

4/2 金



M2608.001300

Machine Learning Fundamentals & Applications

[3: Linear Classification and Regression]

Electrical and Computer Engineering
Seoul National University

© 2018 Sungroh Yoon. This material is for educational uses only. Some contents are based on the material provided by the textbook authors and may be copyrighted by them.

Outline

Linear Classification

Checking Two Basic Criteria for Learning
Non-Separable Data and Pocket Algorithm
Example: Handwritten Digit Recognition

Linear Regression

Introduction
Linear Regression Algorithm
Interpretation via Hat Matrix

Summary

Readings

- *Learning from Data* by Abu-Mostafa, Magdon-Ismail, and Lin
 - ▶ Chapter 3: The Linear Model (Sections 3.1 & 3.2)

references on linear algebra (optional):

- *Linear Algebra (MIT 18.06)* [▶ Link](#)
 - ▶ a well-known OCW (undergraduate level)
- *Introduction to Linear Dynamical Systems (Stanford EE263)* [▶ Link](#)
 - ▶ Lecture 3: Linear Algebra Review
 - ▶ Lecture 5: Least-Squares

The linear model

- set of lines
- often a good first choice
 - ▶ small VC dimension (more on this later) : simple
 - ⇒ generalize well from E_{in} to E_{out}
↳ 하이퍼 파라미터 같은 존재. (조정할 수 있는 것 같다)

Three important problems

1. classification (supervised learning)
2. regression ("")
3. probability estimation (*aka* logistic regression)
↳ 0과 1사이 값 return
 - ▶ come with different but related algorithms

딥러닝

* 3은 1도, 2도 될수

(연속적인 값을 return하는 3을 기준에 따라
차르면 1도 될수 있고,
2도 사용할수도 있다.)

Outline

Linear Classification

Checking Two Basic Criteria for Learning

Non-Separable Data and Pocket Algorithm

Example: Handwritten Digit Recognition

Linear Regression

Introduction

Linear Regression Algorithm

Interpretation via Hat Matrix

Summary

The linear model for binary classification

- uses a hypothesis set of linear classifiers
 - ▶ each h has the form

$$\underbrace{h(\mathbf{x})}_{\text{binary output}} = \text{sign}(\mathbf{w}^T \mathbf{x})$$

- $\mathbf{w} \in \mathbb{R}^{d+1}$: weight column vector
- d : input dimension
- $x_0 = 1$: corresponds to bias w_0

- note

- ▶ we use h and \mathbf{w} interchangeably to refer to the hypothesis
 \mathbf{w} : 무게, weight
(learning의 대상)

Recall: two basic criteria for learning

1. can we make sure $E_{\text{out}}(g) \approx E_{\text{in}}(g)$? ○
 - ▶ what we have learned in sample will generalize out of sample
 2. can we make $E_{\text{in}}(g) \approx 0$? ○
 - ▶ what we have learned in sample is a good hypothesis
- will these two criteria be fulfilled by a linear model? ○

Computational learning theory

- a quick summary (more details will follow in later lectures):

▶ (VC dimension) d_{VC} of linear model: only $d + 1$

▶ VC generalization bound: for any tolerance $\delta > 0$

$$E_{\text{out}} \leq E_{\text{in}} + \sqrt{\frac{8}{N} \ln \frac{4m_{\mathcal{H}}(2N)}{\delta}} \quad (1)$$

with probability $\geq 1 - \delta$

M : hypothesis set의 크기
m : effectively 한 M

▶ the bound on the growth function

$$m_{\mathcal{H}}(N) \leq N^{d_{VC}} + 1 \quad (2)$$

M은 한대라도 m을 bound 가능하다 : learning 가능

First criterion: $E_{\text{out}} \approx E_{\text{in}}$

- based on the above: with high probability

$$E_{\text{out}}(g) = E_{\text{in}}(g) + O\left(\sqrt{\frac{d}{N} \ln N}\right) \quad (3)$$

- for sufficiently large N
 - ▶ $\underbrace{E_{\text{out}} \approx E_{\text{in}}}_{\text{the first criterion: fulfilled}}$
- side note: a recurring pattern
 - ▶ generalization error typically contains $\frac{(d_{vc})}{N}$ term
(*i.e.* effective DOF over # samples)

Second criterion: $E_{in} \approx 0$

- linearly separable data

- ▶ there is always some hypothesis \mathbf{w}^* with $E_{in}(\mathbf{w}^*) = 0$

e.g. perceptron learning algorithm (PLA)

1. start with an arbitrary weight vector $\mathbf{w}(0)$
2. then, at every time step $t \geq 0$
 - 2a. select *any* misclassified data point $(\mathbf{x}(t), y(t))$
 - 2b. update \mathbf{w} :

$$\mathbf{w}(t+1) = \mathbf{w}(t) + y(t)\mathbf{x}(t)$$

- linearly inseparable data

- ▶ we can often make E_{in} very small

e.g. (pocket algorithm), support vector machine (SVM)

M이 무한히 커라도

m은 유한해



A4용지에 구멍

풀는것으로 생각 가능.

Outline

Linear Classification

Checking Two Basic Criteria for Learning

Non-Separable Data and Pocket Algorithm

Example: Handwritten Digit Recognition

Linear Regression

Introduction

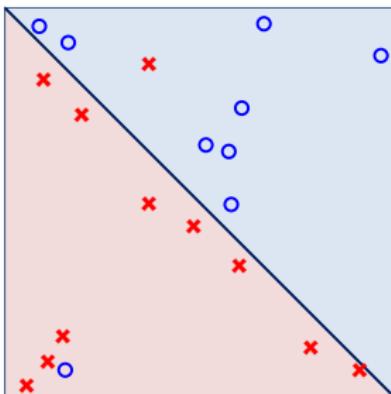
Linear Regression Algorithm

Interpretation via Hat Matrix

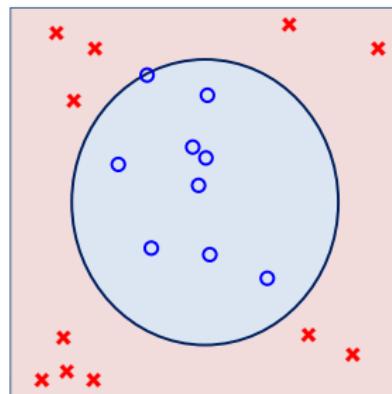
Summary

Data sets that are not linearly separable

2) 2종류 (a) et (b)



(a) a few noisy points

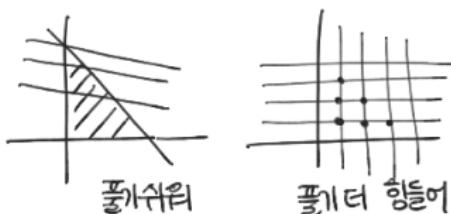


(b) nonlinearly separable

Figure 1: two types of linearly inseparable data

- (a) data becomes linearly separable after removing just 2 examples
(noisy examples or outliers)
- (b) data can be separated by a circle rather than a line

Applicability of linear model



- Figure 1(a)
 - ▶ it seems appropriate to stick with a line
 - ▶ but to somehow tolerate noise
 - ▶ and output a hypothesis with a small E_{in} (not necessarily 0)
 - Figure 1(b)
 - ▶ the linear model does not seem to be correct
 - ▶ may need a technique called non/linear transformation
- ←
때로는 integer 문제가
floating point 문제보다 어려워

Hardness

- the situation in Figure 1(a): encountered very often in practice
 - a linear classifier seems appropriate
 - but data may not be linearly separable due to outliers/noise
- to find a hypothesis with the minimum E_{in}
 - need to solve combinatorial optimization

$$\min_{\mathbf{w} \in \mathbb{R}^{d+1}} \underbrace{\frac{1}{N} \left[\text{sign}(\mathbf{w}^T \mathbf{x}_n) \neq y_n \right]}_{E_{in}(\mathbf{w})} \quad (4)$$

\hookrightarrow 이득(discrete)

- Eq. (4): difficult to solve (in general, NP-hard) than continuous
 - due to the discrete nature of both $\text{sign}(\cdot)$ and $[\cdot]$
- \Rightarrow approximately minimize E_{in}

Approximately minimizing E_{in}

- one approach:
 - ▶ extend PLA through a simple modification
 - ▶ socket algorithm with ratchet
(Gallant, 1990)
 - $E_{in}(w)$ 을 각 Iteration 따라 계산
- in essence
 - ▶ the pocket algorithm keeps 'in its pocket' the best w encountered up to iteration t in PLA
 - ▶ at the end, the best w will be reported as the final hypothesis



The pocket algorithm

- 1: set the pocket weight vector $\hat{\mathbf{w}}$ to $\mathbf{w}(0)$ of PLA;
- 2: **for** $t = 0, \dots, T - 1$ **do**
 - 3: run PLA for one update to obtain $\mathbf{w}(t + 1)$;
 - 4: evaluate $E_{in}(\mathbf{w}(t + 1))$; /* time consuming */
 - 5: **if** $\mathbf{w}(t + 1)$ is better than $\hat{\mathbf{w}}$ in terms of E_{in}
 - 6: **then** set $\hat{\mathbf{w}}$ to $\mathbf{w}(t + 1)$;
 - 7: **return** $\hat{\mathbf{w}}$;

예제에서 더해지게.

Comparison



in each iteration:

- PLA
 - ▶ only checks *some* examples using $\mathbf{w}(t)$ to identify $(\mathbf{x}(t), y(t))$
- pocket algorithm
 - ▶ needs an additional step ([line 4](#))
 - ▶ evaluates *all* examples using $\mathbf{w}(t + 1)$ to get $E_{in}(\mathbf{w}(t + 1))$
 - ⇒ much *slower* than PLA
 - ⇒ no guarantee for how fast it can converge to a good E_{in}

Outline

Linear Classification

Checking Two Basic Criteria for Learning
Non-Separable Data and Pocket Algorithm

Example: Handwritten Digit Recognition

Linear Regression

Introduction

Linear Regression Algorithm

Interpretation via Hat Matrix

Summary

Handwritten digit recognition

- data: sampled digits from the US Postal Service zip code DB
 - ▶ each image has 16×16 pixels

- goal: recognize each digit
 - ▶ non-trivial task (even human E_{out} is about 2.5%)
 - ▶ common confusion: digits $\{4, 9\}$ and $\{2, 7\}$

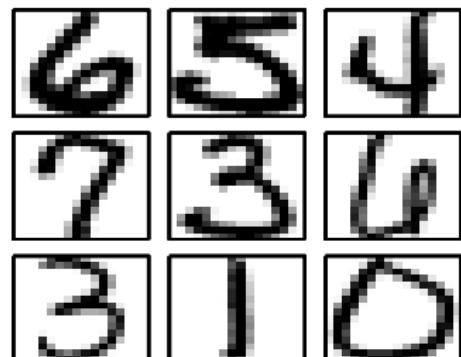


Figure 2: digit images

Solution components

1. decomposition from *multiclass* to *binary* classification

- ▶ a commonly used approach for multiclass setup

2 data의 주요 feature를 뽑아 feature extraction

- ▶ a human approach to determining the digit of an image: to look at the shape (or other properties) of black pixels
- ▶ rather than carrying all the info in 256 pixels, summarize the info contained in the image into a few features.
- ▶ again a very common/essential approach in machine learning

Image features

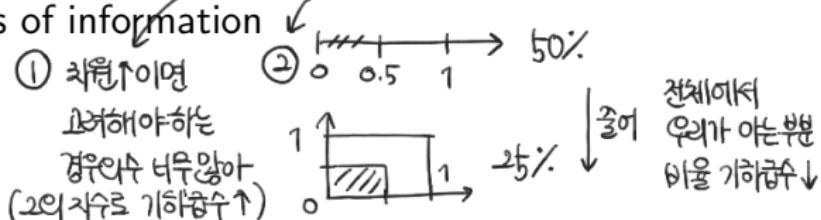
- two important features here: intensity and symmetry
 - ▶ digit 5 usually occupies more black pixels than digit 1
 - ▶ digit 1 is symmetric while digit 5 is not
 - ▶ these two features will define the *feature space*



- definitions
 - ▶ asymmetry: the average absolute difference between an image and its flipped version
 - ▶ symmetry: the negation of asymmetry

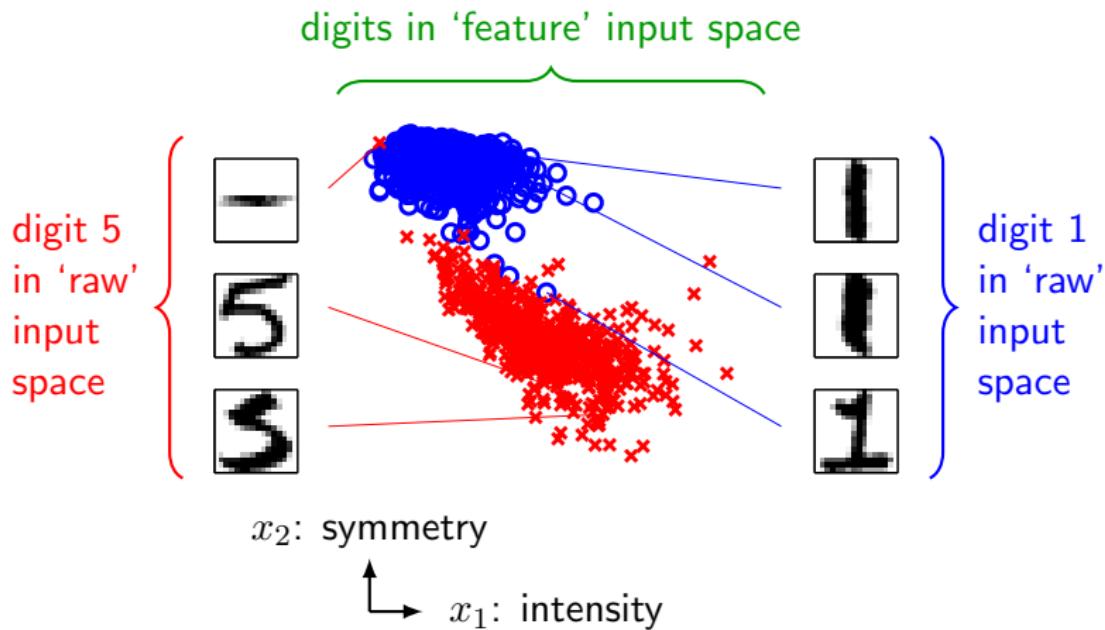
Input representation

- in 'raw' input space \mathbb{R}^{256+1} : $\mathbf{x} = (x_0, x_1, x_2, \dots, x_{256})$
 - ▶ linear model: $\mathbf{w} = (w_0, w_1, w_2, \dots, w_{256})$
- in 'feature' input space \mathbb{R}^{2+1} : $\mathbf{x} = (x_0, x_1, x_2)$
 - ▶ linear model: $\mathbf{w} = (w_0, w_1, w_2)$
- introducing feature space (preprocessing 에서 차원 줄이는 것 매우 중요)
 - ▶ reduces the number of dimensionality (e.g. from 256 to 2)
 - ⇒ allows us to better handle the curse of dimensionality
 - ▶ but can cause loss of information



Example

- a scatter plot for these intensity and symmetry features:
 - ▶ $\mathbf{x} = (x_0, x_1, x_2)$



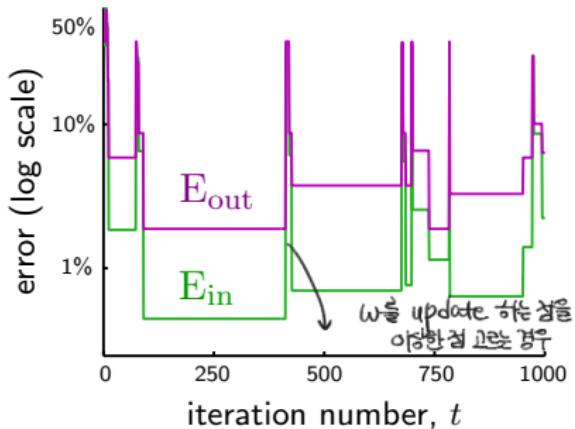
Classification in (intensity, symmetry) space

출인후 공간

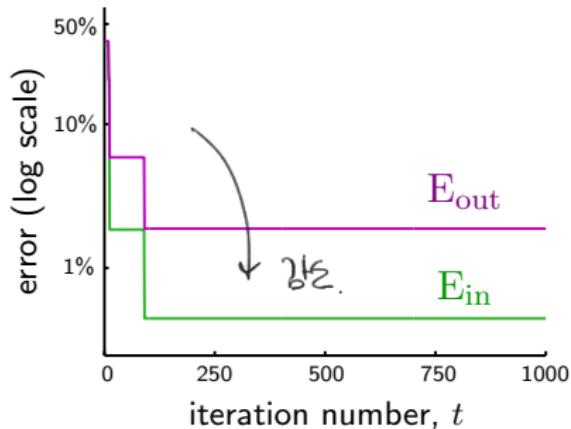
- digits can be roughly separated by a line in the feature plane
 - ▶ but poorly written digits prevent a perfect line separation
 - ▶ e.g. the '5' depicted in the top-left corner
- we now run PLA & pocket algorithm
 - ▶ with data points represented in the feature space

- error versus iteration number

- ▶ PLA: $E_{in} = 2.24\%$ and $E_{out} = 6.37\%$ (quit at iteration 1000)
- ▶ pocket: $E_{in} = 0.45\%$ and $E_{out} = 1.89\%$



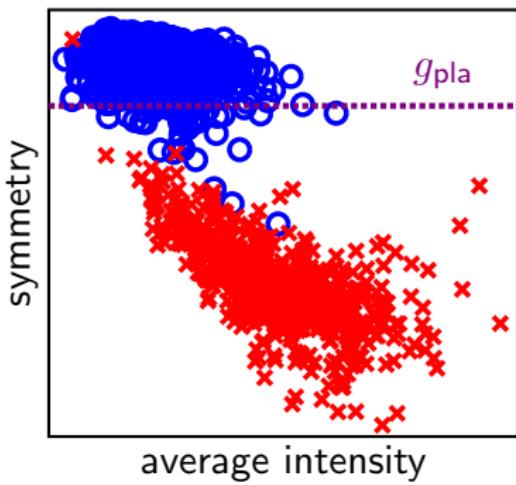
(a) PLA



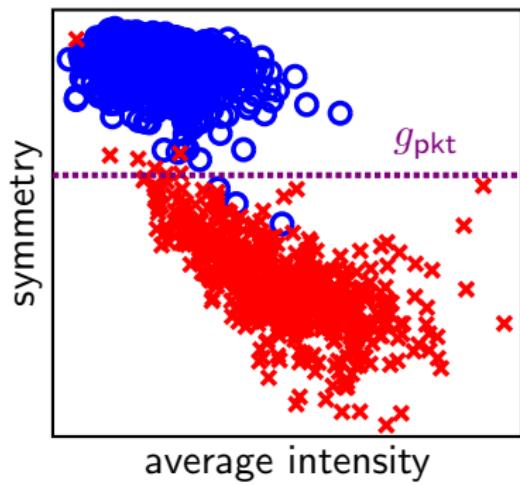
(b) pocket algorithm

Figure 3: comparison of two linear algorithms for separating digits 1 and 5

- learned hypothesis g



(a) PLA



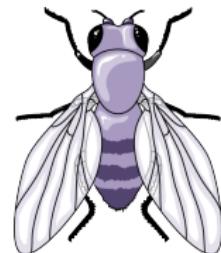
(b) pocket algorithm

MNIST dataset

Link

- 70k handwritten digits
 - ▶ 28×28 pixels
 - ▶ training set: 60k
 - ▶ test set: 10k
- best performing classifiers
 - ▶ linear: 7.6% error
 - ▶ Gaussian SVM: 1.4%
 - ▶ CNN: 0.23%
- “drosophila of machine learning”
 - ▶ Geoffrey Hinton

0
1
2
3
4
5
6
7
8
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9



Outline

Linear Classification

Checking Two Basic Criteria for Learning
Non-Separable Data and Pocket Algorithm
Example: Handwritten Digit Recognition

Linear Regression

Introduction

Linear Regression Algorithm
Interpretation via Hat Matrix

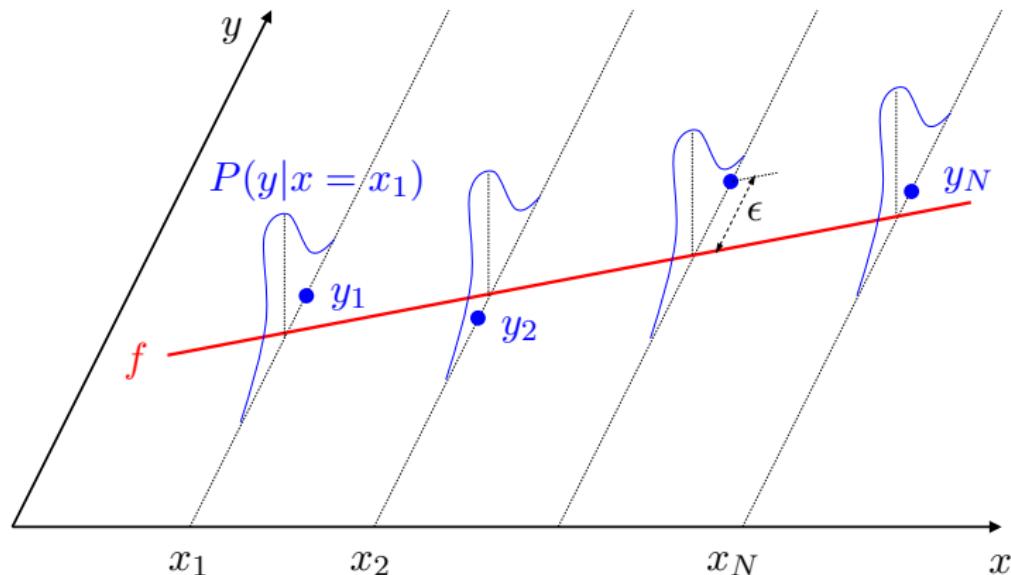
Summary

Regression

통계학적인

- a statistical method to study relationship between x and y
 - covariate/predictor variable/independent variable/feature
 - y response/dependent variable ↪ 다른 이는 다른
- training data $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$
 - ▶ noise ϵ is added to target: $y_n = f(\mathbf{x}_n) + \epsilon$
 - ⇒ $y \sim P(y|\mathbf{x})$ instead of $y = f(\mathbf{x})$ ↪ $y \sim P(y|\mathbf{x})$
- our goal
 - ▶ find a model $g(\mathbf{x})$ that approximates y_n

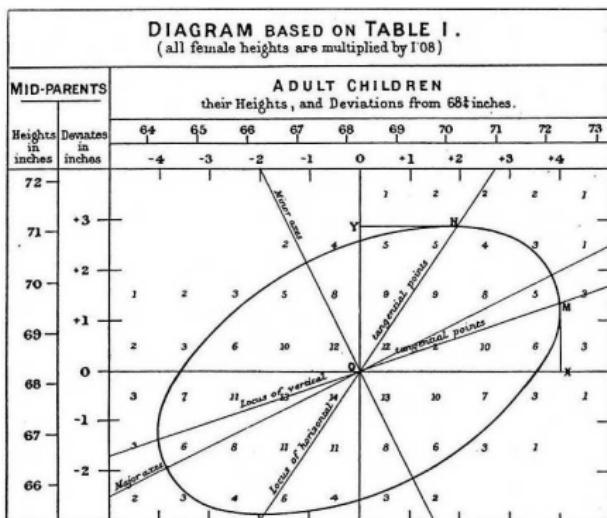
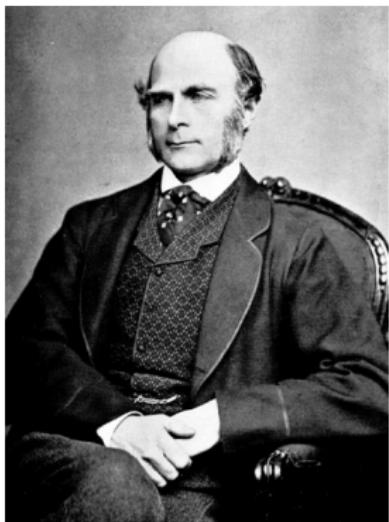
- illustration (simple linear regression):
각지점 차원



각지점 똑같은 bell shape 쓰겠다.

assumption: homoscedasticity (homogeneity of variance; the same finite variance for each value of x)

- etymology: Sir Francis Galton (1822-1911)
 - ▶ re (back)+gression (going): going back from data to formula
 - ▶ regression towards the mean: “tall and short men tend to have sons with heights closer to the mean”
 - data와 formula를 찾자.



Linear regression

y is continuous

- popular linear model for predicting a quantitative response
 - ▶ applies to real-valued target functions (y)
 - ▶ long history in statistics, social and behavioral sciences
 - ▶ “regression” means y is real-valued (inherited from statistics)
- simple versus multiple
 - ▶ $d = 1$: simple linear regression (one predictor)
 - ▶ $d \geq 2$: multiple linear regression (multiple predictors)
- we discuss “simple linear regression” from a learning perspective
 $d=1$

Parameter estimation for linear regression

복잡성 complexity도
같이 줄여주자는 것.

(모델 너무 복잡하면
예측까지 다 learn
하기 때문에)

- **least squares**

- ▶ ordinary least squares (OLS) 최소제곱하는 방법
 - ▷ minimizes the sum of squared residuals
 - ▷ leads to a closed-form expression
- ▶ generalized least squares, iteratively reweighted least squares

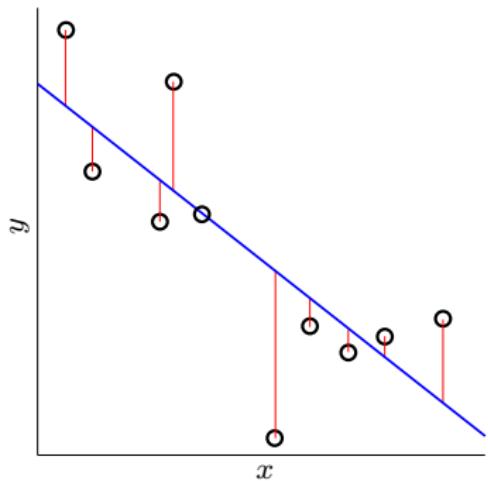
- **maximum likelihood**

- ▶ ridge (Tikhonov regularization)/lasso regression: regularized
- ▶ least absolute deviation regression

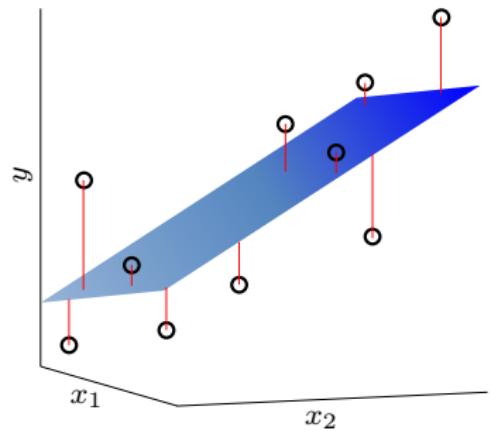
- **other techniques**

- ▶ Bayesian linear regression, principal component regression

Linear regression in 1D and 2D



(a) 1D (line)



(b) 2D (hyperplane)

Figure 4: the solution hypothesis (in blue) of the linear regression algorithm in 1D and 2D. in OLS, the sum of squared error is minimized.

Credit approval example revisited

ശ്രൂ മറ്റ് പര്യായ ചെന്നി.

- consider a regression problem (rather than classification)
 - ▶ e.g. set a credit limit for each customer
- the bank uses historical records to build a data set \mathcal{D}
 - ▶ $\mathcal{D}: (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$
 - ▶ \mathbf{x}_n : customer information
 - ▶ y_n : credit limit (set by human experts)
- y_n is now a real number
 - ▶ e.g. \$123,000 for Alice and \$57,500 for Bob

- the bank wants to use learning to find a hypothesis g
 - ▶ g replicates how human experts determine credit limits
- our target:
 - ▶ not a deterministic function $y = f(\mathbf{x})$
 - ▷ there is more than one human expert
 - ▷ each expert may not be perfectly consistent
 - ▶ instead: noisy target
 - ▷ formalized as a distribution of random variable y

Target distribution rather than a function

- regression label y_n comes from some distribution $P(y|\mathbf{x})$
 - ▶ instead of a deterministic function $f(\mathbf{x})$
- nonetheless, the nature of the problem is not changed
 - ▶ given: *unknown* distribution $P(\mathbf{x}, y)$ generating each (\mathbf{x}_n, y_n)
 - ▶ objective: find a hypothesis g that minimizes the error between $g(\mathbf{x})$ and y wrt the distribution P

Outline

Linear Classification

Checking Two Basic Criteria for Learning
Non-Separable Data and Pocket Algorithm
Example: Handwritten Digit Recognition

Linear Regression

Introduction

Linear Regression Algorithm

Interpretation via Hat Matrix

Summary

The linear regression algorithm

- based on minimizing the squared error between $h(\mathbf{x})$ and y

$$E_{\text{out}}(h) = \mathbb{E} [(h(\mathbf{x}) - y)^2]$$

- ▶ expected value: taken wrt joint pdf $P(\mathbf{x}, y)$
- goal
 - ▶ find a hypothesis h that achieves a smallest $E_{\text{out}}(h)$
- issue: E_{out} cannot be computed
 - $\because P(\mathbf{x}, y)$ is unknown

- thus, we resort to in-sample error instead:

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2 \text{ (정의 OLS)}$$

- in *linear* regression, h takes the form of
 - a linear combination of the components of \mathbf{x} :

$$h(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$$

- $\mathbf{w}^T \mathbf{x}$: also called "signal"

Matrix representation: data

- data matrix $X \in \mathbb{R}^{N \times (d+1)}$ $N \left[\begin{array}{c} d+1 \\ \vdash \end{array} \right] \ni X$

► rows: inputs \mathbf{x}_n as row vectors

- target vector $y \in \mathbb{R}^N$ $N \left[\begin{array}{c} \vdash \\ \vdash \end{array} \right] \ni y$

► components: target values y_n

- example: $d = 1, N = 4$

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ 1 & x_4 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

Matrix representation: in-sample error って何?

- in-sample error is a function of \mathbf{w} and data X, \mathbf{y} :

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

$$= \frac{1}{N} \left\| \begin{bmatrix} \mathbf{x}_1^T \mathbf{w} - y_1 \\ \mathbf{x}_2^T \mathbf{w} - y_2 \\ \vdots \\ \mathbf{x}_N^T \mathbf{w} - y_N \end{bmatrix} \right\|^2$$

$$\begin{aligned} &= \frac{1}{N} \| X\mathbf{w} - \mathbf{y} \|^2 \\ &= \frac{1}{N} (\mathbf{w}^T X^T X \mathbf{w} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \end{aligned} \quad (5)$$

$$= \frac{1}{N} \| X\mathbf{w} - \mathbf{y} \|^2 \quad (\because \|x\|^2 = x^T x) \quad (6)$$

$$= \frac{1}{N} (\mathbf{w}^T X^T X \mathbf{w} - 2\mathbf{w}^T X^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \quad (7)$$

- ▶ $\|\cdot\|$: Euclidean norm of a vector
- ▶ scalar $\mathbf{y}^T X \mathbf{w} = (\mathbf{w}^T X^T \mathbf{y})^T = \mathbf{w}^T X^T \mathbf{y}$

- example: $d = 1, N = 4$

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ 1 & x_4 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

$E_{in}(\mathbf{w})$ = $\frac{1}{4} \sum_{i=1}^4 (\mathbf{w}^T \mathbf{x}_n - y_n)^2$

E_{in} 을 최소화

하는 \mathbf{w} 찾자.

$$= \frac{1}{4} \left\| \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ 1 & x_4 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \right\|^2$$

Getting the solution \mathbf{w}_{lin}

- \mathbf{w}_{lin} : the solution to linear regression : E_{in} 을 최소화하는 \mathbf{W} 찾기.
 - ▶ derived by minimizing $E_{\text{in}}(\mathbf{w})$ over all possible $\mathbf{w} \in \mathbb{R}^{d+1}$
다음을 최소화하는 $\arg \min \mathbf{W}$ 찾기.

$$\mathbf{w}_{\text{lin}} = \underset{\mathbf{w} \in \mathbb{R}^{d+1}}{\operatorname{argmin}} E_{\text{in}}(\mathbf{w}) \quad (8)$$

$$= \underset{\mathbf{w} \in \mathbb{R}^{d+1}}{\operatorname{argmin}} \frac{1}{N} \|X\mathbf{w} - \mathbf{y}\|^2 \quad (9)$$

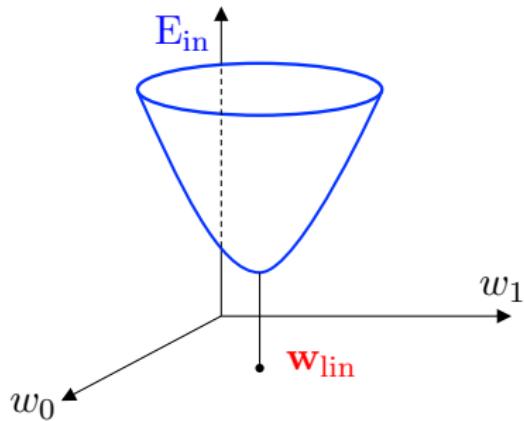
Minimization of $E_{in}(w)$

- Eq. (7) implies
 - ▶ $E_{in}(w)$ is continuous, differentiable, and convex \rightarrow 이쁘다.
 - ⇒ we can use standard matrix calculus to find w that minimizes $E_{in}(w)$ by requiring

$$\nabla E_{in}(w) = 0$$

E_{in} (objective fn) 을
최소로 하는 w 찾자.

- general optimization techniques
 - ▶ e.g. gradient descent



- recall: gradient identities¹

$$\nabla_{\mathbf{w}}(\mathbf{w}^T A \mathbf{w}) = (A + A^T)\mathbf{w}$$

$$\nabla_{\mathbf{w}}(\mathbf{w}^T \mathbf{b}) = \mathbf{b}$$

- scalar w

$$E_{in}(w) = aw^2 - 2bw + c$$

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

- vector \mathbf{w}

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

$$\nabla E_{in}(\mathbf{w}) = (A + A^T)\mathbf{w} - 2\mathbf{b}$$

¹matrix analog of ordinary differentiation of quadratic and linear functions

The solution

- from Eq. (7)

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

 앞에 곱해진 것 같은 듯 하면

- ▶ both \mathbf{w} and $\nabla E_{in}(\mathbf{w})$ are column vectors
- finally, one should solve for \mathbf{w} that satisfies the *normal equations*?

$$X^T X \mathbf{w} = X^T \mathbf{y}$$

X^T 가 invertible: $X\mathbf{w} = \mathbf{y} \Leftrightarrow \forall \varepsilon = 0$

Two scenarios

$$\underline{(X^T X)w = X^T y} : \text{normal eq-}$$

1. if $X^T X$ is invertible, $\boxed{w = X^{-1}y}$: w 의 값 1개로 결정가능

- ▶ $X^\dagger = \boxed{(X^T X)^{-1} X^T}$ is pseudo-inverse of X
- ▶ resulting w is the unique optimal solution to (8)
inverse가 매우 불안정해.
실제로는 $X^T X$ 가 invertible 아님

2. if $X^T X$ is not invertible ◀ 이 책은 사용해

- ▶ a pseudo-inverse can still be defined, but no unique solution
- ▶ there will be many solutions for w that minimizes E_{in}
- ▶ e.g. $w_{\text{lin}} = \boxed{V \Gamma^{-1} U^T y}$ is a solution to (8)²
예외 습득 ↓



² $X = U \Gamma V^T$ is the singular value decomposition (SVD) of X (see page 55)

In practice

- $X^T X$ is invertible in most cases
 - ▶ since N is often much bigger than $d + 1$
 - ⇒ there will likely be $d + 1$ linearly independent vectors \mathbf{x}_n
- when $X^T X$ is (almost) singular
 - ▶ use a well-implemented \tilde{X}^\dagger routine instead of $(X^T X)^{-1} X^T$
 - ▶ this is for numerical stability → 오류 amplify ↓
→ 예전 적도를 디자인한 X^\dagger 을 빌고 사용해.
- we have derived the following *linear regression algorithm*

The linear regression algorithm

1. construct matrix X and vector y as follows:

- ▶ use data set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$
- ▶ each \mathbf{x} includes $x_0 = 1$ bias coordinate

$$X = \underbrace{\begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}}_{\text{input data matrix}} \quad y = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}$$

2. compute pseudo-inverse X^\dagger of X ; if $X^T X$ is invertible,

$$X^\dagger = (X^T X)^{-1} X^T$$

← invertible 인 경우
SVD 등 대체로
방법이 존재한다.

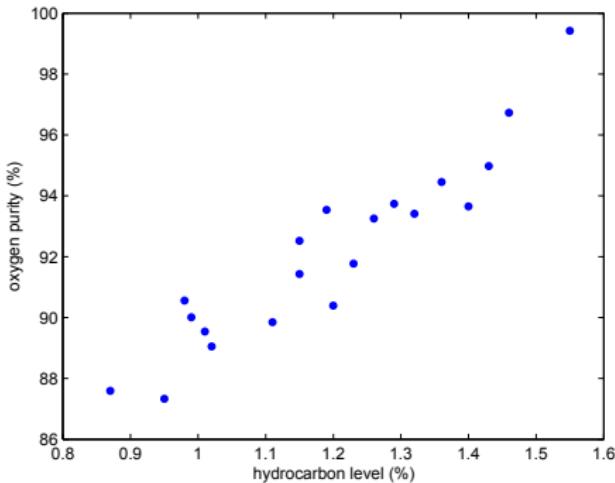
3. return $\mathbf{w}_{\text{lin}} = X^\dagger y$

Example: oxygen and hydrocarbon levels

no.	x	y
1	0.99	90.01
2	1.02	89.05
3	1.15	91.43
4	1.29	93.74
5	1.46	96.73
6	1.36	94.45
7	0.87	87.59
8	1.23	91.77
9	1.55	99.42
10	1.40	93.65
11	1.19	93.54
12	1.15	92.52
13	0.98	90.56
14	1.01	89.54
15	1.11	89.85
16	1.20	90.39
17	1.26	93.25
18	1.32	93.41
19	1.43	94.98
20	0.95	87.33

y: purity (%) of oxygen produced in a chemical distillation process

x: hydrocarbons (%) that are present in the main condenser of the distillation unit



- training data:

data matrix $X =$

$$\left[\begin{array}{cc} 1 & 0.99 \\ 1 & 1.02 \\ 1 & 1.15 \\ 1 & 1.29 \\ 1 & 1.46 \\ 1 & 1.36 \\ 1 & 0.87 \\ 1 & 1.23 \\ 1 & 1.55 \\ 1 & 1.40 \\ 1 & 1.19 \\ 1 & 1.15 \\ 1 & 0.98 \\ 1 & 1.01 \\ 1 & 1.11 \\ 1 & 1.20 \\ 1 & 1.26 \\ 1 & 1.32 \\ 1 & 1.43 \\ 1 & 0.95 \end{array} \right]$$

12

, target vector $y =$

$$\left[\begin{array}{c} 90.01 \\ 89.05 \\ 91.43 \\ 93.74 \\ 96.73 \\ 94.45 \\ 87.59 \\ 91.77 \\ 99.42 \\ 93.65 \\ 93.54 \\ 92.52 \\ 90.56 \\ 89.54 \\ 89.85 \\ 90.39 \\ 93.25 \\ 93.41 \\ 94.98 \\ 87.33 \end{array} \right]$$

- $X^T X$ is invertible

$$X^T X = \begin{bmatrix} 20.00 & 23.92 \\ 23.92 & 29.29 \end{bmatrix} \Rightarrow (X^T X)^{-1} = \begin{bmatrix} 2.15 & -1.76 \\ -1.76 & 1.47 \end{bmatrix}$$

- $(X^T X)^{-1} X^T \mathbf{y}$ yields

$$\mathbf{w}_{\text{lin}} = \begin{bmatrix} 74.28 \\ 14.95 \end{bmatrix}$$

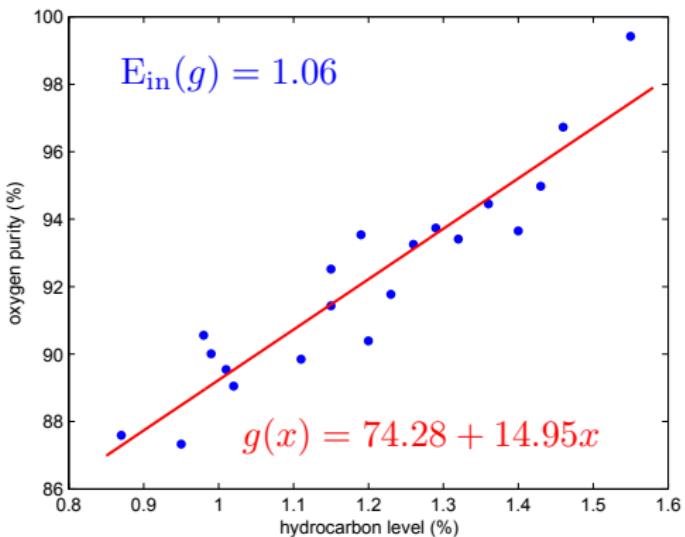
- the learned model:

$$\begin{aligned} g(x) &= \mathbf{w}_{\text{lin}}^T \mathbf{x} \\ &= 74.28 + 14.95x \end{aligned}$$

- error (see Eq. 3):

$$E_{\text{in}}(g) = 1.06$$

$$E_{\text{out}}(g) \approx 1.45$$



Remarks

- the linear regression algorithm on page 50
 - ▶ also called ordinary least squares (OLS)
 - ▶ provides the best linear unbiased estimator (BLUE)
- linear regression does not really look like learning
 - ▶ compared with PLA
 - ▶ hypothesis w_{lin} comes from an analytic solution (matrix inversion/multiplications) rather than iterative learning steps
- but learning *has* occurred (“one-step” learning)
 - ▶ as long as hypothesis w_{lin} has a decent E_{out}

포켓등과 다르게
수렴의 iteration으로
찾을수 있다는 것

- linear regression: a rare case
 - ▶ analytic formula for learning that is easy to evaluate
 \Rightarrow “very popular”
- there are methods for computing the pseudo-inverse X^\dagger directly
 - ▶ without inverting a matrix
 - ▶ numerically more stable than matrix inversion
- assume: the SVD of input matrix X is $X = U\Gamma V^T$
 - ▷ $U \in \mathbb{R}^{N \times \rho}$ satisfies $U^T U = I_\rho$,
 - ▷ $V \in \mathbb{R}^{(d+1) \times \rho}$ satisfies $V^T V = I_\rho$,
 - ▷ $\Gamma \in \rho \times \rho$ is a positive diagonal matrix with $\rho = \text{rank}(X)$
 - ▶ then $\mathbf{w}_{\text{lin}} = V\Gamma^{-1}U^T\mathbf{y}$ satisfies $X^T X \mathbf{w}_{\text{lin}} = X^T \mathbf{y}$

Outline

Linear Classification

Checking Two Basic Criteria for Learning
Non-Separable Data and Pocket Algorithm
Example: Handwritten Digit Recognition

Linear Regression

Introduction
Linear Regression Algorithm
Interpretation via Hat Matrix

Summary

Hat matrix H (aka projection matrix)

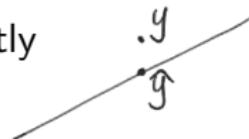
- \textcircled{H} $\rightarrow E_{in}$ 과 E_{out} 의 관계 내용이
 H relates to in-sample and out-of-sample errors
 - ▶ one of analysis tools developed in statistics
- H maps the vector of **observed** values to the vector of **fitted** values
 - ▶ if \mathbf{y} and $\hat{\mathbf{y}}$ denote **observed** and **fitted** values, respectively

$$\hat{\mathbf{y}} = H\mathbf{y} \quad (10)$$

- \textcircled{H} : observed values \mapsto fitted values
- ▶ matrix H ‘puts a hat’ on \mathbf{y} , hence the name
yet estimator: \hat{y}

Hat matrix in linear regression

- the linear regression weight vector \mathbf{w}_{lin} 값을 무너뜨리기 가장 작은 점
 - an attempt to map inputs X to outputs \mathbf{y}
- however, \mathbf{w}_{lin} does not produce \mathbf{y} directly
 - but produces an estimate



$$\hat{\mathbf{y}} = X\mathbf{w}_{\text{lin}}$$

which differs from \mathbf{y} due to in-sample error

- substituting the expression for \mathbf{w}_{lin}
 - we get (assuming $X^T X$ is invertible)

$$\hat{\mathbf{y}} = \underbrace{X(X^T X)^{-1} X^T}_{H \in \mathbb{R}^{N \times N}} \mathbf{y} = X \mathbf{w}_{\text{lin}} = H\mathbf{y}$$

- therefore, estimate \hat{y} is a linear transformation of actual y
 - ▶ through matrix multiplication with H where

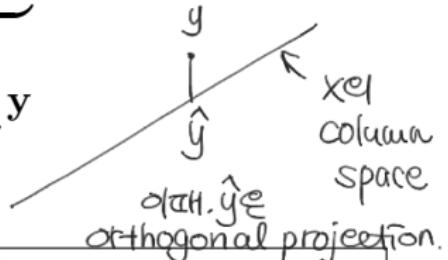
$$H = \underline{X} (\underline{X^T X})^{-1} \underline{X^T} \quad (11)$$

- bottom line:

projection : $X^T X$ 인데

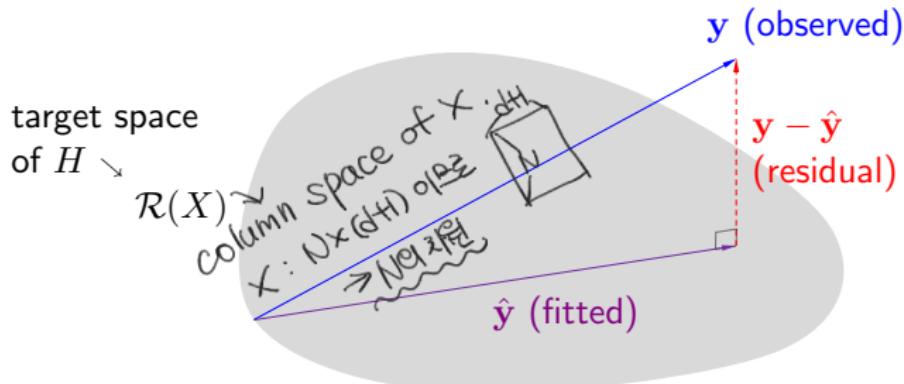
가운데는 normalize 하는 것.

$$\begin{aligned}\hat{y} &= X w_{\text{lin}} \\ &= X \underbrace{(\underline{X^T X})^{-1} \underline{X^T} y}_{w_{\text{lin}}} \\ &= \underbrace{X (X^T X)^{-1} X^T}_{H} y\end{aligned}$$



∴ \hat{y} : orthogonal projection of y onto the column space of X

Geometric interpretation



to minimize $\|y - \hat{y}\|$, \hat{y} should be the orthogonal projection of y onto $\mathcal{R}(X)$

$$\hat{y} = X\mathbf{w}_{\text{lin}} = H\mathbf{y}$$

$$\mathbf{y}, \hat{\mathbf{y}}, \mathcal{R}(X) \in \mathbb{R}^N$$

$$H = X(X^T X)^{-1} X^T$$

$$H \in \mathbb{R}^{N \times N}$$

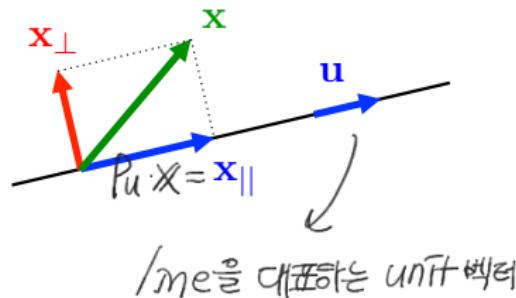
$$\mathcal{R}(X) = \text{range (column span) of } X$$

$$X \in \mathbb{R}^{N \times (d+1)}$$

$$\mathbf{w}_{\text{lin}} \in \mathbb{R}^{d+1}$$

Intuitive/informal explanation of projection

- a simple case: orthogonal projection onto a line
 - ▶ let \mathbf{u} denote a unit vector on the line
 - ▶ the projection matrix is then given by $P_{\mathbf{u}} = \mathbf{u}\mathbf{u}^T$
- to see why, decompose an arbitrary \mathbf{x} into $\mathbf{x} = \mathbf{x}_{||} + \mathbf{x}_{\perp}$
 - ▶ $\mathbf{x}_{||}$: component of \mathbf{x} on the line
 - ▶ \mathbf{x}_{\perp} : component of \mathbf{x} perpendicular to the line



- applying projection to \mathbf{x} gives

$$\begin{aligned}
 \underbrace{P_{\mathbf{u}} \mathbf{x}}_{=} &= \mathbf{u} \mathbf{u}^T \mathbf{x}_{||} + \mathbf{u} \mathbf{u}^T \mathbf{x}_{\perp} \\
 &= \mathbf{u} |\mathbf{x}_{||}| + \mathbf{u} \mathbf{0} \\
 &= \mathbf{x}_{||} = \text{P}_{\mathbf{u}} \mathbf{x} \text{ 입을 증명.}
 \end{aligned}$$

which confirms that $P_{\mathbf{u}}$ works

- generalization (orthonormal basis):

\downarrow vector
여러개로.

- ▶ let $\mathbf{u}_1, \dots, \mathbf{u}_k$ be an orthonormal basis of subspace U
- ▶ let \mathbf{A} denote the matrix whose columns are $\mathbf{u}_1, \dots, \mathbf{u}_k$
- ▶ then the projection matrix onto U is $\boxed{P_A = \mathbf{A} \mathbf{A}^T}$

- generalization (non-orthonormal basis): *orthonormal only*

- let $\mathbf{u}_1, \dots, \mathbf{u}_k$ be a (not necessarily orthonormal) basis of U
- let A denote the matrix whose columns are $\mathbf{u}_1, \dots, \mathbf{u}_k$
- then the projection matrix onto U is $P_A = A(A^T A)^{-1} A^T$
- note that $(A^T A)^{-1}$ is a “normalizing factor”

$$P_A = A \underbrace{(A^T A)^{-1}}_{\downarrow} A^T$$

- for vector \mathbf{y} , $P_A \mathbf{y} = A(A^T A)^{-1} A^T \mathbf{y} = \hat{\mathbf{y}}$
- projects \mathbf{y} orthogonally onto $\mathcal{R}(A)$, the column span of A

Properties of projection/hat matrices

let $H \in \mathbb{R}^{N \times N}$ (*i.e.*, H is a square matrix):

- H is called a projection/hat matrix iff $H^2 = H$
- two basic properties of hat matrix H
 - ▶ symmetric: $H^T = H$
 - ▶ idempotent: $H^2 = H$ (no effect on vectors already on space)
- for hat matrix H and identity matrix $I \in \mathbb{R}^{N \times N}$
 - ▶ $I - H$ is also a hat matrix:
$$(I - H)^2 = I - 2H + H^2 = I - H$$
 - ▶ $H^T(I - H) = 0 \Rightarrow$ target spaces are orthogonal

- for hat matrix $H = X(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ where $X \in \mathbb{R}^{N \times (d+1)}$
 - ▶ H : orthogonal projection onto the column space of X
 - ▶ range $\mathcal{R}(X)$ is thus called the target space of H
 - ▶ $\text{trace}(H) = \underline{d+1}$

$$\begin{aligned}\text{trace}(H) &= \text{trace}(X(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) \\ &= \text{trace}(\mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}) \\ &= \text{trace}(I) \\ &= \underline{d+1}\end{aligned}$$
 - ▷ note: $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{(d+1) \times (d+1)}$ and thus $I \in \mathbb{R}^{(d+1) \times (d+1)}$
- recall: $\text{trace}(H)$ is sum of diagonal elements in H
 - ▶ property: $\text{trace}(AB) = \text{trace}(BA)$ for two matrices A and B

Generalization of linear regression

정의하면: Input X 로 분산 y 를 예측한다.
이때, best X 하이퍼플레인 만들자.
우리는 error을 OLS로 선택.

~~ 그렇지만 이게 어떤 측정에서 좋지 않은가?

- linear regression looks for optimal w in terms of E_{in}

- does this guarantee decent E_{out} ?

yes!

1) 원시방법:

$$w_{lin} = (X^T X)^{-1} X^T y.$$

2) H: hat matrix 방법:

투영과 앤탈 통해 차이가
최소 되도록

- theoretical analysis shows

- generalization bound of linear regression:

이 역시 일의 것 같다.

$$E_{out}(g) = E_{in}(g) + O\left(\frac{d}{N}\right)$$

→ 결국 같은
이상가요?

마치 같은 classification처럼.

⇒ comparable to classification (3): $E_{out}(g) = E_{in}(g) + O\left(\sqrt{\frac{d}{N} \ln N}\right)$

assumption 이용하면

- we can even derive exact formulas for E_{in} & E_{out} (see textbook)

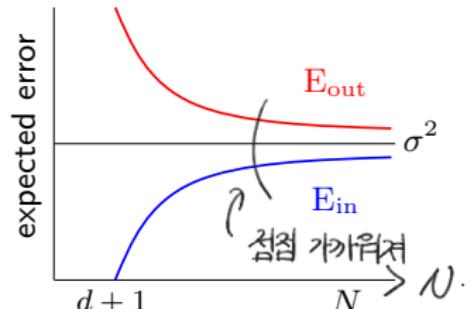
- assume: $\hat{y} = y + \epsilon$ (noise ϵ has zero mean and σ^2 variance)
 $y = f + \varepsilon$

→ 모든 ϵ 을 평균,
전체 power를 토대로 한다.

- the best possible linear fit has the expected error σ^2 (noise level)

$$E_{in} = \sigma^2 \left(1 - \frac{d+1}{N} \right) < \sigma^2$$

$$E_{out} = \sigma^2 \left(1 + \frac{d+1}{N} \right) > \sigma^2$$



- both converge to σ^2 as $N \rightarrow \infty$
- $E_{in} < \sigma^2$ because regression fits in-sample noise as well
- $E_{out} > \sigma^2$ reflects drift in w_{lin} due to fitting in-sample noise
- expected generalization error: $O\left(\frac{2(d+1)}{N}\right)$
- learning has 'happened' by linear regression



Outline

Linear Classification

Checking Two Basic Criteria for Learning
Non-Separable Data and Pocket Algorithm
Example: Handwritten Digit Recognition

Linear Regression

Introduction
Linear Regression Algorithm
Interpretation via Hat Matrix

Summary

Summary

- linear models
 - ▶ used for classification, regression, probability estimation
 - ▶ have small d_{VC} and generalize well \Rightarrow first practical choice
 - ▶ use the 'signal': $\sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$ (곱하기가 아니라 $x_0 = 1$)
 - linear classification
 - linear regression
- $\mathcal{Y} = \{-1, +1\}$ $\mathcal{Y} = \mathbb{R}$
- $h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$ $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- $e(\hat{y}, y) = [\hat{y} \neq y]$ $\begin{cases} 1 & \text{같으면 } 0 \\ 0 & \text{다르면 } 1 \end{cases}$ $e(\hat{y}, y) = (\hat{y} - y)^2$
- ▶ NP-hard in general ▶ efficient closed-form solution

- linear regression algorithm (1-step learning): $\boxed{\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y}}$
- 1 step learning.
+ 일반화도 가능하다