

Statistical Modelling

Cheat Sheet · v1.0 · 2021

Mael Macuglia

D-MATH
ETH Zürich

1 Classical Linear Model

Model: Let $\mathcal{Y} := (Y_1, \dots, Y_n)$ ind. copies of Y
 $\mathcal{Y} = f + \epsilon$ with $EY = f, f \in \mathcal{F} \subseteq \mathbb{R}^n$
Linear Model: $\mathcal{F} := \{f_i = a + \sum_{j=1}^p x_{ij}\beta_j, a \in \mathbb{R}, b \in \mathbb{R}^p\}$
Goal: find an estimate for the systematic component
 $EY = f = X\beta$

1.1 Modelling Effect of Covariates

- Continuous Covariates:** z has non-linear effect $\beta_1 f(z)$ with known f : Linear model $\rightarrow y_i = \beta_0 + \beta_1(x_i)$ where $x_i := f(z_i) - \bar{f}$. z has approximately polynomial effect $\beta_0 + \beta_1 z + \dots + \beta_l z^l \rightarrow x_{ij} = z_i^j$ (apply centering) (use partial residuals to check the effect of modelling)
- Categorical Covariates** Let the covariate $x_i \in \{1, \dots, c\}$ then apply Dummy Encoding: $x_{i1} = \mathbb{1}(x_i = 1), \dots, x_{ic-1} = \mathbb{1}(x_{c-1} = 1)$ where the category omitted (needed for uniqueness of LS) is called the reference category
- Interactions** catcat/catcont: $y = \beta_0 + \beta_1 x + \beta_2 z + \beta_3 xz + \epsilon$ Think: x, z not need to be correlated has interaction has effect on y | Think: β_1 effect of having A $\beta_2 =$ effect of having B $\beta_3 =$ additional effect of having A and B | Think interaction when effect of cat. cov. do not translate.

1.2 Parameter Estimation β

- Assumptions** Homoscedasticity and Uncorrelated errors: $cov(\epsilon) = \sigma^2 I_{n \times n} \rightarrow$ same precision for all measurements | $rank(X) = p$.
- RSS** := $\|Y - X\beta\|^2$
- minLS:** $\hat{\beta} = \min \|Y - X\beta\|^2 = \min_{\beta} \epsilon^T \epsilon$
- LSE:** $\hat{\beta} = (X^T X)^{-1} X^T Y$
- Prediction** $EY = \hat{Y} = X\hat{\beta} = HY = \mathcal{Y}P\mathcal{X}$
- Hat Matrix** $SVD(X) = P\Phi V^T \Rightarrow H = PP^T$ Trace: $tr(H) = p$ | Anti projection: $Q = (1 - H)$ Trace: $tr(Q) = n - p$ | both idempotent and symmetric
- Residuals** $\hat{\epsilon} = Y - X\hat{\beta} = QY = \mathcal{Y}\mathcal{A}\mathcal{X}$

1.3 Parameter Estimation σ

- MLE:** $MLE(\sigma^2, \beta = \hat{\beta}_{LS}) = \sigma_{MLE}^2 = \frac{1}{n} \|\hat{\epsilon}\|^2$ (biased)
- Unbiased Estimator:** $E\|\hat{\epsilon}\|^2 = (n - p)\sigma^2 \Rightarrow \hat{\sigma}^2 = \frac{1}{n - p} \|\hat{\epsilon}\|^2$ (unbiased)
- Residual Standard Error** $RSE := \hat{\sigma} = \sqrt{\frac{1}{n - p} \|\hat{\epsilon}\|^2}$

1.4 Geometric Properties

- Normal Eq:** $\mu = EY = X\beta \in \text{Span}(1, x_1, \dots, x_{p-1}) = p\text{-dim Vector Space}$ Minimize vector of residuals: $X^T(Y - X\hat{\beta}) = 0$

- Orthonormal Design:** In the regression of y onto (x_0, x_1, \dots, x_k) the coefficient of x_k is $\hat{\beta}_k = \frac{\langle z_k, y \rangle}{\langle z_k, z_k \rangle}$ where z_k is the residual from regressing x_k onto $(x_0, \dots, x_{k-1}) \Rightarrow z_k^T x_j = 0 \forall j \in \{0, 1, \dots, k - 1\}$
- Interpretation:** The multiple regression coeff. $\hat{\beta}_k$ represents the additional contribution of x_k on Y after x_k has been adjusted for x_0, \dots, x_{k-1}

1.5 Analysis of Variance

- Sample Variance:** $S^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2$
- Sample Pred. Var.:** $S_Y^2 = \frac{1}{n} \|\hat{y} - \bar{y}\|^2$
- Variance Decomp.:** $\|y - \bar{y}\|^2 = \|y - \hat{y} + \hat{y} - \bar{y}\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - \bar{y}\|^2 + 2 \langle y - \hat{y}, \hat{y} - \bar{y} \rangle$ where we note that in case of LS the last scalar product is = 0
- $SS_{tot} = SS_{reg} + SS_{res}$
- Residual variance is the variance we couldn't explain
- R Coeff:** $R^2 = \frac{SS_{reg}}{SS_{tot}} = \frac{\|\hat{y} - \bar{y}\|^2}{\|y - \bar{y}\|^2} = 1 - \frac{SS_{res}}{SS_{tot}}$
- Interp.:** R^2 proportion of variance explained by the model

Model Comparison using R^2 only meaningful if following 3 statements are fulfilled:

- Every model has the same response y (not transform)
- Every model has the same dimension of parameter space (nb of p param.)
- Every model has an intercept

Note: By LLN the sample variance converges in probability to the variance as $n \rightarrow \infty$

1.6 Statistical Properties

- $E\hat{\beta} = \beta$ and $cov(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$
- $Var(\hat{\beta}_j) = \frac{\sigma^2}{\langle z_j, z_j \rangle} \rightarrow$ correlation between the covariates inflate the variance.
- Est. Var.** $cov(\hat{\beta}) = \hat{\sigma}^2 (X^T X)^{-1} | \hat{\sigma}^2 = \frac{1}{n - p} \epsilon^T \epsilon$
- Estimated Standard Error:** (Std.error in R output) $se_j = (cov(\hat{\beta}))_{jj} = Var(\hat{\beta}_j)^{\frac{1}{2}}$
- If $\epsilon \rightarrow N(0, \sigma^2 I)$ then $Y \rightarrow N(X\beta, \sigma^2 I) | \hat{\beta} \rightarrow N(\beta, \sigma^2 (X^T X)^{-1})$

Theorem 1: Gauss Markov

$Y = X\beta + \epsilon$ $E[\epsilon] = 0$ $Cov(\epsilon) = \sigma^2 I$ $rank[X] = p$
Furthermore let $c \in \mathbb{R}^p$ and $\hat{\beta}$ the LS estimator. Then $c^T \hat{\beta}$ has minimal variance among all linear unbiased estimators of $c^T \beta \rightarrow$ BLUE: Best Linear Unbiased Estimator.
Assume: ϵ has $N(0, \sigma^2 I)$. Then $c^T \hat{\beta}$ has minimal variance among all unbiased estimator \rightarrow UMVU

1.7 Asymptotic Properties of LSE

$Y_n = X_n \beta + \epsilon_n$ $n :=$ nb of observations Assumptions:

- $(X_n^T X_n)^{-1} \rightarrow 0$ as $n \rightarrow \infty$
- $\max_j H_{jj} = \max_j x_j (X^T X)^{-1} x_j^T \rightarrow 0$ as $n \rightarrow \infty$

Under Assumptions:

- $\hat{\beta}_{nLS} \xrightarrow{P} \beta | \hat{\sigma}_n^2, \hat{\sigma}_{nML}^2 \xrightarrow{P} \sigma^2$ (Consistency)
- $(X_n^T X_n)^{\frac{1}{2}} (\hat{\beta} - \beta) \xrightarrow{D} N(0, \sigma^2 I)$

1.8 Statistical Properties of Residuals

- $E\hat{\epsilon} = 0$
- $cov(\hat{\epsilon}) = \sigma^2 (I - H)$
- $Var(\hat{\epsilon}_j) = \sigma^2 (1 - H_{jj})$ (Heteroscedasticity)
- $cov(\hat{\epsilon}_j, \hat{\epsilon}_i) \neq 0$ (Correlated)
- $\epsilon \rightarrow N(0, \sigma^2 Q)$
- Res.Sum.Square: $\frac{\epsilon^T \epsilon}{\sigma^2} \rightarrow \chi_{n-p}^2$
- $\epsilon^T \hat{\epsilon}$ Independent of $\hat{\beta}_{LS}$

2 Hypotheses Testing

2.1 Exact F test

Basics of testing: level α : $P(\phi(y, x, \beta_0) = 1) \leq \alpha$ Def P-value: Prob. to get a value as extrem as at least as extreme as result actually observed during the test, assuming that H_0 is correct

- General linear Hypotheses:** $H_0 : C\beta = d, rank(C) = r \leq k$
- Assumptions:** Gaussians Errors
- Pivot Statistic:** $F = \frac{n-p}{r} \frac{\Delta SSE}{SSE} \sim F_{r, n-p}$
- Computation of Pivot: $F = \frac{1}{r} (C\hat{\beta} - d)' (\hat{\sigma} C (X^T X)^{-1} C')^{-1} (C\hat{\beta} - d)$
- Reject H_0 :** $F > F_{r, n-p}(1 - \alpha)$
- $\Delta SSE = SSE_0 - SSE$ to compute $SSE_0 \rightarrow$ constrained Least square
- ΔSSE get smaller the less the error the more likely to accept H_0

2.2 F-test for Specific Problems

- Test of significance:** $H_0 : \hat{\beta}_j = 0$ (given all other predictors)
- Pivot: $t_j = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{Var(\hat{\beta}_j)}} = \frac{\hat{\beta}_j}{\hat{\sigma}} \|z_j\|^2 \sim t_{n-p}$
- Reject H_0 (significance): $|t_j| > t_{n-p}(1 - \frac{\alpha}{2})$
- Composite test:** $\beta_{sub} = (\beta_1, \dots, \beta_r)^T | H_0 : \beta_{sub} = 0 | H_1 : \exists (\beta_{sub})_i \neq 0$
- Pivot: $F = \frac{1}{r} \hat{\beta}_{sub}^T cov(\hat{\beta}_{sub}) \hat{\beta}_{sub} \sim F_{r, n-p}$
- Note: $cov(\hat{\beta}_{sub})$ are sub element of full model $\hat{\sigma}^2 (X^T X)^{-1}$
- Global F-test:** $H_0 : \beta_{full} = 0, \beta_{full} \in \mathbb{R}^k | H_1 \exists \beta_j \neq 0$
- Pivot: $F = \frac{n-p}{k} \frac{R^2}{1 - R^2}$
- Under $H_0 : \beta_0 = \bar{y}$

2.3 Confidence Regions

- CI $H_0 : \beta_j = d_j$** with Pivot $t_j = \frac{\hat{\beta}_j - \beta_j}{se_j}$
- $(1 - \alpha)CI = [\hat{\beta}_j + -se_j * t_{n-p}(1 - \frac{\alpha}{2})]$
- C El.** $\{\beta : \frac{1}{r} (\hat{\beta} - \beta)' cov(\hat{\beta}) (\hat{\beta} - \beta) \leq F_{r, n-p}\}$
- CI for $Ey_0 = \mu_0$:** $y_0 =$ future observation at location x_0
- $[x_0' \hat{\beta} + -t_{n-p}(1 - \frac{\alpha}{2}) \hat{\sigma} (x_0' (X^T X)^{-1} x_0)^{\frac{1}{2}}]$
- Prediction Interval:** we want an interval for future value y_0
- $[x_0' \hat{\beta} + -t_{n-p}(1 - \frac{\alpha}{2}) \hat{\sigma} (1 + x_0' (X^T X)^{-1} x_0)^{\frac{1}{2}}]$
- Note: Pred.interval bigger than Ci for mean

3 Multiple Testing

3.1 Problem Formulation

- Hypotheses:** $H_0 = H_0^{(1)} \cap \dots \cap H_0^{(m)}$
- Discovery:** reject H_0 (at least accept one alternative)
- False Discovery := $Type I - Error$
- P-Value:** $p := \mathbb{P}_{H_0}(T(X) > T(x)) \in [0, 1]$ for $\phi = \mathbb{1}_{(T(x) \notin [q_l, q_r])}$
- P-value $\sim Uni[0, 1]$
- Proced.:** given $\{p^i\}_{i \in \{1, \dots, m\}} \xrightarrow{Procedure} \{H_0, H_A\}$
- Level:** $\mathbb{P}_{H_0}(falsedisc.) = \mathbb{P}_{H_0}(p \leq \alpha) \leq \alpha$
- Naive approach: test separately $\mathbb{P}_{H_0}(\phi^i = 1; \exists i \in \{1, \dots, m\}) = 1 - (1 - \alpha)^m$ as m grows last eq. goes to 1 \Rightarrow get false discovery by chance.

3.2 Bonferroni Holmes

- Metrics** $FWER = \mathbb{P}(V \geq 1)$ with $V :=$ nb false discovery
- Procedure BF:** Reject $H_0^{(i)}$ if $p^{(i)} \leq \frac{\alpha}{m}$
- Thm:** Bonferroni controls FWER with level α : $\mathbb{P}(V \geq 1) \leq \alpha$
- Procedure BF Holmes:** given ordered $\{p^{(1)} \leq \dots \leq p^{(m)}\}$
- Step1:** Reject if $p^{(1)} \leq \frac{\alpha}{m}$ else accept $\{H_0^{(1)} \dots H_0^{(m)}\}$
- Step2:** Reject if $p^{(2)} \leq \frac{\alpha}{m+1}$ else accept $\{H_0^{(2)} \dots H_0^{(m)}\}$
- Stepi:** Reject if $p^{(i)} \leq \frac{\alpha}{m-i+1}$ else accept $\{H_0^{(i)} \dots H_0^{(m)}\}$
- Thm:** bf Holmes controls FWER at level α and is UMP for simple Bf procedure

3.3 Permutation Test

- **Idea:** non parametric test (not knowing about the dist.)
- (1) Sample $S_{xy} = \{x_1, ..., x_n, y_1, ..., y_m\}$
- (2) Randomly permute S_{xy} and split into \tilde{S}_x, \tilde{S}_y
- (3) Compute statistic $T(\tilde{S}_x, \tilde{S}_y)$
- repeat B times (1,2,3): $\{T^{(1)}, ..., T^{(B)}\}$
- $\hat{p}_{val} = \frac{\sum_{l=1}^B \mathbb{1}(|T^{(l)}| \geq |T(S_x, S_y)|) + 1}{B+1}$
- $T(S_x, S_y) := \text{observed statistic}$
- Decision: Reject if $\hat{p}_{val} \leq \alpha$

3.4 Benjamini Hochberg

- **Idea:** Controlling FWER is very restrictive, hence low power. Gain in power by expecting $(1 - FDR)$ of the discoveries to be true.
- $FDP = \frac{V}{|R|} = \frac{|H_0 \cap R|}{|R|}$
- $FDR := E[FDP]$
- **Procedure:** Given ordered pval $\{p^{(1)} \leq ... \leq p^{(m)}\}$
- $max_k(p^{(k)} \leq \frac{\alpha k}{m}) \rightarrow R := \{p^{(1)}, ..., p^{(k)}\}$
- **Thm:** if pval are independant then BH_α controls FDR at level α

4 Model Selection

4.1 Model Specification

- **Missing Variables:** $\hat{\beta}^M$ is biased | $Var(\hat{\beta}^M) \leq Var(\beta_{true})$
- **Irrelevant Varibales:** $\hat{\beta}^M$ is biased | $Var(\hat{\beta}^M) \geq Var(\beta_{true})$

4.2 Metrics

- **SMSE:** sum of mean squared errors \hat{Y}_i^M from true value μ_i
- $SMSE(M) = E[\|\hat{Y}^M - \mu\|^2] = |M|\sigma^2 + \sum_i^n (E[\hat{Y}_i^M] - \mu_i)^2 = Var + bias^2$
- $\Gamma_p^M = \frac{SMSE}{\sigma^2}$
- **SPSE:** sum of prediction errors; \hat{Y}_i^M pred. of new variable $Y_{n+i} = \mu_i + \epsilon_{n+i}$
- $SPSE(M) = \sum_i^n E[(Y_{n+i} - \hat{Y}_i^M)^2] = n\sigma^2 + SMSE(M)$
- **SSE:** sum of squared errors
- $SSE(M) = E[\|Y - \hat{Y}^M\|^2] = SPSE(M) - 2|M|\sigma^2$
- Note: for all metrics μ_i, σ^2 are unknwn \rightarrow need estimates
- Note: Minimizing one of the above metrics results in minimizing all of them

4.3 Model Choice Criteria

- **Overall Goal:** select covariates that minimize the expected SPSE
- **Approach 1:** Split data into $\{train, validation\}$. Model on train and $\widehat{SPSE} = E[\|Y_{val} - \widehat{Y}_{train}^M\|^2]$ with validation
- **Approach 2:** not enough data for appr.1 compute $\hat{\beta}^M$ with all data
- $\widehat{SPSE}(M) = \sum_i^n (y_i - \hat{Y}_i^M)^2 + 2|M|\hat{\sigma}^2$ ($\hat{\sigma}^2$ based on full model)
- **Mallow's CP:** $\hat{\Gamma}^M = C_p = \frac{\sum_i^n (y_i - \hat{Y}_i^M)^2}{\hat{\sigma}^2} - n + 2|M|$
- **AIC:** $AIC = -2l(\hat{\beta}^M, \hat{\sigma}_{ML}^2) + 2(|M| + 1)$ wish to minimize AIC
- **BIC:** $BIC = -2l(\hat{\beta}^M, \hat{\sigma}_{ML}^2) + \log(n)(|M| + 1)$
- Note: BIC stronger penalization for complex model
- Better fit if small
- **Adj. R coef** $\bar{R}^2 = 1 - \frac{n-1}{n-p}(1 - R^2)$
- **Cross Validation:** $CV_{score} = \sum_j^r \widehat{SPSE}_j^{-j}$ for data splitted in r chunk
- Special case Leave one out: $CV_{score} = \frac{1}{n} \sum_i^n (Y_i - \hat{Y}_i^{-i, M})^2 = \frac{1}{n} \sum_i^n \frac{Y_i - \hat{Y}_i^M}{1 - H_{ii}}$
- big computational cost avoided as one only need to compute H^M once for each M and not one for all CV chunks

4.4 Model Selection Procedures

- Forward selection
- Backward selection
- Stepwise selection
- Best subset selection

4.5 Inference after Model Selection

- After model sel. the data has already been used \implies Statistical Inference rules are broken
- Solution: Randomly split data into 2 parts, apply model selection on first and use second for statistical inference (like CI, PI, etc...)

5 Model Diagnostic

5.1 Model Assumptions

- **Zero mean Errors:** ($\hat{\beta}$ biased | test, CI not valid)
- Transformations, included omitted covariates
- **Homo. Errors:** test and CI
- Transformation of response, general linear model
- **Uncorrelated Errors:** test and CI
- General linear model
- **Normal Ass.** test and CI only asymp. correct

5.2 Residual Plots

- **Stand. Residuals:** $r_i = \frac{\hat{\epsilon}_i}{\hat{\sigma}\sqrt{1-H_{ii}}}$
- **Studentized Res:** $\tilde{r}_i = \frac{\hat{\epsilon}_{(i)}}{\hat{\sigma}^{(-i)}\sqrt{1-H_{ii}}} \sim t_{n-p-1}$
- design matrix with i row $X^{-i} \rightarrow \hat{\beta}^{(-i)} \rightarrow \hat{\sigma}^{(-i)}$ and $\hat{\epsilon}_{(i)} = y_i - x_i' \hat{\beta}^{(-i)}$
- **Tuckey Anscombe:** Plot $\hat{\epsilon}_i$ VS $\hat{y}_i \rightarrow$ check zero means. Even if residuals not uncorr still good visual approx.
- fluctuations \rightarrow non-linearity | omitted covariate
- **Stand Res VS Covariates:** if rdn spread around zero line model explained well effects of covariates on mean
- **Scale Location Plot :** $\sqrt{|r_i|}$ VS y_i to check Homosced. (if well spreaded around zero line)

5.3 Transformations

- **Goal:** Achieve mean function approx linear in the transformed scale
- **Power Fam:** $\Psi(u, \lambda) = \begin{cases} u^\lambda & \lambda \neq 0 \\ \log(u) & \lambda = 0 \end{cases}$
- $\lambda \in \{-1, 0, 1/3, 1/3, 1\}$
- **Scaled Power Fam:** (transforming only covariates) $\Psi(x, \lambda) = \begin{cases} x^\lambda & \lambda \neq 0 \\ \log(x) & \lambda = 0 \end{cases}$
- Choose λ visually, ie model expression best linear relationship (ex: distr. not scewed) AND minimizing the SSE (res sum square)
- **Box Cox Procedure**

5.4 Collinearity Analysis

- **Goal:** Check whether covariates have lin dep. (correlation)
- **Variance Inflation Factor:** $VIF_j := \frac{1}{1-R_j^2}$
- Qtf. $var(\hat{\beta}_j) | R_j :=$ R coeff of $X_j \sim X^{(-j)}$
- Threshold: (> 10)

5.5 Outliers

- **Stud. Res:** $\tilde{r}_i \sim t_{n-p-1}$ | Test all \tilde{r}_i simultaneously (BF correction)
- $\tilde{r}_i > q(\alpha/2)_{t_{n-p-1}} \rightarrow$ outlier
- **Leverage:** "Outlier in x-direction" $lev := H_{ii} \in [\frac{1}{n}, 1]$
- Influential data if : $var(\hat{\epsilon}_i) = \sigma^2(1 - H_{ii})$ if $H_{ii} \rightarrow 1$ (big) then Hyper plane passes through $(x_i, y_i) \rightarrow$ influential data.
- Rule of Thumb: $H_{ii} > \frac{2p}{n}$
- **Cook's Distance:** $D_i = \frac{\frac{1}{p\hat{\sigma}^2} \frac{\hat{\epsilon}^{(-i)2}}{1-H_{ii}} \frac{H_{ii}}{1-H_{ii}}}{p\hat{\sigma}^2} = \frac{\|\hat{y}^{(-i)} - \hat{y}\|^2}{p\hat{\sigma}^2}$
- Euclidean dist bwt Hyper plane Full \hat{y} and omitted i $\hat{y}^{(-i)}$
- Check eq. \rightarrow plot residuals vs leverage
- If $D_i > 0.5$

6 General Linear Model

6.1 Weighted Least Square

- **Gener. Model:** $Y = X\beta + \epsilon | E[\epsilon] = 0 | cov(\epsilon) = \sigma^2 W^{-1}$
- W assumed to be Pos. Definite.
- $W = P^T DP, \exists Bs.t W = BB^T$ and $W^{1/2} W^{1/2} = W$
- **OLS gener. Model:** $\hat{\beta} = (X^T X)^{-1} X^T Y$
- $E[\hat{\beta}] = 0$ but $cov(\hat{\beta}) = \sigma^2 (X'X)^{-1} X'W^{-1} X (X'X)^{-1}$
- (CI and Test no longer valid) and GaussMarkov ass. broken $cov(\epsilon) = \sigma^2 I$
- Hence OLS no longer UMVU
- **WLS :** $W^{1/2} Y = W^{1/2} X\beta + W^{1/2} \epsilon = \tilde{Y} = \tilde{X}\beta + \tilde{\epsilon}$
- $\hat{\beta}_{WLS} = (\tilde{X}'\tilde{X})^{-1} \tilde{X}'\tilde{Y} = (X'WX)^{-1} X'WY$
- $cov(\tilde{\epsilon}) = \sigma^2 I \rightarrow$ GaussMarkov $\hat{\beta}_{WLS}$ is UMVU
- REML: $\hat{\sigma}^2 = \frac{1}{n-p} \tilde{\epsilon}' W \tilde{\epsilon}$ and $\hat{\epsilon} = Y - X\hat{\beta}_{WLS}$
- if $W = diag(w_1, ..., w_n)$ then $WLS(\beta) = (Y - X\beta)' W (Y - X\beta) = \sum_i^n w_i (y_i - x_i' \beta)^2$
- $w_i \propto \frac{1}{var(\epsilon_i)}$ weights bigger impact on objective when var is small
- **Grouped Data (known Heter.):** $y = (\bar{y}_1, ..., \bar{y}_G)'$ hence $\epsilon = (\bar{\epsilon}_1, ..., \bar{\epsilon}_G)'$
- $cov(\bar{\epsilon}_i) = \frac{\sigma^2}{n_i}$
- $cov(\epsilon) = \sigma^2 W^{-1} \rightarrow W = diag\{\frac{1}{n_1}, ..., \frac{1}{n_G}\}$
- **Unknown Heter.** solution: transformations $\sqrt{y}, \log(y)$
- **Solution 2 Stage LS:**
- $var(\epsilon_i) = E[\epsilon_i^2] = \sigma_i^2 \xrightarrow{lin. model} \epsilon_i^2 = \sigma_i^2 + \nu_i = z_i' \alpha + \nu_i$
- Note: $z_i :=$ cov. that affect the errors (usually X)
- errors unknown $\rightarrow \hat{\epsilon}_i^2 \approx z_i' \alpha + \nu_i$
- $\hat{E}[\hat{\epsilon}_i^2] \approx \hat{E}[\epsilon_i^2] = z_i' \hat{\alpha}_{LS} \rightarrow \sigma_i^2 \approx z_i' \hat{\alpha}_{LS}$
- WLS: $w_i := \frac{1}{z_i' \hat{\alpha}_{LS}}$
- **White Esti.** estimation of $cov(\hat{\beta})$ used for correction of stand. errors CI and tests for asymptotics
- $\widehat{cov}(\hat{\beta}) = (X'X)^{-1} X' diag\{\hat{\epsilon}_1^2, ..., \hat{\epsilon}_n^2\} X (X'X)^{-1}$

7 Robust Regression

8 Generalized Linear Models

8.1 Binary Regression

- $y_i \in \{0, 1\}$ and $y_i \sim \text{Ber}(\pi)$
- $E[y_i] = P(y_i = 1) = \pi_i$ and $\text{var}[y_i] = \pi_i(1 - \pi_i)$
- Response Func:** $\pi_i = h(\eta_i) \in [0, 1]$
- Link Func:** $h^{-1}(\pi_i) = g(\pi_i) = \eta_i = x_i' \beta$
- Logit Model:** $\log\left(\frac{\pi_i}{1-\pi_i}\right) = x_i' \beta \Leftrightarrow \pi_i = \frac{\exp(\eta_i)}{1 + \exp(\eta_i)}$
- Probit Model:** $\Phi^{-1}(\pi) = x' \beta \Leftrightarrow \pi = \Phi(\eta)$
- Latent Variable:** $y_i = \mathbb{1}(z_i \geq 0)$ (observed)
- $z_i = x_i' \beta - \epsilon_i \rightarrow E[y_i] = P(y_i = 1) = P(z_i \geq 0) = P(\epsilon_i < x_i' \beta) = F(x_i' \beta)$
- if $\epsilon_i \sim N(0, \sigma) \rightarrow F = \Phi$ Probit model
- if $\epsilon_i \sim \text{logistic} \rightarrow F = \text{logit}$ model
- Note: Logistic variance $= \pi^2/3 \rightarrow \tilde{\beta} = \frac{\pi}{\sqrt{3}} \beta \rightarrow \pi(\eta) = h(x' \tilde{\beta}) = \Phi(x' \tilde{\beta})$
- resulting prob. are \approx equal for adjusted model
- Note: to compare $\hat{\beta}_{\text{logit}}$ vs $\hat{\beta}_{\text{probit}} \rightarrow$ scale: $\hat{\beta}_{\text{probit}} \frac{\pi}{\sqrt{3}}$
- Interpretation:** $\log \text{ ratio } \frac{P(y_i=1)}{P(y_i=0)} = \exp(\beta_0) \exp(x_1 \beta_1) \dots$ (ex: if risk covariate categorical then see the effect of risk present $\frac{P(y_i=1)}{P(y_i=0)} = \exp(1\hat{\beta}_{\text{risk}}) * \dots$ vs $= \exp(0\hat{\beta}_{\text{risk}}) * \dots$)
- Grouped Data:** $\{n_i | \bar{y}_i = \frac{1}{n_i} \sum_{i=1}^{n_i} y_i | x_i\}, i \in \{1, \dots, G\}$ data grouped for identical covariates in G groups
- $y_i \sim \text{Ber}(\pi_i) \rightarrow \sum y_i \sim \text{Bin}(n_i, \pi_i)$
- $P_{\pi_i}(n_i \bar{y}_i) = P_{\pi_i}(\sum_{j=1}^{n_i} y_j) = (n_i, \sum y_j) \pi_i^{\sum y_j} (1 - \pi_i)^{n_i - \sum y_j}$
- $E[\bar{y}_i] = \pi_i \rightarrow$ logit probit model
- $\text{var}(\bar{y}_i) = \frac{\pi(1-\pi_i)}{n_i}$
- Overdispersion:** unobserved Heterogeneity or positive correlation (ex: when people sample comes from same cluster)
- Empirical Var: $s = \frac{\bar{y}_i(1-\bar{y}_i)}{n_i} > \frac{\hat{\pi}_i(1-\hat{\pi}_i)}{n_i} \rightarrow$ Overdispersion
- Solution: adjust model with overdispersion parameter
- $\text{var}(\bar{y}_i) = \phi \frac{\pi(1-\pi_i)}{n_i}$

8.2 Max. Likelihood Estimation

- Score Func:** $s_{\beta}(y) = \frac{d}{d\beta} \log(p_{\beta}(y))$
- M Est:** $s_{\beta}(y)|_{\hat{\beta}} = 0$
- Observed Fisher Info:** $H(\beta) = \frac{d^2}{d^2\beta} l(\beta)$
- Fisher Information:** $F(\beta) = E[s(\beta)s(\beta)'] = E[-\frac{d}{d\beta} s(\beta)]$
- Z Estimator:** $\frac{d}{d\beta} \hat{R}(\beta) = \frac{1}{n} \sum_{i=1}^n s_{\beta}(y) = s_{\beta}(\bar{y})$
- Alg. Newton :** $\beta^{t+1} = \beta^t + H(\beta^t)^{-1} s_{\beta^t}(\bar{y})$
- Alg. Fisher Scoring:** $\beta^{t+1} = \beta^t + F(\beta^t)^{-1} s_{\beta^t}(\bar{y})$
- Note: $s_{\beta}(\bar{y}) = \sum_i^n s_{\beta}(y_i)$
- Note: $\mathbb{F}(\beta) = nF(\beta)$
- special case:** $y_i \sim \text{Ber}(\pi_i) \rightarrow s_{\beta}(\bar{y}) = X'(Y - \pi) \rightarrow \mathbb{F}(\beta) = \sum_i^n x_i x_i' \pi_i (1 - \pi_i) = X' V X$

8.3 Statistical Inference

- Asymp. Theory :** $\hat{\beta}_{MLE} \xrightarrow{D} N(\beta, F^{-1}(\beta))$
- $\widehat{\text{cov}}(\hat{\beta}) = F^{-1}(\hat{\beta})$
- $(F^{(\hat{\beta})})_{ii} = \widehat{\text{var}}(\hat{\beta}_i) = \text{se}_i^2$ standard error for $\hat{\beta}_i$
- Testing:** Likelyhood ratio test $lr = 2(l(\hat{\beta}) - l(\tilde{\beta})) \sim \chi_r^2$
- $H_0 : \beta_1 = 0 \rightarrow \hat{\beta}$
- $H_A : \beta_1 \neq 0 \rightarrow \tilde{\beta}$
- General Hyp:** Wald Statistic $w = (C\hat{\beta} - d)'(CF^{-1}(\hat{\beta})C')(C\hat{\beta} - d) \sim_a \chi_r^2$ Note $C = r \times p$ with $\text{rank}(C) = r < p$
- Sign test:** wald stat $w = z_j^2 = \frac{\hat{\beta}_j^2}{F^{-1}(\hat{\beta})_{jj}} \sim_a N(0, 1)$
- Model fit criteria :** Deviance $D(\hat{\pi}) = -2l(\hat{\pi})$ explicitly compare the fit with the perfect fit $l(\hat{\pi} = 0)$ Like kullback leiber information. The smaller the better the fit
- Deviance for grouped data: Theoretical MLE known: $\hat{\pi}_i = \bar{y}_i$ compare num to theoretical: $D = -2 \sum_i^G (l_i(\hat{\pi}_i) - l_i(\bar{y}_i))$
- Overdispersion estimation: $\hat{\phi} = \frac{1}{n-p} D$
- Note: deviance good for comparing nested models

8.4 Count Data Regression

- $y_i \in \{0, 1, 2, \dots\}$ and $f_{\beta}(y_i) = \frac{\lambda^{y_i}}{y_i!} \exp(-\lambda_i)$
- $E[y_i] = \lambda_i = \text{var}[y_i]$
- Log Linear model:** $\log(\lambda_i) = \eta_i = x_i' \beta$
- Overdispersion: variance higher in data than in model \rightarrow adjust model $\text{var}[y_i] = \phi \lambda_i$
- $\hat{\phi} = \frac{1}{n-p} D$

8.5 Unified Framework for GLM's

- Exp Family:** $f_{\theta}(y) = \exp[\frac{y\theta - b(\theta)}{\phi} w + C(\phi, w, y)]$
- $\theta :=$ natural parameter
- Canonical Link Function:** $f_{\gamma} = \exp[c(\gamma)T(y) - d(\gamma)]h(y) \rightarrow c(\gamma) := \theta$ c canonical link function
- Note: under natural paramter: $E[y] = \frac{d}{d\theta} b(\theta)$
- Summary:** natural paramter modelled as linear $c(\gamma) = \theta = \eta = x' \beta$

8.6 Classification Metrics

- Goal:** use predicted probabilities to classify (predict) new data: $y_{\text{new}} = \mathbb{1}(\hat{p}_i > t) \in \{0, 1\}$ with t being the threshold
- $\hat{\pi} = \frac{\exp(x_{\text{new}} \hat{\beta})}{1 + \exp(x_{\text{new}} \hat{\beta})}$
- if $t = 0$ no mistake on oservation =1
- Accuracy:** $\frac{TP+TN}{n}$
- Sensitivity:** $TPR = \frac{TP}{TP+FN}$ (Recall)
- Specificity:** $TNR = \frac{TN}{TN+FP}$
- $FPR = 1 - \text{spec} = \frac{FP}{TN+FP}$
- ROC curve:** TPR vs FPR goal: reach point (0,1) and think threshold goes $1 \rightarrow 0$ left to right

9 Penalized Regression

- Context:** If model is well specified the OLS estimator is unbiased but may have high variance \rightarrow think subset of coeff that are close to zero bring the same variance as other coeff into the model. Overall idea is to trade some bias to reduce variance
- Case: if there is strong collinearity between covariates $\rightarrow X'X$ instable
- Case: High dimensional Regression $p > n$
- Goal:** reduce SPSE by trading bias for some variance
- Penalized Regr. :** $PLS(\beta) = \|Y - X\beta\|^2 + \lambda \text{pen}(\beta)$
- Ridge:** $\text{pen}(\beta) = \|\beta\|^2 \rightarrow \hat{\beta}_{\text{ridge}} = (X'X + I)^{-1} X'Y$
- SVD:** $X = UDV'$ with $\text{col}(U) = \text{Span}(X)$ and $U'U = I$ with $U \in R^{n \times p}$
- $X\hat{\beta}_{\text{ridge}} = \sum_{j=1}^p u_j \frac{d_j}{d_j + \lambda} u_j' Y$
- Interpretation:** Ridge estimator shrinks the smaller principal component which correspond to small sample variances.
- $\text{cov}[\hat{\beta}_{\text{ridge}}] \leq \text{cov}[\hat{\beta}_{OLS}]$ BUT biased
- Ridge need covariates and response to be scaled and centered (centered because no penalty on intercept and scale because constrain on β so β_j must be of the same scale)
- Amount of shrinkage controlled by λ
- choosing $\lambda \rightarrow$ CV minimizing the SPSE and choose the simplest model
- LASSO:** $\text{pen}(\beta) = |\beta|_1 \rightarrow \|Y - X\beta\|^2 + \lambda \sum_i^p |\beta_i|$
- Note: no close form
- Graphical Overview :** $LS(\beta) = (\beta - \hat{\beta})' X'X (\beta - \hat{\beta}) + \hat{\epsilon}' \hat{\epsilon} \rightarrow$ ELLIPSOID (plot contour form and constraint on β)
- LASSO shrink coefficient to zero (more shrinkage for small coeff but less for big coeff)
- Comparison of shrinkage:** Assume $X'X = I$ orhtogonal design and $\hat{\beta}_j$ is OLS estimate
- Best subset selection: $\hat{\beta}_j \mathbb{1}(|\hat{\beta}_j| > \sqrt{\lambda})$
- Ridge: $\hat{\beta}_j / (1 + \lambda)$
- LASSO: $\text{sgn}(\hat{\beta}_j) \max\{|\hat{\beta}_j| - \lambda/2\}$
- Note on LASSO:** it performs soft thresholding and is like automatic variable selection as the coefficients are set to zero.

10 Maths Tricks Tutorials

Sub Differentials :

- Convex Subset :** A subset $X \subseteq \mathbf{R}^k$ is called a convex subset iff $\forall x_1, x_2 \in X$, the segment that link them is contained in X
- Convex Function:** Let X be a convex subset and $f : X \rightarrow \mathbf{R}$ a function. Then f is convex iff $\forall x_1, x_2 \in X$ and $\forall t \in [0, 1]$

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$

- Subdifferentials :** Let $f : I \rightarrow \mathbf{R}$ a real valued convex function on the open interval I . The sub-differential of f at $x_0 \in I$ is the set $[a, b]$ where

$$a = \lim_{x \rightarrow x_0^-} \frac{f(x) - f(x_0)}{x - x_0}$$

$$b = \lim_{x \rightarrow x_0^+} \frac{f(x) - f(x_0)}{x - x_0}$$

- If f differential at x_0 then $a = b \implies [a, b] = \{f'(x_0)\}$
- Proposition:** $x_0 \in I$ is a global minimum of f iff zero is contained in the subdifferential at x_0