

Linear Regression & Gradient Descent

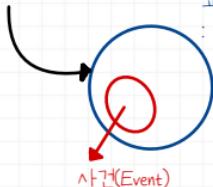
Date

1. Static Review

1) 평균, 분산, 표준 편차

실험(Experiment)

: 결과물을 만드는 행위



표본공간(Sample Space)

: 가능한 모든 결과

- : Sample space의 부분 집합
- . 특정 결과를 가르킴

평균 (mean)

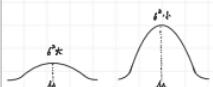
: 기대값(Expectation), μ

$$E(\bar{X}) = \sum X_i P(X_i) : \text{Discrete Data}$$

$$\int x P(x) : \text{Continuous Data}$$

분산(Variance)

: Data들의 산포된 정도, σ^2



$$\sigma^2 = E(\bar{X}^2) - [E(\bar{X})]^2$$

$$\sigma^2 = \frac{\text{분산}}{\text{표준}}$$

예제) X 1 2 3			
P	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$
$E(\bar{X})$	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{3}{4}$
$E(\bar{X}^2)$	$\frac{1}{2}$	1	$\frac{9}{4}$
$\rightarrow \mu = \frac{1}{2}$	$\sigma^2 = \frac{33}{16}$	$\sigma = \sqrt{\frac{33}{16}}$	

2) Random Variable

: 랜덤 변수, 이름에 혼혹 X \rightarrow 랜덤 변수 = 함수 (function)

Random Variable

Discrete RV

\rightarrow 이산적 Data : PMF

Continuous RV

\rightarrow 연속적 Data : CDF, PDF

PMF (Probability Mass Function)

\rightarrow 확률질량함수

$$Ex) x = 1, 2, 3$$

$$P(X=x) = \frac{n}{6}$$

$$P(X=1) = \frac{1}{6}$$

$$P(X=2) = \frac{1}{6}$$

$$P(X=3) = \frac{1}{3}$$

CDF (Continuous Distribution Function)

: 누적 분포 함수

PDF (Probability Density Function)

: 확률 밀도 함수

$$P(X < a) = F(a)$$

$$P(X < b) = F(b)$$

$$P(a < X < b) = F(b) - F(a)$$

CDF
 $P(X \leq x) = F(x)$

3) Gaussian Distribution

→ Normal Distribution

4) Joint Probability

$$P(A, B) = P(A \cap B)$$

$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

= Gaussian PDF

$$= f_{xy} = P(X=x, Y=y)$$

= A, B가 동시에 일어나는 확률

$$\Phi(x|\mu, \sigma^2) = \int_{-\infty}^x N(x|\mu, \sigma^2) dx$$

5) Conditional Probability

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Given B이 대한 A의 확률

6) Bayes Rule

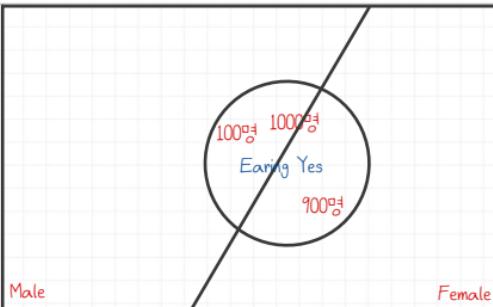
$$P(\theta|X) = \frac{\text{Likelihood} \ P(X|\theta) \ P(\theta)}{\text{Posterior} \ P(X)}$$

$$P(X=x|Y=y) = \frac{P(X=x)P(X=x|Y=y)}{\sum P(X=x)P(Y=y|X=x)}$$

$$P(\theta|X) \propto P(X|\theta)P(\theta)$$

Posterior Likelihood Prior

Ex)



$$\frac{P(\text{남}, \text{구})}{P(\text{구})} = P(\text{남}| \text{구}) = \frac{P(\text{구} | \text{남})P(\text{남})}{P(\text{구})}$$

$$P(\text{구}) = P(\text{구} | \text{남})P(\text{남}) + P(\text{남} | \text{구})P(\text{구})$$

나중에 분모 고려 X

: posterior 확률마다 분모는 동일하기에 고려 안함.

	남	남	Total
남	100	6900	7000
여	900	2100	3000
Total	1000	9000	10000

2. Machine Learning의 시작

1) Generative Classifier

$$P(\theta|X) \propto P(X|\theta)P(\theta)$$

Posterior Likelihood Prior

Posterior Maximization: MAP(Maximum a posterior)

Likelihood Maximization: MLE(Maximum likelihood estimation)

(Likelihood가 training set일 가능성이 많음(알고있을 가능성이 큼))

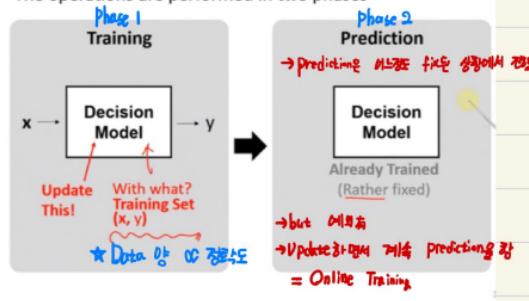
2) Machine Learning의 분류

Supervised Learning VS Unsupervised Learning

정답 있음	정답 없음
Classification	Regression
이산 data 분류	연속 data 분류

In Supervised Learning,

- The operations are performed in two phases



Machine Learning

data를 Y (알고 싶은 것), X (알고 있는 것)으로 나눈다

$Y = f(X) + n$ 관계라 할 때

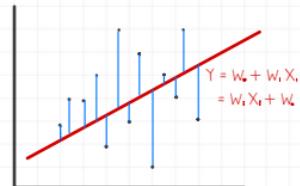
f(Decision Model) 를 추정하는 것이 machine learning이다.



Programming ($X \rightarrow Y$)

X 로부터 Y 를 도출

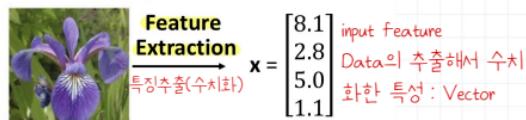
3. Linear Regression



1) Machine Learning에서 matrix를 사용하는 이유

Input feature: Data의 특징들을 수치화해서 카테고리화 시킴

→ This 4D data can be expressed with a vector



Feature를 여러 개 뽑아내면 무조건 행렬로 처리해야 함!!

Linear Regression

2) Weight & Design Matrix

Weight = 가중치, 비중

각각의 input feature에 대해서 영향력/ 가중치를 평가

$$\text{Ex) } \underline{X} = 4D = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \quad \underline{W} = 7 = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ w_4 \end{bmatrix}$$

$$\underline{W}^T \underline{X} = \begin{bmatrix} w_1 & w_2 & w_3 & w_4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4$$

Weight 종류가 여러개

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \\ w_{41} & w_{42} & w_{43} & w_{44} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

3) Design Matrix

$$\underline{X} = \begin{bmatrix} 1.1 \\ 2.1 \\ 3.1 \\ 4.1 \end{bmatrix} \quad \begin{bmatrix} 1.2 \\ 2.2 \\ 3.2 \\ 4.2 \end{bmatrix} \quad \begin{bmatrix} 1.3 \\ 2.3 \\ 3.3 \\ 4.3 \end{bmatrix} \quad \begin{bmatrix} 1.4 \\ 2.4 \\ 3.4 \\ 4.4 \end{bmatrix}$$

$$\text{Design Matrix} = \underline{X} = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \underline{x}_3^T \\ \underline{x}_4^T \end{bmatrix} = \begin{bmatrix} 1.1 & 1.2 & 1.3 & 1.4 \\ 2.1 & 2.2 & 2.3 & 2.4 \\ 3.1 & 3.2 & 3.3 & 3.4 \\ 4.1 & 4.2 & 4.3 & 4.4 \end{bmatrix}$$

Data의 차원수

$$\underline{X} = N \times D$$

Data 수

4) Model parameter

Machine learning을 할 때, 가장 중요한 것?

Model 설정!!!!

Weight = 중요도.

Weight는 input data에 곱하면??

Vector 내적!!!

$$\hat{y} = \underline{W}^T \underline{X} = W_0 + W_1x_1 + W_2x_2 + W_3x_3 + \dots + W_Dx_D \Rightarrow \begin{cases} \hat{y}: \text{Predicted value} \\ \text{Predicted value} = \text{내적 } \underline{w} \cdot \underline{x} \\ \text{내적} \end{cases}$$

$$\therefore \hat{y} = h_{\theta}(\underline{x}) = \underline{W}^T \underline{x}$$

$\Theta = \underline{W}$ = Model parameter vector

=> 우리가 결정해야 하는 model
(Called Hyperparameter)

5) Residual Sum of Squares

-> 정확성을 판단하는 방법이 될까?

-> Cost function

$$\rightarrow \text{RSS}(W) = \text{Cost func} = \sum_{i=1}^N (y_i - \underline{W}^T \underline{x}_i)^2$$

-> RSS = 0 이면 Best!!

6) Mean square Error

-> Square Error의 평균!! 즉, RSS(W)의 평균

$$\rightarrow \text{MSE}(W) = \frac{1}{N} \text{RSS}(W)$$

$$= \frac{1}{N} (y_i - \underline{W}^T \underline{x}_i)^2$$

Finding optimal weight

4. Linear regression - Analytical way

1) Ordinary Least Squares

$$\text{MSE}(\underline{w}) = \frac{1}{N} \sum_{i=1}^N (\underline{y}_i - \underline{w}^T \underline{x}_i)^2 \rightarrow \nabla \text{MSE} = -\frac{2}{N} \sum_{i=1}^N (\underline{y}_i - \underline{w}^T \underline{x}_i) \cdot \underline{x}_i = \underline{0}$$

$$\therefore \sum (\underline{w}^T \underline{x}_i) \cdot \underline{x}_i = \sum \underline{y}_i \underline{x}_i ; \text{Vector 4 } \underline{\underline{W}}$$

→ [같은 행렬로 어떻게 표현하지?]

Idea 1) Design matrix Idea 2) (스칼라) = 스칼라

$$\underline{X} = \begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \\ \vdots \\ \underline{x}_N \end{bmatrix} = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix} \quad \underline{W}^T \underline{x}_i = \text{벡터 내적} \\ \quad \quad \quad = \text{스칼라} \\ \therefore (\underline{W}^T \underline{x}_i)^T = (\underline{W}^T \underline{x}_i)^T = \underline{x}_i^T \underline{W} \\ \underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \quad (\underline{W}^T \underline{x}_i) \underline{x}_i = (\underline{x}_i^T \underline{W}) \underline{x}_i \\ \quad \quad \quad = (\underline{x}_i^T \underline{x}_i) \underline{W}$$

Idea 3) Concatenation

In python : vstack, hstack, dstack, column_stack

→ 이걸 Design matrix!!!!

$$\# \sum (\underline{W}^T \underline{x}_i) \underline{x}_i = \sum \underline{y}_i \underline{x}_i$$

$$\sum (\underline{x}_i^T \underline{W}) \underline{x}_i = \sum \underline{y}_i \underline{x}_i \quad \text{By inverse matrix}$$

$$\sum (\underline{x}_i^T \underline{W}) \underline{W} = \sum \underline{y}_i \underline{x}_i \quad * \text{ Ordinary Least Squares}$$

$$\downarrow \quad \quad \quad \underline{W} = \underline{W}^T = \underline{W}_{OLS}$$

$$\underline{W}_{OLS} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

Stack in python(Concatenation)

```
import numpy as np
import pandas as pd

a = [1, 2, 3, 4, 5]
b = [6, 7, 8, 9, 10]
c = [11, 12, 13, 14, 15]

Row1 = np.r_[a,b,c]           [ 1  2  3  4  5  6  7  8  9 10 11 12 13 14 15]

Row2 = np.hstack([a,b,c])     [ 1  2  3  4  5  6  7  8  9 10 11 12 13 14 15]

Column1 = np.c_[a,b,c]        [[ 1  6 11]
                           [ 2  7 12]
                           [ 3  8 13]
                           [ 4  9 14]
                           [ 5 10 15]]

Column2 = np.dstack([a,b,c])  [[[ 1  6 11]
                           [ 2  7 12]
                           [ 3  8 13]
                           [ 4  9 14]
                           [ 5 10 15]]]

print(Row1)
print(Row2)
print(Column1)
print(Column2)
```

Is There any Numerical way?

* Numerical Way↑ 왜 필요해?

→ 차원이 커지면 Computational Complexity up !!

→ 계산 복잡도 up !!

→ So, Closed-form Solution으로 풀기 불가능!!

5. Linear Regression - Numerical way

차원(N)이 커지면 Computational Complexity up!!!

즉 계산 복잡도 매우 큼!!

많은 경우 Analytical 한 방법으로 풀기 불가능

1) Numerical Way Solution

* Gradient Descent

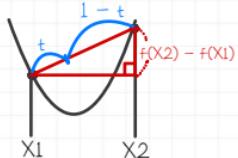
* Newton - Raphson method

* Least-Squares Solution

2) Gradient Descent

Convex 개념

* 볼록함수

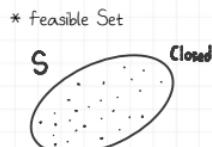


Convexity \Leftrightarrow 모든 $X_1, X_2 \in S$ 정의역, $t = [0,1]$
 $\therefore f(tX_1 + (1-t)X_2) \leq tf(X_1) + (1-t)f(X_2)$

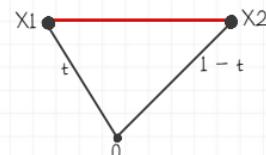
내분점

내분점의 합수값

* 즉, 구간 내에서 함수보다 위에 있다!!



S 안에 존재하는 모든 X_1, X_2 에 대해
 $tX_1 + (1-t)X_2 \in S$



Multi Variable function

: function은 scalar & Input = Vector, Output = Scalar

Scalar를 Vector로 미분
 $=$ Vector (Jacobian)

Vector를 Vector로 미분
 $=$ Matrix (Hessian)

$\underline{z} = \text{Vector}, H = \text{Hessian}$

$$\underline{z}^T H \underline{z} = \begin{matrix} \underline{z}^T \\ \underline{z} \end{matrix} \begin{matrix} H \\ \underline{z} \end{matrix} = \text{Scalar}$$

$D \times 1 \quad D \times D \quad 1 \times D$
 $= 1 \times 1$

Convex가 왜 중요해?

: feasible set에서 convex하면

Local min = Global min

Gradient: Vector로 미분

* ∇ : Gradient = Gradient의 방향 = 가파른 방향

* But, 내가 원하는 것은 미분이 0이 되는 Local point 지점

그럼 그 반대 방향으로 가면
 ∇MSE 가 최소가 되겠나?!!

Gradient Descent Algorithm

* MLE 관점: $P(W|D)$

→ Model parameter를 update

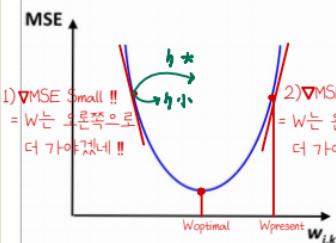
$$MSE(w) = \sum_{i=1}^N \frac{1}{N} (y_i - w^T x_i)^2 : (\text{진짜 Data의 } y\text{좌표} - \text{예측 Data의 } y\text{좌표})$$

$$\nabla MSE(w) = \sum_{i=1}^N -2(y_i - w^T x_i)x_i : \text{의 제곱의 평균} - RSS의 평균$$

Gradient Descent

$$W_{next} = W_{present} - \eta \nabla MSE(W)$$

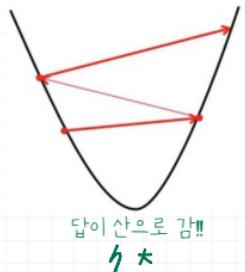
§ Gradient Descent는 초기값이 중요!! §

만약 $W_{present} \neq W_{optimal}$ 이면

W 를 optimal W 쪽으로 update해야겠네!!
 (∇MSE 의 반대방향 !!)

Learning Rate의 중요성 !!

Big learning rate



Small learning rate



Data Set으로 update 하는 방법 !!

$$\rightarrow MAP = P(D|W)$$

Step 1. Set up the initial value

Step 2. $X_{next} = X_{present} - (\text{Learning rate}) \times (\text{Gradient MSE})$

$$X_{next} = X_{present} - \eta \nabla MSE$$

Step 3. Iteration

Step 1) Initial value 설정하기

* Gradient Descent는 초기값이 매우 중요

$$\text{방향} = -\nabla = -\text{Gradient} = -\frac{\partial f}{\partial x}$$

* 얼마나?: Step size = Learning rate

Step 2)

* Scalar를 Vector로 미분

Objective Function : RSS = $\|y - Ax\|^2$ A: Weight matrix

$$f(a, b) = df = \frac{\partial f}{\partial a} da + \frac{\partial f}{\partial b} db$$

$$\rightarrow f_{ab} \quad \frac{\partial f}{\partial a} \frac{\partial f}{\partial b} = \left[\begin{array}{cc} \frac{\partial f}{\partial a} & \frac{\partial f}{\partial b} \end{array} \right] \left[\begin{array}{c} da \\ db \end{array} \right]$$

$$= \frac{\partial f}{\partial a} \cdot da$$

$\rightarrow df$ 를 구한 다음에 da 앞에 있는 그 Gradient

$$f = (y - Ax)^T (y - Ax)$$

$$df = \left[y - A(x+dx) \right]^T \left[y - A(x+dx) \right] - (y - Ax)^T (y - Ax)$$

$$= \left[(y - Ax) - A dx \right]^T \left[(y - Ax) - A dx \right] - (y - Ax)^T (y - Ax)$$

$$= (y - Ax)^T (y - Ax) - (A dx)^T (y - Ax) - (y - Ax)^T A dx + (A dx)^T A dx - (y - Ax)^T (y - Ax)$$

$$= -(A dx)^T (y - Ax) - (y - Ax)^T A dx$$

$$= -d x^T A^T (y - Ax) - (y - Ax)^T A dx$$

$$\boxed{dx} \quad \boxed{A^T} \quad \boxed{y - Ax} = \text{Scalar}$$

$$= -d x^T A^T (y - Ax) - (y - Ax)^T A dx$$

$$= -(y - Ax)^T A dx - (y - Ax)^T A dx$$

$$= -2(y - Ax)^T A dx$$

$$= -2(y - Ax)^T A dx$$

$$\rightarrow \frac{\partial f}{\partial a} \cdot da ; \quad \frac{\partial f}{\partial b} = -2(y - Ax)^T b$$

$$\therefore \frac{\partial f}{\partial a} = -2(y - Ax)^T b = \nabla MSE$$

* X 가 Gradient와 반대방향으로 가면 도달하네!!!

* $+2(y - Ax)^T A$ 방향!!

* 곱의 미분으로도 풀수 있다!!

$$\therefore X_{\text{next}} = X_{\text{present}} + h \cdot 2(y - Ax)^T b$$

$$(f_g)' = \text{앞미그두} + \text{앞그두미}$$

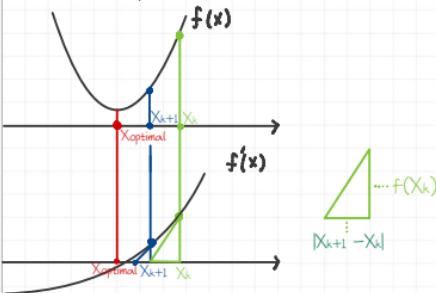
$$= f'_g + f_g'$$

$$\rightarrow f = (y - Ax)^T (y - Ax)$$

$$df = (-A dx)^T (y - Ax) + (y - Ax)^T (-A dx) = -2(y - Ax)^T A dx$$

$$\begin{matrix} \text{Scalar} \\ \text{Transpose} \end{matrix}$$

3) Newton - Raphson Method

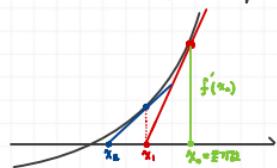


Newton - Raphson method는 Basically "Zero-Finding"

* Let Objective function = $f(\underline{x}) = (\underline{y} - \underline{A}\underline{x})^T(\underline{y} - \underline{A}\underline{x})$

* Gradient Descent의 결과 : Scalar를 Vector로 미분

$$\underline{x}_{k+1} = \underline{x}_k + 2(\underline{y} - \underline{A}\underline{x}_k)^T \underline{A}^{-1}$$



접선의 기울기 = 미분 \rightarrow 미분 역수 $\times f(x_0) = t$

$$|\underline{x}_{\text{next}} - \underline{x}_{\text{present}}| = |\Delta \underline{x}|$$

$$\frac{\Delta \underline{x}}{\Delta \underline{y}} = \frac{t}{f'(\underline{x}_{\text{present}})}, \Delta \underline{y} = f'(\underline{x}_{\text{present}})$$

$$\underline{x}_{\text{next}} = \underline{x}_{\text{present}} - \frac{f'(\underline{x}_{\text{present}})}{f''(\underline{x}_{\text{present}})}$$

Gradient Descent에서 Scalar를 Vector로 미분 \rightarrow Jacobian (Vector)

Newton-Raphson Method

* Scalar를 Vector로 두번 미분 = Vector를 Vector로 미분 = Hessian (Matrix)

* 두번 미분하고 나눈것은 역행렬을 곱한 것과 같다.

$$\frac{\partial}{\partial \underline{x}^T} \left(\frac{\partial f}{\partial \underline{x}^T} \right) = \frac{\partial}{\partial \underline{x}^T} [2(\underline{y} - \underline{A}\underline{x})\underline{A}] = \text{Scalar}$$

• $[\text{Scalar}]^T = \text{Scalar}$

$$\frac{\partial}{\partial \underline{x}^T} [2(\underline{y} - \underline{A}\underline{x})]^T = \frac{\partial}{\partial \underline{x}^T} [2\underline{A}^T(\underline{y} - \underline{A}\underline{x})] = 2\underline{A}^T \underline{A} = \text{Hessian Matrix}$$

$$\cdot [\text{Hessian}]^{-1} = \frac{1}{\underline{\underline{s}}^{nn}}$$

$$\therefore \underline{x}_{\text{next}} = \underline{x}_{\text{present}} + (2\underline{A}^T \underline{A})^{-1} 2\underline{A}^T (\underline{y} - \underline{A}\underline{x}_{\text{present}})$$

4) Least Squares Solution

만약에 Weight matrix $\underline{A}^T \underline{A}$ 가 역행렬이 있다면? invertable 하다면?

$$(2\underline{A}^T \underline{A})^{-1} 2\underline{A}^T (\underline{y} - \underline{A}\underline{x})$$

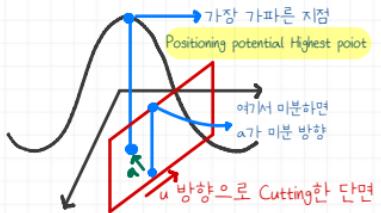
~~$$= \frac{1}{2} \cancel{\underline{A}^T \underline{A}} \cancel{\underline{A}^T} (\underline{y} - \underline{A}\underline{x})$$~~

$$\rightarrow \underline{x}_{\text{next}} = \underline{x}_{\text{present}} + \underline{A}^T (\underline{y} - \underline{A}\underline{x}_{\text{present}})$$

$$= \underline{x}_{\text{present}} + \underline{B}^T \underline{y} - \cancel{\underline{x}_{\text{present}}}$$

$$\therefore \underline{x}_{\text{next}} = \underline{B}^T \underline{y}$$

6. 방향도함수와 Gradient



Directional Derivative

$$\lim_{h \rightarrow 0} \frac{f(\underline{x} + h\vec{u}) - f(\underline{x})}{h} = \nabla f \cdot \vec{u}$$

\vec{u} 방향으로의 방향 미분 $|u| = 1$

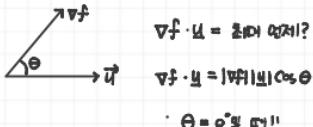
$$* f(\underline{x}) = f\left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}\right) = f\left(\begin{pmatrix} x_1 + hu_1 \\ x_2 + hu_2 \end{pmatrix}\right)$$

→ 이라면 x_1, x_2 는 h 의 함수가됨 (u_1, u_2 는 내가 정한 방향 : 상수)

$$\frac{df}{dh} = \frac{df}{dx_1} \cdot \frac{du_1}{dh} + \frac{df}{dx_2} \cdot \frac{du_2}{dh} = \nabla f \cdot \vec{u}$$

따라서 Gradient = \vec{u} 방향으로의 미분

* 그럼 \vec{u} 를 바꿔 가면서 어느 방향으로 제일 크나? $|u| = 1$



즉, Gradient 벡터는 미분할 때 가장 가깝다 !!

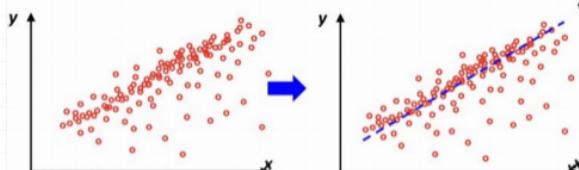
7. Linear Regression 보충

Phase 1. Training에서 MLE를 써보자 !!

$$W_{\text{next}} = W_{\text{present}} - \eta \nabla \text{MSE} : P(W|D)$$

→ MAP임 !!

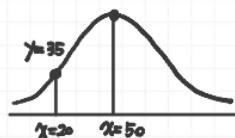
Phase 2. MLE??



$$\text{MLE} = P(W|D)$$

* 파란색 선을 찾을 때 MLE를 쓴다 !!

* Linear Regression이 좀 더 meaningful 하려면 Probability 써야 함 !!



∴ "X = 20에서 Y = 35" 이거보단

↳ "X = 20일 때 평균에서 얼마나큼 떨어졌고 확률이 얼마야 ?!"

∴ 이 훨씬 meaningful !!

Linear Regression || Gaussian Model

가장 표준편차는 fix이고 x에 따라 평균만 change y

$$P(y | x, \theta) = N(y | \underline{w}_x^T \theta)$$

* 계산 평균: $10.9 \pm \frac{1}{2} \text{요}(ln)$

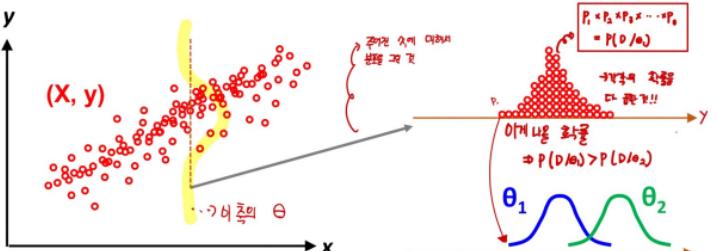
모든 Argument $\ln P(D|Q)$

$$l(\theta) = \log P(D|\theta) = \sum_{i=1}^N \log P(x_i|\theta)$$

Meaning of $p(D|\theta)$?

f p(D|θ)? $\underbrace{p(\theta|D)}_{\text{이거 구하기 힘들어}} \propto \underbrace{p(D|\theta)p(\theta)}_{\substack{\text{이걸} \\ \text{구해} \\ \text{야趾}}}$ $\xrightarrow{\theta \text{가} \text{다른} \text{게}} \text{Uniform}$

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y|\mathbf{w}^T \mathbf{x}, \sigma^2)$$



→ Can you compare $p(D|\theta_1)$ and $p(D|\theta_2)$?

Note that probability is probability \propto PDF

→ Can you relate this with Number Game case? Occam's razor?

→ But, PDP+ 글수를 확장해 더 높아질 예상 가능

MLE 계산 유도

* Analytic한 방법으로 구한 (MAP) wOLS 와 결과 같음

But, why MLE?

$\rightarrow P(D|W)$ 는 구하기 어려우.

-> Model은 눈으로만 보면 알기 어렵다 ; 계산 필요

$$P(W|D) \propto P(D|W)P(W)$$

$$\ell(\theta) = \sum_{i=1}^N \log \left[\left(\frac{1}{2\pi\sigma^2} \right)^{\frac{1}{2}} \exp \left(-\frac{1}{2\sigma^2} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \right) \right]$$

$$= \frac{-1}{2\sigma^2} RSS(\mathbf{w}) - \frac{N}{2} \log(2\pi\sigma^2)$$

1. (a)의 최대 = RSS(b)의 최소

$\Rightarrow (\gamma - \beta_{\text{fit}})^2$ 은 최소화 = MLE를 최대화

Minimizing MSE approach for Comparing it w/ MLE or MAP
Minimizing MSE it to minimize the following

$$\text{MSE}(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

* MAP is to maximize $p(W|D)$

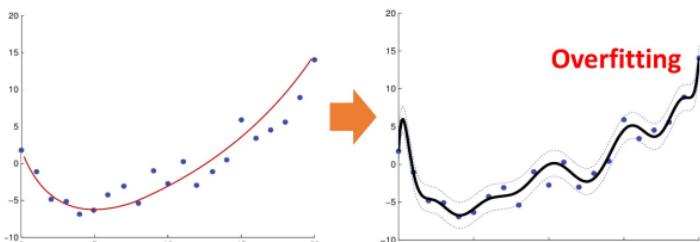
* MI E is to maximize $p(D|w)$

For now, we can roughly say that decreasing error is equivalent to increase probability
(MLE 와 MSE는 비슷한 단위로 결과는 같다!)

8. Overfitting

: Dataset에 매몰되어 Model을 정함

-> Solution 1) Prior 2) Abundant Data 3) Regularization



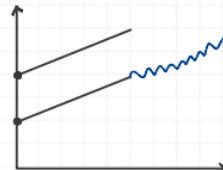
$$P(W|D) = P(D|W)P(W)$$

* 만약 MLE를 위해 $P(D|W)$ 만 maximize하면 overfitting 생길 수 있음

* data 수가 적으면 특히 많이 발생

* Prior 를 이용하거나 더 많은 data를 fitting!!

Performing MAP



$$\begin{aligned} W_0 &\leftarrow (W_1 x_1 + W_2 x_2 + W_3 x_3 + W_4 x_4 + \dots); \\ W_0 &\quad W^T x \end{aligned}$$

W_0 는 Zero mean 이면 안됨

-> W_0 는 $\text{data } \mathbf{x}_i$ 에 대해 중심점을 잡아주는 것 같아 !!

$$p(D|W) = \mathcal{N}(y|w_0 + w^T x, \sigma^2) \text{ and } p(W) = \prod_j \mathcal{N}(w_j|0, \tau^2)$$

- This leads to maximize log likelihood:

$$\underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i=1}^N \log \mathcal{N}(y_i | w_0 + \mathbf{w}^T \mathbf{x}_i, \sigma^2) + \sum_j^D \log \mathcal{N}(w_j | 0, \tau^2)$$

▶ 다만 최대화하면 Overfitting !!

* 목표는 전체 최대화 !!

-> 기준이랑 Different, why?

-> 기준에는 $P(W) = \text{Prior } \mathcal{N}(w|0, I)$

- That is equivalent to minimize

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - (w_0 + \mathbf{w}^T \mathbf{x}_i))^2 + \lambda \|\mathbf{w}\|_2^2$$

where

$$\lambda \triangleq \sigma^2 / \tau^2 \text{ and } \|\mathbf{w}\|_2^2 = \sum_j w_j^2 = \mathbf{w}^T \mathbf{w}$$

- That is equivalent to minimize

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - (\underline{w_0} + \mathbf{w}^T \mathbf{x}_i))^2 + \lambda \|\mathbf{w}\|_2^2$$

where

$$\lambda \triangleq \sigma^2 / \tau^2 \text{ and } \|\mathbf{w}\|_2^2 = \sum_j w_j^2 = \mathbf{w}^T \mathbf{w}$$

- * w_0 는 offset $\rightarrow w_0$ 는 regularized 안됨; 함수의 높이에만 영향을 줌
- * 다시 말해 bias는 regularized X \rightarrow complexity에 영향을 줌
- * lambda = Regularizer = λ

Optimal w with Regularization

- The resultant optimal \mathbf{w} is

$$\hat{\mathbf{w}}_{\text{ridge}} = (\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

* Ridge Regression = Penalized Least Squares

* Prior $\stackrel{\text{은}}{\sim}$ uniform distribution 이 아닌! Gaussian prior term으로 더해주면 이를 Regularization 또는 Weight decay라고 한다.

* Regularization의 가장 큰 목적은 OVERFITTING을 줄이는 것!!!!!!

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import math as m

plt.close("all")

dfLoad = pd.read_csv("https://raw.githubusercontent.com/hanwoolJeong/lecture1/master/linearregression.csv", sep="\s+")
xxRaw = dfLoad["xx"]
yyRaw = dfLoad["yy"]
yyRawNP = np.array(yyRaw)
plt.plot(xxRaw, yyRaw, ".r")

# # 이 상황에서 linear regression 모델 찾아보겠다이 말이야
Ndata = len(xxRaw)
X = np.c_[np.ones([100, 1]), xxRaw]
# wOLS = X.T.dot(X) 행렬곱까지

# wOLS = np.linalg.inv(X.T.dot(X)) 역행
wOLS = np.linalg.inv(X.T.dot(X)).dot(X.T).dot(yyRawNP.reshape(Ndata, 1))
# wOLS 구해짐
xPredict = np.linspace(0, 2, num=101)
xPredictPadding = np.c_[np.ones([101, 1]), xPredict]
yPredict = wOLS.T.dot(xPredictPadding.T)

# plt.plot(xPredict.reshape(1, 101), yPredict, "b.")
# 직선 만들기
eta = 0.1
Niteration = 10
# 20번 돌렸을 때는 차이가 멀어 안나서 200번으로 바꿈(Niteration), print page 27, 28

wGD = np.array([0, 0]).reshape(2, 1)
for iteration in np.arange(Niteration):
    gradients = -(2/Ndata)*(X.T.dot(yyRawNP.reshape(Ndata, 1))-X.dot(wGD)))
    wGD = wGD - gradients*eta
    print(iteration)
    yGD = wGD.T.dot(xPredictPadding.T)
    plt.plot(xPredict.reshape(1, 101), yGD, "b.")

```