1/26

$\hat{p}$ is a value calculated from observed data

⎰ A statistic

⎰ probability distribution is called sampling distribution

If the mean $\hat{p} = p$ it is unbaised

Sampling without Replacement

$$E(\hat{p}) = p \quad \text{still}$$

$$Var[\hat{p}] = \langle \hat{p}^2 \rangle - \langle \hat{p} \rangle^2$$

$$\hat{p} = \left( \frac{X_1 + X_2 + \cdots X_n}{n} \right)$$

$$\langle \hat{p}^2 \rangle = \frac{1}{n^2} \left\langle X_1^2 + \cdots X_n^2 + 2X_1 X_2 + \cdots 2X_{n-1} X_n \right\rangle$$

$$= \frac{1}{n^2} \left\langle n X_1^2 + \binom{n}{2} \langle 2X_1 X_2 \rangle \right\rangle$$

↑ Expected Values are same for each $X_i$

$$= \frac{1}{n^2} \cdot \langle X_1^2 \rangle + \frac{n-1}{n} \langle 2X_1 X_2 \rangle$$

‖        ‖

$\langle X_1 \rangle$     $P[X_1 = 1 \text{ and } X_2 = 1]$

since $X_i = 0$ or $1$    $= P[X_1 = 1] \cdot P[X_2 = 1 \mid X_1 = 1]$

$$= p \cdot p\frac{N-1}{N-1}$$

$$Var(\hat{p}) = \langle \hat{p}^2 \rangle - \langle \hat{p} \rangle^2$$

$$= \frac{1}{n} \cdot p + \frac{n-1}{n} \cdot p \cdot p\frac{N-1}{N-1}$$

$$= \frac{p(1-p)}{n} \left( 1 - \frac{n-1}{N-1} \right)$$

⟅ Correction factor

Some stuff on Central Limit Theorem

Types of Inference Questions

- Hypothesis Testing : Asking a yes or no question about the distribution

- Estimation : Determining the distribution or some characteristic

- Confidence Interval

## 1128

Random Variables

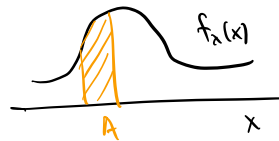Discrete : finite our countably infinite number of possible values

Probability Mass Function : $f_X(x) = \mathbb{P}[X=x]$

For a set of values A $\quad \mathbb{P}[X \in A] = \sum_{x \in A} f_X(x) = 1$

when A is all possible

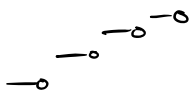Continuous : takes any real value

Probability Density Function :

$$\mathbb{P}[X \in A] = \int_A f_X(x) dx$$

$$\int_R f_X(x) dx = \mathbb{P}[X \in R] = 1$$



Cumulative Distribution Function (CDF): $F_X(x) = \mathbb{P}[X \leq x]$

Discrete CDF

$$F_X(x) = \sum_{y : y \leq x} f_X(y)$$

Continuous case : $F_X(x) = \int_{-\infty}^{\infty} f_X(y) dy$

Properties of CDF $F_X(x)$

· Increasing : For any $X \leq Y$ always $F_X(x) \leq F_X(y)$

As $x \to -\infty$ $\quad F_X(x) \to 0$ $\quad$ as $x \to \infty$ $f_X(x) \to 1$

If $F_X(x)$ is strictly increasing and continuous, then $F_X(x)$ has an inverse function $F_X^{-1}$

$: (0,1) \to \mathbb{R}$ called quantile function of X

$F_X^{-1}(t)$ satisfies $\mathbb{P}[X \leq x] = t$ i.e. $F_X^{-1}(t)$ is the $t^{th}$ quantile of X

Expected Value and Variance

$$\mathbb{E}[X] = \sum_{x \in X} x \cdot f_X(x) \qquad \mathbb{E}[X] = \int_R x f(x) dx$$

$$g: \mathbb{R} \to \mathbb{R}$$

$$\mathbb{E}[g(x)] = \sum_{x \in X} g(x) f_X(x)$$

$$\mathbb{E}[g(x)] = \int_{\mathbb{R}} g(x) f_X(x) dx$$

Expectation is linear!!!

$$\text{Var}[X] = \mathbb{E}\left[(X - \mathbb{E}(x))^2\right] = \mathbb{E}[X^2] - \mathbb{E}(x)^2$$

Translation invariant: $\text{Var}[X+c] = \text{Var}[X]$

Multiply by constant: $\text{Var}[cX] = c^2 \text{Var}[X]$

Independent RV: $\text{Var}[X_1 + X_2] = \text{Var}[X_1] + \text{Var}[X_2]$

Dependent RV: $\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X,Y]$

The standard deviation of $X$ is $SD(x) = \sqrt{\text{Var}(X)}$

A Bernoulli r.v. $X \sim \text{Bernoulli}(p)$ takes all possible values 0 and 1

The PMF is $f_X(1) = \mathbb{P}[X=1] = p$

$$f_X(0) = \mathbb{P}[X=0] = 1-p$$

Mean: $\mathbb{E}(x) = p \cdot 1 + (1-p) \cdot 0 = p$

Variance: $\text{Var}[X] = \mathbb{E}(X^2) - \mathbb{E}(x)^2$

$$= E(X) - E(X^2) \quad \text{since } X^2 = X \text{ when } X=0 \text{ or } 1$$

$$= p - p^2 = p(1-p)$$

A Binomial distribution is repeated bernoulli trials
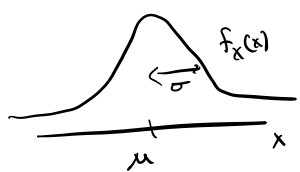
$$X = X_1 + X_2 + \ldots X_n \sim \text{Binomial}(n,p)$$

PMF is $f_X(x) = \binom{n}{x} p^x (1-p)^{n-x}$ for $x \in \{0,1,\ldots n\}$

Mean: $\mathbb{E}(x) = \mathbb{E}(X_1 + \ldots + X_n) = \mathbb{E}(X_1) + \ldots + \mathbb{E}(X_n) = np$

Variance: $\text{Var}[X] = \text{Var}[X_1 + \ldots X_n] = \text{Var}[X_1] + \ldots \text{Var}[X_n] = np(1-p)$

A normal r.v. $X \sim N(\mu, \sigma^2)$ is a continuous r.v. on $\mathbb{R}$,

PDF $\quad f_X(x) = \dfrac{1}{\sqrt{2\pi\sigma^2}} \, e^{-\frac{(x-\mu)^2}{2\sigma^2}}$



$$E(X) = \int_{-\infty}^{\infty} x \cdot \frac{1}{\sqrt{2\pi\sigma^2}} \, e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \mu$$

$$\text{Var}[X] = \int_{-\infty}^{\infty} (x-\mu)^2 \cdot \frac{1}{\sqrt{2\pi\sigma^2}} \, e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \sigma^2$$

A Gamma r.v. $X \sim \text{Gamma}(\alpha, \beta)$ (for $\alpha, \beta > 0$) is continuous

$$f_X(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} \, x^{\alpha-1} \, e^{-\beta x} \qquad \text{for } x > 0$$

$$\Gamma(\alpha) = \int_{-\infty}^{x} x^{\alpha-1} \, e^{-x} \, dx$$

For integer $n \geq 1$, $\quad \Gamma(n) = (n-1)!$

# Joint Distributions

Discrete case: Joint PMF $\quad f_{X_1 \cdots X_n}(x_1 \ldots x_n) = P[X_1 = x_1 \ldots X_k = x_k]$

Continuous Case: Joint PDF $\quad f_{X_1 \cdots X_k}(x_1 \ldots x_k)$

$$P[(X_1 \ldots X_n) \in A] = \iint_A f_{X_1 \cdots X_n}(x_1 \ldots x_n) \, dx_1 \ldots dx_k$$

Something about multi-nomial
$\quad \mapsto$ generalized binomial

$$f_{X_1 \cdots X_n}(x_1 \ldots x_n) = \binom{n}{x_1 \cdots x_n} P_1^{x_1} P_2^{x_2} \cdots P_n^{x_k}$$

$$\frac{n!}{x_1! \, x_2! \ldots x_n!}$$

Random Variable $X_1 \ldots X_k$ are independent

$$f_{X_1 \cdots X_n}(x_1 \ldots x_n) = f_{X_1}(x_1) \times f_{X_2}(x_2) \ldots \times f_{X_k}(x_n)$$

Properties of independent r.v

i) $\mathbb{P}[X_1 \in A_1 \text{ and } X_2 \in A_2 \text{ and } \dots X_k \in A_k]$

$$= \mathbb{P}[X_1 \in A_1] \times \mathbb{P}[X_2 \in A_2] \times \dots \times \mathbb{P}[X_k \in A_k]$$

ii) For any function $g_1 \dots g_k : \mathbb{R} \to \mathbb{R}$

$$\mathbb{E}[g_1(X_1) \cdot g_2(X_2) \dots g_k(X_k)] = \mathbb{E}[g_1(X_1)] \dots \mathbb{E}[g_k(X_k)]$$

E.g. If $X, Y$ are independent, then $\mathbb{E}[XY] = E(X) \, E(Y)$

Covariance between r.v.'s $X, Y$

$$\text{Cov}[X,Y] = \mathbb{E}[(X - E(X))(Y - E(Y))] = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y)$$

$$\text{Cov}[X, X] = \text{Var}[X]$$

Translationally Invariant: $\text{Cov}[X+a, Y+b] = \text{Cov}[X,Y]$

Bilinear: $\text{Cov}[aX, bY] = ab \, \text{Cov}[X,Y]$

If $X_i$ and $X_j$ are independent

$$\text{Cov}[X_i, X_j] = \mathbb{E}(X_i Y_j) - \mathbb{E}(X_i)\mathbb{E}(X_j)$$
$$= E(X_i) E(X_j) - \mathbb{E}(X_i)\mathbb{E}(X_j) = 0$$

Independent $\to$ Cov $= 0$

Correlation

$$\text{Corr}(X,Y) = \frac{\text{Cov}[X,Y]}{\sqrt{\text{Var}(X)} \sqrt{\text{Var}[Y]}}$$

1/31

Def: The moment generating function (MGF) of a random variable $X$ is given by

$$M_X(t) = \mathbb{E}[e^{tX}]$$

ex. Normal MGF

Let $X \sim \mathcal{N}(0,1)$. What is MGF of $M_X(t)$

$$M_X(t) = \mathbb{E}[e^{tX}] = \int_{-\infty}^{\infty} e^{tx} f_X(x) \, dx$$

$$= \int_{-\infty}^{\infty} e^{tx} \frac{1}{\sqrt{2\pi}} e^{\frac{-x^2}{2}} \, dx$$

$$= e^{-t^2/2} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-t)^2} = e^{-t^2/2}$$

If $M_X(t)$ is finite in a small interval around 0, then it uniquely determines the distribution of X

Theorem: Let X and Y be different random variables where MGF's are finite in an interval $(-h, h)$ around 0, and

$M_X(t) = M_Y(t)$ for all $t \in (-h, h)$. Then X and Y have the same distribution

If $X_1 \dots X_n$ are independent r.v.'s then

$$M_{X_1 + \dots + X_n}(t) = \mathbb{E}\left[e^{t(X_1 + \dots X_n)}\right]$$

$$= \mathbb{E}\left[e^{tX_1}\right] \times \dots \times \mathbb{E}\left[e^{tX_n}\right] = M_{X_1}(t) \times \dots M_{X_n}(t)$$

Multivariate Normal Distribution

In k dimensions it is a continuous distribution for $(X_1 \dots X_k) \in \mathbb{R}^k$

It is specified by

mean vector: $\mu \in \mathbb{R}^k$

covariance matrix: $\Sigma \in \mathbb{R}^{k \times k}$

k-dimensional generalization of the normal $\mathcal{N}(\mu, \sigma^2)$

Def: $(X_1 \dots X_k)$ are multivariate normal if for any constants $a_1 \dots a_k \in \mathbb{R}$, the linear combination $a_1 X_1 + a_2 X_2 + \dots a_k X_k$ has a

normal distribution. More specifically, $(X_1 \dots X_k) \sim \mathcal{N}(\mu, \Sigma)$ if in addition h this property

$$\mathbb{E}[X_i] = \mu_i \text{ and } Var[X_i] = \Sigma_{ii} \quad \text{for each } i = 1 \dots k$$

$$Cov[X_i, X_j] = \Sigma_{i,j} \quad \text{for } i \neq j$$

<span style="color:orange">Note: This implies each individual r.v. is normal</span>

You can prove this via expanding moment generating functions

$Cov(X, Y) = 0$ implies independence for bivariate normal

theorem: If $\mathbb{X} = (X_1 \dots X_k)$ is multivariate normal and $\mathbb{X}_1$ and $\mathbb{X}_2$ are two subvectors that are in pairwise uncorrelated, then $\mathbb{X}_1$ and $\mathbb{X}_2$ are independent.

Sampling Distribution of Statistics

For data $X_1 \dots X_n$ a statistic $T(X_1 \dots X_n)$ is any value computed from this data

· Sample mean: $\bar{X} = \frac{1}{n}(X_1 + \dots + X_n)$

· Sample variance: $s^2 = \frac{1}{n-1}\left((X_1 - \bar{X})^2 + \dots + (X_n - \bar{X})^2\right)$

· Sample Range: $R = \max(X_1 \dots X_n) - \min(X_1 \dots X_n)$

A statistic is any function of data

Randomness of data induces a sampling distribution

## Chi-Squared Distribution

$$X_1, \ldots, X_n \overset{IID}{\sim} N(0,1)$$

Sampling distribution of $X_1^2 + \cdots + X_n^2 \sim \chi_n^2$ with $n$ degrees of freedom

$$M_{X_1^2 + \cdots + X_n^2} = M_{X_1^2}(t) \times \cdots \times M_{X_n^2}(t)$$

$$M_{X_i^2}(t) = \mathbb{E}\left[e^{tX_i^2}\right] = \int_{-\infty}^{\infty} e^{tx^2} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{(t-1/2)x^2} dx$$

$t \geq 1/2$ blows up

$$M_{X_i^2}(t) = \frac{1}{\sqrt{1-2t}} \int_{-\infty}^{\infty} \sqrt{\frac{1-2t}{2\pi}} e^{-\frac{1}{2}(1-2t)x^2} dx$$

Just MGF of Gamma $\left(\frac{1}{2}, \frac{1}{2}\right)$   So $X_i^2 \sim$ Gamma $\left(\frac{1}{2}, \frac{1}{2}\right)$

$$M_{X_1^2 + \cdots + X_n^2}(t) = \begin{cases} \infty & t \geq 1/2 \\ (1-2t)^{-n/2} & t < 1/2 \end{cases} \quad \leftarrow \text{Gamma} \left(\frac{n}{2}, \frac{1}{2}\right)$$

Distributions that are difficult to study will be studied by approximation

Simulate via R

Asymptotic approximations

1) Faster

2) Theoretical Understanding

Weak Law of Large Numbers

Suppose $X_1, \ldots X_n$ are IID, with $\mathbb{E}[X_i] = \mu$ and $\text{Var}[X_i] < \infty$. Let

$$\bar{X} = \frac{1}{n}(X_1 + \cdots + X_n)$$

$$\bar{X} \to \mu \quad \text{as} \quad n \to \infty$$

## Central Limit Theorem

Suppose $X_1 \ldots X_n$ are iid with $\mathbb{E}[X_i] = \mu$ and $\text{var}[X_i] = \sigma^2$

Let $\bar{X} = \frac{1}{n}(X_1 + \ldots + X_n)$ Then

$$\sqrt{n}\left(\frac{\bar{X} - \mu}{\sigma}\right) \to N(0,1) \quad \text{is distribution as } n \to \infty$$

## Normal Distribution Approximation accuracy depends on
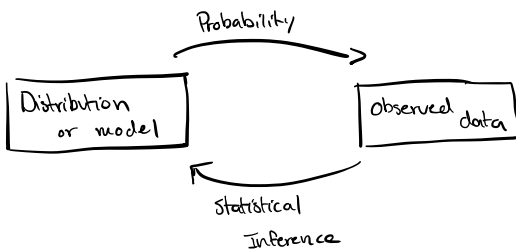
i) Sample size

ii) Skewness

iii) Heavyness of tails

## Continuous Mapping Theorem

Let $g(x)$ be a continuous function of $x$. As $n \to \infty$

i) If $S_n \to Z$ in distribution, then $g(S_n) \to g(Z)$ in dist

ii) Analogous for probability

2/7



Probability

Distribution or model → Observed data

Statistical Inference

## Einstein's Theory of Brownian Motion

Suppose the particle is at position $P_t \in \mathbb{R}^2$ at time $t$. Then at time $t + \Delta t$, the position $P_{t+\Delta t}$ is random

and has a bivariate normal distribution around $P_t$

$$P_{t+\Delta t} - P_t \sim N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{pmatrix}\right)$$

$$\sigma^2 = \frac{RT}{3\pi\eta r N_A}(\Delta t) \quad \leftarrow \text{how tf he derive this}$$

Hypothesis test is a binary question about the distribution of the data

Accept/reject null hypothesis $H_0$ in favor of alternative hypothesis $H_1$

For the magicians die. Null hypothesis: die is fair

$$H_0 : (X_1 \ldots X_6) \sim \text{Multinomial}\left(n, (\tfrac{1}{6} \ldots \tfrac{1}{6})\right)$$

$$H_1 : (X_1 \ldots X_6) \sim \text{Multinomial}\left(n, (P_1 \ldots P_6)\right)$$

$$P_1 \ldots P_6 \neq (\tfrac{1}{6}, \ldots \tfrac{1}{6})$$

Setting up notation for brownian motion experiment

$$(X_1, Y_1) = P_1 - P_0$$
$$(X_2, Y_2) = P_2 - P_1$$
$$\vdots$$
$$(X_n, Y_n) = P_n - P_{n-1}$$

displacements are measured every 30 sec

$$H_0 : (X_1, Y_1) \cdots (X_n, Y_n) \overset{iid}{\sim} N\left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 2.23e{-}7 & 0 \\ 0 & 2.23e{-}7 \end{pmatrix} \right)$$ ← Einstein's Theory

$$H_1 : (X_1, Y_1) \cdots (X_n, Y_n) \overset{iid}{\sim} N\left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{pmatrix} \right)$$ ← Variance is wrong

Many other possibilities

## Neyman Pearson Paradigm

Hypothesis testing is a binary test

Does that provide sufficiently strong evidence to reject $H_0$, in favor of $H_1$?

Default assumption is that $H_0$ is true

A test statistic $T$ is a statistic computed from data that provides evidence against $H_0$ when extreme valued

1) How can we design a test statistic
2) How can you use it?

1) How can we design a test statistic?

For the brownian motion example:

$$\bar{R} = \frac{1}{n}(R_1 + R_2 + \cdots R_n) \qquad R_i = X_i^2 + Y_i^2 \quad \text{average distance in 30 sec interval}$$

$$\mathbb{E}(\bar{R}) = \mathbb{E}(R_i) = E(X_i^2) + \mathbb{E}(Y_i^2) = 4.46e{-}7$$

Extreme values of $\bar{R}$ reject null hypothesis

Can also compare against $(2.23e{-}7)\, \chi_2^2$
    Plot expected distribution against histogram

Hanging histogram plots $O_i - E_i$ for each histogram bin where $O_i$ is the observed count in bin $i$ and $E_i$ is the theoretical expected count

Test statistic $T = \sum_i (O_i - E_i)^2$
    Large $T$ indicates rejection of null hypothesis

You can stabilize the variance by plotting $\dfrac{O_i - E_i}{\sqrt{E_i}} = \dfrac{O_i - E_i}{\sqrt{n p_i}}$ so $\mathbb{E}\left[ \left( \dfrac{O_i - E_i}{\sqrt{n p_i}} \right)^2 \right] = 1$

Called hanging chi-gram

$T = \sum_i \dfrac{(O_i - E_i)^2}{E_i}$ is the pearson chi-squared statistic for goodness of fit

Another alternative is to consider $\sqrt{O_i} - \sqrt{E_i}$

$$\sqrt{O_i} - \sqrt{E_i} \simeq \frac{O_i - E_i}{2\sqrt{E_i}} \quad \text{via taylor expansion}$$

Tukey's hanging histogram

QQ plot

plots sorted values $R_1, \dots R_n$ against $\frac{1}{n}, \frac{2}{n} \dots \frac{n}{n}$ quantiles of their hypothesized distribution

Values deviating from $y=x$ provide evidence against hypothesized distribution

Let $R_{(1)} < \dots < R_{(n)}$ be the sorted values of $R_1 \dots R_n$. the maximum vertical deviation is

$$T = \max_{i=1}^{n} \left| R_{(i)} - F^{-1}\left(\frac{i}{n}\right) \right|$$

$F^{-1}$ is the quantile function $\leftarrow$ inverse cdf

more deviation is higher quantiles (not as scrunched together)
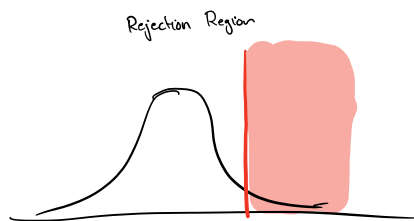
Stabilize Quantile Spacing

$$T = \max_{i=1}^{n} \left| F(R_{(i)}) - \frac{i}{n} \right|$$

Kolmogorov-Smirnov Statistic

2) What values of $T$ would allow us to reject $H_0$?

Null distribution of $T$ is the distribution of $T$ under the assumption the null hypothesis is true

We divide space of possible $T$ values into acceptance and rejection regions

Rejection Region



Type 1 Error: probability that we reject $H_0$ when $H_0$ is true

In Neyman Pearson paradigm, we choose rejection error s.t.

Type 1 error $\leq \alpha$

for a specified significance level $\alpha \in (0,1)$

Generally don't want to specify a significance level

p-value is the smallest significance level that we would reject $H_0$

probability that the null distribution assigns values $\geq t_{obs}$ $\leftarrow$ specifically for large $T$

2/9

Recap:

$H_0$ : Null hypothesis

$H_1$ : Alternative hypothesis

① Designing a good test statistic

② Choosing rejection region

Type I error: $\mathbb{P}_{H_0}[\text{reject } H_0] \leq \alpha$ ← significance level

What is the best choice of test statistic?

## Simple hypothesis

Def: A hypothesis is simple if it completely specifies its distribution

no unknown parameters

## Neyman Pearson Lemma

Simple $H_0$ vs. simple $H_1$

Def: The type II error for a simple alternative hypothesis $H_1$ is

$$B = \mathbb{P}_{H_1}[\text{accept } H_0]$$

Power of the test is $1-B$

$$1-B = \mathbb{P}_{H_1}[\text{reject } H_0]$$

Goal of hypothesis testing is to maximize power against $H_1$ while constraining type I error $\leq \alpha$

Example: Consider $X \in \{1,2,3,4,5\}$ Two hypotheses

| $X$ | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| $f_0(x)$ | 0.2 | 0.2 | 0.2 | 0.2 | 0.2 |
| $f_1(x)$ | 0 | 0.1 | 0.2 | 0.3 | 0.4 |

Let's say we want a test w/ significance level $\alpha = 0.4$

A hypothesis is defined by a rejection region $R \subseteq \{1,2,3,4,5\}$

If $X \in R$ : Reject $H_0$

If $X \notin R$ : Accept $H_0$

Type I Error: $\mathbb{P}_{H_0}[X \in R]$

Power $= \mathbb{P}_{H_1}[X \in R]$. To minimize power, set $R = \{4,5\}$. Then power is 0.7

More generally: Suppose the data is

$$X = (X_1, \dots X_n)$$

Takes values $x = (x_1, \dots x_n) \in X$ is finite

$H_0: X$ has joint pmf $f_0(x)$
$H_1: X$ has joint pmf $f_1(x)$

To define a test, we need to define a region $R \subseteq X$

If $X \in R$ : Reject $H_0$

If $X \notin R$ : Accept $H_0$

Want to find points of high $f_1$ and low $f_0$

Intuition: $R$ should contain points $X$ that have the smallest values of

$$L(X) = \frac{f_0(x)}{f_1(x)} \quad \leftarrow \text{ increase in Type I error per unit increase in power}$$

$L(X)$ is the likelihood ratio. Test rejects $H_0$ for small $L(X)$ is the likelihood ratio test

Analogous for continuous distributions

Neyman - Pearson Lemma: Let $H_0$ and $H_1$ be simple hypotheses. The significance level $\alpha \in (0,1)$. Suppose there is $c > 0$ s.t. the likelihood ratio test

rejects $H_0$ when $L(X) < c$

accepts $H_1$ when $L(X) \geq c$

has type I error of exactly $\alpha$. Then for any other test of type I error $\leq \alpha$, its power is at most the power of LRT.

## Proof : Consider the discrete case. Let $R = \{X: L(X) < c\}$ be the rejection region of LRT.

Then $R$ maximizes

$$\sum_{X \in R} (c f_1(x) - f_0(x)) \quad \text{among all subsets of } X$$

b/c if $X \in R$ then $\frac{f_0(x)}{f_1(x)} < c$ so $c f_1(x) - f_0(x) > 0$

if $X \notin R$ then $\frac{f_0(x)}{f_1(x)} \geq c$ so $c f_1 - f_0 \leq 0$

Consider another test w/ rejection region $R'$.

$$\sum_{X \in R} c f_1(x) - f_0(x) \geq \sum_{X \in R'} c f_1(x) - f_0(x)$$

$\Rightarrow$
$$c \left( \underbrace{\sum_{X \in R} f_1(x)}_{\substack{\text{power of} \\ \text{LRT}}} - \underbrace{\sum_{X \in R'} f_1(x)}_{\substack{\text{power of} \\ \text{other test}}} \right) \geq \underbrace{\sum_{X \in R} f_0(x)}_{\substack{\text{type I error} \\ \text{of LRT}}} - \underbrace{\sum_{X \in R'} f_0(x)}_{\substack{\text{type I error} \\ \text{of other test}}}$$

$$\geq 0$$

power of LRT $\geq$ power of other test

## Examples:

$X_1 \dots X_n$ are normal

$H_0 : X_1 \dots X_n \overset{iid}{\sim} N(0,1)$

$H_1 : X_1 \dots X_n \overset{iid}{\sim} N(\mu,1)$

for some fixed $\mu > 0$.

Let's apply N-P lemma

$$f_0(X) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{x_i^2}{2}} = \left(\frac{1}{\sqrt{2\pi}}\right)^n e^{-\left(\frac{x_1^2}{2} + \dots \frac{x_n^2}{2}\right)}$$

$$f_1(X) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{(X_i - \mu)^2}{2}} = \left(\frac{1}{\sqrt{2\pi}}\right)^n e^{-\left(\frac{(X_1-\mu)^2}{2} + \dots \frac{(X_n - \mu)^2}{2}\right)}$$

$$L(X) = \frac{f_0(X)}{f_1(X)} = \frac{e^{-\left(\frac{x_1^2}{2} + \dots \frac{x_n^2}{2}\right)}}{e^{-\left(\frac{(X_1-\mu)^2}{2} + \dots + \frac{(X_n - \mu)^2}{2}\right)}} = e^{-\left(\frac{x_1^2}{2} + \dots \frac{x_n^2}{2}\right) + \left(\frac{(X_1-\mu)^2}{2} + \dots \frac{(X_n - \mu)^2}{2}\right)}$$

$$= e^{\frac{n\mu^2}{2} - \mu(X_1 + \dots X_n)}$$

By the N-P lemma : Want to pick $c > 0$ s.t.

$$\mathbb{P}_{H_0}[L(X) < c] = \alpha$$
$\hookleftarrow$ type I error

Then LRT which

rejects the null hypothesis when $L(X) < c$

accepts the null hypothesis when $L(X) \geq c$

is the most powerful test

Two observations:

$$L(X) = e^{\frac{n\mu^2}{2} - \mu(X_1 + \dots + X_n)} =: f(\bar{x})$$

where $\bar{x} = \frac{(X_1 + \dots X_n)}{n}$ and $f(y) = e^{\frac{n\mu^2}{2} - n\mu y}$

This function is decreasing in $y$

so $L(X) = f(\bar{x}) < c \iff \bar{x} > f^{-1}(c)$

we want to pick $f^{-1}(c)$ s.t.

$$\mathbb{P}_{H_0}[\bar{x} > f^{-1}(c)] = \alpha$$

but we know under the null the $X_1 \dots X_n \overset{iid}{\sim} N(0,1)$, $\bar{x} \sim N(0, 1/n)$. So $f^{-1}(c)$ should be the $(1-\alpha)^{th}$ quantile

of $N(0, 1/n)$. i.e. $\frac{1}{\sqrt{n}} z(\alpha)$ where $z(\alpha)$ is the upper portion of $N(0,1)$

So the most powerful test is:

$$\text{Reject } H_0 \text{ if } \overline{X} > \frac{1}{\sqrt{n}} Z(\alpha)$$

$$\text{Accept } H_0 \text{ if } \overline{X} \leq \frac{1}{\sqrt{n}} Z(\alpha)$$

② No dependence of $M$. The most powerful test is this same test for every $M > 0$

## Lecture 7 (2/14/22)

Recap: Test $H_0 \sim f_0$

$\quad\quad\quad H_1 \sim f_1$

Most powerful test is to:

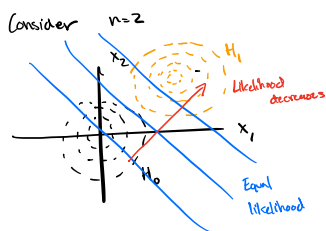Compute likelihood: $L(X) = \frac{f_0(X)}{f_1(X)}$

Pick a number $c > 0$ s.t. $P_{H_0}[L(X) < c] = \alpha$

Reject $H_0$ when $L(X) < c$, accept $H_0$ when $L(X) \geq c$

### Example:

$$H_0 : X_1 \ldots X_n \overset{iid}{\sim} N(0,1)$$

$$H_1 : X_1 \ldots X_n \overset{iid}{\sim} N(M,1) \quad \mu > 0 \text{ fixed and known}$$

Consider $n = 2$



Most powerful test is equivalent to

Compute $\overline{X} = \frac{X_1 + \cdots + X_n}{n}$

Pick $c \in \mathbb{R}$ s.t. $P_{H_0}[\overline{X} > c] = \alpha$

$$c = \frac{1}{\sqrt{n}} Z(\alpha)$$

Reject $H_0$ if $\overline{X} > c$, accept $H_0$ if $\overline{X} \leq c$

EX: $H_0 : X_1 \ldots X_n \overset{iid}{\sim} \text{Bernoulli}(\frac{1}{2})$

$\quad\quad H_1 : X_1 \ldots X_n \overset{iid}{\sim} \text{Bernoulli}(p)$

$$L(X) = \frac{f_0(X)}{f_1(X)} = \left(\frac{1}{2(1-p)}\right)^n \cdot \left(\frac{1-p}{p}\right)^{X_1 + \cdots X_n}$$

Observation: $L(X)$ only depends on $X_1 \ldots X_n$ via their sum

For $p > \frac{1}{2}$ $\quad \frac{1-p}{p} < 1$. So $L(X)$ is decreasing in $X_1 + \cdots X_n$

Likelihood procedure is equivalently:

- Compute $S = X_1 + \cdots X_n$

- Under $H_0$ $S \sim \text{Binomial}(n, \frac{1}{2})$

  pick $c$ as the upper $\alpha$ point $b_n(\alpha)$ of binomial $(n, \frac{1}{2})$

- Reject $H_0$ when $S > c$, accept when $S \leq c$

## Composite Hypothesis and pivotal statistic

A hypothesis (either $H_0$ or $H_1$) that is not simple is composite

Ex. Let $n = 250$ be the number of students in STAT 242

Does 242 improve student knowledge of Statistics?

Suppose each student takes

  — diagnostic exam at start

  — final exam at end

Let $X_i$ be the difference (final - diagnostic) for student $i$.

Two formulations:

$$H_0 : X_1, \ldots X_n \overset{iid}{\sim} N(0, \sigma^2)$$
$$H_1 : X_1, \ldots X_n \overset{iid}{\sim} N(\mu, \sigma^2) \quad \text{for some } \mu, \sigma^2 > 0$$

Both $H_0$ and $H_1$ are composite

$H_0$: $X_1, \ldots X_n$ are iid from some pdf $f$ w/ median $0$

$H_1$: $X_1, \ldots X_n$ are iid from some pdf $f$ w/ median $> 0$

Both $H_0$ and $H_1$ are composite

When testing composite null $H_0$:

- We want to ensure probability of Type I error is $\leq \alpha$ for <u>every</u> possible data distribution described $H_0$

- Simplify design of the test by picking test statistic $T$ whose distribution is the same under very distribution in $H_0$

  $\uparrow$ pivotal or distribution free test statistic

When testing composite Alternative $H_1$:

- Oftentimes no single test that maximizes power against all possible Alternatives
- We often design a test to balance the power against different alternatives

One-sample t-test

Setup: $X_1, \ldots X_n \overset{iid}{\sim} N(\mu, \sigma^2)$ where both $\mu, \sigma^2$ are unknown

$H_0 : \mu = 0 \qquad H_1 : \mu > 0$

Q: What if $\sigma^2$ is known

In this case, we can first standardize our data: $Z_i = \frac{X_i}{\sigma}$

Then $Z_i \sim N(\frac{\mu}{\sigma}, 1)$. The above is equivalent to testing $H_0 : \frac{\mu}{\sigma} = 0$ vs. $H_1 : \frac{\mu}{\sigma} > 0$ based on $Z_1, \ldots Z_n$

The Neyman Pearson Lemma implies Most powerful level-$\alpha$ test is to reject if $\bar{Z} = \frac{\bar{X}}{\sigma} > \frac{1}{\sqrt{n}} Z(\alpha)$

Idea: If sigma is unknown, let's estimate it from the data

Let $S^2 = \frac{1}{n-1} \left[ (X_1 - \bar{X})^2 + \cdots + (X_n - \bar{X})^2 \right]$

Use instead

$$T = \frac{\sqrt{n} \, \bar{X}}{S} \quad \leftarrow \text{One sample t-statistic}$$

Pivotal Statistic!

Thm: Let $X_1, \ldots X_n \overset{iid}{\sim} N(\mu, \sigma^2)$. Let $\bar{X}$ and $S^2$ be the sample mean and variance. Then

1) $S^2$ is independent of $\bar{X}$

2) $S^2 \sim \frac{\sigma^2}{n-1} \cdot \chi^2_{n-1}$

Proof: WLOG $\mu = 0$

1) It suffices to show $\bar{X}$ is independent of $(X_1 - \bar{X}, \ldots X_n - \bar{X})$

$\uparrow$ $S_n$ is a simple function of $(X_1 - \bar{X}, \ldots, X_n - \bar{X})$

$(\bar{X}, X_1 - \bar{X}, \dots X_n - \bar{X})$ are all linear combinations of $X_1 \dots X_n \overset{iid}{\sim} N(u, \sigma^2)$

$\Rightarrow (\bar{X}, X_1 - \bar{X}, \dots X_n - \bar{X})$ is multivariate

So to show $\bar{X}$ is independent of $(X_1 - \bar{X}, \dots X_n - \bar{X})$ its enough $Cov[\bar{X}, X_1 - \bar{X}] = 0$

We have:

$$Cov[\bar{X}, X_i] = Cov\left[\frac{1}{n}\sum_{j=1}^{n} X_j, X_i\right] = \frac{1}{n}\sum_{j=1}^{n} Cov[X_j, X_i] = \frac{1}{n} Cov[X_j, X_i] = \frac{1}{n} Var[X_i] = \frac{\sigma^2}{n}$$

$\underbrace{\qquad}_{=0 \text{ for } i \neq j}$

$$Cov[\bar{X}, X_i] = Var[\bar{X}] = \frac{\sigma^2}{n}$$

$$\Rightarrow Cov[\bar{X}, X_i - \bar{X}] = \frac{\sigma^2}{n} - \frac{\sigma^2}{n} = 0 \qquad \text{so} \quad \bar{X} \text{ and } S^2 \text{ are independent}$$

2) For the distribution of $S^2$, write $Z_i = \frac{X_i}{\sigma} \sim N(0,1)$

$$S^2 = \frac{1}{n-1}\left[(X_i - \bar{X})^2 + \dots + (X_n - \bar{X})^2\right]$$

$$\Rightarrow \underbrace{\frac{(n-1)}{\sigma^2} S^2}_{U} = (Z_1 - \bar{Z})^2 + \dots + (Z_n - \bar{Z})^2$$

$$= (Z_1^2 - 2Z_1\bar{Z} + \bar{Z}^2) + \dots + (Z_n^2 - 2Z_n\bar{Z} + \bar{Z}^2)$$

$$= (Z_1^2 + \dots Z_n^2) - 2\bar{Z}\underbrace{(Z_1 + \dots + Z_n)}_{n\bar{Z}} + n\bar{Z}^2$$

$$= \underbrace{(Z_1^2 + \dots Z_n^2)}_{W} - \underbrace{(\sqrt{n}\,\bar{Z})^2}_{V}$$

This shows that $W = U + V$

we know $\bar{X}$ is independent of $S^2$

So $V$ is independent of $U$

Here: $W \sim \chi^2_n$ , $V \sim \chi^2_1$ b/c $\sqrt{n}\bar{Z} \sim N(0,1)$

$\Rightarrow U \sim \chi^2_{n-1}$ .

Thus, $\frac{(n-1)}{\sigma^2} S^2 \sim \chi^2_{n-1} \Rightarrow S^2 \sim \frac{\sigma^2}{n-1} \cdot \chi^2_{n-1}$

Returning to $T = \frac{\sqrt{n}\bar{X}}{S} = \frac{\sqrt{n}\, \bar{X}/\sigma}{S/\sigma}$

· Under $H_0$, $\frac{\sqrt{n}\bar{X}}{\sigma} \sim N(0,1)$ b/c $X_i \sim N(0, \sigma^2)$

· $\frac{S^2}{\sigma^2} \sim \frac{1}{n-1}\chi^2_{n-1}$

· And $\frac{\sqrt{n}\bar{X}}{\sigma}$ is independent of $S/\sigma^2$

**Def:** If $Z \sim N(0,1)$, $U \sim \chi_n^2$ and $Z$ and $U$ are independent, then the distribution of $\dfrac{Z}{\sqrt{\frac{1}{n}U}}$ is called the t-distribution w/ $n$ degrees of freedom

**Remark:** The preceeding thm explains why we use $\frac{1}{n-1}$ to define $S^2$:

$$\mathbb{E}[S^2] = \mathbb{E}\left[\frac{\sigma^2}{n-1}\chi_{n-1}^2\right], \text{ hence } \mathbb{E}[\chi_{n-1}^2] = n-1 \text{ so } \mathbb{E}[S^2] = \sigma^2 \text{ so } S^2 \text{ is unbiased for } \sigma^2$$

## Lecture 8 (2/16/22)

**Recap:** Data $X_1 \dots X_n$ iid samples. Test whether distribution of $X_i$'s is "centered around 0"

Formulation in Normal Setting:

$$H_0 : X_1 \dots X_n \overset{iid}{\sim} N(0,\sigma^2)$$
$$H_1 : X_1 \dots X_n \overset{iid}{\sim} N(\mu,\sigma^2), \quad \mu > 0$$

One-sample t-statistic

$$T = \frac{\sqrt{n}\,\bar{X}}{S} \text{ where } \bar{X} \text{ is the average and } S^2 = \frac{1}{n-1}\left((X_1-\bar{X})^2 + \dots + (X_n-\bar{X})^2\right)$$

What is distribution of $t$ under $H_0$?

- $\dfrac{\sqrt{n}\,\bar{X}}{\sigma} \sim N(0,1)$

- $\dfrac{S^2}{\sigma^2} \sim \dfrac{1}{n-1}\chi_{n-1}^2$

- These two statistics are independent

So $T \sim t_{n-1}$, the t-distribution with $n-1$ degrees of freedom

**T-test:** Reject $H_0$ when $T > t_{n-1}$

Formulation in a non-parametric: $X_1 \dots X_n \overset{iid}{\sim} f$

$H_0$ : Median of distribution is 0

$H_1$ : Median is greater than 0

Consider the sign statistic $S = \sum_{i=1}^{n} \mathbb{1}\{X_i > 0\}$ ← Indicator

$$\mathbb{1}\{X_i > 0\} = \begin{cases} 1 & \text{if } X_i > 0 \\ 0 & \text{if } X_i \leq 0 \end{cases}$$

Under the distribution $f$ described by $H_0$, $S \sim \text{Binomial}(n, \frac{1}{2})$

Let $b_n(\alpha)$ be upper-$\alpha$ point of binomial $(n, \frac{1}{2})$. Then the test that rejects for $S > b_n(\alpha)$ is the sign test

**Remark:** For large $n$ we can approximate this test via a normal approximation

- By the CLT: $\dfrac{\frac{S}{n} - \frac{1}{2}}{\sqrt{\frac{1}{4}}} \cdot \sqrt{n} \longrightarrow N(0,1)$

$$\frac{S}{n} - \frac{1}{2} \approx N\left(0, \frac{1}{4n}\right) \Rightarrow S \approx N\left(\frac{n}{2}, \frac{n}{4}\right)$$

- The $b_n(\alpha) =$ upper alpha point of $N\left(\frac{n}{2}, \frac{n}{4}\right)$

$$= \frac{n}{2} + \sqrt{\frac{n}{4}} \cdot z(\alpha)$$

↑ upper alpha point

Approximate sign test by rejecting $H_0$

$$S > \frac{n}{2} + \sqrt{n/4} \ z(\alpha) \iff \sqrt{4n} \left( \frac{S}{n} - \frac{1}{2} \right) > z(\alpha)$$

The type I error is not **exactly** $\alpha$ but under $H_0$

$$P[\text{Reject } H_0] \xrightarrow{n \to \infty} P[N(0,1) > z(\alpha)] = \alpha$$

type I error approaches $\alpha$


# Two-sample Tests

Ex. $X_1 \ldots X_n$ are differences in exam scores for $n = 250$ students

Q: What if the exams are not equally difficult?

A: Take a separate group of $m = 100$ students not in S&DS 242

Let $Y_1 \ldots Y_m$ be the differences in scores

Informally, test whether the distribution of $X_i$'s is larger than the $Y_i$'s

Normal Model Function:

Suppose $X_1 \ldots X_n \overset{iid}{\sim} N(\mu_x, \sigma^2)$
$Y_1 \ldots Y_n \overset{iid}{\sim} N(\mu_y, \sigma^2)$

$H_0: \mu_x = \mu_y$

$H_1: \mu_x > \mu_y$

More non-parametric formulation

Suppose $X_1 \ldots X_n \overset{iid}{\sim} f$, $Y_1 \ldots Y_n \overset{iid}{\sim} g$

$H_0: f = g$

$H_1: f$ stochastically dominates $g$

↰ IF $x \sim f$ and $y \sim g$ then for any $x \in \mathbb{R}$ $P[X \geq x] \geq P[Y \geq x]$

Two-sample $t$-test

Intuition: Look at difference of sample means $\bar{X} - \bar{Y}$

Reject $H_0: \mu_x = \mu$ in favor $H_1: \mu_x > \mu_y$

if $\bar{X} - \bar{Y}$ is "large enough"

What is the null distribution of $\bar{X} - \bar{Y}$

$$\bar{X} = \frac{X_1 + \ldots X_n}{n} \sim N\left(\mu_x, \frac{\sigma^2}{n}\right)$$

$$\bar{Y} = \frac{Y_1 + \ldots Y_m}{m} \sim N\left(\mu_y, \frac{\sigma^2}{m}\right)$$

$$\bar{X} - \bar{Y} \sim N\left(\mu_x - \mu_y, \ \frac{\sigma^2}{n} + \frac{\sigma^2}{m}\right)$$

Under $H_0: \mu_x = \mu_y$

$$\bar{X} - \bar{Y} \sim N\left(0, \ \sigma^2\left(\frac{1}{n} + \frac{1}{m}\right)\right)$$

$$\frac{\bar{X} - \bar{Y}}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim N(0,1)$$

If $\sigma^2$ were known, then reject $H_0$ if $\dfrac{\bar{X} - \bar{Y}}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}} > z(\alpha)$

When $\sigma^2$ is unknown, estimate $\sigma^2$ from the data

$$S^2_{pooled} = \frac{1}{m+n-2} \left( \sum_{i=1}^{n} (X_i - \bar{X})^2 + \sum_{j=1}^{m} (Y_j - \bar{Y})^2 \right)$$

Test Statistic

$$T = \frac{\bar{X} - \bar{Y}}{S_{pooled} \sqrt{\frac{1}{n} + \frac{1}{m}}} \qquad \longleftarrow \text{pooled two-sample } t\text{-statistic}$$

Statistic is pivotal under $H_0$: Let $\mu = \mu_x = \mu_y$

Write $X_i = \mu + \sigma Z_i$ where $Z_i$'s and $W_j$'s are now $N(0,1)$

$\qquad Y_j = \mu + \sigma W_j$

Substituting into $T$: $\mu$ cancels from $\bar{X} - \bar{Y}$

$\qquad\qquad\qquad\quad \sigma$ cancels from $\dfrac{(\bar{X} - \bar{Y})}{S_{pooled}}$

What is the distribution of $T$?

$\left. \begin{array}{l} \cdot \dfrac{\bar{X}}{\sigma} \sim N\left(\dfrac{\mu}{\sigma}, \dfrac{1}{n}\right) \\[3mm] \cdot \dfrac{\bar{Y}}{\sigma} \sim N\left(\dfrac{\mu}{\sigma}, \dfrac{1}{m}\right) \end{array} \right\} \quad \dfrac{\bar{X} - \bar{Y}}{\sigma} \sim N\left(0, \dfrac{1}{n} + \dfrac{1}{m}\right)$

$\cdot \dfrac{1}{\sigma^2} \sum_{i=1}^{n} (X_i - \bar{X})^2 \sim \chi^2_{n-1}$

$\cdot \dfrac{1}{\sigma^2} \sum_{j=1}^{m} (Y_j - \bar{Y})^2 \sim \chi^2_{m-1}$

$\qquad\qquad\qquad\qquad\qquad\qquad \overbrace{\chi^2_{n-1}} \qquad\qquad \overbrace{\chi^2_{m-1}}$

$\Rightarrow \dfrac{S_{pooled}^2}{\sigma^2} = \dfrac{1}{m+n-2} \left( \dfrac{1}{\sigma^2} \sum_{i=1}^{n} (X_i - \bar{X})^2 + \dfrac{1}{\sigma^2} \sum_{j=1}^{m} (Y_j - \bar{Y})^2 \right)$

$$\sim \dfrac{1}{m+n-2} \cdot \chi^2_{m+n-2}$$

This implies $\mathbb{E}\left[ S_{pooled}^2 \right] = \sigma^2$, so this is unbiased

$\text{So } T = \dfrac{\bar{X} - \bar{Y}}{S_{pooled} \sqrt{\frac{1}{n} + \frac{1}{m}}} = \dfrac{\frac{1}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}} (\bar{X} - \bar{Y})}{S_{pooled} / \sigma} = \dfrac{N(0,1)}{\sqrt{\frac{1}{m+n-2} \chi^2_{m+n-2}}}$

$\qquad\qquad\qquad \sim t_{m+n-2}, \text{ the } t\text{-distribution w/ } m+n-2 \text{ degrees of freedom}$

The two sampled $t$-test rejects $H_0 = \mu_x = \mu_y$ in favor of $H_1 : \mu_x > \mu_y$ when $T > t_{m+n-2} (\alpha)$

Remark: Assumed $X_i$ and $Y_j$ have the same variance

$\qquad$ It's common to see applications where variances are not identical

$\qquad\qquad$ Suppose $X_1 \dots X_n \overset{iid}{\sim} N(\mu_x, \sigma_x^2)$

$\qquad\qquad\qquad\qquad Y_1 \dots Y_n \overset{iid}{\sim} N(\mu_y, \sigma_y^2)$

$\qquad\qquad$ Test $H_0 : \mu_x = \mu_y$ vs. $H_1 : \mu_x > \mu_y$

Idea: Look at $\bar{X} - \bar{Y}$. Under $H_0$:

$\qquad\qquad \bar{X} - \bar{Y} \sim N\left(0, \dfrac{\sigma_x^2}{n} + \dfrac{\sigma_y^2}{n}\right)$

$\qquad$ Estimate $\sigma_x^2$ by sample var. $S_x^2 = \dfrac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X})^2$

$\qquad\qquad\qquad \sigma_y^2 \qquad \dots \qquad\qquad S_y^2$

Define $T_{welch}$ : $\dfrac{\bar{X} - \bar{Y}}{\sqrt{\frac{1}{n} \cdot S_x^2 + \frac{1}{m} S_y^2}} \qquad \longleftarrow \text{Not exactly pivotal under } H_0 \text{ and not exactly } t\text{-distributed}$

## Mann-Whitney-Wilcoxon Rank-Sum Test

$X_1 \ldots X_n \overset{iid}{\sim} f$ and $Y_1 \ldots Y_m \overset{iid}{\sim} g$

Test $H_0 : f = g$ vs. $H_1 : f$ stochastically dominates $g$

The M-W-W rank-sum statistic $T$ is defined as:

① Pool all $X_i$'s and $Y_j$'s and rank them

   Let smallest have rank 1

   Largest have rank $m+n$

② Define $T$ as the sum of ranks of only the $Y_j$'s

## Lecture 9 (2/21/22) — Permutation Tests

$X_1 \ldots X_n \overset{iid}{\sim} f$    $Y_1 \ldots Y_m \overset{iid}{\sim} g$

$H_0 : f = g$        $H_1 : f$ stochastically dominates $g$

   ↑ one-sided alternative
   two-sided : $f \neq g$

Mann-Whitney-Wilcoxon Rank Sum

① Pool together $X_1 \ldots X_n$, $Y_1 \ldots Y_m$ and rank them smallest to largest

② $T = $ sum of ranks $Y_1 \ldots Y_m$

Under $H_0$, $X_1 \ldots X_n$, $Y_1 \ldots Y_m$ are all iid

By symmetry, equally likely for $Y_1 \ldots Y_m$ to be any set of $m$ ranks among $\{1, 2 \ldots m+n\}$

**Thm:** Under $H_0 : f = g$

a) $\mathbb{E}[T] = \dfrac{m(m+n+1)}{2}$ and $Var[T] = \dfrac{mn(m+n+1)}{12}$

b) $\dfrac{T - \mathbb{E}[T]}{\sqrt{Var[T]}} \longrightarrow N(0,1)$ in distribution as $n, m \to \infty$

Proof of part a

   Let $N = m+n$. Define

   $$I_k = \begin{cases} 1 & \text{rank } k \text{ observation is a } Y \\ 0 & \text{rank } k \text{ observation is a } X \end{cases}$$

   $$T = \sum_{k=1}^{N} k I_k$$

   Have the values $k$ when $I_k = 1$ is a simple random sample from $\{1 - N\}$ under $H_0$

   $$\mathbb{E}[I_k] = \mathbb{P}[I_k = 1] = \frac{m}{N}$$

   $$\mathbb{E}[T] = \sum_{k=1}^{N} k \cdot \mathbb{E}[I_k] = \sum_{k=1}^{N} k \cdot \frac{m}{n} = \frac{N(N+1)}{2} \cdot \frac{m}{N} = \frac{m(m+n+1)}{2}$$

   $$Var[T] = \mathbb{E}[T^2] - \left(\mathbb{E}[T]\right)^2$$

   $$\mathbb{E}[T^2] = \mathbb{E}\left[\left(\sum_{k=1}^{N} k I_k\right)\left(\sum_{k=1}^{N} j I_j\right)\right]$$

   $$= \sum_{k=1}^{N} jk \, \mathbb{E}[I_k]\mathbb{E}[I_j]$$

   $$= \sum_{k=1}^{N} k^2 \underbrace{\mathbb{E}[I_k^2]}_{\frac{m}{N}} + 2 \sum_{j \neq k} jk \underbrace{\mathbb{E}[I_k I_j]}_{\frac{m}{n} \cdot \frac{m-1}{N-1}}$$

Apply: $\sum_{k=1}^{N} k^2 = \frac{N(N+1)(2N+1)}{6}$

$2 \sum_{j \neq k} jk = \left( \sum_{k=1}^{N} k \right)^2 - \sum_{k=1}^{N} k^2$

$\qquad = \left( \frac{N(N+1)}{2} \right)^2 - \frac{N(N+1)(2N+1)}{6}$

$\Rightarrow \mathbb{E}[T^2] = \frac{m}{N} \cdot \frac{N(N+1)(2N+1)}{6} + \frac{m}{N} \cdot \frac{m-1}{N-1} \left[ \left( \frac{N(N+1)}{2} \right)^2 - \frac{N(N+1)(2N+1)}{6} \right]$

$Var[T] = \mathbb{E}[T^2] - \mathbb{E}(T)^2$

$\qquad = \frac{mn(m+n+1)}{12}$

To preform an asymptotic test

For large $n, m$ distribution of $T$ is $\approx N \left( \frac{m(m+n+1)}{2}, \frac{mn(m+n+1)}{12} \right)$

To test against $H_1: f$ stochastically dominates $g$

The $T$ takes smaller value under $H_1$

Reject $H_0$ when $T < \frac{m(m+n+1)}{2} - \sqrt{\frac{mn(m+n+1)}{12}} \cdot z(\alpha)$

To test against $H_1: f \neq g$

Reject when $\left| T - \frac{m(m+n+1)}{2} \right| > \sqrt{\frac{mn(m+n+1)}{12}} \cdot z\left(\frac{\alpha}{2}\right)$

Permutation Tests

$X_1 \dots X_n \overset{iid}{\sim} f$, $Y_1 \dots Y_m \sim g$

Test $H_0: f = g$ vs. $H_1: f \neq g$

Symmetry Underlying rank-sum test

Pool all observations as $\{Z_1 \dots Z_{n+m}\}$

ignoring their ordering

Given these pooled values $\{Z_1 \dots Z_{n+m}\}$ under $H_0$, its equally likely for $X_1 \dots X_n \ Y_1 \dots Y_n$ to be any permutation of these values

Idea: For any test statistic $T(X_1, \dots X_n, Y_1 \dots Y_m)$

The permutation null distribution of $T$ is its distribution upon:

1) Take a random permutation $X_1^* \dots X_n^*, Y_1^* \dots Y_m^*$

2) Compute $T(X_1^*, \dots X_n^*, Y_1^* \dots Y_m^*)$

$\quad$ uniform distribution on all $(m+n)!$ possible possibilities

Can approximate this by simulation:

· Sample $B$ uniformly random permutations

· Look at the distribution of the $B$ values

A test of $H_0$ using $T$ and this permutation null distribution is a permutation test

Permutation test is a conditional test

pivotal on $H_0$ <u>conditional</u> on $\{Z_1 \dots Z_{m+n}\}$

The test guarantees

$$\mathbb{P}_{H_0}\left[\text{reject } H_0 \mid \{Z_1 \dots Z_{m+n}\}\right] \leq \alpha \quad \text{for any set of } \{Z_1 \dots Z_{m+n}\}$$

This also holds unconditionally

$$\mathbb{P}_{H_0}\left[\text{reject } H_0\right] \leq \alpha$$

Example: Suppose $X_1 \dots X_n$ are objects and $Y_1 \dots Y_n$ are a second sample

Suppose we have definite distance $d(x, x')$ on $X$

Here are 3 different statistics

· Average Distance Statistic

$$T_1 = 2 \frac{1}{mn} \sum_{i=1}^{n} \sum_{j=1}^{m} d(x_i, y_j) - \frac{1}{\binom{n}{2}} \sum_{\substack{i=1 \\ j=1}}^{n} d(x_i, x_j) - \frac{1}{\binom{m}{2}} \sum_{\substack{j=1 \\ i=1}}^{m} d(y_i, y_j)$$

· $k$-nearest neighbors

How many of its $k$-nearest numbers from its sample

$T_2 = $ avg across all samples

· Freedmann – Rafely minimum spanning tree statistics

· Constructing a minimum spanning tree from the pooled observations

· Remove edges connected $X$ with $Y$

$T_3 = $ number of connected components

$T_1, T_2, T_3$ are not necessarily pivotal under $H_0$
Their full null distributions are complicated


<u>Fisher's exact test</u>

Randomly shuffle atrames in table


Lecture 10 (2/23/22) - Effect size, power, and experimental design

Steps of a Scientific Study

1. Identify and formulate question of interest
2. Design experiment to answer this question
3. Visualize and explore collected data
4. Apply statistical procedure

Questions
- Predict in advance whether the study will succeed
- Predict size of study
- Why does experimental design influence our ability to ID this effect

Case Study: Stanford Peer Grading Experiment
   Divide course into peer grading and control groups

Predict the power
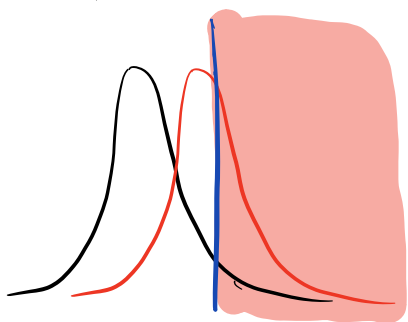   Rejecting $H_0$ at desired level of significance

Power in the one-sample $Z$-test
   $X_1 \dots X_n \overset{iid}{\sim} N(\mu, \sigma^2)$

   $H_0: \mu = 0$   vs.   $H_1: \mu > 0$   assume $\sigma^2$ is known

Neyman-Pearson lemma tells us that the most powerful test rejects $H_0$ for large values of $\overline{X}$ $\leftarrow$ $Z$ test

   Null: $\overline{X} \sim N(0, \frac{\sigma^2}{n})$      Alternative: $\overline{X} \sim N(\mu, \frac{\sigma^2}{n})$



Analytically: $Z$-test rejects $H_0$ when $\frac{\sqrt{n}}{\sigma} \overline{X} > Z(\alpha)$

   this ensures type 1 error $= \alpha$

Under $H_1$: $\overline{X} \sim N(\mu, \sigma^2/n)$

   $\overline{X} = \mu + \frac{\sigma}{\sqrt{n}} Z$ where $Z \sim N(0,1)$

   Power: $\mathbb{P}_{H_1}\left[\frac{\sqrt{n}}{\sigma} \overline{X} > Z(\alpha)\right]$

   $= \mathbb{P}\left[\frac{\sqrt{n}}{\sigma}\left(\mu + \frac{\sigma}{\sqrt{n}} Z\right) > Z(\alpha)\right]$

   $= \mathbb{P}\left[Z > Z(\alpha) - \sqrt{n}\frac{\mu}{\sigma}\right]$

   $= \Phi\left(\sqrt{n}\frac{\mu}{\sigma} - Z(\alpha)\right)$

Power in comparing two samples
   Consider $X_1 \dots X_n \sim N(\mu_x, \sigma^2)$   $Y_1 \dots Y_m \sim N(\mu_y, \sigma^2)$

   $H_0: \mu_x = \mu_y$      $H_1: \mu_x > \mu_y$

   Pooled two-sample

$$T = \frac{\bar{X} - \bar{Y}}{S_{pooled}\sqrt{\frac{1}{n} + \frac{1}{m}}}$$

$$\text{Power} = \mathbb{P}_{H_1}\left[T > t_{m+n-2}(\alpha)\right]$$

For large $n, m$

$$\approx \mathbb{P}_{H_1}\left[\frac{\bar{X} - \bar{Y}}{\sigma\sqrt{\frac{1}{n} + \frac{1}{m}}} > z(\alpha)\right]$$

Under $H_1$: $\bar{X} - \bar{Y} \sim N\left(\mu_X - \mu_Y, \sigma^2\left(\frac{1}{n} + \frac{1}{m}\right)\right)$

$$\bar{X} - \bar{Y} = \mu_X - \mu_Y + \sigma\sqrt{\frac{1}{n} + \frac{1}{m}} \cdot Z \quad \text{where } Z \sim N(0,1)$$

$$\text{Power} = \mathbb{P}\left[Z + \underbrace{\frac{1}{\sqrt{\frac{1}{n} + \frac{1}{m}}} \cdot \frac{\mu_X - \mu_Y}{\sigma}}_{d} > z(\alpha)\right] = \phi(d - z(\alpha))$$

Power is increasing in

$$d = \frac{1}{\sqrt{\frac{1}{n} + \frac{1}{m}}} \cdot \frac{\mu_X - \mu_Y}{\sigma} \quad \Rightarrow \quad \text{decreasing in } \frac{1}{n} + \frac{1}{m} \quad \text{so} \quad \frac{1}{n} + \frac{1}{m} \text{ is minimized by}$$

$$n = m = \frac{N}{2}$$

Predicting typical p-value

$$\text{p-value} = \mathbb{P}_{H_0}\left[T > t_{obs}\right] \approx 1 - \phi(t_{obs})$$

Under $H_1$

$$t_{obs} \approx \frac{\bar{X} - \bar{Y}}{\sigma\sqrt{\frac{1}{n} + \frac{1}{m}}} = Z + \frac{1}{\sqrt{\frac{1}{n} + \frac{1}{m}}} \cdot \frac{\mu_X - \mu_Y}{\sigma}$$

Median value of $t_{obs}$ under $H_1$ is roughly $d$.

For $d = 0.95$, this p-value is $0.17$

Paired Design

split course into 2 units and swap groups between units

Consider paired differences

$$D_i = X_i - Y_i$$

If $X_i, Y_i$ is bivariate normal then $D_i$ has a normal distribution

$$\mathbb{E}(D_i) = \mu_X - \mu_Y$$

$$\text{Var}(D_i) = \text{Cov}(X_i - Y_i, X_i - Y_i) = 2\sigma^2(1 - \rho)$$

Reduces problem to one-sample testing problem

Level $\alpha$ t-test

$$\frac{\sqrt{n}}{s} \bar{D} > t_{n-1}(\alpha)$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^{n} (D_i - \bar{D})^2 \quad \text{is sample variance}$$

Difference between paired and unpaired is a

$$\frac{1}{1-p} \quad \text{term}$$

$1-p$ is the relative efficiency of the unpaired design to the paired design

paired test with $n$ pairs has the same power as an unpaired design with sample size $\frac{n}{1-p}$ per group
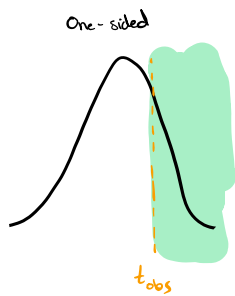
Confounding variables

- systematic biases
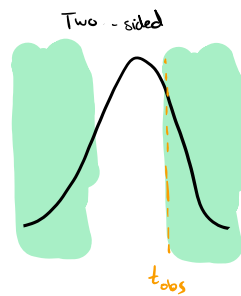- Inflate variance

Lecture 11: Testing Multiple Hypotheses

If I test $n$ null hypotheses at level $\alpha$, all of which are true, then on average I'll falsely reject $\alpha n$ of them

Most statistical multiple-testing procedures take p-values as inputs

p-value is the smallest significance value at which the test rejects the null hypothesis

One-sided



$t_{obs}$

p-val is 0.036   (smallest value that rejects null)

Two-sided



$t_{obs}$

p-value is 0.072

For a one-sided test with continuous test statistic $T$, we reject $H_0$ when $T$ exceeds upper-$\alpha$ point of its null distribution

$$P = \mathbb{P}_{H_0}[T > t_{obs}] = 1 - F(t_{obs})$$

For a two-sided test, we reject $H_0$ when $T$ is larger than upper-$\frac{\alpha}{2}$ or less than lower-$\frac{\alpha}{2}$

$$P = 2 \cdot \min\left(F(t_{obs}), 1 - F(t_{obs})\right)$$

p-value is the probability under $H_0$ of observing a value of $T$ that is more extreme then $t_{obs}$.

p-value can be considered as a test statistic itself

reject if $P \leq \alpha$

$P$ is uniform $(0,1)$ under $H_0$

Bonferroni Method:

For $n$ different null hypotheses, reject at $\alpha/n$ instead of $\alpha$

Justification:

$$\mathbb{P}[\text{reject any null hypothesis}] = \mathbb{P}\left[\{\text{reject } H_0^{(1)}\} \cup \ldots \cup \{\text{reject } H_0^{(n)}\}\right]$$

$$\leq \mathbb{P}\left[\text{reject } H_0^{(1)}\right] + \ldots + \mathbb{P}\left[\text{reject } H_0^{(n)}\right] \leq \frac{\alpha}{n} + \ldots + \frac{\alpha}{n} = \alpha$$

Family Wise Error Rate

$n$ null hypotheses, $n_0$ are true nulls

$$\text{FWER} = \mathbb{P}[\text{reject any true null hypothesis}]$$

Controls FWER at level $\alpha$ guarantees FWER $\leq \alpha$

Bonferroni Method controls FWER at $\alpha$

False Discovery Rate

$$\text{False Discovery Proportion} = \frac{\text{number of true null hypotheses rejected}}{\text{number of total null hypotheses rejected}} = \frac{V}{R} \quad \leftarrow 0 \text{ when } V=0 \text{ and } R=0$$

Controls FDR at level $\alpha$ if FDR $\leq \alpha$

Controlling FWER is more appropriate if the consequences of a single Type I error is high
result will be interpreted as truth

Controlling FDR is more appropriate if the test IDs candidate discoveries for further study
if false discoveries are acceptable as long most of the discoveries are correct

Estimating $\text{FDP} = \frac{V}{R}$ but we don't know $V$

We can estimate $V$ since p-values are uniformly distributed $(0,1)$.

For a rejection value of $t$, we can expect $t n_0$ of the true nulls to have $p \leq t$

$$V \approx t n_0$$

We don't know $n_0$ so

$$\widehat{\text{FDP}} = \frac{t n}{R} \quad \leftarrow \text{estimate}$$

Control $\widehat{\text{FDP}} = \frac{t n}{R(t)} \leq \alpha$

$\quad \uparrow$ number of rejections

goal is to find maximum $t$ that satisfies this relation

For $r$ rejections, $t = P_{(r)}$, the $r^{th}$ smallest p-value and find largest $r$ s.t.

$$\frac{P_{(r)}}{r} \cdot n \leq \alpha \iff P_{(r)} \leq \frac{\alpha r}{n} \quad \leftarrow \text{Benjamini-Hochberg}$$

1. Sort $n$ p-values from smallest to largest

2. Find largest $r$ s.t. $P_{(r)} \leq \frac{\alpha r}{n}$

3. Reject null hypothesis corresponding to $P_{(1)} \ldots P_{(r)}$

Notice: $P_{(1)}$ is compared to $\frac{\alpha}{n}$, $P_{(2)} = \frac{2\alpha}{n}$, etc.

## 3/2 : Parametric models and the method-of-moments

Def: A parametric model is a family of distributions indexed by a small number of unknown parameters

e.g. $N(u, \sigma^2)$

Notation: Vector of parameters $\Theta \in \mathbb{R}^k$

PDF / PMF as $f(x|\Theta)$

<span style="color:orange">↑ not conditional</span>

The set of allowable parameters is the parameter space

How to choose model to fit?

- What the data represents
- How the data arose?
- Visual examination of the data
- Considerations of complexity

Suppose we observe $X_1 \ldots X_n \overset{iid}{\sim} f(x|\Theta)$

· How can we estimate $\Theta$
· How can we quantify our uncertainty

## The Method of Moments

Suppose $\Theta \in \mathbb{R}$ is a single unknown parameter

Pick $\Theta$ so that the mean of $f(x|\Theta)$ matches sample mean

$$\bar{x} = \frac{x_1 + \ldots x_n}{n}$$

Ex. The poisson distribution for parameter $\lambda > 0$ is a discrete distribution over non-negative integer counts

$$\text{PMF: } f(x|\lambda) = \frac{e^{-\lambda} \lambda^x}{x!} \quad \text{for } x \in \{0, 1, 2, \ldots\}$$

This distribution has mean $\lambda$. The M-O-M estimate is $\hat{\lambda} = \bar{x}$

Ex. The exponential distribution $\lambda > 0$

$$f(x|\lambda) = \lambda e^{-\lambda x} \qquad \text{mean}: \frac{1}{\lambda}$$

$$\text{equate} \qquad \frac{1}{\hat{\lambda}} = \overline{X}$$

More generally, suppose the parameters are $\Theta \in \mathbb{R}^k$. Equating the theoretical mean w/ sample avg. gives 1 equation

To get $k$ equations, consider the first $k$ moments of $X \sim f(x|\theta)$

$$\left.\begin{array}{l} \mu = \mathbb{E}[x] \\[6pt] \mu_2 = \mathbb{E}[X^2] \\[6pt] \quad\vdots \\[6pt] \mu_k = \mathbb{E}[X^k] \end{array}\right\} \quad \text{depend on } \Theta$$

M-O-M estimate

① Compute $\mu_1, \ldots \mu_k$ in terms of $\Theta$

② Equate theoretical moments w/ sample moments

③ Solve for $\Theta$


Ex. Let $X_1 \ldots X_n \overset{iid}{\sim} N(\mu, \sigma^2)$

$$\mu_1 = \mathbb{E}[X] = \mu$$

$$\mu_2 = \mathbb{E}[X^2] = \text{Var}[X] + (\mathbb{E}X)^2 = \sigma^2 + \mu^2$$

$$\text{Set } \hat{\mu} = \frac{1}{n}(X_1 + \ldots X_n)$$

$$\tilde{\sigma}^2 + \hat{M}^2 = \frac{1}{n}\left(X_1^2 + \ldots X_n^2\right)$$

Ex. $X_1 \ldots X_n \overset{iid}{\sim} \text{Gamma}(\alpha, B)$. Recall Mean in $\frac{\alpha}{B}$

Variance is $\frac{\alpha}{B^2}$

$$M_1 = \frac{\alpha}{B} \qquad M_2 = \frac{\hat{\alpha} + \hat{\alpha}^2}{\hat{B}}$$

## Bias, Variance, and mean-squared error

Any estimate $\hat{\Theta}$ for $\Theta \in \mathbb{R}^k$ is a statistic.

The bias of $\hat{\Theta}$ is $\mathbb{E}_\Theta[\hat{\Theta}] - \Theta$, where $\mathbb{E}_\Theta$ means expectation when the parameter is $\Theta$.

"Whats the difference between average value of $\hat{\Theta}$ and $\Theta$"

Standard error is just standard deviation of $\hat{\Theta}$

$$\sqrt{\text{Var}_\Theta(\hat{\Theta})}$$

"How far is $\hat{\Theta}$ typically from its average value"

Mean-squared-error of $\hat{\Theta}$ is $\mathbb{E}\left[(\hat{\Theta} - \Theta)^2\right]$

MSE combines bias and standard error

$$MSE = \text{Variance} + \text{Bias}^2$$

Usually, MSE, bias, and variance all depend on $\Theta$

An estimator is said to be unbiased if bias $= \mathbb{E}_\Theta[\hat{\Theta}] - \Theta = 0$ for every possible $\Theta$ in the parameter space

Ex. Consider $X_1 \ldots X_n \overset{iid}{\sim} \text{Poisson}(\lambda)$

Recall M-O-M estimate for $\hat{\lambda} = \bar{X}$

① Bias: $\mathbb{E}_\lambda(\hat{\lambda}) = \mathbb{E}_\lambda\left[\frac{1}{n}(X_1 + \ldots + X_n)\right] = \frac{1}{n}\left(\mathbb{E}(X_1) + \ldots \mathbb{E}(X_n)\right) = \lambda$

so bias $\mathbb{E}_\lambda(\hat{\lambda}) - \lambda = 0$ $\hat{\lambda}$ is unbiased for $\lambda$

② Standard error: $\text{Var}_\lambda[\hat{\lambda}] = \text{Var}_\lambda\left[\frac{1}{n}(X_1 + \ldots X_n)\right] = \frac{\lambda}{n}$

$$\text{s.e.} = \sqrt{\frac{\lambda}{n}}$$

③ $MSE = (s.e.)^2 + (bias)^2 = \frac{\lambda}{n}$

3/7: Maximum Likelihood Estimator

Recap: Parametric model $f(X|\Theta)$ parameterized by $\Theta \in \mathbb{R}^k$

   Estimate $\Theta$ via method-of-moments

   Compute $M_j = \mathbb{E}[X^j]$ in terms of $\Theta$ for $j = 1 \dots k$

   Equate values to sample values and solve

Def: The joint PMF or PDF of data $X_1 \dots X_n$ viewed as a function of the parameter $\Theta \in \mathbb{R}^k$ is the likelihood function $lik(\Theta)$

E.g. If $X_1 \dots X_n \overset{iid}{\sim} f(X|\Theta)$ then

$$lik(\Theta) = f(X_1|\Theta) \times \dots \times f(X_n|\Theta)$$

Recall from NP

$$H_0: X \sim f_0 \quad \text{vs.} \quad H_1: X \sim f_1$$

Likelihood ratio statistic $L(X) = \dfrac{f_0(X)}{f_1(X)}$

In the context of parametric models, if

$$f_0(X_1 \dots X_n) = f(X_1|\Theta_0) \times \dots \times f(X_n|\Theta_0)$$
$$f_1(X_1 \dots X_n) = f(X_1|\Theta_1) \times \dots \times f(X_n|\Theta_1)$$

then $L(X) = \dfrac{lik(\Theta_0)}{lik(\Theta_1)}$

Maximum Likelihood Estimator of $\Theta$ is the value of $\Theta$ in the parameter space that maximizes $lik(\Theta)$

Examples

How to compute MLE $\hat{\Theta}$?

Note: It's equivalent to maximize log likelihood

$$\ell(\Theta) = \log \, lik(\Theta)$$
$$= \sum_{i=1}^{n} \log f(x_i|\Theta) \quad \text{if } X_1 \dots X_n \text{ are } iid$$

Let $X_1 \dots X_n \overset{iid}{\sim}$ Poisson $(\lambda)$, Parametric space: $\lambda \in (0,\infty)$

The PMF is $f(X|\lambda) = \dfrac{e^{-\lambda} \lambda^x}{X!}$

Log Likelihood function

$$\ell(\Theta) = \sum_{i=1}^{n} \log f(x_i|\lambda) = \sum_{i=1}^{n} -\lambda + x_i \log(\lambda) - \log(X_i!)$$

$$= -n\lambda + \log(\lambda) \sum_{i=1}^{n} x_i - \sum_{i=1}^{n} \log(x_i!)$$

To maximize $\lambda$, we want $0 = \ell'(\lambda)$

$$\ell'(\Theta) = -n + \frac{1}{\lambda} \sum_{i=1}^{n} x_i = 0$$

$$\lambda = \frac{1}{n} \sum_{i=1}^{n} x_i = \overline{X}$$

First check: $\hat{\lambda}$ is in the parameter space

$$\hat{\lambda} = \bar{X} \quad \text{is in} \quad (0,\infty) \quad \text{not all} \quad X_i = 0$$

Second check: $\hat{\lambda}$ is a maximizer of the function

$$0 > \ell'(\lambda) \quad \lambda > \bar{X} \quad \text{and} \quad 0 < \ell'(\lambda)$$

$$\Updownarrow \qquad\qquad\qquad \Updownarrow$$

$$\hat{\lambda} > \hat{\lambda} = X \qquad\qquad \lambda < \hat{\lambda} = \bar{X}$$

$\ell(\lambda)$ is decreasing left of $\bar{X}$ and increasing right of $\bar{X}$

Alternatively check if $\ell''(\lambda) < 0$

$$\ell''(\lambda) = -\frac{1}{\lambda^2} \sum_{i=1}^{n} X_i$$

This is negative for all $\lambda \in (0,\infty)$ so $\ell(\lambda)$ is concave

If all $X_i$'s are $0$, our estimate is $\hat{\lambda} =$

Example: Let $X_1, \dots X_n \overset{iid}{\sim} N(\mu, \sigma^2)$. The log likelihood function is

$$\ell(\mu, \sigma^2) = \sum_{i=1}^{n} \log f(X_i | \mu, \sigma^2)$$

$$= \sum_{i=1}^{n} \log \left( \frac{1}{\sqrt{2\pi \sigma^2}} e^{-\frac{(X_i - \mu)^2}{2\sigma^2}} \right)$$

$$= \sum_{i=1}^{n} -\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{(X_i - \mu)^2}{2\sigma^2}$$

$$= -\frac{n}{2} \log 2\pi - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (X_i - \mu)^2$$

Maximize: $\ell(\mu, \sigma^2)$

First order condition

$$\ell'(\theta) = 0$$

$$\frac{\partial \ell}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^{n} (X_i - \mu) \qquad \Longleftrightarrow \qquad 0 = \sum_{i=1}^{n} X_i - n\mu \implies \hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} X_i = \bar{X}$$

$$\frac{\partial \ell}{\partial \sigma^2} = \frac{-n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^{n} (X_i - \mu)^2 \qquad \Longleftrightarrow \qquad 0 = \frac{-n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^{n} (X_i - \bar{X})^2 \implies \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{X})^2$$

Verifying Maximization

$$0 < \frac{\partial \ell}{\partial \mu} \Longleftrightarrow \mu < \bar{X} \qquad \text{and} \qquad 0 > \frac{\partial \ell}{\partial \mu} \Longleftrightarrow \mu > \bar{X}$$

Similarly for $\sigma^2$ by subbing $\mu = \bar{X}$

Example: Let $X_1, \ldots X_n \overset{iid}{\sim}$ Gamma $(\alpha, \beta)$    $a, \beta > 0$

Log Likelihood

$$\ell(\alpha, \beta) = \sum_{i=1}^{n} \log f(X_i | \alpha, \beta) = \sum_{i=1}^{n} \log\left(\frac{\beta^{\alpha}}{\Gamma(\alpha)} X_i^{\alpha-1} e^{-\beta X_i}\right) = \sum_{i=1}^{n} \alpha \log \beta - \log \Gamma(\alpha) + (\alpha-1) \log(X_i) - \beta X_i$$

$$= n\alpha \log \beta - n \log(\Gamma(\alpha)) + (\alpha-1) \sum_{i=1}^{n} \log(X_i) - \beta \sum_{i=1}^{n} X_i$$

Maximize

$$\frac{\partial \ell}{\partial \alpha} = n \log \beta - \frac{n \Gamma'(\alpha)}{\Gamma(\alpha)} + \sum_{i=1}^{n} \log X_i$$
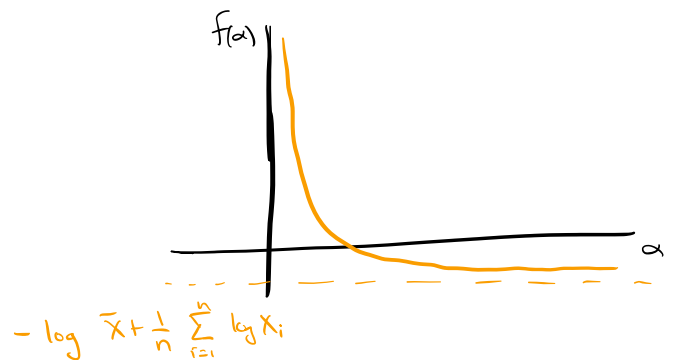
$$\frac{\partial \ell}{\partial \beta} = \frac{n\alpha}{\beta} - \sum_{i=1}^{n} X_i$$

Solve 2nd equation in $\beta$ so $\hat{\beta} = \hat{\alpha}/\bar{X}$

Substituting into first equation

$$0 = n \log\left(\frac{\hat{\alpha}}{\bar{X}}\right) - \frac{n \Gamma'(\hat{\alpha})}{\Gamma(\hat{\alpha})} + \sum_{i=1}^{n} \log(X_i)$$

$$0 = \underbrace{\log(\hat{\alpha}) - \frac{\Gamma'(\hat{\alpha})}{\Gamma(\hat{\alpha})}}_{f(\hat{\alpha}} - \log \bar{X} + \frac{1}{n} \sum_{i=1}^{n} \log X_i$$



$$-\log \bar{X} + \frac{1}{n} \sum_{i=1}^{n} \log X_i$$

Observe $f(\alpha)$ is decreasing

$$-\log(\bar{X}) + \frac{1}{n} \sum_{i=1}^{n} \log(X_i) < 0$$

by Jensen's Inequality

So there must be a unique root to $0 = f(\hat{\alpha})$

Then this $\hat{\alpha}$ and $\hat{\beta} = \frac{\hat{\alpha}}{\bar{X}}$ are the MLE's

you can check that these are the maximizers

Usually no explicit form for MLE, Instead compute via optimization algorithm

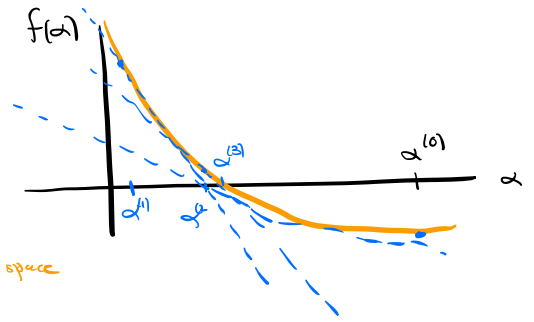Newton - Raphson method is a common approach

1. Initialize with a guess

   Often use MOM estimate

2. linearize $\alpha^{[t]}$ using a Taylor expansion

   $$0 \approx f(\alpha^t) + (\hat{\alpha} - \alpha^t) f'(\alpha^t)$$

   $$\hat{\alpha} = \alpha^t - \frac{f(\alpha^t)}{f'(\alpha^t)} \quad \longleftarrow \text{reset } \alpha^t \text{ if } \alpha^t \text{ falls outside of the parameter space}$$

3. Use linearized estimate for next iteration

This is different from M-O-M estimator

Example: Let $X_1 \dots X_K \sim$ Multinomial $(n, (p_1 \dots p_K))$

$X_1 \dots X_K$ are not IID, they sum to $n$

Log likelihood function is the log PMF for $(X_1 \dots X_K)$

$$\ell(p_1 \dots p_K) = \log\left[ \binom{n}{X_1, X_2, \dots X_K} p_1^{X_1} \cdot p_2^{X_2} \cdot \dots \cdot p_K^{X_K} \right]$$

$$= \log \binom{n}{X_1 \dots X_K} + \sum_{n=1}^{K} X_i \log p_i$$

The parameter space is the set $(p_1 \dots p_K)$ where $0 \leq P_i \leq 1$ and $P_1 + \dots + P_K = 1$

MLE maximizes $\ell(p_1 \dots p_K)$ subject to these constraints

Easiest way to do this is via Lagrange Multiplier Method

$$L(p_1 \dots p_K, \lambda) = \ell(p_1, \dots, p_K) + \lambda(p_1 + \dots + p_K - 1)$$

— Lagrange Multiplier

$$= \log \binom{n}{X_1 \dots X_K} + \sum_{i=1}^{K} X_i \log(p_i) + \lambda(p_1 + \dots p_K - 1)$$

Set all partials of $L$ to $0$

$$0 = \frac{\partial L}{\partial p_i} = \frac{X_i}{P_i} + \lambda$$

$$0 = \frac{\partial L}{\partial \lambda} = P_1 + \dots + P_K - 1$$

Solve 1st eq to find $P_i = -\frac{X_i}{\lambda}$

Substitute into 2nd eqn:

$$0 = \frac{-X_1}{\lambda} - \frac{X_2}{\lambda} - \dots \frac{X_K}{\lambda} - 1$$

$$\lambda = -(X_1 + \dots X_K) = -n$$

$$\hat{P_i} = \frac{X_i}{n} \quad \text{(fraction of observations in class } i)$$

Rationale for Lagrange Multiplier

1) Fix any $\lambda$. Maximizing $L(p_1 \dots p_K, \lambda)$ subject to $p_1 + \dots p_K = 1$ is equivalent to maximizing $\ell(p_1 \dots p_K)$ under the same constraint

   $\lambda$ term cancels to $0$

2) The unconstrained maximizer of $L(p_1, \ldots p_K, \lambda)$ over $(p_1, \ldots p_K)$ is $p_i = -\frac{X_i}{\lambda}$ as calculated above

3) Choosing $\lambda = -(X_1 + \ldots X_n) = -n$, the unconstrained max satisfies the constraint

Implies that the solved $p_i$'s are also constrained maximizers

from i) we have that this solution maximizes $\ell(p_1, \ldots p_n)$

Example: The genotypes $AA, Aa, aa$ at a locus satisfy Hardy-Weinberg Equilibrium

Occur with probabilities

$$(1-\Theta)^2, \quad 2\Theta(1-\Theta) \text{ and } \Theta^2$$

$\Theta = [0,1]$ is the frequency of $a$

Count occurrences of $AA, Aa, aa$ in $n$ samples can be modeled as multinomial

$$(X_1, X_2, X_3) \sim \text{Multinomial}\left(n, \left((1-\Theta)^2, 2\Theta(1-\Theta), \Theta^2\right)\right)$$

Similar to previous example with single parameter $\Theta = [0,1]$

Compute MLE for $\Theta$

$$\ell(\Theta) = \log\left[\binom{n}{X_1, X_2, X_3}\left((1-\Theta)^2\right)^{X_1}\left(2\Theta(1-\Theta)\right)^{X_2}\left(\Theta^2\right)^{X_3}\right]$$

$$= \log\binom{n}{X_1, X_2, X_3} + (2X_1 + X_2)\log(1-\Theta) + (X_2 + 2X_3)\log\Theta$$

Maximize over $\Theta$ (constraint already accounted for)

$$0 = \ell'(\Theta) = -\frac{2X_1 + X_2}{1-\Theta} + \frac{X_2 + 2X_3}{\Theta}$$

$$\hat{\Theta} = \frac{2X_3 + X_2}{2n}$$

# 3/9 : Normal Approximation, confidence Interval

Example: Poisson Model

$X_1 \ldots X_n \overset{iid}{\sim} \text{Poisson}(\lambda)$. Let $\lambda_0$ be the true parameter

Both M-O-M and MLE find $\hat{\lambda} = \bar{X}$

From lecture 12: $\mathbb{E}_{\lambda_0}[\hat{\lambda}] = \lambda_0$   unbiased estimator

$\text{Var}_{\lambda_0}[\hat{\lambda}] = \frac{\lambda_0}{n}$   standard error of $\sqrt{\frac{\lambda_0}{n}}$

By LLN: $\hat{\lambda} = \bar{X} \to \lambda_0$ in probability as $n \to \infty$

We say $\hat{\lambda}$ is consistent for $\lambda$

By CLT: $\sqrt{n}(\hat{\lambda} - \lambda_0) \longrightarrow N(0, \lambda_0)$ in distribution as $n \to \infty$

Informally, $\hat{\lambda}$ has approximate distribution $N\left(\lambda_0, \frac{\lambda_0}{n}\right)$ for large $n$

This allows us to make a confidence interval for $\lambda_0$

random interval containing $\lambda_0$ w/ a pre-specifies probability, $1 - \alpha \in (0,1)$

Let $Z(\alpha/2)$ be the upper-$\alpha$ point of $N(0,1)$

CLT implies $\mathbb{P}\left[-\sqrt{\frac{\lambda_0}{n}}\, z\left(\alpha/2\right) \leq \hat{\lambda}-\lambda_0 \leq \sqrt{\frac{\lambda_0}{n}}\, z\left(\alpha/2\right)\right] \simeq 1-\alpha$

Substituting $\sqrt{\hat{\lambda}/n}$ for $\sqrt{\lambda_0/n}$ we find

$\mathbb{P}\left[-\sqrt{\frac{\hat{\lambda}}{n}}\, z\left(\alpha/2\right) \leq \hat{\lambda}-\lambda_0 \leq \sqrt{\frac{\hat{\lambda}}{n}}\, z\left(\alpha/2\right)\right] \simeq 1-\alpha$

We have lower and upper bounds of $\lambda_0$

$\Rightarrow \lambda_0$ belongs to the interval $\hat{\lambda} + \sqrt{\frac{\hat{\lambda}}{n}} \cdot z\left(\alpha/2\right)$  w/ probability $1-\alpha$

More formally: The coverage guarantee is

$$\mathbb{P}_{\lambda_0}\left[\lambda_0 \in \left[\hat{\lambda} - \sqrt{\hat{\lambda}/n} \cdot z(\alpha/2),\ \hat{\lambda} + \sqrt{\hat{\lambda}/n} \cdot z(\alpha/2)\right]\right] \longrightarrow 1-\alpha \text{ as } n\to\infty$$

Because

· $\sqrt{n}\left(\hat{\lambda}-\lambda_0\right) \longrightarrow N(0,\lambda_0)$  by CLT

· $\dfrac{1}{\sqrt{\hat{\lambda}}} \longrightarrow \dfrac{1}{\sqrt{\lambda_0}}$  by LLN

· $\dfrac{\sqrt{n}\left(\hat{\lambda}-\lambda_0\right)}{\sqrt{\hat{\lambda}}} \longrightarrow \dfrac{1}{\sqrt{\lambda_0}} \cdot N(0,\lambda_0) = N(0,1)$  by Slutsky's Lemma

## Asymptotic Normality of the MLE

Thm: Let $f(X|\theta)$ be a parametric model, where $\theta \in \mathbb{R}$. Let $\theta_0$ be the true parameter, let $X_1,\ldots X_n \overset{iid}{\sim} f(X|\theta_0)$

Let $\hat{\theta}$ be the MLE.

Under some smoothness conditions for $f(X|\theta)$, as $n\to\infty$

a) $\hat{\theta}$ is a consistent estimator

$\hat{\theta} \to \theta_0$ in probability as $n\to\infty$

b) $\hat{\theta}$ is asymptotic normal, and $\sqrt{n}\left(\hat{\theta}-\theta_0\right) \longrightarrow N\left(0, 1/I(\theta_0)\right)$ in distribution

$I(\theta)$ has two definitions

Single sample

$$I(\theta) = \mathrm{Var}_\theta\left[\frac{\partial}{\partial\theta}\log f(X|\theta)\right] = -\mathbb{E}_\theta\left[\frac{\partial^2}{\partial\theta^2}\log f(X|\theta)\right]$$

· $I(\theta)$ is the Fisher Information

· $\dfrac{\partial}{\partial\theta}\log f(X|\theta)$ is the score

We will show $\mathbb{E}_\theta\left[\frac{\partial}{\partial\theta}\log f(X|\theta)\right] = 0$

So $I(\theta) = \mathbb{E}_\theta\left[\left(\frac{\partial}{\partial\theta}\log f(X|\theta)\right)^2\right]$

· Implications

- Asymptotically $\hat{\theta}$ is unbiased. The bias is smaller scale than $1/\sqrt{n}$

$\mathbb{E}\left(\sqrt{n}(\hat{\theta}-\theta_0)\right) \to 0$, if bias was larger the constant term of the bias would be $\mathbb{E}$ of expression

- For large $n$, the standard error of the MLE is

$$S.E. \ \hat{\theta} = \sqrt{1/nI(\theta)} \quad . \quad \text{Is on the order of } 1/\sqrt{n}$$

$$MSE = (bias)^2 + (std. error)^2 \text{ is dominated by std. error for large } n$$

- Distribution of $\hat{\theta}$ is $\approx N(\theta_0, \frac{1}{nI(\theta_0)})$

- Confidence interval for $\theta_0$ with coverage $1-\alpha$

$$\hat{\theta} \pm \sqrt{\frac{1}{nI(\hat{\theta})}} \cdot z(\%_2) \quad , \quad \sqrt{\frac{1}{I(\hat{\theta})}} \text{ estimates } \sqrt{\frac{1}{I(\theta_0)}}$$

**Example:** Consider again $X_1, \dots X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda_0)$

MLE is $\hat{\lambda} = \overline{X}$

$$\log f(X|\lambda) = \log \frac{e^{-\lambda} \lambda^x}{x!} = -\lambda + \log \lambda - \log(x!)$$

$$\frac{\partial}{\partial \lambda} \log f(X|\lambda) = -1 + \frac{X}{\lambda} \quad \leftarrow \text{ Score}$$

$$\frac{\partial^2}{\partial \lambda^2} \log f(X|\lambda) = -\frac{X}{\lambda^2}$$

$$\mathbb{E}_{\lambda_0}\left[\frac{\partial}{\partial \lambda} \log f(X|\lambda)\right] = \mathbb{E}_{\lambda_0}\left[-1 + \frac{X}{\lambda}\right] = 0$$

$$I(\lambda_0) = \text{Var}_{\lambda_0}\left[-1 + \frac{X}{\lambda_0}\right] = \text{Var}\left[\frac{X}{\lambda_0}\right] = \frac{1}{\lambda_0^2} \text{Var}[X] = \frac{1}{\lambda_0}$$

Alternatively,

$$I(\lambda_0) = -\mathbb{E}_{\lambda_0}\left[-X/\lambda_0^2\right] = -\frac{1}{\lambda_0^2} \mathbb{E}[-X] = \frac{1}{\lambda_0}$$

Theorem shows $\sqrt{n}(\hat{\lambda} - \lambda_0) \to N(0, \frac{1}{I(\lambda_0)}) = N(0, \lambda_0)$

**Proof Sketch**

Consistency: The MLE $\hat{\theta}$ maximizes

$$\frac{1}{n} \sum_{i=1}^{n} \log f(X|\theta)$$

Suppose $\theta_0$ is the true parameter. For any fixed $\theta$, this is the average of $n$ IID random variables.

As $n \to \infty$, by LLN

$$\frac{1}{n} \sum_{i=1}^{n} \log f(X|\theta) \to \mathbb{E}_{\theta_0}\left[\log f(X_i|\theta)\right]$$

Implies under some conditions that the maximizer of LHS maximizes RHS

Maximizer of RHS is $\theta_0$

$$\mathbb{E}_{\theta_0}\left[\log f(X_i|\theta)\right] - \mathbb{E}_{\theta_0}\left[\log f(X_i|\theta_0)\right] = \mathbb{E}_{\theta_0}\left[\log \frac{f(X_i|\theta)}{f(X_i|\theta_0)}\right] \leq \log \mathbb{E}_{\theta_0}\left[\frac{f(X_i|\theta)}{f(X_i|\theta_0)}\right]$$

*PMF*   *Jensen's Inequality since log is concave*

$$= \log \int \frac{f(x|\theta)}{f(x|\theta_0)} \cdot f(x|\theta_0) \, dx = \log \int f(x|\theta) \, dx = \log 1 = 0$$

$$\Rightarrow \quad \mathbb{E}_{\Theta_0}\left[\log f(x_i|\Theta)\right] - \mathbb{E}_{\Theta_0}\left[\log f(x_i|\Theta_0)\right] \leq 0$$

So this is maximized over $\Theta$ at $\Theta = \Theta_0$

## Fisher Information

$$\int f(x|\Theta)\, dx = 1$$

$$\Rightarrow \quad 0 = \frac{\partial}{\partial \Theta} \int f(x|\Theta)\, dx = \int \frac{\partial}{\partial \Theta} f(x|\Theta)\, dx$$

$$\frac{\partial}{\partial \Theta} \log f(x|\Theta) = \frac{\frac{\partial}{\partial \Theta} f(x|\Theta)}{f(x|\Theta)} \quad \Longleftrightarrow \quad \frac{\partial}{\partial \Theta} f(x|\Theta) = \left(\frac{\partial}{\partial \Theta} \log f(x|\Theta)\right) f(x|\Theta)$$

$$\Rightarrow \quad 0 = \int \left(\frac{\partial}{\partial \Theta} \log f(x|\Theta) \cdot f(x|\Theta)\right) dx = \mathbb{E}_\Theta\left[\frac{\partial}{\partial \Theta} \log f(x|\Theta)\right]$$

So the score has mean 0.

## Differentiate a second time

$$0 = \frac{\partial}{\partial \Theta} \int \left(\frac{\partial}{\partial \Theta} \log f(x|\Theta)\right) f(x|\Theta)\, dx$$

$$\underset{\text{product rule}}{=} \int \left(\frac{\partial^2}{\partial \Theta^2} \log f(x|\Theta)\right) \cdot f(x|\Theta)\, dx + \int \left(\frac{\partial}{\partial \Theta} \log f(x|\Theta)\right)\left(\frac{\partial}{\partial \Theta} f(x|\Theta)\right) dx$$

$$\underbrace{\qquad\qquad\qquad}$$
$$\frac{\partial}{\partial \Theta} \log f(x|\Theta) \cdot f(x|\Theta)$$

$$= \mathbb{E}_\Theta\left[\frac{\partial^2}{\partial \Theta^2} \log f(x|\Theta)\right] + \mathbb{E}_\Theta\left[\left(\frac{\partial}{\partial \Theta} \log f(x|\Theta)\right)^2\right]$$

$$\Rightarrow -\mathbb{E}_\Theta\left[\frac{\partial^2}{\partial \Theta^2} \log f(x|\Theta)\right] = \mathbb{E}_\Theta\left[\left(\frac{\partial}{\partial \Theta} \log f(x|\Theta)\right)^2\right] = \text{Var}_\Theta\left[\frac{\partial}{\partial \Theta} \log f(x|\Theta)\right]$$

## Asymptotical Normality

$$\sqrt{n}\,(\hat{\Theta} - \Theta_0) \longrightarrow N\left(0, \frac{1}{I(\Theta_0)}\right)$$

The MLE $\hat{\Theta}$ maximizes $\ell(\Theta) = \sum_{i=1}^{n} \log f(x|\Theta)$

So $0 = \ell'(\hat{\Theta}) \approx \ell'(\Theta_0) + (\hat{\Theta} - \Theta) \cdot \ell''(\Theta_0) \quad \longleftarrow$ Taylor Expansion

$$\hat{\Theta} - \Theta_0 \approx -\frac{\ell'(\Theta_0)}{\ell''(\Theta_0)}$$

$$\sqrt{n}\,(\hat{\Theta} - \Theta_0) \approx -\sqrt{n}\,\frac{\ell'(\Theta_0)}{\ell''(\Theta_0)} = \frac{-\frac{1}{\sqrt{n}} \ell'(\Theta_0)}{-\frac{1}{n} \ell''(\Theta_0)}$$

$$-\frac{1}{n} \ell''(\Theta_0) = \frac{1}{n} \sum_{i=1}^{n} \left[-\frac{\partial^2}{\partial \Theta^2} \log f(x_i|\Theta_0)\right]$$

iid w/ mean $-\mathbb{E}_0\left[\frac{\partial^2}{\partial\theta^2}\log f(x_i|\theta)\right] = I(\theta_0)$ in pr

$$\frac{1}{\sqrt{n}}\ell'(\theta_0) = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}\frac{\partial}{\partial\theta}\log f(x_i|\theta_0)$$

iid w/ mean $= 0$

$$\text{Var}_{\theta_0}\left[\frac{\partial}{\partial\theta}\log f(x_i-\theta_0)\right] = I(\theta_0)$$

$$\longrightarrow N\left(0, I(\theta_0)\right) \text{ by CLT}$$

By Slutsky Lemma

$$\sqrt{n}(\hat\theta - \theta) \xrightarrow{L} \frac{1}{I(\theta_0)}\cdot N\left(0, I(\theta_0)\right) = N\left(0, \frac{1}{I(\theta)}\right)$$

---

3/14: Plug-in Estimates, delta Method

Recap: $X_1 \dots X_n \overset{iid}{\sim} f(x|\theta)$, $\theta \in \mathbb{R}^k$

① Method of Moments

Write
$M = h(\theta)$
$$\begin{cases} M_1 = \mathbb{E}_\theta[X] = \frac{1}{n}\sum_{i=1}^{n} X_i \\ \vdots \\ M_k = \mathbb{E}_\theta[X^k] = \frac{1}{n}\sum_{i=1}^{n} X_i^k \end{cases}$$
Solve equations for $\theta$

② Maximum Likelihood Estimation

$$\ell(\theta) = \sum_{i=1}^{n}\log f(x_i|\theta)$$

Maximize $\ell(\theta)$ over the parameter space to get $\hat\theta$

Thm: The MLE $\hat\theta$ is consistent for $\theta$, as $n\to\infty$

$$\sqrt{n}\left(\hat\theta - \theta\right) \to N\left(0, \frac{1}{I(\theta)}\right) \text{ for } \theta\in\mathbb{R}$$

$$I(\theta) = \text{Var}_\theta\left[\frac{\partial}{\partial\theta}\log f(x|\theta)\right] = -\mathbb{E}_\theta\left[\frac{\partial^2}{\partial\theta^2}\log f(x|\theta)\right]$$

Asymptotic Confidence Interval level $1-\alpha$ for $\theta$

$$\hat\theta \pm \sqrt{\frac{1}{n\,I(\theta)}}\cdot z(\%)$$

Estimating a function of $\theta$

First estimate $\theta$ by $\hat\theta$, then use $g(\hat\theta)$ as an estimate for $g(\theta)$

Example: Binomial coin, heads prob is $p$

If heads lose $1
If tails win $X

What value of $X$ makes this game fair
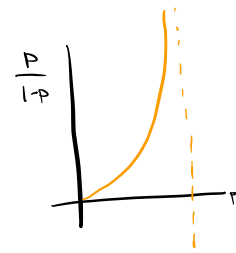
Expected winnings: $p \cdot (-1) + (1-p)x = 0$

$X = \dfrac{p}{1-p}$ is called the odds of this game

Often consider log-odds or logit

Estimate $\log \dfrac{p}{1-p}$ from $X_1, \ldots X_n \overset{iid}{\sim}$ Bernoulli$(p)$

① First estimate $p$, by $\bar{X}$

② Plug-in to get $\log \dfrac{\bar{X}}{1-\bar{X}}$ as an estimate of $\log \dfrac{p}{1-p}$



Example: Pareto Distribution

PDF: $f(X \mid x_0, \Theta) = \Theta x_0^{\Theta} x^{-\Theta-1}$ for $X \geq x_0$

Suppose we know $x_0 = 1$, we don't know $\Theta$, so

$$f(X \mid \Theta) = \Theta x^{-\Theta-1} \text{ for } X \geq 1$$

Mean: $\dfrac{\Theta}{\Theta-1}$   when $\Theta > 1$

Variance: $\dfrac{\Theta}{(\Theta-1)^2 (\Theta-2)}$   when $\Theta > 2$

How to estimate mean from $X_1 \ldots X_n \overset{iid}{\sim} f(X \mid \Theta)$

① Estimate $\Theta$

$$\ell(\Theta) = \sum_{i=1}^{n} \log f(x \mid \Theta) = \sum_{i=1}^{n} \log \Theta - (\Theta+1) \log X_i$$

$$= n \cdot \log \Theta - (\Theta+1) \sum_{i=1}^{n} \log X_i$$

To maximize this over $\Theta$

$$0 = \ell'(\Theta) = \dfrac{n}{\Theta} - \sum_{i=1}^{n} \log X_i$$

$$\Rightarrow \hat{\Theta} = \dfrac{n}{\sum_{i=1}^{n} \log X_i} \text{ is the MLE for } \Theta$$

② Estimate the mean $\dfrac{\Theta}{\Theta-1}$ by $\dfrac{\hat{\Theta}}{\hat{\Theta}-1}$

## Delta Method

Goal: Quantify the uncertainty (asymptotically) for $g(\hat{\Theta})$ based on uncertainty of $\hat{\Theta}$ itself

Thm: If $g: \mathbb{R} \to \mathbb{R}$ is continuously differentiable at $\Theta$, and if $\sqrt{n}(\hat{\Theta} - \Theta) \to N(0, v(\Theta))$ in distribution as $n \to \infty$, then

$$\sqrt{n}(g(\hat{\Theta}) - g(\Theta)) \to N\left(0, g'(\Theta)^2 \cdot v(\Theta)\right) \text{ in distribution as } n \to \infty$$

**Proof sketch:** Apply Taylor expansion of $g(\hat{\theta})$ around $\hat{\theta} = \theta$

$$g(\hat{\theta}) \approx g(\theta) + (\hat{\theta} - \theta) g'(\theta)$$

Then, $\sqrt{n}\left(g(\hat{\theta}) - g(\theta)\right) \approx \sqrt{n}\left(\hat{\theta} - \theta\right) g'(\theta)$

$$\longrightarrow g'(\theta) \, N(0, v(\theta)) = N(0, g'(\theta)^2 \cdot v(\theta))$$

**Example:** Let $X_1, \ldots X_n \overset{iid}{\sim}$ Bernoulli $(p)$

We estimate $\log \frac{p}{1-p}$ by $\log \frac{\bar{X}}{1-\bar{X}}$

Apply Delta method:

① By CLT: $\sqrt{n}(\bar{X} - p) \longrightarrow N(0, p(1-p))$

② Let $g(p) = \log \frac{p}{1-p} = \log p - \log 1-p$

$$g'(p) = \frac{1}{p(1-p)}$$

So, $\sqrt{n}\left(\log \frac{\bar{X}}{1-\bar{X}} - \log \frac{p}{1-p}\right) \longrightarrow N\left(0, \frac{1}{p(1-p)}\right)$

Suppose we toss $n = 100$ coins with $60$ heads

$$\bar{X} = 0.6$$

$$\log \text{ odds} \approx \frac{0.6}{1-0.6} = 0.41$$

$$\text{Standard error} = \sqrt{\frac{1}{n\bar{X}(1-\bar{X})}} \approx 0.2$$

Asymptotic Level $1-\alpha$, confidence interval is $0.41 \pm 0.2 \cdot z\left(\frac{\alpha}{2}\right)$

**Back to Pareto Example**

Let $X_1, \ldots X_n \overset{iid}{\sim}$ Pareto $(1, \theta)$

Recall: MLE is $\hat{\theta} = \dfrac{n}{\sum\limits_{i=1}^{n} \log X_i}$

Plug in estimate for mean $\dfrac{\theta}{\theta - 1}$ was $\dfrac{\hat{\theta}}{\hat{\theta} - 1}$

Apply delta method:

① Understand $\hat{\theta}$
$$\sqrt{n}(\hat{\theta} - \theta) \longrightarrow N\left(0, \frac{1}{I(\theta)}\right)$$

Computing $I(\theta)$:

$$\log f(X|\theta) = \log \theta - (\theta+1)\log X$$

$$\frac{\partial}{\partial\theta} \log f(X|\theta) = \frac{1}{\theta} - \log X$$

$$\frac{\partial^2}{\partial\theta^2} \log f(X|\theta) = -\frac{1}{\theta^2}$$

$$I(\theta) = -\mathbb{E}_\theta\left[-1/\theta^2\right] = \frac{1}{\theta^2}$$

$$\sqrt{n}\left(\hat{\theta} - \theta\right) \longrightarrow N(0, \theta^2)$$

② Set $g(\theta) = \dfrac{\theta}{\theta-1}$

$$g'(\theta) = -\frac{1}{(\theta-1)^2}$$

So, $\sqrt{n}\left(\dfrac{\hat{\theta}}{\hat{\theta}-1} - \dfrac{\theta}{\theta-1}\right) \longrightarrow N\left(0, \left[-\dfrac{1}{(\theta-1)^2}\right]^2 \theta^2\right) = N\left(0, \dfrac{\theta^2}{(\theta-1)^4}\right)$

An alternative estimate for the mean is $\bar{X}$

For $\bar{X}$ estimate, we can apply CLT

$$\sqrt{n}\left(\bar{X} - \frac{\theta}{\theta-1}\right) \longrightarrow N\left(0, \frac{\theta}{(\theta-1)^2(\theta-2)}\right)$$

Notice variance from MLE method is less than $\bar{X}$ estimate

Reduces impact of extreme observations

Standard Error for Method of Moments

Consider $\theta \in \mathbb{R}$. Estimate $\mu = \mathbb{E}_\theta[X]$ by $\bar{X}$

Suppose $\mu = h(\theta)$ for some function $h()$

Let $g$ be the inverse function of $h$ so $\theta = g(\mu)$

MOM estimate for $\theta = g(\mu)$ is $\hat{\theta} = g(\bar{X})$

For Standard Error

① By CLT, $\sqrt{n}(\bar{X} - \mu) \longrightarrow N(0, V(\theta))$

$$V(\theta) = \text{var}_\theta[X_i]$$

② By delta method

$$\sqrt{n}(\hat{\theta} - \theta) = \sqrt{n}\left(g(\bar{x}) - g(h(\theta))\right) \longrightarrow N\left(0, g'(h(\theta))^2 \cdot V(\theta)\right)$$

Standard error of $\hat{\theta}$ is

$$\sqrt{\frac{g'(h(\theta))^2 \cdot V(\hat{\theta})}{n}}$$

Example: Consider $X_1, \dots X_n \overset{iid}{\sim} \text{Exponential}(\lambda)$

PDF: $f(x|\lambda) = \lambda e^{-\lambda x}$ for $x > 0$

mean: $1/\lambda$    var: $1/\lambda^2$

Inverse of $h()$ is $\lambda = g(\mu) = 1/\mu$

$$g(\bar{x}) = 1/\bar{x}$$

Estimate std. error

① By CLT, $\sqrt{n}(\bar{x} - 1/\lambda) \longrightarrow N(0, 1/\lambda^2)$

② $g'(\mu) = -1/\mu^2$ so $g'(1/\lambda) = -\lambda^2$

$$\sqrt{n}\left(\frac{1}{\bar{x}} - \lambda\right) = \sqrt{n}\left(g(\bar{x}) - g(1/\lambda)\right) \longrightarrow N\left(0, g'(1/\lambda)^2 \cdot 1/\lambda^2\right) = N(0, \lambda^2)$$

Therefore std. error is $\approx \sqrt{\frac{\lambda^2}{n}}$ for large $n$, can estimate via $\sqrt{\frac{1}{n\bar{x}^2}}$

## 3/16: Cramer - Rao Bound, Asymptotic Efficiency

Recap: Pareto Model

PDF: $f(x|\theta) = \theta x^{-\theta - 1}$ for $x \geq 1$

$X_1, \dots X_n \overset{iid}{\sim} f(x|\theta)$

MLE $\hat{\theta}$: $\dfrac{n}{\sum_{i=1}^{n} \log x_i}$ ,  $\sqrt{n}(\hat{\theta} - \theta) \longrightarrow N(0, \theta^2)$

Mean in this model: $\mu = \dfrac{\theta}{\theta - 1}$ , $\sqrt{n}\left(\dfrac{\hat{\theta}}{\hat{\theta} - 1} - \dfrac{\theta}{\theta - 1}\right) \longrightarrow N\left(0, \dfrac{\theta^2}{(\theta - 1)^4}\right)$

Method-of-Moments Estimate for $\theta$

Solve $\mu = \dfrac{\theta}{\theta - 1}$ for $\theta$

$\theta = g(\mu) = \dfrac{\mu}{\mu - 1}$

M-O-M Estimator $\qquad \hat{\Theta} = \dfrac{\overline{X}}{\overline{X}-1}$

Delta Method

① By CLT, $\sqrt{n}\left(\overline{X} - \dfrac{\Theta}{\Theta-1}\right) \longrightarrow N\left(0, \dfrac{\Theta}{(\Theta-1)^2(\Theta-2)}\right)$

② $g'(\mu) = \dfrac{1}{\mu-1} - \dfrac{\mu}{(\mu-1)^2} = \dfrac{-1}{(\mu-1)^2}$

$\Rightarrow g'\left(\dfrac{\Theta}{\Theta-1}\right) = -(\Theta-1)^2$

So $\sqrt{n}\left(\hat{\Theta} - \Theta\right) = \sqrt{n}\left(g(\overline{X}) - g\left(\dfrac{\Theta}{\Theta-1}\right)\right) \longrightarrow N\left(0, \dfrac{\Theta}{(\Theta-1)^2(\Theta-2)} \cdot \left[-(\Theta-1)^2\right]^2\right)$

$\qquad\qquad = N\left(0, \dfrac{\Theta(\Theta-1)^2}{\Theta-2}\right)$

Geometry of the Fisher Information

Recall that as $n \to \infty$, fixing the parameter $\Theta_0$,

$\qquad \dfrac{1}{n}\sum_{i=1}^{n}\log f(X_i|\Theta) \to E_{\Theta_0}\left[\log f(X|\Theta)\right]$ for every $\Theta$



Sample $-\log$ likelihood $\qquad\qquad$ Expected log-likelihood

Fisher information is the curvature of $\overline{\ell}(\Theta)$ at $\Theta = \Theta_0$
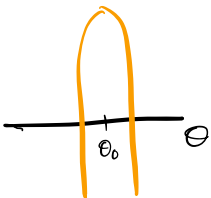
$\qquad \overline{\ell}(\Theta) = E_{\Theta_0}\left[\log f(X|\Theta)\right]$

Reason: $\quad I(\Theta) = -E_{\Theta}\left[\dfrac{\partial^2}{\partial\Theta^2}\log f(X|\Theta)\right]$

For a fixed true parameter $\Theta_0$

$\qquad I(\Theta_0) = -E_{\Theta_0}\left[\dfrac{\partial^2}{\partial\Theta^2}\log f(X|\Theta)\Big|_{\Theta=\Theta_0}\right]$

$\qquad\qquad = -\dfrac{\partial^2}{\partial\Theta^2}\underbrace{E_{\Theta_0}\left[\log f(X|\Theta)\right]}_{\overline{\ell}(\Theta)}\Big|_{\Theta=\Theta_0} \qquad = -\overline{\ell}''(\Theta_0)$

If $I(\Theta_0)$ is large, then the curvature of $\overline{\ell}$ around $\Theta_0$ is big



$\qquad$ If we perturb $\Theta$ a bit from $\Theta_0$, then $\overline{\ell}(\Theta)$ decreases a lot from $\overline{\ell}(\Theta_0)$

$\qquad$ Data contains a lot of information about $\Theta_0$

If $I(\Theta_0)$ is small, then $\bar{\ell}$ would have small curvature around $\Theta_0$

perturbations from $\Theta_0$ do not result in large changes



Data contains less information about $\Theta_0$

## Cramer-Rao Low Bound

Thm: Consider a parametric model $f(x|\Theta)$ where $\Theta \in \mathbb{R}$ is a single parameter. Let $T$ be any unbiased estimator of $\Theta$ based on $n$ observations $x_1, \ldots x_n \overset{iid}{\sim} f(x|\Theta)$. Then (under smoothness conditions)

$$\text{Var}_\Theta [T] \geq \frac{1}{n I(\Theta)}$$

Interpretation: If $T = \hat{\Theta}$ is the MLE, then for large $n$, $\text{Var}[T] \simeq \frac{1}{n I(\Theta)}$

Any unbiased estimator can't achieve a smaller variance

Proof: Let $Z = \frac{\partial}{\partial \Theta} \log f(x_1, \ldots x_n | \Theta) = \sum_{i=1}^{n} \frac{\partial}{\partial \Theta} \log f(x_i | \Theta)$

Recall: $\mathbb{E}_\Theta \left[ \frac{\partial}{\partial \Theta} \log f(x|\Theta) \right] = 0$

$\text{Var}_\Theta \left[ \frac{\partial}{\partial \Theta} \log f(x|\Theta) \right] = I(\Theta)$

So, $\mathbb{E}_\Theta [Z] = 0$ and $\text{Var}_\Theta [Z] = n I(\Theta)$

The correlation between $Z$ and $T$ belongs to $[-1, 1]$

$$\text{Cov}_\Theta [Z, T]^2 \leq \text{Var}_\Theta [Z] \, \text{Var}[T]$$

$$\text{Var}_\Theta [Z] = n I(\Theta)$$

$$\text{Cov}_\Theta [Z, T] = \mathbb{E}_\Theta \left[ (Z - \mathbb{E}_\Theta Z)(T - \mathbb{E}_\Theta T) \right]$$

$$= \mathbb{E}_\Theta \left[ (Z - \mathbb{E}_\Theta Z) T \right]$$

$$= \mathbb{E}_\Theta [ZT] \qquad bc \quad \mathbb{E}_\Theta Z = 0$$

Since $T$ is unbiased for $\Theta$

$$\Theta = \mathbb{E}_\Theta [T] = \int T(x_1, \ldots x_n) f(x_1, \ldots x_n | \Theta) dx_1 \ldots dx_n$$

$$1 = \int T(x_1 \cdots x_n) \cdot \frac{\partial}{\partial \Theta} f(x_1 \cdots x_n | \Theta) \, dx_1 \cdots dx_n$$

Recall: $\left[ \frac{\partial}{\partial \Theta} \log f(x_1 \cdots x_n | \Theta) \right] \cdot f(x_1 \cdots x_n | \Theta) = \frac{\partial}{\partial \Theta} f(x_1 \cdots x_n | \Theta)$

$$= \int T(x_1 \cdots x_n) \cdot \underbrace{\left[ \frac{\partial}{\partial \Theta} \log f(x_1 \cdots x_n | \Theta) \right]}_{Z} \cdot f(x_1 \cdots x_n | \Theta) \, dx_1 \cdots dx_n$$

where $\underbrace{T(x_1 \cdots x_n)}_{T}$

$$= \mathbb{E}_{\Theta}[TZ]$$

Therefore, $\text{Cov}[T, Z] = 1$

$$\text{Var}_{\Theta}[T] \geq \frac{1}{n I(\Theta)} \qquad \leftarrow \text{Plugging everything}$$

$\hat{\Theta}$ is an asymptotically efficient estimator for $\Theta$

if $\quad \sqrt{n}(\hat{\Theta} - \Theta) \longrightarrow N\left(0, \frac{1}{I(\Theta)}\right)$ in distribution

MLE is asymptotically efficient

If two estimators $\hat{\Theta}, \tilde{\Theta}$ satisfy

$$\sqrt{n}(\hat{\Theta} - \Theta) \longrightarrow N(0, M(\Theta))$$

$$\sqrt{n}(\tilde{\Theta} - \Theta) \longrightarrow N(0, V(\Theta))$$

Then, $\frac{V(\Theta)}{M(\Theta)}$ is the asymptotic relative efficiency of $\hat{\Theta}$ relative to $\tilde{\Theta}$

Interpretted as ratio of sample sizes required

$$\text{Var}[\hat{\Theta}] \approx \frac{M(\Theta)}{n} \qquad \text{Var}[\tilde{\Theta}] = \frac{V(\Theta)}{n}$$

For plug-in estimates

$$\sqrt{n}(g(\hat{\Theta}) - g(\Theta)) \longrightarrow N\left(0, g'(\Theta)^2 / I(\Theta)\right)$$

$$\text{Var}[g(\hat{\Theta})] \approx \frac{g'(\Theta)^2}{n I(\Theta)}$$

# Cramer-Rao Bound for plug-in estimates

In a parametric model $f(x|\Theta)$ for $\Theta \in \mathbb{R}$, if $T$ is <u>any</u> unbiased estimator for $g(\Theta)$ based on $X_1 \dots X_n \overset{iid}{\sim} f(x|\Theta)$

then
$$\text{Var}_\Theta [T] \geq \frac{g'(\Theta)^2}{n \cdot I(\Theta)}$$

Analogous proof as for traditional Cramer-Rao lower bound

# Fisher Information for multiple parameters

$\Theta \in \mathbb{R}^k$, the fisher information matrix $I(\Theta) \in \mathbb{R}^{k \times k}$ is the matrix w/

$(i,j)$ entry:
$$I(\Theta)_{ij} = \text{Cov}_\Theta \left[ \frac{\partial}{\partial \Theta_i} \log f(x|\Theta), \frac{\partial}{\partial \Theta_j} \log f(x|\Theta) \right]$$
$$= -\mathbb{E}_\Theta \left[ \frac{\partial^2}{\partial \Theta_i \partial \Theta_j} \log f(x|\Theta) \right]$$

Thm: Let $f(x|\Theta)$ be a parametric model where $\Theta \in \mathbb{R}^k$

Let $\hat{\Theta}$ be the MLE for $\Theta$ based on $n$ observations $X_1 \dots X_n \overset{iid}{\sim} f(x|\Theta)$

then, under smoothness conditions and assuming $I(\Theta)$ is invertible

$$\underbrace{\sqrt{n} (\hat{\Theta} - \Theta)}_{\in \mathbb{R}^k} \to \underbrace{\mathcal{N}(0, I(\Theta)^{-1})}_{\text{k-dimension multivariate normal}}$$

← Inverse matrix

Example: Consider $X_1 \dots X_n \overset{iid}{\sim} \text{Gamma}(\alpha, B)$  $\alpha, B > 0$

We can compute $\alpha, B$ via MLE

Use $I(\alpha, B)$ to quantify their uncertainty

$$\log f(x|\alpha, B) = \alpha \log B - \log \Gamma(x) + (\alpha - 1) \log x - Bx$$

$$\frac{\partial}{\partial \alpha} \log f(X(\alpha, \beta)) = \log \beta - \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} + \log X$$

$$\frac{\partial}{\partial \beta} \log f(X|\alpha, \beta) = \frac{\alpha}{\beta} - X$$

$$\frac{\partial^2}{\partial \alpha^2} \log f(X|\alpha, \beta) = \frac{\Gamma'(\alpha)^2}{\Gamma(\alpha)^2} - \frac{\Gamma''(\alpha)}{\Gamma(\alpha)} = -\psi(\alpha)$$

← tri-gamma function

$$\frac{\partial^2}{\partial \alpha, \partial \beta} \log f(X(\alpha, \beta)) = \frac{1}{\beta}$$

$$\frac{\partial^2}{\partial \beta^2} \log f(X|\alpha, \beta) = \frac{-\alpha}{\beta^2}$$

So
$$I(\alpha, \beta) = - \begin{pmatrix} -\psi(\alpha) & \frac{1}{\beta} \\ \frac{1}{\beta} & -\frac{\alpha}{\beta^2} \end{pmatrix} = \begin{pmatrix} \psi(\alpha) & -\frac{1}{\beta} \\ -\frac{1}{\beta} & \frac{\alpha}{\beta^2} \end{pmatrix}$$

$$I(\alpha, \beta)^{-1} = \frac{1}{\psi(\alpha)\frac{\alpha}{\beta} - \frac{1}{\beta^2}} \begin{pmatrix} \frac{\alpha}{\beta^2} & \frac{1}{\beta} \\ \frac{1}{\beta} & \psi(\alpha) \end{pmatrix}$$

Errors are positively correlated

## 3|28 : Bayesian Inference

Parametric model : Data $X = (X_1, \dots X_n)$

Modeled by some distribution $f(X, \Theta)$ w/ parameter $\Theta$

Frequentist Perspective : Fixed true value we try to estimate

$\Theta$ is non-random

Bayesian Perspective : Treat $\Theta$ as a random variable with a distribution

## Prior and Posterior Distributions

Review joint, marginal, and conditional distributions

Consider two r.v.'s $X + Y$

Joint PDF or PMF $f_{X,Y}(x,y)$

Marginal Distribution of $X$ is given by

$$f_X(x) = \int f_{X,Y}(x,y)\, dy \quad \text{or} \quad f_X(x) = \sum_y f_{X,Y}(x,y)$$

<span style="color:orange">continuous</span>                                           <span style="color:orange">discrete</span>

Conditional distribution of $Y$ given $X = x$ is then

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)}$$

This is a PDF or PMF over values $y$, fixing the value $X$

Similarly define

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$$

The Joint PDF/PMF factors in two ways

$$f_{X,Y}(x,y) = f_{Y|X}(y|x) \cdot f_X(x)$$

$$= f_{X|Y}(x|y) \cdot f_Y(\ )$$

Bayesian Inference: Observed data $\mathbb{X} = (X_1 \ldots X_n)$

    Parametric model for $\mathbb{X}$ with unknown parameter $\textcircled{H}$

        Think about $\textcircled{H}$ as random

        - Interpret the model fo $\mathbb{X}$ as the conditional distribution given $\textcircled{H} = \Theta$

$$f_{\mathbb{X}|\textcircled{H}}(\mathbb{X}|\Theta)$$

        $\textcircled{H}$ is random w/ distribution $f_{\textcircled{H}}(\Theta)$ — the prior distribution

        $\Rightarrow$ This defines a joint distribution over both $\textcircled{H}$ and $\mathbb{X}$

$$f_{\mathbb{X},\Theta}(\mathbb{X},\Theta) = f_{\mathbb{X}|\textcircled{H}}(\mathbb{X}|\Theta) \times f_{\textcircled{H}}(\Theta)$$

    We can factor this joint distribution

$$f_{\mathbb{X},\Theta}(\mathbb{X},\Theta) = f_{\Theta|\mathbb{X}}(\Theta|\mathbb{X}) \times f_{\mathbb{X}}(\mathbb{X})$$

$f_{\Theta}(\theta)$ is the prior distribution for $\Theta$, ie. marginal distribution of $\Theta$

$f_{X|\Theta}(X|\theta)$ is our parametric model for $X$ i.e. the likelihood function

$$f_X(x) = \int f_{X,\Theta}(x,\theta)\,d\theta \quad \text{is the marginal distribution of } X.$$

Distribution of data averaging over $\Theta$

$f_{\Theta|X}(\theta|X)$ is the posterior distribution of $\Theta$

This is our belief about the value of $\Theta$, after seeing the data

Goal: Understand $f_{\Theta|X}(\theta|X)$

$$f_{\Theta|X}(\theta|X) = \frac{f_{X,\Theta}(x,\theta)}{f_X(x)} = \frac{f_{X|\Theta}(X|\theta) \times f_{\Theta}(\theta)}{f_X(x)}$$

$$f_{\Theta|X}(\theta|X) \propto f_{X|\Theta}(X|\theta) \times f_{\Theta}(\theta)$$

"posterior $\propto$ likelihood $\times$ prior"

Example: Coin, heads probability $P$

Consider a prior distribution for $P \sim$ Uniform $(0,1)$

Observe $X_1 \ldots X_n \overset{iid}{\sim}$ Bernoulli $(p)$ given $P=p$

what is the posterior of $P$?

Joint distribution of $X_1 \cdots X_n$, $P$ is

$$f_{X,P}(X,P) = f_{X|P}(X|p) \times f_P(p)$$

$$= \prod_{i=1}^{n} p^{x_i}(1-p)^{1-x_i} \times 1 \quad \leftarrow \text{PDF uniform}$$

$$= p^s(1-p)^{n-s} \quad \text{where} \quad s = X_1 + \cdots + X_n$$

Marginal PMF of $X$ is

$$f_X(X) = \int_0^1 f_{X,P}(X,p) \, dp$$

$$= \int_0^1 p^s(1-p)^{n-s} \, dp$$

This is the beta integral

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1}(1-x)^{\beta-1} \, dx = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$$

$$\Rightarrow f_X(X) = B(s+1, n-s+1)$$

Posterior Distribution of $P$ given $X = x$ is

$$f_{P|X}(P|X) = \frac{f_{X,P}(X,p)}{f_X(X)}$$

$$= \frac{1}{B(s+1, n-s+1)} \cdot p^s(1-p)^{n-s} \quad \leftarrow \text{PDF of Beta}(s+1, n-s+1)$$

**Example:** Can extend to a more general prior for P

Consider Beta $(\alpha, \beta)$ prior for P:

$$f_P(p) = \frac{1}{\beta(\alpha,\beta)} p^{\alpha-1} (1-p)^{\beta-1} \quad \text{for } p \in (0,1)$$

$\alpha = \beta = 1$ gives Uniform prior

Posterior for P satisfies

$$f_{P|X}(p|X) \propto f_{X|p}(X|p) \times f_P(p)$$

$$\propto p^s (1-p)^{n-s} \cdot p^{\alpha-1} (1-p)^{\beta-1}$$

$$= p^{s+\alpha-1} (1-p)^{n-s+\beta-1}$$

This is proportional to Beta $(s+\alpha, n-s+\beta)$

Beta distribution is the posterior for P

**Example:** Observe counts $X_1 \ldots X_n$ model as $X_1 \ldots X_n \overset{iid}{\sim}$ Poisson $(\lambda)$

Treat parameter $\Lambda$ as random

Take the prior $\Lambda \sim$ Gamma $(\alpha, \beta)$. This has PDF $f_\Lambda(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}$ for $\lambda > 0$

What is the posterior distribution of $\Lambda$?

$$f_{\Lambda|X}(\lambda|X) \propto f_{(X|\lambda)}(X|\lambda) \times f_\Lambda(\lambda)$$

$$= \prod_{i=1}^n \frac{\lambda^{x_i} e^{-\lambda}}{x_i!} \times \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}$$

$$\propto \lambda^{\sum_{i=1}^n x_i} e^{-n\lambda} \times \lambda^{\alpha-1} e^{-\beta\lambda}$$

$$= \lambda^{\alpha + \sum_{i=1}^n x_i - 1} \times e^{-(\beta+\alpha)\lambda}$$

This is proportional to the Gamma $\left(\alpha + \sum_{i=1}^n x_i, \beta + \alpha\right)$

So this Gamma distribution is the posterior for $\Lambda$

Example: Observe $X_1 \ldots X_n \overset{iid}{\sim} N(\theta, 1/\xi)$ where $\xi = 1/\sigma^2$ is the precision

① Assume that $\xi$ is known, treat $\Theta$ as random

Consider a normal prior for $\Theta \sim N(\mu_{prior}, 1/\xi_{prior})$

$$f_{\Theta|X}(\theta|x) \propto f_{X|\Theta}(x|\theta) \times f_{\Theta}(\theta)$$

$$= \prod_{i=1}^{n} e^{-\frac{\xi}{2}(x_i - \theta)^2} \times e^{-\frac{\xi_{prior}}{2}(\theta - \mu_{prior})^2}$$

$$\propto \exp\left( \frac{-n\xi}{2}\theta^2 + \xi\left(\sum_{i=1}^{n} x_i\right)\theta - \frac{\xi_{prior}}{2}\theta^2 + \xi_{prior}\mu_{prior}\theta \right)$$

Suppose $\xi_{post} = n\xi + \xi_{prior}$ and $\mu_{post} = \dfrac{\xi\sum_{i=1}^{n}x_i + \xi_{prior}\mu_{prior}}{n\xi + \xi_{prior}}$

$$f_{\Theta|X}(\theta|x)$$

Some more algebra i dont fucking know

$$N\left(\mu_{post}, 1/\xi_{post}\right)$$

② Assume mean $\Theta$ is known, precision $\xi$ is unknown

Consider Gamma prior $\frac{1}{\sigma^2} \sim \text{Gamma}(\alpha, \beta)$

Then

$$f_{\frac{1}{\sigma^2}|X}(\xi|x) \propto \prod_{i=1}^{n} \sqrt{\xi} \, e^{-\frac{\xi}{2}(x_i - \theta)^2} \times \xi^{\alpha - 1} e^{-\beta\xi}$$

$$\propto \xi^{\frac{\alpha + n}{2} - 1} e^{-\left(\beta + \sum_{i=1}^{n} \frac{(x_i - \theta)^2}{2}\right)\xi}$$

GAMMA!

## Bayesian Point Estimates and Credible Intervals

To get a numeric estimate for $\Theta$:

- Mean of the posterior distribution
- Mode of the posterior ("MAP estimate")

To get an interval estimate for $\Theta$:

- Define a Bayesian credible interval
  $I$ w/ coverage level $1-\alpha$

An interval containing $1-\alpha$ portion of the posterior distribution

$$\mathbb{P}[\Theta \in I \mid X = x] = \int_I f_{\Theta|X}(\Theta|X)\, d\Theta = 1-\alpha$$

Common to use: lower $\frac{\alpha}{2}$-point to upper $\frac{\alpha}{2}$-point of the posterior

Example: $X_1, \ldots X_n \overset{iid}{\sim} \text{Bernoulli}(p)$

Prior $P \sim \text{Beta}(\alpha, \beta)$

Recall: Posterior $P|X = x \sim \text{Beta}(S + \alpha, n - S + \beta)$

where $S = X_1 + \cdots + X_n$

Estimate $\hat{p}$ by posterior mean

$$\hat{p} = \frac{S + \alpha}{n + \alpha + \beta}$$

Distinct from $\frac{S}{n}$ as calculated by M-O-M and MLE

$\hat{p}$ is a weighted average of sample and prior means

Credible Interval

lower $-0.05$ point to upper $-0.05$ point of $\text{Beta}(S+\alpha, n-S+\beta)$

## 3/30: Bayesian Inference (cont'd)

Likelihood model: $X = (X_1, \ldots X_n) \sim f_{X|\Theta}(X|\Theta)$

Prior: $f_\Theta(\theta)$

Posterior: $f_{\Theta|X}(\Theta|X) \propto f_{X|\Theta}(X|\Theta) \times f_\Theta(\theta)$

Examples:

① $X_1 \cdots X_n \overset{iid}{\sim} \text{Bernoulli}(p)$, prior $P \sim \text{Uniform}(0,1)$

$\Rightarrow (P|X) \sim \text{Beta}(X_1 + \cdots + X_n + 1, n - (X_1 + \cdots + X_n) + 1)$

② $X_1, \ldots X_n \overset{iid}{\sim} \text{Bernoulli}(p)$, prior $P \sim \text{Beta}(\alpha, \beta)$

$\Rightarrow (P|X) \sim \text{Beta}(X_1 + \cdots + X_n + \alpha, n - (X_1 + \cdots + X_n) + \beta)$

③ $X_1, \ldots X_n \overset{iid}{\sim} \text{Poisson}(\lambda)$, prior $\lambda \sim \text{Gamma}(\alpha, \beta)$

$\Rightarrow (\lambda|X) \sim \text{Gamma}(X_1 + \cdots X_n + \alpha, n + \beta)$

④ $X_1, \ldots X_n \overset{iid}{\sim} N(\Theta, \frac{1}{\xi})$

known $\xi$, prior $\Theta \sim N(\mu_{prior}, \frac{1}{\xi_{prior}})$

$$\Rightarrow (\Theta | \underline{X}) \sim N(\mu_{post}, {}^{1}/_{\mathcal{I}_{post}})$$

$$\mu_{post} = \frac{\sum_{i=1}^{n} X_n + \left(\mathcal{I}_{prior}/_{\mathcal{I}}\right)\mu_{prior}}{n + \left(\mathcal{I}_{prior}/_{\mathcal{I}}\right)} \quad , \quad \mathcal{I}_{post} = n\mathcal{I} + \mathcal{I}_{prior}$$

Posterior means:

③ $\hat{\lambda} = \dfrac{X_1 + \cdots + X_n + \alpha}{n + \beta}$

Different from MLE/MoM which is $\overline{X}$

— As if we had $\beta$ extra samples summing to $\alpha$

— $\hat{\lambda} = \dfrac{n}{n+\beta} \cdot \overline{X} + \dfrac{\beta}{n+\beta} \cdot \dfrac{\alpha}{\beta}$

<span style="color:orange">↑ Sample mean</span>  <span style="color:orange">↑ prior mean</span>

④ Posterior Mean: $\hat{\Theta} = \mu_{post}$

$$\hat{\Theta} = \mu_{post} = \frac{n}{n + \left(\mathcal{I}_{prior}/_{\mathcal{I}}\right)} \cdot \overline{X} + \frac{\mathcal{I}_{prior}/_{\mathcal{I}}}{n + \left(\mathcal{I}_{prior}/_{\mathcal{I}}\right)} \cdot \mu_{prior}$$

<span style="color:orange">↑ Sample mean</span>  <span style="color:orange">↑ prior mean</span>

A 70% Bayesian credible interval for $\Theta$ would be

$$\mu_{post} \pm \sqrt{\frac{1}{\mathcal{I}_{prior}}} \cdot z(0.05)$$

## Conjugate Priors and Improper Priors

A conjugate prior is a prior distribution when the resulting posterior has the same parametric form

- Beta prior $\longrightarrow$ Bernoulli prob
- Gamma prior $\longrightarrow$ Poisson rate
- Normal prior $\longrightarrow$ Normal mean
- Gamma prior $\longrightarrow$ Normal precision parameter

Unfortunately tend to be light-tailed (bias inference towards prior mean)
— Use heavier-tailed, non-conjugate priors for more robust inference

Goal is to use an uninformative prior

E.g. Poisson example $X_1 \ldots X_n \overset{iid}{\sim} Poisson(\lambda)$

Gamma$(\alpha, \beta)$ prior $\Rightarrow$ Posterior mean $\dfrac{X_1 + \cdots + X_n + \alpha}{n + \beta}$

Uninformative prior would mean smaller values of $\alpha, \beta$

At its limit : Gamma$(0,0) \propto \lambda^{-1}$

<span style="color:orange">↑ not a proper probability distribution</span>

Since gamma$(0,0)$ is not an actual distribution, we call it an improper prior

Improper priors can still yield proper posterior distributions

Bayesian vs. Frequentist Coverage Guarantees

Bayesian level-$(1-\alpha)$ credible interval guarantees

$$\mathbb{P}\left[\Theta \in I \mid \underline{X} = \underline{x}\right] = 1 - \alpha$$

Frequentist Confidence Interval

$$\mathbb{P}_{\Theta}\left[\Theta \in I\right] = 1 - \alpha \text{ where } \mathbb{P}_{\Theta} \text{ is over } X_1, \ldots X_n \overset{iid}{\sim} f(X|\Theta)$$

Example: Let $X_1, \ldots X_n \overset{iid}{\sim} N(\Theta, 1)$

MLE/MoM is that $\hat{\Theta} = \overline{X}$

Under parameter $\Theta$, $\overline{X} \sim N(\Theta, 1/n)$ so a frequentist level-$\alpha$ confidence interval is

$$\overline{X} \pm \frac{1}{\sqrt{n}} \cdot z(\alpha/2)$$

This guarantees $\mathbb{P}_\theta\left[\theta \in \bar{X} + \frac{1}{\sqrt{n}} \cdot z(\alpha/2)\right] = 1 - \alpha$

For a bayesian analysis $\quad \theta \sim N\left(0, 1/\xi_{prior}\right)$

$\quad (\theta|X) \sim N\left(\mu_{post}, 1/\xi_{prior}\right)$

$$\mu_{post} = \frac{\sum\limits_{i=1}^{n} X_i}{n + \xi_{prior}} = \frac{n}{n + \xi_{prior}} \bar{X}$$

$$\xi_{post} = n + \xi_{prior}$$

Level $(1-\alpha)$ Bayesian credible interval is

$$\frac{n}{n + \xi_{prior}} \cdot \bar{X} \pm \sqrt{\frac{1}{n + \xi_{prior}}} \cdot z(\alpha/2)$$

This guarantees : $\mathbb{P}\left[\theta \in \frac{n}{n+\xi_{prior}} \bar{X} \pm \sqrt{\frac{1}{n + \xi_{prior}}} \cdot z(\alpha/2) \,\bigg|\, X = x\right] = 1 - \alpha$

Bayesian credible interval does guarantee that the frequentist coverage probability averaged according to the prior for $\theta$

$\quad$ is $\quad 1 - \alpha$

the average coverage probability with prior $f_\theta(\theta)$ is $1 - \alpha$

## Normal Approximation

$\quad$ Frequentist and bayesian approach converge for large $n$

$\quad$ For a fixed prior, as $n \to \infty$, the influence of the prior vanishes
$\quad\quad$ mean and shape of posterior distribution are determined by data

$\quad$ Posterior will approach $N\left(\hat{\theta}, \frac{1}{nI(\theta)}\right)$

$\quad\quad\quad\quad \hat{\theta}$ is MLE $\quad$ and $\quad I(\theta)$ is Fischer Information

$\quad$ Bayesian Credible Interval will be $\hat{\theta} \pm \sqrt{\frac{1}{nI(\theta)}} \cdot z(\alpha/2)$

## Heuristic Explanation

$\quad$ Let $X_1 \dots X_n \overset{iid}{\sim} f(X|\theta) \quad$ prior $f_\theta(\theta)$

$\quad$ Define $\ell(\theta) = \sum\limits_{i=1}^{n} \log\left(f(x_i|\theta)\right)$ the total log likelihood

$\quad$ Posterior :

$$f_{\theta|x}(\theta|x) \propto \text{likelihood} \times \text{prior}$$

$$e^{\ell(\theta)} \times f_{(\theta)}(\theta)$$

Taylor Expand for $\theta$ close to MLE

$$\ell(\theta) = \ell(\hat{\theta}) + (\theta - \hat{\theta})\ell'(\hat{\theta}) + \frac{1}{2}(\theta - \hat{\theta})^2 \ell''(\hat{\theta})$$

$$\ell'(\hat{\theta}) = 0 \quad \text{be max}$$

$$\frac{1}{n}\ell''(\hat{\theta}) \simeq \frac{1}{n}\ell''(\theta_0) = -I(\theta_0) \simeq -I(\hat{\theta})$$

$$\ell(\theta) = \ell(\hat{\theta}) - \frac{1}{2}(\theta - \hat{\theta})^2 \cdot nI(\hat{\theta})$$

$\ell(\hat{\theta})$ depends on $X$ so it can be absorbed into proportionality

## 4/4 : MLE under misspecified models

$$X_1 \dots X_n \overset{iid}{\sim} f(X|\theta_0)$$

Estimate $\theta$, quantify uncertainty

- Bias, variance, MSE
- Consistency, asymptotic normality, efficiency

### Kullback-Leibler divergence

for discrete distributions w/ PMFs $f$ and $g$ on the sample space $X$

$$D_{KL}(g||f) = \sum_{x \in X} g(x) \log \frac{g(x)}{f(x)}$$

For continuous distributions w/ PDFs $f$ and $g$

$$D_{KL}(g||f) = \int g(x) \log \frac{g(x)}{f(x)} dx$$

Equivalent to

$$D_{KL}(g||f) = \mathbb{E}_g\left[\log \frac{g(x)}{f(x)}\right] = \mathbb{E}_g[\log g(x)] - \mathbb{E}_g[\log f(x)]$$

Asymmetric definition

expectation under true distribution $X \sim g(x)$

### Properties:

if $f(x) = g(x)$ for all $x$, then $D_{KL}(g||f) = 0$ since $\log \frac{g(x)}{f(x)} = \log 1 = 0$

For any $f$ and $g$, $D_{KL}(g||f) \geq 0$

Follows from applying Jensen's inequality

$$D_{KL}(g||f) = \mathbb{E}_g\left[-\log \frac{f(x)}{g(x)}\right] \geq -\log \mathbb{E}_g\left[\frac{f(x)}{g(x)}\right] = -\log \int g(x) \cdot \frac{f(x)}{g(x)} dx = -\log 1 = 0$$

$f$ and $g$ don't need to come from the same family

Ex. Let $f$ be $N(\mu_0, \sigma^2)$ and $g$ be $N(\mu_1, \sigma^2)$. What is $D_{KL}(g||f)$?

$$\log \frac{g(x)}{f(x)} = \log\left[\frac{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu_1)^2}{2\sigma^2}}}{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu_0)^2}{2\sigma^2}}}\right] = -\frac{(x-\mu_1)^2}{2\sigma^2} + \frac{(x-\mu_0)^2}{2\sigma^2} = \frac{2(\mu_1-\mu_0)x - (\mu_1^2 - \mu_0^2)}{2\sigma^2}$$

$$D_{KL}(g\|f) = \mathbb{E}_g\left[\log \frac{g(x)}{f(x)}\right] = \frac{2(\mu_1 - \mu_0)\mathbb{E}_g[X] - (\mu_1^2 - \mu_2^2)}{2\sigma^2} \quad \text{where} \quad \mathbb{E}_g[X] = \mu_1$$

$$= \frac{1}{2\sigma^2}(\mu_1 - \mu_2)^2$$

KL divergence is the squared difference of means, normalized by variance

In this example KL divergence is symmetric

Ex. Let $f$ be bernouilli $(p)$    What is $D_{KL}(g\|f)$
    $g$ be Bernoulli $(q)$

Then $\log \dfrac{g(x)}{f(x)} = \begin{cases} \log \ q/p & x=1 \\ \log \ \dfrac{1-q}{1-p} & x=0 \end{cases}$

$$D_{KL}(g\|p) = \mathbb{E}_g\left[\log \frac{g(x)}{f(x)}\right] = q \log \frac{q}{p} + (1-q) \log \frac{1-q}{1-p}$$

If $p \simeq q$, we can approximate this by taylor expansion

$$\log p \doteq \log q + (p-q)\cdot \frac{1}{q} + \frac{1}{2}(p-q)^2 \cdot \left(-\frac{1}{q^2}\right)$$

$$\log(1-p) \doteq \log(1-q) + (p-q)\left(-\frac{1}{1-q}\right) + \frac{1}{2}(p-q)^2 \cdot \left(-\frac{1}{(1-q)^2}\right)$$

$$D_{KL}(g\|f) = q\left(\log q - \log p\right) + (1-q)\left(\log(1-q) - \log(1-p)\right)$$

$$\approx q \cdot \left(-(p-q)\frac{1}{q} + \frac{1}{2}(p-q)^2 \cdot \frac{1}{q^2}\right) + (1-q)\left((p-q)\cdot\frac{1}{1-q} + \frac{1}{2}(p-q)^2 \cdot \frac{1}{(1-q)^2}\right)$$

$$= \frac{1}{2}(p-q)^2 \cdot \left(\frac{1}{q} + \frac{1}{1-q}\right) = \frac{(p-q)^2}{2q(1-q)}$$

Ex. $f$ is binomial $(n, p)$
   $g$ is binomial $(n, q)$

Then $\log \dfrac{g(x)}{f(x)} = \log\left(\binom{n}{x} q^x (1-q)^{n-x} \Big/ \binom{n}{x} p^x (1-p)^{n-x}\right) = x \log \frac{q}{p} + (n-x) \log \frac{1-q}{1-p}$

$$D_{KL}(g\|f) = \mathbb{E}_g\left[\log \frac{g(x)}{f(x)}\right] = \mathbb{E}_g[X] \cdot \log \frac{q}{p} + (n - \mathbb{E}_g[X]) \log \frac{1-q}{1-p}$$

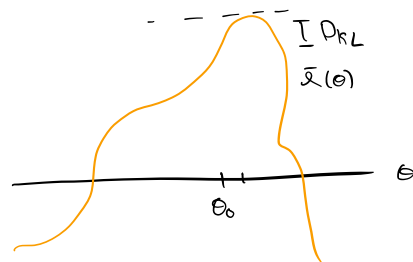$$= n\left(q \cdot \log \frac{q}{p} + (1-q) \log \frac{1-q}{1-p}\right)$$

<span style="color:orange">↑ exactly n times KL divergence in Bernoulli example</span>

For $p \doteq q$
$$D_{KL}(g\|f) \doteq n \cdot \frac{(p-q)^2}{2q(1-q)}$$

Generally, for $f(x) = f(x|\Theta)$ and $g(x) = f(x|\Theta_0)$

$$D_{KL}\left(f(x|\Theta_0) \| f(x|\Theta)\right) = \mathbb{E}_{\Theta_0}\left[\log \frac{f(x|\Theta_0)}{f(x|\Theta)}\right]$$

$$= \mathbb{E}_{\Theta_0}\left[\log f(x|\Theta_0)\right] - \mathbb{E}_{\Theta_0}\left[\log f(x|\Theta)\right]$$

$$= \bar{\lambda}(\Theta_0) - \bar{\lambda}(\Theta)$$

$$\tilde{\ell} = \mathbb{E}_{\Theta_0}\left[\log f(X|\Theta)\right]$$

Using a taylor expansion we find

$$D_{KL}\left(f(X|\Theta_0) \| f(X|\Theta)\right) \approx \frac{1}{2}(\Theta - \Theta_0)^2 \cdot I(\Theta_0)$$

Interpretation of $I(\Theta_0)$: scale factor that relates $(\Theta - \Theta_0)^2$ in parameter space to the theoretical "differences" $D_{KL}$

## Consistency of MLE in misspecified model

Suppose $X_1 \ldots X_n \overset{iid}{\sim} g(x)$, the pdf of the true distribution

We fit a parametric model $f(X|\Theta)$ that doesn't contain $g(x)$

What happens to the MLE?

By definition the mle $\hat{\Theta}$ maximizes

$$\frac{1}{n}\ell(\Theta) = \frac{1}{n}\sum_{i=1}^{\hat{n}} \log f(X_i|\Theta)$$

By LLN, $\frac{1}{n}\ell(\Theta) \longrightarrow \bar{\ell}(\Theta)$ in probability

$$\bar{\ell}(\Theta) = \mathbb{E}_g\left[\log f(X|\Theta)\right]$$

Here, if $g$ doesn't belong to our model

$$\bar{\ell}(\Theta) = \underbrace{\mathbb{E}_g\left[\log g(x)\right]}_{\substack{\text{independent of} \\ \Theta}} - \underbrace{\mathbb{E}_g\left[\log \frac{g(x)}{f(X|\Theta)}\right]}_{D_{KL}(g \| f(X|\Theta))}$$

therefore we maximize $\bar{\ell}(\Theta)$ by maximizing $D_{KL}(g\|f(X|\Theta))$

Suppose $\Theta \mapsto D_{KL}(g(x) \| f(X|\Theta))$ has a unique maximizer $\Theta^*$

Under some smoothness assumptions as $n \to \infty$, the MLE $\hat{\Theta}$ converges to $\Theta^*$ in probability

Example: Suppose we observe $X_1 \ldots X_n \geq 0$ and fit a model Exponential $\lambda$

$$f(X|\lambda) = \lambda e^{-\lambda x} \text{ for } x > 0$$

In reality $X_1 \ldots X_n \overset{iid}{\sim}$ Gamma $(\alpha, 1)$, $g(x) \frac{1}{\Gamma(\alpha)} x^{\alpha-1} e^{-x}$ only exponential if $\alpha = 1$

Thm tells us that maximizing $D_{KL}(g(x) \| f(X|\lambda^*))$ estimates $\lambda^*$

$$\log \frac{g(x)}{f(x|\lambda)} = \log\left(\frac{\frac{1}{\Gamma(\alpha)} x^{\alpha-1} e^{-x}}{\lambda e^{-\lambda x}}\right) = -\log \Gamma(\alpha) + (\alpha-1)\log x - \log \lambda - (\lambda-1) x$$

Minimizing $\lambda$

$$0 = -\frac{1}{\lambda} + \mathbb{E}_g[X] \implies \lambda^* = \frac{1}{\mathbb{E}_g[X]} = \frac{1}{\alpha}$$

MLE $\hat{\lambda} \longrightarrow \frac{1}{\alpha}$ in probability

Could solve this more directly through traditional methods

## Ex. What is the asymptotic variance of $\hat{\lambda}$?

We have explicit $\hat{\lambda} = \frac{1}{\bar{X}}$ so we can apply the delta method

$$\sqrt{n}\,(\bar{X} - \alpha) \longrightarrow N(0, \mathrm{Var}_g[X]) \text{ where } \mathrm{Var}_g[X] = \alpha$$

Apply delta method w/ $h(\alpha) = \frac{1}{\alpha} \Rightarrow h'(\alpha)^2 = \frac{1}{\alpha^4}$

So, $\sqrt{n}\left(\hat{\lambda} - \frac{1}{\alpha}\right) = \sqrt{n}\left(h(\bar{X}) - h(\alpha)\right)$

$$\longrightarrow N\left(0, \alpha \cdot \frac{1}{\alpha^4}\right) = N\left(0, \frac{1}{\alpha^3}\right)$$

Variance $\hat{\lambda}$ for large $n$ is $\frac{1}{n\alpha^3}$

Fisher information estimate would be incorrect

$$I(\lambda) = \frac{1}{\lambda^2}$$

$$\frac{1}{n I(\hat{\lambda})} \quad \text{to estimate variance}$$

$$\hat{\lambda} \approx \lambda^* = \frac{1}{\alpha}$$

$$I(\hat{\lambda}) \approx I(\lambda^*) = \alpha^2 \quad \text{and} \quad \frac{1}{n I(\hat{\lambda})} \approx \frac{1}{n\alpha^2} \quad \text{instead of} \quad \text{correct variance} \quad \frac{1}{n\alpha}$$

## 4/6: The bootstrap

Simulation based approach to quantify uncertainty of statistical estimates

Can estimate standard error or a confidence interval

Given $X_1 \dots X_n \overset{iid}{\sim} f(X|\theta)$, what is the standard error for an estimator $\hat{\theta}$ for $\theta$

Idea of the bootstrap is to simulate new data and compute $\hat{\theta}$ from each dataset

Unfortunately, you can't simulate $f(X|\theta)$ without knowing $\theta$

Bootstrap method involves simulations from an estimate of the true distribution

### Parametric Bootstrap

Assume $X_1 \dots X_n \overset{iid}{\sim} f(X|\theta)$

Estimate $\theta$ by $\hat{\theta}$ and simulate $X_1^* \dots X_n^* \overset{iid}{\sim} f(X|\hat{\theta})$

Analogous to the plug-in principle

### Non-parametric bootstrap

No assumption of a parametric model

Instead, we sample $X_1^* \dots X_n^*$ independently <u>with replacement</u> from the original $X_1 \dots X_n$

Sample size is $n$

Likely to have repeated values

Some samples will be lost (unsampled)

63.2% of samples are expected to be present

### Rationale for the nonparametric bootstrap

Estimated distribution is the empirical distribution

each observation places a mass of $\frac{1}{n}$

Draw new data from this estimated distribution

### Key differences between true distribution and empirical distribution

Empirical distribution is always discrete

Some statistics don't make sense to compare

mode, max value, min value

The empirical CDF very closely resembles the true CDF

$$\overline{F}_n(t) = \frac{1}{n} \sum_{i=1}^{\hat{n}} \mathbb{1}\{x_i \le t\}$$

as $n \to \infty$, this converges to CDF

Mean translates well between empirical and true distributions

## Bootstrap and Misspecified Models

Suppose $X_1, \dots X_n \overset{iid}{\sim} g(x)$. We fit $g(x)$ with Poisson $\lambda$

Fischer Information:

$$I(\lambda) = 1/\lambda$$

$$\sqrt{1/nI(\lambda)} = \sqrt{\frac{\lambda}{n}} = \sqrt{\frac{\overline{x}}{n}}$$

So true standard error is $\sqrt{\frac{\sigma^2}{n}}$

If poisson were correct then $\sigma^2 = \lambda$ so $\sqrt{\frac{\lambda}{n}}$ is an accurate estimate of standard error

Non parametric bootstrap guards against model misspecification

## Bootstrap Confidence Intervals

Let $\hat{\Theta}$ be an estimator of $\Theta$ and $\hat{se}$ be the bootstrap standard error estimate of $\hat{\Theta}$

$$\Theta \pm z(\alpha/2) \hat{se} \quad \text{for} \quad 100-\alpha \quad \text{confidence interval} \quad \textcolor{orange}{\leftarrow \text{Normal}}$$

### Percentile bootstrap interval

Let $\hat{\Theta}_1^*, \dots \hat{\Theta}_B^*$ be the values of $\hat{\Theta}$ computed in B simulations

Let $\hat{\Theta}^{*(\alpha/2)}$ and $\hat{\Theta}^{*(1-\alpha/2)}$ be the empirical $\alpha/2$ and $1-\alpha/2$ quantiles of the simulated values

$$\left[ \hat{\Theta}^{*(\alpha/2)}, \hat{\Theta}^{*(1-\alpha/2)} \right]$$

Requires symmetry?

### Basic Bootstrap Interval

Let $q^{(\alpha/2)}$ and $q^{(1-\alpha/2)}$ be the $\alpha/2$ and $1-\alpha/2$ quantiles of $\hat{\Theta}^* - \hat{\Theta}, \dots \hat{\Theta}_B^* - \hat{\Theta}$

Use this to approximate true distribution of $\hat{\Theta} - \Theta$

$$\left[ \hat{\Theta} - q^{(1-\alpha/2)}, \hat{\Theta} - q^{(\alpha/2)} \right] \quad \text{since} \quad \hat{\Theta} - \Theta \in \left[ q^{(\alpha/2)}, q^{(1-\alpha/2)} \right] \iff$$

Same as percentile confidence interval if it is symmetric

### Advantages and Disadvantages of bootstrap

- Easy to apply
- Obtains estimates valid under misspecification
- Can be computationally prohibitive
- Validity of bootstrap requires analytic proofs

## 4/11 : Generalized Likelihood Ratio Test

① Simple Null Hypothesis
② Sub-model null hypothesis

### Generalized Likelihood Ratio Test for Simple null hypothesis

Observe: $X_1, \dots X_n \overset{iid}{\sim} f(x|\Theta)$

$$H_0 : \Theta = \Theta_0$$
$$H_1 : \Theta \ne \Theta_0$$

The GLRT rejects $H_0$ for small values of $\Lambda$

$$\Lambda = \frac{\text{lik}(\Theta_0)}{\max_\Theta \text{lik}(\Theta)} \quad \text{where} \quad \text{lik}(\Theta) = \prod_{i=1}^{\hat{n}} f(x_i|\Theta)$$

$H_1$ is composite, so you replace $\text{lik}(\theta_1)$ with $\max\text{-lik}(\theta)$

$$\max_\theta \text{lik}(\theta) = \text{lik}(\hat\theta) \quad \text{where} \quad \hat\theta \text{ is the MLE}$$

By definition of MLE

$$\text{lik}(\hat\theta) \geq \text{lik}(\theta_0) \quad \text{so} \quad \Lambda \leq 1$$

If $H_0$ were true, we expect $\hat\theta \approx \theta_0$ for large $n$

Rejecting $H_0$ for small $\Lambda$ is the same as rejecting $H_0$ for large $-2\log\Lambda$

$$-2\log\Lambda = 2\ell(\hat\theta) - 2\ell(\theta_0)$$

$$\ell(\theta) = \log\text{lik}(\theta) = \sum_{i=1}^{n} \log\left(f(x_i|\theta)\right)$$

Since $\Lambda \leq 1$, $\quad -2\log\Lambda \geq 0$

We reject $H_0$ when

$$-2\log\Lambda \geq \chi^2_K(\alpha) \quad \textcolor{orange}{\leftarrow \text{upper } \alpha \text{ point of } \chi^2_K}$$

$K$ is the dimension of the parameter space

**Example:** Let $X_1 \ldots X_n \overset{iid}{\sim} N(\theta,1)$. Test null hypothesis $\theta=0$ and $H_1: \theta \neq 0$

MLE for $\theta$ is $\hat\theta = \bar X$. We compute

$$\text{lik}(0) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}} e^{\frac{-x_i^2}{2}}$$

$$\text{lik}(\hat\theta) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}} e^{-\frac{(x_i-\bar x)^2}{2}}$$

$$\Lambda = \frac{\text{lik}(0)}{\text{lik}(\hat\theta)} = \exp\left(-\sum_{i=1}^{n}\frac{x_i^2}{2} + \sum_{i=1}^{n}\frac{(x_i-\bar x)^2}{2}\right)$$

$$= \exp\left(-\frac{n}{2}\bar x^2\right) \quad \textcolor{orange}{\leftarrow \text{expand square and combine terms}}$$

$$-2\log\Lambda = n\bar x^2 = (\sqrt{n}\bar x)^2$$

Under $H_0: \theta = 0 \quad - \quad \bar X \sim N(0, 1/n) \quad$ so $\quad \sqrt{n}\bar x \sim N(0,1)$

so $\quad -2\log\Lambda = (\sqrt{n}\bar x)^2 \sim \chi^2_1$

For more general non-parametric models

**Thm:** Let $X_1 \ldots X_n \overset{iid}{\sim} f(X|\theta)$, a parametric model with parameter space of dimension $K$. Let $\theta_0$ be in the interior of the parameter space. Under smoothness conditions

$$H_0: \theta = \theta_0 \quad \text{vs.} \quad H_1: \theta \neq \theta_0$$

$$-2\log\Lambda \longrightarrow \chi^2_K \text{ in distribution as } n \to \infty \text{ under } H_0$$

So, the GLRT that rejects $-2\log\Lambda \geq \chi^2_K(\alpha)$ is an asymptotic level-$\alpha$ test

**Proof sketch:** Suppose $K=1$ when $\theta \in \mathbb{R}$ is a single parameter

$$-2\log\Lambda = 2\ell(\hat\theta) - 2\ell(\theta_0)$$

Taylor Expansion: $\ell(\theta_0)$ around $\hat\theta$

$$\ell(\theta_0) \approx \ell(\hat\theta) + (\theta_0 - \hat\theta)\cdot\ell'(\hat\theta) + \frac{1}{2}(\theta_0-\hat\theta)^2\cdot\ell''(\hat\theta)$$

$\textcolor{orange}{\text{Recall: } \ell'(\hat\theta) = 0 \text{ since it is at a maximizing point}}$

$$\textcolor{orange}{\ell''(\hat\theta) \approx -n I(\hat\theta) = -n I(\theta_0)}$$

$\textcolor{orange}{\uparrow \text{from earlier}}$

$$\Rightarrow \ell(\theta_0) \approx \ell(\hat\theta) - \frac{n}{2}I(\theta_0)(\theta_0-\hat\theta)^2$$

$$\Rightarrow -2\log\Lambda = n I(\theta_0)\cdot(\theta_0-\hat\theta)^2$$

Under $H_0: \sqrt{n}(\hat\theta - \theta_0) \longrightarrow N\left(0, 1/I(\theta_0)\right) \quad$ in distribution, by asymptotic normality for MLE

$$\sqrt{n I(\theta_0)}(\hat\theta - \theta_0) \longrightarrow N(0,1)$$

$$n I(\theta_0)(\hat\theta - \theta_0)^2 \longrightarrow \chi^2_1$$

Dimension is the number of parameters minus the number of equality constraints

Example: Consider $(X_1 \ldots X_n) \sim \text{Multinomial}(n, (P_1, P_2, \ldots P_k))$

$$H_0: (P_1, \ldots, P_k) = (P_{0_1}, \ldots P_{0_k})$$
$$H_1: (P_1, \ldots, P_k) \neq (P_{0_1}, \ldots P_{0_k})$$

Likelihood function:

$$\text{lik}(P_1 \ldots P_k) = \binom{n}{X_1 \ldots X_k} P_1^{X_1} P_2^{X_2} \cdots P_k^{X_k}$$

$$\ell(P_1 \ldots P_k) = \log\binom{n}{X_1 \ldots X_k} + X_1 \log P_1 + \cdots + X_k \log P_k$$

The MLE is $(\hat{P}_1, \ldots \hat{P}_k) = \left(\frac{X_1}{n}, \ldots, \frac{X_k}{n}\right)$

$$\Rightarrow -2\log \Lambda = 2\ell(\hat{P}_1 \ldots \hat{P}_k) = -2\ell(P_{0_1}, \ldots P_{0_k})$$
$$= 2X_1\left(\log \frac{X_i}{n} - \log P_{0_1}\right) + \cdots + 2X_k\left(\log \frac{X_k}{n} - \log P_{0_k}\right)$$

Dimension is $k-1$ (1 equality constraint)

GLRT rejects when $-2\log \Lambda \geq \chi^2_{k-1}(\alpha)$

## GLRT for a sub-model

Let $X_1, \ldots X_n \stackrel{iid}{\sim} f(X|\theta)$. Let $\Omega$ be the parameter space.

$$H_0: \theta \in \Omega_0 \qquad \Omega_0 \subset \Omega \text{ is a lower dimensional space in } \Omega$$
$$H_1: \theta \notin \Omega_0$$

Example: Consider $X_1 \ldots X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

$$\Omega = \{(\mu, \sigma^2): \sigma^2 > 0\} \text{ is a region of the plane } \mathbb{R}^2$$

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

Then $\Omega_0 = \{(\mu, \sigma^2): \mu = 0, \sigma^2 > 0\} \Leftarrow$ a line inside $\Omega$

$\Omega_0$ has dimension 1 while $\Omega$ has dimension 2

GLRT rejects $H_0$ for small values of

$$\Lambda = \frac{\max_{\theta \in \Omega_0} \text{lik}(\theta)}{\max_{\theta \in \Omega} \text{lik}(\theta)}$$

In other words, $\Lambda = \frac{\text{lik}(\hat{\theta}_0)}{\text{lik}(\hat{\theta})}$ 
$\hat{\theta}_0$ is MLE in $\Omega_0$
$\hat{\theta}$ is MLE in $\Omega$

Once again we consider $-2\log \Lambda$
GLRT rejects for

$$-2\log \Lambda \geq \chi^2_k(\alpha)$$

where $k$ is the difference in dimension between $\Omega$ and $\Omega_0$

If $\Omega_0 = \Theta_0$ we have dimension $= 0$, so the submodel generalization simplifies to the simple null example

## Example: (One-sample t-test)

$$X_1, \dots X_n \overset{iid}{\sim} N(\mu, \sigma^2)$$

log-likelihood:

$$\ell(\mu, \sigma^2) = \sum_{i=1}^{n} \log \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(X_i - \mu)^2}{2\sigma^2}}$$

$$= -\frac{n}{2} \log(2\pi\sigma^2) - \sum_{i=1}^{n} \frac{(x_i - \mu)^2}{2\sigma^2}$$

Full model MLE:

$$(\hat{\mu}, \hat{\sigma}^2) = \left(\bar{X}, \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2\right)$$

$$\Rightarrow \ell(\hat{\mu}, \hat{\sigma}^2) = -\frac{n}{2}\left(\log(2\pi\hat{\sigma}^2) + 1\right)$$

Sub-model MLE:

$$\mu = 0$$

MLE is $(\mu_0^2, \bar{\sigma}_0^2) = \left(0, \frac{1}{n} \sum_{i=1}^{n} x_i^2\right)$    ← from HW6

$$\Rightarrow \ell(\mu_0^2, \sigma_0^2) = -\frac{n}{2} \log(2\pi\hat{\sigma}^2) - \frac{n}{2}$$

$$-2 \log \Lambda = 2\ell(\hat{\mu}, \hat{\sigma}^2) - 2\ell(\hat{\mu}_0, \hat{\sigma}_0^2)$$

$$= n \log\left(\frac{\hat{\sigma}_0^2}{\hat{\sigma}^2}\right)$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2 \qquad \hat{\sigma}_0^2 = \frac{1}{n} \sum_{i=1}^{n} (x_i)^2 \qquad \text{so} \quad \hat{\sigma}^2 = \hat{\sigma}_0^2 - \bar{X}$$

$$\Rightarrow -2 \log \Lambda = n \log \frac{\hat{\sigma}_0^2}{\hat{\sigma}^2} = n \log\left(1 + \frac{\bar{x}^2}{\hat{\sigma}^2}\right) = n \log\left(1 + \frac{1}{n-1} T^2\right)$$

recall: $T = \dfrac{\sqrt{n}\,\bar{x}}{\frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2}$

This is increasing in $|T|$ so rejecting for large $-2\log\Lambda$ is the same as rejecting large $|T|$

Except, for 2 sided t-test we reject $H_0$ when $|T| \geq t_{n-1}(\alpha/2)$

     for GLRT we reject $H_0$ when $-2\log\Lambda \geq \chi_1^2(\alpha)$

for large $n$    $-2\log\Lambda \approx \chi_1^2$

$$-2\log\Lambda = n \log\left(1 + \frac{1}{n-1} T^2\right)$$

$$\approx n \cdot \frac{1}{n-1} T^2 \approx T^2 = Z^2 \quad \text{where} \quad Z \sim N(0,1)$$

## 4/13: Chi-squared tests for categorical data

### Recap: GLRT

$$H_0 : \Theta \in \Omega_0 \quad \text{vs.} \quad H_1 : \Theta \notin \Omega_0$$

Generalized Neyman Pearson Lemma

$$\Lambda = \frac{\text{lik}(\hat{\Theta}_0)}{\text{lik}(\hat{\Theta})} \quad \text{where:} \quad \hat{\Theta}_0 \text{ the MLE in } \Omega_0$$

$$\hat{\Theta} \text{ the MLE in full model}$$

Test statistic : $-2\log\Lambda$

For large sample distribution is $\approx \chi_k^2$ where $k$ is the difference in dimension between $\Omega$ and $\Omega_0$

## Example: Hardy-Weinberg Equilibrium

Genotypes at single locus $\in \{AA, Aa, aa\}$ in $n$ individuals with observed counts $N_{AA}, N_{Aa}$, and $N_{aa}$

Full model: $(N_{AA}, N_{Aa}, N_{aa}) \sim \text{Multinomial}(n, (P_{AA}, P_{Aa}, P_{aa}))$

$H_0 : P_{AA} = (1-\Theta)^2, \quad P_{Aa} = 2\Theta(1-\Theta), \quad P_{aa} = \Theta^2 \quad$ for some $\Theta \in (0,1)$

$$\text{lik}(P_{AA}, P_{Aa}, P_{aa}) = \binom{n}{N_{AA}, N_{Aa}, N_{aa}} \times P_{AA}^{N_{AA}} \, P_{Aa}^{N_{Aa}} \, P_{aa}^{N_{aa}}$$

$$\Rightarrow \quad \log \binom{n}{N_{AA}, N_{Aa}, N_{aa}} + N_{AA} \log P_{AA} + N_{Aa} \log P_{Aa} + N_{aa} \log P_{aa}$$

$$\Rightarrow \quad -2 \log \Lambda = 2\ell(\hat{P}_{AA}, \hat{P}_{Aa}, \hat{P}_{aa}) - 2\ell\left(\hat{P}_{0\,AA}, \hat{P}_{0\,Aa}, \hat{P}_{0\,aa}\right)$$

$$= 2N_{AA}\left(\log \frac{\hat{P}_{AA}}{\hat{P}_{0\,AA}}\right) + 2N_{Aa}\left(\log \frac{\hat{P}_{Aa}}{\hat{P}_{0\,Aa}}\right) + 2N_{aa}\left(\log \frac{\hat{P}_{aa}}{\hat{P}_{0\,aa}}\right)$$

Full Model MLE: $\hat{P}_{AA} = \dfrac{N_{AA}}{n}$ $\quad \hat{P}_{Aa} = \dfrac{N_{Aa}}{n}$ $\quad \hat{P}_{aa} = \dfrac{N_{aa}}{n}$

Sub-model MLE: $\quad \hat{\theta} = \dfrac{2N_{aa} - N_{Aa}}{2n}$

$$P_{0_{AA}} = (1-\hat{\theta})^2, \quad P_{Aa} = 2\hat{\theta}(1-\hat{\theta}), \quad P_{aa} = \hat{\theta}^2$$

Dimensions : Full model: $2$ ($3$ params w/ $1$ constraint)

Sub-model: $1$ (single param $\theta$)

$k = 1$

Compare $-2\log \Lambda$ to $\chi_1^2(\alpha)$

## Test of Independence

Example: General Social Survey

Random sample of $1972$ people

| | Dem | Repub | Indep |
|---|---|---|---|
| M | 422 | 381 | 273 |
| F | 299 | 365 | 232 |

Want to test whether gender is independent of party affiliation

Model : $\begin{pmatrix} N_{11} & N_{12} & N_{13} \\ N_{21} & N_{22} & N_{23} \end{pmatrix} \sim \text{Multinomial}\left(n, \begin{pmatrix} P_{11} & P_{12} & P_{13} \\ P_{21} & P_{22} & P_{23} \end{pmatrix}\right)$

Let $P_{i\cdot} = \sum_j P_{ij}$ (row sums)

$P_{\cdot j} = \sum_i P_{ij}$ (column sums)

$H_0 : P_{ij} = P_{i\cdot} \times P_{\cdot j}$ for every $i,j$ vs. $H_1:$

Dimensions :

Full model : $5$ ($6$ params, $1$ constraint $\sum_{ij} P_{ij} = 1$)

Sub model : $3$ ($5$ params, $2$ constraints)

Parameters : $P_{1\cdot} \quad P_{2\cdot} \quad P_{\cdot 1} \quad P_{\cdot 2} \quad P_{\cdot 3}$

Constraints : $P_{1\cdot} + P_{2\cdot} = 1 \quad P_{\cdot 1} + P_{\cdot 2} + P_{\cdot 3} = 1$

$k = 2$ so $-2 \log \Lambda$ compare against $\chi_2^2(\alpha)$

Consider more generally,

$$(N_1, \ldots, N_k) \sim \text{Multinomial}(n, (P_1 \ldots P_k))$$

Test $H_0 : (P_1 \ldots P_k) \in \Omega_0$

Multinomial likelihood

$$\text{lik}(P_1, \ldots P_k) = \binom{n}{N_1 \ldots N_k} \times \prod_{i=1}^{k} P_i^{N_i}$$

Let $\hat{P}_i$ and $\hat{P}_{0i}$ be the full model and submodel MLEs.

$$-2 \log \lambda = 2 \log \text{lik}(\hat{P}_1 \ldots \hat{P}_k) - 2 \log \text{lik}(\hat{P}_{0_1} \ldots \hat{P}_{0_k})$$

$$= 2 \sum_{r=1}^{k} N_i \cdot \log \frac{\hat{P}_i}{\hat{P}_{0,i}}$$

Recall: $\hat{P}_i = \frac{N_i}{n}$

Expected counts in sub-model $E_i = n \cdot \hat{P}_{0 \cdot i}$

$$\Rightarrow -2 \log \lambda = 2 \sum_{i=1}^{k} N_i \log \frac{N_i}{E_i}$$

Return to $k=6$ example

$$\text{lik}(P_1., P_2., P_{\cdot 1}, P_{\cdot 2}, P_{\cdot 3}) = \binom{n}{N_{11}, N_{12} \ldots N_{23}} \times \prod_{r=1}^{2} \prod_{j=1}^{3} (P_i \ldots P_j)^{N_{ij}}$$

$$= \binom{n}{N_{11}, N_{12} \ldots N_{23}} \times \prod_{i=1}^{2} P_i.^{N_i.} \times \prod_{j=1}^{3} P_{\cdot j}^{N_{\cdot j}}$$

where $N_i. = N_{i1} + N_{i2} + N_{i3}$

$$N_{\cdot j} = N_{1j} + N_{2j}$$

Take a log and maximize subject to $\sum_i P_i. = 1$, $\sum_j P_{\cdot j} = 1$

$\mathcal{L} = \text{Lagrangian}: \log \binom{n}{N_{11} \ldots N_{23}} + \sum_{i=1}^{2} N_i. \log P_i. + \sum_{j=1}^{3} N_{\cdot j} \log P_{\cdot j} + \lambda \left( \sum_{r=1}^{2} P_i. - 1 \right) + \mu \left( \sum_{j=1}^{3} P_{\cdot j} - 1 \right)$

$$0 = \frac{\partial \mathcal{L}}{\partial P_i.} = \frac{N_i.}{P_i.} + \lambda \quad \Rightarrow \quad P_i. = \frac{-N_i.}{\lambda}$$

$$0 = \frac{\partial \mathcal{L}}{\partial P_{\cdot j}} = \frac{N_{\cdot j}}{P_{\cdot j}} + \mu \Rightarrow P_{\cdot j} = \frac{-N_{\cdot j}}{\mu}$$

$$0 = \frac{\partial \mathcal{L}}{\partial \lambda} = P_1. + P_2. - 1 \Rightarrow -\frac{N_1.}{\lambda} - \frac{N_{2i}}{\lambda} = -\frac{n}{\lambda} = 1$$

$$\Rightarrow \lambda = -n$$

$$\Rightarrow \hat{P}_i. = \frac{N_i.}{n}$$

$$0 = \frac{\partial \mathcal{L}}{\partial \mu} = P_{\cdot 1} + P_{\cdot 2} + P_{\cdot 3} - 1 = -\frac{N_{\cdot 1}}{\mu} - \frac{N_{\cdot 2}}{\mu} - \frac{N_{\cdot 3}}{\mu} = \frac{-n}{\mu} = 1$$

$$\Rightarrow \mu = -n$$

$$\Rightarrow \hat{P}_{\cdot j} = \frac{N_{\cdot j}}{n}$$

$$\Rightarrow \hat{P}_{0_{ij}} = \hat{P}_i. \times \hat{P}_{\cdot j} = \frac{N_i. N_{\cdot j}}{n^2}$$

Expected counts: $E_{ij} = n \cdot \hat{P}_{0_{ij}} = \frac{N_i. \times N_{\cdot j}}{n}$

Plugin into $-2 \log \Lambda = 2 \sum_{r=1}^{2} \sum_{j=1}^{3} N_{ij} \log \frac{N_{ij}}{E_{ij}} = 8.31$

Compare to $\chi_2^2$ so we find a p-value of 0.016

Pearson chi-squared test

Alternative to GLRT

$$\chi^2 = \sum_{r=1}^{k} \frac{(N_i - E_i)^2}{E_i}$$

Test rejects $H_0$ when $\chi^2$ exceeds $\chi^2_{dof}(\alpha)$ value

<span style="color:orange">↗ difference in dimension of the models</span>

$\chi^2$ is similar to $-2 \log \Lambda$

<span style="color:orange">See via a taylor expansion (4/13 lecture 22)</span>

Test of Homogeneity

Setting: $(N_1, \ldots, N_k) \sim$ Multinomial $(n, (p_1, \ldots p_k))$

$(M_1, \ldots M_k) \sim$ Multinomial $(m, (q_1, \ldots q_k))$

$H_0: p_i = q_i$ for all $i = 1 \ldots k$

Example: Jane Austen + Emulator

|          | a   | an | this | that | with | without |
|----------|-----|----|------|------|------|---------|
| Ch. 1 + 6 | 101 | 11 | 15   | 37   | 28   | 10      |
| Ch. 12 + 24 | 83 | 29 | 15   | 22   | 43   | 4       |

Ch 1 + 6 $\sim$ Multinomial $(202, (p_1, \ldots p_6))$

Ch 12 + 24 $\sim$ Multinomial $(196, (q_1, \ldots q_6))$

Test $H_0: p_i = q_i$ for all $i = 1 \ldots 6$

GLRT statistic

$$-2 \log \Lambda = 2 \sum_{r=1}^{k} N_i \log \frac{N_i}{E_i} + M_i \log \frac{M_i}{F_i}$$

$N_i, M_i$ are the observed counts

$E_i = n \cdot \hat{P}_{0i} \qquad F_i = m \cdot \hat{P}_{0i}$

MLEs are $\hat{P}_{0,i} = \frac{N_i + M_i}{n + m}$

Dimension of full model : $5 + 5 = 10$

Dimension of sub model : $5$

$k = 10 - 5 = 5$ so compare against $\chi^2_5(\alpha)$

Find p-value to be 0.0014 <span style="color:orange">← very different!</span>

Data: $\mathcal{Y} = (Y_1 \dots Y_n)$

Parametric model — $f(\mathcal{Y}|\Theta)$ is the joint PDF/PMF of our data dependent on $\Theta \in \mathbb{R}^k$

Log-likelihood function — $\ell(\Theta) = \sum_{i=1}^{n} \log f(Y_i|\Theta)$

MLE — $\hat{\Theta}$ that maximizes $\ell(\Theta)$

Recall: If $Y_1 \dots Y_n \overset{iid}{\sim} f(y|\Theta)$

$$J(\Theta) = \text{Var}_\Theta\left[\frac{\partial}{\partial \Theta} \log f(Y|\Theta)\right] = -\mathbb{E}_\Theta\left[\frac{\partial^2}{\partial \Theta^2} \log f(Y|\Theta)\right]$$

For large $n$, the MLE $\hat{\Theta}$ is approximately distributed as

$$N\left(\Theta_0, \frac{1}{nJ(\Theta_0)}\right)$$

Define $I_{\mathcal{Y}}(\Theta) = nJ(\Theta) = \text{Var}_\Theta\left[\sum_{i=1}^{n}\frac{\partial}{\partial \Theta} \log f(Y_i|\Theta)\right] = \text{Var}_\Theta[\ell'(\Theta)]$

$$= -\mathbb{E}_\Theta\left[\sum_{i=1}^{n}\frac{\partial^2}{\partial \Theta^2} \log f(Y_i|\Theta)\right] = -\mathbb{E}_\Theta[\ell''(\Theta)]$$

Generally for $\mathcal{Y} \sim f(\mathcal{Y}|\Theta)$, define Fisher information of all observations

$$I_{\mathcal{Y}}(\Theta) = \text{Var}_\Theta[\ell'(\Theta)] = -\mathbb{E}_\Theta[\ell''(\Theta)]$$

Under regularity assumptions, the MLE $\hat{\Theta}$ has approximately distribution $N\left(\Theta_0, \frac{1}{I_{\mathcal{Y}}(\Theta_0)}\right)$ ← Estimate when data isn't iid

For multiple parameters $\Theta \in \mathbb{R}^k$, we define $I_{\mathcal{Y}}(\Theta) \in \mathbb{R}^{k \times k}$

$$I_{\mathcal{Y}}(\Theta)_{ij} = \text{Cov}_\Theta\left[\frac{\partial}{\partial \Theta_i}\ell(\Theta), \frac{\partial}{\partial \Theta_j}\ell(\Theta)\right]$$

$$= -\mathbb{E}_\Theta\left[\frac{\partial^2}{\partial \Theta_i \partial \Theta_j}\ell(\Theta)\right]$$

Again, $\hat{\Theta}$ is $\approx N(\Theta_0, I_{\mathcal{Y}}(\Theta_0)^{-1})$

## Bradley Terry Model

Example: NBA has 30 basketball teams

Each team plays 82 games

How can we rank teams?

① Naively: Count number of wins against number of losses

However, each team plays each other team between 2 and 4 times

② Bradley-Terry: Represent strength of team $i$ by $B_i \in \mathbb{R}$

If game played between teams $i$ and $j$, outcome is random and depends on $B_i$ and $B_j$

Outcome — Bernoulli $(P_{ij})$

$$\log \frac{P_{ij}}{1-P_{ij}} = B_i - B_j \implies P_{ij} = \frac{e^{B_i - B_j}}{1 + e^{B_i - B_j}} = \frac{e^{B_i}}{e^{B_i} + e^{B_j}}$$

- Each $B_i$ is only meaningful relative to other $B_j$'s

We can add any constant $C \in \mathbb{R}$ to all $B_i$'s without changing the model

Allows us to select a team as a standard (set to $0$)

in this case $B_j$ is relative strength to the standardized teams

- Can specify order of each game s.t. the first team is always the home team

Incorporate home team advantage by adding an additional intercept

$$\log \frac{P_{ij}}{1-P_{ij}} = \alpha + B_i + B_j$$

# Estimation and Inference

$K = 30$ (number of teams)

$n = 1580$ (number of total games)

Questions: Estimate $\alpha$ and $B_1, B_2, \dots B_K$ (constraining $B_{\text{nets}} = B_1 = 0$)

Test the null hypothesis $H_0 : \alpha = 0$ (no home team advantage)

Obtain confidence interval

We observe $n$ games $(i_1, j_1), (i_2, j_2), \dots (i_n, j_n)$ where $i$ is always the home team

Outcomes: $Y_1, Y_2, \dots Y_n \in \{0, 1\}$

$$Y_m = \begin{cases} 1 & \text{if } i \text{ beat } j \text{ in the } m\text{th game} \\ 0 & \text{otherwise} \end{cases}$$

Log-likelihood:

$$\text{lik}(\alpha, B_2, \dots B_K) = \prod_{m=1}^{n} P_{i_m j_m}^{Y_m} (1 - P_{i_m j_m})^{Y_m} = \prod_{m=1}^{n} (1 - P_{i_m j_m}) \left( \frac{P_{i_m j_m}}{1 - P_{i_m j_m}} \right)^{Y_m}$$

$$\ell(\alpha, B_2 \dots B_K) = \sum_{m=1}^{n} Y_m \underbrace{\log \frac{P_{i_m j_m}}{1 - P_{i_m j_m}}}_{\substack{\text{log - odds} \\ = \alpha + B_{i_m} + B_{j_m}}} + \log(1 - P_{i_m j_m})$$

$$= \sum_{m=1}^{n} Y_m [\alpha + B_{i_m} + B_{j_m}] - \log\left(1 + e^{\alpha + B_{i_m} - B_{j_m}}\right)$$

① Common approach to estimate $\Theta$ $(\alpha, B_2, \dots B_K)$ via MLE

Solve:

$$0 = \frac{\partial \ell}{\partial \alpha} = \sum_{m=1}^{n} \left[ Y_m - \frac{e^{\alpha + B_{i_m} - B_{j_m}}}{1 + e^{\alpha + B_{i_m} - B_{j_m}}} \right]$$

$$0 = \frac{\partial \ell}{\partial B_i} = \sum_{m : i_m = i} \left[ Y_m - \frac{e^{\alpha + B_{i_m} - B_{j_m}}}{1 + e^{\alpha + B_{i_m} - B_{j_m}}} \right] + \sum_{m : j_m = i} \left[ -Y_m + \frac{e^{\alpha + B_{i_m} - B_{j_m}}}{1 + e^{\alpha + B_{i_m} - B_{j_m}}} \right]$$

No closed form so we solve it numerically

Gradient: $\nabla \ell(\Theta) = \left( \frac{\partial \ell}{\partial \alpha}, \dots \frac{\partial \ell}{\partial B_K} \right)$

Hessian: $\nabla^2 \ell(\Theta) = \begin{pmatrix} \frac{\partial^2 \ell}{\partial \alpha^2} & \cdots & \frac{\partial^2 \ell}{\partial \alpha \partial B_K} \\ \vdots & \ddots & \\ \frac{\partial \ell^2}{\partial \alpha \partial B_K} & & \frac{\partial \ell^2}{\partial B_K^2} \end{pmatrix}$

Newton – Raphson : $\Theta^{(t+1)} = \Theta^{(t)} - \left( \nabla^2 \ell(\Theta^{(t)}) \right)^{-1} \nabla \ell(\Theta^{(t)})$

② To test $H_0 : \alpha = 0$

Use the GLRT sub-model where $\alpha = 0$

$$0 = \frac{\partial \ell}{\partial B_i} = \sum_{m : i_m = i} \left[ Y_m - \frac{e^{B_{i_m} - B_{j_m}}}{1 + e^{B_{i_m} - B_{j_m}}} \right] + \sum_{m : j_m = i} \left[ -Y_m + \frac{e^{B_{i_m} - B_{j_m}}}{1 + e^{B_{i_m} - B_{j_m}}} \right]$$

Solve numerically using Newton-Raphson

Compute $-2 \log \Lambda = 2\ell(\hat{\alpha}, \hat{\beta}_2, \dots \hat{\beta}_K) - 2\ell(0, \hat{\beta}_{0_2}, \dots \hat{\beta}_{\theta_K})$

Compare against $\chi^2_1(\alpha)$

③ Confidence Interval for $B_i - B_j$:

Center the interval at $\hat{\beta}_i - \hat{\beta}_j$ with the MLE's

Estimate standard error of $\hat{\beta}_i - \hat{\beta}_j$

$$Var[\hat{\beta}_i - \hat{\beta}_j] = Cov[\hat{\beta}_i - \hat{\beta}_j, \hat{\beta}_i - \hat{\beta}_j]$$

$$= Cov[\hat{\beta}_i, \hat{\beta}_i] - Cov[\hat{\beta}_i, \hat{\beta}_j] - Cov[\hat{\beta}_j, \hat{\beta}_i] + Cov[\hat{\beta}_j, \hat{\beta}_j]$$

$$= Var[\hat{\beta}_i] + Var[\hat{\beta}_j] - 2 Cov[\hat{\beta}_i, \hat{\beta}_j]$$

$$\approx (I_y(\theta)^{-1})_{ii} + (I_y(\theta)^{-1})_{jj} - 2(I_y(\theta)^{-1})_{ij}$$

Holds for large $n$

$$I_y(\theta) = -\mathbb{E}_\theta[\nabla^2 \ell(\theta)] \quad \text{where } \nabla^2 \ell(\theta) \text{ is the hessian}$$

Can estimate standard error of $\hat{\beta}_i - \hat{\beta}_j$ as

$$\hat{se} = \sqrt{(I_y(\theta)^{-1})_{ii} + (I_y(\theta)^{-1})_{jj} - 2(I_y(\theta)^{-1})_{ij}}$$

We expect $\hat{\beta}_i - \hat{\beta}_j$ to be approx. normal for large $n$.

$$\hat{\beta}_i - \hat{\beta}_j \pm \hat{se} \cdot z(\alpha/2)$$

②' Test $H_0: \alpha = 0$

Randomly permute $(c_m, J_m)$ for each game

Compute a test statistic $T$ on the permuted data

$$T = -2\log \Lambda$$

③' Non parametric bootstrap

Resample $(i_m^*, J_m^*)$ with replacement

Estimate $\hat{\beta}_i - \hat{\beta}_j$ using bootstrap samples

Avoid model misspecification

Model for binary classification

Example: Predict whether a user will click on an ad

Have $n$ ad impressions

- Binary response $Y_i = \begin{cases} 1 & \text{if clicked} \\ 0 & \text{otherwise} \end{cases}$

- Other features
  - size of ad
  - user age      } represented by a collection of $p$ covariates
  - etc.

Logistic Regression: $n$ responses are independent

$$Y_i \sim \text{Bernoulli}(p)$$

log odds: $\log \dfrac{P_i}{1-P_i} = \alpha + B_1 X_{i_1} + \cdots B_p X_{i_p}$

$\alpha$ is the intercept or baseline log-odds

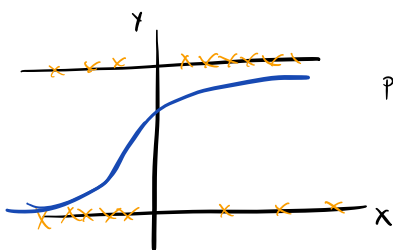$B_j$ represents the change in log odds for a one unit change in $X_{i_j}$

Parameters: $(\alpha, B_1, \dots B_p)$

$$\mathbb{P}[Y_i = 1] = P_i = \dfrac{e^{\alpha + B_1 X_{i_1} + \cdots B_p X_{i_p}}}{1 + e^{\alpha + B_1 X_{i_1} + \cdots B_p X_{i_p}}}$$

Assume one covariate $p=1$ for simplicity

$$B = B_1 \quad \text{and} \quad X_i = X_{i_1}$$

Data $= (X_1, Y_1), \dots (X_n, Y_n)$



$$p(x) = \dfrac{e^{\alpha + Bx}}{1 + e^{\alpha + Bx}}$$

## Estimation and Inference

① Estimating the regression coefficients

② Estimating the "conversion" probability $p(x) = \dfrac{e^{\alpha + Bx}}{1 + e^{\alpha + Bx}}$ for a new impression with covariate $x$

③ Provide a 95% confidence interval for $B$

④ Test a null hypothesis $H_0 : B = 0$

① Estimating $(\alpha, B)$ via MLE

$$\text{lik}(\alpha, B) = \prod_{i=1}^{n} P_i^{Y_i} (1-P_i)^{1-Y_i} \quad \text{where} \quad P_i = \dfrac{e^{\alpha + Bx_i}}{1 + e^{\alpha + Bx_i}}$$

$$= \prod_{i=1}^{n} (1-P_i)\left(\dfrac{P_i}{1-P_i}\right)^{Y_i}$$

$$\Rightarrow \ell(\alpha, B) = \sum_{i=1}^{n} Y_i \underbrace{\log \dfrac{P_i}{1-P_i}}_{\alpha + Bx_i} + \underbrace{\log(1-P_i)}_{-\log(1+e^{\alpha + Bx_i})}$$

$$= \sum_{i=1}^{n} \left(Y_i(\alpha + Bx_i) - \log\left(1 + e^{\alpha + Bx_i}\right)\right)$$

MLE: Set derivatives to 0

$$0 = \dfrac{\partial \ell}{\partial \alpha} = \sum_{i=1}^{n} (Y_i - P_i)$$

$$0 = \frac{\partial \ell}{\partial B} = \sum_{r=1}^{n} (Y_i - P_i) X_i$$

No closed form solution so we once again use Newton–Raphson Method

$$\begin{pmatrix} \alpha^{(t+1)} \\ B^{(t+1)} \end{pmatrix} = \begin{pmatrix} \alpha^{(t)} \\ B^{(t)} \end{pmatrix} - \left( \nabla^2 \ell \left( \alpha^{(t)}, B^{(t)} \right) \right)^{-1} \cdot \nabla \ell \left( \alpha^{(t)}, B^{(t)} \right)$$

$$\nabla \ell (\alpha, B) = \begin{pmatrix} \frac{\partial \ell}{\partial \alpha} \\ \frac{\partial \ell}{\partial B} \end{pmatrix} = \sum_{i=1}^{n} (Y_i - P_i) \begin{pmatrix} 1 \\ X_i \end{pmatrix}$$

$$\nabla^2 \ell (\alpha, B) = \begin{pmatrix} \frac{\partial^2 \ell}{\partial \alpha^2} & \frac{\partial^2 \ell}{\partial \alpha \partial B} \\ \frac{\partial^2 \ell}{\partial \alpha \partial B} & \frac{\partial \ell^2}{\partial B^2} \end{pmatrix} = \sum_{i=1}^{n} -P_i(1 - P_i) \begin{pmatrix} 1 & X_i \\ X_i & X_i^2 \end{pmatrix}$$

*Skipped Algebra*

Interpretation: The update $(\alpha^{(t+1)}, B^{(t+1)})$ solves a least-squares problem

$$\underset{\alpha, B}{\arg\min} \sum_{i=1}^{n} P_i^{(t)} (1 - P_i^{(t)}) \left( z_i^{(t)} - (\alpha + B X_i) \right)^2$$

$$z_i^{(t)} = \alpha^{(t)} + B^{(t)} X_i + \frac{Y_i - P_i^{(t)}}{P_i^{(t)} (1 - P_i^{(t)})}$$

Check: To minimize, we would set derivatives $\alpha, B$ to 0

$$0 = \sum_{i=1}^{n} P_i^{(t)} (1 - P_i^{(t)}) \cdot 2 \left( \alpha + B X_i - z_i^{(t)} \right)$$

$$= 2 \left[ (\alpha - \alpha^{(t)}) \cdot \sum_{i=1}^{n} P_i^{(t)} (1 - P_i^{(t)}) + (B - B^{(t)}) \cdot \sum_{i=1}^{n} P_i^{(t)} (1 - P_i^{(t)}) X_i - \sum_{i=1}^{n} (Y_i - P_i^{(t)}) \right]$$

$$0 = \sum_{i=1}^{n} P_i^{(t)} (1 - P_i^{(t)}) \cdot 2 X_i \left( \alpha + B X_i - z_i^{(t)} \right)$$

$$= 2 \left[ (\alpha - \alpha^{(t)}) \cdot \sum_{i=1}^{n} P_i^{(t)} (1 - P_i^{(t)}) X_i + (B - B^{(t)}) \sum_{i=1}^{n} P_i^{(t)} (1 - P_i^{(t)}) X_i^2 - \sum_{i=1}^{n} (Y_i - P_i^{(t)}) X_i \right]$$

$$\iff \sum_{i=1}^{n} (Y_i - P_i^{(t)}) \cdot \begin{pmatrix} 1 \\ X_i \end{pmatrix} = \left[ \sum_{i=1}^{n} P_i^{(t)} (1 - P_i^{(t)}) \begin{pmatrix} 1 & X_i \\ X_i & X_i^2 \end{pmatrix} \right] \cdot \begin{pmatrix} \alpha - \alpha^{(t)} \\ B - B^{(t)} \end{pmatrix}$$

Also called Iterative Reweighted Least Squares (glm in R)

② Estimate $p(x) = \dfrac{e^{\alpha + Bx}}{1 + e^{\alpha + Bx}}$ : Use plug-in estimates

$$\hat{p}(x) = \frac{e^{\hat{\alpha} + \hat{B}x}}{1 + e^{\hat{\alpha} + \hat{B}x}}$$

③ Confidence Interval for B

Model Based Approach: Compute Fisher information to estimate std error of $\hat{B}$

Fisher Information for all $n$ observations is

$$I_y(\alpha, B) = -\mathbb{E}\left[\nabla^2 \ell(\alpha, B)\right] = \sum_{i=1}^{n} P_i(1-P_i)\begin{pmatrix} 1 & x_i \\ x_i & x_i^2 \end{pmatrix} \qquad P_i = \frac{e^{\alpha + Bx_i}}{1 + e^{\alpha + Bx_i}}$$

Covariance of $(\hat{\alpha}, \hat{B})$ is approximately $I_y(\alpha, B)^{-1}$ for large $n$

Estimate $P_i$ by $\hat{P}_i = \dfrac{e^{\hat{\alpha} + \hat{B}x_i}}{1 + e^{\hat{\alpha} + \hat{B}x_i}}$ and estimate $I_y(\alpha, B)$ by

$$\hat{I}_y(\hat{\alpha}, \hat{B}) = \sum_{i=1}^{n} \hat{P}_i(1-\hat{P}_i)\begin{pmatrix} 1 & x_i \\ x_i & x_i^2 \end{pmatrix}$$

$$\hat{se} = \sqrt{\left(I_y(\alpha, B)^{-1}\right)_{22}}$$

95% Confidence interval is $\hat{B} \pm z(0.025) \cdot \hat{se}$


④ To test $H_0: B = 0$

GLRT: Need to compute MLE $\hat{\alpha}_0$ for $\alpha$ in the null model where $B = 0$

log-likelihood:

$$\ell(\alpha) = \sum_{i=1}^{n} Y_i \underbrace{\log \frac{P_i}{1-P_i}}_{\alpha} + \underbrace{\log(1-P_i)}_{-\log(1+e^{\alpha})}$$

$$= \sum_{i=1}^{n} \left(Y_i \alpha - \log(1 + e^{\alpha})\right)$$

$$0 = \frac{\partial \ell}{\partial \alpha} = \sum_{i=1}^{n}\left(Y_i - \frac{e^{\alpha}}{1 + e^{\alpha}}\right)$$

$$\Longleftrightarrow \hat{\alpha}_0 = \log \frac{\bar{Y}}{1 - \bar{Y}}$$

GLRT test statistic is $-2\log \Lambda = 2\ell(\hat{\alpha}, \hat{B}) - 2\ell(\hat{\alpha}_0, 0)$

Compare against $\chi^2_1$ null distribution


⑤ Diagnostic of model-fit is based on residual

Pearson Residual $\dfrac{Y_i - \hat{P}_i}{\sqrt{\hat{P}_i(1 - \hat{P}_i)}}$

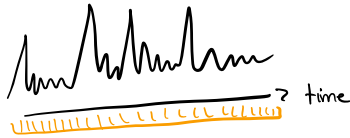If logistic regression is correct, we expect

    mean 0

    variance 1 ← Common to not follow this

    uncorrelated w/ covariate $x_i$

Overdispersion is a common mis-specified model outcome variance greater than 1

Example: Neuron spiking overtime



→ time

Basic Model: Poisson Process

# spikes in the $i^{th}$ bin = $Y_i \sim$ Poisson $(\lambda_i, \Delta)$ independent for $i=1,2,\ldots n$

- WLOG $\Delta = 1$
- $\lambda_i$ is the spiking rate in the $i^{th}$ bin (Depends on external Stimuli)

- Encode stimuli using $p$ covariates

Poisson Log-linear model

$$\log \lambda_i = \alpha + B_1 X_{i_1} + \ldots B_p X_{i_p}$$

$\alpha$ is the intercept, $B_1 \ldots B_p$ are the regression coefficients

Equivalently,

$$\lambda_i = e^{\alpha + B_1 X_{i_1} + B_p X_{i_p}}$$

The responses are then $Y_i \sim$ Poisson $(\lambda_i)$ independent for $i=1,\ldots,n$

To simplify notation, $P=1$

Data: $(X_1, Y_1), \ldots, (X_n, Y_n)$

Consider fixed X's with random responses $Y_i$

Model: $Y_i \sim$ Poisson $(\lambda_i)$  $\lambda_i = e^{\alpha + B X_i}$

Estimation and Inference

- Estimate regression coefficients
- Provide 95% Confidence interval for B

Likelihood Function:

$$lik(\alpha, B) = \prod_{i=1}^{n} \frac{\lambda_i^{Y_i} e^{-\lambda_i}}{Y_i!} \quad \text{where} \quad \lambda_i = e^{\alpha + B x_i}$$

$$\ell(\alpha, B) = \sum_{i=1}^{n} Y_i \underbrace{\log(\lambda_i)}_{\alpha + B X_i} - \lambda_i - \log(Y_i!)$$

$$= \sum_{i=1}^{n} Y_i(\alpha + B X_i) - e^{\alpha + B X_i} - \log(Y_i!)$$

Solving MLE via derivatives

$$0 = \frac{\partial \ell}{\partial \alpha} = \sum_{i=1}^{n} \left( Y_i - e^{\alpha + B X_i} \right) = \sum_{i=1}^{n} (Y_i - \lambda_i)$$

$$0 = \frac{\partial \ell}{\partial B} = \sum_{i=1}^{n} (Y_i - \lambda_i) X_i$$

No closed form so we can apply Newton-Raphson Method

Gradient of $\ell(\alpha, \beta)$:

$$\nabla \ell(\alpha, \beta) = \begin{pmatrix} \frac{\partial \ell}{\partial \alpha} \\ \frac{\partial \ell}{\partial \beta} \end{pmatrix} = \sum_{i=1}^{n} (Y_i - \lambda_i) \begin{pmatrix} 1 \\ X_i \end{pmatrix}$$

Hessian of $\ell(\alpha, \beta)$

$$\nabla^2 \ell(\alpha, \beta) = \begin{pmatrix} \frac{\partial^2 \ell}{\partial \alpha^2} & \frac{\partial^2 \ell}{\partial \alpha \partial \beta} \\ \frac{\partial^2 \ell}{\partial \alpha \partial \beta} & \frac{\partial^2 \ell}{\partial \beta^2} \end{pmatrix} = - \sum_{r=1}^{n} w_i \begin{pmatrix} 1 & X_i \\ X_i & X_i^2 \end{pmatrix} \qquad w_i = \lambda_i = e^{\alpha + \beta X_i}$$

$\Rightarrow$ Newton Raphson Iterations

$$\begin{pmatrix} \alpha^{(t+1)} \\ \beta^{(t+1)} \end{pmatrix} = \begin{pmatrix} \alpha^{(t)} \\ \beta^{(t)} \end{pmatrix} - \left( \nabla^2 \ell(\alpha^{(t)}, \beta^{(t)}) \right)^{-1} \cdot \nabla \ell(\alpha^{(t)}, \beta^{(t)})$$

Newton - Raphson update solves a weighted least- squares regression

Also called iterative reweighted least square

To make a model based 95% Confidence Interval for B

- Fischer Information of all $n$ observations

$$\mathcal{I}_y(\alpha, \beta) = - \mathbb{E}\left[ \nabla^2 \ell(\alpha, \beta) \right]$$

<span style="color:orange">$\leftarrow$ full log-likelihood</span>

<span style="color:orange">(can't use $n \cdot \mathcal{I}$ since each is different)</span>

$$= \sum_{r=1}^{n} \lambda_i \begin{pmatrix} 1 & X_i \\ X_i & X_i^2 \end{pmatrix}$$

Use plug-in estimate

$$\mathcal{I}_y(\hat{\alpha}, \hat{\beta}) = \sum_{r=1}^{n} \hat{\lambda}_i \begin{pmatrix} 1 & X_i \\ X_i & X_i^2 \end{pmatrix}$$

Estimate standard error of $\hat{B}$ by

$$\hat{se} = \sqrt{\left( \mathcal{I}_y(\hat{\alpha}, \hat{\beta}) \right)^{-1}_{22}}$$

Model Diagnostics:

- Based on normalized Residuals

  Pearson Residual : $\dfrac{Y_i - \hat{\lambda}_i}{\sqrt{\hat{\lambda}_i}}$

  If model is correctly specified we expect

  mean 0

  variance 1

  uncorrelated with $X_i$

# Generalized Linear Model

A GLM is a model for responses $Y_i \ldots Y_n$ where $Y_i \sim f(Y|\Theta)$, $\Theta \in \mathbb{R}$

$\Theta_i$ is the value for observation $i$

For observation $i$ we observe $p$ covariates $x_i \ldots x_{ip}$

GLM assumes

$$g(\Theta_i) = \alpha + B_1 x_{i_1} + \ldots B_p x_{ip}$$

$g: \mathbb{R} \to \mathbb{R}$ is the link function

In logistic regression, $g \to$ log odds

poisson log-linear, $g \to \log \lambda$

Alternative link functions

$$g(p) = \Phi^{-1}(p)$$

$\Phi$ is the standard normal quantile function

Choice of link function

- Goodness of model fit to data

- Interpretation of model / parameters

- Mathematical convinience

Change of variable $\Theta \to \eta(\Theta)$ so pdf/pmf has the form

$$f(y|\eta) = e^{\eta y - A(y)} \cdot h(y) \qquad \leftarrow \text{Canonical or natural link}$$

- For Bernoulli $(p)$: The PMF is

$$f(y) = p^y (1-p)^{1-y} = (1-p) \left( \frac{p}{1-p} \right)^y$$

$$= e^{\log\left(\frac{p}{1-p}\right) y + \log(1-p)}$$

Set $\eta(p) = \log \frac{p}{1-p}$, $A(\eta) = -\log(1-p) = \log\left(1 + e^\eta\right)$ and $h(y)=1$

$$p = \frac{e^\eta}{1+e^\eta}$$

- For poisson : The PMF is

$$f(y) = \frac{\lambda^y e^{-y}}{y!} = e^{(\log \lambda) y - \lambda} \cdot \frac{1}{y!}$$

Set $\eta(\lambda) = \log \lambda$  $A(\eta) = \lambda$  and $h(y) = \frac{1}{y!}$

$$\lambda = e^\eta$$

$$f(y|\eta) = e^{\eta y - A(\eta)} h(y) \quad \text{is called the exponential family form of the model}$$

$\eta$ is the natural parameter of the model

For a GLM based on a parametric model $f(Y|\Theta)$, the natural/canonical link is the choice $g(\Theta) = \eta(\Theta)$ where $\eta$ is natural parameter

If we use the natural link:

Log-likelihood function for $(X_1, Y_1) \ldots (X_n, Y_n)$ is

$$\ell(\alpha, \beta) = \log \prod_{i=1}^{n} e^{\eta_i y_i - A(\eta_i)} h(y_i)$$

$$= \sum_{i=1}^{n} Y_i \eta_i - A(\eta_i) + \log(h(y_i))$$

$\uparrow$

$\alpha + \beta x_i$ when we use the natural link

$$= \sum_{i=1}^{n} Y_i(\alpha + \beta x_i) - A(\alpha + \beta x_i) + \log h(y_i)$$

Computing MLE:

$$0 = \frac{\partial \ell}{\partial \alpha} = \sum_{i=1}^{n} Y_i - A'(\eta_i)$$

$$\eta_i = \alpha + \beta x_i$$

$$0 = \frac{\partial \ell}{\partial \beta} = \sum_{i=1}^{n} Y_i x_i - x_i A'(\eta_i)$$

$$\Rightarrow \quad 0 = \sum_{i=1}^{n} (Y_i - A'(\eta_i)) \begin{pmatrix} 1 \\ x_i \end{pmatrix}$$

What is $A'(\eta)$

$$1 = \sum_Y f(y|\eta) = \sum_Y e^{\eta y - A(\eta)} h(y)$$

differentiate w.r.t $y$

$$0 = \sum_Y (Y - A'(\eta)) \cdot \underbrace{e^{\eta y - A(\eta)} h(y)}_{f(y|\eta)}$$

$$= E[Y] - A'(\eta)$$

$$\Rightarrow A'(\eta) = E[Y]$$

For GLM, $A'(\eta) = E[Y_i]$ is the model prediction for mean of $Y_i$ and $Y_i - A'(\eta)$ is a residual

Fisher Information is

$$I_y(\alpha, B) = -\mathbb{E}\left[\nabla^2 \ell(\alpha, B)\right]$$

$$= \sum_{i=1}^{n} w_i \begin{pmatrix} 1 & x_i \\ x_i & x_i^2 \end{pmatrix} \quad \text{where } x_i = A''(\eta_i)$$

## 4/27 : Proportional Hazards Model

Example: Clinical trial studying the effect of a cancer drug produces data

$$T_1, T_2, \ldots T_n$$

where $T_i$ is the time of remission

For the $i^{th}$ patient we have $p$ covariates

ex. Treatment vs. Control

Stage of Cancer

Age

Family history

Goal: Model $T_i$ via $x_{i_1} \ldots x_{i_p}$

Let $T$ be a continuous variable

Hazard Function

$$\lambda(t) = \lim_{\delta \to 0} \frac{1}{\delta} \mathbb{P}\left[T \in [t, t+\delta] \mid T \geq t\right] \quad \leftarrow \text{Instantaneous Risk}$$

$$\lambda(t) = \lim_{\delta \to 0} \frac{\frac{1}{\delta}\mathbb{P}\left[T \in [t, t+\delta] \mid T \geq t\right]}{\mathbb{P}[T \geq t]} = \frac{f(t)}{1 - F(t)}$$

$\underset{\text{Conditional probability}}{\smile}$

Example: when $T \sim$ Exponential $(\Theta)$

$$f(t) = \Theta e^{-\Theta t}$$
$$F(t) = 1 - e^{-\Theta t}$$

$$\lambda(t) = \frac{\Theta e^{-\Theta t}}{1 - (1 - e^{-\Theta t})} = \Theta$$

Hazard function is constant in time $\quad \leftarrow \text{Memoryless Property!!}$

In general the hazard may depend on time and covariates for each patient

Model $T_i$ via



$$\lambda_i(t) = \lambda(t) \exp\left(B_1 x_{i_1} + \cdots B_p x_{i_p}\right)$$
$\qquad\qquad\quad \uparrow \qquad\qquad\quad \uparrow$
$\qquad\qquad$ time $\qquad\qquad$ covariates
$\qquad$ dependence

Same time dependence for each patient

Shape of hazard function will have the same shape over time but it will be scaled based on the covariates

Origin of the name proportional hazards

This model is semi-parametric

parametric component: Regression Coefficient

Non-parametric component: hazard function

Usually we care more about $B_1 \ldots B_p$ than $\lambda(t)$

Idea: Condition on the set of all observed reoccurrence terms

$$t_{(1)} < t_{(2)} < \cdots < t_{(n)}$$

Fixes times at which the $n$ reoccurence events occured but not the patient

Inference for $B_1, \ldots B_p$ will then be based on only the information about which patient had reoccurance

    Akin to permutation tests

For each $t_{(k)}$, let $R_{(k)}$ be the at risk set immediately before time $t_{(k)}$ (Patients still in remission before $t_k$)

Conditional on some patient in $R_{(k)}$ having reoccurance at time $t_{(k)}$, the probability it is patient $i$ for $i \in R_{(k)}$

$$\frac{\lambda_i(t_{(k)})}{\sum\limits_{j \in R_{(k)}} \lambda_j(t_{(k)})}$$  ← Ratio of the instantaneous rate of reoccurance in patient $i$ to the sum of risk for all patients

The baseline $\lambda(t)$ cancels and we find

$$\frac{\lambda_i(t_{(k)})}{\sum\limits_{j \in R_{(k)}} \lambda_j(t_{(k)})} = \frac{\exp(B_1 x_{i_1} + \cdots B_p x_{ip})}{\sum\limits_{j \in R_{(k)}} \exp(B_1 x_{j_1} + \cdots B_p x_{jp})}$$

For each time $t_{(k)}$ $\quad k = 1, \ldots, n$

    Partial likelihood function

$$\text{plik}(B_1 \ldots B_p) = \prod_{k=1}^{n} \frac{\exp(B_1 x_{k_1} + \cdots B_p x_{k_p})}{\sum\limits_{j \in R_{(k)}} \exp(B_1 x_{j_1} + \cdots B_p x_{jp})}$$

    Use partial likelihood in place of actual likelihood

Note: Typically the responses $T_1 \ldots T_n$ are right-censored

    If the $i^{th}$ patient is still in remission by the end of the trial we do not know $T_i$

        just that $T_i \geq l_i$

      $l_i$ is the duration of the trial

      If cancer never occurs, then $T_i = \infty$

  When some responses are right-censored, the at risk set $R_{(k)}$ is defined as the set of patients

      — Still in remission

      — Not right censored

Maximum Likelihood Estimation

    Assume $p=1$ for simplicity

      Data: $(X_1, T_1) \ldots (X_n, T_n)$

    log-partial-likelihood

$$\ell(B) = \log \prod_{k=1}^{n} \frac{\exp(B X_{i_k})}{\sum\limits_{j \in R_{(k)}} \exp(B X_j)} = \sum_{k=1}^{n} \left( B X_{i_k} - \log \sum_{j \in R_{(k)}} \exp(B X_j) \right)$$

$$\hat{B} = \arg \max \ell(B)$$

Solve via Newton–Raphson

Variance is estimated by

$$I(\hat{B})^{-1} \simeq \left( -\frac{\partial^2 \ell(\hat{B})}{\partial B^2} \right)^{-1}$$

Test of $H_0 : B = 0$ based on GLRT statistic

$$-2 \log \Lambda = 2\ell(\hat{B}) - 2\ell(0)$$

w/ distribution $\chi_1^2$

# Course Review

Hypothesis Testing: Deciding whether a particular null hypothesis about the underlying distribution is true/false

Estimation: Estimating quantities/parameters related to this distribution

# Hypothesis Testing

Binary decision to accept/reject $H_0$ based on data $X_1 \ldots X_n$ using a test statistic $T(X_1 \ldots X_n)$

Step 1: How to choose $T$?

Step 2: Deciding whether $H_0$ is true/false based on $T$?

Neyman-Pearson Lemma: Maximize power via likelihood ratio statistic

If monotonic increasing/decreasing we can use its input in its place

Simple Alternatives

Composite Alternatives: