

① Classical Linear Model $Y = X\beta + \epsilon$ (X full rank)

Key Assumptions:

- Linearity of cov. effects \rightarrow transform
- Additivity of errors \rightarrow can apply log. transformation for multiplicative
- Error mean zero assumption: $E[\epsilon_i] = 0$
- Heteroscedastic error var: $\text{Var}[\epsilon_i] = \sigma^2 \neq \text{const}$
- Uncorrelated errors: $\text{Cov}(\epsilon_i, \epsilon_j) = 0 \quad i \neq j$
- Gaussian errors: $\epsilon_i \sim N(0, \sigma^2 I)$

Estimation: $\min RSS \quad z^T \epsilon$

$$- \|z\|_2^2 = y^T y - 2\beta^T z^T y + \beta^T z^T X \beta \quad (\text{same equations via MLE})$$

$$- \text{Ruler for diff: } \frac{\partial}{\partial \beta} \beta^T A \beta = 2AB$$

Pred. values and Residuals: $\hat{y} = X\hat{\beta} = X(X^T X)^{-1} X^T y = Hy$

Successive Orthogonalization (Gram-Schmidt Orthogonalization)

- Initial step: $z^0 = x^0 = y$
- Regress x_j on z^0, z_1, \dots, z_{j-1} (for $j=1, k$) to produce

$$\hat{y}_{ij} = \frac{x_j^T \cdot x_i}{x_i^T x_i}, \quad i=0, \dots, j-1 \quad \text{and} \quad z_i = x_i - \sum_{l=0}^{j-1} \hat{y}_{il} z^l$$

- Regress y on residuals z^k to give $\hat{\beta}_k = \frac{z^k^T \cdot y}{z^k^T z^k}$

\rightarrow QR decomp: Write output of step 2 as

$$X = Z D^T Q^T \Gamma \rightarrow \hat{\beta} = R^{-1} Q^T \gamma$$

Error Variance: $\hat{\sigma}^2 = \frac{1}{n-p} \hat{\epsilon}^T \hat{\epsilon} \rightarrow E[\hat{\sigma}^2] = \sigma^2$

Distribution of Parameter (without distri./Normality assumptions)

$$E[\hat{\beta}] = \beta \rightarrow \text{still unbiased}$$

$$\text{Cov}[\hat{\beta}] = \sigma^2 (X^T X)^{-1} \rightarrow \text{Cov}[\hat{\beta}] = \hat{\sigma}^2 (X^T X)^{-1} \text{ with } \hat{\sigma}^2 = \frac{1}{n-p} \hat{\epsilon}^T \hat{\epsilon}$$

$$\text{Var}[\hat{\beta}_j] = \sigma^2 / (z_j^T z_j) \rightarrow \text{high corr} \rightarrow \hat{\epsilon} \perp \hat{\epsilon} \rightarrow \text{Var} \uparrow$$

Gauss-Markov: LSE vs BLUE (without normality, others still hold)

$$\text{Var}(\hat{\beta}_j) \leq \text{Var}(\beta_j^{\text{true}}) \rightarrow \text{smallest var of all lin. unbiased}$$

UNIV: Ctr. Ass. + normality of error

\rightarrow LSE has min. var among all unbiased (not only linear) estimators
 \rightarrow attains the Cramér-Rao lower bound

$$\hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$$

Asymptotics: Consistency and asympt. normality of LSE \rightarrow 2 assump.

$$1. \lim_{n \rightarrow \infty} (X_n^T X_n)^{-1} = \lim_{n \rightarrow \infty} (\sum_{i=1}^n x_i x_i^T)^{-1} = 0 \quad \text{"covariate onto T as n \uparrow"}$$

$$2. \lim_{n \rightarrow \infty} \max_{i=1, \dots, n} x_i^T (X_n^T X_n)^{-1} x_i = 0 \quad \text{"influence of each x_i is negligible relative to the entire into X_n^T X_n"}$$

Properties of residuals

$$\hat{\epsilon} = y - \hat{y} = (I - H)y \rightarrow \text{residuals are orthogonal to X: } X^T \hat{\epsilon} = 0$$

$$E[\hat{\epsilon}] = 0, \text{Cov}[\hat{\epsilon}] = \sigma^2 (I - H), \text{Var}[\hat{\epsilon}_i] = \sigma^2 (1 - h_{ii})$$

$$\hat{\epsilon}^T \hat{\epsilon} = 0 \quad (\text{residuals are pred. values are orthogonal + independent})$$

② Hypothesis Testing

General linear hypothesis: $C\beta = d$ vs. $C\beta \neq d$

rank matrix = rk(C) = r < p ($r = \# \text{constraints}$)

F-stat: $\frac{n-p}{r} \frac{\text{ASSE}}{\text{SSE}} \sim F_{r, n-p}$ (derivation based on the normality of errors)

\rightarrow Wald Test: $F = [((\hat{\beta} - d)^T \text{Cov}(\hat{\beta}))^{-1} ((\hat{\beta} - d))] / r \rightarrow (W = r) F$

③ Test of sign. (t-test): $t_{ij} = \frac{\hat{\beta}_j}{\text{sej}} \sim t_{n-p}$ as $F = \frac{\hat{\sigma}^2}{\text{Var}(\hat{\beta}_j)} \sim F_{r, n-p}$

④ Composite test of subvector (Partial F-test): $\hat{\beta}_1 = (\beta_1, \dots, \beta_r)^T$

$$H_0: \hat{\beta}_1 = 0 \quad \text{vs.} \quad H_1: \hat{\beta}_1 \neq 0$$

$$\rightarrow F = \frac{1}{r} \hat{\beta}_1^T \text{Cov}(\hat{\beta}_1)^{-1} \hat{\beta}_1 \sim F_{r, n-p} \quad \text{where } \text{Cov}(\hat{\beta}_1) \text{ would be cov. elements of } X^T X^{-1} X^T$$

⑤ Test for sign. of regression (Global F-test)

$$H_0: \beta_i = 0 \quad \forall i \in k, \quad H_1: \text{at least one} > 0 \rightarrow \frac{n-p}{k} \frac{R^2}{1-R^2} \sim F_{k, n-p}$$

2. Confidence and Prediction Intervals

$$⑥ \text{CI based on } t_{ij}: P(\hat{\beta}_j - t_{n-p, 1-\alpha/2} \text{sej} < \beta_j < \hat{\beta}_j + t_{n-p, 1-\alpha/2} \text{sej}) = \alpha \rightarrow \text{can extract } \hat{\beta}_j$$

⑦ For multiple parameters (subvector β_1) \rightarrow Confidence Ellipsoid

$$CE: [P_1: \frac{1}{r} (\hat{\beta}_1 - \beta_1)^T \text{Cov}(\hat{\beta}_1)^{-1} (\hat{\beta}_1 - \beta_1) \leq F_{r, n-p} (1-\alpha)]$$

⑧ CI for $\mu_0 = E[y_0]$ (for future obs. y_0 at known x_0)

$$[x_0^T \hat{\beta} \pm t_{n-p, 1-\alpha/2} \hat{\sigma} (1 + x_0^T (X^T X)^{-1} x_0)^{1/2}]$$

⑨ Prediction Intervals (for prediction $\hat{y}_0 = x_0^T \hat{\beta}$ for new/future obs. x_0)

$$[x_0^T \hat{\beta} \pm t_{n-p, 1-\alpha/2} \hat{\sigma} (1 + x_0^T (X^T X)^{-1} x_0)^{1/2}] \quad \text{contains future obs. } y_0 \text{ with } Pr \geq 1-\alpha$$

3. Multiple Testing

$$⑩ p\text{-value: } = P_{H_0} (|T| \geq |t_{\text{obs}}|) \sim U(0,1) \quad (\text{under } H_0)$$

Errors:	H ₀	H ₁	Residuals
Decision	H ₀ Type I = α	correct	
	H ₁ correct	Type II = β	<ul style="list-style-type: none"> discovery = rejecting H₀ false discovery = type I error power = $1 - \text{Type II error} (1 - \beta)$

\rightarrow as $n \uparrow$ or eff. size $\uparrow \rightarrow \beta \downarrow \rightarrow (1-\beta) \uparrow$ but \propto same

$$⑪ \text{FWER: } P(V \geq 1) = 1 - (1-\alpha)^m \quad (m = \text{total # of hypotheses})$$

* for independent test = α
* for pos. dependent tests: $\text{FWER} < \alpha$ (ind.)
* for neg. dependence $\rightarrow \text{FWER} > \alpha$ (ind.)

- Bonferroni: reject $H_0^{(i)}$ when $p_i \leq \frac{\alpha}{m}$; can also be used to construct CI

$\hookrightarrow \alpha^* = \alpha/m, \quad \text{OR } p^* = p \cdot m$

\hookrightarrow also needs no assumption on dependence structure of $H^{(i)}$ \rightarrow more powerful than BF

- Permutation test: permute test statistic $T(S_x, S_y) \rightarrow P_H(|T| \geq |T(S_x, S_y)|)$

\rightarrow calc p-values after permutation: $p\text{-value} = \frac{1}{m} \sum_{i=1}^m I(|T^{(i)}| \geq |T(S_x, S_y)|) + 1$

* non-parametric approach that works with finite sample
* comp. expensive + only a global test

- FDR: weaker error control than FWER (FDR also controls FDR)

$$FDR = E[FDP] = E[V/R] \quad (\text{exp. false disc. / all disc.})$$

- BH_α: order p-values (same as for BH₀)

* R be the largest k s.t. $p_k \leq \frac{\alpha}{k}$ \rightarrow R will be $\#$

* Reject all null hypotheses for $p_i \leq \frac{\alpha}{m}$ which $p_i \leq p_{k,R}$

\hookrightarrow only works well for independent and pos. dependent (PROS) p-values

\hookrightarrow under arbitrary dependence: $\text{FDR} \approx \alpha \cdot \log(m)$ and FDP high variance!

4. Model Specification

Truth	Reduced model ($\hat{\beta}_1$)	Full model ($\hat{\beta}_2$)
-------	-----------------------------------	--------------------------------

Reduced model ($\hat{\beta}_1$): $E[\hat{\beta}_1] = \beta_1$ both unbiased
 $\text{Cov}(\hat{\beta}_1) = \sigma^2 (X^T X)^{-1}$

Full model ($\hat{\beta}_2$): $E[\hat{\beta}_2] = \beta$

$\text{Cov}(\hat{\beta}_2) = \sigma^2 (X^T X)^{-1}$

③ Prediction Quality: \rightarrow Consider SPSE for assessment

$$\text{SPSE} = \sum_{i=1}^n (y_{ni} - \hat{y}_{ni})^2 \quad \hat{y}_{ni} \text{ is pred. based on a model}$$

$$= n \cdot \sigma^2 + \text{IM}(\sigma^2 + \sum_{i=1}^n E[(\hat{y}_{ni} - y_{ni})^2])$$

(I) SMSE (II) MSPE (III) MSE

I) Irreducible prediction error \rightarrow just depends on σ^2

II) Variance: depends on $(M \neq \# \text{ variables included in model})$

III) Squared bias: smaller as more variables are included (as they have expl. power)

5. Model Choice Criteria: need to use estimated SPSE:

a) Estimate using independent data \rightarrow split into training + validation

b) Use all existing data to estimate $\hat{\beta}_M \rightarrow$ correct training error $\rightarrow C_p$

$$④ \text{Hartman's } C_p: C_p = \sum_{i=1}^n (y_i - \hat{y}_{ni})^2 - n + 2(M+1) \rightarrow \text{Min. } \sigma^2 \text{ est. on full model}$$

$$⑤ \text{AIC: } AIC = -2 \ell(\hat{\beta}_M, \hat{\sigma}_{ML}^2) + 2(C_p+1) \rightarrow \text{Min.}$$

$$= n \cdot \log(\hat{\sigma}_{ML}^2) + 2(M+1) \text{ for lin. model with } \epsilon \sim N(\mu, \sigma^2)$$

$$⑥ \text{BIC: } BIC = -2 \ell(\hat{\beta}_M, \hat{\sigma}_{ML}^2) + \log(n)(M+1)$$

$$= n \cdot \log(\hat{\sigma}_{ML}^2) + \log(n)(M+1) \text{ for lin. model with } \epsilon \sim N(\mu, \sigma^2)$$

⑦ Adjusted Coeff. of Determination (R^2): always smaller than mult. R^2

$$R^2 = \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{\sum \hat{\epsilon}_i^2}{\sum (y_i - \bar{y})^2} \rightarrow \bar{R}^2 = 1 - \frac{(n-p)}{(n-p)} (1-R^2)$$

⑧ Cross Validation: can address overfit/selection bias by splitting up data

- Partition data set in k sets of similar size

- take set i as validation set, $k-1$ sets left for estimation

\rightarrow set SPSE for all k subsets

\rightarrow CV = $\sum \text{SPSE}^{(i)}$ \rightarrow use model with smallest sum of SPSE

⑨ LOOCV: use all observations except one data point

$$- CV = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_{-i})^2 \quad (= \frac{1}{n} \sum_{i=1}^n \frac{(y_i - \hat{y}_{-i})^2}{n-1}) \text{ for LSE}$$

\hookrightarrow select model using the one-se rule

6. Model Selection Procedures: (if no pre-selection exists)

a) Complete model selection \rightarrow all start with H₀ (null model)

1. go through every model for $k=1, 2, \dots, p \rightarrow (P_k)$ possible models

2. Pre-select from models with same p using RSS (less diff. specification)

3. Select single best model based on model choice criteria

\rightarrow can apply leaps & bounds alg.

b) Forward Selection

* consider all $p-k$ models that add an add. predictor to M_k

* as above \rightarrow variable which offers greatest reduction in the pre-selected model choice criteria is chosen

* also applicable if $p > n$ (unlike back ward selection)

c) Backward Selection: as above but starts with full model

d) Stepwise Selection: combination of b) & c), every step addition or deletion of a variable is considered

* Model Selection based on sign. not useful as t-test would not $\sim t$

\rightarrow as $n \uparrow$, any variable C $\neq 0$ would become sign.

7. Model Diagnostics: assessing validity of the model

- look at residuals (but can be heteroscedastic & correlated)

\rightarrow a) Standardized residuals (not useful as unknown dist.).

b) Studentized residuals (based on LS-estimator)

$$t^* = \frac{\hat{\epsilon}_i}{\hat{\sigma}_{\text{est}}} \frac{1}{\sqrt{1 + (\hat{\epsilon}_i^2 / \hat{\sigma}_{\text{est}}^2) \cdot x_i^T X^{-1} x_i}} = \frac{\hat{\epsilon}_i}{\hat{\sigma}_{\text{est}} \sqrt{1 + \hat{\epsilon}_i^2 / \hat{\sigma}_{\text{est}}^2 \cdot x_i^T X^{-1} x_i}} \sim t_{n-p-1}$$

A + B) Mean zero & Homoscedasticity

- 1) Tukey-Anscombe Plot: raw residuals $\hat{\epsilon}_i$ vs. $\hat{\epsilon}_i$ plot
 - 2) Scale-Location plot: (standardized residuals) $\hat{\epsilon}_{i:n}$ vs. $\hat{\epsilon}_i$ plot
- cf heteroscedastic errors do seem to occur: Transformation OR WLS
- i) fix non-linearity → pol. transform x and y
 - ii) stabilize the var → transform y
- often use particular transformation family:
- Power family:** x^λ for $\lambda=0$ and $\log(x)$ (for $\lambda=0$)
 - Transforming only predictor:** Scaled power
- $\psi_\lambda(x_i) = \begin{cases} x_i^{-1} & \text{for } \lambda \neq 0 \\ \log(x_i) & \text{for } \lambda = 0 \end{cases}$ preserves direction of assoc. of x and y (not as power family)

- Transforming response:** Box-Cox Procedure
 - need to have $y > 0$ (otherwise add constant)
 - automatically chooses a transformation (estimates λ) from the (scaled) power family

C) Variance-Stabilizing Transformations: most common $\log(y)$, also \sqrt{y}

D) Auto-correlated errors: $\varepsilon \sim N(0, \Sigma) \rightarrow \beta \sim N(\beta, (X^T X)^{-1} \Sigma X (X^T X)^{-1})$

- AR1 process: $\varepsilon_i = \rho \varepsilon_{i-1} + u_i, -1 < \rho < 1$ size of corr. depends on X & Σ
- Then $E[\varepsilon] = 0, \text{Var}[\varepsilon] = \frac{\sigma^2}{1-\rho^2} \begin{pmatrix} 1 & \rho & \dots & \rho^{n-1} \\ \vdots & \ddots & \ddots & \vdots \end{pmatrix} = \sigma^2 W^{-1}$
- decay of correlations between ε_i and ε_{i-j} as $j \uparrow$
 - for $\rho > 0 \rightarrow$ geometric decay
 - for $\rho < 0 \rightarrow$ corr. decreases with alternating signs

$$- ACF(\varepsilon_j) = \frac{\text{Cov}(\varepsilon_i, \varepsilon_{i+j})}{\sqrt{\text{Var}[\varepsilon_i]}} = \rho^j \rightarrow \text{for AR1: } j=1 \text{ can both be estimated}$$

$$- PACF(j) \text{ is } j\text{th reg. coef. in } \varepsilon_i = \alpha_0 \varepsilon_{i-1} + \dots + \alpha_j \varepsilon_{i-j} + u_i \rightarrow \text{diagnose by looking at residual plots over time}$$

E) Normality of errors: use Q-Q / Normal Plot

F) Linear model → multicollinearity issue

$$\begin{aligned} - \text{Var}(\hat{\beta}_i) &= \frac{\sigma^2}{\hat{\varepsilon}_i^T \hat{\varepsilon}_i} = \frac{\sigma^2}{(1-R_i^2) \sum (x_{ii}-\bar{x}_i)^2} & \bullet \hat{\varepsilon}_i \downarrow \text{with corr. } X \\ - VIF &= \frac{1}{1-R_i^2}, \text{ problem of } > 10 & \bullet R_i^2 \uparrow \text{with higher dependence} \\ \text{where } R_i^2 &\text{ is } R^2 \text{ from reg. on } x_i \end{aligned}$$

$$- \text{Scaled Residuals: } \tilde{\varepsilon}_i = \frac{y_i - x_i^T \hat{\beta}}{\|y_i - x_i^T \hat{\beta}\|_2} = \frac{(I-H)x_i^T \varepsilon_i}{\|x_i^T \varepsilon_i\|_2} \text{ scaled by L2 norm}$$

and $\tilde{\varepsilon}_i \sim N(0, 1) \rightarrow$ simulate under H_0 and compare with the observed scaled residuals

- **Residual Predictor Test:** $H_0: y = x^T \beta + \varepsilon$ vs $H_1: y = f(x) + \varepsilon$ → compare scaled residuals for both models → can get p-values as we know that under $H_0 \tilde{\varepsilon}_i \sim N(0, 1)$

8. Outlier Detection under $E[\varepsilon_i | x_i] = x_i^T \beta + \eta$

- may affect both parameter estimates (part. if not out)
- look for large residuals (\rightarrow see studentized residuals)
- observations with large r^* are pot. outliers (can also test)
- **Leverage:** outliers in x direction (extreme values of x)
 - diagonal elements of H : $\frac{1}{n} \leq h_{ii} \leq 1 \rightarrow h_{ii} > \frac{2p}{n}$ examine!
 - large leverage $\rightarrow \text{Var}(\hat{\varepsilon}_i)$ small
 - for only one covariate x_i : $h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_j - \bar{x})^2}$
 - lever ↑ as x_i further away from \bar{x}
 - obs. with large lever → large influence on analysis

Cooks Distance: compare $\hat{\varepsilon}_i$ vs. $\hat{\varepsilon}_{(i)}$ → look at each distance

$$D_i = \frac{1}{p} \frac{\hat{\varepsilon}_i^2}{\hat{\varepsilon}_{(i)}^2 (1-h_{ii})} \cdot \frac{h_{ii}}{1-h_{ii}} \text{ depends on both levs and residuals}$$

$$= (\hat{\varepsilon}_{(i)} - \hat{\varepsilon}_i)^T (\hat{\varepsilon}_{(i)} - \hat{\varepsilon}_i) \quad D_i > 0.5 \text{ attention} \\ p \cdot \hat{\varepsilon}_{(i)}^2 \quad D_i > 1 \text{ always examine}$$

G) General Linear Model (with heteroscedastic and/or corr. errors)

- Now: $\text{Cov}(\varepsilon) = \sigma^2 W^{-1}$ ($\varepsilon \neq 0$), for hetero var but no correlation: $\text{Cov}(\varepsilon_i) = \frac{\sigma^2}{w_i}$
- Aitken/WLS Estimator:** $\hat{\beta}_{WLS} = (X^T W X)^{-1} X^T W Y$ (\Rightarrow MLE for $\sigma^2 W$)
- from minimizing: $\text{WLS}(\beta) = (Y - X\beta)^T W (Y - X\beta)$ for hetero + uncor.
- Properties: $E[\hat{\beta}_{WLS}] = \beta$, $\text{Cov}(\hat{\beta}_{WLS}) = \sigma^2 (X^T W X)^{-1}$, GLM holds
- with $\text{err} \sim N(0, \sigma^2 W^{-1})$: $\hat{\sigma}_{ML}^2 = \frac{1}{n} \sum \hat{\varepsilon}_i^2 W_i \rightarrow$ biased estimator
- \rightarrow REML estimator: $\frac{1}{n-p} \hat{\varepsilon}^T W \hat{\varepsilon}$ GLS residual

Two-stage Least Squares (Simultaneous estimation of β , σ^2 , w_i) → want to estimate the var. key: $\hat{\varepsilon}^2 = \sigma^2 + v \rightarrow \hat{\sigma}^2 = \hat{\varepsilon}^2 - v$

1. Obtain $\hat{\beta}$ through OLS → get residuals $\hat{\varepsilon}_i$ ($Z^T \hat{\varepsilon}_i = x_i^T$)
2. Obtain $\hat{\sigma}^2$ from uncorrected reg. Cols $\rightarrow \hat{\varepsilon}_i \sim Z_i$
- fit GLM with $\hat{w}_i = \frac{1}{\hat{\sigma}^2} \hat{\varepsilon}_i^2 \rightarrow \hat{\sigma}^2 = \exp(\hat{\varepsilon}_i^2/2)$ to ensure $\hat{\varepsilon}_i^2 > 0$

$$\text{White Estimator: } \text{Cov}(\hat{\beta}) = (X^T X)^{-1} X^T \text{diag}(\hat{\varepsilon}_1^2, \dots, \hat{\varepsilon}_n^2) X (X^T X)^{-1}$$

Two-stage Least Squares for AR(1): as above but in step 2:

$$\hat{\beta} = \sum_{i=1}^n \hat{\varepsilon}_i \hat{\varepsilon}_{i-1} / (\sqrt{\sum \hat{\varepsilon}_i^2} \sqrt{\sum \hat{\varepsilon}_{i-1}^2})$$

3. Insert $\hat{\beta}$ into weight matrix $W \rightarrow \hat{W} \rightarrow$ apply GLS

10. Robust Regression (when we have obs. with high influence and/or leverage)

Influence function: local concept → measure sensitivity of $\hat{\beta}$ w.r.t. one outlier

- Version of canonical IF for $n \rightarrow \infty$: consider cont. dist. of $\hat{\beta}$ $F = (1-\varepsilon) F + \varepsilon \delta_x$ (as $\varepsilon \rightarrow 0$) measures sensitivity of $\hat{\beta}$ to an small change ε
- Shape of IF for LSE: $\text{IF}_{LS}(\varepsilon) = 2\varepsilon$ (prop. to ε)

Breusch-Pagan: global concept → how much contaminated data until we lost

- Finite sample BP: fraction of data that can be given arbitrary values without making estimator arb. bad
↳ cannot be $> 3\%$
- asymptotic BP: limit of the finite sample BP as $n \rightarrow \infty$
e.g. for OLS (with mean) the BP = $1/n$; median would be $2/\sqrt{n}$

M-estimator: method of estimation; pays less attention to obs. with large

- 1) min. obj. function: $\sum \rho(\varepsilon_i) = \sum \rho(y_i - x_i^T \beta)$ $\rho(\cdot)$ is obs. function

- 2) \rightarrow min. means: $0 = \frac{\partial}{\partial \beta} \sum \rho(\varepsilon_i) = \sum \psi(y_i - x_i^T \beta) \cdot x_i^T$

- 3) Set $\omega(\varepsilon) = \psi(\varepsilon)/\varepsilon \rightarrow \sum \omega_i (y_i - x_i^T \beta) x_i^T = 0$ ↳ influence aware!
- 4) \rightarrow like WLS → solve using IRLS

Objective functions: → the PCG for M-Estimation

- a) LS estimator: $\rho_{LS}(\varepsilon) = \varepsilon^2, \omega_{LS}(\varepsilon) = 1$
- b) LAD estimator: $\rho_{LAD}(\varepsilon) = \varepsilon_{+}, \omega_{LAD}(\varepsilon) = 1$

- c) Huber loss: combines LS and LAD for large $|\varepsilon|$ $\rho_{Huber}(\varepsilon) = \begin{cases} \frac{1}{2} \varepsilon^2 & \text{if } |\varepsilon| \leq c \\ c(|\varepsilon| - c/2) & \text{if } |\varepsilon| > c \end{cases}$

- d) Biweight loss: for large $|\varepsilon|$ \rightarrow weighting becomes constant / levels off

$$\rho_{Biweight}(\varepsilon) = \begin{cases} \frac{c}{2} \varepsilon^2 - [1 - (\frac{\varepsilon}{c})^2]^2 & \text{if } |\varepsilon| \leq c \\ \frac{c^2}{2} & \text{if } |\varepsilon| > c \end{cases}$$

with $c \approx 7 \cdot \text{MAD}$

11. Generalized Linear Models (GLMs) - Binary $y \sim \text{Bin}(n)$

- **Logit Model:**

$$\begin{aligned} \text{Logistic response: } \pi &= h(\eta) = \exp(\eta) / (1 + \exp(\eta)) \\ \text{Logit link: } g(\pi) &= \log\left(\frac{\pi}{1-\pi}\right) = \log\left(\frac{P(y=1)}{P(y=0)}\right) = \eta = x^T \beta \end{aligned}$$

\rightarrow effect on odds

- **Probit Model:** Probit response: $\pi = \Phi(\eta) = \exp(-\eta^2/2) / (1 + \exp(-\eta^2/2))$

$$\text{Var: } \sigma^2 = \pi(1-\pi) \rightarrow \text{to compare logit and probit } (\sigma^2 = \frac{\pi}{n})$$

- **Latent linear model:** $y_i = \sum_{i=1}^n \varepsilon_i \leq 0$ Assume $\varepsilon_i = x_i^T \beta - \eta_i$

$$\pi_i = P(y_i=1) = P(\varepsilon_i \leq 0) = P(x_i^T \beta - \eta_i \leq 0) \rightarrow \text{Prob}(\text{given } x_i \text{ often round} \rightarrow \text{Prob about } \eta_i)$$

- for $\varepsilon_i \sim N(0, \sigma^2)$ → only identifiable to $\eta = \beta^T x_i$ due to std. → look at β_i/σ

- for $\varepsilon_i \sim \text{Logistic}$ → Logit

- **For Grouped Data:** rel. freq. (per group): $\bar{Y}_i \sim \text{Bin}(n_i, \pi_i) / n_i$

$$\cdot E[Y_i] = \pi_i, \text{Var}[Y_i] = [\pi_i(1-\pi_i)] / n_i = [\bar{Y}_i(1-\bar{Y}_i)] / n_i \text{ (estimate)}$$

- **Overdispersion:** correct with $\theta > 1$: $\text{Var}[Y_i] = \theta / n_i$

- **Maximum Likelihood Estimation**: when there is no closed form solution to $S(\beta) = 0$

a) **Newton-Raphson method**

$$1. \text{ If } \theta \text{ save } \theta \text{ approx. theory 1st order Taylor expansion: } S(\theta^{(0)}) + \frac{1}{2} S''(\theta^{(0)}) (\theta - \theta^{(0)})$$

$$2. \text{ Improve approx. iteratively: } \theta^{(t+1)} = \theta^{(t)} + I(\theta^{(t)})^{-1} \cdot S(\theta^{(t)})$$

b) **Fisher Scoring algorithm**

$$2. \text{ use } J(\theta^{(t)}) \text{ instead obs. T: } \hat{\theta}^{(t+1)} = \hat{\theta}^{(t)} + J(\theta^{(t)})^{-1} \cdot S(\theta^{(t)})$$

$$\text{with exp. F.T. } J(\theta) = E[C(S(\theta))] = E[S(C(\theta)) \cdot S(C(\theta))]$$

→ when estimating β in logit: $\hat{\beta}^{(t+1)} = \beta^{(t)} + (X^T U^{(t)})^T X^T (y - \pi^{(t)})$

3. Stop when convergence criterion is met ↳ need X full rank for converge

- **Testing linear hypotheses**: same setup as in classical linear model with diff. statistics:

a) **Likelihood Ratio Test:** $\text{LR} = 2(l(\hat{\beta}) - l(\beta)) \stackrel{n \rightarrow \infty}{\sim} \chi^2_r$

b) **Wald test:** $\text{Wald} = (C\hat{\beta} - d)^T (C(C\hat{\beta} - d)^T)^{-1} (C\hat{\beta} - d) \stackrel{n \rightarrow \infty}{\sim} \chi^2_r$ ↳ reject if LR, Wald are "large"

- **Deviance**: measures goodness-of-fit for GLMs

• compares perfect fit / saturated model to training model

a) For ord. data: $D(\hat{\beta}) = -2 \ell(\hat{\beta})$

b) For grouped data: $D(\hat{\beta}) = 2(l(\hat{\beta}) - l(\beta))$ with $\hat{\beta}_i = \frac{y_i - \bar{y}_i}{n_i}$ ↳ if model are nested = LR

- **Count Data Regression**: often based on Poisson distribution

a) **Log-Linear Poisson model**: $\lambda_i = \exp(\eta_i)$ with $\eta_i = x_i^T \beta$

b) **Linear Poisson model**: $\lambda_i = x_i^T \beta$ (if x have a linear effect on λ)

- **Unified framework for GLMs**: $\mu = E[y] = h(\eta)$

• part of univariate exp. family: $f(y|\eta) = \exp\left\{\frac{y\theta - b(\theta)}{\phi} + c(y, \theta, \phi)\right\}$ θ : canonical parameter (of interest), ϕ : dispersion parameter (variance), c : $c(y, \theta, \phi)$

- **Penalized Regression**: column LSE unstable (collinearity or $p > n$)

- Penalized LS objective: $\text{PLS}(\beta) = (y - X\beta)^T (y - X\beta) + \lambda \text{pen}(\beta)$

- **Ridge Regression**: $\text{pen}(\beta) = \| \beta \|^2 = \beta^T \beta \rightarrow \hat{\beta}_{\text{Ridge}} = (X^T X + \lambda I)^{-1} X^T y$

• $E[\hat{\beta}_{\text{Ridge}}] = (X^T X + \lambda I)^{-1} X^T y$ ↳ biased

• $\text{Cov}[\varepsilon_i] = \sigma^2 I$; $\text{Cov}[\hat{\beta}_{\text{Ridge}}] = \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}$ ↳ smaller var than $\hat{\beta}_{\text{OLS}}$ for $\lambda > 0$

• For orthogonal covariates: $\hat{\beta}_{\text{Ridge}} = n/(n+\lambda) \beta$ $\hat{\beta}_{\text{OLS}}$ are $\text{Cov}(\hat{\beta}_{\text{OLS}}) = \sigma^2 (I/n)^{-1}$

• high Var / Corr. variables are penalized less

- **LASSO Regression**: $\text{pen}(\beta) = \lambda \|\beta\|_1 \rightarrow \hat{\beta}_{\text{LASSO}} = \arg \min_{\beta} (y - X\beta)^T (y - X\beta) + \lambda \sum_{i=1}^p 1(\beta_i \neq 0)$

• chsl. model selection → can shrink $\beta_i = 0$ ↳ # non-zero parameters

• also biased but with $\text{Var} \leq \text{Var}(\hat{\beta}_{\text{OLS}})$

- **Best subset Selection**: $\arg \min_{\beta} (y - X\beta)^T (y - X\beta) + \lambda \sum_{i=1}^p 1(\beta_i \neq 0)$