

① Understand conventional regression model (personal perspective)

Result: the basic of Mean Regression is exactly our model setting!

Firstly, we have

{ dependent var: Y
 independent var: X

→ when X is fixed, consider $\underline{Y|X=x}$ r.v.,

it can be decomposed into $\underline{Y|X=x = E[Y|X=x] + \varepsilon_x}$

generally speaking, the decomposition guarantees that $\underline{E[\varepsilon_x] = 0}$

Issue:

- { ① ε_x 's variance might be different
- ② ε_x 's distribution might be different

Necessary condition

$\hat{f}(x) \sim \dots$

equivalent formulation

$$Y|X=x = f(x) + \varepsilon_x \text{ where } E[\varepsilon_x] = 0$$

Model Setting for Linear Reg

then, **a necessary condition is that** $f(x) = E[Y|X=x]$

further stronger assumption (prior)

- ① $\text{Var}[Y|X=x] = \sigma^2$ (independent w.r.t x !)
- ② $Y|X=x_i$ is independent w.r.t different x_i
- ③ $Y|X=x \sim \text{Gaussian}$

Since Under this Model Setting, $f(x) = \underline{E[Y|X=x]}$

Our aim is to construct an Estimator $\hat{f}(x)$ to approximate $f(x)$

the given method is to use MLE under our assumed prob model

⇒ finally, what we have done is **Regression towards Mean**

Linear Regression

① Model :

$$y = X\beta + \varepsilon \quad \left\{ \begin{array}{l} 1. \varepsilon \sim N(0, \sigma^2 I_n) \\ 2. \text{rank}(X) = p+1 \end{array} \right.$$

$$\textcircled{2} \quad \text{MLE} \Rightarrow \left\{ \begin{array}{l} \hat{\beta} = (X^T X)^{-1} X^T y \Rightarrow \hat{y} = X \hat{\beta} = X(X^T X)^{-1} X^T y \\ \hat{\sigma}^2 = \frac{1}{n} y^T (I - H) y \end{array} \right. := H \cdot y$$

③ Property of estimator :

$$\left\{ \begin{array}{l} E[\hat{\beta}] = (X^T X)^{-1} X^T \cdot E[y] = \beta \\ \text{var}[\hat{\beta}] = \text{var}[(X^T X)^{-1} X^T y] = (X^T X)^{-1} X^T \text{var}[y] X (X^T X)^{-1} \\ = \sigma^2 (X^T X)^{-1} \\ \hat{e} = y - \hat{y} = (I - H)y, \quad \text{cov}(\hat{\beta}, \hat{e}) \\ = \text{cov}((X^T X)^{-1} X^T y, (I - H)y) \end{array} \right.$$

$$\text{var}(\hat{e}) = \sigma^2 (I - H) = (X^T X)^{-1} X^T \text{var}(y) (I - H)$$

$$E[\text{SSE}] = \text{tr}(\text{var}[\hat{e}]) = (n-p-1)\sigma^2 = \sigma^2 \cdot 0 = 0$$

④ Best Linear Unbiased Estimator $\Rightarrow \hat{y} \xrightarrow{\text{BLUE}} X^T \hat{\beta}$

$$\rightarrow \text{for } \forall C, \quad C^T \hat{\beta} \xrightarrow{\text{BLUE}} C^T \beta$$

Pf Sketch : Suppose $\alpha^T y$ is an estimator for $C^T \beta$

$$\text{then } E[\alpha^T y] = \alpha^T X \beta = C^T \beta \quad \forall \beta$$

$$\Rightarrow \alpha^T X = C^T \Leftrightarrow X^T \alpha = C$$

$$\begin{aligned}
 & \text{Consider } \text{Var}[\alpha^T \gamma] - \text{Var}[\beta^T \hat{\beta}] \\
 & = 6^2 \alpha^T \alpha - 6^2 \beta^T (X^T X)^{-1} \beta \\
 & = 6^2 \cdot \alpha^T \underbrace{\left(I - X(X^T X)^{-1} X^T \right)}_{\text{幂等 } (I-H)} \alpha \geq 0
 \end{aligned}$$

⑤ Hypothesis Testing

$$H_0: \beta_1 = \dots = \beta_p = 0$$

a) $\frac{\frac{SSR}{p}}{\frac{SSE}{n-p-1}} \sim F(p, n-p-1)$

$H X = X \Rightarrow H \cdot 1 = 1$

$\text{tr}(\cdot) = p+1-1=p$ $\beta_0 \cdot 1$

$\left\{ \begin{array}{l} SSR = y^T \left[H - \frac{1}{n} 1 \cdot 1^T \right] y \quad (H - \frac{1}{n} 1 1^T) X \begin{pmatrix} \beta_0 \\ 0 \end{pmatrix} = 0 \\ SSE = y^T [I - H] y \quad (I - H) \cdot X \beta = 0 \quad \text{tr}(\cdot) = n-p-1 \end{array} \right.$

SSE $\perp\!\!\!\perp$ SSR

b) $\hat{\beta} \sim \text{Gaussian}(\beta, 6^2 (X^T X)^{-1})$

$\text{Cov}(\hat{\beta}, \hat{\epsilon}) = 0 \Rightarrow \hat{\beta} \perp\!\!\!\perp \hat{\epsilon} \Rightarrow \hat{\beta} \perp\!\!\!\perp SSE$

$\Rightarrow \frac{\frac{\hat{\beta}_j}{6 \sqrt{e_{jj}}}}{\frac{SSE}{n-p-1}} = \frac{\hat{\beta}_j}{\hat{6}_{\text{修正}} \cdot \sqrt{e_{jj}}} \sim t(n-p-1)$

$H_0: \beta_j = 0$

c) $R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$



Quantile Regression → A different perspective

① First Problem-Setting

$$Y|X=x = g_z(x) + \varepsilon_z | X=x$$



$$\underbrace{Q_z(y|X=x)}_{\Downarrow} = g_z(x)$$

$$\text{where } \underbrace{Q_z(\varepsilon_z|X=x)}_{\Downarrow} = 0$$

We want to seek for an estimator for $\underbrace{g_z(x)}_{\text{function}}$

$\boxed{\text{ReLU function}}$

② Ground-Truth (Amazing observation)



$$P_z(\cdot) = \begin{cases} 0, & x < 0 \\ zx, & x \geq 0 \end{cases}$$

$$Q_{Y|X=x}(z) = \inf \{y : F_{Y|X=x}(y) \geq z\}$$

$$\text{then } Q_{Y|X=x}(z) = \underset{\theta}{\operatorname{argmin}} \mathbb{E}[P_z(Y|X=x - \theta)]$$



$$Q_{Y|X}(z) = \underset{g \in G}{\operatorname{argmin}} \mathbb{E}[P_z(Y - g(X))] \quad (X, Y) \sim P$$

$$= \underset{g \in G}{\operatorname{argmin}} T(P, g)$$

③ Plug-in estimator

$$\hat{Q}_{Y|X}(z) = \underset{g \in G}{\operatorname{argmin}} T(\hat{P}_n, g)$$

(X, Y) random

$$= \underset{g \in \mathcal{G}}{\operatorname{argmin}} \sum_{i=1}^n p_i (Y_i - g(X_i))^2$$

for fixed X , only random Y

similar derivation also holds!

Generally, $\mathcal{G} = \{g : g: x \rightarrow \beta^T x, \forall \beta\}$

Why Mean Regression?

(Another interpretation)

Observation: Optimality of conditional expectation w.r.t MSE

$$\mathbb{E}[Y|X=x] := g(x) = \underset{g \in \mathcal{G}}{\operatorname{argmin}} \mathbb{E}[(Y-g(x))^2]$$

should be similar!

plug-in estimator

More related to

Non-parametric Framework

(like quantile regression)

$$\hat{g} = \underset{g \in \mathcal{G}}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - g(X_i))^2$$

choose $\mathcal{G} = \{g : g(x) = \beta^T x, \beta \in \mathbb{R}^n\}$,

then this is exactly Linear Regression framework

Update at 11/18

(view)

→ plug-in estimator is just one understanding of Linear Regression Algo

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (\gamma_i - x_i^\top \beta)^2 = \underline{(x^\top x)^{-1} x^\top y}$$

$$\tilde{g}(x) = \mathbb{E}[y|x=x] \leftarrow \tilde{g} = \underset{g \in G}{\operatorname{argmin}} \mathbb{E}_{(x,y)}[(y - g(x))^2]$$

→ However, it is more reasonable to use Parametric Estimation Framework

(under Gaussian assumption, infer Gaussian $(x^\top \beta, \sigma^2)$ ~ $\hat{\beta}$)

$\hat{y} \rightarrow X\hat{\beta}$ estimator (when we achieve $y \sim \text{Gaussian}(x^\top \hat{\beta}, \sigma^2)$, estimate y !) #

Previously, mainly talk about: (LR)

① derivation of $\hat{\beta} = (x^\top x)^{-1} x^\top y$

② Basic property of $\hat{\beta}$ (estimator of β)

1) unbiased

$$2) \operatorname{Var}[\hat{\beta}] = (x^\top x)^{-1} x^\top \operatorname{Var}(y) x (x^\top x)^{-1}$$

$$= \sigma^2 (x^\top x)^{-1} \star$$

$$\Rightarrow \hat{\beta} \sim N(\beta, \sigma^2 (x^\top x)^{-1})$$

⇒ t test for $\hat{\beta}_i$

$$3) \hat{e} = y - X\hat{\beta} = (I - H)y \quad H = X(x^\top x)^{-1} x^\top$$

$$\mathbb{E}[\hat{e}] = 0$$

$$\operatorname{Var}[\hat{e}] = \sigma^2 (I - H)$$

$$\begin{aligned} \operatorname{tr}(H) &= \operatorname{tr}(I_{p+1}) \\ &= p+1 \end{aligned}$$

$$\Rightarrow \mathbb{E}[\hat{e}^\top \hat{e}] = \operatorname{tr}[\operatorname{Var}[\hat{e}]]$$

$$= \operatorname{tr}[\sigma^2 (I - H)]$$

$$= \sigma^2 (n - p - 1)$$

⇒ $\hat{e}^\top \hat{e} / (n - p - 1)$ is an unbiased estimator for σ^2

$$4) \hat{\beta} \perp \hat{e} \Leftarrow \operatorname{cov}(\hat{\beta}, \hat{e})$$

$$= \operatorname{cov}((x^\top x)^{-1} x^\top y, (I - H)y)$$

(characterization of
chi-square dist)

$$5) \quad ① \text{ t-test : } H_0: \hat{\beta}_i = 0 \quad H_a: \hat{\beta}_i \neq 0$$

KEY POINT is:

$$Y^T A Y \sim \chi^2(p)$$

$$\Leftrightarrow \begin{cases} A \cdot E[Y] = 0 \\ A \text{ 对称零异} \\ \text{tr}(A) = p \end{cases}$$

$$SSE = \hat{e}^T \hat{e} = y^T (I - H) y \sim \sigma^2 \chi^2(n-p-1)$$

$$② \text{ SSR} = \underline{y^T (X(X^T X)^{-1} X^T - \frac{1}{n} \mathbf{1} \mathbf{1}^T) y}$$

(also need Cochran theorem
to characterize the independency)

对称零异

if $\beta_1 = \dots = \beta_p = 0$

$$\Rightarrow F\text{-test} \Rightarrow H_0: \beta_1 = \dots = \beta_p = 0$$

Notation

$$(y_i - \bar{y}) \quad (\hat{y}_i - \bar{y}) \quad (y_i - \hat{y}_i)$$

$$SST = SSR + SSE$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$TSS \quad ESS \quad RSS$$

$$H_2: 0 | \omega$$

$$\frac{SSR}{P}$$

$$\frac{SSE}{n-p-1}$$

$$\sim F(p, n-p-1)$$

③ A more general statement: (Constrained LR)

→ Problem Setting (Model)

$$\begin{cases} y = X\beta + \varepsilon \\ \varepsilon \sim N(0, \sigma^2 I_n) \end{cases}$$

$H\beta = c$ → Linear Constraints

→ Hypothesis testing:

$$\begin{cases} H_0: H\beta = c \\ H_a: H\beta \neq c \end{cases}$$

$$\Rightarrow \frac{\frac{SSE_H - SSE}{g}}{\frac{SSE}{n-m-1}} \sim F(g, n-m-1)$$

Rank: Previous

$$\frac{SST - SSE}{SSE}$$

is a special case of this result!

Here, we summarize on the Scenarios that violate assumptions

1) Briefly discuss

① different variances \rightarrow [e.g.] { poor
rich } Consumption

② co-variance $\neq 0 \rightarrow$ observations are co-related

③ Multi-collinearity \rightarrow features are too many

Property: (with respect to Design Matrix X)

Row-wise Issue: ① + ② \rightarrow observation-wise

Column-wise Issue: ③ \rightarrow attribute-wise

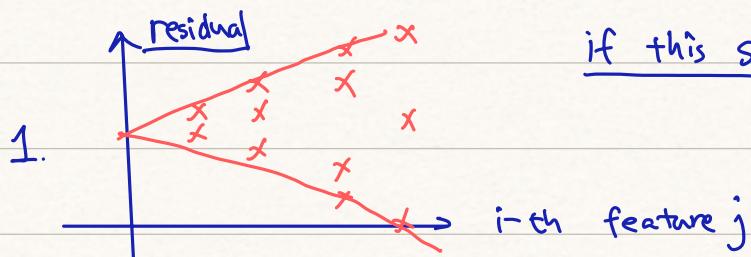
2) Details understanding (TBC 11/14)

12/12

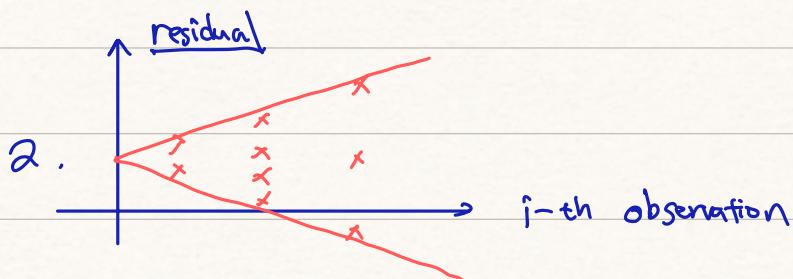
① Heterogeneity (different Variance)

$$\text{var}[\varepsilon_i] \neq \sigma^2$$

a) Residual-plot



if this situation exists for some x_i



b) Rank-test around same variance

(Non-parametric Statistical Testing)

Spearman Test

c) Recipe → WLS (weighted least square)

we should estimate the variance (weight)

Power Function

$$w_i = \frac{1}{6^2} \approx \frac{1}{k x_{ji}^2}$$

Here, $x_{ji} \rightarrow j\text{-th attribute for } i\text{-th observation}$

② Cor-related issue → $\text{Cov}(\varepsilon_i, \varepsilon_j) \neq 0$

a) Residual-Plot



b) Auto-correlation Coefficient → first-order correlation

$$\rho = \frac{\sum_{t=2}^n \varepsilon_t \cdot \varepsilon_{t+1}}{\sqrt{\sum_{t=2}^n \varepsilon_t^2} \sqrt{\sum_{t=2}^n \varepsilon_{t+1}^2}} \in [-1, 1]$$

Recap: if $E[X] = E[Y] = 0$

then $\text{corr}(X, Y)$

$$= \frac{E[X] \cdot E[Y]}{\sqrt{E[X^2]} \cdot \sqrt{E[Y^2]}} \in [-1, 1]$$

$$\varepsilon_t \leftarrow e_t$$

c) DW test $\rightarrow \varepsilon_t = \rho \varepsilon_{t-1} + u_t$

$$\left\{ \begin{array}{l} H_0: \rho = 0 \\ H_A: \rho \neq 0 \end{array} \right.$$

consider $DW = \frac{\sum_{t=2}^n (\varepsilon_t - \varepsilon_{t-1})^2}{\sum_{t=2}^n \varepsilon_t^2} \approx 2(1-\rho)$

Rmk: Under Null Hypothesis, $DW = \frac{\sum_{t=2}^n (u_t - u_{t-1})^2}{\sum_{t=2}^n u_t^2}$

$\{u_t\}_{t=1}^n \rightarrow$ i.i.d Gaussian

interpretation of the effectiveness

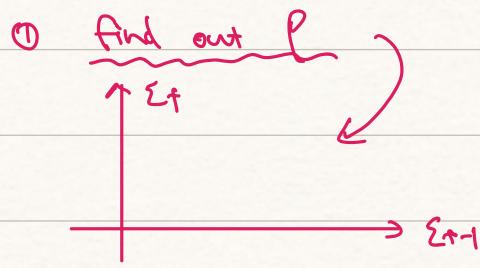
d) Recipe

- 1. add attributes
- 2. Assume it is first-order correlated,
then apply some models to handle this issue!

keep on iterating until pass DW Test!

2. one example:

Assume $\varepsilon_t = \rho \varepsilon_{t-1} + u_t$ (suppose it is correct)



although in real-application,

it is highly likely that this
is incorrect!

② modify the Regression Model:

$$Y_t - \rho Y_{t-1} = (\beta_0 - \rho \beta_0) + \beta_1 (x_t - \rho x_{t-1}) + (\varepsilon_t - \rho \varepsilon_{t-1})$$



y'_t

x'_t

$$\boxed{Y_t - \rho Y_{t-1}} = (1 - \rho) \beta_0 + \beta_1 \boxed{x_t - \rho x_{t-1}} + u_t$$

(adjusted model)

$$\left\{ \begin{array}{l} x'_t = x_t - \rho x_{t-1} \\ y'_t = y_t - \rho y_{t-1} \end{array} \right.$$

③ DW-test for Residual term

③ Col-linearity Issue

($X^T X$ is not full-rank)

1. Judgement

方差扩大量子 VIF

eigenvalue of $X^T X$ (Condition Number)

the measure of ill-condition for one matrix

2. Recipe → Generally speaking, 2 ways

$\left\{ \begin{array}{l} \uparrow \text{sample size} \\ \text{feature selection} \end{array} \right.$

Details:

a) PCA + LR

[idea]: Thru PCA, select those 'important' principal component

and conduct LR for them (subset of all PCs)

make sure the Design Matrix is well-conditioned

Model: ① Standard LR:

$$Y = \beta_0 \cdot \mathbf{1}_n + X\beta + \varepsilon$$

② PCA LR:

$$Y = \beta_0 \cdot \mathbf{1}_n + X \cdot Q \cdot Q^T \beta + \varepsilon$$

$$:= \beta_0 \mathbf{1}_n + Z \cdot \hat{\alpha} + \varepsilon$$

Then $Z^T Z = Q^T X^T X Q = \text{diag}(\lambda_1, \dots, \lambda_p)$

Idea: if $(Z^T Z)_{ii} = \lambda_i$ is small.

then it means the i-th attribute is almost constant

for all observation $\Rightarrow \hat{\alpha}_i$ can be set to 0

Conclusion: $\hat{\alpha} = \begin{pmatrix} \hat{\alpha}_1 \\ 0 \end{pmatrix} \Rightarrow \hat{\beta}^* = Q \hat{\alpha}$

$$= (Q_1 \ Q_2) \begin{pmatrix} \hat{\alpha}_1 \\ 0 \end{pmatrix}$$

assume $\hat{\alpha}_1 \in \mathbb{R}^r$

$$Z = (Z_1, Z_2)$$

$$Z_1 \in \mathbb{R}^{n \times r}$$

$$= Q_1 \hat{\alpha}_1$$

$$= Q_1 (Z_1^T Z_1)^{-1} Z_1^T Y$$

$$= Q_1 \Lambda_r^{-1} Z_1^T Y$$

$$= Q_1 \Lambda_r^{-1} (X \cdot Q_1)^T Y$$

$$= Q_1 \Lambda_r^{-1} Q_1^T \underline{X^T Y}$$

$$= Q_1 \Lambda_r^{-1} \underline{Q_1^T X^T X \hat{\beta}}$$

$$= Q_1 \Lambda_r^{-1} Q_1^T \underline{Q^T \hat{\beta}}$$

$$= Q_1 \boxed{\Lambda_r^{-1} (I_r \circ) \Lambda} Q^T \hat{\beta}$$

$$= Q_1 Q_1^T \hat{\beta} \Rightarrow \underline{(I_r \circ)}$$

Recap: $\hat{\beta} = (X^T X)^{-1} X^T Y$

$$\Rightarrow (X^T X) \cdot \hat{\beta} = X^T Y$$

Collary { 1. $\|\beta^*\| < \|\hat{\beta}\|$
 2. $\exists r.$ s.t. $\underline{\text{MSE}(\beta_r)} < \text{MSE}(\hat{\beta})$

b) Ridge Regression

idea: 1. bias - variance Trade-off

LR estimator \rightarrow { unbiased
 extremely large variance if $X^T X$
 is ill-conditioned

Surrogation \rightarrow estimator that { small bias
 small variance

2. Small component-wise value prior

\rightarrow we prefer the estimator that $\|\beta\|$ is small

Details: $\min_{\omega} \|y - X\beta\|_2^2 + \underbrace{\lambda \|\beta\|_2^2}_{\text{penalty term}}$

$$\Rightarrow \hat{\beta}_{RR}(\lambda) = (X^T X + \lambda I)^{-1} X^T y.$$

- Property:
- ① $\|\hat{\beta}_{RR}(\lambda)\| \leq \|\hat{\beta}\| \quad \text{for } \forall \lambda \geq 0$
 - ② $\hat{\beta}_{RR}(\lambda)$ is a biased estimator for β
 - ③ $\exists \lambda \text{ s.t. } \text{MSE}(\hat{\beta}_{RR}(\lambda)) < \text{MSE}(\hat{\beta})$

Rank: The Proof for Ridge Regression Property is strongly associated with the canonical form of Regression:

That is,

Standard form

$$y = \beta_0 \mathbf{1}_n + X\beta + \varepsilon$$

Canonical form

$$\text{Consider } Z = X \cdot Q$$

$$\alpha = Q^T \beta$$

$$\text{where } \underline{Q^T (X^T X) Q = \text{diag}(\lambda_1, \dots, \lambda_p)}$$

$$y = \beta_0 \mathbf{1}_n + Z \alpha + \varepsilon$$

$$\text{Here } ① \hat{\alpha}_{RR}(x) = Q^T \hat{\beta}_{RR}(x)$$

$$\begin{aligned} ② \hat{\alpha}_{RR}(\lambda) &= (Z^T Z + \lambda I)^{-1} Z^T y \\ &= (Z^T Z + \lambda I)^{-1} Z^T Z \cdot \hat{\alpha} \\ &= (\lambda + \lambda I)^{-1} \lambda \hat{\alpha} \end{aligned}$$

#