

# <보통 data 복제>

$\lambda = X_1 + \dots + X_p$ ,  $p=100$  (차원)  
 $Y \sim \text{Pois}(\lambda)$   
 $\lambda = \exp(\hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p)$

① Data  $\rightarrow \hat{\beta}$   
 ②  $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\lambda}_i)^2$

$\hat{\lambda}_i = \exp(\hat{\beta}_0^* + \dots + \hat{\beta}_p^* X_{ip})$   
 $\rightarrow$  가장 data 얼마나 잘 설명하는지 평가하기 X.

good  $\rightarrow$  비교  
 $\rightarrow$  무나 더 좋은지 알수 없음

bad  $\rightarrow$  가장 data 성능 얼마나 나쁜가?

비교 : original MSE better 항상  $\hat{\beta}$ 로 리미트됨

$\hat{\beta}$  찾을 때 :  $\hat{\beta}^* = \arg\min_{\hat{\beta}} \rightarrow$  original data X, Y 관계 잘 설명해야 되는데 가장 data가 bad.

(MSE)  $\frac{1}{n} \sum_{i=1}^n (y_i - \exp(\hat{\beta}_0^* + \hat{\beta}_1^* X_{i1} + \dots + \hat{\beta}_p^* X_{ip}))^2$   
 $\star \rightarrow$  더 나쁜 경우 없다.  $\hat{\beta} = \arg\min_{\hat{\beta}} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\lambda}_i)^2$

• Test용 set에서 X, Y 발생, 복제해보기

Test data  
 $y_{n+1}, x_{n+1}$   
 $y_{n+T}$

$\frac{1}{T} \sum_{i=n+1}^{n+T} (y_i - \exp(\hat{\beta}_0^* + \hat{\beta}_1^* X_{i1} + \dots + \hat{\beta}_p^* X_{ip}))^2$   
 ||  
 MSE replicated

MSE Original =  $\frac{1}{T} \sum_{i=n+1}^{n+T} (y_i - \exp(\hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \dots + \hat{\beta}_p X_{ip}))^2$

(original data)  $\rightarrow \hat{\beta}$   
 복제된 data  $\rightarrow \hat{\beta}^*$

$p=100$

$\vec{X}_i \sim \text{MVN}(\vec{0}, \Sigma_p)$ ,  $[\beta_0, \dots, \beta_p \sim \text{unif}(-1, 1)] \rightarrow$  고정  $\sigma^2=1$  고정

$y_i \sim N(u_i, \sigma^2)$ ,  $u_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}$

$= \beta_0 + \beta_1^* X_i$  for  $i=1, \dots, n$

VAE  
 $\begin{bmatrix} y_1 & \vec{X}_1 \\ \vdots & \vdots \\ y_n & \vec{X}_n \end{bmatrix} \rightarrow \begin{bmatrix} y_1^* & \vec{X}_1^* \\ \vdots & \vdots \\ y_n^* & \vec{X}_n^* \end{bmatrix}$

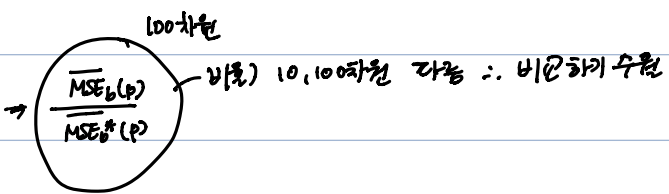
Test MSE (Test data 사용)  
 $MSE = \frac{\sum_{i=n+1}^{n+T} (y_i - (\hat{\beta}_0 + \hat{\beta}_1^* X_i))^2}{T}$   
 $MSE^* = \frac{\sum_{i=n+1}^{n+T} (y_i - (\hat{\beta}_0^* + \hat{\beta}_1^* X_i))^2}{T}$

$\frac{\sum_{b=1}^B MSE_b}{B} < \frac{\sum_{b=1}^B MSE_b^*}{B}$   
 더 좋음 (original 강행)

Test  
 $\begin{bmatrix} y_{n+1} & \vec{X}_{n+1} \\ y_{n+T} & \vec{X}_{n+T} \end{bmatrix} \Rightarrow$  Iteration  $b=1, \dots, B \therefore \beta$  항상 바뀜

\* 차원 ↑ 예측력 ↓ 차원(p) ↓ 예측력 ↑

ex. iteration 10번 (p=10) 계산  $\overline{MSE_b}$  (평균),  $\overline{MSE_s}$



p = 10 20 ... 100 시도 → 어느 순간부터 예측력 나빠지는지

가장

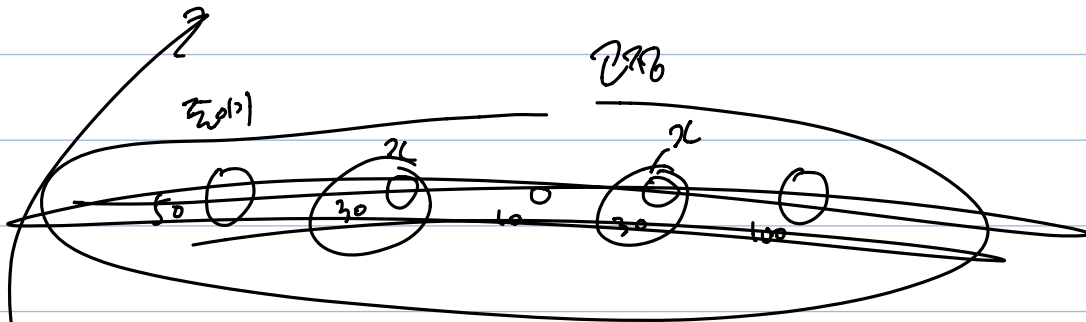
문제: GAN, AE, VAE → 복제 잘 안되서. y 예측에서 data와 MSE 안맞음.

↳ 보통은 data 10 ~ 100 개 랜덤하게 뽑기 (h/w)

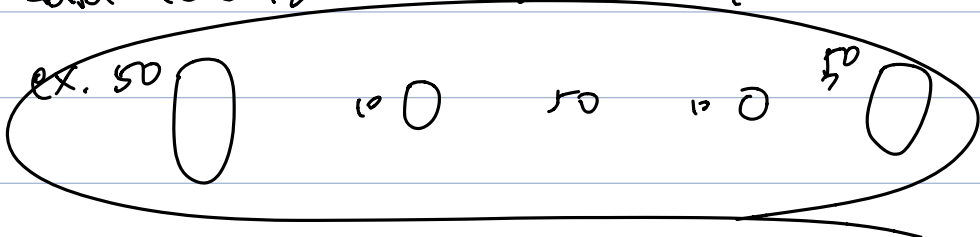
B=10 X 차원 10, 20 ... 100 만지

1개 data = 10000개 정도

<AE>



data 10만개도 하기! 100 20 10 20 100 정도는 해보기



param 개수 신경쓰지 않기

10:1 비율로 작아지게 하기

다들) 발표...?

(여유) discrete GAN model 어떻게 구현 가능하  
(커먼의 논문 참고)  
+ 부정확함 알기, 분석

1. 2. 3. 4. 5. 6. 7.