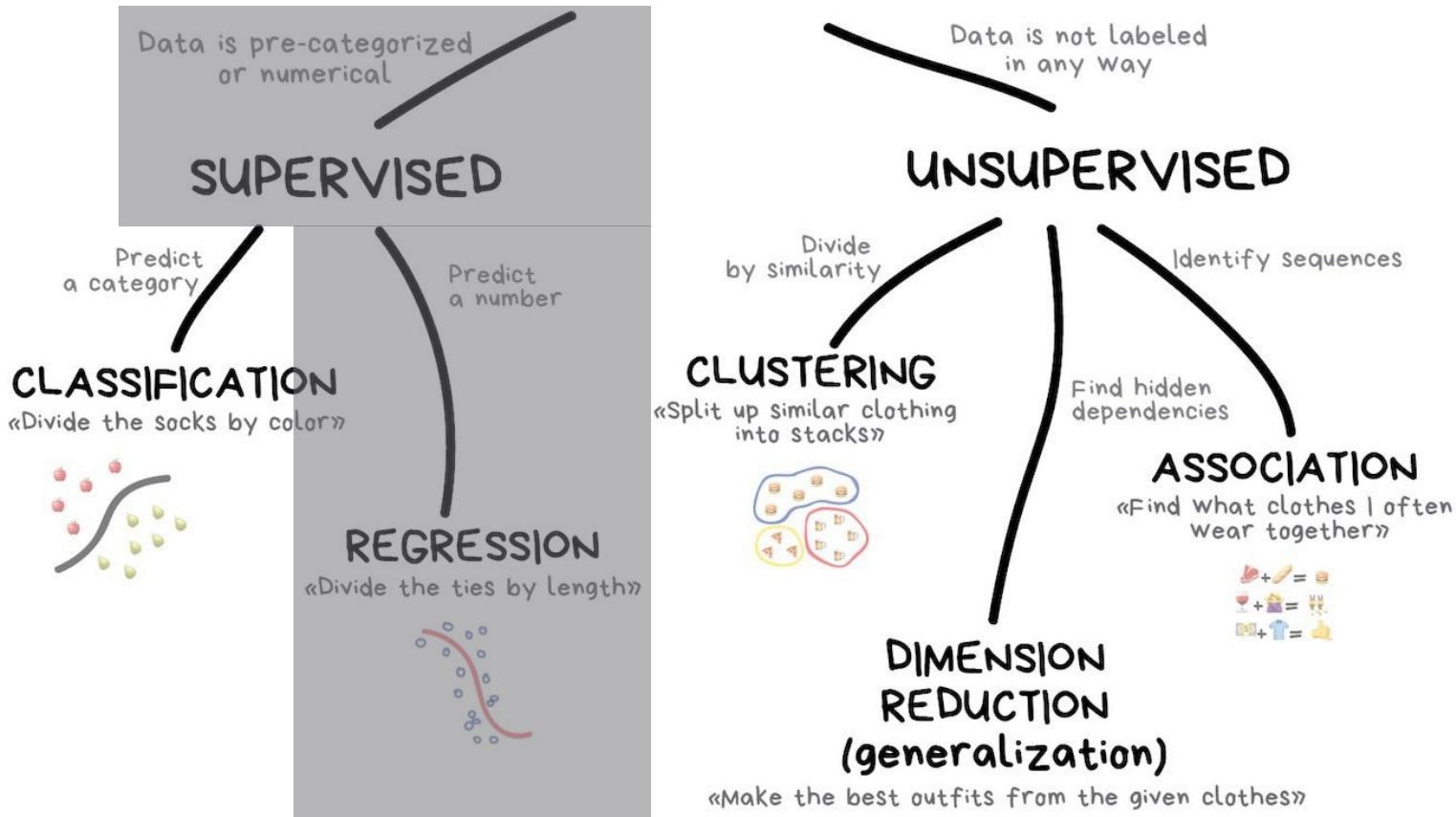


LINEAR REGRESSION

Week04

Topics Covered in This Class

CLASSICAL MACHINE LEARNING



Linear Regression

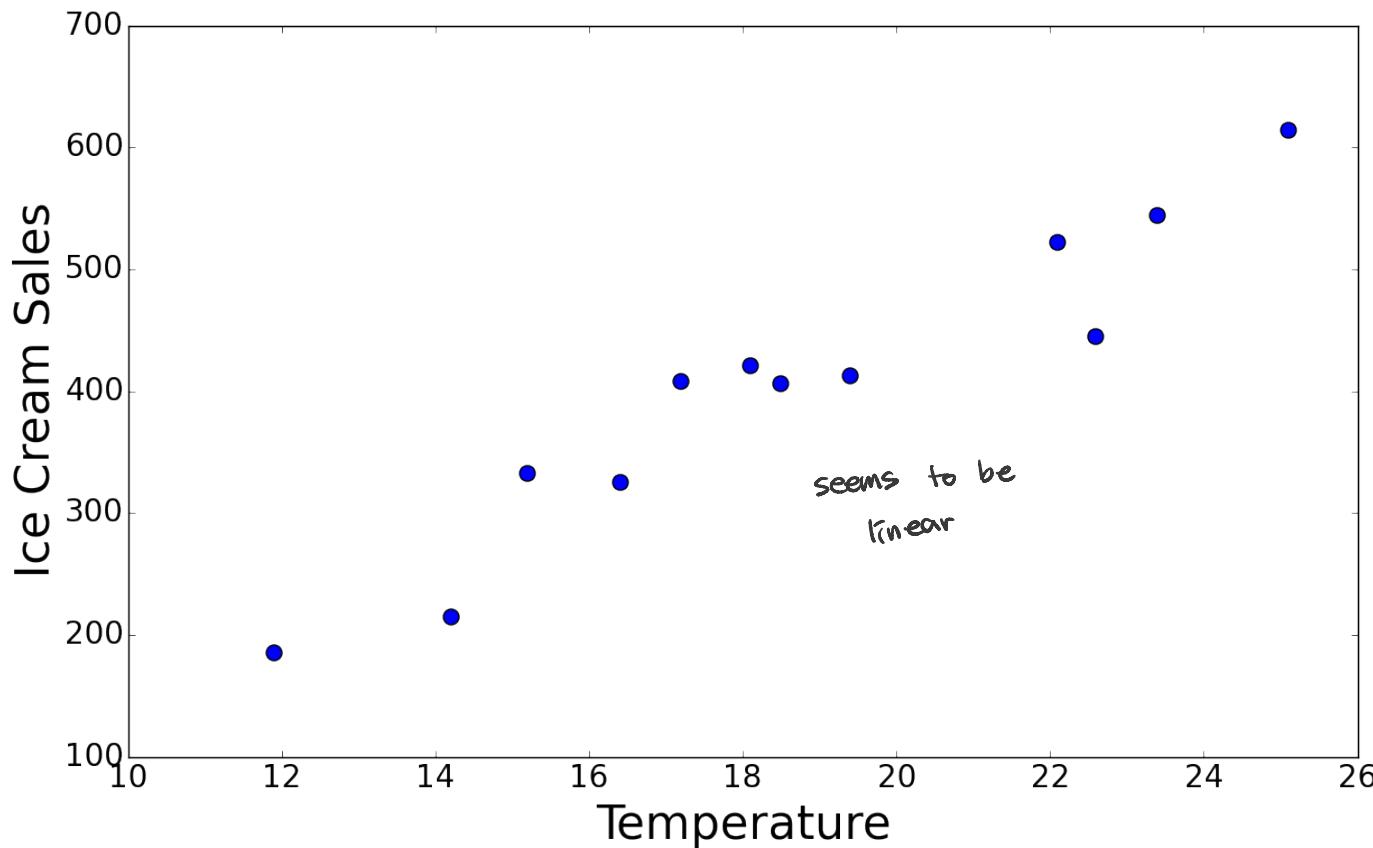
Supervised Learning: Regression

설정변수와 output 변수가
시정변수로 보일 때 적용

- Prediction ice cream sales over given temperature

$$\text{ice cream sales} = f(\text{temperature})$$

predict ice cream sales ← new temperature



Temperature (°C)	Ice Cream Sales (\$)
14.2	215
16.4	325
11.9	185
15.2	332
18.5	406
22.1	522
19.4	412
25.1	614
23.4	544
18.1	421
22.6	445
17.2	408

Linear Regression

- Linear regression
 - ▣ Based on the assumption that the relationship between a scalar dependent variable y and explanatory(independent) variables X is linear
 - ▣ $X = [x_1, x_2, x_3, \dots, x_n]$
Explanatory variables: print run(x_1), page number(x_2)

x_1	x_2
2800	22
2670	14
2800	37
2784	15
2800	38

→ X

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \epsilon$$

$\epsilon \sim N(0, \sigma^2)$

noise follow normal

distribution

regardless of sample input variable.

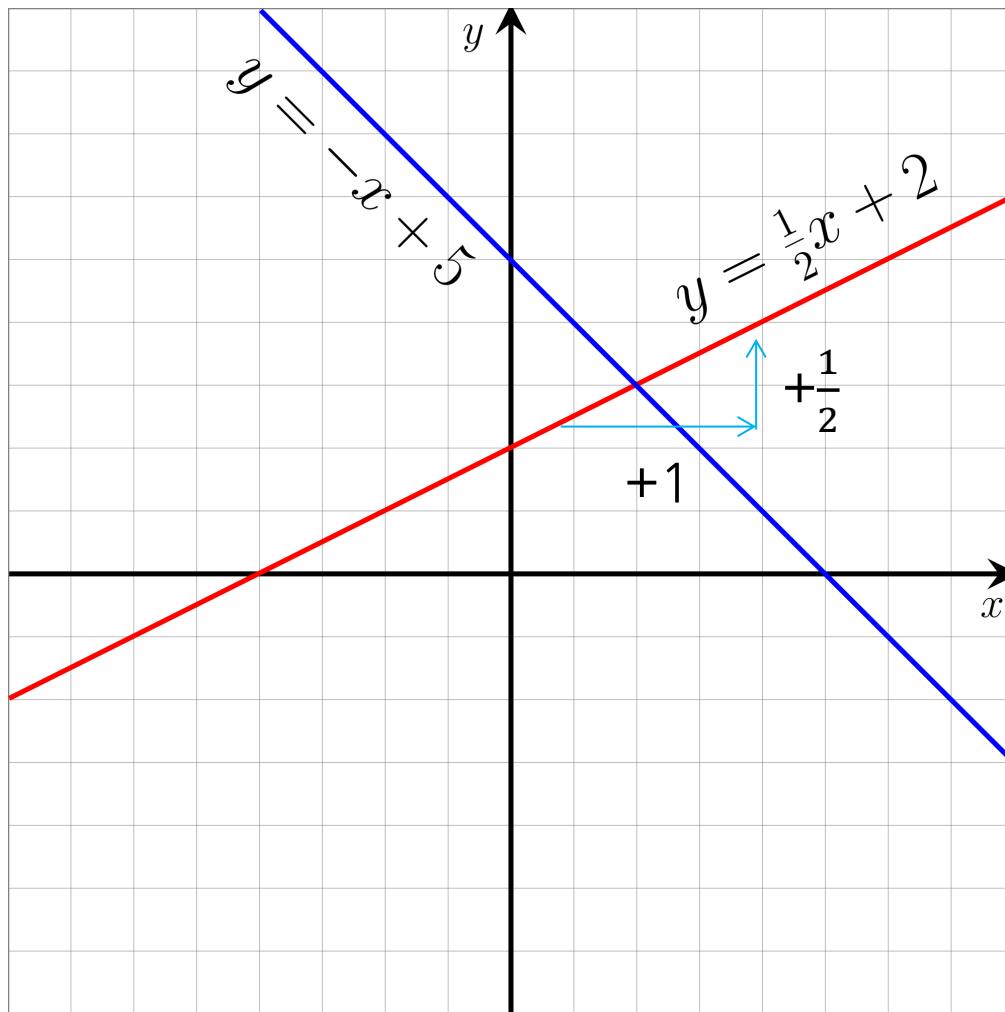
한정된 y에 미치는 영향이 한정적이라고
설정한 관계

describe certain
noise error

p number of the input variables
only one scalar explain the relation
between input, output

* Linear function

- You studied linear function when you are high school student!



$$y = \frac{1}{2}x + 5$$

slope
→ coefficient
HIP
linear regression

intercept

Main Assumptions of Linear Regression

선형 회귀 분석이 제대로 적용되려면
매개지 통계적 가정이 충족되어야 함

- Linear regression analysis makes several key assumptions

input variable has

- Linear relationship with target variable

$$E(Y) = \beta_1 x + \beta_0$$

독립(인연) 변수와 종속 변수 간에는 선형적인 관계가 있어야 함
오차의 분산이 모든 양의 값에 대해 동일, 대비로의 어느 지점에서든 오차의 표준(분산)이 일정해야 함

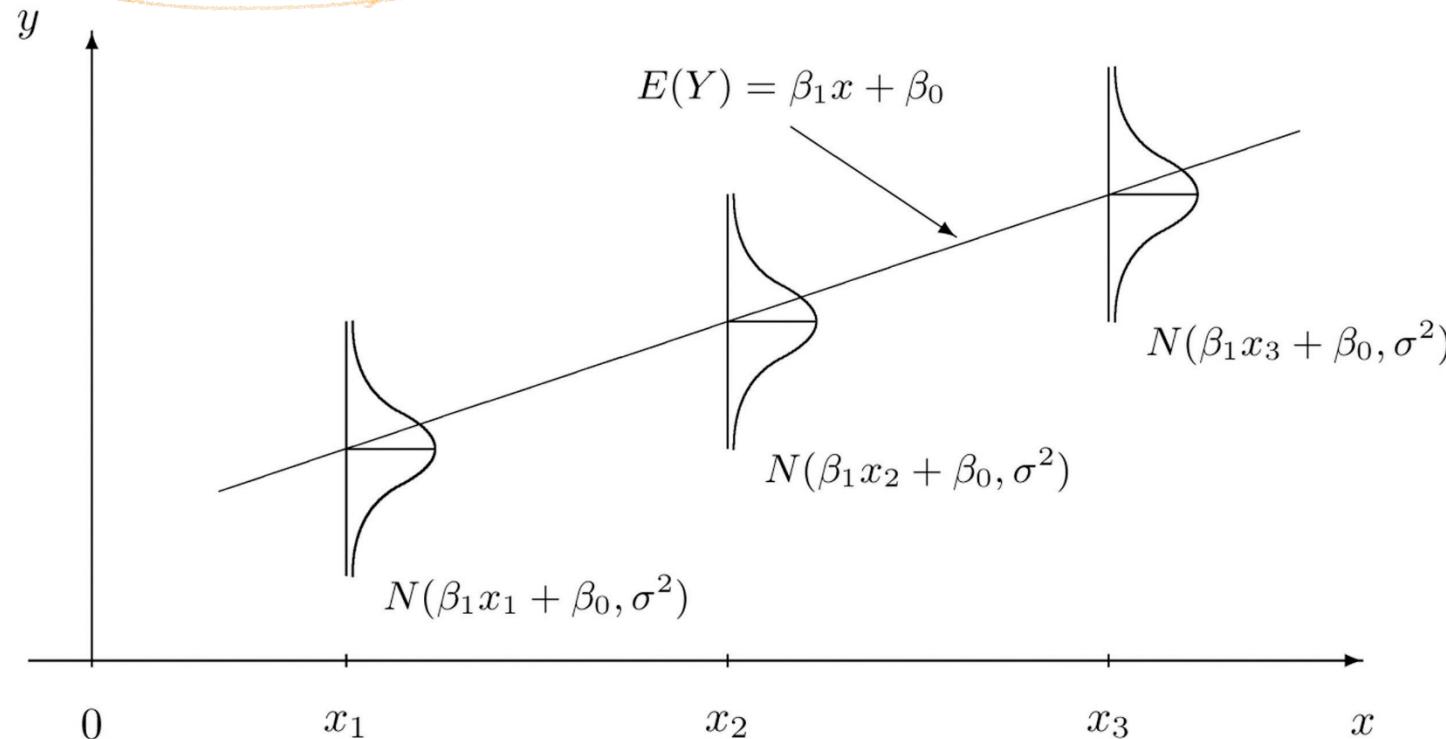
$$N(\beta_1 x + \beta_0, \sigma^2)$$

- Homoscedasticity → error term ~ normal distribution with fixed variance

- Normality 정규성 오차의 정규분포를 따른다는 가정

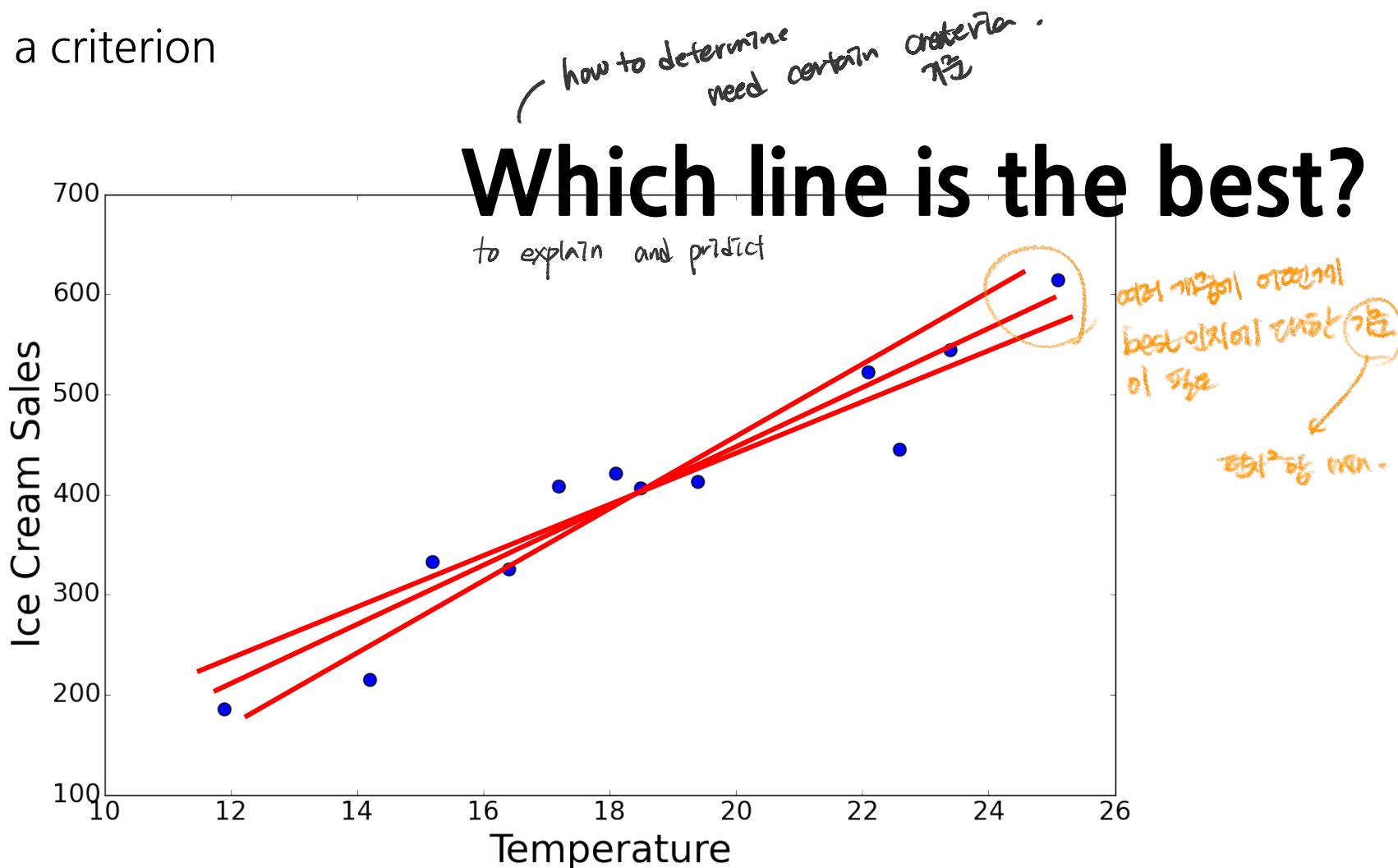
- No or little multicollinearity

noise term in linear regression



How to Determine Relationship between x and y

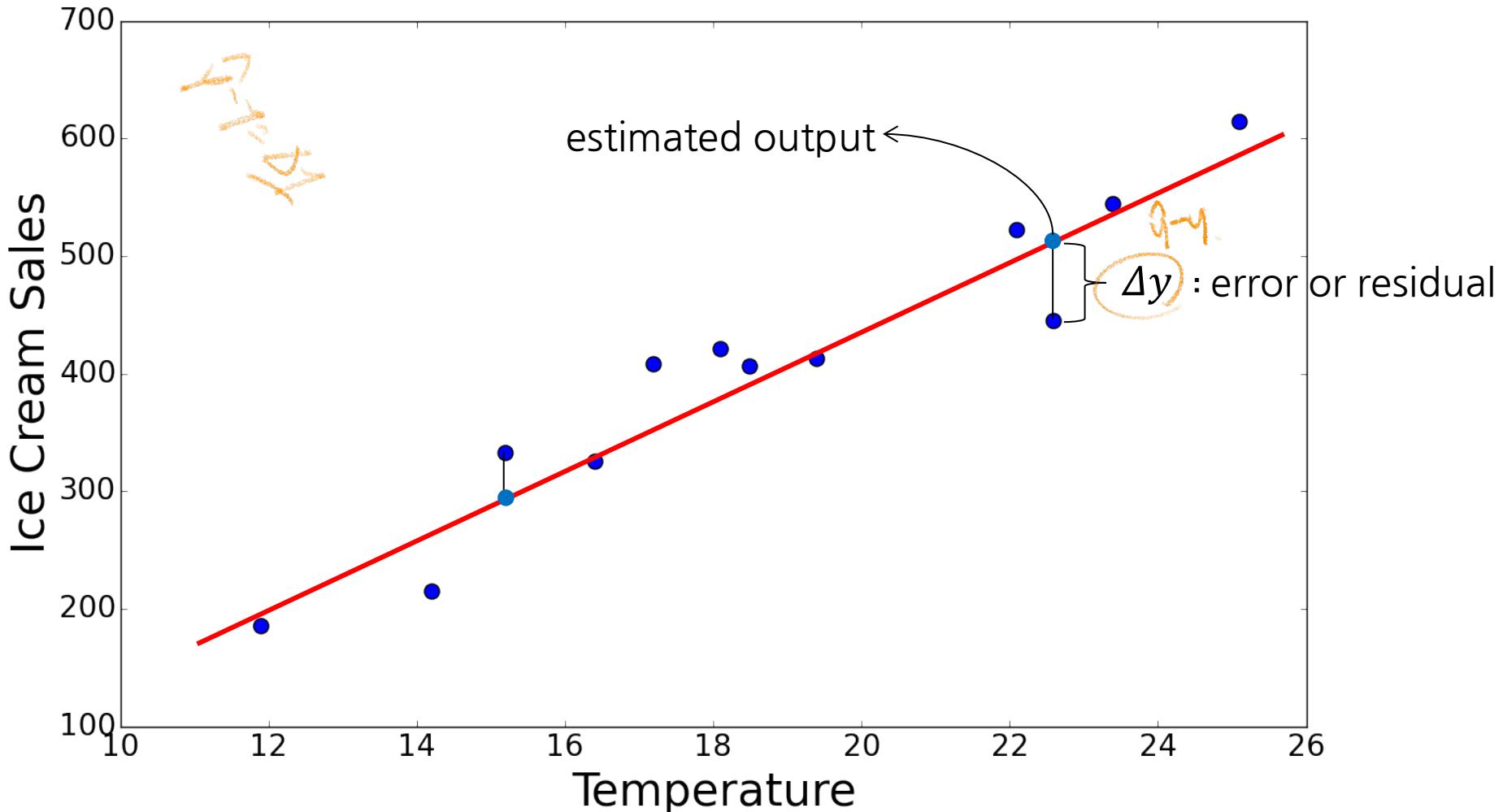
- Need a criterion



Least Square Method

creatin

- Minimize summation of squared error
 - Squared error=(estimated output-real output)²=error²



Least Square Method

- Minimize summation of squared error

$$\sum_i (y_i - \hat{y}_i)^2$$

Sum for all data points in train set

how can we change the sign of error

estimated output

real output

- In the simple case: Only one independent variable
 - Estimated output $\hat{y}_i = \beta_0 + \beta_1 x_i$

$$\min \sum_i (y_i - \beta_0 - \beta_1 x_i)^2$$

Unknown values

Known values

※ Summary for Notation

- Hat, ($\hat{}$)
 - ▣ Represents estimation
 - ▣ β_1 is unknown true value, $\hat{\beta}_1$ is estimation for β_1 through model learning
 - ▣ y_i is known output value of i -th sample, \hat{y}_i is estimated output by learned model
- Bar, ($\bar{}$)
 - ▣ Represents sample mean
 - ▣ Arithmetic average of the observed values of variable

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

- ▣ If the number of input variables is more than one, elements of sample mean vector consist of average of each variable

$$\bar{\mathbf{x}} = \left(\frac{\sum_{i=1}^n x_{1i}}{n}, \frac{\sum_{i=1}^n x_{2i}}{n}, \dots, \frac{\sum_{i=1}^n x_{pi}}{n} \right) = \frac{\sum_{i=1}^n \mathbf{x}_i}{n}$$

- Bold character usually represents vector

Least Square Method

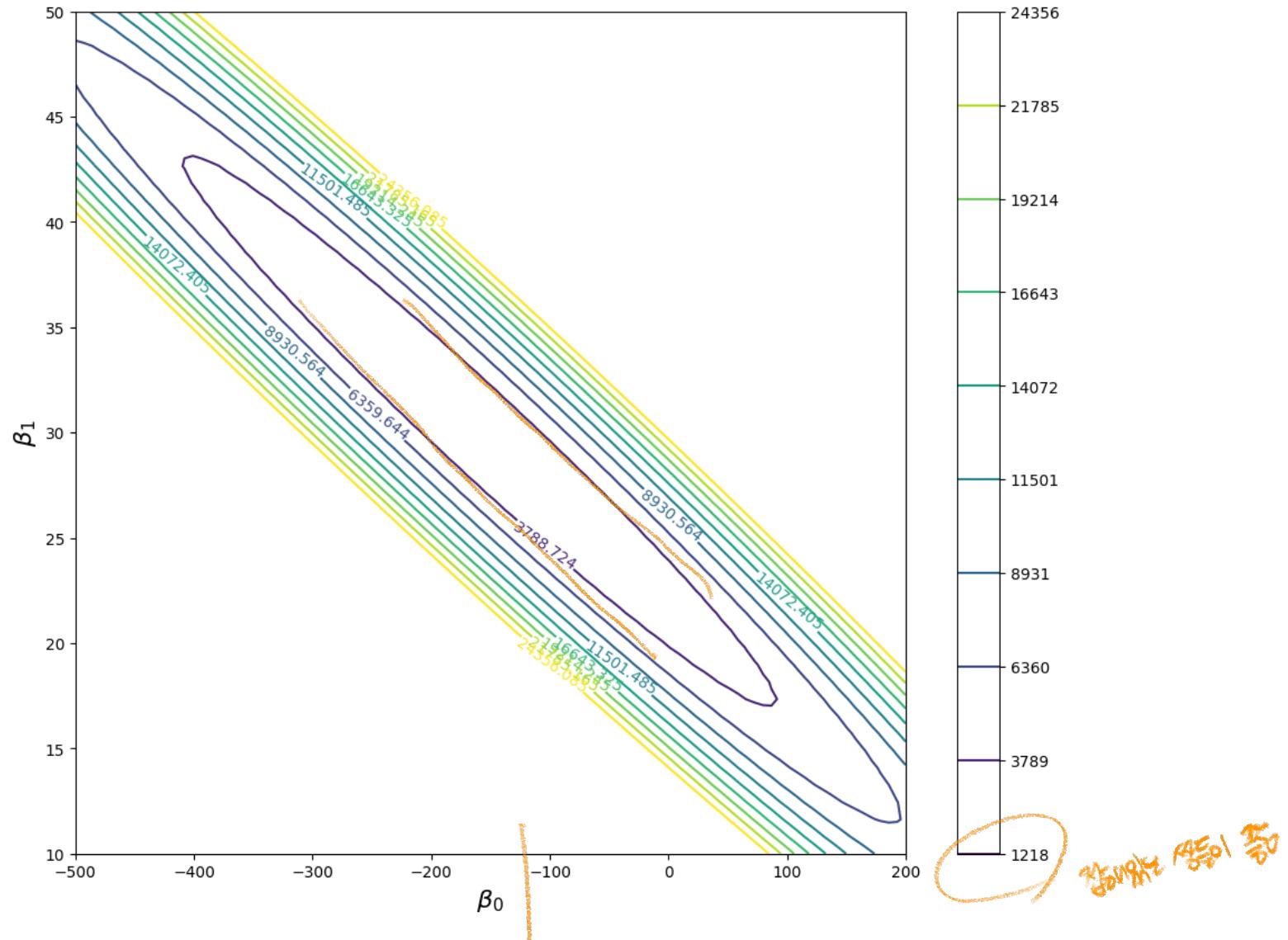
- Which one is better?

Temperature (°C)	Ice Cream Sales (\$)	$\beta_1 = 20, \beta_0 = -80$		$\beta_1 = 30, \beta_0 = -160$	
		Estimated Sales	Squared error	Estimated Sales	Squared error
14.2	215	204	121	266	2601
16.4	325	248	5929	332	49
11.9	185	158	729	197	144
15.2	332	224	11664	296	1296
18.5	406	290	13456	395	121
22.1	522	362	25600	503	361
19.4	412	308	10816	422	100
25.1	614	422	36864	593	441
23.4	544	388	24336	542	4
18.1	421	282	19321	383	1444
22.6	445	372	5329	518	5329
17.2	408	264	20736	356	2704
sum		174901	>	14594	

이거 더 높다

Least Square Method

- Summation of squared error with different β_0 and β_1



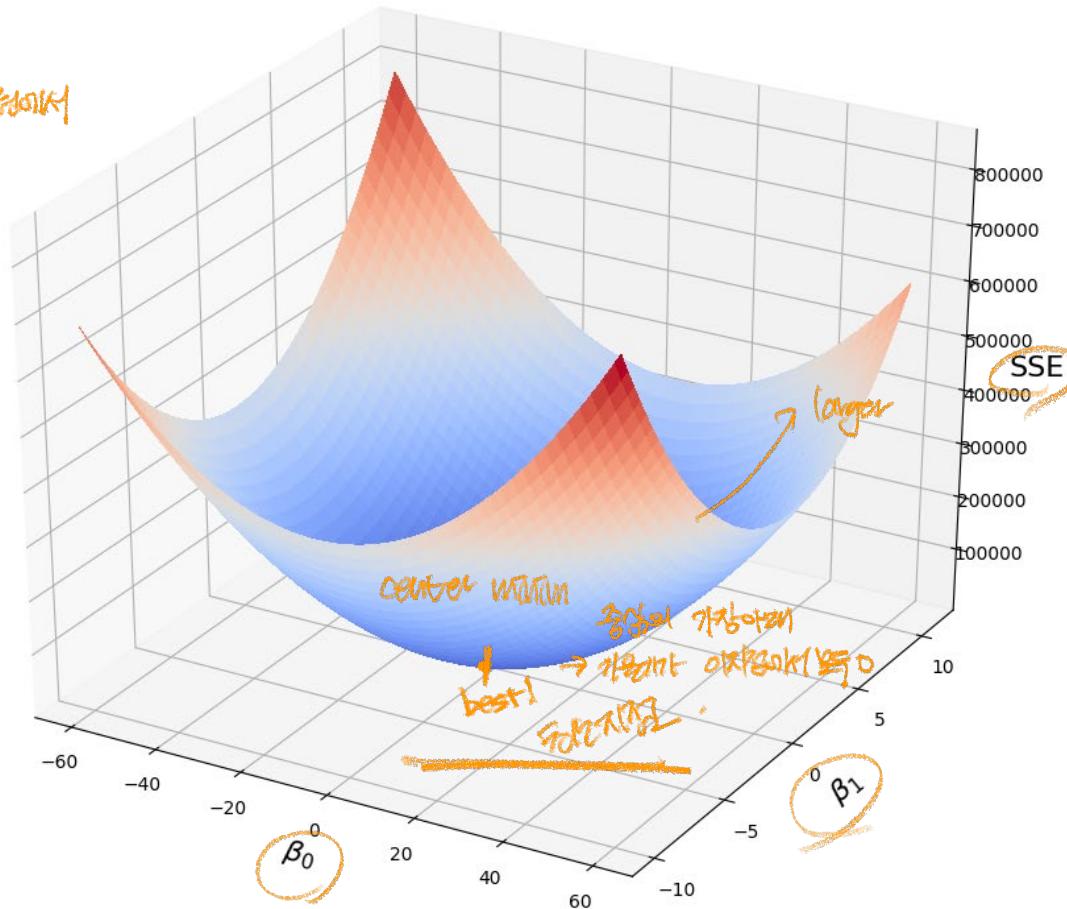
Least Square Method

- Summation of squared error with different β_0 and β_1 for simulated data from $y = x + 1$

sum of square error

beta의 서로 다른 값들에서

인



Optimization for Linear Regression

- Variables to be determined

β_0, β_1

- Objective function

$$\min f(\beta_0, \beta_1) = \min \sum_i (y_i - \beta_0 - \beta_1 x_i)^2$$

- Constraints

- No constraint

Optimization for Linear Regression

- Solution

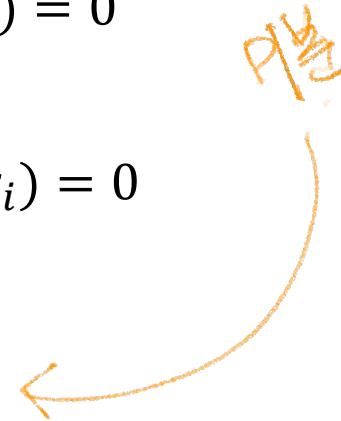
- ▣ Calculate partial derivatives with respect to β_0, β_1

$$\frac{\partial f(\beta_0, \beta_1)}{\partial \beta_0} = \sum_{i=1}^n -2(y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\frac{\partial f(\beta_0, \beta_1)}{\partial \beta_1} = \sum_{i=1}^n -2x_i(y_i - \beta_0 - \beta_1 x_i) = 0$$

- ▣ Solve linear equations

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$



$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Multiple Input Variables

- More than one input variable
 - Want to predict consumption of petrol

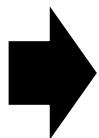
Petrol Tax(\$)	Average Income (\$)	Paved Highways (miles)	Proportion of population with driver's license	Consumption of petrol (M of gallons)
9	3571	1976	0.525	541
9	4092	1250	0.572	524
9	3865	1586	0.58	561
7.5	4870	2351	0.529	414
8	4399	431	0.544	410
10	5342	1333	0.571	457
8	5319	11868	0.451	344
8	5126	2138	0.553	467
8	4447	8577	0.529	464
7	4512	8507	0.552	498
...

Multiple Input Variables

- Estimation based on petrol tax, average income, length of paved highways, proportion of population with driver's license

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \epsilon$$

- y = consumption of petrol
- x_1 = petrol tax
- x_2 = average income
- x_3 = length of paved highways
- x_4 = proportion of population with driver's license
- ϵ is random error which follows Gaussian distribution with 0 mean, σ^2 variance



$$\min \sum_i (y_i - \hat{y}_i)^2$$

Same as the simple case!

Optimization for Linear Regression: Multivariate

- multivariate linear regression

$$\min f(\beta_0, \dots, \beta_p) = \min \sum_i (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_p x_{pi})^2$$

- Estimated parameters are obtained by setting partial derivatives zero

$$\frac{\partial f(\beta_0, \dots, \beta_p)}{\partial \beta_0} = \sum_{i=1}^n -2(y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_p x_{pi}) = 0$$

$$\frac{\partial f(\beta_0, \dots, \beta_p)}{\partial \beta_1} = \sum_{i=1}^n -2x_{1i}(y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_p x_{pi}) = 0$$

⋮

$$\frac{\partial f(\beta_0, \dots, \beta_p)}{\partial \beta_p} = \sum_{i=1}^n -2x_{pi}(y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_p x_{pi}) = 0$$

Eq 8 PM

Multiple Input Variables

- Matrix approach to multiple regression model

$$y_1 = \beta_0 \cdot 1 + \beta_1 x_{11} + \beta_2 x_{21} + \cdots + \beta_p x_{k1}$$

$$y_2 = \beta_0 \cdot 1 + \beta_1 x_{12} + \beta_2 x_{22} + \cdots + \beta_p x_{k2}$$

⋮

$$y_n = \beta_0 \cdot 1 + \beta_1 x_{1n} + \beta_2 x_{2n} + \cdots + \beta_p x_{kn}$$

n samples, p input variables

*this is linear regression
diagram*

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{21} & \cdots & x_{p1} \\ 1 & x_{12} & x_{22} & \cdots & x_{p2} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{pn} \end{bmatrix}$$

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{21} & \cdots & x_{p1} \\ 1 & x_{12} & x_{22} & \cdots & x_{p2} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{1n} & x_{2n} & \cdots & x_{pn} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\Rightarrow \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad \boldsymbol{\epsilon} \sim MVN(0, \sigma^2 \mathbf{I})$$

$$E = \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2$$

Optimization for Linear Regression: Multivariate

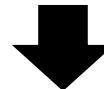
$$\min E = \min \|y - X\beta\|^2$$

$$\begin{aligned} & \min \|y - X\beta\|^2 \\ &= (y - X\beta)^T (y - X\beta) \end{aligned}$$

$(x - As)^T w (x - As)$
 $x \rightarrow Y$
 $A \rightarrow X$
 $s \rightarrow \beta$
 $w \rightarrow L$

- Solution is obtained by setting $\frac{\partial E}{\partial \beta} = 0$

$$\frac{\partial (x - As)^T W (x - As)}{\partial s} = -2A^T W (x - As)$$



$$\frac{\partial (y - X\beta)^T (y - X\beta)}{\partial \beta} = -2X^T (y - X\beta) = 0$$



$$X^T X \beta = X^T y$$

$$X : n \times (p+1) \\ X^T : (p+1) \times n$$

$$X^T X : (p+1) \times (p+1)$$

$$X^T X \beta - X^T y = 0$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

을 통해 얻어지는
는 선형 방정식의 해

- Reference

- Matrix Cookbook

- <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>

Estimation of Regression Coefficients: Multivariate

- Use least square methods as same as simple linear regression

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

양기파여

- \mathbf{X}^T : transpose matrix of \mathbf{X}
- \mathbf{X}^{-1} : inverse matrix of \mathbf{X}

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\beta} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H} \mathbf{y}$$

- Residual(error) terms

$$\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{H}) \mathbf{y} = [e_1 \ e_2 \ \cdots \ e_n]^T$$

- Covariance

$$\begin{aligned} \text{Cov}[\mathbf{e}] &= \text{Cov}[\mathbf{y} - \hat{\mathbf{y}}] = \text{Cov}[(\mathbf{I} - \mathbf{H}) \mathbf{y}] = \text{Cov}[(\mathbf{I} - \mathbf{H})(\mathbf{X} \beta + \epsilon)] \\ &= \text{Cov}[(\mathbf{I} - \mathbf{H}) \epsilon] = (\mathbf{I} - \mathbf{H}) \text{Cov}[\epsilon] (\mathbf{I} - \mathbf{H})^T = \sigma^2 (\mathbf{I} - \mathbf{H}) (\mathbf{I} - \mathbf{H})^T = \sigma^2 (\mathbf{I} - \mathbf{H}) \end{aligned}$$

- SSE

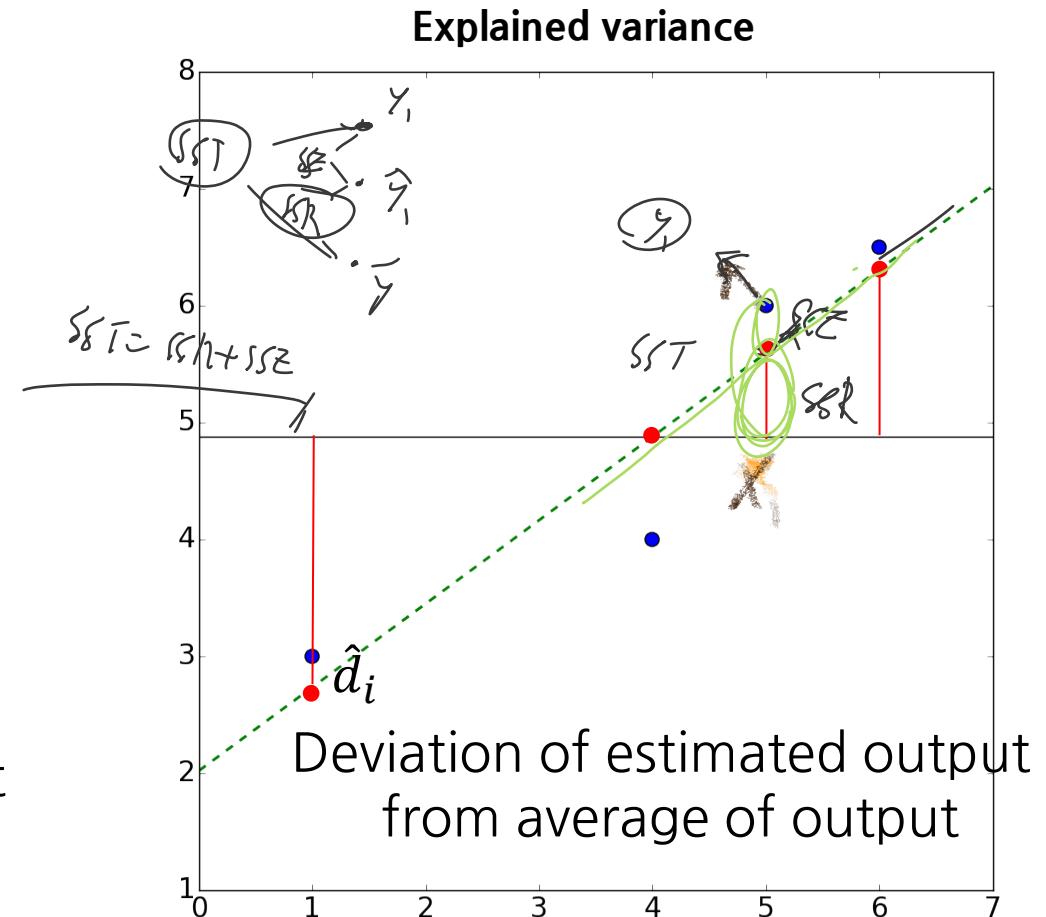
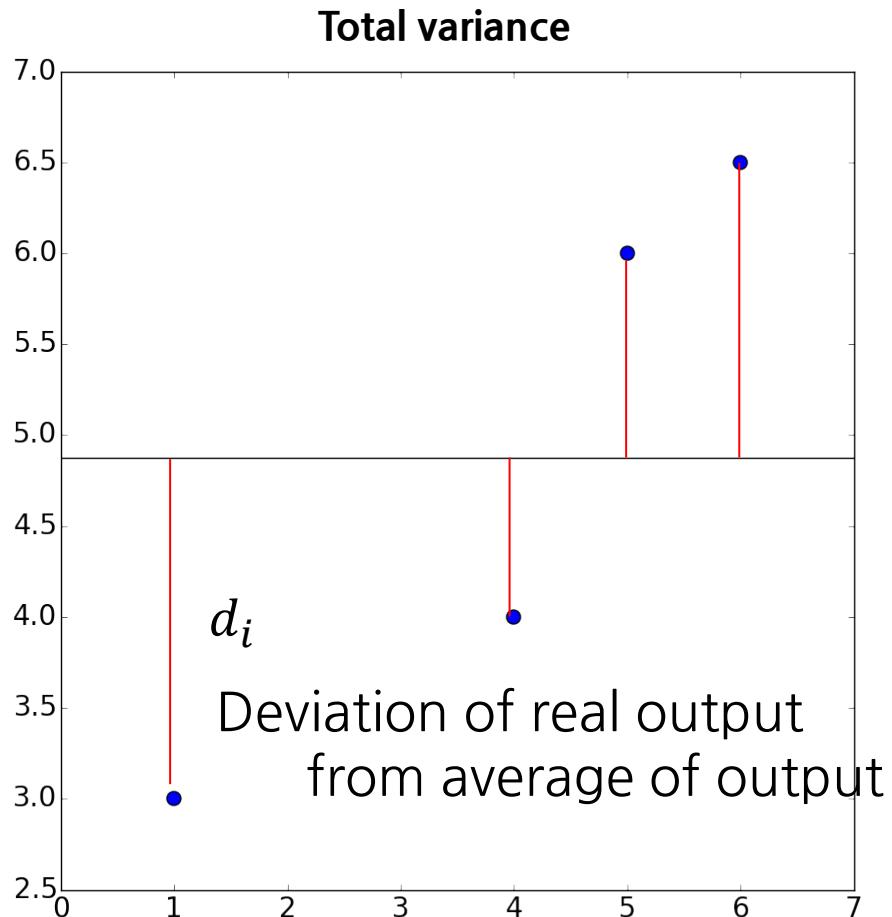
$$SSE = \sum_i (y_i - \hat{y}_i)^2 = \sum_i e_i^2 = \mathbf{e}^T \mathbf{e}$$

Is the Regression Model Significant?

- Modeling learning is not the end of the analysis
 - ▣ Check overall significance in regression models
 - Whether the regression model is overall significant for predicting a target
 - ▣ Check significance of regression coefficients
 - Whether the specific variable is significant for predicting a target
- In the case of simple linear regression, testing overall significance of the model is the same as testing significance of regression coefficients
 - ▣ Because only one explanatory variable is used

Sum of Square

- Explained variance(SSR) and Total variance(SST)



Sum of Square



양수(여기)

이 가지를 가지고 f-test를 한
↓

유의적정성이 얼마나
유의미한지 -

- Total variance: the total sum of squares

$$SST = \sum_i (y_i - \bar{y})^2$$

- Explained variance: the regression sum of squares, also called the explained sum of squares

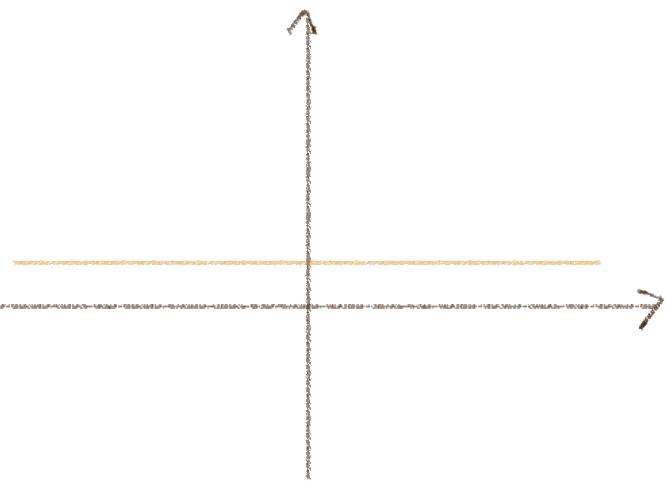
$$SSR = \sum_i (\hat{y}_i - \bar{y})^2$$

- Residual variance: the sum of squares of residuals, also called the residual sum of squares

$$SSE = \sum_i (y_i - \hat{y}_i)^2$$

- Relationship among three values

$$SST = SSR + SSE$$
$$= \sum (y_i - \bar{y})^2$$



Test of Model Significance

- F -test for general regression models

- Hypothesis

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_p = 0 \quad \text{가설이 뜻하는 것은 } X \text{와 } Y \text{ 사이에 선형관계가 없다!}$$

$$H_1: \text{not all } \beta_i (i = 1, 2, \dots, p) \text{ equal zero} \quad \rightarrow \text{선형관계가 있는지를 주장하는 것!}$$

- Test statistic

$$F^* = \frac{\text{MSR}}{\text{MSE}}$$

- F follows F -distribution with $(p, n - p - 1)$ degree of freedom

- Decision rule

If $F^* \leq F(1 - \alpha; p, n - p - 1)$, conclude H_0

If $F^* > F(1 - \alpha; p, n - p - 1)$, conclude H_1

- α : significance level

* Statistical Test

통계적 검증

- A statistical test provides a mechanism for making quantitative decisions about a process or processes
 - ▣ The intent is to determine whether there is enough evidence to "reject" a conjecture or hypothesis about the process
 - ▣ The procedure is based on how likely it would be for a set of observations to occur if the null hypothesis were true
- Null hypothesis ~~추정~~
 - ▣ A general statement or default position that there is no relationship between two measured phenomena, or no association among groups
기본가설, 전제가설, 무연관설정, 관찰가설
- Alternative hypothesis ~~대립~~
 - ▣ It is the hypothesis used in hypothesis testing that is contrary to the null hypothesis

차이가설

* Statistical Test

- Steps in testing for statistical significance

State the research hypothesis

State the null hypothesis & alternative hypothesis

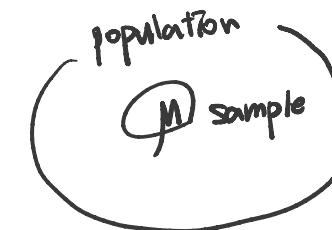
Select a probability of error level (alpha level)

Select and compute the test for statistical significance

Interpret the results

※ Statistical Test

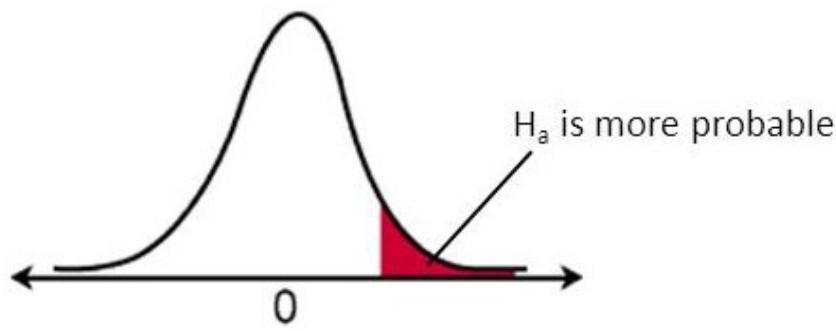
- Consider 20 first year resident female doctors drawn at random from one area
 - ▣ resting systolic blood pressures measured using an electronic sphygmomanometer
 - Sample mean = 130.05
 - ▣ Research hypothesis is that a resting systolic blood pressure of 120 mm Hg is predicted as the population mean
 - ▣ Null Hypothesis
- $$H_0: \mu = 120$$
- $$H_1: \mu \neq 120$$



$H_0: \mu = 120$
 $H_1: \mu < 120$
 $H_1: \mu > 120$

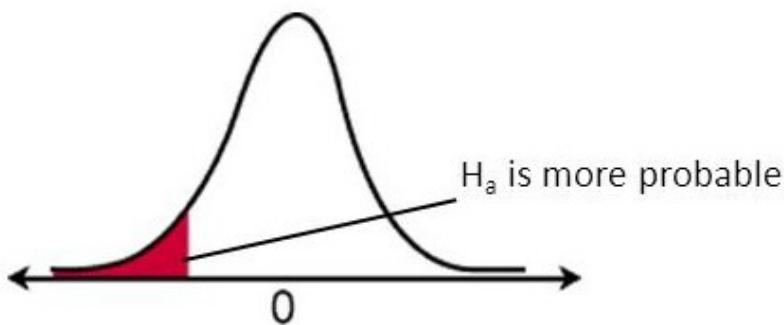
- ▣ Set significance level as 0.05
- ▣ Determine test statistics and underlying distribution
 - $t = \frac{\bar{x} - \mu}{\sqrt{s^2/n}}$
 - t follows t -distribution with the degree of freedom as $n - 1$

* Statistical Test



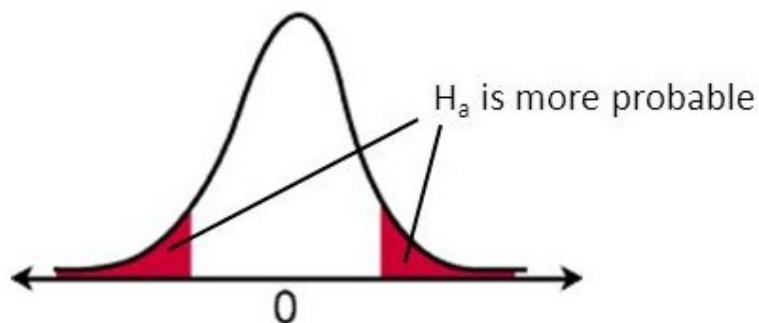
Right-tail test

$$H_a: \mu > \text{value}$$



Left-tail test

$$H_a: \mu < \text{value}$$



Two-tail test

$$H_a: \mu \neq \text{value}$$

Test of Model Significance

- ANOVA table for multiple regression model with p input variables

Factor	Sum of square	Degree of freedom	Mean square	F-value	p-value
Model	SSR	p	$MSR = SSR/p$	$F_0 = MSR/MSE$	$P\{F_{p,n-p-1} > F_0\}$
Residual	SSE	$n - p - 1$	$MSE = SSE/(n - p - 1)$		soil depth soil texture
Total	SST	$n - 1$			

- Analysis of Variance (ANOVA)

Degree of Freedom

- The number of degrees of freedom is the number of values in the final calculation of a statistic that are free to vary
 - ▣ The number of independent ways by which a dynamic system can move, without violating any constraint imposed on it
- Sample variance

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

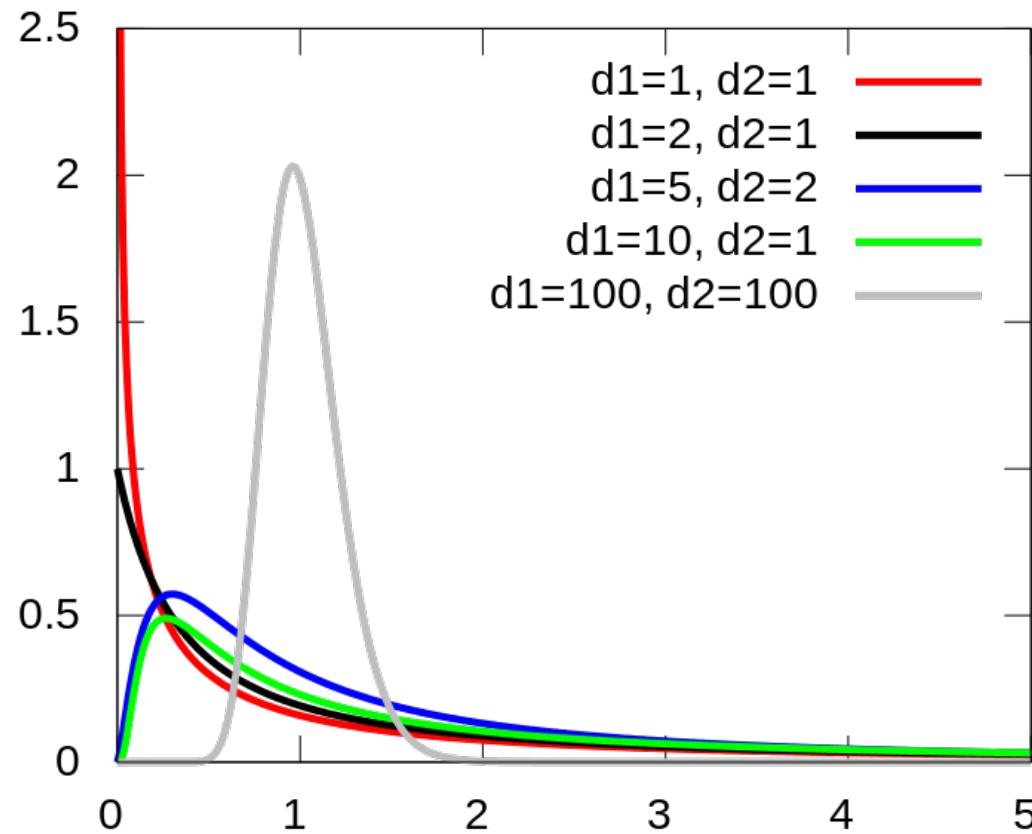
- ▣ The reason that denominator is $n - 1$ is that degree of freedom of sample mean, \bar{y} is $n - 1$
- ▣ Another reason is that in the case of that denominator is $n - 1$, S^2 is unbiased estimator of variance of population
- Mean squared error for simple linear regression

$$MSE = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- ▣ The reason that denominator is $n - 2$ is that \hat{y}_i is calculate from $\hat{\beta}_0 + \hat{\beta}_1 x_i$ and it depends on two estimators $\hat{\beta}_0, \hat{\beta}_1 \rightarrow$ Decrease two degrees of freedom

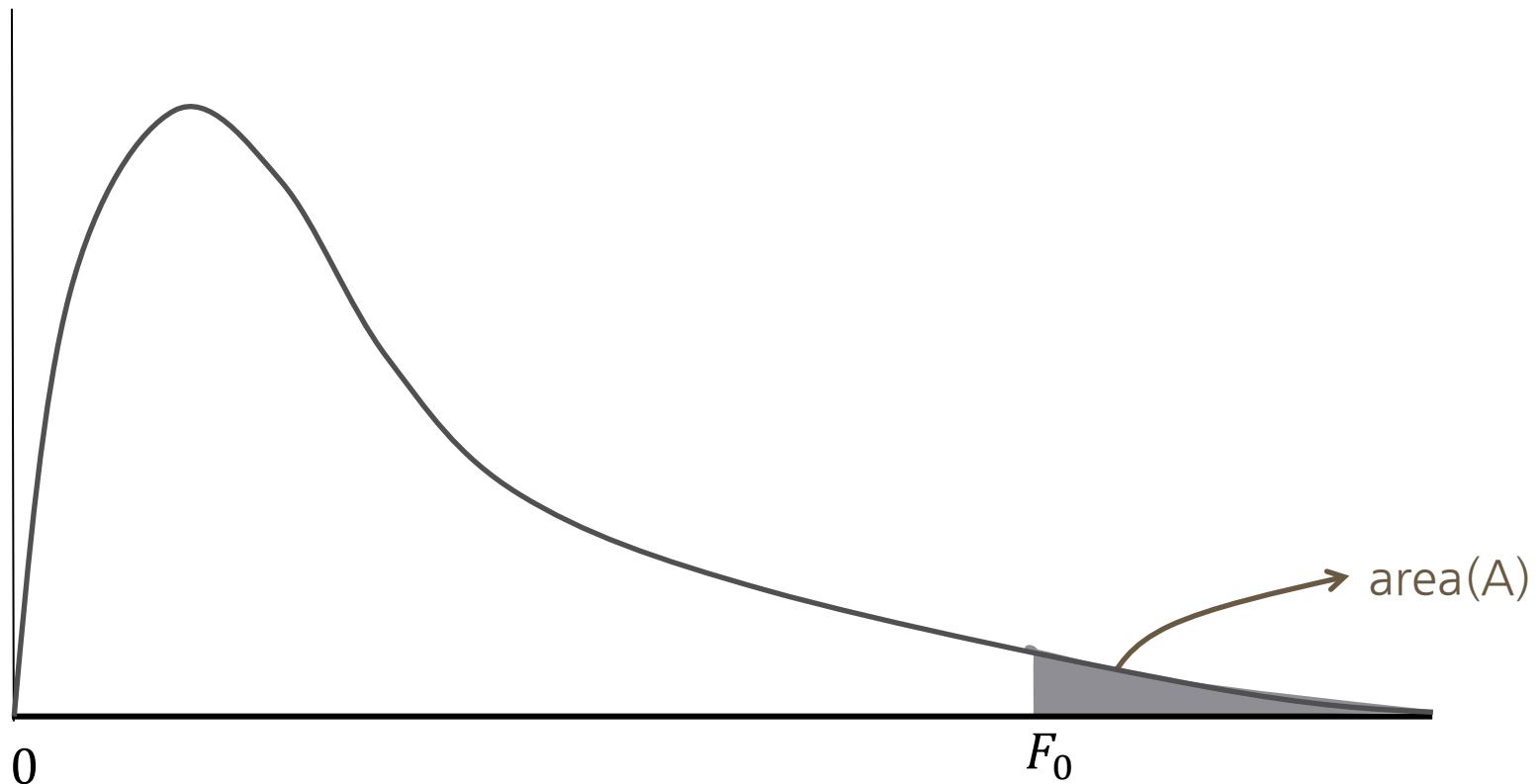
※ F distribution

- F statistics follows F distribution with $(p, n - p - 1)$ degree of freedom
 - Probability density function of F distribution with different parameters
 - F distribution is determined by two parameters



Test of Model Significance

- If (area under density function from F_0 to ∞) $< \alpha$
→ Reject null hypothesis → not all $\beta_i (i = 1, 2, \dots, p)$ equal zero
 - ▣ α is significance level
 - ▣ significance level is usually set to 0.1, 0.05, or 0.01
 - The higher significance level, the higher probability to reject null hypothesis



※ Table for Distributions

F-Distribution, Continued
Upper 0.01 Critical Points

		$F_{0.01}(r_1, r_2)$								
		r_1								
r_2		10	15	20	25	30	40	60	120	∞
1	6055.9	6157.3	6208.7	6239.8	6260.7	6286.8	6313.0	6339.4	6365.9	
2	99.40	99.43	99.45	99.46	99.47	99.47	99.48	99.49	99.50	
3	27.23	26.87	26.69	26.58	26.50	26.41	26.32	26.22	26.13	
4	14.55	14.20	14.02	13.91	13.84	13.75	13.65	13.56	13.46	
5	10.05	9.72	9.55	9.45	9.38	9.29	9.20	9.11	9.02	
6	7.87	7.56	7.40	7.30	7.23	7.14	7.06	6.97	6.88	
7	6.62	6.31	6.16	6.06	5.99	5.91	5.82	5.74	5.65	
8	5.81	5.52	5.36	5.26	5.20	5.12	5.03	4.95	4.86	
9	5.26	4.96	4.81	4.71	4.65	4.57	4.48	4.40	4.31	
10	4.85	4.56	4.41	4.31	4.25	4.17	4.08	4.00	3.91	
11	4.54	4.25	4.10	4.01	3.94	3.86	3.78	3.69	3.60	
12	4.30	4.01	3.86	3.76	3.70	3.62	3.54	3.45	3.36	
13	4.10	3.82	3.66	3.57	3.51	3.43	3.34	3.25	3.17	
14	3.94	3.66	3.51	3.41	3.35	3.27	3.18	3.09	3.00	
15	3.80	3.52	3.37	3.28	3.21	3.13	3.05	2.96	2.87	
16	3.69	3.41	3.26	3.16	3.10	3.02	2.93	2.84	2.75	
17	3.59	3.31	3.16	3.07	3.00	2.92	2.83	2.75	2.65	
18	3.51	3.23	3.08	2.98	2.92	2.84	2.75	2.66	2.57	
19	3.43	3.15	3.00	2.91	2.84	2.76	2.67	2.58	2.49	
20	3.37	3.09	2.94	2.84	2.78	2.69	2.61	2.52	2.42	
21	3.31	3.03	2.88	2.79	2.72	2.64	2.55	2.46	2.36	
22	3.26	2.98	2.83	2.73	2.67	2.58	2.50	2.40	2.31	
23	3.21	2.93	2.78	2.69	2.62	2.54	2.45	2.35	2.26	
24	3.17	2.89	2.74	2.64	2.58	2.49	2.40	2.31	2.21	
25	3.13	2.85	2.70	2.60	2.54	2.45	2.36	2.27	2.17	
26	3.09	2.81	2.66	2.57	2.50	2.42	2.33	2.23	2.13	
27	3.06	2.78	2.63	2.54	2.47	2.38	2.29	2.20	2.10	
28	3.03	2.75	2.60	2.51	2.44	2.35	2.26	2.17	2.06	
29	3.00	2.73	2.57	2.48	2.41	2.33	2.23	2.14	2.03	
30	2.98	2.70	2.55	2.45	2.39	2.30	2.21	2.11	2.01	
40	2.80	2.52	2.37	2.27	2.20	2.11	2.02	1.92	1.80	
60	2.63	2.35	2.20	2.10	2.03	1.94	1.84	1.73	1.60	
120	2.47	2.19	2.03	1.93	1.86	1.76	1.66	1.53	1.38	
∞	2.32	2.04	1.88	1.77	1.70	1.59	1.47	1.32	1.00	