

Functions of a random variable (i.e. finding $f_Y(y)$ from $f_X(x)$ where $Y=g(X)$)
 - generally, you find $F_X(x)$, then $F_Y(y) = P(Y \leq y) = P(X \leq g^{-1}(y)) = F_X(g^{-1}(y))$
 if g is strictly monotonic increasing

Magic formula: If X is cts ranvar with density $f_X(x)$,
 s.t. $f_X(x) = 0$ if $x \notin I$ interval,
 and g is differentiable & strictly monotonic on I ,
 then $Y=g(X)$ has density $f_Y(y) = f_X(g^{-1}(y)) \cdot \left| \frac{d}{dy} g^{-1}(y) \right|$

Joint distribution: $F(x,y) = P(X \leq x, Y \leq y)$

$$F(x_0, y_0) = \int_{-\infty}^{x_0} \int_{-\infty}^{y_0} f(x,y) dy dx$$

$$f(x_0, y_0) = \frac{\partial^2}{\partial x \partial y} F(x_0, y_0)$$

$$f_X(x_0) = \int_{-\infty}^{+\infty} f(x_0, y) dy \leftarrow \text{marginal distribution}$$

Independence: joint c.d.f. factors into product of marginal c.d.f.s:

$$F(x_1, \dots, x_n) = F_{X_1}(x_1) \dots F_{X_n}(x_n) \text{ for all } x_1, \dots, x_n$$

or equiv. the joint density function factors:

$$f(x_1, \dots, x_n) = f_{X_1}(x_1) \dots f_{X_n}(x_n) \text{ for all } x_1, \dots, x_n$$

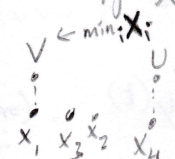
Quotient of cts ranvars:

$f_{X,Y}(x,y)$ is given, and $Z = Y/X$

$$\text{Then } f_Z(z) = \int_{-\infty}^{+\infty} |x| f(x, xz) dx$$

if X & Y are indep, then $f_Z(z) = \int_{-\infty}^{+\infty} |x| f_X(x) f_Y(xz) dx$
 where X_i are indep. and each have c.d.f. F and density f .

Extrema



$$F_U(u) = P(U \leq u) = P(X_1 \leq u) \dots P(X_n \leq u) = F(u)^n$$

$$\therefore f_U(u) = \frac{d}{du} F(u)^n = n f(u) F(u)^{n-1}$$

$$F_V(v) = 1 - (1 - F_V(v))^n$$

$$\therefore f_V(v) = n f(v) (1 - F(v))^{n-1}$$

E.g. if $X_i \sim \text{Exp}(\lambda)$, then $f_V(v) = n \cdot \lambda e^{-\lambda v} \cdot (e^{-\lambda v})^{n-1} = n \lambda e^{-n \lambda v}$

E.g. if $X_i \sim \text{Uniform}(0, \theta)$, then $f_U(u) = n f(u) \cdot F(u)^{n-1} = \frac{n u^{n-1}}{\theta^n}$ (where $0 \leq u \leq \theta$)

Expectation $E(X) = \int_{-\infty}^{+\infty} x f(x) dx$

- of a function $Y = g(X) : E(Y) = \int_{-\infty}^{+\infty} g(x) f(x) dx$
- of a function over jointly distrib. ranvar $Y = g(X_1, \dots, X_n) : E(Y) = \int \dots \int g(x_1, \dots, x_n) f(x_1, \dots, x_n) dx_1 \dots dx_n$
- if X, Y indep. ranvar : $E[XY] = E[X] \cdot E[Y]$
- if X, Y indep. ranvar & g, h are any functions, then $E[g(X)h(Y)] = E[g(X)] \cdot E[h(Y)]$
- linearity of expectation: $E(a + \sum_{i=1}^n b_i X_i) = a + \sum_{i=1}^n b_i E(X_i)$

Variance $Var(X) = E[(X - E[X])^2] = E[X^2] - (E[X])^2$

- $Var(a + bX) = b^2 Var(X)$
- Chebyshev's ineq: $P(|X - \mu| > t) \leq \frac{\sigma^2}{t^2}$
(if σ^2 is small then X should not deviate too much from μ)
- If X_i are indep then $Var(\sum_{i=1}^n X_i) = \sum_{i=1}^n Var(X_i)$

Moment generating function: $M(t) = E[e^{tx}] = \int_{-\infty}^{\infty} e^{tx} f(x) dx$ (cts)

- m.g.f. uniquely determines p.d.f. if $M(t) = \sum_i e^{tx_i} p(x_i)$ (discrete)
- $M^{(r)}(0) = E[X^r]$
↑
rth differentiation
- $M_{a+bX}(t) = e^{at} M_X(bt)$ (moment of a linear change of variable)
- if X, Y indep and $Z = X + Y$, then $M_Z(t) = M_X(t) M_Y(t)$ (on the common interval where M_X and M_Y exist)

Limit theorems

- Weak Law of Large Numbers (WLLN):
If X_1, \dots, X_n, \dots are a sequence of indep. ranvar with $E(X_i) = \mu$ and $Var(X_i) = \sigma^2$ (not necessarily normal),
Then $\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$ satisfies:
 $\forall \epsilon > 0, P(|\bar{X}_n - \mu| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$ (Pf: using Chebyshev)
(in other words, taking more indep. readings gets us arbitrarily close to the mean)
- "Converge in probability to α ": $P(|Z_n - \alpha| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$
- If we don't know σ^2 , we can estimate it with $\frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2$

"Converges in distribution"

If X, X_1, X_2, \dots is a sequence of r.v.s with c.d.f.s F, F_1, F_2, \dots resp, then

X_n converges in distribution to X if $F_n(x) \rightarrow F(x)$ as $n \rightarrow \infty$

ptwise convergence of the c.d.f.

Central Limit Theorem (CLT):

If X_1, X_2, \dots are indep r.v.s with mean μ and variance σ^2 with c.d.f. F (and m.g.f. is defined in a neighbourhood of 0), then:

$$S_n := \sum_{i=1}^n X_i$$

$$P\left(\frac{S_n - n\mu}{\sigma\sqrt{n}} \leq x\right) \rightarrow \Phi(x) \text{ as } n \rightarrow \infty$$

std. normal c.d.f.

(i.e. $\frac{S_n}{n}$ fluctuates around μ in a normal dist. with variance $\frac{\sigma^2}{n}$)

- skewed dist. or with large tails will need larger n for a good approx.
- sometimes, we are given $E(\hat{\alpha})$ and $Var(\hat{\alpha})$ that are already obtained as an estimate from a large sample. Then,

$$P\left(\frac{\hat{\alpha} - E(\hat{\alpha})}{\sqrt{Var(\hat{\alpha})}} \leq x\right) \rightarrow \Phi(x) \text{ as } n \rightarrow \infty$$

(or in other words, $\hat{\alpha} \rightarrow N(E(\hat{\alpha}), Var(\hat{\alpha}))$)

t Distribution : Given $Z \sim N(0,1)$ (i.e. stdnorm) and $U \sim \chi_n^2$ } indep

then $T = \frac{Z}{\sqrt{U/n}} \sim t_n$ (t_n is "Cauchy distribution")

• $f(t) = \frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2})\sqrt{n\pi}} (1 + \frac{t^2}{n})^{-\frac{n+1}{2}}$ t distribution with n dof.

- note: it is symmetric about 0, since $f(t) = f(-t)$
- As $n \rightarrow \infty$, $t_n \rightarrow N(0,1)$ (tails become lighter)

F Distribution : Given $U \sim \chi_m^2$ and $V \sim \chi_n^2$ } indep

then $W = \frac{U/m}{V/n} \sim F_{m,n}$

• $f(w) = \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} (\frac{m}{n})^{\frac{m}{2}} w^{\frac{m}{2}-1} (1 + \frac{m}{n}w)^{-\frac{m+n}{2}}$ F distribution with m and n degrees of freedom.

• $E(W) = \frac{n}{n-2}$ for $n > 2$

• $(t_n)^2 \sim F_{1,n}$

• If $X, Y \sim \text{Exp}(\lambda=1)$, then $\frac{X}{Y} \sim F_{2,2}$ (indep)

On normal distributions & sample variance

• If $X \sim N(\mu_x, \sigma_x^2)$ and $Y \sim N(\mu_y, \sigma_y^2)$ (indep)

then $X+Y \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$

• let X_1, \dots, X_n be indep $N(\mu, \sigma^2)$ r.v.s:

• $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ (sample mean) and $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$

• $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ (sample variance)

↑ note the (n-1):

• \bar{X} and $(X_1 - \bar{X}, \dots, X_n - \bar{X})$ are indep.

• \bar{X} and S^2 are indep.

• $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$ (or $\frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-1}^2$, where $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$)

• $\frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$

Estimation

• given data points X_1, \dots, X_n (indep, same dist)

• want to estimate θ (parameters for dist); the indiv. distrib is $f(x|\theta)$, joint distrib is $f(x_1|\theta) \dots f(x_n|\theta)$

Method of Moments

• k^{th} moment is $\mu_k := E(X^k)$

• want observed $\hat{\mu}_k$ to be equal to expected μ_k .

$$\hat{\mu}_k = \frac{1}{n} \sum_{i=1}^n X_i^k$$

• Steps: (1) Express low-order moments (starting from $k=1$) in terms of parameters
(e.g. $\mu_1 = E(X) = \mu$; $\mu_2 = E(X^2) = \mu^2 + \sigma^2$)

→ if moment does not depend on params, then ignore. e.g. μ_1, σ^2

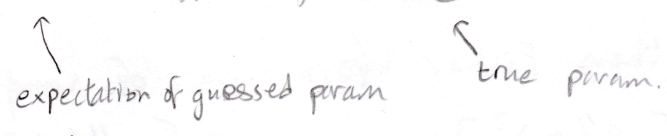
(2) Rearrange expressions to make parameters the subjects
(e.g. $\mu = \mu_1$; $\sigma^2 = \mu_2 - \mu_1^2$)

(3) Calculate and substitute sample moments

(e.g. calculate $\hat{\mu}_1, \hat{\mu}_2$ from the given data, sub in μ_1, μ_2 to find μ and σ^2)
e.g. $\hat{\mu} = \bar{X}$; $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$

Bias:

• How much $E(\hat{\theta})$ differs from θ



$$(e.g. E(\hat{\alpha}) = 3 E[\bar{X}] = 3 \frac{1}{n} \sum_i E[X_i] = 3 \frac{1}{n} \cdot n \cdot \frac{\alpha}{3} = \alpha)$$

Annotations: $E(\hat{\alpha})$ from MOM estimate; $E[X_i]$ L.O.E.; $\frac{1}{n} \cdot n \cdot \frac{\alpha}{3}$ from the mean (based on distribution param)

Consistency of an estimate $\hat{\theta}_n$ of θ

• $\hat{\theta}_n$ is consistent in probability if $\hat{\theta}_n$ converges in probability to θ as $n \rightarrow \infty$.
i.e. $\forall \epsilon > 0, P(|\hat{\theta}_n - \theta| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$.

• MOM estimates are consistent (due to WLLN)
• so (assuming functions are continuous), estimates will converge.

Method of Maximum Likelihood

• Random vars X_1, \dots, X_n ; parameter to estimate: θ

• density or freq: $f(x_1, \dots, x_n | \theta) = f(x_1 | \theta) \dots f(x_n | \theta)$

$$lik(\theta) = f(x_1, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta) \text{ (if i.i.d.)}$$

if indep. \nmid i.i.d.

want to find θ that maximises this.

• equiv. to finding θ that maximises $\ell(\theta) := \log(lik(\theta)) = \sum_{i=1}^n \log(f(x_i | \theta))$
• differentiate $\ell(\theta)$... use partial derivative if more than one parameter

Large sample theory for MLE (p. 38)

- Fischer information: $I(\theta) = E \left(\left[\frac{\partial}{\partial \theta} \log f(x|\theta) \right]^2 \right)$ (indiv. sample)
- if sufficiently smooth: $I(\theta) = -E \left(\frac{\partial^2}{\partial \theta^2} \log f(x|\theta) \right)$
 - ↑ over x given θ
 - ↑ natural log
 - ↑ may not hold if support of f depends on θ.
- if sufficiently smooth, and i.i.d.

(a) MLE $\hat{\theta}$ is consistent

(b) $\sqrt{nI(\theta_0)} (\hat{\theta} - \theta_0) \rightarrow N(0, 1)$

↑ estimate ↑ true value

(i.e. $\hat{\theta}$ is approx $N(\theta_0, \frac{1}{nI(\theta_0)})$)
 asymptotic variance

Fischer info of i.i.d. sample of size n is $nI(\theta)$

If $I(\theta_0)$ is unknown, use $I(\hat{\theta})$ instead.

If not i.i.d., then Fisher information of sample is $E[l'(\theta)^2]$ or $-E[l''(\theta)]$

and the asymptotic variance is $\frac{1}{E[l'(\theta)^2]}$ or $-\frac{1}{E[l''(\theta)]}$

Confidence interval:

For normal dist:

- 100(1-α)% confidence interval for μ is $\bar{X} \pm \frac{s}{\sqrt{n}} t_{n-1}(\frac{\alpha}{2})$ (since $\frac{\bar{X}-\mu}{s/\sqrt{n}} \sim t_{n-1}$)
- 100(1-α)% confidence interval for σ² is $(\frac{n\hat{\sigma}^2}{\chi^2_{n-1}(\frac{\alpha}{2})}, \frac{n\hat{\sigma}^2}{\chi^2_{n-1}(1-\frac{\alpha}{2})})$ (since $\frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-1}$)

Large sample theory approx:

$P \left(-z(\frac{\alpha}{2}) \leq \sqrt{nI(\hat{\theta})} (\hat{\theta} - \theta_0) \leq z(\frac{\alpha}{2}) \right) \approx 1 - \alpha$

so 100(1-α)% confidence interval for θ is $\hat{\theta} \pm \frac{z(\frac{\alpha}{2})}{\sqrt{nI(\hat{\theta})}}$ (right tail cumulative)

Efficiency

- Mean squared error (MSE): $E[(\hat{\theta} - \theta_0)^2] = \underbrace{\text{Var}(\hat{\theta})}_{\text{variance}} + \underbrace{(E(\hat{\theta}) - \theta_0)^2}_{\text{squared bias}}$
- Efficiency of $\hat{\theta}$ relative to $\tilde{\theta}$: $\text{eff}(\hat{\theta}, \tilde{\theta}) = \frac{\text{Var}(\tilde{\theta})}{\text{Var}(\hat{\theta})}$
 - two estimates of θ
 - how concentrated the estimate is

• Cramér-Rao lower bound: for any unbiased estimate of θ (from n i.i.d. r.v.s)

$$\text{Var}(T) \geq \frac{1}{nI(\theta)}$$

- gives the best possible variance
- An efficient unbiased estimate attains this lower bound
 - if it is biased, then we don't say that it is efficient.

Sufficiency: finding a statistic that contains all the info in a sample about θ

- A statistic $T=T(X_1, \dots, X_n)$ is sufficient for θ if $P(X_1=x_1, \dots, X_n=x_n | T=t)$ does not depend on θ (for any t)

** Equiv: $T(\tilde{x})$ is sufficient for θ iff $f(\tilde{x} | \theta) = g(T(\tilde{x}), \theta) \cdot h(\tilde{x})$

(i.e. f factorizes into g and h)

e.g. for $X_i \sim N(\mu, \sigma^2)$

$$f(\tilde{x} | \mu, \sigma) = \frac{1}{\sigma^n (2\pi)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2} \left(\sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2\right)\right)$$

so $\sum_{i=1}^n x_i$ and $\sum_{i=1}^n x_i^2$ are (together) sufficient statistics

• Exponential family of dists: those with one-parameter

$$f(x|\theta) = \begin{cases} \exp(c(\theta)T(x) + d(\theta) + S(x)), & x \in A \\ 0 & x \notin A \end{cases}$$

↑
can be a vector

where A does not depend on θ .

(i.i.d.) $f(\tilde{x}|\theta) = \prod_{i=1}^n f(x_i|\theta)$ → These factorize into $\exp(c(\theta) \sum_{i=1}^n T(x_i) + nd(\theta)) \exp(\sum_{i=1}^n S(x_i))$

so $\sum_{i=1}^n T(x_i)$ is a sufficient statistic.

k-parameters: $f(x|\theta) = \begin{cases} \exp\left[\sum_{j=1}^k c_j(\theta) T_j(x) + d(\theta) + S(x)\right], & x \in A \\ 0 & x \notin A \end{cases}$

The parameters $\sum_{i=1}^n T_j(x_i)$ ($1 \leq j \leq k$) are sufficient statistics

Rao-Blackwell thm : (sufficient statistics make the best estimators)

· If $\hat{\theta}$ is an estimator of θ with $E(\hat{\theta}^2) < \infty$ for all θ ,

and T is sufficient for θ ,

and $\tilde{\theta} = E(\hat{\theta} | T)$,

then $\forall \theta, \quad E[(\tilde{\theta} - \theta)^2] \leq E[(\hat{\theta} - \theta)^2]$

(where equality holds iff $\hat{\theta} = \tilde{\theta}$)

Hypothesis Testing

- H_0 : null hypothesis
- H_1 : alternative hypothesis
- Type I error: rejecting H_0 when it is true
- $P(\text{Type I error}) = \alpha$: significance level of the test
- Type II error: accepting H_0 when it is false
- $P(\text{Type II error}) = \beta$
- $P(H_0 \text{ rejected when it is false}) = 1 - \beta$: power of the test
- test statistic: some value X s.t. we decide to accept/reject H_0 based on X .
- acceptance region: set of values for test statistic s.t. we accept H_0
- rejection region: set of values for test statistic s.t. we reject H_0
- null distribution: distribution of test statistic under H_0
- simple hypothesis: H_0 and H_1 are completely specified distributions (i.e. no unfixed variables)
(e.g. $H_0: X \sim \text{binomial}(10, 0.5)$)

Neyman-Pearson lemma: Given that H_0 and H_1 are simple hypotheses, for any significance level α , the likelihood ratio test has the most power.

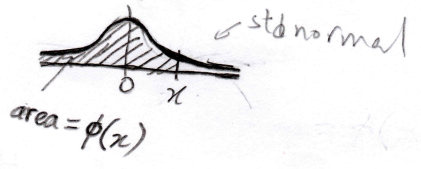
Likelihood ratio test: test statistic = $\frac{P(x|H_0)}{P(x|H_1)}$ → Want to reject H_0 when likelihood ratio $< c$ (for some c to get the desired significance level α)

usually, we can map this back to something like $\bar{X} > \kappa_0$ or $\bar{X} < \kappa_0$.

see p. 49 for ~~an~~ example when α is given

p-value: smallest significance level under which H_0 will be rejected (i.e. $P(T \geq t_{\text{obs}} | H_0)$)

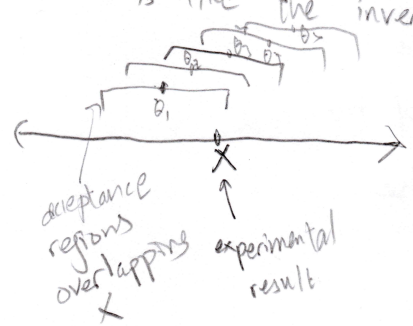
↑ depends on direction ↑ the observed value



Uniformly most powerful test: The rejection region depends on the null distribution and H_0 but not on H_1 's variables (e.g. μ_1) apart from the direct

Confidence intervals: p. 52, Ex 34.

= the set of values for μ_0 for which $H_0: \mu = \mu_0$ is accepted.
(so it is like the inverse of a hypothesis test)



confidence interval is the set of all θ where X is accepted.

$C(X) = \{\theta \mid X \in A(\theta)\}$

100(1- α)% confidence region for θ .
acceptance region of θ at level α

$A(\theta_0) = \{X \mid \theta_0 \in C(X)\}$

— p. 53, Ex 35.

Generalized likelihood ratio tests for non-simple (i.e. composite) hypotheses:

- not generally optimal, but one of the best ways since no optimal tests exist.
- Let Ω be the set of all possible values of θ , and $\omega_0 \subseteq \Omega$ be the subset where H_0 is valid (i.e. $H_0: \theta \in \omega_0$)

Then $\Lambda = \frac{\max_{\theta \in \omega_0} \text{lik}(\theta)}{\max_{\theta \in \Omega} \text{lik}(\theta)}$

rejection region: $\Lambda \leq \lambda_0$ (i.e. it is not likely that $\theta \in \omega_0$)
↑ for some λ_0 .

~~...~~
~~...~~
• $-2 \log \Lambda \rightarrow \chi^2$ with $df = \dim \Omega - \dim \omega_0$
always natural degrees of freedom where $\dim S$ is the number of free params under S .

e.g. for $H_0: \mu = \mu_0$ and $H_1: \mu \neq \mu_0$ where X_1, \dots, X_n i.i.d. from normal dist with variance σ^2 , then $\dim \Omega = 1$ (μ is free variable), $\dim \omega_0 = 0$ (no free params), so $-2 \log \Lambda \rightarrow \chi_1^2$ but not σ^2 (no free params),
(actually, it is equal for normal dist)

Likelihood ratio tests for ~~multinomial~~ multinomial distribution (p.55, Ex. 36.) (11)

- the vector $\tilde{p} = \tilde{p}(\theta)$ where $\theta \in \omega_0$ for some θ
- for Ω , the \tilde{p} are free such that they are valid (i.e. $p_i \geq 0$ and $\sum p_i = 1$)

$$\text{likelihood ratio } \Lambda = \frac{\max_{\theta \in \omega_0} \text{lik}(\theta)}{\max_{\theta \in \Omega} \text{lik}(\theta)} \rightarrow \max_{\theta \in \omega_0} \left(\frac{n!}{x_1! \dots x_m!} p_1(\theta)^{x_1} \dots p_m(\theta)^{x_m} \right)$$

where x_1, \dots, x_m are observed cell counts

$$= \frac{\left(\frac{n!}{x_1! \dots x_m!} p_1(\hat{\theta})^{x_1} \dots p_m(\hat{\theta})^{x_m} \right)}{\left(\frac{n!}{x_1! \dots x_m!} \hat{p}_1^{x_1} \dots \hat{p}_m^{x_m} \right)}$$

where $\hat{\theta}$ is the mle

and $\hat{p}_i = \frac{x_i}{n}$ is the unrestricted mle

$$= \prod_{i=1}^m \left(\frac{p_i(\hat{\theta})}{\hat{p}_i} \right)^{x_i}$$

$$\therefore -2 \log \Lambda = 2 \sum_{i=1}^m O_i \log \left(\frac{O_i}{E_i} \right) \text{ where } O_i = n \hat{p}_i = x_i \text{ (observed count)}$$

$\dim \Omega = m-1$ (since there is a constraint that $\sum p_i = 1$)
 $k := \dim \omega_0$ is the dimensionality of θ (possibly a vector)

So $-2 \log \Lambda \rightarrow \chi^2_{m-k-1}$

Pearson's χ^2 statistic: $\chi^2 = \sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i}$

(Ex. 36)

$$\chi^2 \rightarrow -2 \log \Lambda \rightarrow \chi^2_{m-k-1}$$

not equal, but asymptotically equivalent under H_0

Comparing two samples

If $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu_X, \sigma^2)$

and $Y_1, \dots, Y_m \stackrel{iid}{\sim} N(\mu_Y, \sigma^2)$

we want to see if $\mu_X = \mu_Y$ or not.

If σ^2 is known, then $Z = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim N(0, 1)$

so a $100(1-\alpha)\%$ confidence interval for $\mu_X - \mu_Y$ is $(\bar{X} - \bar{Y}) \pm z\left(\frac{\alpha}{2}\right) \sigma \sqrt{\frac{1}{n} + \frac{1}{m}}$

If σ^2 is unknown, then estimate it from pooled sample variance:

$$s_p^2 = \frac{(n-1)S_X^2 + (m-1)S_Y^2}{m+n-2}$$

where $S_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$

and $S_Y^2 = \frac{1}{m-1} \sum_{i=1}^m (Y_i - \bar{Y})^2$

$S_{\bar{X}-\bar{Y}} = s_p \sqrt{\frac{1}{n} + \frac{1}{m}}$ (est. std. error of $\bar{X} - \bar{Y}$)

$t = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{S_{\bar{X}-\bar{Y}}} \sim t_{m+n-2}$

so a $100(1-\alpha)\%$ confidence interval for $\mu_X - \mu_Y$ is $(\bar{X} - \bar{Y}) \pm t_{m+n-2}\left(\frac{\alpha}{2}\right) S_{\bar{X}-\bar{Y}}$

For hypothesis testing:

test statistic $t = \frac{\bar{X} - \bar{Y}}{S_{\bar{X}-\bar{Y}}}$

- $H_0: \mu_X = \mu_Y$
- $H_1: \mu_X \neq \mu_Y \iff |t| > t_{m+n-2}\left(\frac{\alpha}{2}\right)$
- $H_2: \mu_X > \mu_Y \iff t > t_{m+n-2}(\alpha)$
- $H_3: \mu_X < \mu_Y \iff t < -t_{m+n-2}(\alpha)$

If variances of X_i 's and Y_i 's are not equal, then: $t = \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}}}$ (?)

test statistic $t = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}}}$

where $df = \frac{\left(\frac{S_X^2}{n} + \frac{S_Y^2}{m}\right)^2}{\left(\frac{S_X^2}{n}\right)^2/(n-1) + \left(\frac{S_Y^2}{m}\right)^2/(m-1)}$

nonnormality: for large sample size, it is approximately valid