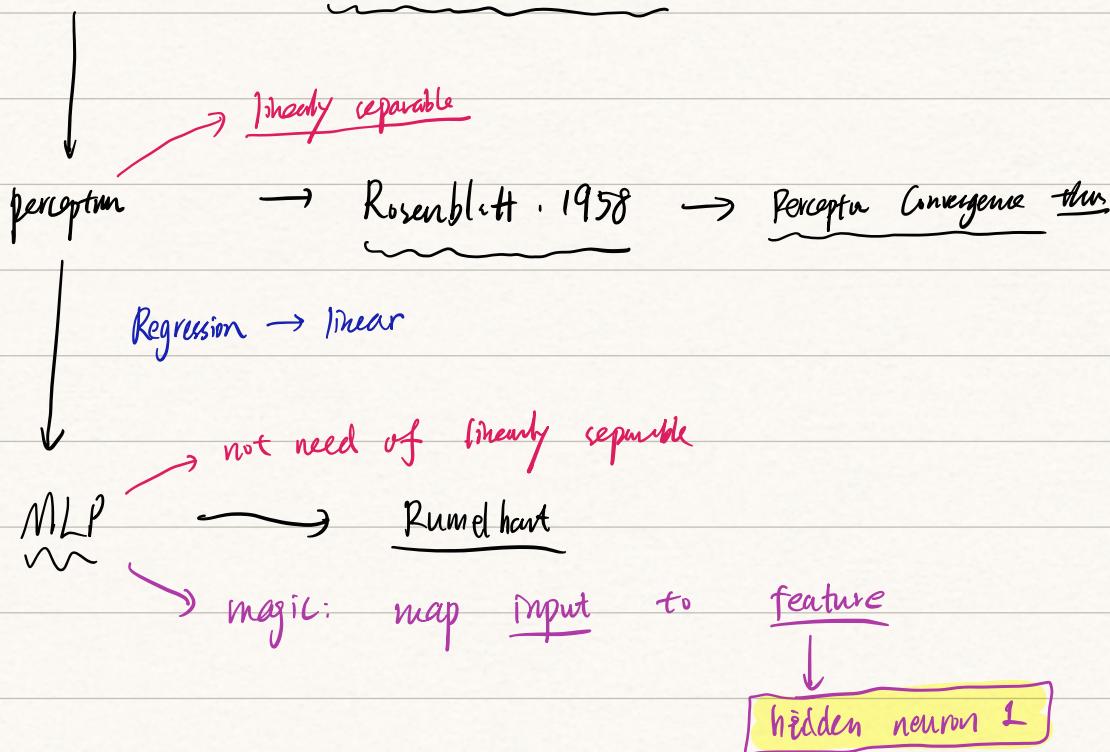


EE5904 lec5

Mathematics Neuron \rightarrow McCulloch & Pitts, 1943



Learning Algorithm: Stochastic Gradient

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

output err input signal

Back propagation \rightarrow easy to compute!

How to design?

① # of hidden layer $\left\{ \begin{array}{l} x-x-x \\ x-x-x-x \end{array} \right.$

② # of hidden neuron

③ activation f° (hidden)

(4) activation f^2 (out) { regression
PR

(5) pre-process \rightarrow input \Rightarrow Normalization

(6) { Batch \rightarrow try first { fast!
But local minimum
Sequential \rightarrow { slow
better results!

(7) overfitting { minimal structure
regularization

MLP { regression
PR

universal approximator
A way of approximation! \Leftrightarrow Real model is too complicated!

desired output
Recognition \rightarrow users-defined
Regression \rightarrow observation

WEAKNESS

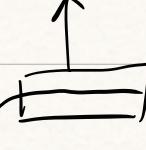
just approximation \Leftrightarrow cannot explain! (no intuitive explanation)

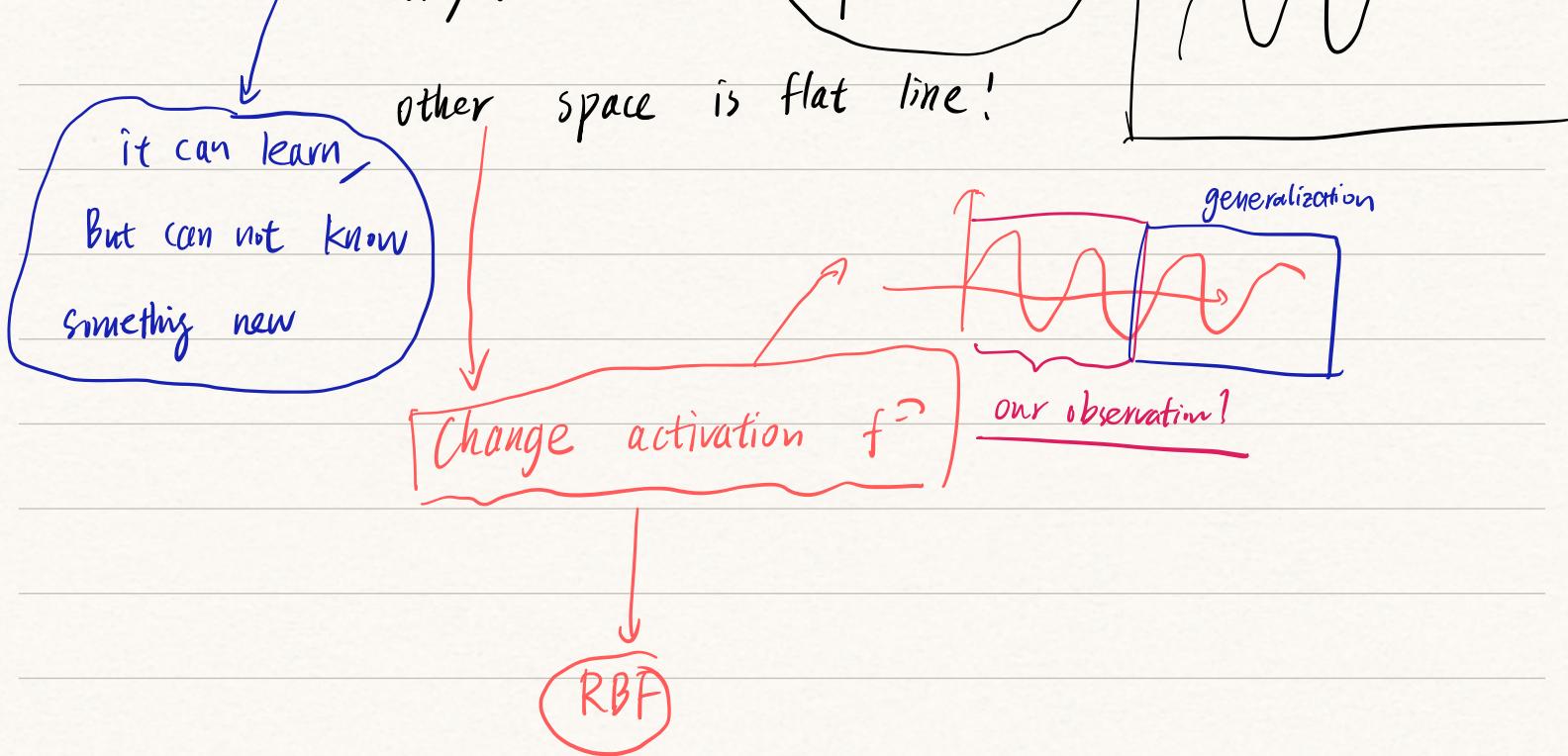
Black box

WEAKNESS 2

only focus on the input area

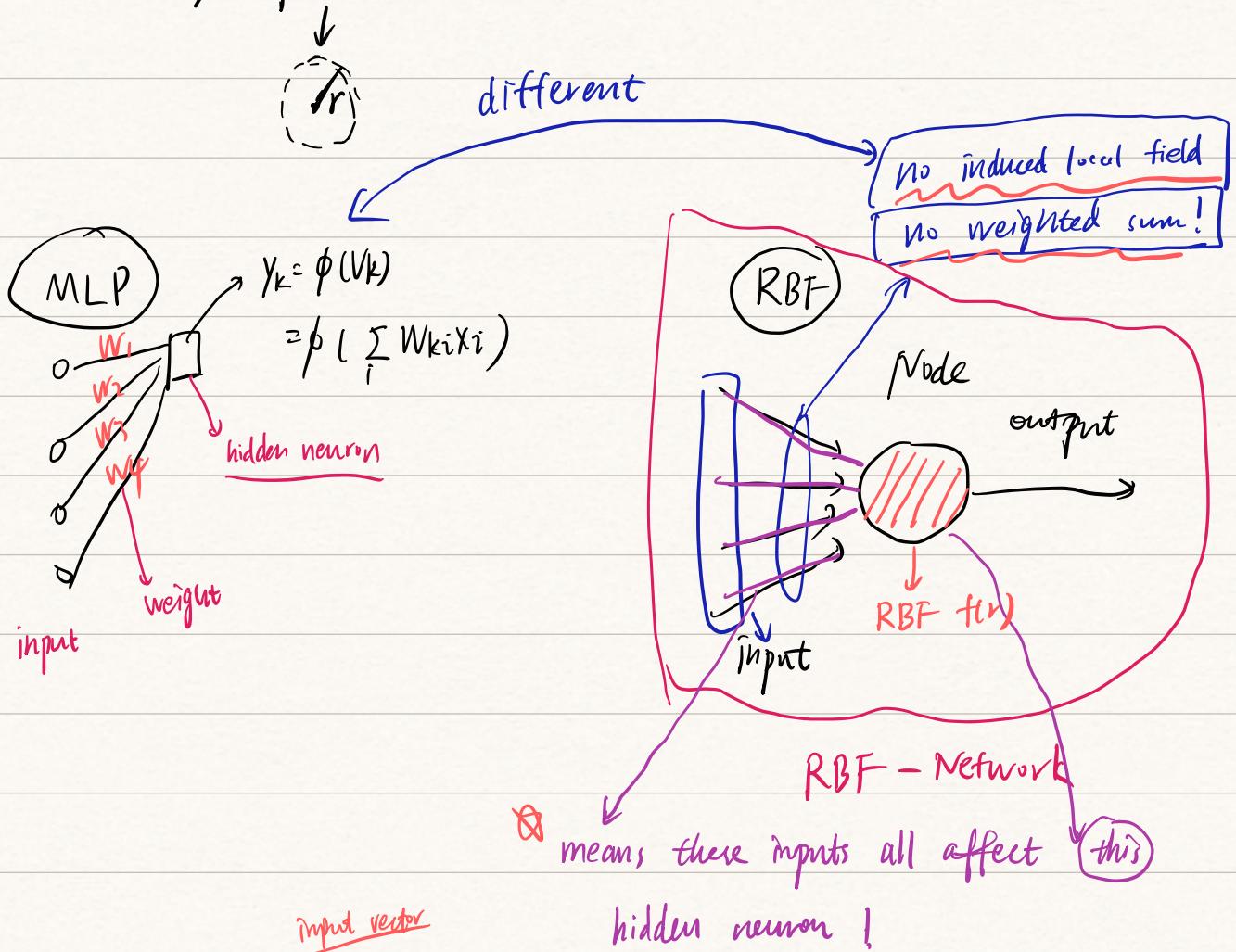
BAD





Radial Basis Function \longleftrightarrow function approximation
(RBF)

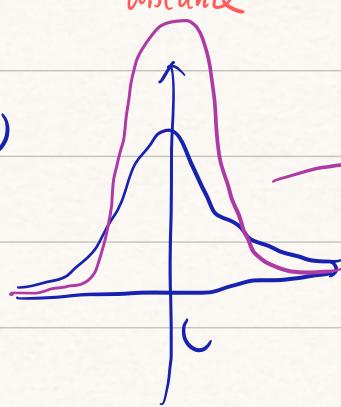
RBF: f^* only depends on the radial distance



Process

$$① \quad g(x) = g(\underbrace{\|\underline{x} - \underline{c}\|}_\text{distance}) = g(r)$$

$$② \quad g(\cdot)$$



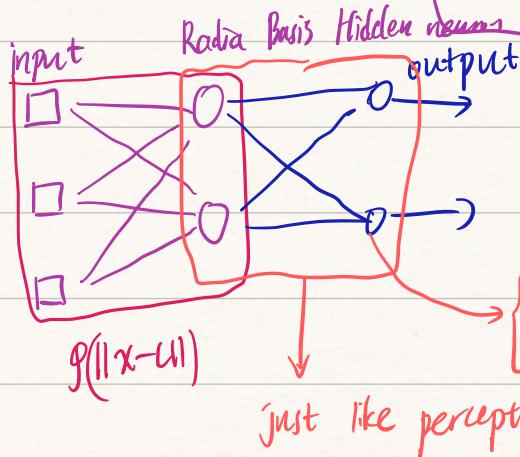
Gaussian-like function

our goal: if x match c

\downarrow
 $\|\underline{x} - \underline{c}\|$ small

then $g(\|\underline{x} - \underline{c}\|)$ BIG!

Note



- ① weighted sum \rightarrow induced local field
- ② activation f^c

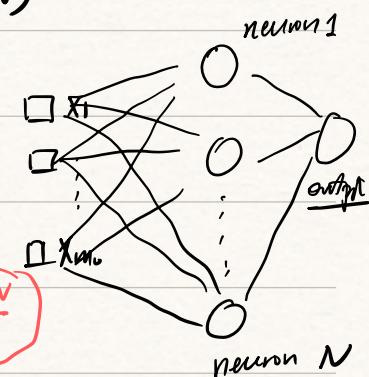
just like perception

$\left\{ \begin{array}{l} \# \text{ of hidden neurons : } \# \text{ of data samples } (N) \\ \text{center : } \underline{x}_i \end{array} \right.$

$$f(x) = \sum_{i=1}^N w_i \phi(\|\underline{x} - \underline{x}_i\|)$$

$$\underline{x}(1), \dots, \underline{x}(N)$$

$$\underline{x}(i) \in \mathbb{R}^{m_x}$$



$$\Rightarrow f(\underline{x}_j) = \sum_{i=1}^N w_i \phi(\|\underline{x}_j - \underline{x}_i\|) = d_j$$

$$\Leftrightarrow \begin{pmatrix} d_1 \\ \vdots \\ d_N \end{pmatrix} = \begin{pmatrix} \phi_{11} & \cdots & \phi_{1N} \\ \vdots & \ddots & \vdots \\ \phi_{N1} & \cdots & \phi_{NN} \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_N \end{pmatrix}$$

$$\uparrow \quad \phi_{ij} = \phi(\|\underline{x}_i - \underline{x}_j\|)$$

$$\psi_{ji} = \phi(\|\mathbf{x}_j - \mathbf{x}_i\|)$$

\downarrow

$$\underline{d} = \underline{\Phi} \underline{w} \Rightarrow w = \underline{\Phi}^{-1} \underline{d} \quad \text{if } \underline{\Phi} \text{ is invertible}$$

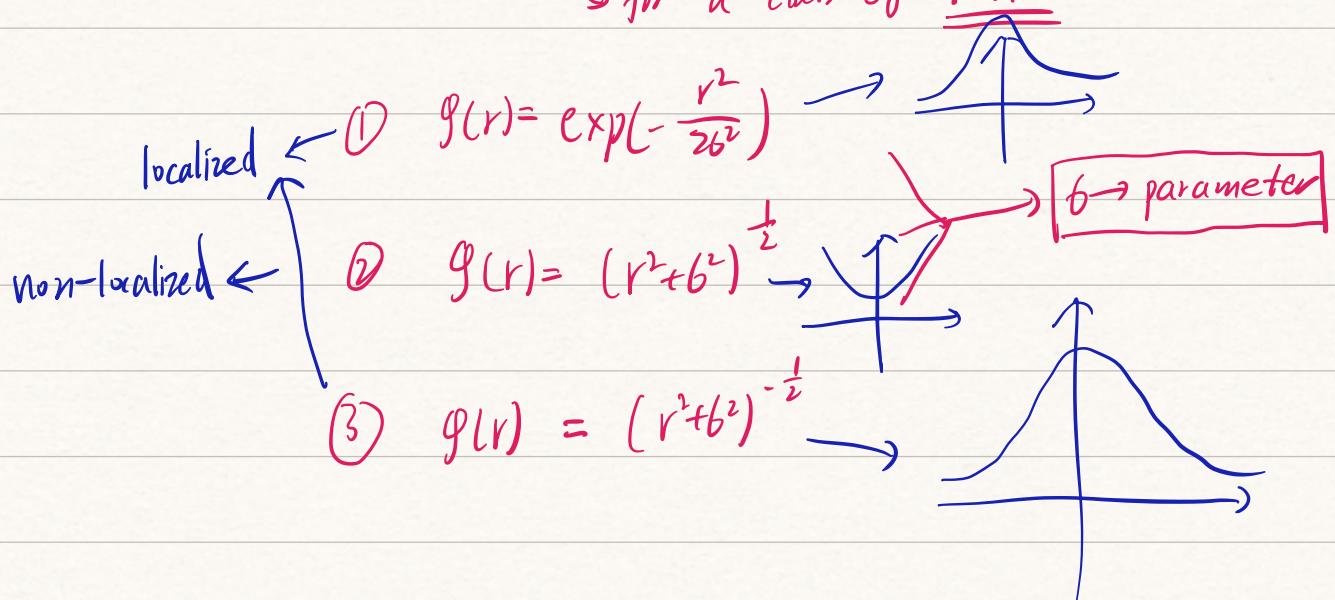
However, if $\underline{\Phi} \rightarrow$ noninvertible?

\downarrow

thm: if $\underbrace{\mathbf{x}(1), \dots, \mathbf{x}(n)}$ in \mathbb{R}^{m_o}
distinct points

then $\psi_{ji} = \phi(\|\mathbf{x}_j - \mathbf{x}_i\|)$ non-singular

for a class of RBF



the algorithm is quite easy!

$\downarrow \underline{\Phi}, \underline{d} \Rightarrow w = \underline{\Phi}^{-1} \underline{d}$

Limitation!

↓
over-fitting! \leftrightarrow noisy data!



Strategy: shrink the hidden neuron!

How to improve?

- ↓
- { not fix the # of hidden units
 - { not fix the center of RBF \leftarrow learn!
 - spread parameter? \leftarrow learn

we can also add Bias term w_0

Gaussian RBF

$$\varphi_i(\|x - \text{M}_i\|) = \exp\left(-\frac{\|x - \text{M}_i\|^2}{2\sigma_i^2}\right)$$

$\begin{cases} \text{M}_i \\ \sigma_i \\ w_i \end{cases}$ } \rightarrow How to find?

$$y(x) = \sum_{i=1}^M w_i \varphi_i(\|x - \text{M}_i\|)$$

And it also has Universal Approximation Theorem for RBFN

↓
just like MLP

only existence

$E(w) = \sum_{i=1}^N e(i)^2 = \sum_{i=1}^n (d(i) - y(i))^2$

① Steepest descent

↑
calculate $\frac{\partial E(w)}{\partial M_i}$ $\frac{\partial E(w)}{\partial b_i}$ $\frac{\partial E(w)}{\partial w_i}$

↑
we don't have BP Alg. here!

↓
no efficient algorithm

to calculate these!

② our method

↓
by PFF train

- ① M_i & b_i
- ② w_1, \dots, w_M, w_o

according to M_i & b_i

Algorithm:

Stage 1: { para of hidden layer
 $M \rightarrow$ # of hidden neurons \Rightarrow unsupervised
 M_i
 b_i \Rightarrow weights

para between hidden & output

Stage 2 { w_1, \dots, w_M
 w_o \Rightarrow supervised learning

Firstly, stage 2 fix $\begin{cases} w \\ u_1, \dots, u_M \\ b_1, \dots, b_M \end{cases}$

$$y(x_j) = \sum_{i=0}^M w_i \phi_i(x_j) = d_j$$

$$\Rightarrow \begin{pmatrix} \phi_0(x_1) & \dots & \phi_M(x_1) \\ \phi_0(x_N) & \dots & \phi_M(x_N) \end{pmatrix}_{N \times (M+1)} \begin{pmatrix} w_0 \\ \vdots \\ w_M \end{pmatrix} = \begin{pmatrix} d_1 \\ \vdots \\ d_N \end{pmatrix}$$

$$\Rightarrow \phi w = d$$

~~if $N > M+1$~~ ~~hidden neuron~~ \nmid

if $\underbrace{N > M+1}_{\text{hidden neuron}} \rightarrow$ we can not find f

pass through all data samples!

$$\downarrow \\ \text{MSE}$$

$$E(w) = (\phi w - d)^T (\phi w - d)$$

$$= w^T \phi^T \phi w - 2 d^T \phi w$$

optimality condition: $\nabla E(w) = 0$

$$\Rightarrow \nabla E(w) = 2 \phi^T \phi w - 2 \phi^T d = 0$$

$$\Leftrightarrow \boxed{w = (\phi^T \phi)^{-1} \phi^T d}$$

LLS solution

Stage 1.

① fix centers randomly

② clustering algorithm → k-means

supervised learning of RBFN

↓

well result
great cost

↓ like MLP

$$E = \frac{1}{2} \sum_{i=1}^N (e_i)^2$$

$$g_n = \exp\left(-\frac{\|\underline{x}_i - \underline{\mu}_n\|^2}{2\sigma_n^2}\right)$$

$$e_i = d_i - y(\underline{x}_i) = d_i - \sum_{n=1}^M w_n g_n(\underline{x}_i; \underline{\mu}_n, \sigma_n)$$

Gradient Descent → second order method is also OK

$$\begin{cases} \Delta w_i = -\eta_w \frac{\partial E}{\partial w_i} \\ \Delta \mu_i = -\eta_\mu \frac{\partial E}{\partial \mu_i} \\ \Delta \sigma_i = -\eta_\sigma \frac{\partial E}{\partial \sigma_i} \end{cases}$$

Q: How to choose # of hidden neurons?



- 1. minimal structure
- 2. regularization

Goal → overcome overfitting



restrict the value of weight

$$F = E_D + \lambda E_W$$

$$\frac{1}{2} \sum_j e_j^2$$

penalty on complexity (smoothness)

Question: How to measure the smoothness?

derivatives (gradient)!

$$f(x) = \sum_i w_i g(\|x - \mu_i\|)$$

$$g_i(x) = \exp\left(-\frac{\|x - \mu_i\|^2}{2\sigma_i^2}\right)$$

$$\frac{\partial f}{\partial x} = \sum_{i=0}^M w_i \frac{\partial g(\|x - \mu_i\|)}{\partial x}$$

Influence the gradient!

w_i big \Rightarrow less smoothness!

a measurement
of SMOOTHNESS

the motivation of minimizing weights!

$$\frac{\partial \hat{e}^T e}{\partial w} = \frac{\partial \hat{e}^T e}{\partial e} \frac{\partial e}{\partial w}$$

$$F(w) = \frac{1}{2} \sum_i e_i^2 + \frac{1}{2} \lambda w^T w$$

$$= 2e^T \cdot \Phi$$

$$= \frac{1}{2} (\Phi w - d)^T (\Phi w - d) - \frac{1}{2} \lambda w^T w$$

$$\Phi = \begin{pmatrix} g_0(x_1) & \dots & g_M(x_1) \end{pmatrix} \quad w = \begin{pmatrix} w_0 \\ \vdots \\ w_M \end{pmatrix} \quad = 2(\Phi w - d)^T \Phi$$

$$= 2w^T \Phi^T \Phi - 2d^T \Phi$$

$$g_0(x_N) \dots g_M(x_N)$$

$$\nabla F(w) = \Phi^T \Phi w - \Phi^T d + \lambda w$$

$$= (\Phi^T \Phi + \lambda I) w - \Phi^T d$$

$$\frac{\partial F}{\partial w} = \left(\frac{\partial e^T e}{\partial w} + \frac{\partial \lambda w^T w}{\partial w} \right)^{\frac{1}{2}}$$

$$\Rightarrow w = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T d$$

$$= 2 e^T \Phi + 2 \lambda w^T$$

{ MLP - Regularization \rightarrow H/SVD
 RBFN \leftarrow Reg \rightarrow EASY

for classification

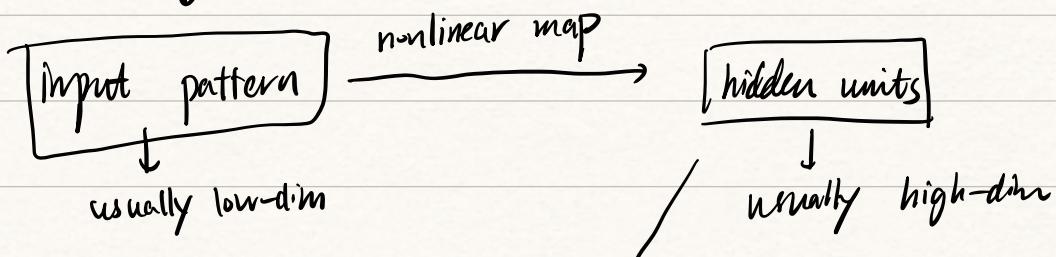
\downarrow
easy to design output

\downarrow
k-class \Rightarrow k-output $\begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \dots \begin{bmatrix} 0 \\ 1 \end{bmatrix}_{k \times 1}$

Note:

All output neurons \rightarrow linear neuron

$\cancel{\text{All power lies in}}$ hidden neurons \rightarrow RBF



linear map

Cover's theorem

output

A complex pattern classification

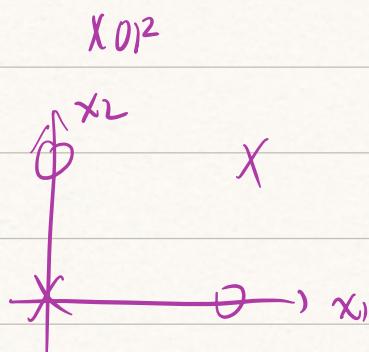
cast in a high-dim space. (nonlinearly)

then it is more likely to be linearly separable than
in low-dim space!

then we can solve it easily!

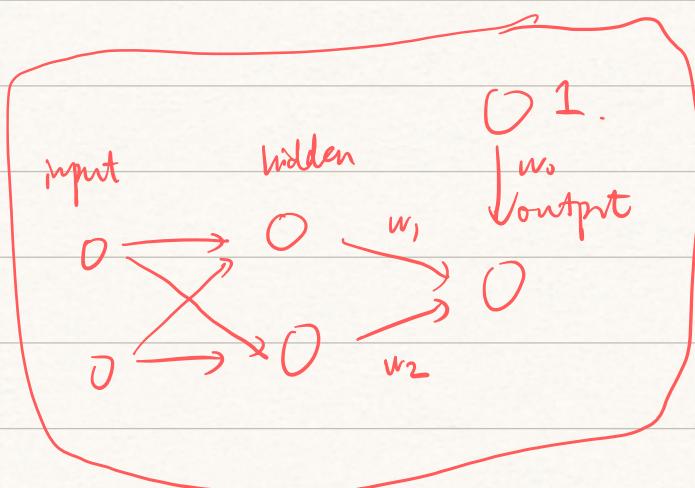
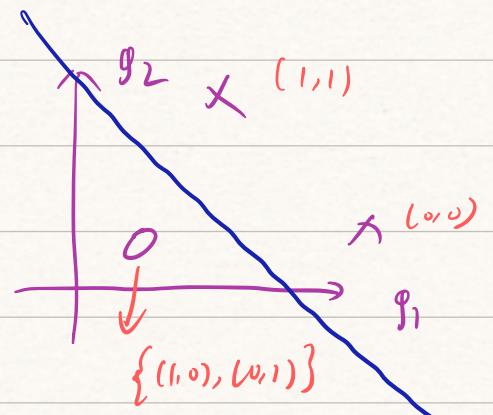
$$\mu_1 = (1, 1, 0)$$

$$\mu_2 = (1, 1, 1)$$



$$g_1(x) = \exp(-\|x - \mu_1\|^2)$$

$$g_2(x) = \exp(-\|x - \mu_2\|^2)$$



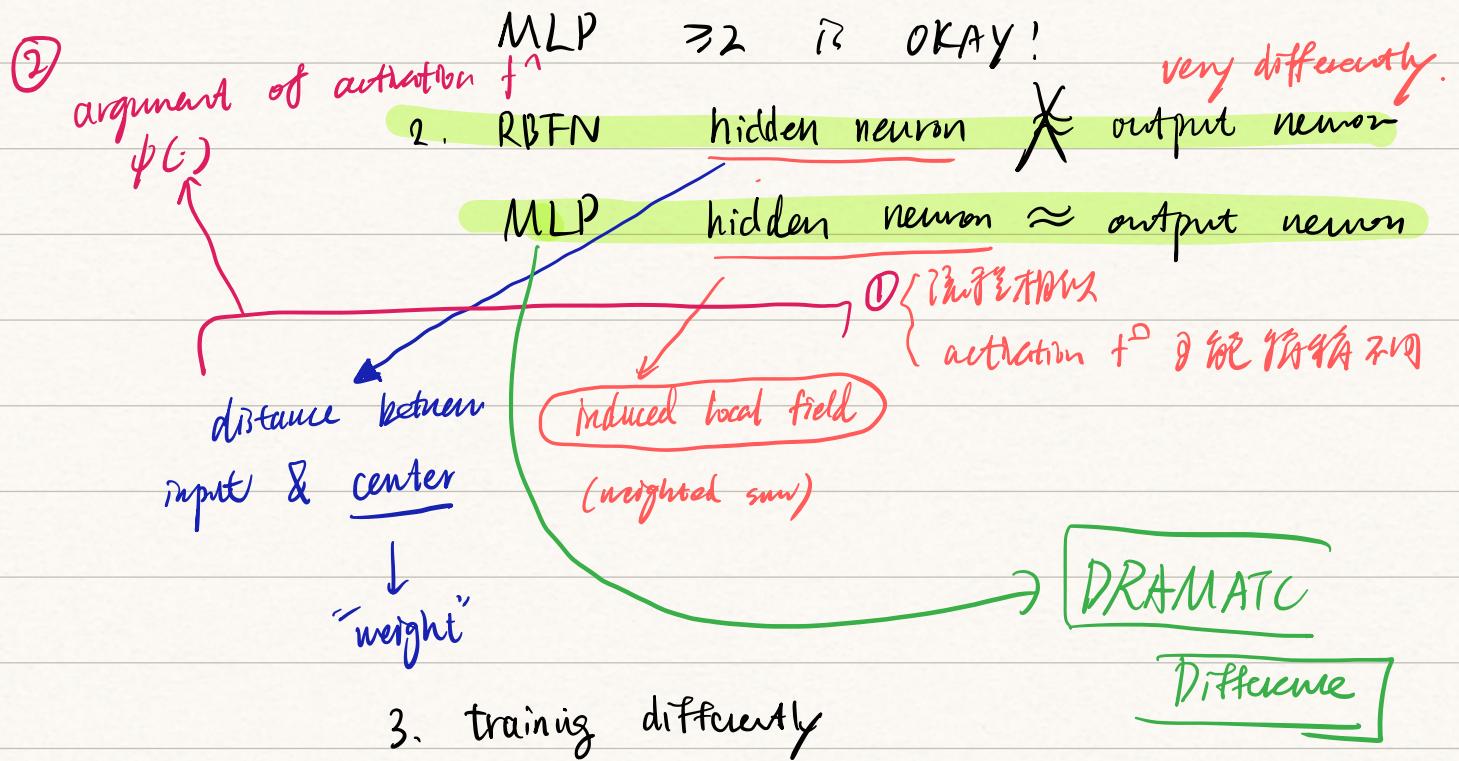
RBFN vs MLP

Similarity ① non-linear feed-forward network

② universal approximators

③ in similar application areas

difference : 1. RBFN cannot have 2 (or more) hidden neurons

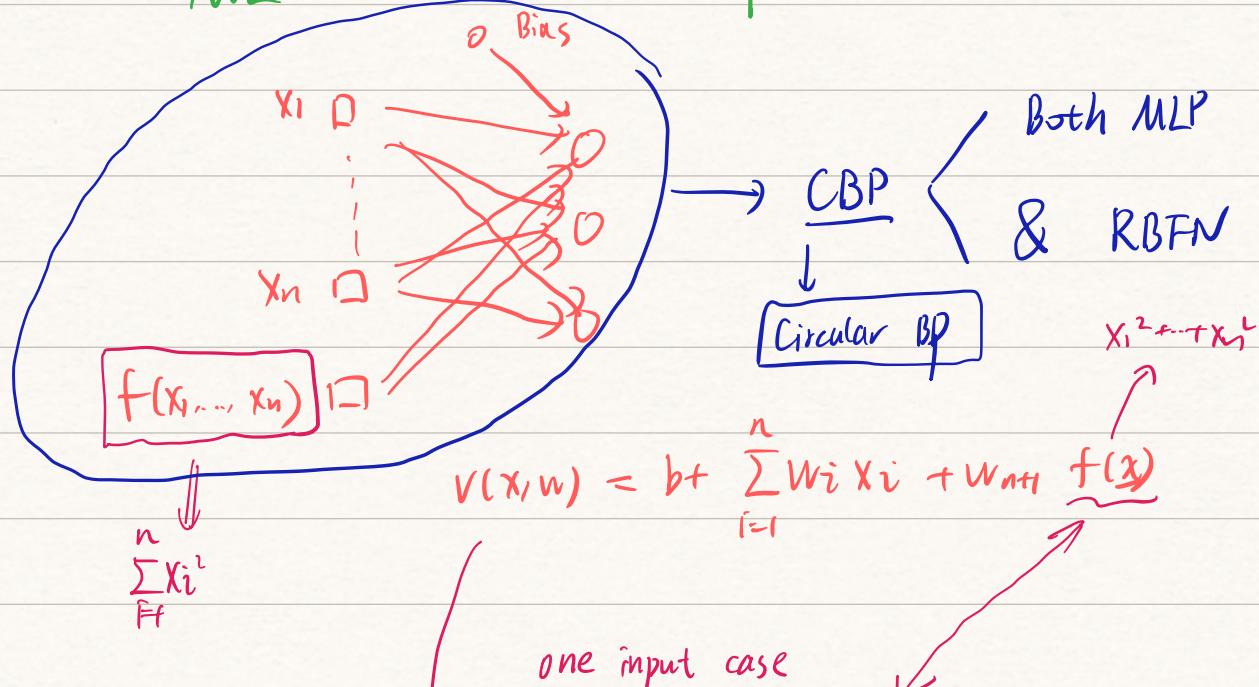


MLP \rightarrow BP

RBFN \rightarrow { 2-stage
GD

Faster

MLP with additional input \rightarrow RBFN



$$b + w_i x_i + \underbrace{(x_i - c)^2 + d}_{\approx}$$

$$\approx (x_i - c)^2 := V^2 \quad \underline{V = \|x - c\|} \\ = g(\|x - c\|)$$

$$V(x, \underline{w}) = w_{n+1} \sum_{i=1}^n \left(x_i + \frac{w_i}{2w_{n+1}} \right)^2 + \infty$$

$$\boxed{V(x) = g(\|x - c\|^2 + \theta)}$$

$$\text{where } g = w_{n+1} \quad \underline{c} = \begin{pmatrix} 0 \\ \vdots \\ c_n \end{pmatrix} \quad \underline{i} = -\frac{\underline{w}}{2w_{n+1}}$$

$$\theta = \frac{b}{w_{n+1}} - \sum_{i=1}^n \frac{w_i^2}{4w_{n+1}}$$

$$\text{let } r = \|x - c\|$$

RBF

$$\leftarrow V(x) = g(\|x - c\|^2 + \theta) = gr^2 + a \rightarrow \text{just find one homogeneous } f$$

$$\phi(V(x)) \rightarrow RBF$$