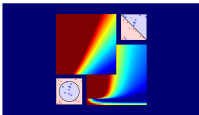


Machine Learning Foundations

(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)



Roadmap

- 1 When Can Machines Learn?
- 2 Why Can Machines Learn?

Lecture 8: Noise and Error

learning can happen
with **target distribution** $P(y|\mathbf{x})$ and **low** E_{in} **w.r.t. err**

- 3 **How** Can Machines Learn?

Lecture 9: Linear Regression

- Linear Regression Problem
- Linear Regression Algorithm
- Generalization Issue
- Linear Regression for Binary Classification

- 4 How Can Machines Learn Better?

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

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

(ideal credit **limit** formula)

training examples

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

(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)

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

Linear Regression Hypothesis

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

- For $\mathbf{x} = (x_0, x_1, x_2, \dots, x_d)$ 'features of customer', approximate the **desired credit limit** with a **weighted** sum:

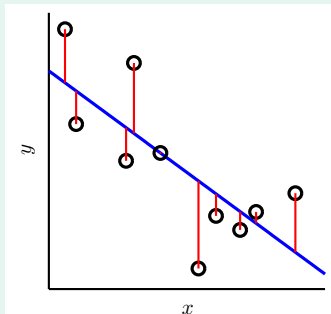
$$y \approx \sum_{i=0}^d w_i x_i$$

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

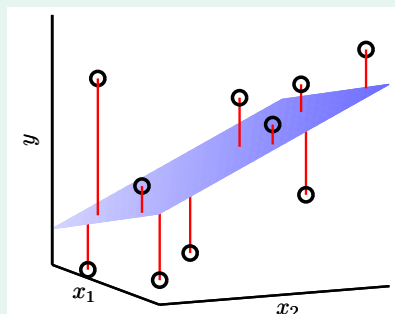
$h(\mathbf{x})$: like **perceptron**, but without the **sign**

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:

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

in-sample

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

out-of-sample

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

有noise的意思是x和y都是从某个分布中抽出来的，其中y的分布和抽取出来的x有关

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

Reference Answer: ②

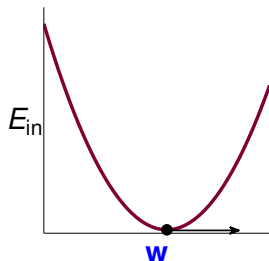
Customers with higher monthly income should naturally be given a higher credit limit, which is captured by the positive weight on the 'monthly income' feature.

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

w实际上就是那个系数

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

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



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

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

—not possible to 'roll down'

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

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

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{N} \left(\underbrace{\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}}_A - 2 \underbrace{\mathbf{w}^T \mathbf{X}^T \mathbf{y}}_b + \underbrace{\mathbf{y}^T \mathbf{y}}_c \right)$$

one w only

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

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

simple! :-)

vector \mathbf{w}

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

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

similar (**derived by definition**)

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

Optimal Linear Regression Weights

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

invertible $\mathbf{X}^T \mathbf{X}$ 如果可以求逆

- easy!** unique solution

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

$\mathbf{X}: N \times d+1$

- often the case because

$$N \gg d + 1$$

为何通常可逆???

singular $\mathbf{X}^T \mathbf{X}$ 如果不可逆

- many** optimal solutions
- one of the solutions

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

by defining \mathbf{X}^\dagger in other ways

practical suggestion:

use **well-implemented** \dagger **routine**

instead of $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

for numerical stability when **almost-singular**

Linear Regression Algorithm

- 1 from \mathcal{D} , construct **input matrix X** and **output vector y** by

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

- 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**

Fun Time

After getting \mathbf{w}_{LIN} , we can calculate the predictions $\hat{y}_n = \mathbf{w}_{\text{LIN}}^T \mathbf{x}_n$. If all \hat{y}_n are collected in a vector $\hat{\mathbf{y}}$ similar to how we form \mathbf{y} , what is the matrix formula of $\hat{\mathbf{y}}$?

- ① \mathbf{y}
- ② $\mathbf{X}\mathbf{X}^T \mathbf{y}$
- ③ $\mathbf{X}\mathbf{X}^\dagger \mathbf{y}$
- ④ $\mathbf{X}\mathbf{X}^\dagger \mathbf{X}\mathbf{X}^T \mathbf{y}$

Reference Answer: ③

Note that $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}_{\text{LIN}}$. Then, a simple substitution of \mathbf{w}_{LIN} reveals the answer.

Is Linear Regression a ‘Learning Algorithm’?

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

No!

- analytic (**closed-form**) solution, ‘instantaneous’
- not improving E_{in} nor E_{out} iteratively

Yes!

- good E_{in} ?
yes, optimal!
- good E_{out} ? 只有d+1个变数? 可操作的量
yes, finite d_{VC} like perceptrons
- improving iteratively?
somewhat, within an iterative pseudo-inverse routine

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


刚才给胡的是 E_{in} 最小化

Benefit of Analytic Solution: 'Simpler-than-VC' Guarantee

为什么取平均？

$$\overline{E_{\text{in}}} = \mathcal{E}_{\mathcal{D} \sim P^N} \left\{ E_{\text{in}}(\mathbf{w}_{\text{LIN}} \text{ w.r.t. } \mathcal{D}) \right\} \stackrel{\text{to be shown}}{=} \text{noise level} \cdot \left(1 - \frac{d+1}{N}\right)$$

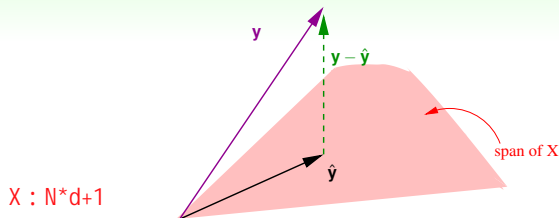
一定是错误的量



$$\begin{aligned} E_{\text{in}}(\mathbf{w}_{\text{LIN}}) &= \frac{1}{N} \|\mathbf{y} - \underbrace{\hat{\mathbf{y}}}_{\text{predictions}}\|^2 = \frac{1}{N} \|\mathbf{y} - \mathbf{X} \underbrace{\mathbf{X}^\dagger \mathbf{y}}_{\mathbf{w}_{\text{LIN}}}\|^2 \\ &= \frac{1}{N} \|(\underbrace{\mathbf{I}}_{\text{identity}} - \mathbf{X} \mathbf{X}^\dagger) \mathbf{y}\|^2 \end{aligned}$$

call $\mathbf{X} \mathbf{X}^\dagger$ the **hat matrix** \mathbf{H}
because it **puts \wedge on \mathbf{y}**

Geometric View of Hat Matrix



in \mathbb{R}^N

这里的w和X的列乘开后的系数k, w实际上是任意的, 我们在里头找一个最吼的w

- $\hat{y} = X\mathbf{w}_{\text{LIN}}$ within the span of X columns
- $\mathbf{y} - \hat{y}$ smallest: $\mathbf{y} - \hat{y} \perp \text{span}$
- H : project \mathbf{y} to $\hat{y} \in \text{span}$
- $I - H$: transform \mathbf{y} to $\mathbf{y} - \hat{y} \perp \text{span}$

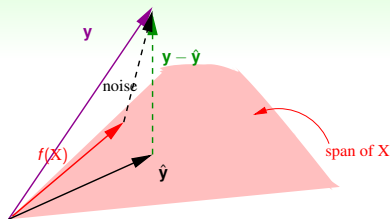
$H = XX^+$ 是个投影矩阵, 线性变换

$$\text{trace}(AB) = \text{trace}(BA)$$

claim: $\text{trace}(I - H) = N - (d + 1)$. Why? :-)

N维向量投影到一个最多d+1维的空间X的span, 自由度减少d+1 (物理意义)

An Illustrative 'Proof'



- if **y** comes from some ideal $f(X) \in \text{span}$ plus **noise**
- noise** transformed by $I - H$ to be $y - \hat{y}$

都是投影，想象noise是一个新的y那就好理解了

$$\begin{aligned}
 E_{\text{in}}(\mathbf{w}_{\text{LIN}}) &= \frac{1}{N} \|\mathbf{y} - \hat{\mathbf{y}}\|^2 = \frac{1}{N} \|(I - H)\mathbf{noise}\|^2 \\
 &= \frac{1}{N} (N - (d + 1)) \|\mathbf{noise}\|^2
 \end{aligned}$$

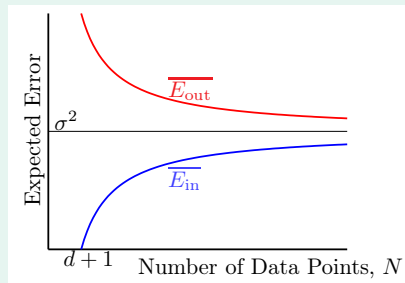
??? 这一步怎么来的?

$$\begin{aligned}
 \overline{E}_{\text{in}} &= \text{noise level} \cdot \left(1 - \frac{d+1}{N}\right) \\
 \overline{E}_{\text{out}} &= \text{noise level} \cdot \left(1 + \frac{d+1}{N}\right) \text{ (complicated!)}
 \end{aligned}$$

The Learning Curve

$$\overline{E}_{\text{out}} = \text{noise level} \cdot \left(1 + \frac{d+1}{N}\right)$$

$$\overline{E}_{\text{in}} = \text{noise level} \cdot \left(1 - \frac{d+1}{N}\right)$$



- both converge to σ^2 (**noise level**) for $N \rightarrow \infty$
- expected generalization error: $\frac{2(d+1)}{N}$
 —**similar to worst-case guarantee from VC**

linear regression (LinReg):
learning ‘happened’!

Fun Time

Which of the following property about H is not true?

- ① H is symmetric
- ② $H^2 = H$ (double projection = single one)
- ③ $(I - H)^2 = I - H$ (double residual transform = single one)
- ④ none of the above

Reference Answer: ④

You can conclude that ② and ③ are true by their physical meanings! :-)

Linear Classification vs. Linear Regression

Linear Classification

$$\mathcal{Y} = \{-1, +1\}$$

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

$$\text{err}(\hat{y}, y) = \mathbb{I}[\hat{y} \neq y]$$

NP-hard to solve in general

Linear Regression

$$\mathcal{Y} = \mathbb{R}$$

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

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

efficient analytic solution

$\{-1, +1\} \subset \mathbb{R}$: linear regression for classification?

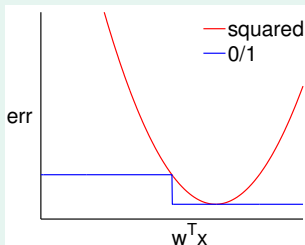
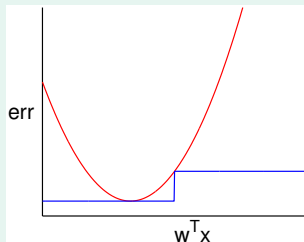
好想法，直接用Regression考虑classification的问题

- ① run LinReg on binary classification data \mathcal{D} (**efficient**)
- ② return $g(\mathbf{x}) = \text{sign}(\mathbf{w}_{\text{LIN}}^T \mathbf{x})$

but explanation of this **heuristic**?

Relation of Two Errors

$$\text{err}_{0/1} = \left[\text{sign}(\mathbf{w}^T \mathbf{x}) \neq y \right] \quad \text{err}_{\text{sqr}} = (\mathbf{w}^T \mathbf{x} - y)^2$$

desired $y = 1$ desired $y = -1$ 

\mathbf{w} 怎么取？取 \mathbf{x} 为一维咯，一般情况呢？

$$\text{err}_{0/1} \leq \text{err}_{\text{sqr}}$$

Linear Regression for Binary Classification

$$\text{err}_{0/1} \leq \text{err}_{\text{sqr}}$$

$$\begin{aligned} \text{classification } E_{\text{out}}(\mathbf{w}) &\stackrel{\text{VC}}{\leq} \text{classification } E_{\text{in}}(\mathbf{w}) + \sqrt{\dots\dots\dots} \\ &\leq \text{regression } E_{\text{in}}(\mathbf{w}) + \sqrt{\dots\dots\dots} \end{aligned}$$

R是C的一个upper bound，说实话，如果R做好的话，C也会好一些，即 E_{in}

- (loose) upper bound err_{sqr} as $\hat{\text{err}}$ to approximate $\text{err}_{0/1}$
- trade **bound tightness** for **efficiency**

\mathbf{w}_{LIN} : useful baseline classifier,
or as **initial PLA/pocket vector**

给出PLA的初始向量？！好想法，可能可以从一个好的解开始！可以试一试

Fun Time

Which of the following functions are upper bounds of the pointwise 0/1 error $\mathbb{I}[\text{sign}(\mathbf{w}^T \mathbf{x}) \neq y]$ for $y \in \{-1, +1\}$?

- ① $\exp(-y\mathbf{w}^T \mathbf{x})$ 0/1 上限
- ② $\max(0, 1 - y\mathbf{w}^T \mathbf{x})$
- ③ $\log_2(1 + \exp(-y\mathbf{w}^T \mathbf{x}))$
- ④ all of the above

Reference Answer: ④

Plot the curves and you'll see. Thus, all three can be used for binary classification. In fact, all three functions connect to very important algorithms in machine learning and we will discuss one of them soon in the next lecture.

Stay tuned. :-)

Summary

- 1 When Can Machines Learn?
- 2 Why Can Machines Learn?

Lecture 8: Noise and Error

- 3 **How** Can Machines Learn?

Lecture 9: Linear Regression

- Linear Regression Problem
use hyperplanes to approximate real values
- Linear Regression Algorithm
analytic solution with pseudo-inverse
- Generalization Issue
$$E_{\text{out}} - E_{\text{in}} \approx \frac{2(d+1)}{N} \text{ on average}$$
- Linear Regression for Binary Classification
0/1 error \leq squared error

• **next: binary classification, regression, and then?**

- 4 How Can Machines Learn Better?