

sample Pearson: $\gamma = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}}$ Two Ordinal Variables: (not int in excel, but direction)

estimation: $V = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$ y low (low), high (stronger tendency)

Measures of variability: standard deviation & range: $x_{(n)} - x_{(1)}$ & interquartile range (IQR): $Q_3 - Q_1$, first, second, third quartiles.

Measures of shape: positive/rights skewed, negative/left skewed
skewness = $\frac{E[(X - \bar{X})^3]}{\sigma^3}$ (sample) = $\frac{\sum_{i=1}^n (X_i - \bar{X})^3}{(n-1)(n-2) \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{S^2}}$
kurtosis = $\frac{E[(X - \bar{X})^4]}{\sigma^4} - 3$, measures the heaviness of tails, compared to a normal distribution.

sample = $\frac{n(n+1)}{(n-1)(n-2)(n-3)} \sum_{i=1}^n \frac{(X_i - \bar{X})^4}{S^4} - \frac{3(n-1)^2}{(n-2)(n-3)}$.

3 Basic Hypothesis Testing

3.1 Basic concepts

Type I error(α): reject H_0 when H_0 is actually the truth. Type II error(β): fail to reject H_0 when H_1 is the truth.
 $\alpha = \Pr(\mathbf{x} \in RR|H_0)$ is true, $\beta = \Pr(\mathbf{x} \notin RR|H_1)$ is true).
The p-value is the probability of test result at least as extreme as the result actually observed during the test, assuming H_0 is true.
Due to the randomness of the observed data, p-value is r.v., which $\sim U[0, 1]$
Statistic Power: the prob of rejecting H_0 when H_1 is true, power = $1 - \beta$
Sampling dist: dist of the point estimate based on samples of a fixed size from a population.

Interpretation of CI: having numerous sample datasets and the 95% CI is computed for each sample dataset, then the fraction of computed CI that encompass the true parameter would tend toward 95%.

3.2 Hypothesis Testing for Categorical Variables

One-sample z test: $\hat{p} = \frac{\sum_{i=1}^n X_i}{n}$, SE of estimate $SE(\hat{p}) = \sqrt{p(1-p)/n}$
test statistic: $Z = \frac{\hat{p} - p_0}{\sqrt{p_0(1-p_0)/n}} \sim_{asy} N(0, 1)$.

Wald interval: $\hat{p} \pm z_{\alpha/2} \times \sqrt{\hat{p}(1-\hat{p})/n}$
Need for sample size: $CI(n) \hat{p} \geq 10 \& n(1-\hat{p}) \geq 10$ (test on p) $np_0 \geq 10 \& n(1-p_0) \geq 10$
There are other types of intervals, e.g., the Wilson (or score) interval, the Clopper-Pearson (or exact) interval, etc

Two-sample z test: $H_0: p_1 = p_2 = 0$, we have $\hat{p}_1 = \frac{\sum X_{1i}}{n_1}$, $\hat{p}_2 = \frac{\sum X_{2i}}{n_2}$
 $SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$, under H_0 ,

$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{p(1-p)(\frac{1}{n_1} + \frac{1}{n_2})}$
test statistic: $Z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{p(1-p)(\frac{1}{n_1} + \frac{1}{n_2})}} \sim_{asy} N(0, 1)$ where p is estimated by

pooled proportion $\hat{p} = \frac{n_1 \hat{p}_1 + n_2 \hat{p}_2}{n_1 + n_2}$
Need for sample size: $CI(n) \hat{p}_1 \geq 10 \& n_1(1-\hat{p}_1) \geq 10$ (test on p) $n_1 \hat{p}_1 \geq 10 \& n_1(1-\hat{p}_1) \geq 10$ for \hat{p}_1

Test for contingency table: Pearson's chi-square test can be used to assess: Goodness of fit/Homogeneity&Independence
For Got: $H_0: p_1 = p_{01}, \dots, p_m = p_{0m}$. Pearson's chi-square test statistic is $\chi^2 = \sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i} \sim_{asy} \chi^2_{m-1}$, O_i is the obs freq of the i th category, $E_i = np_{0i}$ is the expected freq under H_0 .

For homog&indep: groups=r, category=c, $\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$
 $\sim_{asy} \chi^2_{(r-1)(c-1)}$ under hom or indep ass, where E_{ij} is the expected freq of cell (i, j) assuming indep: $E_{ij} = np_{i.}p_{.j}$, when $r=c=2$, χ^2 equiv two sample z test for binary variables.
Ass for appay this test: (1) obs X_{1i} and X_{2i} are indep samples (2) sample size is enough (cell counts greater or equal to 10)

When sample is small, apply exact tests to compute p-values: 1. Fisher's exact test: with the same margins (same row and col sums) 2. Barnard's exact test: only the row margins are fixed, more powerful than the Fisher's McNemar's Test for paired samples: $H_0: p_1 = p_2$

3.3 Hypothesis testing for Continuous Variable

One-sample t test: (for normal population) $H_0: \mu = \mu_0$, test statistic

$$T = \frac{\sqrt{n}(\bar{X} - \mu_0)}{S} \sim t_{n-1} \text{ (exact)}$$

The larger the d.f., the more closely the dist approximates $N(0, 1)$

By CLT, T asy follows $N(0, 1)$, under H_0 , the t-test provides an exact test.
Two-sample t test: (for 2 indep normal populations) $H_0: \mu_1 = \mu_2$, if assuming that $\sigma_1 = \sigma_2 = \sigma$, then pooled sample standard deviation $S_p = \sqrt{\frac{1}{n_1 + n_2 - 2} (\sum_{i=1}^{n_1} (X_{1i} - \bar{X}_1)^2 + \sum_{i=1}^{n_2} (X_{2i} - \bar{X}_2)^2)}$ =
 $\sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}}$, test statistic: $T = \frac{\bar{X}_1 - \bar{X}_2}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1 + n_2 - 2}$ (exact)

F test: $H_0: \sigma_1^2 = \sigma_2^2$, test statistic: $F = \frac{S_1^2}{S_2^2} \sim F(n_1 - 1, n_2 - 1)$
When the variances are not equal, $SE(\bar{X}_1 - \bar{X}_2)$ is better estimated by $\sqrt{S_1^2/n_1 + S_2^2/n_2}$, thus the test statistic is $T_s = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{S_1^2/n_1 + S_2^2/n_2}} \sim_{asy} t_{v_t}$, where $v_t = \frac{(S_1^2/n_1 + S_2^2/n_2)^2}{\frac{(S_1^4/n_1^2)}{n_1 - 2} + \frac{(S_2^4/n_2^2)}{n_2 - 2}}$ is called Satterthwaite/Welch's t-test.

Nonparametric test for Means/Median: **Sign test:** $H_0: m = m_0$, let N^+ be the number of positive signs obtained upon calculating $X_i - m_0$ for $i = 1, \dots, n$, under H_0 , $N^+ \sim \text{Bin}(n, p)$ with $p = 0.5$, take one-sample z-test. (rout)

Wilcoxon signed-rank test: compute $\{X_i - m_0\}_{i=1}^n \rightarrow$ order $\{|X_i - m_0|\}_{i=1}^n$ and assign ranks \rightarrow sums of ranks (positive) = S^+
When sample size is large (> 20), by CLT, under H_0 , we have $W = S^+ \sim_{asy} N(\frac{n(n+1)}{4}, \frac{n(n+1)(n+2)}{24})$

3.4 Multiple Comparisons

Type I error/alpha inflation. To control the family-wise error rate (FWER), i.e., the prob of incorrectly rejecting at least one H_0 .

Bonferroni: $\hat{p}_i = \min\{m \times p_i, 1\}$ or $\hat{\alpha}_i = \alpha/m$ is conservative when m is large or the tests are highly positively correlated.

Holm adjustment: step 1: if $p_{(1)} \leq \alpha/m$, reject $H_{(1)0}$ and continue, else stop. \dots step m : if $p_{(m)} \leq \alpha/m$, reject $H_{(m)0}$, with larger threshold (more powerful), $\hat{\alpha}_i = \alpha/(m - i + 1)$ & adjusted p-value: $\hat{p}_{(i)} = \{1, \max\{(m - i + 1)p_{(i)}, \hat{p}_{(i-1)}\}\}$

4 Linear Regression (Model) Model Fitting

4.1 The Multiple Linear Regression

Regression analysis: describe the mean of the distribution of one variable (response) as a function of other variables (explanatory): $E(Y|X) = f(X)$.
Regression Model: $\mathbf{Y} = \mathbf{X}\beta + \epsilon$; \mathbf{X} : design matrix, β : vector of parameter.
Least Square Method: assumptions: 1. All explanatory variables X_i are fixed 2. random errors are uncorrelated with $E(\epsilon) = 0$ & $\text{Var}(\epsilon) = \sigma^2$.

Two results: $\frac{\partial \hat{\beta}}{\partial \mathbf{a}} = \frac{\partial \hat{\beta}}{\partial \mathbf{a}} \frac{\partial \mathbf{a}}{\partial \mathbf{a}} = \mathbf{A}$ and $\frac{\partial (\hat{\beta}^T \mathbf{A} \hat{\beta})}{\partial \hat{\beta}} = (\mathbf{A} + \mathbf{A}^T) \hat{\beta}$
 $SSE(\hat{\beta}) = \|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2 \rightarrow \hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ - fitted value (orthogonal projection): $\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H} \mathbf{y}$, \mathbf{H} : hat matrix (projection matrix)
 $E(\hat{\beta}) = \beta$, $\text{Var}(\hat{\beta}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$, $\hat{\sigma}^2 = \frac{RSS}{n-p-1}$ is a unbiased estimator.

Maximum Likelihood Estimation: Assumptions: 1. All explanatory variables X_i are fixed 2. random errors are i.i.d. $N(0, \sigma^2)$
Under ass, have important results: 1. $\hat{\beta} \sim N(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$ 2.

$\frac{n \hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-p-1}$ 3. $\hat{\beta}$ and $\hat{\sigma}^2$ are independent. need has two cons
If \mathbf{A}, \mathbf{B} (scalars matrices) and $\hat{\sigma} \sim N(\mu, \Sigma)$, then 1. $\mathbf{A} \mathbf{y} \sim N(\mathbf{A} \mu, \mathbf{A} \Sigma \mathbf{A}^T)$ 2. $\mathbf{A} \mathbf{y}$ and $\mathbf{B} \mathbf{y}$ are indep iff $\mathbf{A} \Sigma \mathbf{B}^T = 0$

4.2 Testing the Regression Coefficients

Single: $H_0: \beta_i = 0$. Denote $(i+1)$ th diagonal element of $(\mathbf{X}^T \mathbf{X})^{-1}$ as c_{ii} , as $\hat{\beta} \sim N(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$, have $\hat{\beta}_i \sim N(\beta_i, c_{ii} \sigma^2)$, test statistic: $T = \frac{\hat{\beta}_i - 0}{\sqrt{c_{ii} \hat{\sigma}^2}} \sim t_{n-p-1}$ under H_0

Sever: $H_0: \beta_k + 1 = \dots = \beta_p = 0$. Def two models, full model: $\mathbf{y} = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \dots + \beta_p X_p + \epsilon$ and reduced/restricted model ($k < p$): $\mathbf{y} = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \epsilon \rightarrow RSS_R \geq RSS_F$, test statistic: $F = \frac{(RSS_R - RSS_F)/(p-k)}{RSS_F/(n-p-1)} \sim F(p-k, n-p-1)$ under H_0

Overall significance, $H_0: \beta_1 = \dots = \beta_p = 0 \rightarrow$ ANOVA table SST: total sum of squares, SSM: explained sum of squares of the model, SSE: residual sum of squares

R Squared/Coefficient of determination: $R^2 = SSM/SST$, represents

the proportion of variance in the response variable that is explained by the explanatory variables, the remaining can be attributed to unknown variables or inherent variability.

Interactive effects: two exp vars are said to interact if the effect that one of them has on the mean response depends on the value of the other.

Gauss-Markov Theorem: in linear regression model, if $\epsilon_1, \dots, \epsilon_n$ satisfy: $E(\epsilon_i) = 0$ and $\text{Var}(\epsilon_i) = \sigma^2 < \infty$; $\text{Cov}(\epsilon_i, \epsilon_j) = 0$, $\forall i \neq j$, then $\hat{\beta}_{LS}$ has the lowest sampling variance within the class of linear unbiased estimators, termed the BLUE.

5 Linear Regression: Model Selection and Diagnosis

5.1 Model Selection

With less variables: overfitting, simplicity/Interpretation
Model-fitting criterion: $R^2_{adj} = 1 - \frac{1}{n} \sum_{i=1}^n SSE_k = 1 - (1 - R^2) \frac{n-1}{n-k-1}$, largest R^2_{adj} is equiv to choose model smallest MSE
Mallows's C_p : $C_p = \frac{SSE_k}{SSE_{p-1}} - (n - 2k - 2)$, under full model, $k=p$, $C_p = k + 1 = p + 1$, can be proven that $E(C_p) = k + 1$, choose the model with C_p closest to $k+1$ and k is small.
AIC = $-2 \log(L) + 2(k+1)$, **BIC** = $-2 \log(L) + \log(n(k+1))$, when $n \geq 8$, BIC imposes heavier penalty on k than AIC.

$l(\beta, \sigma^2) = -\frac{n}{2} \log \sigma^2 - \frac{(\mathbf{y} - \mathbf{X} \beta)^T (\mathbf{y} - \mathbf{X} \beta)}{2\sigma^2} + c$, and $\hat{\sigma}^2 = \frac{(\mathbf{y} - \mathbf{X} \beta)^T (\mathbf{y} - \mathbf{X} \beta)}{n}$

AIC = $n \log(\frac{SSE_k}{n}) + 2(k+1) + c$, **BIC** = $n \log(\frac{SSE_k}{n}) + (k+1) \log n + c$
Sequential Selection: Begin with the current model, sequentially add and/or drop one explanatory variable at a time based on whether the resulting model is superior.

forward selection/ backward elimination/ stepwise selection

Shrinkage method:

Ridge, minimize $SSE(\beta, \lambda) = \sum_{i=1}^n (y_i - \beta_0 - \dots - \beta_p x_{ip})^2 + \lambda \sum_{j=1}^p \beta_j^2$ is equiv $SSE(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \dots - \beta_p x_{ip})^2$ subject $\lambda \sum_{j=1}^p \beta_j^2 \leq t$
 $\hat{\beta}^{ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$, addition of $\lambda \mathbf{I}_p$ makes nonsingular.
Lasso: Because of the nature of the constraint, making sufficiently small (especially large) will cause some of the coefficients to be exactly zero.
Cross validation: to test the model's ability to predict new data to obtain an insight on how the model will generalize to an unknown dataset.
Leave-one-out CV & k-fold CV: Shuffle the data randomly and split the data into groups of approximately equal size.

5.2 Model Diagnosis

Linearity & Homoscedasticity & Independence & Normality, lin_i^2 hom, nor Residual plots: $y_i = \epsilon_i$, first check the linear and hom by Fitted(X) versus Residual Plot(Y), i.e., scatterplots of the residuals against the fitted. Noraml Quantile-Q plot: check the normality ass, we expect that the points in the Q-Q plot will closely lie on a straight line, or histogram.

Hypothesis: **Hom:** Breusch-Pagan test and White test

BP: auxiliary regression model: $r_i^2 = \gamma_0 + \gamma_1 \epsilon_i + \dots + \gamma_k z_{ki} + \epsilon_i$, $H_0: \gamma_1 = \dots = \gamma_k = 0$, using F-statistic. **White:** All explain vars, all square vars, all inter terms are included. another form: $r_i^2 = \gamma_0 + \gamma_1 \epsilon_i + \gamma_2 \epsilon_i^2 + \epsilon_i$
Normality: Shapiro-Wilk test and Kolmogorov-Smirnov test

SW: test statistic: $W = \frac{(\sum_{i=1}^n a_i r_i)}{\sqrt{\sum_{i=1}^n (r_i - \bar{r})^2}}$, $0 \leq W \leq 1$ and small values of W lead to rejection of normality.

Dist of W under norm has no closed form, only applied when $n \leq 2000$.
KS: based on empirical (edf), $F_n(x) = \frac{1}{n} \sum_{i=1}^n I(r_i \leq x)$, test statistic: $D = \sup_x |F_n(x) - F(x)|$, F is normal cdf, $D \sim_{asy}$ Kolmogorov dist under normality, K-S test requires a relatively large ($n > 2k$) to take proper cdf
Independence: **Durbin-Watson Test**, we can judge whether it is reasonable to assume independence based on the nature of how the data were collected.

for time series data, test statistic: $DW = \frac{\sum_{i=1}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2}$, detect the first order auto-corr (as $e_i = \rho e_{i-1} + u_i$), $H_0: \rho = 0$, $1 \leq DW \leq 4$, and $DW=2$ indicates no auto-corr, $DW < 1$ means strong positive auto-corr, $DW > 3$ means strong negative auto-corr.

5.3 Unusual Observation

6 Analysis of Variance

Source	DF	Sum of Squares	Mean Squares	F Value	P > F
Model	P	SSM	MSM = $\frac{SSM}{P}$	F = $\frac{MSM}{MSE}$	
Error	n - p - 1	SSE	MSE = $\frac{SSE}{n - p - 1}$		
Total	n - 1	SST			

[illegible]

