

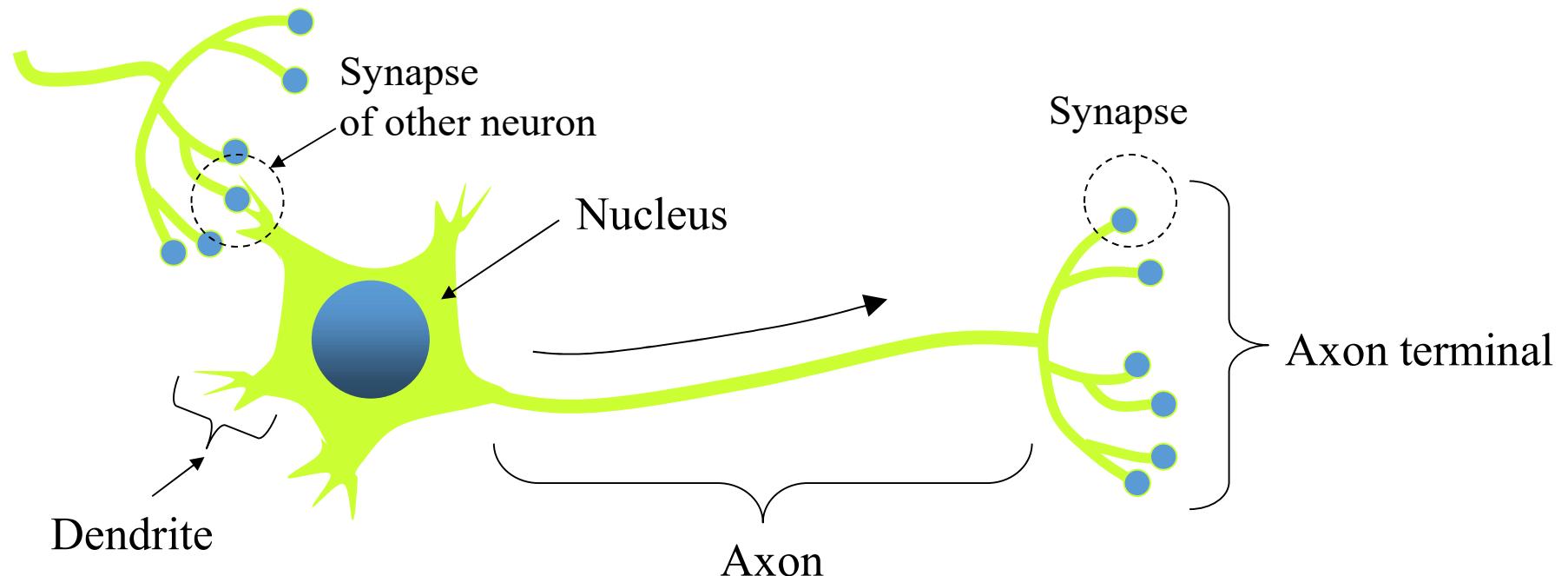
Topics in data engineering

Artificial neural network

Masaomi Kimura

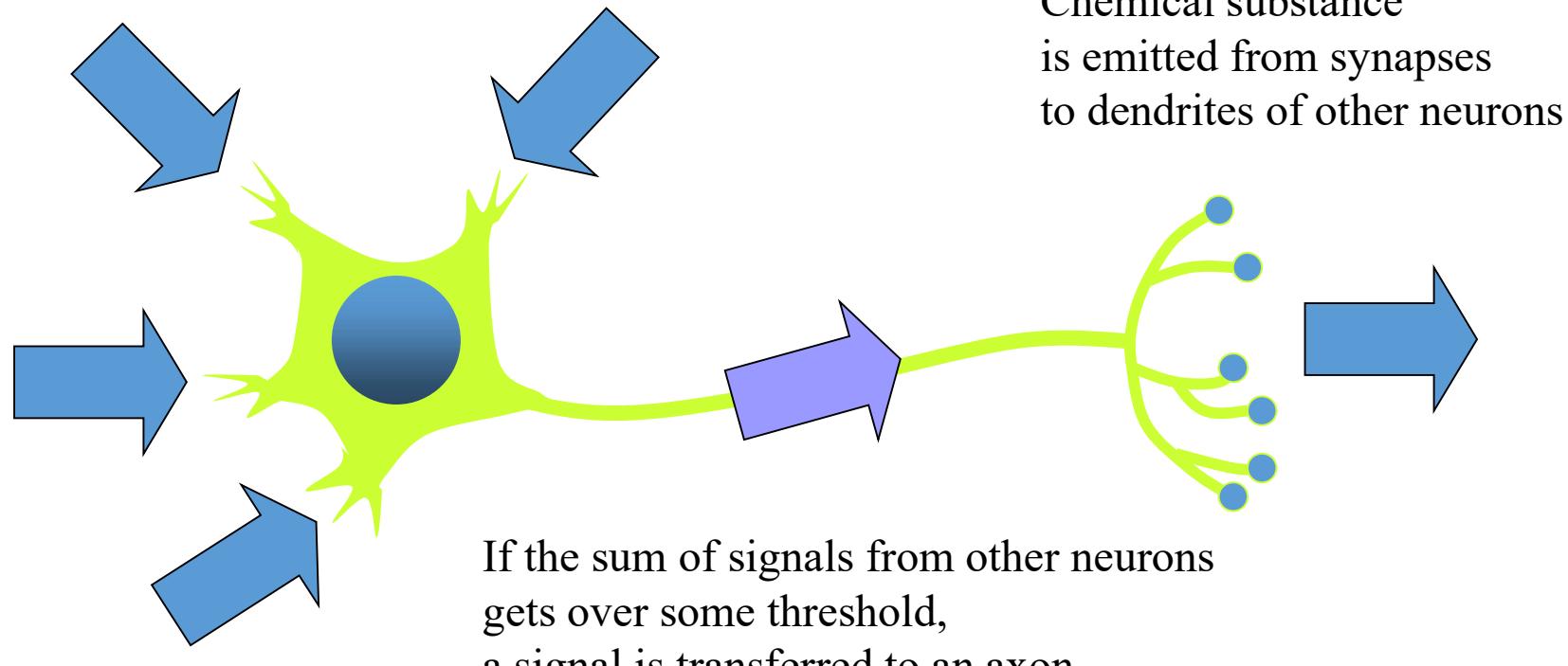
Neuron

- cells in nervous system
 - contains dendrites, a nucleus, an axon
 - Synapses at an axon terminal send a signal to dendrites in other neurons

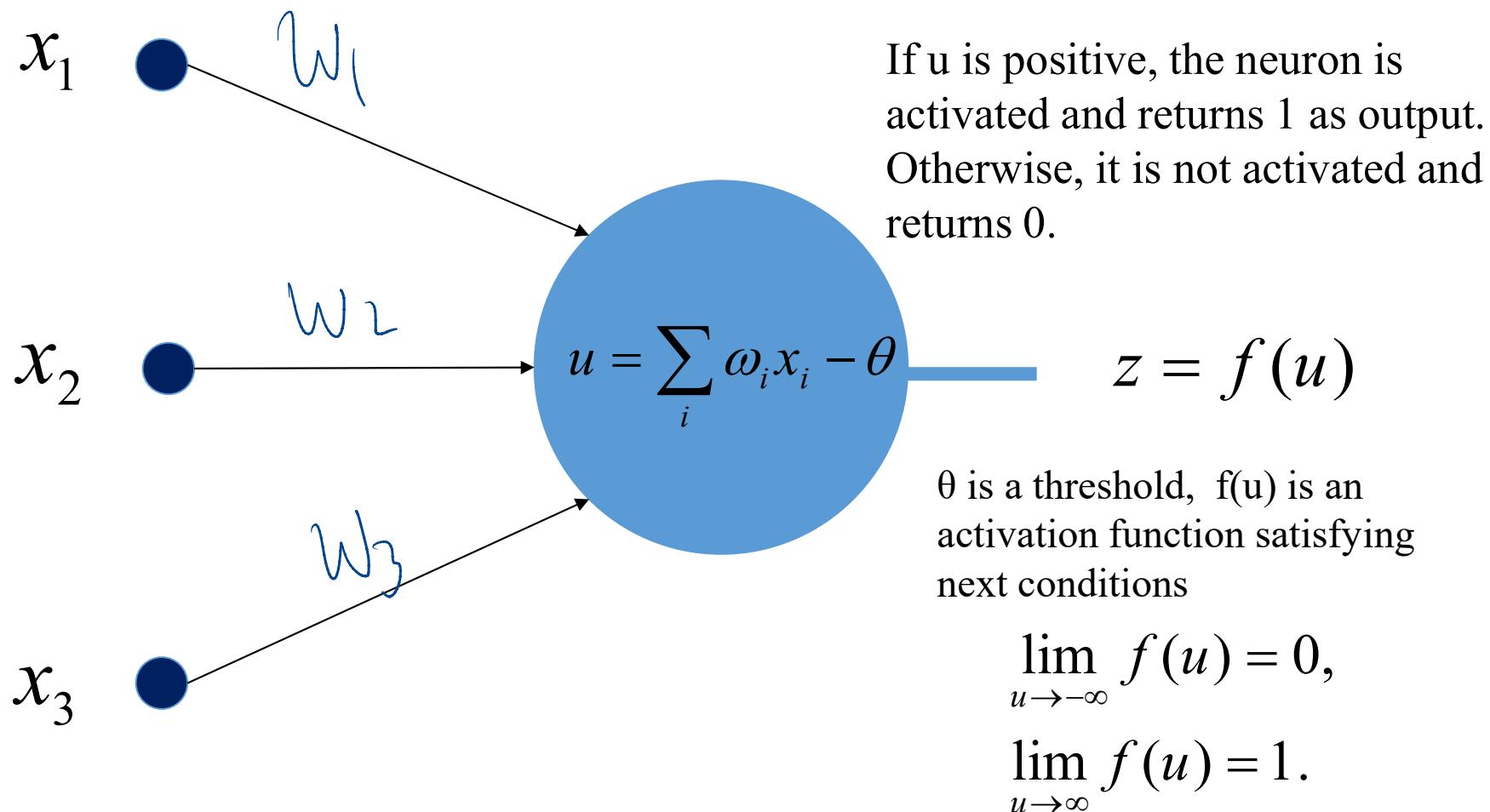


Signal transfer

Signals from synapses of other neurons



Perceptron



Activation function $f(u)$

- Step function

$$f(u) = \begin{cases} 0 & (u < 0) \\ 1 & (u \geq 0) \end{cases}$$

- Sigmoid function

$$\sigma(u) = \frac{1}{1 + e^{-u}}$$

- Hyperbolic tangent function

$$\tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$$

$$\sigma(u) = \frac{1}{1 + e^{-u}}$$

$$\frac{d\sigma}{du} = \frac{d}{du} \left(\frac{1}{1 + e^{-u}} \right)$$

$$= \frac{d}{du} \left(\frac{1}{1 + e^{-u}} \right)^{-1}$$

$$= \frac{d\sigma}{dv} \frac{dv}{du}$$

$$= \left(\frac{d}{dv} v^{-1} \right) \left(\frac{d}{du} (1 + e^{-u}) \right)$$

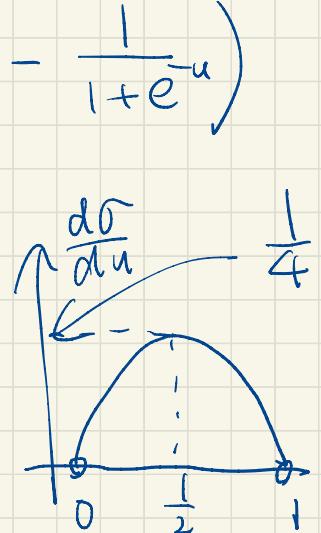
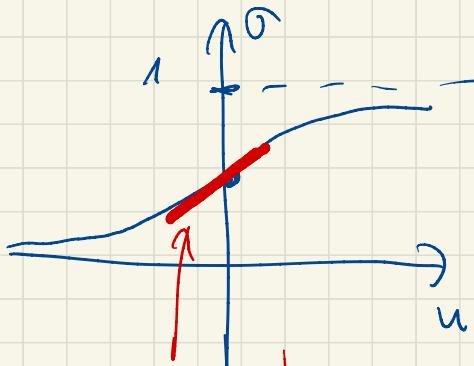
$$= -v^{-2} \times (-e^{-u})$$

$$= \left(\frac{1}{1 + e^{-u}} \right)^2 e^{-u} \quad \text{f}(1 - \frac{1}{1 + e^{-u}})$$

$$= \frac{1}{1 + e^{-u}} \cdot \frac{e^{-u}}{1 + e^{-u}}$$

$$= \sigma(u) (1 - \sigma(u))$$

$$0 < \sigma(u) < 1$$



Exercise

- Let's find the behavior of a sigmoid function $\sigma(x)$.
 - If we let x get larger ($x \rightarrow \infty$), explain the behavior of $\sigma(x)$.
 - If we let x get smaller ($x \rightarrow -\infty$), explain the behavior of $\sigma(x)$.
 - Write a graph of $\sigma(x)$ in the range $-5 < x < 5$.
- Calculate the derivative $\sigma'(x)$
 - Find the range of $\sigma'(x)$

Exercise

- Let the function

$$z(x) = \sigma(ax - 4) = \frac{1}{1+e^{-ax+4}}.$$

- Draw graphs of $f(x)$ for

- $a = 1$
- $a = 2$
- $a = 4$

and compare them.

Training of perceptrons

- A training set of an input vector $\{x_i\}$ and its expected output value z_i^* is used to determine weights $\{\omega_k\}$.
- adjust the weights $\{\omega_k\}$ in order to minimize the mean squared error R based on Gradient descent
(z_i is an output for the an input vector x_i , z_i^* is a correct output, ε is a constant.)

$$R = \frac{1}{2} \sum_i (z_i - z_i^*)^2$$

$$\delta \omega_i = -\varepsilon \frac{\partial R}{\partial \omega_i} = \varepsilon \sum_k (z_k^* - z_k) \frac{\partial z_k}{\partial \omega_i} = \varepsilon \sum_k (z_k^* - z_k) x_k^i \frac{\partial f(u_k)}{\partial u}$$

$$Z = \sigma(w \cdot x - \theta) \quad \text{-- perceptron}$$

$$\begin{matrix} d_1 & \tilde{z}_1 \\ d_2 & \tilde{z}_2 \\ \vdots & \vdots \\ d_n & \tilde{z}_n \end{matrix} \quad L = \frac{1}{2} \sum_i \left(\tilde{z}_i - \sigma(w \cdot \underline{x}_i - \underline{\theta}) \right)^2$$



$$w' = w - \varepsilon \frac{\partial L}{\partial w}$$

$$\theta' = \cancel{\theta} - \varepsilon' \frac{\partial L}{\partial \theta}$$

$$w \cdot x - \theta = w_1 d_1 + w_2 d_2 + \dots + \theta (-1)$$

$$= (w_1, w_2, \dots, \theta) \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ -1 \end{pmatrix}$$

$$= \hat{w} \cdot \hat{x}$$

$$L = \frac{1}{2} \sum_i (\hat{z}_i - \sigma(wx_i))^2$$

$$\frac{\partial L}{\partial w} = \sum_i (\hat{z}_i - \sigma(wx_i)) \left(-\frac{\partial}{\partial w} \sigma(wx_i) \right)$$

$$w' = w - \epsilon \boxed{\frac{\partial L}{\partial w}}$$
$$= \sum_i (\sigma(wx_i) - \hat{z}_i) \sigma(wx_i) (1 - \sigma(wx_i)) \times x_i$$

Exercise

- Let $z(x) = \frac{1}{1+e^{-wx+b}}$ and $R(w) = \frac{1}{2} \left(z(2) - \frac{1}{3} \right)^2$.
 - Calculate $\frac{\partial}{\partial w} z(2)$
 - Calculate $\frac{\partial}{\partial w} R(w)$

Exercise

- Let $w' = w - \epsilon \frac{\partial}{\partial w} R(w)$, where ϵ is a small and positive constant.
- Which is larger, $R(w)$ or $R(w')$?
 - Explain why?

Exercise

- Let $z(x) = \sigma(wx - b)$.
- Train the perceptron with the training data:

X	Z*
1.0	0.1
2.0	0.9

- Draw a graph of the perceptron.

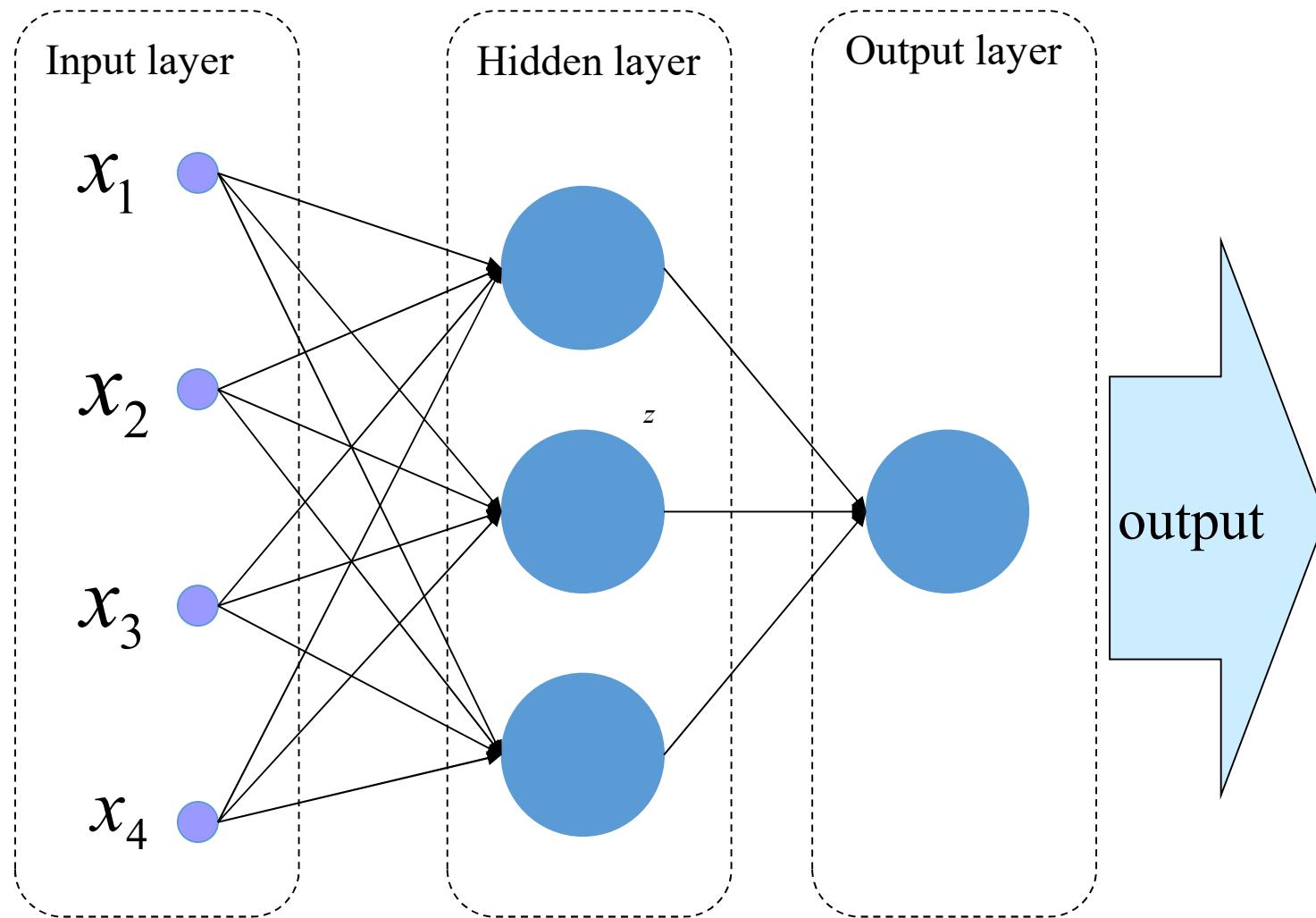
Exercise

- Let $z(x_1, x_2) = \sigma(w_1x_1 + w_2x_2 - b)$.
- Discuss whether we can make a perceptron which satisfies the following table:

x_1	x_2	z
0.0	0.0	0.0
1.0	0.0	1.0
0.0	1.0	1.0
1.0	1.0	0.0

- This is called XOR problem.

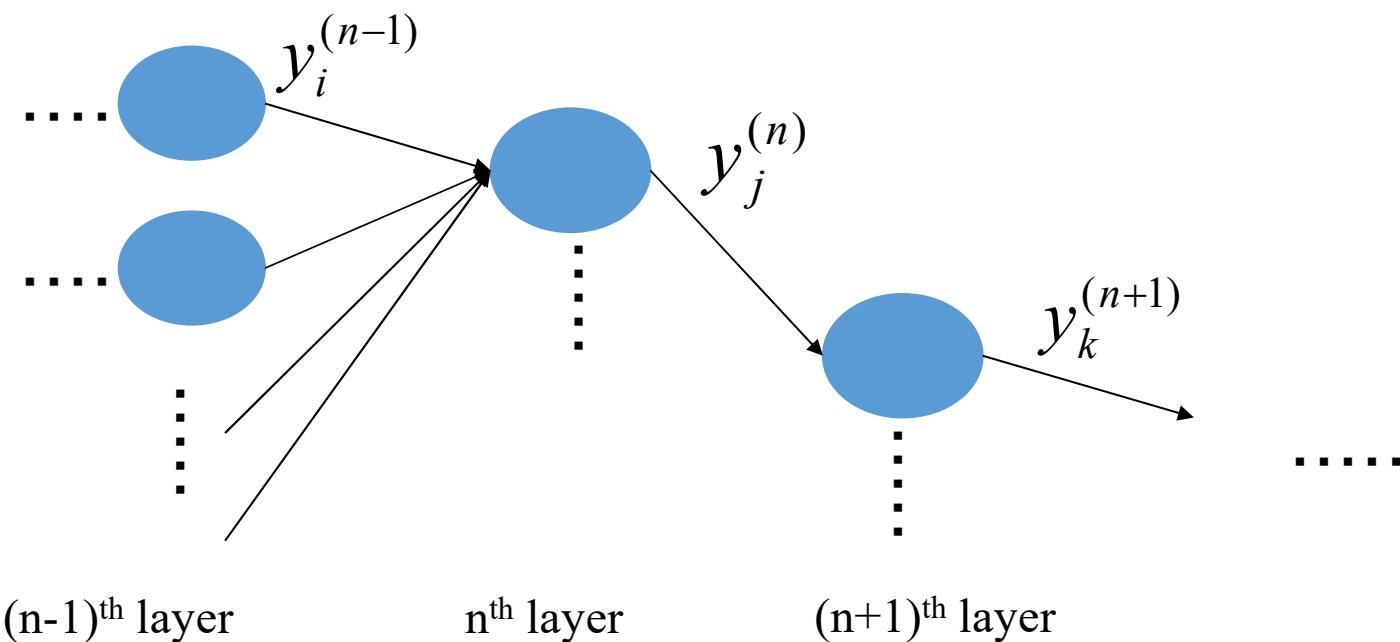
Multi layer perceptron



“Signal” propagation in NN

$$y_j^{(n)} = f\left(\sum_i \omega_{ji}^{(n)} y_i^{(n-1)} - \theta_j^{(n)}\right)$$

$$y_k^{(n+1)} = f\left(\sum_k \omega_{kj}^{(n+1)} y_j^{(n)} - \theta_k^{(n+1)}\right)$$



Exercise

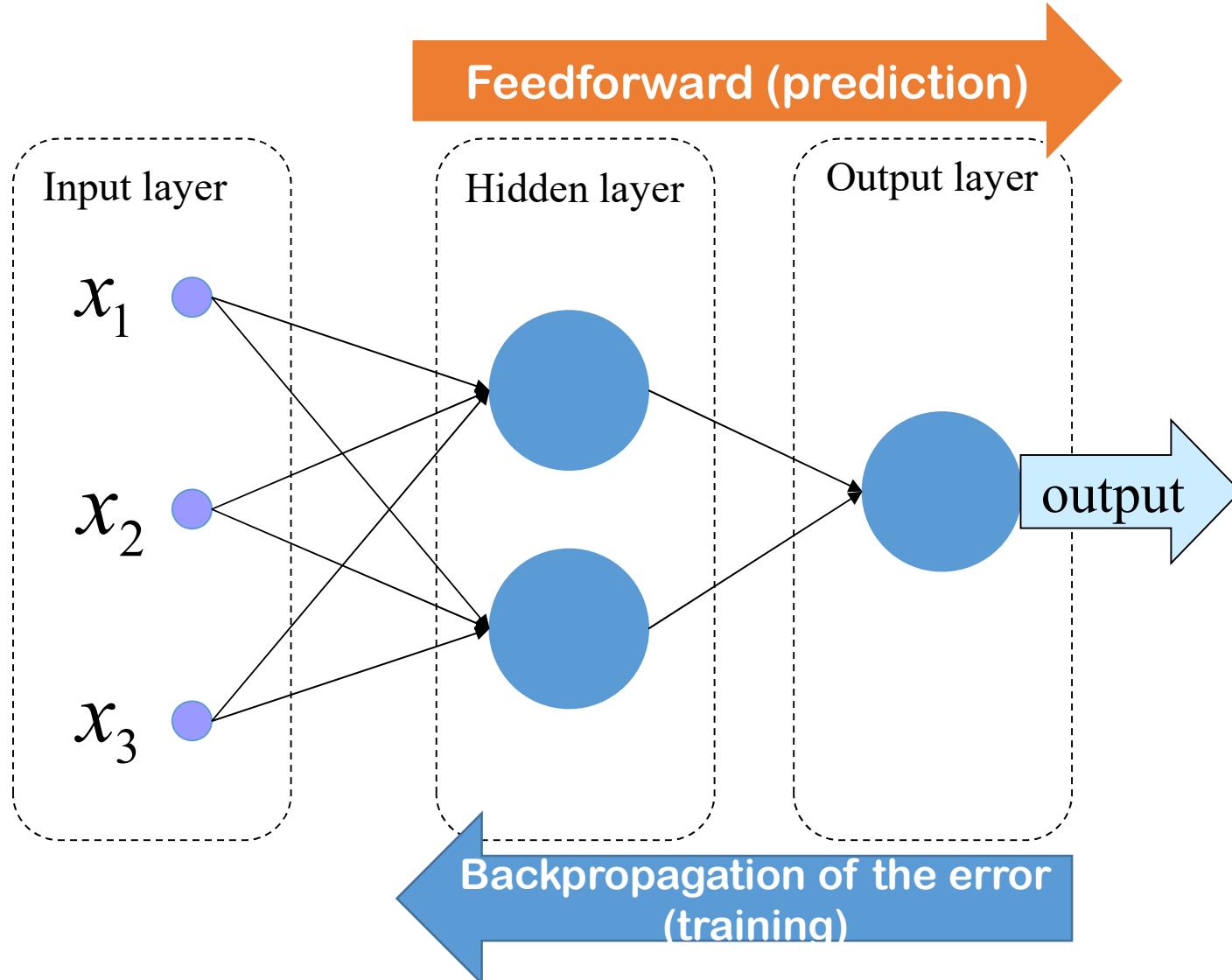
- Let

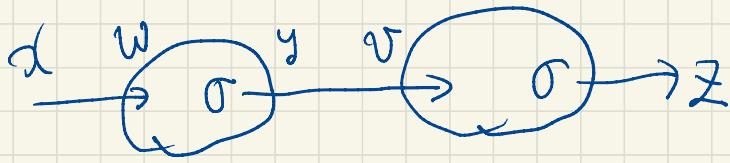
$$z(x_1, x_2) = \sigma \left(10\sigma \left(x_1 + x_2 - \frac{1}{2} \right) - 10\sigma \left(x_1 + x_2 - \frac{3}{2} \right) \right).$$

- Draw a graph of this function.
- Show this approximately solve XOR problem.

Gradient Vanishing
Problem

Feedforward and backpropagation





$$\begin{array}{l}
 x_1 \quad \tilde{x}_1 \\
 x_2 \quad \tilde{x}_2 \\
 \vdots \quad \vdots \\
 x_n \quad \tilde{x}_n
 \end{array} \quad L = \frac{1}{2} \sum_i \left(\tilde{x}_i - z_i \right)^2$$

$$\Downarrow = \frac{1}{2} \sum_i \left(\tilde{x}_i - \sigma(v\sigma(wx_i)) \right)^2$$

$$z = \sigma(vy) \quad y = \sigma(wx)$$

$$z_i = \sigma(v\sigma(wx_i))$$

$$\frac{\partial L}{\partial v} = \sum_i \left\{ \left(\sigma(v\sigma(wx_i)) - \tilde{x}_i \right) \times \sigma(v\sigma(wx_i)) \left(1 - \sigma(v\sigma(wx_i)) \right) \times \sigma(wx_i) \right\}$$

$$= \sum_i (z_i - \tilde{x}_i) z_i (1 - z_i) y_i$$

$$V' = V - \alpha \frac{\partial L}{\partial w}$$

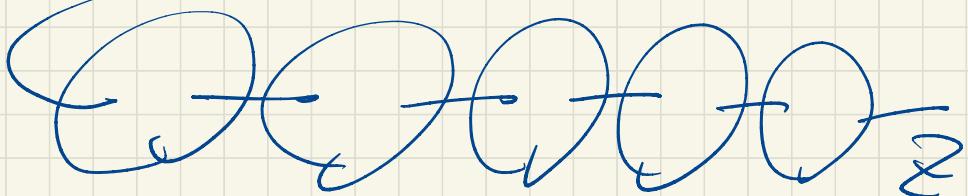
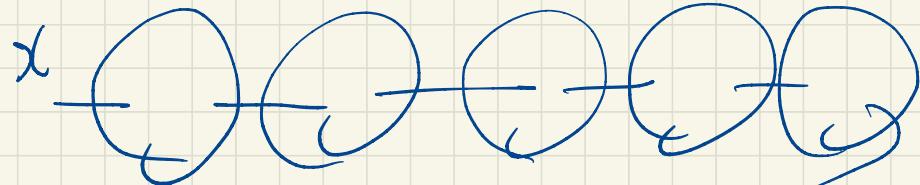
$$\frac{\partial L}{\partial w} = \frac{\partial}{\partial w} \frac{1}{2} \sum_i \left\{ \sigma(v \sigma(wx_i)) - \hat{z}_i \right\}^2$$

$$= \sum_i \left(\sigma(v \sigma(wx_i)) - \hat{z}_i \right) \times \frac{\partial \sigma(v \sigma(wx_i))}{\partial (v \sigma(wx_i))} \cdot \frac{\partial (v \sigma(wx_i))}{\partial w x_i} \frac{\partial w x_i}{\partial w}$$

$$= \sum_i (\hat{z}_i - \hat{z}_i) \underbrace{\sigma(vy_i)}_{\times v \sigma(wx_i)} \underbrace{(1-\sigma(vy_i))}_{(1-\sigma(wx_i))x_i}$$

$$= \sum_i (\hat{z}_i - \hat{z}_i) \hat{z}_i ((1-\hat{z}_i) \times \underbrace{v y_i}_{w} \underbrace{(1-y_i)}_{(1-x_i)} x_i)$$

$$\frac{\partial L}{\partial v} = \sum_i (\hat{z}_i - \hat{z}_i) \hat{z}_i ((1-\hat{z}_i) \underbrace{y_i}_{y_i})$$



$$\frac{1}{4}x \sim \frac{1}{4}x \frac{1}{4}x \overbrace{\frac{1}{4}x \frac{1}{4}x \frac{1}{4}}$$

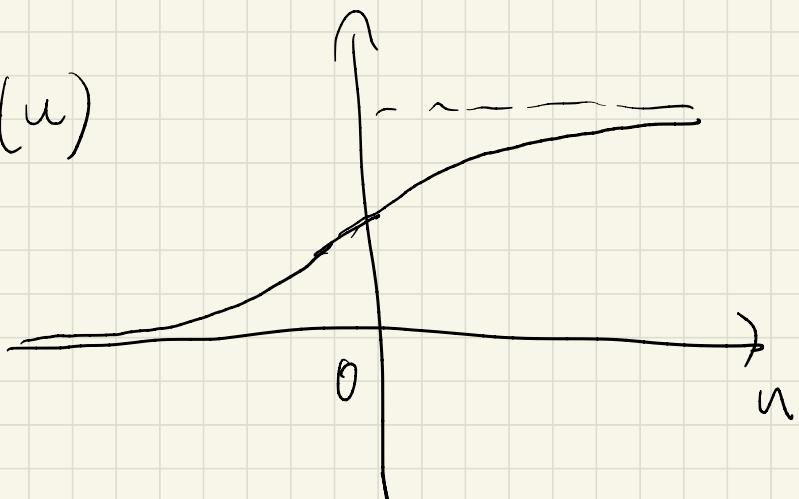
$$\left(\frac{1}{4}\right)^{10} = \left(\frac{1}{2}\right)^{20} = \frac{1}{2^{20}}$$

$$2^{10} = 1024 \approx 10^3$$

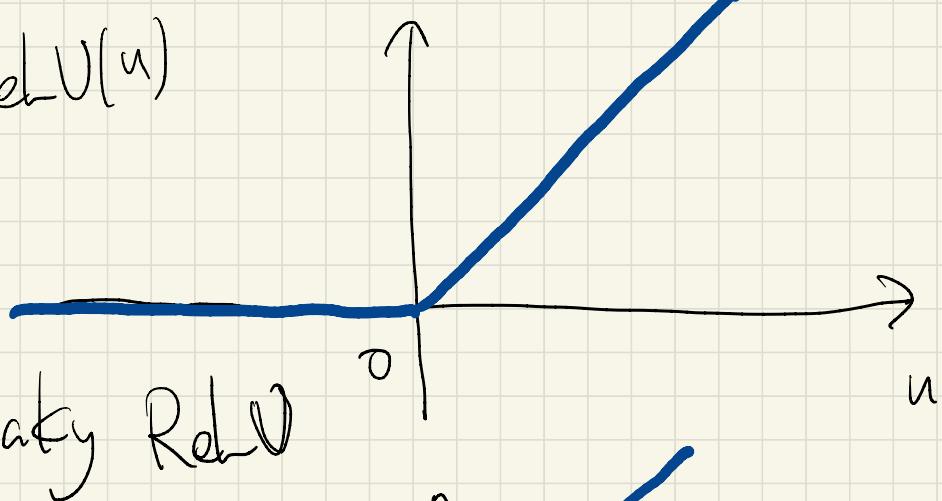
$$2^{20} \approx (10^3)^2 = 10^6$$

$$\approx \frac{1}{10^6}$$

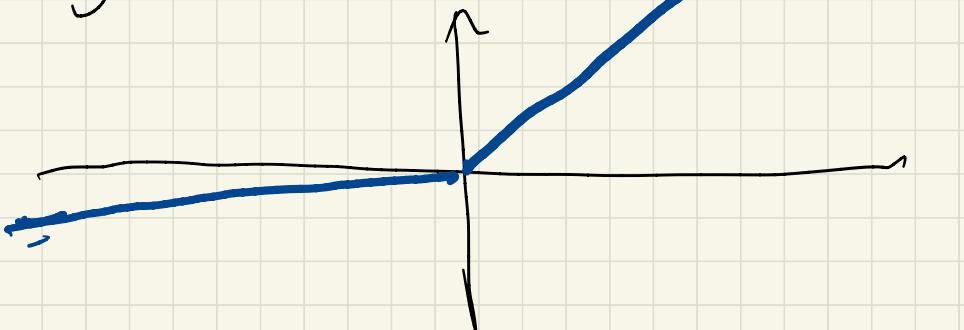
$\delta(u)$



$\text{ReLU}(u)$



Leaky ReLU



Back propagation

- adjusts the weights $\{\omega_k\}$ in order to minimize the mean squared error R based on Gradient descent

$$R = \frac{1}{2} \sum_k (z_k - z_k^*)^2 \quad \begin{aligned} y_j^{(n)} &= f(u_j^{(n)}) \\ u_j^{(n)} &= \sum_i \omega_{ji}^{(n)} y_i^{(n-1)} - \theta_j^{(n)} \end{aligned}$$

$$\omega'_{ji}^{(n)} = \omega_{ji}^{(n)} - \varepsilon \frac{\partial R}{\partial \omega_{ji}^{(n)}} = \omega_{ji}^{(n)} - \varepsilon \frac{\partial R}{\partial y_j^{(n)}} \frac{\partial y_j^{(n)}}{\partial u_j^{(n)}} \frac{\partial u_j^{(n)}}{\partial \omega_{ji}^{(n)}} = \omega_{ji}^{(n)} - \varepsilon O_j^{(n)} y_i^{(n-1)}$$

$$O_j^{(n)} = \frac{\partial R}{\partial y_j^{(n)}} \frac{\partial y_j^{(n)}}{\partial u_j^{(n)}} = \frac{\partial R}{\partial y_j^{(n)}} f(u_j^{(n)}) (1 - f(u_j^{(n)}))$$

Back propagation (cont'd)

$$O_i^{(n-1)} = \sum_j \frac{\partial R}{\partial y_j^{(n)}} \frac{\partial y_j^{(n)}}{\partial y_i^{(n-1)}} \frac{\partial y_i^{(n-1)}}{\partial u_i^{(n-1)}}$$

$$= \sum_j \frac{\partial R}{\partial y_j^{(n)}} \frac{\partial y_j^{(n)}}{\partial u_j^{(n)}} \boxed{\frac{\partial u_j^{(n)}}{\partial y_i^{(n-1)}}} \frac{\partial y_i^{(n-1)}}{\partial u_i^{(n-1)}}$$

$$= \sum_j O_j^{(n)} \omega_{ji}^{(n)} f(u_i^{(n-1)}) \{1 - f(u_i^{(n-1)})\}$$

where $y_i^{(n-1)} = f(u_i^{(n-1)})$ and $u_j^{(n)} = \sum_i \omega_{ji}^{(n)} y_i^{(n-1)} - \theta_j^{(n)}$

Exercise

- Let $z(x) = \sigma(v \sigma(wx - b) - d)$ and imagine we get training data (x_i, z_i) for $i = 1, 2, \dots, n$. Define a quadratic loss function as

$$R(v, w, b, d) = \frac{1}{2} \sum_i (z(x_i) - z_i)^2,$$

- Calculate
 - $\frac{\partial}{\partial v} z(x_i)$
 - $\frac{\partial}{\partial w} z(x_i)$ and express it with $\frac{\partial}{\partial v} z(x_i)$
 - $\frac{\partial}{\partial w} R$ and express it with $\frac{\partial}{\partial v} R$

Hint: Note that x_i and z_i should be dealt as a constant.

Exercise

- Update rules (**gradient descent**) of w and v are

$$w' = w - \epsilon \frac{\partial}{\partial w} R$$
$$v' = v - \epsilon \frac{\partial}{\partial v} R$$

- Which should be calculated first to reduce calculation time?

Input data

- Numerical data
 - are usually standardized as:

$$z = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

- Categorical data
 - are converted to a set of flags:

Male→(1,0), Female→(0,1), unknown→(0,0)

Examples

Estimation of fellowship grants

Organization	# of Division	# of Researchers	Total grants (actual)
团体A	29	3409	246374
团体B	499	15556	230978
团体C	73	17760	1020195
团体D	1041	133989	1367747
团体E	227	22266	188409
团体F	1961	134892	1717810



Normalization

$$z = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Organization	# of Division	# of Researchers	Total grants (actual)
团体A	0.000	0.000	0.038
团体B	0.243	0.092	0.028
团体C	0.023	0.109	0.544
团体D	0.524	0.993	0.771
团体E	0.102	0.143	0.000
团体F	1.000	1.000	1.000

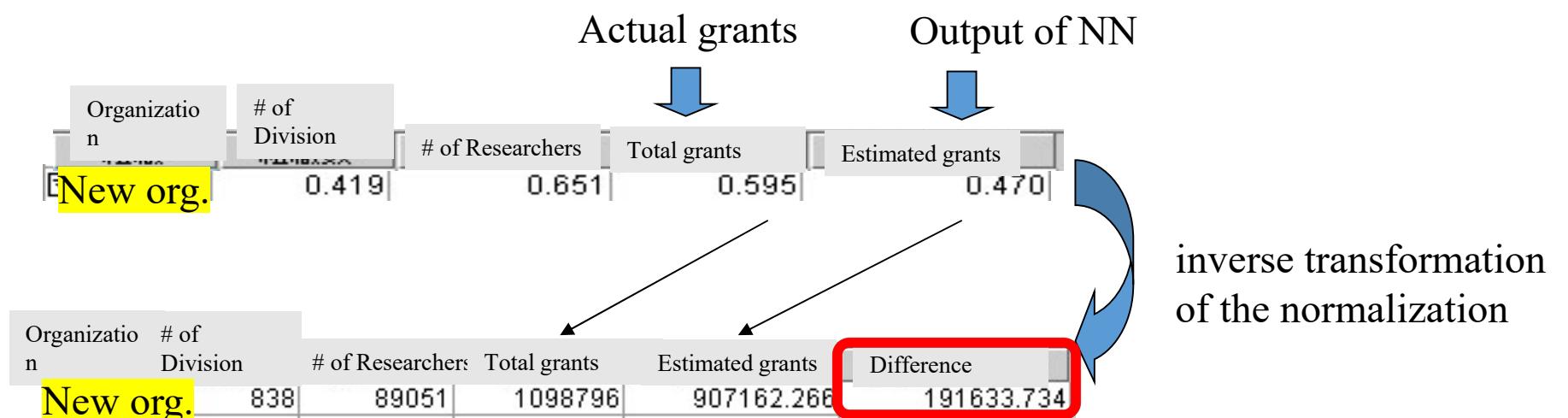
Estimation of fellowship grants (cont'd)

Estimated by NN

Organization	# of Division	# of Researchers	Total grants (actual)	Estimated grants
团体A	0.000	0.000	0.038	0.040
团体B	0.243	0.092	0.028	0.000
团体C	0.023	0.109	0.544	0.541
团体D	0.524	0.993	0.771	0.771
团体E	0.102	0.143	0.000	0.024
团体F	1.000	1.000	1.000	0.994

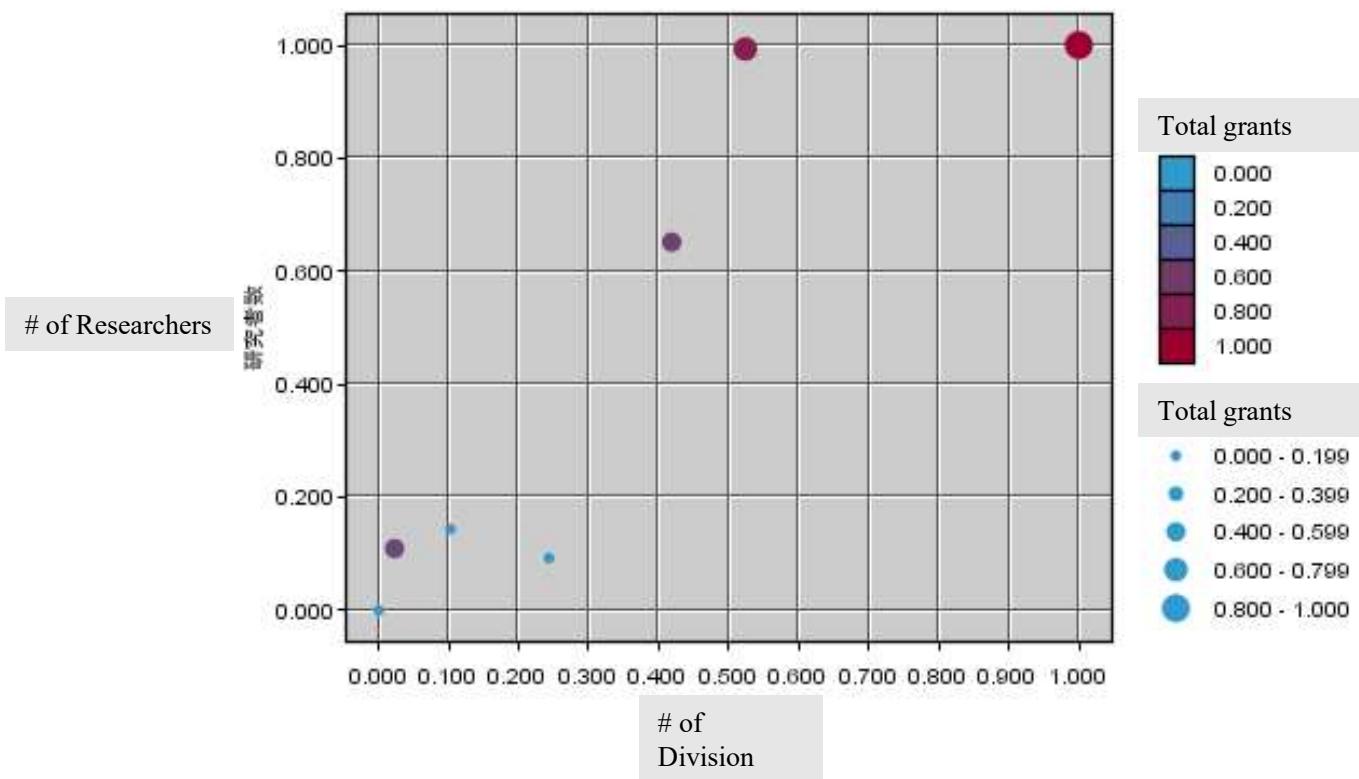
Estimation of fellowship grants (cont'd)

- Compare the estimated grants for a new organization with the actual grants



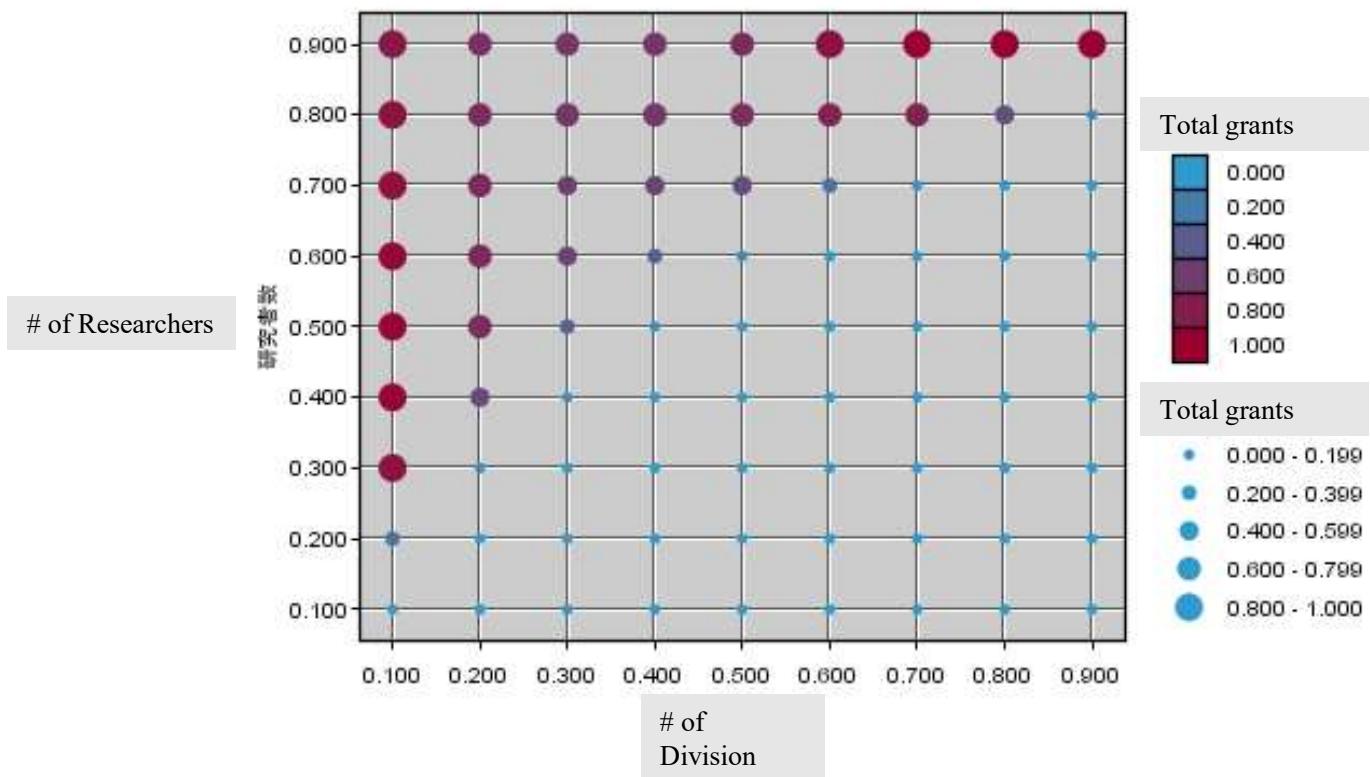
The result shows this organization receives more grants than other organizations

Distribution of the input data



The distribution of NN after learning

NN can be regarded as a non-linear interpolation function



Advantage and disadvantage

- Advantage
 - Ability to realize a highly non-linear complex function
 - Generally high precision of outputs
 - No restriction on input variables
- Disadvantage
 - Hard to understand the reason why the output is obtained
 - Long training time
 - Categorical data need to be converted to numerical data