# Machine Learning Foundations
## (機器學習基石)

Lecture 9: Linear Regression

### Hsuan-Tien Lin (林軒田)

htlin@csie.ntu.edu.tw

Department of Computer Science
& Information Engineering

National Taiwan University
(國立台灣大學資訊工程系)

# Credit **Limit** Problem

| age | 23 years |
|---|---|
| gender | female |
| annual salary | NTD 1,000,000 |
| year in residence | 1 year |
| year in job | 0.5 year |
| current debt | 200,000 |

credit limit? **100,000**

unknown target function
$f: \mathcal{X} \to \mathcal{Y}$

*(ideal credit **limit** formula)*

training examples
$\mathcal{D}: (\mathbf{x}_1, y_1), \cdots, (\mathbf{x}_N, y_N)$

*(historical records in bank)*

learning
algorithm
$\mathcal{A}$

final hypothesis
$g \approx f$

*('learned' formula to be used)*

hypothesis set
$\mathcal{H}$

*(set of candidate formula)*

$\mathcal{Y} = \mathbb{R}$: **regression**

# Linear Regression Hypothesis

| age | 23 years |
|---|---|
| annual salary | NTD 1,000,000 |
| year in job | 0.5 year |
| current debt | 200,000 |

- For $\mathbf{x} = (x_0, x_1, x_2, \cdots, x_d)$ 'features of customer', approximate the desired credit limit with a weighted sum:
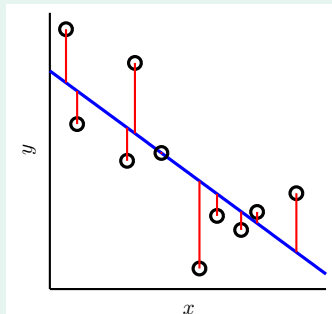
$$y \approx \sum_{i=0}^{d} w_i x_i$$

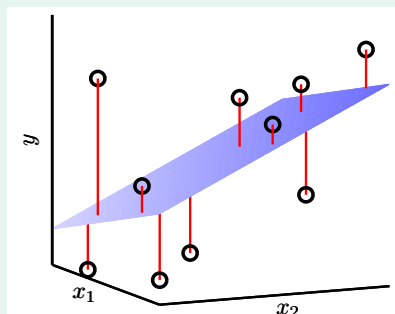- linear regression hypothesis: $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$

$h(\mathbf{x})$: like **perceptron**, but without the sign

# Illustration of Linear Regression

$\mathbf{x} = (x) \in \mathbb{R}$



$\mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$



linear regression:
find lines/hyperplanes with small residuals

# The Error Measure

popular/historical error measure:

$$\text{err}(\hat{y}, y) = (\hat{y} - y)^2$$

### in-sample

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^{N} (\underbrace{h(\mathbf{x}_n)}_{\mathbf{w}^T \mathbf{x}_n} - y_n)^2$$

### out-of-sample

$$E_{\text{out}}(\mathbf{w}) = \underset{(\mathbf{x}, y) \sim P}{\mathcal{E}} (\mathbf{w}^T \mathbf{x} - y)^2$$

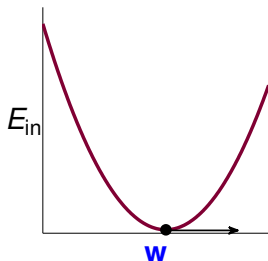next: how to minimize $E_{\text{in}}(\mathbf{w})$?

# Fun Time

Consider using linear regression hypothesis $h(\mathbf{x}) = \mathbf{w}^T\mathbf{x}$ to predict the credit limit of customers $\mathbf{x}$. Which feature below shall have a positive weight in a **good hypothesis** for the task?

1. birth month
2. monthly income
3. current debt
4. number of credit cards owned

# Matrix Form of $E_{\text{in}}(\mathbf{w})$

$$
\begin{aligned}
E_{\text{in}}(\mathbf{w}) &= \frac{1}{N}\sum_{n=1}^{N}(\mathbf{w}^T\mathbf{x}_n - y_n)^2 = \frac{1}{N}\sum_{n=1}^{N}(\mathbf{x}_n^T\mathbf{w} - y_n)^2 \\
&= \frac{1}{N}\left\|\begin{array}{c} \mathbf{x}_1^T\mathbf{w} - y_1 \\ \mathbf{x}_2^T\mathbf{w} - y_2 \\ \dots \\ \mathbf{x}_N^T\mathbf{w} - y_N \end{array}\right\|^2 \\
&= \frac{1}{N}\left\|\left[\begin{array}{c} --\mathbf{x}_1^T-- \\ --\mathbf{x}_2^T-- \\ \dots \\ --\mathbf{x}_N^T-- \end{array}\right]\mathbf{w} - \left[\begin{array}{c} y_1 \\ y_2 \\ \dots \\ y_N \end{array}\right]\right\|^2 \\
&= \frac{1}{N}\|\underbrace{\mathrm{X}}_{N\times d+1}\ \underbrace{\mathbf{w}}_{d+1\times 1} - \underbrace{\mathbf{y}}_{N\times 1}\|^2
\end{aligned}
$$

$$\min_{\mathbf{w}} E_{\text{in}}(\mathbf{w}) = \frac{1}{N}\|\mathrm{X}\mathbf{w} - \mathbf{y}\|^2$$



- $E_{\text{in}}(\mathbf{w})$: continuous, differentiable, **convex**
- necessary condition of 'best' $\mathbf{w}$

$$\nabla E_{\text{in}}(\mathbf{w}) \equiv \begin{bmatrix} \frac{\partial E_{\text{in}}}{\partial w_0}(\mathbf{w}) \\ \frac{\partial E_{\text{in}}}{\partial w_1}(\mathbf{w}) \\ \dots \\ \frac{\partial E_{\text{in}}}{\partial w_d}(\mathbf{w}) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

—not possible to 'roll down'

task: find $\mathbf{w}_{\text{LIN}}$ such that $\nabla E_{\text{in}}(\mathbf{w}_{\text{LIN}}) = \mathbf{0}$

# The Gradient $\nabla E_{\text{in}}(\mathbf{w})$

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N}\|\mathrm{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{N}\left(\underbrace{\mathbf{w}^T\mathrm{X}^T\mathrm{X}\mathbf{w}}_{\mathrm{A}} - \underbrace{2\mathbf{w}^T\mathrm{X}^T\mathbf{y}}_{\mathbf{b}} + \underbrace{\mathbf{y}^T\mathbf{y}}_{c}\right)$$

**one $w$ only**

$$E_{\text{in}}(w) = \frac{1}{N}\left(aw^2 - 2bw + c\right)$$

$$\nabla E_{\text{in}}(w) = \frac{1}{N}\left(2aw - 2b\right)$$

**simple! :-)**

**vector $\mathbf{w}$**

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N}\left(\mathbf{w}^T\mathrm{A}\mathbf{w} - 2\mathbf{w}^T\mathbf{b} + c\right)$$

$$\nabla E_{\text{in}}(\mathbf{w}) = \frac{1}{N}\left(2\mathrm{A}\mathbf{w} - 2\mathbf{b}\right)$$

similar (**derived by definition**)

$$\nabla E_{\text{in}}(\mathbf{w}) = \frac{2}{N}\left(\mathrm{X}^T\mathrm{X}\mathbf{w} - \mathrm{X}^T\mathbf{y}\right)$$

# Optimal Linear Regression Weights

task: find $\mathbf{w}_{\text{LIN}}$ such that $\frac{2}{N}\left(X^TX\mathbf{w} - X^T\mathbf{y}\right) = \nabla E_{\text{in}}(\mathbf{w}) = \mathbf{0}$

### invertible $X^TX$

- **easy!** unique solution

$$\mathbf{w}_{\text{LIN}} = \underbrace{\left(X^TX\right)^{-1}X^T}_{\text{pseudo-inverse } X^\dagger}\ \mathbf{y}$$

- often the case because $N \gg d + 1$

### singular $X^TX$

- **many** optimal solutions
- one of the solutions

$$\mathbf{w}_{\text{LIN}} = X^\dagger\mathbf{y}$$

by defining $X^\dagger$ in other ways

practical suggestion:
use **well-implemented** † **routine**
instead of $\left(X^TX\right)^{-1}X^T$
for numerical stability when **almost-singular**

# Linear Regression Algorithm

1. from $\mathcal{D}$, construct input matrix $X$ and output vector $\mathbf{y}$ by

$$X = \underbrace{\begin{bmatrix} --\mathbf{x}_1^T-- \\ --\mathbf{x}_2^T-- \\ \cdots \\ --\mathbf{x}_N^T-- \end{bmatrix}}_{N \times (d+1)} \quad \mathbf{y} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \cdots \\ y_N \end{bmatrix}}_{N \times 1}$$

2. calculate pseudo-inverse $\underbrace{X^{\dagger}}_{(d+1) \times N}$

3. return $\underbrace{\mathbf{w}_{\text{LIN}}}_{(d+1) \times 1} = X^{\dagger}\mathbf{y}$

simple and efficient
with **good** $\dagger$ **routine**

# Is Linear Regression a 'Learning Algorithm'?

$$\mathbf{w}_{\text{LIN}} = \mathbf{X}^{\dagger}\mathbf{y}$$

## No!

- analytic (**closed-form**) solution, 'instantaneous'
- not improving $E_{\text{in}}$ nor $E_{\text{out}}$ iteratively

## Yes!

- good $E_{\text{in}}$?
  **yes, optimal!**
- good $E_{\text{out}}$?
  **yes, finite $d_{\text{VC}}$ like perceptrons**
- improving iteratively?
  **somewhat, within an iterative pseudo-inverse routine**

if $E_{\text{out}}(\mathbf{w}_{\text{LIN}})$ is good, **learning 'happened'**!