# Linear Regression

$$RSS = e_1^2 + e_2^2 + \cdots + e_n^2$$
$$= (y_1 - \hat{y}_1)^2 + \cdots + (y_n - \hat{y}_n)^2$$

Single Linear Regression : $y = \hat{\beta}_0 + \hat{\beta}_1 x$ model 에서

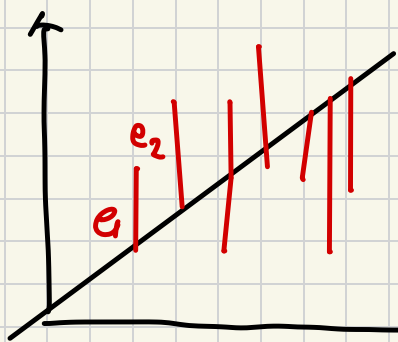최적의 $\hat{\beta}_0$ 과 $\hat{\beta}_1$ 을 구하기 위해서 RSS를 minimize 한다.

$$RSS = \sum_{i=1}^{n} \left( y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right)^2$$

이제 RSS를 minimize 하는 $\hat{\beta}_0$ 과 $\hat{\beta}_1$ 을 각각으로 RSS를 미분한다

$$\frac{\partial RSS}{\partial \hat{\beta}_0} = 0 \quad , \quad \frac{\partial RSS}{\partial \hat{\beta}_1} = 0 \quad \text{이 되도록 해보자. (이를 만족하면}$$
$$\text{max 아니면 min 이다 )}$$

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

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

\# Single Linear Regression에서 최적의 $\hat{\beta}_0$ 과 $\hat{\beta}_1$ 를 얻기 위해서
RSS를 각각 $\hat{\beta}_0$ 과 $\hat{\beta}_1$ 로 미분해서 구해준다.



〈평가방식〉

RSS $\Rightarrow e^2$의 sum \# data ↑ error ↑

RSE, RMSE $\Rightarrow$ data수로 나눈다.

$\sqrt{\dfrac{RSS}{n-2}} \quad \sqrt{\dfrac{RSS}{n}}$

R-square : $1 - \dfrac{RSS}{TSS}$ (이면 1에가까) 

→ 두의 model의 그차의합
→ 실제 분산의 합
$\sum_{i=1}^{n} (y_i - \bar{y})^2$
→ model이 이저 설명까지 몇만 변동.

# Maximum likelihood

Why 곱이 ( )^x 을 사용해서? 그냥 abs()는 사용해도 되는거 아닌가?

IID : Independent and Identically distribution
⇒ 모든 사건이 독립이고, 동일한 확률분포를 따른다.

우리 data들이 어느 분포에서 얻어졌는지 모으기에 명확한 반대 증거가 없는한
IID 를 따른다고 가정한다.

Likelihood : 이미 관측된 data가 있을 때, 그 data를 가장 잘 설명해 주는 분포 parameter θ를
찾기 위한 값음.

일반확률은 θ가 주어진 경우 data가 나올 확률 $p(data | \theta)$ 이면
가능도 함수는 θ를 확률 변수로 두고 $p(\theta | data)$ 를 의미한다.

그래서 $L(\theta ; \underline{x_1, x_2 \cdots x_n})$ 이 가능도이다.


n개의 data가 관측된 경우

θ값

⇓
IID가정여기에

$L(\theta ; x_1, x_2, \cdots, x_n) = p_\theta(x_1) \cdot p_\theta(x_2) \cdots \cdot p_\theta(x_n)$ 이 된다.

결국 모든 확률값의 곱을 maximize하는 θ를 구하는게 가능도 함수의 목적이다.

\* Distribution의 모형은 이미 가정 (예를 들어 Gaussian, 나들바시즌 ---)

그래서 Gaussian인 경우 $\mu$ 와 $\gamma$ 를 찾게 되는것이다.

# Log-likelihood

: 여러의 maximize, minimize를 구할거기에



를 곱해서 해당 value는 변화X.
monotonus increasing 하기에.
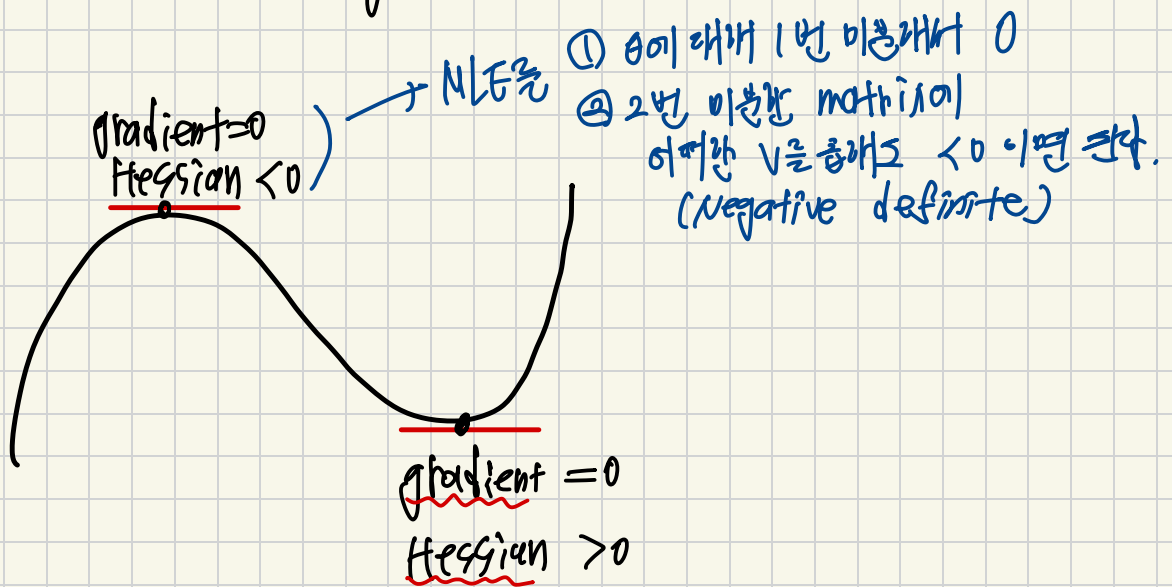
$$\log \prod_{i=1}^{n} p_\theta(x_i) = \sum_{i=1}^{n} \log p_\theta(x_i)$$

$$\arg\max_x f(x) = \arg\max_x \log f(x)$$

각각의 gradient 구하기에 용이

# MLE : Maximum Liklihood Estimation

$$\hat{\theta}(x_1, x_2 \cdots . x_n) = \underset{\theta}{\arg\max} L(\theta; x_1, x_2 \cdots x_n)$$

$$= \underset{\theta}{\arg\max} l(\theta; x_1, x_2 \cdots . x_n)$$

$$= \underset{\theta}{\arg\max} \sum_{i=1}^{n} \log P_\theta(x_i)$$

gradient=0
Hessian <0 $\rightarrow$ MLE은

① $\theta$에 대해 1번 미분해서 0
② 2번 미분한 matrix이
어떠한 V를 곱해도 <0 이면 된다.
(Negative definite)

gradient = 0
Hessian >0

## 〈MLE 예시 : Bernoulli distribution〉

$$\underset{\theta}{\arg\max} \sum_{i=1}^{n} \log P_\theta(x_i) = \underset{\theta}{\arg\max} \sum_{i=1}^{n} \log\left(\theta^{x_i}(1-\theta)^{(1-x_i)}\right)$$

$$= \underset{\theta}{\arg\max} \log\theta \sum_{i=1}^{n} x_i + \log(1-\theta)\sum_{i=1}^{n}(1-x_i)$$

$$\frac{\partial\left(\log\theta\sum_{i=1}^{n}x_i + \log(1-\theta)\sum_{i=1}^{n}(1-x_i)\right)}{\partial\theta} = 0 \; 을 \; 풀면$$

$$\theta = \frac{1}{n}\sum_{i=1}^{n}x_i \; 가 \; 된다. \quad 자연스럽게 \; 평균값이 \; 나오게 된다.$$

2번 미분해서 Hessian을 구하면

$$\frac{\partial^2 l(\theta)}{\partial \theta^2} = -\left(\frac{1}{\theta^2}\right) \sum_{i=1}^{n} x_i - \frac{1}{(1-\theta)^2} \sum_{i=1}^{n}(1-x_i) \quad \therefore 0보다 작다.$$

## ⟨ MLE 예시 : Gaussian distribution ⟩

$\theta = \{ \mu, \gamma \}$ 이다. // $\mu$ 따 $\gamma$에 대해 모두 미분했다.

$$\hat{\theta}(x_1, x_2, \cdots x_n) = \underset{\theta}{argmax} \sum_{i=1}^{n} \log \frac{1}{\gamma\sqrt{2\pi}} \cdot e^{-\frac{1}{2}\left(\frac{x_i - \mu}{\gamma}\right)^2}$$

$$= \underset{\theta}{argmax} \sum_{i=1}^{n} \log \frac{1}{\gamma\sqrt{2\pi}} + \sum_{i=1}^{n} -\frac{1}{2}\left(\frac{x_i - \mu}{\gamma}\right)^2$$

$\therefore \mu$에 대해 1차 미분시 $\mu = \frac{1}{n}\sum_{i=1}^{n} x_i$ 가 된다.

$\gamma$의 "  $\gamma = \frac{1}{n}\sum_{i=1}^{n}(x_i - \mu)^2$ 가 된다.

※ 정규분포는 가시안 가정시 MLE를 통해서 나온 결과이다.

# MLE for Linear Regression

기본적인 Linear Regression은 $y = \beta_0 + \beta_1 x_1 + \cdots + \beta_n x_n$ 표현가능

이를 $y = \beta^T x + \underline{\varepsilon}$ 로 표현이 가능하다.

<span style="color:red">└→ 애가 $\sim N(0, \gamma^2)$을 따른다고 가정하자.</span>

$\Rightarrow$ 결국 Linear Regress은 $p(y|x) \sim N(\underline{\beta^T x}, \gamma^2)$

<span style="color:red">└ 우리가 구한 선형을 따르되<br>특정 $\gamma$만큼 분산이 존재한다.</span>

이를 MLE를 구하기 위해 가능도를 구하게 되면

$$\arg\max_\beta \sum_{i=1}^{n} \log \frac{1}{\gamma\sqrt{2\pi}} \cdot e^{-\frac{1}{2}\left(\frac{y - (\beta^T x)}{\gamma}\right)^2}$$

$$\simeq \arg\max_\beta \textcircled{\small$-$} \sum_{i=1}^{n} (y_i - \beta^T x)^2$$

$$\arg\min_\beta \sum_{i=1}^{n} (y_i - \beta^T x_i)^2 \quad \left(\underline{\text{오차제곱의 합이다}}\right)$$

<span style="color:red">$e$를 제곱하는 이유이다.</span>

행렬에 대해서 $\frac{\partial RSS(\beta)}{\partial \beta}$를 구하게 되면 $\beta = \underline{(X^T X)^{-1} X^T y}$ 가 된다.

<span style="color:blue">이게 정사형 Matrix이다.</span>

<span style="color:blue">∴ $y$가장 잘 대변하는 $\beta$</span>



<span style="color:blue">X가 span</span>    <span style="color:red">이게 $\beta$이다.</span>

# Bias, Variance and MSE

$$\text{Bias}(\hat{\theta}) = E[\hat{\theta}] - \theta$$

$$\text{Var}(\hat{\theta}) = E[(\hat{\theta} - E[\hat{\theta}])^2]$$

$$\text{MSE} = E\left(\sum_{i=1}^{d}(\hat{\theta}_i - \theta_i)^2\right)$$

$$\text{MSE} = \text{tr}(\text{Var}) + \|\text{Bias}\|^2$$

$$\mathbb{E}\left[\sum_{j=1}^{d}(\hat{\theta}_j - \theta_j)^2\right] \overset{?}{=} \sum_{j=1}^{d}\mathbb{E}\left[(\hat{\theta}_j - \mathbb{E}(\hat{\theta}_j))^2\right] + \sum_{j=1}^{d}(\mathbb{E}(\hat{\theta}_j) - \theta_j)^2$$

For simplicity, let's suppose d = 1 without loss of generality.

$$\mathbb{E}\left[(\hat{\theta} - \theta)^2\right] \overset{?}{=} \mathbb{E}\left[(\hat{\theta} - \mathbb{E}(\hat{\theta}))^2\right] + (\mathbb{E}(\hat{\theta}) - \theta)^2$$

- Using the Bias-Variance decomposition, we can compute the MSE of linear regression:

$$\text{MSE}(\hat{\theta}|\mathbf{X}) = \boxed{\text{Var}(\hat{\theta}|\mathbf{X})} + \boxed{\text{Bias}(\hat{\theta}|\mathbf{X})^2} = \sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}$$

$$\sigma^2(\mathbf{X}^\top\mathbf{X})^{-1} \qquad 0$$

  - Linear regression is unbiased (bias = 0).
  - With infinite data, variance also converges to 0.

- It can be also proved that the MLE is the best unbiased estimator.
  - That is, no other $\theta$ has lower variance than the one found by MLE. (Gauss-Markov Theorem)