

STAT 231 - S25
Statistics
3-page Cheat Sheet

With Prof James Huang

This is my final exam cheat sheet. They don't cover everything, but hopefully they are still helpful.

[Josiah Plett](#)

Unit: Individual measurable item.

Population: Static collection of units.

Process: Produced collection of units.

Variate: Any varying value.

↳ continuous, discrete, categorical, ordinal

Attribute: Function of a variate.

Skewness:  < 0  = 0  > 0

Kurtosis:  1.8  3  5

IQR: inter-quartile range $q(0.75) - q(0.25) = \text{IQR}$

A		B	
y_{11}	y_{12}	y_{21}	y_{22}
Relative Risk of A in B (vs B):		$RR = \frac{y_{11}/(y_{11}+y_{12})}{y_{21}/(y_{21}+y_{22})}$	

Summary: $\{y_{(1)}, q(0.25), q(0.5), q(0.75), y_{(n)}\}$

Sample Variance:

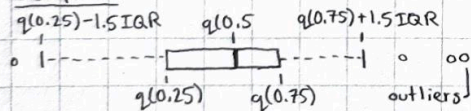
$$s^2 = \frac{1}{(n-1)} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{(n-1)} \sum_{i=1}^n (y_i^2 - n\bar{y}^2)$$

\uparrow
Covariance: $\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$

Correlation: $\text{corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$

Choose: $\binom{n}{x} = \frac{n!}{x!(n-x)!}$

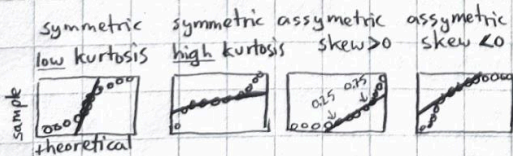
Boxplot:



Sample Correlation: $r = \frac{S_{XY}}{\sqrt{S_{XX} S_{YY}}}$

$$S_{XY} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

QQPlot: Gaussian \rightarrow linear plot



Central Limit Theorem (CLT):

$$Z_n = \frac{\sqrt{n}(\bar{Y} - \mu)}{\sigma} \quad Z_n \sim G(0, 1) \quad n \rightarrow \infty$$

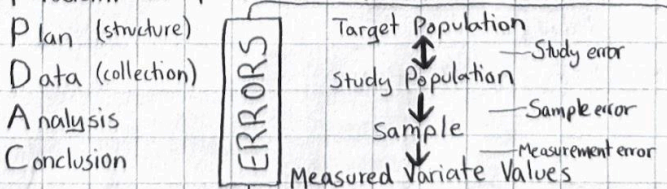
(Y_1, \dots, Y_n IID)

Maximum Likelihood Estimate (MLE):

$$L(\theta) = \prod_{i=1}^n f_Y(y_i; \theta) \quad \hat{\theta} = \underset{\theta \in \Theta}{\text{argmax}} L(\theta) \quad R = \frac{L(\theta)}{L(\hat{\theta})}$$

Model	Params, PDF, E[x], Var[x]	MLE	$R(\theta) = \frac{L(\hat{\theta})}{L(\theta)}$	(G) Pivotal Quantity
Poisson	λ $\frac{\lambda^x e^{-\lambda}}{x!}$ $E[x] = \lambda$ $\text{Var}[x] = \lambda$	\bar{y}	$\left(\frac{\theta}{\hat{\theta}}\right)^{n\hat{\theta}} e^{n(\hat{\theta}-\theta)}$ $\theta > 0$	$\frac{\bar{y} - \theta}{\sqrt{\theta/n}}$
Binomial	n, p $\binom{n}{x} p^x (1-p)^{n-x}$ $E[x] = np$ $\text{Var}[x] = np(1-p)$	$\frac{y}{n}$	$\left(\frac{\theta}{\hat{\theta}}\right)^y \left(\frac{1-\theta}{1-\hat{\theta}}\right)^{n-y}$ $0 < \theta < 1$	$\frac{\bar{y} - \theta}{\sqrt{\theta(1-\theta)/n}}$
Geometric	p $(1-p)^{x-1} p$ $E[x] = 1/p$ $\text{Var}[x] = 1/p^2$	$\frac{1}{1+\bar{y}}$	$\left(\frac{\theta}{\hat{\theta}}\right)^n \left(\frac{1-\theta}{1-\hat{\theta}}\right)^{\bar{y}}$ $0 < \theta < 1$	$\frac{\bar{y} - \theta}{\hat{\theta}/\sqrt{n}}$
Exponential	λ $\lambda e^{-\lambda x}, x \geq 0$ $E[x] = 1/\lambda$ $\text{Var}[x] = 1/\lambda^2$	\bar{y}	$\left(\frac{\hat{\theta}}{\theta}\right)^n e^{n(\hat{\theta}-\theta)}$ $\theta > 0$	$\frac{\bar{y} - \theta}{\hat{\theta}/\sqrt{n}}$
Gaussian (plain)	μ, σ $\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ $E[x] = \mu$ $\text{Var}[x] = \sigma^2$	$\hat{\mu} = \bar{y}$ $\hat{\sigma} = \sqrt{s^2}$	$e^{-\frac{n}{2\sigma^2}(\mu - \bar{x})^2}$	$\frac{\bar{y} - \mu}{\sigma/\sqrt{n}} \Rightarrow \bar{x} - \mu$
Simple Linear Regression	$Y_i \sim G(\alpha + \beta x_i, \sigma)$ $i=1, 2, \dots, n$ independent	Least Squares Estimates (MLE): $\hat{\beta} = \frac{S_{xy}}{S_{xx}} = r \cdot \frac{s_y}{s_x}$ $\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$ $\hat{\sigma}^2 = \frac{1}{n} (S_{yy} - \hat{\beta} S_{xy})$		
	Unknown	Pivotal Quantity	100p% Confidence Interval	
$G(\mu, \sigma)$ σ known	μ	$\frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \sim G(0, 1)$	$\bar{y} \pm a\sigma/\sqrt{n}$ $a = q_{\text{norm}}(\frac{1+p}{2})$	
$G(\mu, \sigma)$ σ unknown	μ	$\frac{\bar{Y} - \mu}{S/\sqrt{n}} \sim t(n-1)$	$\bar{y} \pm bs/\sqrt{n}$ $b = qt(\frac{1+p}{2}, n-1)$	
$G(\mu, \sigma)$ μ, σ unknown	Y	$\frac{Y - \bar{Y}}{S/\sqrt{1+1/n}} \sim t(n-1)$	$\bar{y} \pm bs\sqrt{1+1/n}$ $b = \uparrow$	
$G(\mu, \sigma)$ μ unknown	$\sigma^{(2)}$	$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$	$\left[\sqrt{\frac{(n-1)s^2}{d}}, \sqrt{\frac{(n-1)s^2}{c}} \right]$ c, d below	
a : $P(Z \leq a) = \frac{1+p}{2}$		b : $P(T \leq b) = \frac{1+p}{2}$, $T \sim t(n-1)$		
c, d : $P(W \leq c) = \frac{1-p}{2} = P(W > d)$, $W \sim \chi^2(n-1)$		$c = q_{\chi^2}(\frac{1-p}{2}, n-1)$ $d = q_{\chi^2}(\frac{1+p}{2}, n-1)$		
Exponential(θ)	θ	$\frac{2n\bar{y}}{\theta} \sim \chi^2(2n)$	$\left[\frac{2n\bar{y}}{d_1}, \frac{2n\bar{y}}{c_1} \right]$ $c_1, d_1: n-1 \rightarrow 2n$	

Problem: Descriptive, Causative, Predictive Survey: Sample, observational, experimental



(approximate)

p-value based on Gaussian distribution:

$$2P(Z \geq |\text{pivotal quantity}|), \quad Z \sim G(0, 1)$$

$$\downarrow$$

$$2 \cdot (1 - \text{pnorm}(d))$$

use $\hat{\theta} \pm Z \cdot \frac{1}{\text{pivotal quantity without numerator}}$

$$SSE = \sum (y_i - \hat{y}_i)^2$$

$$S_e = \sqrt{\frac{SSE}{n-2}}$$

$$\hat{y}_i = y_i - \hat{\alpha} - \hat{\beta}x_i$$

$$\hat{y}_i^* = \hat{y}_i / S_e$$

Binomial: $H_0: \theta = \theta_0$
 $D = |y - n\theta_0|$

Poisson: $H_0: \theta = \theta_0 \rightarrow$ average per interval
 $D = |y - t\theta_0|$

Exponential: $H_0: \theta = \theta_0 \rightarrow$ average wait time
 $D = \frac{2n\bar{y}}{\theta_0}$

Binomial: Calculate sample Size

$$n \geq \left(\frac{a}{m}\right)^2 \hat{\theta}(1-\hat{\theta})$$

$a = Z$: value from p confidence level.

$m = (U(y) - L(y))/2$: $\frac{1}{2}$ confidence interval.

Likelihood Interval \leftrightarrow Confidence Interval

100p% L.I. is defined as $\{\theta: R(\theta) \geq p\}$

100p% L.I. is an approx 100q% C.I. where:

$$q = P(W \leq -\ln p), W \sim \chi^2$$

Gaussian Confidence Intervals

$$\frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \sim G(0,1) \quad \frac{\bar{Y} - \mu}{S/\sqrt{n}} \sim t(n-1) \quad \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$$

Equal-Tailed 100p% C.I.: $P(-a \leq \text{pivot quantity} \leq a) = p \rightarrow [\bar{y} \pm a \overbrace{S/\sqrt{n}}^{\text{denominator}}]$

One-sided 100p% C.I.: $P(-a \leq p.q.) = p$ or $P(p.q. \leq a) = p$
upper lower

χ^2 Distribution

$$f_k(x) = \frac{x^{k/2-1} e^{-x/2}}{2^{k/2} \Gamma(k/2)}$$

↑
degrees of freedom

$$F_k \sim G(k, 2k)$$

$$\Gamma(n) = (n-1)! \text{ for } n \in \mathbb{N}$$

$$\Gamma(n) = (n-1)\Gamma(n-1) \text{ for } n > 1$$

$$\Gamma(1/2) = \sqrt{\pi}$$

$$Y_1, \dots, Y_n \sim G(0,1) \Rightarrow \sum_{i=1}^n Y_i^2 \sim \chi^2(n)$$

$$F_2 \sim \text{Exp}(2)$$

$$W_1, \dots, W_n \sim \chi^2(k_i) \Rightarrow \sum W_i \sim \chi^2(\sum k_i)$$

Student t Distribution

Let $Z \sim G(0,1)$, $Y \sim \chi^2(k)$ be independent: $T = \frac{Z}{\sqrt{Y/k}} \sim t(k)$

Simple Linear Regression Intervals

$$\hat{\beta} \pm a \text{Se}/\sqrt{S_{xx}}$$

$$[a]: P(T \leq a) = \frac{1+p}{2}$$

where $T \sim t(n-2)$

$$\hat{\alpha} \pm a \text{Se} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}$$

$$\hat{\mu} \pm a \text{Se} \sqrt{\frac{1}{n} + \frac{(x-\bar{x})^2}{S_{xx}}}$$

Constant variance, Independence, Normality, Linearity

Residual Plots



standardized residual is uncorrelated with x

Simple Linear Regression Question

Given: $Y_i = \alpha + \beta x_i + R_i$, $R_i \sim G(0, \sigma)$, $i=1, \dots, n$ independently

a) Show $E[\hat{\alpha}] = \alpha$ (note $E[\hat{\beta}] = \beta$)

$$E[\hat{\alpha}] = E[\bar{Y} - \hat{\beta}\bar{x}] \text{ (definition of } \hat{\alpha})$$

$$= \frac{1}{n} \sum_{i=1}^n (\alpha + \beta x_i) - \bar{x}\beta = \alpha$$

Prediction Interval > Confidence Interval

d) Give 95% C.I. for β

$$\hat{\beta} \pm c \text{SE}(\hat{\beta}) \text{ where } c = qt(0.975, n-2)$$

Two-variate Analysis

$$H_0: P(x \cap y) = P(x) \cdot P(y)$$

① Table is averages of real totals:

	high	low		high	low		
tall	20	5	25	tall	17.5	7.5	25
short	15	10	25	short	17.5	7.5	25
	35	15	same	35	15		

$$\lambda = 2 \sum_{j=1}^4 y_j \ln\left(\frac{y_j}{e_j}\right) \quad y_j = \text{observed} \quad e_j = \text{expected}$$

$$\text{③ } p\text{-value} = P(\chi^2 \geq \lambda) \quad (H_0: \lambda \sim \chi^2_1)$$

Multinomial

Test $H_0: \mu_1 = \mu_2 = \dots$

assuming same s.d.

$$T = \frac{\bar{y}_1 - \bar{y}_2}{\text{SE} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t(n_1 + n_2 - 2)$$

$$S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1 + n_2 - 2}$$

$$\rightarrow p\text{-value} = 2 \cdot (1 - P(T \leq t))$$

Pearson's χ^2 Goodness of Fit (dice)

$$D = \sum_{j=1}^6 \frac{(Y_j - E_j)^2}{E_j} \sim \chi^2(6-1) \text{ under } H_0$$

d = input the values

$$p\text{-value} = P(D \geq d | H_0) \approx P(W \geq d), W \sim \chi^2(5)$$

Expected Frequency

$$e_{\#} = n(\hat{\theta})e^{-\hat{\theta}} \leftarrow \text{Poisson}$$

Graphical Summaries

Q1 $\bar{y} = 1.58$ $s^2 = \frac{\sum (y_i - \bar{y})^2}{n-1} = 1.30$
 $\sum_{i=1}^n y_i^2 = (n-1)s^2 + n\bar{y}^2 = 378.34$

MLE

Given: $f_Y(y; \lambda) = \frac{e^{-\lambda} \lambda^y}{y!}$

STEPS:

- $L(\lambda) = \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{y_i}}{y_i!} = \text{constant} \cdot e^{-n\lambda} \lambda^{n\bar{y}}$
- $\ell(\lambda) = \ln(L(\lambda)) = \text{constant} - n\lambda + n\bar{y} \ln(\lambda)$
- solve $\ell'(\lambda) = 0 \rightarrow \ell'(\lambda) = -n + \frac{n\bar{y}}{\lambda}$

t, χ^2 Properties

Z_1, \dots, Z_{n+1} are IID $\sim G(0, 1)$

$\chi^2(n) \sim Z_1^2 + \dots + Z_n^2$

$t(n) \sim \frac{Z}{\sqrt{U/n}} = \frac{Z_{n+1}}{\sqrt{\sum_{i=1}^n Z_i^2/n}}$

$G(0, n) \sim nZ_1$

Likelihood Interval

40% likelihood interval \rightarrow C.I.

$\alpha: e^{-\alpha^2/2} = 0.4 \rightarrow \alpha = \sqrt{-2 \log(0.4)}$

$q: q = 2P(Z \leq \alpha) - 1 \rightarrow 100q\%$ C.I.

General 100p% C.I.

- $P(a \leq \text{pivotal quantity} \leq b) = p$
 $P(W \leq a) = P(W \geq b) = \frac{1-p}{2}$
 where $W \sim \text{table}$
- Rearrange!

Hypothesis Test (G)

Given: $\bar{y} = 101.7$ $s = 13.5$ $Y_i \sim G(\mu, \sigma)$

Q: Run a test to determine if average could be 105.

base case

① $H_0: \mu = 105$ $H_1: \mu \neq 105$

② $D: \frac{|\bar{Y} - \mu_0|}{s/\sqrt{n}}$

③ $2P(T \geq D) = 2P(T \geq 0.843) \approx 0.417$

student t since we have s not σ $n-1$ d.f.

(BONUS) 95% C.I.: $\bar{y} \pm t_{0.025, 11}^{qt} \cdot s/\sqrt{n}$

"Two-sided" hypothesis test

Testing σ^2 with $Y_i \sim G(\mu, \sigma)$

$D = \frac{(n-1)s^2}{\sigma_0^2} \rightarrow \min\{2P(U \leq D), 2P(U \geq D)\}$

Approx n based on C.I.

- $l = \text{range of } 100p\% \text{ C.I.}$
denominator of r.q.
- l should be $2 \cdot (\alpha \cdot \sigma/\sqrt{n})$

Asymptotic Pivotal Q-Proof

Given: 100 coin flips, 58 heads.

Assume $Y \sim \text{Binomial}(100, \theta)$, recall $(\theta = \frac{\text{heads}}{n})$
 $E[Y] = n\theta$, $\text{Var}(Y) = n\theta(1-\theta)$, $\tilde{\theta} = Y/n$

Prove $Z_n = \frac{\tilde{\theta} - \theta}{\sqrt{\theta(1-\theta)/n}}$ is asymptotic p.q. for θ .

- $E[\tilde{\theta}] = E[\frac{Y}{n}] = \frac{1}{n} E[Y] = \theta$
- $\text{Var}[\tilde{\theta}] = \text{Var}[\frac{Y}{n}] = \frac{1}{n^2} \text{Var}(Y) = \frac{\theta(1-\theta)}{n}$
- Note $Y = \sum_{i=1}^{100} X_i$, where X_i IID $\sim \text{Bernoulli}(\theta)$
- By CLT, big n: $Z_n = \frac{\ln(\bar{X} - E[X_i])}{\sqrt{\text{Var}[X_i]}} \sim G(0, 1)$
- Since we know distribution of Z_n , it's a p.q.

Likelihood Ratio Test Statistic TEST

Given: $n=25$ Poisson(θ) $\rightarrow \hat{\theta} = \bar{y}$, $L(\theta) \propto \theta^{n\bar{y}} e^{-n\theta}$
 Use $H_0: \theta = 2.5$, $H_1: \theta \neq 2.5$ $\bar{y} = 3.16$

- $\Lambda(\mu_0) = -2 \ln\left(\frac{L(\theta_0)}{L(\hat{\theta})}\right) = -2 \ln(R(\theta)) = 4.016$ ★
- p-value $\sim P(W \geq 4.016)$, $W \sim \chi^2(1)$ (pchisq)

↑ I should've had more examples of this!!!!!!

Linear Regression Model

$\hat{\alpha}$: average y when $x=0$.

$\hat{\beta}$: mean increase in y per x .

Test linear relationship ($\hat{\beta} = 0$)

$H_0: \beta = 0$ $H_1: \beta \neq 0$

$t: \frac{|\hat{\beta} - \beta_0|}{s_e/\sqrt{S_{xx}}} = 7.037$ df: n-2 pF

p-value: $2P(T \geq 7.037) = 0.0005$

95% Confidence Interval for $\hat{\mu}$ (given x)

$\hat{\mu}(12) = \hat{\alpha} + \hat{\beta}x = 2285.324$

★ $\hat{\alpha} + \hat{\beta}x \pm \alpha s_e \sqrt{\frac{1}{n} + \frac{(x-\bar{x})^2}{S_{xx}}}$ ★

R: TRUE: $W \leq d$ FALSE: $W \geq d$