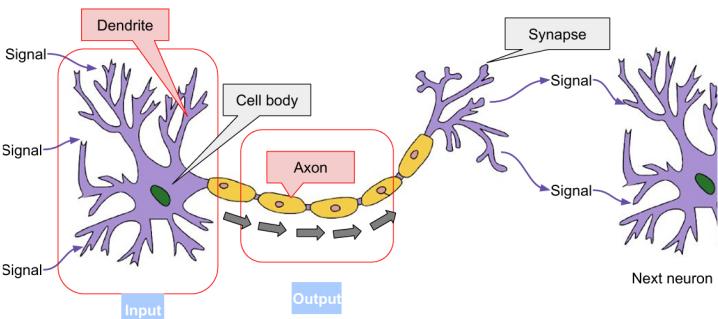


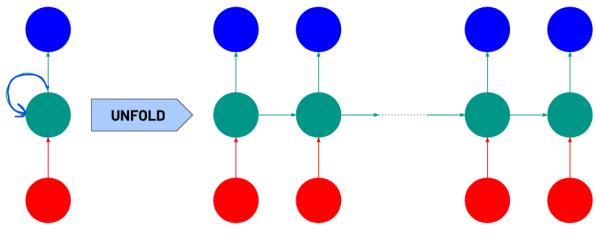
McCulloch-Pitts (MCP) neuron, in 1943

- McCulloch and Pitts described a nerve cell as a simple logic gate with binary outputs:
- multiple signals arrive at the dendrites,
 - they are then integrated into the cell body, and,
 - if the accumulated signal exceeds a certain threshold, an output signal is generated that will be passed on by the axon.

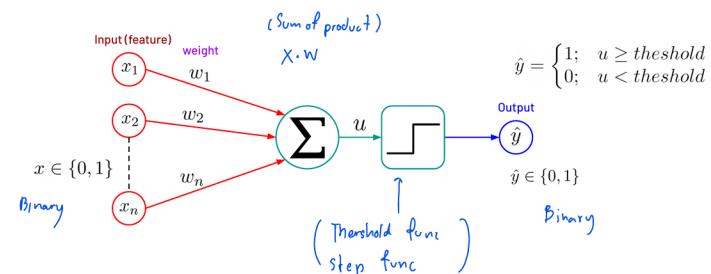
Biological neuron



Recurrent Neural Network



McCulloch-Pitts Neuron



• Two neural organs

• man Warren McCulloch (neuro-sci)
Walter Pitts (logician)

• McCulloch-Pitts neurons implement information processing in the brain

• Input / output functions

78%

1 is an logic gate or not

2 known as weight (w)

3 unknown \leq input instance

4 unknown threshold func

* return *

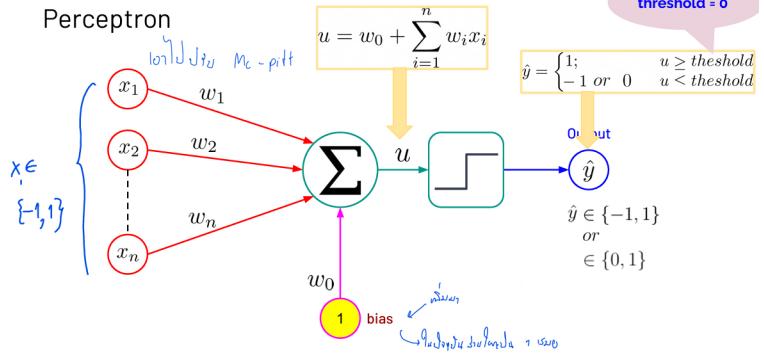
now set weight & threshold w_0 } model of 1st stage function
set $y = 1$ then

Linearly separable

Input Mac-pitts one time weight, threshold

$$f(x,y) = -y + mx + c$$

Perceptron



• Frank Rosenblatt based on McCulloch-Pitts

↳ Model \hat{y}

• weight simultaneous with training iteration with output y (convergence theory)

• learning rate (η -etc) $\eta \in (0, 1]$

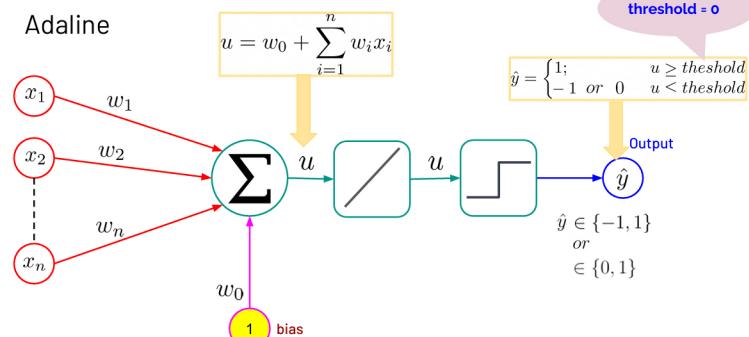
$$\Gamma w_{i+1} = w_{i+1} + \eta \cdot (y - \hat{y}) \cdot x_i$$

gradient

1 (sum) weight w_{i+1}
dinner η Learning rate $\eta \in (0, 1]$
dinner η w_{i+1}

2 \hat{y} is the instance
answer u without threshold $y \in \{-1, 1\}$
during weights

Adaline



• man Bernard Widrow & Ted Hoff

• Based on McCulloch-Pitts neurons
For output via only linear activation func ($f(x) = x$) instead of step func

• return Perceptron ↳

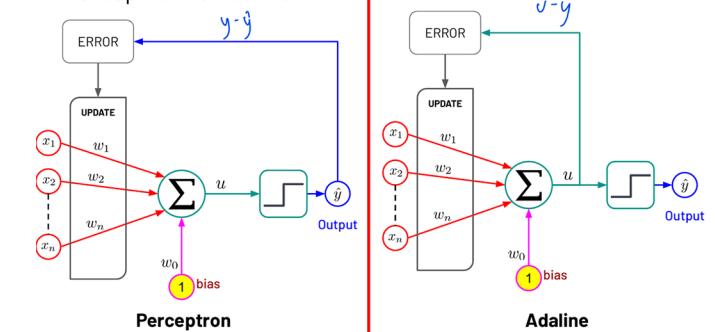
step func \hat{y} is used to predict y

89% works Perceptron

most difficult in u is \hat{y} which is not y for threshold func
↓ w weight

↳ what is the threshold func

Perceptron vs Adaline



Behind the Adaline Learning

Function Approximation

- We have a function $y = f(x, w)$ weight w sum an error J of function
- target function $y = \hat{y} = f(x, w)$ x : data w : weight
- $J = \|y - \hat{y}\|$ $\|\cdot\|$ denotes

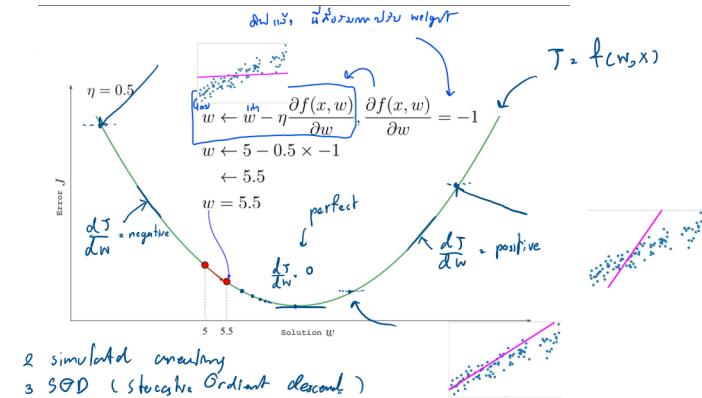
, Gradient-descent

$$w_i \leftarrow w_i - \eta \frac{\partial J}{\partial w_i} \quad \{ \text{minimizes vector } w_i \text{ to minimize sum } J \}$$

Gradient-descent

, Gradient-descent

$$w_i \leftarrow w_i - \eta \frac{\partial J}{\partial w_i} \quad \{ \text{minimizes vector } w_i \text{ to minimize sum } J \}$$



Delta Learning Rule

- gradient descent Widrow & Hoff (original Adaline)

↳ Delta Learning Rule, Widrow & Hoff Learning Rule, Least Mean Square (LMS) algorithm

- Given target vs the output of activation function (\hat{y}) → error
- To minimize error \hat{y} uses gradient descent shrinks weight

error function, loss function → MSE, RMSE, MAE etc

Adaline → Q: Delta Learning → Q: Gradient descent → derivative of error → derivative of activation function

Delta Learning Rule

- Given

$$E_j = y_j - \hat{y}_j \quad \text{is the error for the j-th sample}$$

$$\hat{y}_j = \phi \left(w_0 + \sum_{i=1}^n w_i x_{j,i} \right) \quad \text{is the output for the j-th sample}$$

$$J_j = E_j^2 \quad \text{is the squared error for the j-th sample}$$

The learning of delta rule use the Gradient Descent Learning of E_j

$$w_i \leftarrow w_i - \eta \frac{\partial J_j}{\partial w_i} \quad \text{on loss function}$$

Loss function is: Square Error

$$w_i \leftarrow w_i + \eta (y_j - \hat{y}_j) \frac{\partial \phi(u)}{\partial u} x_{j,i}$$

when η is the learning rate which its value is $\eta = (0, 1]$

$\frac{\partial \phi(u)}{\partial u}$ is the derivative of activation function used in the neural model

Compute the derivative

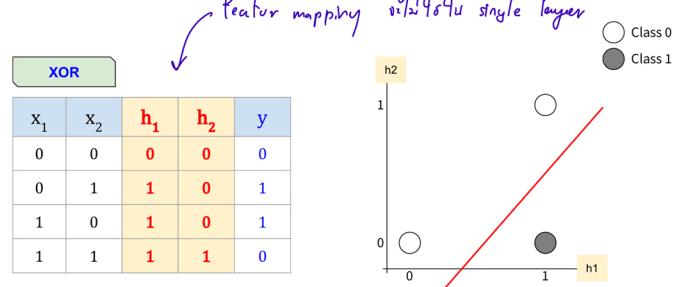
$$\begin{aligned} \phi(u) &= \frac{1}{1 + e^{-u}} \\ \frac{\partial \phi(u_j)}{\partial u_j} &= \phi(u_j)(1 - \phi(u_j)) \\ \frac{\partial J_j}{\partial w_i} &= -(y_j - \hat{y}_j) \frac{\partial \phi(u)}{\partial u} x_{j,i} \end{aligned}$$

Weight adjusting

$$\begin{aligned} w_i &\leftarrow w_i - \eta \frac{\partial J_j}{\partial w_i} \\ &\leftarrow w_i + \eta (y_j - \hat{y}_j) \frac{\partial \phi(u)}{\partial u} x_{j,i} \\ &\leftarrow w_i + \eta (y_j - \hat{y}_j) \phi(u_j)(1 - \phi(u_j)) x_{j,i} \end{aligned}$$

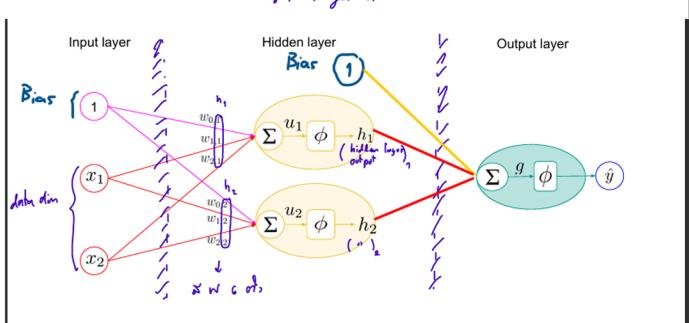
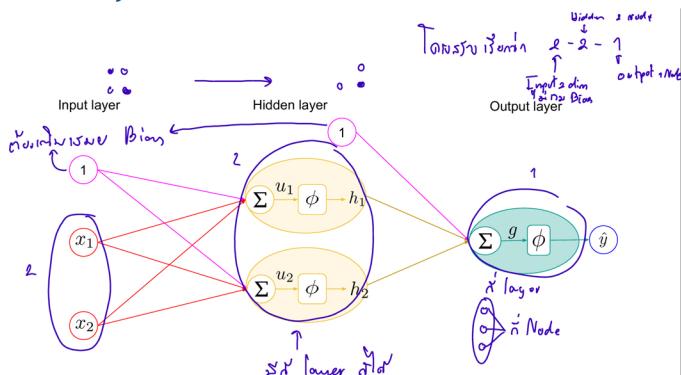
Multi-Layer Perceptron (MLP)

Map the original space (x_1, x_2) to the new space (h_1, h_2)



1 Target vs. 1/2/2 2 - 2 - 1

on. Node hidden layer
↓
2 - 2 - 1
↑
Input Output 1 Node
2 dimension
1/2/2 Bias



→ Feed-forward (compute output)
← Backpropagation (weight tuning)

→ Error stückweise output → Layer → Layer

Step 1 : train Parameter

- define Network architecture
 - input hidden layer
 - input hidden node
- define Weight/Bias
- define activation function
- define learning rate
- define epoch

Step 2 :

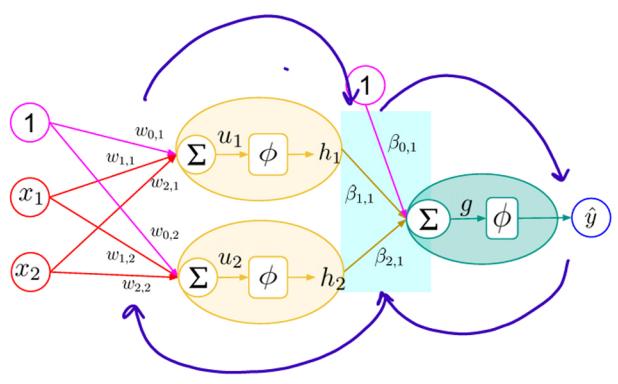
Feed - Forward

- 2.1 forward output of hidden layer h_1
- 2.2 reverse output of output layer c_y - hats

Back propagation

- 2.3 JWD weight output layer
- 2.4 JWD weight input layer

Feed forward



Back propagation

Recall:

$$E(t) = \frac{1}{2} e^2(t)$$

$$e = y - \hat{y}$$

$$\hat{y} = \phi(g)$$

$$g = \beta_{0,1} + \sum_{k=1}^l \beta_{k,1} h_k$$

$$u_1 = w_{0,1} + \sum_{j=1}^n w_{j,1} x_j$$

$$u_2 = w_{0,2} + \sum_{j=1}^n w_{j,2} x_j$$

$$h_1 = \phi(u_1)$$

$$h_2 = \phi(u_2)$$

$$\beta = \beta + \eta e(\hat{y}(1 - \hat{y}))h$$

$$w = w + \eta e(\hat{y}(1 - \hat{y}))\beta(h(1 - h))x$$

Step 2: Train model

- For each data point (x)

Feed-forward

- Step 2.1: Compute outputs of hidden layer
- Step 2.2: Compute outputs of output layer

Backpropagation

- Step 2.3: Adjust the weights of output layer
- Step 2.4: Adjust the weights of input (hidden) layer

$$u_1 = w_{0,1} + \sum_{j=1}^n w_{j,1} x_j$$

$$u_2 = w_{0,2} + \sum_{j=1}^n w_{j,2} x_j$$

$$h_1 = \phi(u_1)$$

$$h_2 = \phi(u_2)$$

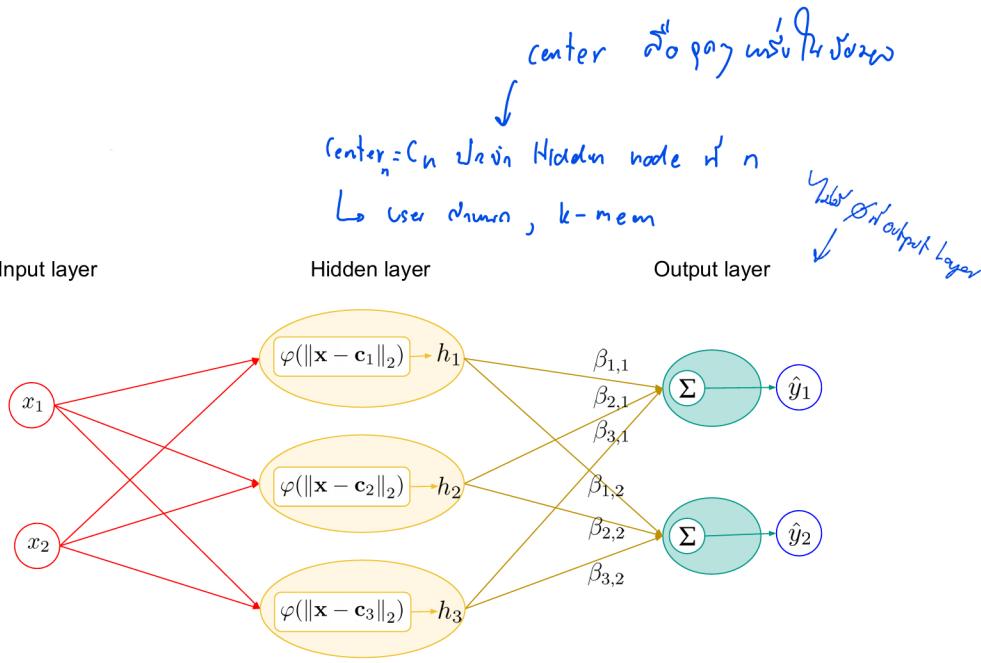
Activation function:

$$\text{Sigmoid} \quad \hat{y} = \phi(g) = \frac{1}{1 + e^{-u(t)}}$$

$$\beta = \beta - \eta \nabla_\beta E$$

$$w = w - \eta \nabla_w E$$

Radial Basis Function (RBF) Neural Network

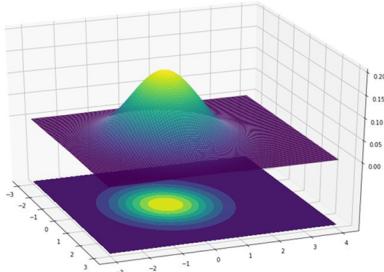


The **Gaussian Functions** are generally used for Radial Basis Function.

Euclidean distance

$$\text{Given: } r = \|\mathbf{x} - \mathbf{c}\|_2$$

$$\varphi(r) = \exp\left(-\frac{r^2}{2\sigma^2}\right)$$



Gaussian RBF:

$$\varphi(\|\mathbf{x} - \mathbf{c}\|_2) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}\|_2^2}{\sigma^2}\right)$$

1. Hidden layer

From center vector To hidden node

- no training set
- no k-means clustering in

2. Output weight tuning

- gradient descent

$$J = \frac{1}{2} (y - \hat{y})^2 = \text{square error}$$

$$\begin{aligned} \beta_{l,m} &\leftarrow \beta_{l,m} - \eta \frac{\partial J_i}{\partial \beta_{l,m}} \\ &\leftarrow \beta_{l,m} + \eta (y_{i,m} - \hat{y}_{i,m}) h_{i,l} \end{aligned}$$

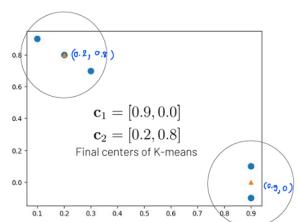
- generalized pseudo inverse

$$\boldsymbol{\beta} = (\mathbf{H}^\top \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{Y}$$

- 1st output layer, 1 Hidden Layer, output layer
- 1st input weight
- activation function φ Radial-Basis function
- Gaussian Node of Hidden node within "prototype" \Rightarrow "center"
- tuning output by LMS, GD
- Generalized pseudo inverse

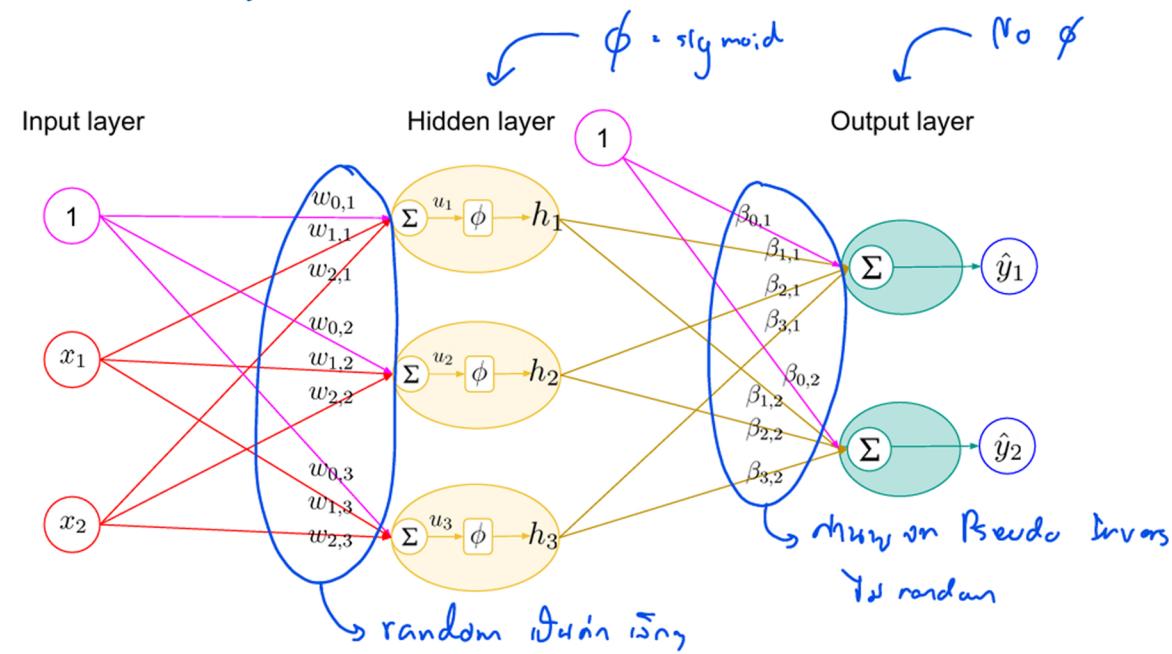
Example

x_1	x_2	y
0.1	0.9	-1
0.2	0.8	-1
0.3	0.7	-1
0.9	0.1	1
0.9	-0.1	1



Random Neural Network

- 1 hidden layer ไม่มี Backprob แต่ทำ Pseudo Inverse แทน
- จด input weight
- Out weight หาด้วย Pseudo inverse
- Fast learning model (เรียนเร็ว)



data

$$\mathbf{X} \in \underset{N}{\underbrace{\left[\begin{array}{c|cccc} \text{bias + data} & \\ \hline 1 & x_1 & \dots & x_N \\ 1 & x_1 & \dots & x_N \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1 & \dots & x_N \end{array} \right]}}_{(d+1)} \quad N \times (d+1)$$

Single Hidden Node = L

$$W \in \underset{(d+1) \times L}{\underbrace{\left[\begin{array}{c|cccc} 1 & 2 & \dots & L \\ \hbar & \hbar & \dots & \hbar \\ \hbar & \hbar & \dots & \hbar \\ \vdots & \vdots & \ddots & \vdots \\ \hbar & \hbar & \dots & \hbar \end{array} \right]}}$$

$$H = \phi(X \cdot W) = \underset{N \times (L+1)}{\underbrace{\left[\begin{array}{c|cccc} 1 & h_1 & h_2 & \dots & h_L \\ \hbar & \hbar & \hbar & \dots & \hbar \\ \hbar & \hbar & \hbar & \dots & \hbar \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \hbar & \hbar & \hbar & \dots & \hbar \end{array} \right]}}$$

Output Weight

Output Node = m

$$\underset{(L+1) \times M}{\underbrace{\left[\begin{array}{c|cccc} 1 & 2 & \dots & m \\ \beta & \beta & \dots & \beta \\ \beta & \beta & \dots & \beta \\ \vdots & \vdots & \ddots & \vdots \\ \beta & \beta & \dots & \beta \end{array} \right]}}$$

$$\hat{Y} = H \cdot \beta$$

$$\underset{N \times M}{\underbrace{\left[\begin{array}{c|cccc} y_1 & y_2 & \dots & y_m \\ y_1 & y_2 & \dots & y_m \\ y_1 & y_2 & \dots & y_m \\ \vdots & \vdots & \ddots & \vdots \\ y_1 & y_2 & \dots & y_m \end{array} \right]}}$$

$$\boxed{\mathbf{H}\beta = \mathbf{Y}}$$

$$\mathbf{H}^T \mathbf{H}\beta = \mathbf{H}^T \mathbf{Y}$$

$$(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{H}\beta = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{Y}$$

$$\mathbf{I}\beta = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{Y}$$

`numpy.linalg.pinv()` $\beta = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{Y}$ ←

Deep Learning Timeline

