

Computer Age Statistical Inference. CHAPTER 01

★ Statistics: Science about learning from data.

- i) Inference - measuring the uncertainty around those estimates
- ii) Algorithm - following a set of data ~~processes~~ ~~steps~~ to produce estimates

e.g. Say you have observations x_1, \dots, x_n

▷ The mean $\bar{x} = \frac{1}{n} \sum x_i$ summarizes x_1, \dots, x_n in a number.

▷ $\hat{\sigma}_x = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2 / (n(n-1))} = \frac{1}{\sqrt{n}} \sigma$

the standard error of the mean \bar{x} . aka. \hat{s}_e

- σ is the standard deviation of the sample.

* Here, averaging is the algorithm, while the standard error provides an inference of the algorithm's accuracy.

▷ \hat{s}_e is itself an algorithm, which could be subject to further inferential analysis.

e.g. Linear Regression

▷ both linear regression & the mean minimize the error.

↳ however, linear regression is a conditional mean.

$$\bar{y} = E(y) \quad \downarrow \text{expectation of } y$$

$\hat{y} = E(y|x)$ → conditional mean

unconditional mean.

e.g.: $y = \hat{\beta}_0 + \hat{\beta}_1 x$ → is the regression model

↳ Age

↳ Liver.

▷ Least-Squares Algorithm: $\sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$

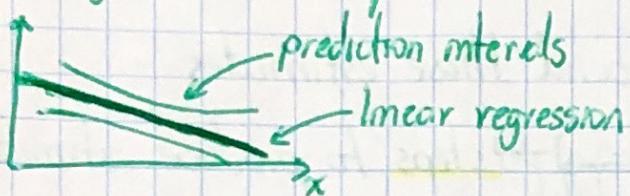
y_i are the actual observed y 's

Since we are estimating two parameters from the data

▷ Inference is:

$$s_{\text{LSQ}}^2 = \left[\frac{1}{n-2} \sum_{i=1}^n (x_i - \bar{x})^2 \right] \frac{1}{2}$$

Q: Why aren't the prediction intervals the same width across the data?



A) i) the local variance is not the same
ii) we don't have the same data points at every x^*

e.g. Lowess (aka Local Regression) Algorithm

This method uses a weighting function ~~with a window~~ with the effect that the influence of a neighboring value on the smoothed value at a certain position decreases with their distance to that position.

Step 1: The range of the points, to be included is determined. The number n of these points is specified.

- the larger n is, the smoother the adapted curve will be.
- the range shall be exactly n values, including the selected point. and the selected point is in the center of the selected range.
- it's possible that the number of points to the left and right may differ; i.e. the 1st value only has neighboring point to the right.

Step 2: Establish the weighting for the locally regression smoothing.

$$w(x_k) = \left(1 - \left|\frac{x_i - x_k}{d_i}\right|^3\right)^3 \quad \text{for } k = 1, \dots, n \quad \begin{matrix} \text{(Proximity)} \\ \text{(Weighting)} \end{matrix}$$

d_i is the distance from x_i to the k^{th} ~~neighboring~~ neighboring point
The largest value is attained at point x_c (the focal point)
The value 0 is at the range limits.

Step 3: A linear regression function based on the least squares method is estimated for the LOESS procedure. The LOESS procedure uses a quadratic function

The parameters are calculated in such a way the the following function is minimized.

$$\text{LOWESS: } \hat{y}_k = a + b x_k$$

$$\text{LOESS: } \hat{y}_k = a + b x_k + c x_k^2$$

Step 4: (optional) It is possible that the estimated regression function is influenced by outliers.

Therefore, robust weightings can be determined.

$$G(x_k) = \begin{cases} \left(1 - \left(\frac{|y_k - \hat{y}_k|}{6 \text{ median}(|y_i - \hat{y}_i|)}\right)^2\right)^2 & \left|\frac{|y_k - \hat{y}_k|}{6 \text{ median}(|y_i - \hat{y}_i|)}\right| \leq 1 \\ 0 & \left|\frac{|y_k - \hat{y}_k|}{6 \text{ median}(|y_i - \hat{y}_i|)}\right| \geq 1 \end{cases}$$

The robust weighting is 0 if a residual is greater than or equal to $6m$ ($m = \text{median residuals}$). This eliminates the outliers' influence.

The robust weightings are multiplied with the proximity weightings and reused to re-estimate a regression function within the individual ranges.

$$\text{LOWESS: } \sum_k w(x_k) G(x_k) (y_k - a - b x_k)^2$$

$$\text{LOESS: } \sum_k w(x_k) G(x_k) (y_k - a - b x_k - c x_k^2)^2$$

Step n: A series of new smoothed values, is the result. The procedure can be repeated several times - the higher the number of iterations, the more precise the curve fitting.

> Bootstrapping (aka sampling with replacement) - INFERENCE
(Bootstrap Confidence Intervals)

16 Aug 2019

- Under the bootstrap principle: $\sigma_{(\text{sampling w/replacement})}^2 \sim \sigma_{(\text{across samples})}^2$

- let B (usually in the hundreds or thousands) be the number of bootstrap iterations.

- If our original sample is $(x_i, y_i)_{i=1}^N$, at each bootstrap iteration, we'll have:

$$(x_{j(b)}, y_{j(b)})_{j \in \mathbb{I}} \text{ for } b = 1, \dots, B \quad \& \quad I = \{1, \dots, N\}$$

j is the index that is randomly sampled from the i indexes.

- For each b , compute $\hat{Y}_{j(b)} = E[Y_{j(b)} | x_{j(b)}]$

↳ using our LOWESS model (or any other model)

- The purpose is to output a table such as the following:

j	b	1	2	3	...	B
1		$\bar{Y}_{1(1)}$	$\bar{Y}_{1(2)}$	$\bar{Y}_{1(3)}$...	$\bar{Y}_{1(B)}$
2		$\bar{Y}_{2(1)}$	$\bar{Y}_{2(2)}$	$\bar{Y}_{2(3)}$...	$\bar{Y}_{2(B)}$
:		:	:	:	:	:
N		$\bar{Y}_{N(1)}$	$\bar{Y}_{N(2)}$	$\bar{Y}_{N(3)}$...	$\bar{Y}_{N(B)}$

→ we are calculating \bar{Y} (the mean) since j can be repeated

→ also, due to replacement, ~~the~~ j is not continuous.

- For each j in the table, compute \hat{S}_{boot} or $\hat{\sigma}_{j,boot}$

$$\hat{\sigma}_{j,boot} = \sqrt{\left(\bar{Y}_j - \bar{\bar{Y}}_j\right)^2 / (B-1)}$$

} Then compute the confidence intervals assuming a t distribution.

⇒ then sort by value of $\bar{Y}_{j(b)}$ from min to max

↳ then take the 5th & 95th values to get a 90% confidence interval.

↳ ~~the~~ sorting does not assume any stat. distribution, but assumes that the min and max values are indeed the min & max that we can ever get.

⇒ the OLS confidence intervals work asymptotically, that means they assume that the number of ~~the~~ available observations is infinite

↳ but it assumes normality.

⇒ In LOWESS, n is not infinite, but it doesn't assume any distribution.

⇒ The OLE confidence intervals ~~the~~ width are minimized in the mean.

★ Big Data Implications on HYPOTHESIS TESTING:

e.g. following the example from the book,

72 leukemia patients, where:

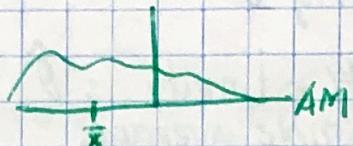
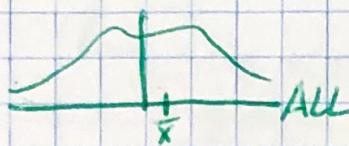
47 have ALL (acute lymphoblastic leukemia) - aka bad leukemia

25 have AML (acute myeloid leukemia) - aka worse leukemia

↳ measured 7128 genes \Rightarrow find whose genes are relevant to the diagnosis.

\rightarrow let's take gene 136.

\rightarrow let's say the distribution of gene 136 looks like:



So let's find out if the mean of the two diagnosis is different.

\rightarrow the standard classical answer is to use a two-tailed t-test for hypothesis

$$H_0: \mu_{ALL} \neq \mu_{AML}$$

\rightarrow The t-statistic is:

$$t = \frac{\bar{Y}_{136}^{AML} - \bar{Y}_{136}^{ALL}}{\hat{s}_e} ; \hat{s}_e = \hat{\sigma} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)^{1/2} ; n_1 = 47 \text{ & } n_2 = 25$$

if we assume same variances.

The degrees of freedom ν is given by

$$\nu = n_1 + n_2 - 2 = 70$$

Computer Age Statistical Inference CHAPTER 02

- For Frequentist Inference, we assume that the observed sample comes from an probability distribution F

- we call our data vector $\mathbf{x} = (x_1, \dots, x_n)$, which is a realization of a vector of random variables $\mathbf{X} = (X_1, \dots, X_n)$, which are distributed according to F .

X is language to say "a sample" & \mathbf{x} is the sample's value

PROBABILITY

So: F is a probability distribution

X are individual draws of F

x are the very particular values for each of the X 's we got from our sample.

▷ We want to learn properties from F , based on x (vector of realizations)

Usually the expectation of a single random draw X_i from F .

→ If I get a number from that distribution, what can I expect that number to be, and how much it deviates (se)

→ let's call the expectation property θ

$\theta = E_F(X_i)$ / The obvious ~~WW~~ estimate is $\hat{\theta} = \bar{x}$
equal to the sample average.

θ = the true expectation value of any draw. X_i

$\hat{\theta}$ = is the best guess of what θ is

\bar{x} = is the mean.

\bar{x} is the algorithm; $\hat{\theta}$ is the estimation; θ is what is being estimated

$\hat{\theta}$ is calculated from x following an algorithm $\hat{\theta} = t(x)$

for instance, $t(x) = \sum x_i / n = \bar{x}$

↳ t is any algorithm (like sample mean \bar{x} or sample std dev)

▷ $\hat{\theta}$ is sample specific and is a realization of $\hat{\theta} = t(x)$

computing the algorithm to any sample

$\hat{\theta}$ is a number

$\hat{\theta}$ is accurate only to a degree which we want to quantify

▷ Define $\mu = E_F(\hat{\theta})$, the expected value ~~WW~~ (across samples) of producing an estimate using algorithm $t(x)$ when samples come from distribution F

➤ Bias & variance of estimate $\hat{\theta}$ of parameter θ to be:

$$\text{bias} \stackrel{d}{=} \mu - \theta \quad \& \quad \text{var} \stackrel{d}{=} E_F \{ (\hat{\theta} - \mu)^2 \}$$

→ **variance**: if we follow algorithm $t(c)$, how much are the values of Θ going to vary with respect to the expected value, μ

→ bias: expected value - true value

↳ given an algorithm) is that algorithm estimating the parameter properly.

> Bias-variance trade-off:

- ↳ • The estimator can be very precise, but it can change a lot from sample to sample (no bias, high variance)
• The estimator can be a simple algorithm that does not varie a lot from sample to sample. But are not able to describe all the data points (biased, low variance)

- Frequentist usually define quantities (parameters) with respect to an infinite sequence of trials.

In other words: Hypothetical datasets $X^{(1)}, X^{(2)}, \dots$ generate infinite $\Theta^{(1)}, \Theta^{(2)}, \dots$ samples

→ So frequentist calculate properties of estimator $\hat{\theta} = t(x)$

↳ However we don't know F...

- Frequentist principles to conduct inference.

i) Plugin principle $\Rightarrow \text{se}(\bar{x}) = [\text{var}_f(X)/n]^{1/2}$ / $\text{se}(\bar{x}) = \sigma_{\bar{x}}$

$$\hookrightarrow \bar{x} = \sum x_i / n$$

However, $\text{var}_F(x) = \hat{\text{var}}_F = \left[\frac{\sum (x_i - \bar{x})^2}{n(n-1)} \right]$ without bias

aka. relate the sample standard error with the true variance

Recall: θ is the true parameter of F

$\hat{\theta}$ is the estimate of θ

$\hat{\theta}^t$ is the estimate of θ , using algorithm $t(\cdot)$

If F is $N(\mu, \sigma^2)$ aka a normal distribution

↳ θ would be μ or σ^2

→ for a gamma distribution $T(\alpha, \beta)$. θ would be α or β .

» The purpose of Frequentist Inference is to assess the inference of the estimates. $\hat{\theta}$. So we need the se of $\hat{\theta}$

↳ However the se calculation requires many samples $\hat{\theta}^{(1)}, \hat{\theta}^{(2)}, \dots$ but we only have one...

» Plug-in principle: (in human words)

if we took one x from the distribution F , I can link the variance of that x to the variance of my sample mean.

→ the variance around the mean of a sample $\hat{\theta}$ is ~~the~~ an estimator of the variance of F .

$$\hat{var}_F = \left[\frac{\sum (x_i - \bar{x})^2}{n-1} \right] \Leftrightarrow \sigma_{\bar{x}} = se(\bar{x}) = \left[\frac{var_F(x)}{n} \right]^{1/2} \dots \textcircled{1}$$

↳ (-1) because we already estimated \bar{x} using the data

» Taylor Series Approximations:

relate $t(x)$ by local linear approximations.

$$\text{e.g. if } \hat{\theta} = \bar{x}^2 \rightarrow \frac{d\hat{\theta}}{d\bar{x}} = 2\bar{x}$$

Since \bar{x} is constant, $\sigma_{\bar{x}^2} = 2|\bar{x}| \sigma_{\bar{x}}$

we can compute the standard error of the transformed estimator.

In other words you can get the standard error of $\hat{\theta}$ by taking $d\hat{\theta}/d\bar{x}$ so that ~~the~~:

$$\sigma_{\hat{\theta}} = \sigma_{\bar{x}} \frac{d\hat{\theta}}{d\bar{x}}$$

To compute the $se(\bar{x}^2)$ of an estimator that is not \bar{x} and the estimator a differential function of the original estimator,

3 Parametric Families & MLE (maximum likelihood estimate)

e.g. → assume family $N(\mu, \sigma^2)$, where $\sigma^2 = 5$
 → given x , find μ such that the probability of $N(\hat{\mu}, 5)$ is maximized.

where $x = [x_1, x_2, \dots, x_n]$

$$P(x_i | N(\hat{\mu}, 5)) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \Rightarrow \text{probability to observe } x_i$$

If each element x was sample independently, we can multiply their probabilities.

$$\hookrightarrow P(x_1, x_2 | N(\mu, 5)) = P(x_1 | N(\mu, 5)) \cdot P(x_2 | N(\mu, 5))$$

↓

$$P(x | N(\mu, \sigma^2)) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \prod_{i=1}^n e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

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

↪ The Likelihood Function.

$L(x)$ is the probability of observe the whole sample x .

4 MLE → find $\hat{\mu}$, such that $L(x)$ is maximized. for normal dist
 Tune $\hat{\mu}$ (in this case μ)

$$\max_{\mu} L(x) \Rightarrow \max_{\mu} \left[\frac{1}{\sqrt{2\pi\sigma^2}} \right]^n e^{\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}} \Rightarrow \hat{\mu}^{\text{MLE}}$$

4 Simulation & Bootstrap

Implement a repeated sequence of trials

→ Start with an estimate \hat{F} of F , then simulate values from \hat{F} to get a prior sample $\hat{\Theta}^{(k)} = t(x^{(b)})$ by computing the algorithm $t(\cdot)$ on each of the bootstrap iterations, for $b = 1, \dots, B$

→ Then the empirical standard deviation of the $\hat{\Theta}$'s is the frequentist estimate for $\sigma_{\hat{\Theta}}$

Y6 Pivotal Statistics.

e.g. say we got two samples $x_1 = [x_{11}, x_{12}, \dots, x_{1n}]$ and $x_2 = [x_{21}, x_{22}, \dots, x_{2n}]$.

→ Assume $x_{1i} \stackrel{\text{ind}}{\sim} N(\mu, \sigma^2) \quad i = 1, 2, \dots, n,$

$$x_{2i} \stackrel{\text{ind}}{\sim} N(\mu_2, \sigma^2) \quad i=1, 2, \dots, n_2$$

→ Test the null hypothesis $H_0: \mu_1 = \mu_2$

$$\hat{\theta} = \bar{x}_2 - \bar{x}_1 \rightarrow \text{difference of means test-statistic}$$

$$\bar{X}_2 - \bar{X}_1 \sim N(\mu^2, \sigma^2) ; \text{ where } \mu^2 = 0, \sigma^2 = \sigma^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$$

$$\sigma_{\bar{x}_2}^2 = se(\bar{x}_2) = \frac{\sigma^2}{n_2}$$

* σ^2 is unknown

So, estimate σ^2 using some sample variances.

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^{n_1} (x_{1i} - \bar{x}_1)^2 + \sum_{i=1}^{n_2} (x_{2i} - \bar{x}_2)^2}{(n_1 + n_2 - 2)}$$

→ However, if we have more occurrences of this problem, say x_3, x_4, \dots for each hypothesis test iteration,

$\hat{\sigma}^2$ changes (every time) for each sample x_j

→ The issue. In hypothesis testing we compute a sample statistic, and we compare that value to a critical value on a theoretical distribution.

↳ The theoretical distribution will be different for each x_j set of samples

→ The ~~unknown~~ statistician Student tried to solve that...

▷ t-student

$$t = \frac{\bar{x}_2 - \bar{x}_1}{\hat{s}_d}$$

\hat{s}_d is the estimation of the standard deviation.

$$\hat{s}_d = \hat{\sigma} \left(\frac{1}{n_1} + \frac{1}{n_2} \right) \dots \text{use ①}$$

Using t-student, the estimator does not depend on the other parameters of the distribution

it allows to test $\bar{x}_1 - \bar{x}_2$ (and come with confidence intervals) without the estimation of $\sigma_{\bar{x}_1 - \bar{x}_2}^2 = \sigma^2$

t-test statistic is a pivotal statistic as it does not depend on other parameters the distribution might have.

★ So. Frequentist use pivotal statistics whenever they are available to conduct stat tests.

27 Aug 2019

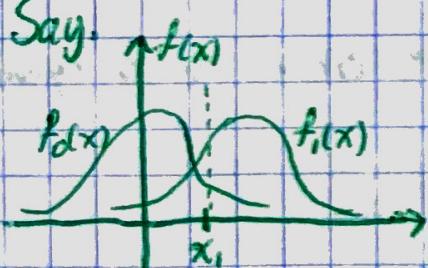
▷ Frequentist Optimality

- Neyman-Pearson lemma provides an optimum hypothesis-testing algorithm.

→ NP lemma states:

- 1) if we are to decide between 2 possible density functions for observed data x
→ a null hypothesis density $f_0(x)$
→ and an alternative density $f_1(x)$

Say:



for the simple case
 $x = x_1$

So, with what confidence level can we say
 x_1 does not belong to f_0 (rejecting H_0)?

Determine the hypothesis test that is powerful at:

rejecting, but not rejecting when it's the appropriate thing.
(no false positives)

ii) A testing rule $t(x)$ says which choice (0 or 1) we'll make given x .

For any decision we make, there will be two associated errors α and β .

$\alpha = \Pr_{f_0} \{t(x) = 1\} \Rightarrow$ the probability of rejecting the null hypothesis when x belongs to the null hypothesis (incorrectly rejecting f_0)

$\beta = \Pr_{f_1} \{t(x) = 0\} \Rightarrow$ choosing f_0 when actually f_1 generated x (probability of β incorrectly failing to reject)

iii) Let $L(x)$ be the Likelihood ratio

$$L(x) = \frac{f_1(x)}{f_0(x)}$$

what is the likelihood of the alternative hypothesis $f_1(x)$ giving the data x

vs. what is the likelihood of the null hypothesis $f_0(x)$ giving your data x

iv) Testing rule $t_c(x)$ by:

$$t_c(x) = \begin{cases} 1 & \text{if } \ln(L(x)) \geq c \\ 0 & \text{if } \ln(L(x)) < c \end{cases} \quad \text{As only rules of the } t_c(x) \text{ form can be optimal}$$

There is one rule such rule for each choice of the cutoff c .

→ that is $\alpha_c < \alpha$ & $\beta_c < \beta$ for some c and where α & β are generated by another rule.

if $t_c(x)$ has an α_c with a c such that the α error α_c is equal to α , then $\beta_c < \beta$ (and vice versa)

e.g.

We want to get c such that α is equal to Type I error rate and Type II error β is minimized.

$$f_0 \sim N(0, 1) \quad \not\sim f_1 \sim N(\gamma_2, 1)$$

$$f_0(x_1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x_1^2}{2}} \quad f_1(x_1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x_1 - \mu_0)^2}{2}}$$

for normal dists
density functions.

$$f_0(x) = \left[\frac{1}{\sqrt{2\pi}} \right]^n \prod_{i=1}^n e^{-\frac{x_i^2}{2}} = \left[\frac{1}{\sqrt{2\pi}} \right]^n e^{-\frac{1}{2} \sum_{i=1}^n x_i^2}$$

$$f_1(x) = \left[\frac{1}{\sqrt{2\pi}} \right]^n e^{-\frac{1}{2} \sum_{i=1}^n (x_i - \mu_1)^2}$$

$$L(x) = \frac{e^{-\frac{1}{2} \sum_{i=1}^n (x_i - \mu_1)^2}}{e^{-\frac{1}{2} \sum_{i=1}^n x_i^2}}$$

$$L(x) = e^{-\frac{1}{2} \left[\sum_{i=1}^n x_i + \frac{n}{4} \right]}$$

$$L(x) = e^{-\frac{1}{2} \left[n\bar{x} + \frac{n}{4} \right]}$$

$$L(x) > c \Rightarrow e^{-\frac{1}{2} \left[n\bar{x} + \frac{n}{4} \right]} > c \Rightarrow \bar{x} \text{ shall be greater than some constant.}$$

$$-\frac{1}{2} \left[n\bar{x} + \frac{n}{4} \right] > c$$

$$n\bar{x} + n/4 > c_2$$

$$\bar{x} > c_3 \rightarrow \text{only the mean depends on the sample.}$$

➤ If two likelihood functions come from a normal distribution, then the only thing we need for the hypothesis testing is the sample mean \bar{x}

➤ Most powerful hypothesis test at any Type I error rate α is to compare \bar{x} to a constant.

→ let's relate $c, \alpha & \beta$...

$$\alpha = P(\bar{x} > c | \theta = 0) \text{ where } \theta = \mu = 0$$

↳ probability under $f_0 \sim N(0, \sigma^2)$ where $\sigma^2 = 1$

$$\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n} = \frac{1}{n} \Rightarrow \sigma_{\bar{x}}^2 = \frac{1}{\sqrt{n}}$$

$$\text{Then } \frac{\bar{x} - 0}{1/\sqrt{n}} = \frac{\bar{x} - \mu}{\sigma_{\bar{x}}^2} = \bar{x}\sqrt{n} \sim N(0, 1)$$

normalize

$$\bullet \alpha = P(\bar{x} > c | \mu = 0)$$

$$\alpha = P(\bar{x}\sqrt{n} > c\sqrt{n} | \mu = 0)$$

$$\alpha = 1 - P(\bar{x}\sqrt{n} \leq c\sqrt{n} | \mu = 0)$$

$$\alpha = 1 - \Phi(c\sqrt{n})$$

↳ function for the CDF of a $N(0,1)$ distribution.

$$\Phi(c\sqrt{n}) = 1 - \alpha$$

$$c\sqrt{n} = \Phi^{-1}(1 - \alpha)$$

In general:

$$c = \frac{1}{\sqrt{n}} \Phi^{-1}(1 - \alpha)$$

$$c = \mu_0 + \frac{1}{\sqrt{n}} \Phi^{-1}(1 - \alpha)$$

for two normal distributions with the same variance

$\triangleright 1 - \beta$: Power of hypothesis test

↳ Probability of correctly ~~rejecting~~ rejecting

- A hypothesis test for which $1 - \beta$ is maximized (β is minimized) at level α is called a Most Powerful Test (MPT)

On the other hand:

$$\bullet \beta = P(\bar{x} \leq c | \mu = \frac{1}{2})$$

↳ fail to reject the null hypothesis when the alternative hypothesis is true.

$$\frac{\bar{x} - \frac{1}{2}}{\frac{1}{\sqrt{n}}} \sim N(0, 1) \quad \left. \right\} \text{normalize to a } N(0, 1) \text{ distribution}$$

$$\beta = P[(\bar{x} - \frac{1}{2}) \cdot \sqrt{n} \leq (c - \frac{1}{2}) \cdot \sqrt{n} | \mu = \frac{1}{2}]$$

$$\beta = \Phi((c - \frac{1}{2})\sqrt{n})$$

Smallest TypeII error ratio for hypothesis test with α TypeI error rate is given by:

$$\beta = \mathbb{E}[(e - \frac{1}{2}) \sqrt{n}]$$

30 Aug 2019

Computer Age Statistical Inference CHAPTER 03

→ Bayesian Inference

$$\text{Bayes Rule: } P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

- like in Frequentist Inference, the fundamental unit of inference is a family of probability densities.

$$F = \{f_{\mu}(x); x \in X, \mu \in \Omega\}$$

x is a point in the sample space X
 μ is an unobserved point in the parameter space Ω

→ we observe x from $f_{\mu}(x)$ & infer μ .

→ Bayesian Inference also assumes the knowledge of a prior density $g(\mu)$, $\mu \in \Omega$

e.g. if μ is the average height of a room, I would need to state a priori (before observing the sample x)

→ Bayesian works for cases where we don't have a lot of data but we have expert knowledge.

↳ in frequentist, with few data the confidence intervals would be very thick (wide)

- Bayes Theorem (aka Bayes Rule) is a rule to combine prior knowledge in $g(\mu)$ with current evidence in x .

→ How to consistently update our belief of μ ?

Let $g(\mu|x)$ denote the posterior density of μ

↳ after observing the data vector X

Bayes Rule posterior distribution

$$g(\mu|x) = \frac{g(\mu) \cdot f_{\mu}(x)}{f(x)} \quad \text{where } f(x) \text{ is the marginal density of } x$$

$$f(x) = \int_{\Omega} f_{\mu}(x) \cdot g(\mu) d\mu$$

prior distribution

likelihood that μ is the true parameter given x

probability of having observed the vector (data) x after averaging across all the possible values that μ can take.

Frequentists assumes that there is a true distribution where the data comes from. Also assumes that just like the sample you got, there are an infinite number of them

In Bayesian, x is fixed (only one sample)
The only thing that changes is the belief of μ

Bayes Rule can be written as:

$$g(\mu|x) = c_x L_x(\mu) g(\mu) \quad \text{if } f_{\mu}(x) = L_x(\mu) = f(x|\mu)$$

constant that only depends on x

aka. normalization constant of the posterior dist.

as any probability distribution integrates to 1

for any two μ_1, μ_2 values on Ω , the ratio of posterior densities is given by:

$$\frac{g(\mu_1|x)}{g(\mu_2|x)} = \frac{g(\mu_1)}{g(\mu_2)} \cdot \frac{f_{\mu_1}(x)}{f_{\mu_2}(x)} \Rightarrow \text{The posterior odds ratio is the prior odds ratio times the Likelihood ratio.}$$

The relative probability (aka. the odds ratio) of μ_1 being the right value of the parameter μ , given the sample x

is equal to the relative odds $\frac{\mu_1}{\mu_2}$ of the prior belief times the odds ratio of the likelihood function from the observed data x .

e.g. An engineer knows she's having twins. She asks what's the probability that they'll be identical. The doctor says: $\frac{1}{3}$ of twin births are identical } prior belief

Say x is a sonogram result (either same sex or opposite sex) and same sex is observed.

(identical twins always are of the same sex, while fraternals have 0.5 probability of same or different sex)

- How does the sonogram modify our prior of $1/3$ probability on the twins being identical?

$$\frac{g(I, S)}{g(F, S)} = \frac{g(I)}{g(F)} \cdot \frac{f_I(S)}{f_F(S)} \quad \begin{matrix} I = \text{identical twins.} \\ S = \text{same sex.} \\ F = \text{fraternal.} \end{matrix}$$

$$= \frac{\frac{1}{3}}{1 - \frac{1}{3}} \cdot \frac{1}{2} \rightarrow \begin{matrix} \text{if they are identical, what's the prob. of being same sex} \\ \text{" fraternal, "} \end{matrix} \quad \begin{matrix} \text{same sex} \\ \text{same sex} \end{matrix}$$

$$= 1$$

↳ $g(I, S)$ & $g(F, S)$ are equally likely.

⇒ Fraternal & Identical are equally likely

	S	D	joint probability distribution
I	$\frac{1}{3}$	0	$\frac{1}{3} = g(I)$
F	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{2}{3} = 1 - \frac{1}{3} = g(F)$

e.g. The binomial distribution counts the number of successes or failures in n trials

⇒ If a random variable is distributed binomially

$$x \sim B(n, p) \Rightarrow P(X = k | p, n)$$

↳ number of successes being k

p = probability of success

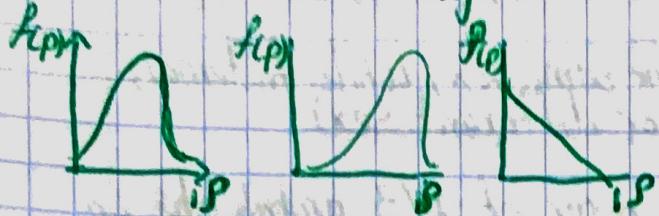
n = number of observations (trials)

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} = \frac{T(n+1)}{T(k+1) \cdot T(n-k+1)}$$

$$P(X=k|n,p) = \binom{n}{k} p^k (1-p)^{n-k}$$

Probability that the number of successes is exactly k \hookrightarrow observing $n-k$ failures
 \hookrightarrow observing k successes

> let's assume that before observing our data, we belief p to look either of these ways:



\Rightarrow all of these can be modeled using the Beta family of distributions.

So let's model our prior distribution on p as:

$$\underset{\text{prior}}{\Pi(p|\alpha, \beta)} \sim \text{Beta}(\alpha, \beta) = \frac{p^{\alpha-1} (1-p)^{\beta-1}}{B(\alpha, \beta)} ; 0 \leq p \leq 1$$

where $B(\alpha, \beta)$ is the Beta function

$$B(\alpha, \beta) = \frac{\Gamma(\alpha) \Gamma(\beta)}{\Gamma(\alpha+\beta)} \text{ where } \Gamma \text{ is the gamma function}$$

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx ; z > 0 \quad \hookrightarrow \frac{1}{B(\alpha, \beta)} \text{ is the normalizing constant}$$

\hookrightarrow is the ^{continuous} generalization of the factorial function.

$$\Gamma(n+1) = n!$$

03 Sep 2019

$$\text{Also: } B(\alpha, \beta) = \int_0^1 p^{\alpha-1} (1-p)^{\beta-1} dp$$

> if $\Pi(p|\alpha, \beta) \sim \text{Beta}(\alpha, \beta)$, then the mean of the Beta distribution is given by:

$$\underset{\text{The mean}}{E(p|\alpha, \beta)} = \frac{\alpha}{\alpha + \beta}$$

For the following analysis is convenient to parameterize our Beta distribution such that the mean can be represented using only one parameter... So, let...

$$\mu = \frac{\alpha}{\alpha + \beta} ; M = \alpha + \beta \Rightarrow \text{prior } \Pi(\rho | \mu, M) = \text{Beta}(M\mu, M(1-\mu))$$

↳ under this parametrization the expected value E is.

$$E(\rho | \mu, M) = \mu ; V(\rho | \mu, M) = \mu(1-\mu) / M+1$$

↳ expected value of ρ

↳ variance of ρ

▷ Using this new prior, we can proceed to get the posterior distribution.

$$P(\rho | k) \propto (\text{aka. is proportional to}) \ l(k | \rho) \cdot \Pi(\rho | \mu, M)$$

↳ posterior dist.
our belief about ρ
after observing k successes

↳ likelihood

↳ prior function

* → we are applying the Bayes rule, without the normalizing constant $1/B(\alpha, \beta)$
because it doesn't alter the distribution "shape".
- Bayes rule only needs the proportionality.

$$\text{▷ then } P(\rho | k) \propto (\rho^k (1-\rho)^{n-k}) (\rho^{M\mu-1} (1-\rho)^{M(1-\mu)-1})$$

↳ without $\binom{n}{k}$ because the same reason as *

$$\propto \rho^{(k+M\mu)-1} (1-\rho)^{n-k+M(1-\mu)-1}$$

$$\propto \text{Beta}[k+M\mu, n-k+M(1-\mu)] \Rightarrow \text{probability density function}$$

▷ if the prior for ρ is Beta & the likelihood is Binomial, then the posterior for ρ is also Binomial.

↳ that is a conjugate prior (aka. the posterior is of the same family as the prior)

▷ the expected value of the posterior given we have observed k success is:

$$E(\rho, k) = \frac{k+M\mu}{n+M}$$

↳ size of the sample

$$x = [0, 0, 1, 0, 1, \dots, 0, 1]$$

n

$$k = \sum x_i$$

k = number of successes

M, μ = the two prior parameters, which is equivalent to have observed k successes in the past (before observing the sample)

M = is equivalent to having had a sample of size n in the past. (aka. the size of the prior sample)

μ is the expectation for p

> we can say our prior assumed having seen $M\mu$ successes in M trials before running the current experiment.

> What is our prediction of the number of successes given our posterior for p ?

- Using conditional probability, we can write:

$$m(k|\mu, M, n) = \int_0^1 \ell(k|p) \cdot P(p|\mu, M, k) \cdot dp$$

↓
marginal distribution of the # of successes in n future trials given parameters μ, M & n

$$= \frac{T'(M)}{T'(M\mu) \cdot T'(M(1-\mu))} \binom{n}{k} \int_0^1 p^{k+M\mu-1} (1-p)^{n-k+M(1-\mu)-1} dp$$

This is the Beta Function

$$= \frac{T'(M)}{T'(M\mu) \cdot T'(M(1-\mu))} \cdot \frac{T'(n+1)}{T'(k+1) \cdot T'(n-k+1)}.$$

$$\frac{T(k+M\mu) \cdot T(n-k+M(1-\mu))}{T(n+M)} \Rightarrow \text{probability mass function}$$

= The probability distribution for the number of successes you are expecting to observe.

successes

$$\alpha = k - M\mu$$

= is the posterior predicted distribution

as the random variable (number of successes k) is discrete.

$$\beta = n - k + M(1-\mu)$$

failures

> look for "Conjugate Prior" tables for more examples.

▷ Flaws in Frequentist Inference

I.e.g. (in frequentist inference) Say that an ongoing experiment is being run each month. \downarrow observation at month t

An independent normal variate is observed $X_t \sim N(\mu, 1)$

You're planning to run an hypothesis test to test $H_0: \mu = 0$ vs. $H_a: \mu > 0$

$$Z_i = \frac{\sum_{t=1}^i X_t / i - 0}{\sqrt{i}}$$

up to month i

$$Z_i = \sum_{t=1}^i X_t / \sqrt{i}$$

is the z-score based on data up to month i

→ Say that at month 30 (the scheduled end of the experiment)

$$Z_{30} = 1.66 > 1.645$$

one sided 95%
for a $N(0, 1)$ distribution, Then reject H_0 .

$t_1(x)$

→ So in this case, our hypothesis testing algorithm is collect data for 30 months, once you have all the data points, compute Z_{30} and compare to a critical value.

hypothesis testing

→ let's say we change the algorithm, to: collect data for 20 months & check Z_{20} , if it is significant stop (if $Z_{20} > 1.645$, stop) if it's not, then collect data from 10 more months and compare Z_{30} to a critical value.

→ In Frequentist, if the algorithm changes, even though the data points stay exactly the same, the significance level is different for each algorithm.

↳ which is a flaw.

→ for $t_1(x)$, $\alpha = 5\%$

- For $t_2(x)$, $\alpha_2 = P(\text{rejecting } H_0 \text{ incorrectly})$

$$\alpha_2 = P(Z_{20} > 1.645 \cup Z_{30} > 1.645 \mid \mu = 0)$$

$$\alpha_2 = P(Z_{20} > 1.645 \mid \mu = 0) + P(Z_{30} > 1.645 \mid \mu = 0)$$

$$- P(Z_{20} > 1.645 \cap Z_{30} > 1.645 \mid \mu = 0)$$

→ Since Z_{20} happens before Z_{30} ; use conditional probability to write the α_2 .

$$\alpha_2 = 0.05 + 0.05 - [1 - P(Z_{30} \leq 1.645 \mid Z_{20} > 1.645, \mu = 0)]$$

$$P(Z_{20} > 1.645 \mid \mu = 0)$$

$$\alpha_2 = 0.1 - [1 - P(Z_{30} \leq 1.645 \mid Z_{20} > 1.645, \mu = 0)] 0.05$$

$$\alpha_2 = 0.1 - [1 - 0.7] 0.05$$

$$\alpha_2 = 0.085$$

Under the algorithm $t_2(x)$, Frequentist would have said the result was not significant

but significant under algorithm $t_1(x)$

On Bayesian inference, we always use the same algorithm called Bayes Rule in which the likelihood function of $\mathbf{x} = (x_1, x_2, x_3, \dots, x_{20}, x_{30})$ is always

HTHLY $L_{\mathbf{x}}(\mu) = \prod_{i=1}^{30} e^{-\frac{1}{2}(x_i - \mu)^2}$ regardless of stopping the experiment

early or not. Stopping early does not affect on the posterior as it only depends on \mathbf{x} through the likelihood function.

II e.g. We have data of prostate cancer comparing 52 patients with 50 healthy controls. Each subject got genetic activity measured for a panel of $n = 6033$ genes.

For each gene, the two-sample t-statistic is

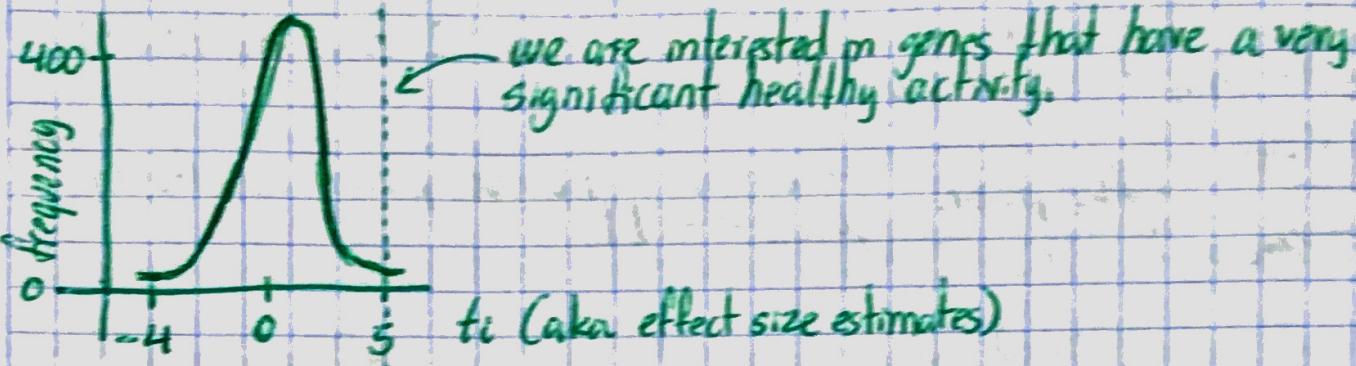
$$t_i = \frac{\bar{x}_{ip} - \bar{x}_{ic}}{s_e^2} \quad \begin{cases} \text{patients' mean for gene } i = \bar{x}_{ip} \\ \text{control's } " = \bar{x}_{ic} \end{cases}$$

Then, (since large sample size) $t_i \sim N(\mu_i, 1) \quad i = 1, \dots, 6033$

where μ_i is the true effect size for gene i

↳ from F

Let's say: the following ~~histogram~~ histogram



Say gene 610 with $t_i = t_{610} = 5.29$ (which is the largest t_i)

If I try to estimate μ_{610} , Is t_{610} OK?

↳ 610 showed a very large value in this case, and that's the reason we picked it to estimate μ ; However if the experiment is conducted again 610 might not be the largest

↳ We biased the estimate of μ by choosing the largest effect size. Since in our sample 610 was particularly high.

↳ ("picking the highest value and take it as significant" is part of the algorithm)

↳ & that creates "selection bias"

↳ 610 might be an overestimate of μ_{610}

↳ selecting genes randomly would ease the overestimation/bias.

Even though gene 610 is individually unbiased for μ_{610} , frequentists would worry that there is an upward bias on $t_{610} = 5.29$.

↳ As frequentist assumes this large value was obtained by choice.

↳ It's likely frequentists would recommend a procedure to correct the bias.

On Bayesian Inference would ignore whether t_{010} was picked because it was large.

However Bayesian Inference is sensitive to the choice of the prior

↳ for instance, the choice of a flat prior (aka "I have no idea"), the result might be $\mu_{010} = 5.29$.

$\pi(\mu_{010})$ } flat prior. However a normal prior results in
 $\mu_{010} = 4.11$

★ Attention shifts from

choosing an algorithm $t(x)$ in frequentist inference to choosing a prior $\pi(\cdot)$ in bayesian inference.

- is much easier to create biases in frequentist inference as any change in the algorithm has an effect on the statistical significance
- ↳ in Bayesian, the prior may be subjective, but it's the only thing that can be biased.

▷ Bayesian vs. Frequentist.

BAYESIAN

- 1) operates only in one sample with the whole parameter space
- 2) requires a prior distribution (past experience)

FREQUENTIST

- 1) operates with one parameter in many samples.
- 2) replaces the choice of a prior with the choice of a method (algorithm $t(x)$)
- 3) is much flexible than Bayes inf. as we can come up with many algorithms.

BAYESIAN

- 4) Bayesian analyses answers all possible questions at once, (because the posterior is a distribution)

FREQUENTIST

- 4) usually only computes the expected value and the variance. (each characteristic requires a specific algorithm.)

10 Sep 2019

Computer Age Statistical Inference CHAPTER 04

→ Fisherian Inference and Maximum Likelihood Estimation (MLE)

- ▷ For a family of probability densities $f_{\theta}(x)$, where θ is a vector of parameters, the log-likelihood function is defined as:

$$l_x(\theta) = \log \{ f_{\theta}(x) \} \quad \text{for a fixed } x \text{ and variable } \theta.$$

(only one sample)
Get the most likely parameters that would have generated the sample we have.

- ▷ the MLE is the value of $\theta \in \Omega$ that maximizes the likelihood function.
parameter space
parameter vector

$$\text{MLE: } \hat{\theta} = \underset{\theta \in \Omega}{\operatorname{argmax}} \{ l_x(\theta) \}$$

- we can also provide MLE estimates for a function $\hat{\theta} = T(\theta)$ using $\hat{\theta} = T\hat{\theta}$. (aka. estimate functions of the true parameter)

- MLE properties

i) it's automatic: likelihood function $\xrightarrow{\text{data}} \text{MLE} \rightarrow \hat{\theta}^{\text{MLE}}$

ii) excellent frequentist properties (good bias & variance)

- bias: $\mu - E(\hat{\mu})$

Expected value of the estimation

→ true value of the parameter

- unbiased estimator \rightarrow bias = 0

↳ [next page...]

↳ expected deviation from the true value.

- variance - is the deviation of the expected value from the true value

$$\text{Variance} = \sum_{i=1}^I (\hat{\mu}^{(i)} - E(\hat{\mu}))^2 \quad / \quad I = \text{number of samples} \\ \hat{\mu}^{(i)} = \mu \text{ of sample } i$$

for a normal distribution

$$= E_F \{ (\mu^{(i)} - E(\hat{\mu}))^2 \}$$

for any distribution's expected value for the MLE's density probability function F .

iii) has a reasonable Bayesian justification.

$$P(\theta|x) = Cx \pi(\theta) e^{\ell(x|\theta)}$$

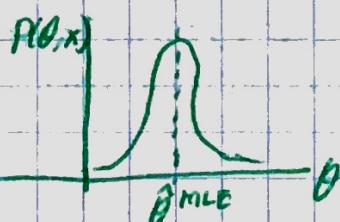
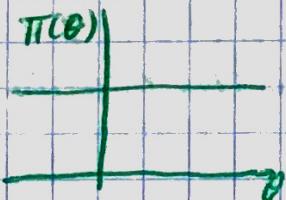
maximum likelihood estimation is the log likelihood function

$$P(\theta|x) \rightarrow \text{posterior} \quad / \quad \theta^{\text{MLE}} \text{ is a maximizer of } P(\theta|x) \\ Cx \rightarrow \text{a constant} \quad / \quad \text{if } \pi(\theta) \text{ is flat (aka. unknown)} \\ \pi(\theta) \rightarrow \text{prior}$$

* - If we assume we know nothing about the parameter we are to estimate, then the MLE will be the highest point of the posterior distribution.

- Remember that in Bayesian we have a posterior distribution of θ .

↳ If $\pi(\theta)$ is flat, then $P(\theta|x)$ shall be as follows.



➢ as the algorithm (described in *) does not change, then the significance level is not affected by unexpected changes in the algorithm (as in frequentist).

e.g. we'll use the Glomerular Filtration Rate data:

- we'll consider 2 potential families

i) Normal: let $\theta = (\mu, \sigma)$, then

$$f_{\theta}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} \Rightarrow \text{density function}$$

ii) Gamma (shifted): let $\theta = (\lambda, \sigma, \nu)$, then

$$f_{\theta}(x) = \frac{(x-\lambda)^{\nu}}{\sigma^{\nu} \Gamma(\nu)} e^{-(x-\lambda)/\sigma} \text{ for } x \geq \lambda, 0 \text{ otherwise}$$

↳ we are shifting the density distribution with λ .

Since

$$f_{\theta}(x) = \prod_{i=1}^n f_{\theta}(x_i) \Rightarrow \text{The Likelihood Function.}$$

- Under iid sampling, we got.

$$l_x(\theta) = \sum_{i=1}^n \log f_{\theta}(x_i) = \sum_{i=1}^n l_{x_i}(\theta)$$

↳ log-likelihood function

- For Normal: $\hat{\mu}^{\text{MLE}} = \bar{x}$

$$\hat{\sigma}^{\text{MLE}} = \left[\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n} \right]^{1/2}$$

iid = independent and identically distributed.

identically = every observation comes from the same distribution.

independent = the value of a previous observation does not influence the next observation.

But $\hat{\mu}^{\text{MLE}}$ and $\hat{\sigma}^{\text{MLE}}$ have a closed form solution.

- For Gamma; there is no closed form solution, so it needs to be numerically maximized in the computer.
(Jupyter Lecture 7 part 1)

MLE can cause overfitting identification problems in high dimensions.

★ ↳ if we fit a lot of parameters in θ that fit would become very specific to our current data set. (& would not represent the population)

- However regularized versions of MLE such as Lasso WML are a workaround to this issue.

~~⇒ Permutation & Randomization~~

> Permutation & Randomization

- Previously we conducted a hypothesis test on the activity of gene 136 on 2 groups of leukemia patients.
- We computed a 2-sample t-test to assess the significance of the effect on gene 136 of ALL diagnosis vs. AML diagnosis.
 - ↳ whenever we conduct a t-test, we assume that the data samples come from a Normal distribution.
 - ↳ however, as the sample is small, it follows a different distribution
- So Fisher suggested the use of randomization to avoid the Normality assumption

↳ Randomization is: taking groups from the data that are ^{random} of the same size as the tested groups.
(in our case $n_1 = 47$ & $n_2 = 25$). Computing the t-statistic for each randomly sampled pair of groups & get their histogram

Utilizing random generated groups, we expect the t values not to be very high, so we can construct an empirical distribution of t values.

13 Sep 2019

~~> Fisher Information and the MLE~~

- We'll go over the univariate case of Fisher Information with one parameter family of densities

$$F = \{ f_\theta(x), \theta \in \Omega, x \in X \}$$

↳ sample space
↳ parameter space
↳ parameter density function

- we'll consider the case of continuous random variables.

► The Log-likelihood is defined as:

$$l_x(\theta) = \log f_{\theta}(x)$$

The derivative of $l_x(\theta)$ with respect to θ is the score function. $i_x(\theta)$

$$i_x(\theta) = \frac{\partial}{\partial \theta} \log f_{\theta}(x) = \frac{f'_\theta(x)}{f_\theta(x)}$$

→ How higher or lower the likelihood value of our sample gets as θ varies.

① Let's compute the expectation of the score function.

$$E(x) = \int_x x f(x) dx$$

density function, so...

$$\begin{aligned} E[i_x(\theta)] &= \int_x i_x(\theta) f_{\theta}(x) dx = \int_x \frac{f'_\theta(x)}{f_\theta(x)} dx \\ &= \int_x \frac{\partial}{\partial \theta} f_{\theta}(x) dx \end{aligned}$$

if $f_{\theta}(x)$ is continuous & continuously differentiable,

$$= \frac{\partial}{\partial \theta} \int_x f_{\theta}(x) dx = \frac{\partial}{\partial \theta} \cdot 1 = 0$$

$$\underline{E[i_x(\theta)] = 0}$$

② Let's compute the variance of $i_x(\theta)$

$$V(x) = \int_x [x - E(x)]^2 f(x) dx$$

$$\begin{aligned} V[i_x(\theta)] &= \int_x [i_x(\theta) - E(i_x(\theta))]^2 f_{\theta}(x) dx \\ &= \int_x [i_x(\theta)]^2 f_{\theta}(x) dx \end{aligned}$$

The Fisher Information is defined as: the variance of the score function

$$I_\theta \triangleq \mathbb{V}[\dot{l}_x(\theta)] = \int_X [\dot{l}_x(\theta)]^2 f_\theta(x) dx$$

③ The MLE estimator of θ : $\hat{\theta}^{\text{MLE}}$

$$\hat{\theta}^{\text{MLE}} \text{ approx } N(\theta, \frac{1}{I_\theta})$$

$\xrightarrow{\text{mean}}$
 $\xrightarrow{\text{variance}}$

► Proof of ③

$$\text{let } \ddot{l}_x(\theta) = \frac{\partial^2}{\partial \theta^2} \log f_\theta(x)$$

$$\begin{aligned} &= \frac{\partial}{\partial \theta} \left[\frac{\dot{l}_\theta(x)}{f_\theta(x)} \right] \\ &= \frac{\ddot{f}_\theta(x) f_\theta(x) - \dot{f}_\theta(x) \dot{f}_\theta(x)}{(f_\theta(x))^2} \\ &= \frac{\ddot{f}_\theta(x)}{f_\theta(x)} - \left(\frac{\dot{f}_\theta(x)}{f_\theta(x)} \right)^2 \end{aligned}$$

$\ddot{l}_x(\theta)$ has expectation

$$\begin{aligned} E[\ddot{l}_x(\theta)] &= \int_X \frac{\ddot{f}_\theta(x)}{f_\theta(x)} f_\theta(x) dx = \int_X \left(\frac{\dot{f}_\theta(x)}{f_\theta(x)} \right)^2 f_\theta(x) dx \\ &= \int_X \ddot{f}_\theta(x) dx - \int_X [\dot{l}_x(\theta)]^2 f_\theta(x) dx \\ &= \frac{\partial^2}{\partial \theta^2} \int_X \int_X f_\theta(x) dx dx - I_\theta \end{aligned}$$

$$E[\ddot{l}_x(\theta)] = -I_\theta$$

► Now suppose $\mathbf{x} = (x_1, x_2, \dots, x_n)$ is a sample from $f_\theta(x)$

Then $\dot{l}_x(\theta) = \sum_{i=1}^n \dot{l}_{x_i}(\theta)$, as $\dot{l}_{x_i}(\theta)$ is in Log space.

$$\ddot{l}_x(\theta) = \sum_{i=1}^n \ddot{l}_{x_i}(\theta)$$

Since $\hat{\theta}^{\text{MLE}}$ for the full sample x satisfies the maximizing condition

$\dot{l}_x(\theta) = 0 \rightarrow$ as we are maximizing the likelihood function, and every maximum has a 1st derivative of 0.

Then, we can get a Taylor Series approximation.

↳ (local linearization)

$$\dot{l}_x(\hat{\theta}^{\text{MLE}}) = 0 \approx \dot{l}_x(\theta) + \ddot{l}_x(\theta)(\hat{\theta}^{\text{MLE}} - \theta)$$

Solve for $\hat{\theta}^{\text{MLE}}$

$$\hat{\theta}^{\text{MLE}} = \frac{\theta \ddot{l}_x(\theta) - \dot{l}_x(\theta)}{2 \ddot{l}_x(\theta)}$$

$$\hat{\theta}^{\text{MLE}} = \theta + (\dot{l}_x(\theta) / \ddot{l}_x(\theta))$$

④

$\dot{l}_x(\theta) = \sum_{i=1}^n \dot{l}_{x_i}(\theta)$ is a sum of a function of random variables.

- Means of random variables are preferable over Sums of random variables.

↳ because means (are random variables) of very large number of random variables follow asymptotically a Normal distribution.

↳ which is the central limit theorem (CLT)

So, by the CLT, $\frac{\dot{l}_x(\theta)}{n} \sim N(E[\dot{l}_x(\theta)], V[\dot{l}_x(\theta)]^{\frac{1}{2}}/n)$
 $\sim N(0, I_{\theta}^{\frac{1}{2}}/n)$ $\quad \frac{I_{\theta}}{n} = \text{variance}$

$$\Rightarrow \text{then } \hat{\theta}^{\text{MLE}} = \theta + \frac{\dot{L}_X(\theta)/n}{\ddot{L}_X(\theta)/n}$$

As $n \rightarrow \infty$, the numerator $\sim N(0, I_d/n)$

and the mean of the function will converge to the expected value

, the denominator

$$\frac{\ddot{E}_x(\theta)}{n} \rightarrow \frac{E_0[\ddot{E}_x(\theta)]}{n} = \cancel{I_{\text{ext}}} - I_0$$

$$\checkmark V(cx) = c^2 V(x)$$

$$\frac{\hat{I}_X(\theta)/n}{\hat{I}_X(\theta)/n} \sim N\left(0, \frac{I_0/n}{I_0^2}\right)$$

$$\sim N\left(0, \frac{1}{n I_0}\right)$$

thus $\hat{\theta}^{\text{MLE}} \sim N(\theta, \frac{1}{n} I_{\theta})$

I_a is the Fisher Information for n iid observations

$$E\left[\theta + \frac{\hat{L}_x(\theta)}{\hat{L}_x(\theta)}\right] = E(\theta) + E\left(\frac{\hat{L}_x(\theta)}{\hat{L}_x(\theta)}\right)$$

$$\stackrel{a}{=} \theta - \theta = \theta$$

So, $\hat{\theta}^{\text{MLE}} \sim N(\theta, \frac{1}{I(\theta)})$ End of Proof. ③

Fisher Information Implementation.

to compute the Fisher Information I_θ ...

Let $x_i \stackrel{\text{ind}}{\sim} N(\theta, \sigma^2)$, where σ^2 is known.

→ we compute Log-likelihood $l_x(\theta)$

~~density function for a normal distribution~~ $f_\theta(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\theta)^2}{2\sigma^2}}$

normal likelihood function $f_\theta(x) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i-\theta)^2}{2\sigma^2}}$



$$\underbrace{l_\theta(x)}_{\sim} = -\frac{1}{2} \sum_{i=1}^n \frac{(x_i - \theta)^2}{\sigma^2} - \frac{n}{2} \log(2\pi\sigma^2)$$

→ the score function

$$\begin{aligned} \underbrace{\hat{l}_\theta(x)}_{\sim} &= -\frac{1}{2} \frac{1}{\sigma^2} 2(-1) \sum_{i=1}^n (x_i - \theta) \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \theta) \end{aligned}$$

→ $\ddot{l}_\theta(x)$

$$\underbrace{\ddot{l}_\theta(x)}_{\sim} = \frac{1}{\sigma^2} (-1) \sum_{i=1}^n 1 = -\frac{n}{\sigma^2} \quad \text{as } E[\underbrace{\ddot{l}_\theta(x)}_{\sim}] = -I_\theta$$

→ Fisher Information

$$I_\theta = \frac{n}{\sigma^2} \quad \text{as } E[\hat{l}_x(\theta)] = 0$$

→

$$E[\hat{l}_x(\theta)] = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \theta) = 0 \quad \text{such that}$$

$$\sum_{i=1}^n x_i = n\theta \Rightarrow \hat{\theta}^{\text{MLE}} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}$$

→ Finally; for a large enough n ,

$$\hat{\theta}^{\text{MLE}} \approx N(\theta, \frac{1}{I_\theta}) \Rightarrow \hat{\theta}^{\text{MLE}} \sim N(\theta, \frac{\sigma^2}{n})$$

↳ standard error of the MLE.

→ To finish our discussion about the properties of MLE in a 1-dimensional setting, suppose that

$\hat{\theta} = t(x)$ $\hat{\theta}$ is an estimator computed by some algorithm $t(x)$,
based on iid sample $x = (x_1, \dots, x_n)$ from $f_\theta(x)$.
↳ that is: $E_\theta\{t(x)\} = \theta$

► Then the Cramér-Rao lower bound says that the variance of $\hat{\theta}$ exceeds the Fisher Information bound

$$\text{Var}_\theta(\hat{\theta}) \geq \frac{1}{n I_\theta}$$

→ MLE Variance.

- In other words, MLE has variance at least as small as the best unbiased estimate of θ .

- Note that MLE is not unbiased in general, but its bias is small. (in the order of $1/n$) → making the comparison with unbiased estimates appropriate.

► Conditional Inference

Say we have an iid sample $x_i \stackrel{\text{iid}}{\sim} N(\theta, 1)$ which has produced estimate $\hat{\theta} = \bar{x}$.

However people conducting the sampling initially disagreed on what the sample size would be. So they flipped a coin to decide:

$$n = \begin{cases} 25 & \text{prob } 1/2 \\ 100 & \text{prob } 1/2 \end{cases}$$

and $n=25$ was decided

Q: what is $\text{Var}_{\bar{x}}$?

Classical Frequentist rationale could have resulted in

$$\left[\frac{1}{2} \frac{\sigma^2}{100} + \frac{1}{2} \frac{\sigma^2}{25} \right]^{1/2} = 0.158 = \sigma_{\bar{x}} = \text{sd}(\bar{x})$$

↳ variance of \bar{x} with $n=25$
↳ variance of \bar{x} with $n=100$

However Conditional Inference would resulted in

$$\left[\frac{\sigma^2}{25} \right]^{1/2} = 0.158 \cdot 0.2 = \sigma_{\bar{x}}$$

- we have the likelihood function without the prior

↳ based on the observations
- "just take the sample you have"

} it's a mix of Freq.
and Bayes thinking

▷ Fisher's arguments for conditional inference:

i) More relevant inferences

(as the inferences only have to do with what really happened)

ii) Simpler inferences

(no need to assess any correlation between the result and the sample size selection)

e.g. Observed Fisher Information.

Rather than using $\hat{\theta} \sim N(\theta, 1/nI_0)$, Fisher suggested using
where $I(x)$ is the observed Fisher Information

$$\hat{\theta} \sim N(\theta, \frac{1}{I(x)}) \quad I(x) \triangleq -\ddot{\ell}_x(\hat{\theta}^{\text{MLE}}) = -\frac{\partial^2}{\partial \theta^2} \ell_x(\theta) \Big|_{\hat{\theta}^{\text{MLE}}}$$

$E[I(x)] = nI_0$, in large samples, the observed Fisher Information is the same as the Fisher Information

↳ However Fisher suggested that in smaller samples $I(x)$ gives a better idea of θ 's accuracy.

► We can confirm Fisher's suggestion by sampling from a distribution.
 - For instance, let's use the Cauchy distribution.
 ↳ doesn't have an expected value, no variance

- 1) Compute the log-likelihood
- 2) get its 1st derivative
- 3) then the 2nd one
- 4) get the Fisher Information
- 5) get the variance of the estimate

even if the distribution doesn't have a variance or an expected value

$$f_{\theta}(x) = \frac{1}{\pi} \frac{1}{1+(x+\theta)^2} \Rightarrow \text{Cauchy distribution's density function}$$

We are to estimate θ .

→ Compute the log likelihood function $\ell_{\theta}(x)$

$$\ell_{\theta}(x) = \log\left(\frac{1}{\pi}\right) + \log(1) - \log(1+(x+\theta)^2)$$

$$\hat{\ell}_{\theta}(x) = \frac{2(x-\theta)}{1+(x+\theta)^2}$$

$$\hat{\ell}_{\theta}(x) = \frac{-2[1+(x-\theta)^2] - [2(x-\theta)(-1)(+2)(x-\theta)]}{[1+(x-\theta)^2]^2}$$

$$\hat{\ell}_{\theta}(x) = \frac{-2[1+(x-\theta)^2] + 4(x-\theta)^2}{[1+(x-\theta)^2]^2}$$

$$I(x) = \cancel{\ell_x(\hat{\theta}^{\text{MLE}})}$$

→ 10000 samples of size $= 20$ drawn with $\theta = 0$
 We computed $\frac{1}{I(x)}$ for each sample.

→ Grouped the 10000 $\hat{\theta}^{\text{MLE}}$ values according to quantiles of $\frac{1}{I(x)}$
 and calculated the empirical variance for each group.

↳ Rough estimate of conditional variance of $\hat{\theta}^{\text{MLE}}$ given $\frac{1}{I(x)}$

for all samples, the unconditional variance $\frac{1}{n} I_0$ is the same.

↳ Specifically $I_0 = 1/2$ for a single cauchy observation.

Thus $\frac{1}{20 \cdot (1/2)} = \frac{1}{10} = 0.1$

↳ $\frac{1}{n} I_0 \Rightarrow$ all the samples have the same Fisher Information because they are of the same size.

↳ The standard Fisher Information is defined as an expected value.

On the other hand, the observed Fisher Information will vary from sample to sample. Because the $\hat{\theta}^{\text{MLE}}$ are different for each sample. (jupyter Lecture 11)

- The observed Fisher Information is related to the variance

20 Sep 2019

➤ Permutation & Randomization

→ permutation (shuffling data)

- recall our dataset where we had data for activity of 7128 genetic markers on Leukemia patients, with 2 different types of diagnosis, AML & ALL.

↳ 47 patients

↳ 25 patients.

- we first attempted obtaining t -values for the group means of each marker. then we got an histogram of the 7128 t -values.

$$t_i = \frac{\bar{X}_{\text{ALL}} - \bar{X}_{\text{AML}}}{\hat{\sigma}_i} \rightarrow \text{assumes that the AML \& ALL are distributed Normally.}$$

- an early conclusion we got was that the empirical distribution of t -values didn't match the t -student, with 70 degrees of freedom it was supposed to follow (as it had fatter tails). (jupyter Lecture 2 part 2)

→ The permutation method we will not need to assume normality.

- one major assumption to use the two-sample t-test that compared our t-scores vs. the t-student with 70 degrees of freedom was that our observations of gene activity come from a normal distribution
- To get rid of the normality assumption, Fisher suggested:
 - to take random permutations (of a single gene's activity from our 72 individuals).
 - pick the first 72 permuted individuals as "having ALL" and the next 25 individuals as "having AML"
 - The intention is to make any differences in the mean activity for each gene across the two groups be completely random
 - Take gene 136 & let's compute 10 000 random permutations of the 72 observations of gene 136.
(jupyter lecture 12 part 1)
 - the t-value distribution tells us if that t-value could or couldnt be generated by chance.

Hypothesis

- In the classical approach, we compare the t-values of all the 7128 genes against the same t-distribution with 70 degrees of freedom
- ↳ using permutation, we have a distribution for each gene.

→ Randomization

Randomly assigning experimental units (individuals) to possible treatment groups.

In the AML / ALL study, there might have existed other reasons (variables) that may influence the diagnosis.

for instance, from an homogeneous group, choose test samples and control samples randomly.

Computer Age Statistical Inference CHAPTER 05

► Univariate Parametric Models

Philosophical differences aside, the three classical approaches we covered relied on low-dimensional (the number of parameters we are to learn / fit) parametric models.

A brief table of useful univariate families

Name, Notation	Density	x Sample Space	Ω Parameter Space	Expectation
Normal $N(\mu, \sigma^2)$	$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}}$	$\mathbb{R}^{(0)}$	$\mu \in \mathbb{R}^{(0)}$ $\sigma^2 \in \mathbb{R}^+$	μ σ^2
Poisson $Poi(\lambda)$	$\frac{e^{-\lambda} \lambda^x}{x!}$	\mathbb{N}_0	$\lambda \in \mathbb{R}^+$	λ λ

Useful to model quantities that take positive and negative continuous values • if the distribution is ~~symmetric~~ or • if there aren't many extreme values

Poisson
 $Poi(\lambda)$

↳ not a density function but a probability mass function, as $Poi(\lambda)$ is discrete (not continuous)

(if the mean grows/shrinks, the variance grows/shrinks proportionally)

λ must stay positive,
 λ is mean interval of time of an exponential distribution, which is continuous → meaning that the expected number of successes can have decimals.

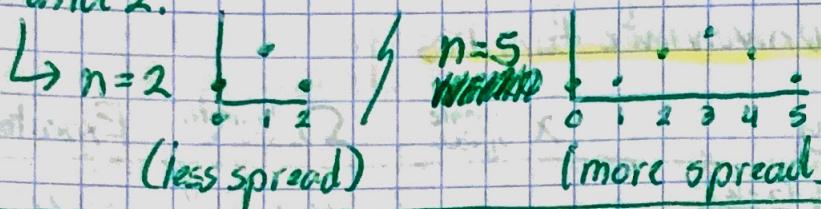
Use to model a quantity that is discrete, it's the number of counts of something. It's not very flexible, as only has one parameter to tweak.

Note, in the Normal we have 2 parameters to tweak, but they are very specific on what they describe. μ is the location parameter \leftrightarrow and σ^2 \leftrightarrow

Name, Notation	Density	x	Σ	Expectation, Variance
Binomial Bin(n, θ)	$\binom{n}{x} \theta^x (1-\theta)^{n-x}$	$\{0, 1, \dots, n\}$	$0 \leq \theta \leq 1$	$n\theta$ $n\theta(1-\theta)$

We also count the number of successes, however we know the number of trials (unlike the Poisson)

The n is constraining what the x can be, giving some flexibility.
→ if n is 2 the probability needs to be concentrated between 0 and 2.



Gamma	$\frac{x^{\nu-1} e^{-x/\sigma}}{\sigma^\nu \Gamma(\nu)}$	\mathbb{R}^+	$\nu > 0$	σ
Ga(ν, σ)			$\sigma > 0$	$\sigma^2 \nu$

$\nu-1$ shall stay positive & $-x/\sigma$ shall stay negative

The Gamma is used to model positive quantities. It's common to use the inverse Gamma to model variances.

Beta	$\frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)}$	$0 \leq x \leq 1$	$\alpha > 0$	α/β
Be(α, β)			$\beta > 0$	*

As the x goes from 0 to 1, it's mostly used to talk about probabilities (as a probability distribution)

$$* = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$$

- Both the Gamma & Beta have 2 parameters that convey some degree of flexibility.
- Gamma is flexible but not as flexible as the Beta
- The Binomial with a large n and small probability can approximate the Poisson distribution and viceversa