

> t-student

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

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

$$\hat{s}_d^2 = \hat{\sigma}^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right) \quad \dots \text{use ①}$$

→ 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.

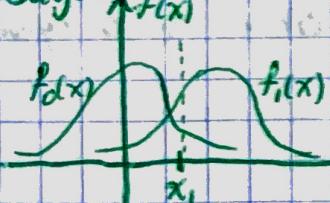
→ NP lemma states:

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 rerejecting 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:

$$t_c(x) = \begin{cases} 1 & \text{if } L(x) \geq c' \\ 0 & \text{if } L(x) < c' \end{cases}$$

As only rules of the $t_c(x)$ form can be optimal

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

\rightarrow 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 viceversa)

e.g.

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

$$f_0 \sim N(0, 1) \quad \nmid \quad f_1 \sim CN(\gamma, 1)$$

$$f_0(x_i) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x_i^2}{2}} \quad f_1(x_i) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{2}} \quad \text{for normal dists}$$

density functions.

$$f_0(\mathbf{x}) = \left[\frac{1}{\sqrt{2\pi}} \right]^n \prod_{i=1}^n e^{-\frac{x_i^2}{2}} = \left[\frac{1}{\sqrt{2\pi}} \right]^n e^{-\frac{1}{2} \sum_{i=1}^n x_i^2} \quad \left. \begin{array}{l} \text{the likelihood} \\ \text{functions} \end{array} \right\}$$

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

$$L(\mathbf{x}) = \frac{e^{-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2}}{e^{-\frac{1}{2} \sum_{i=1}^n x_i^2}} \quad \left. \begin{array}{l} \text{likelihood} \\ \text{ratio} \end{array} \right\}$$

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

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

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

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

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

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

> 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 $\sigma^2 = 1$

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

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

normalize

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

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

$$\alpha = 1 - P(\bar{x}\sqrt{n} \leq c\sqrt{n} | \mu = 0) \quad / \text{ } \Phi \text{ is the cumulative density function (CDF) of a normal distribution.}$$

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

\hookrightarrow 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:

$$c = \frac{1}{\sqrt{n}} \Phi^{-1}(1 - \alpha) \quad / \quad c = \mu_0 + \frac{1}{\sqrt{n}} \Phi^{-1}(1 - \alpha) \quad \text{for two normal distributions with the same variance}$$

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

\hookrightarrow Probability of correctly ~~rejecting~~ rejecting

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

On the other hand:

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

\hookrightarrow 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 a 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 ~~knowledge~~
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 on 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:

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

↓ 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 ~~all~~ 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?

$\text{MMD} = \text{different twins}$.

$$\frac{g(I, S)}{g(F, S)} = \frac{g(I)}{g(F)} \cdot \frac{f_I(S)}{f_F(S)}$$

$I = \text{identical twins}$.
 $S = \text{same sex}$.
 $F = \text{fraternal}$,

$$= \frac{1/3}{1/2} \cdot \frac{1}{1/2} \rightarrow \text{if they are identical, what's the prob. of being same sex}$$

Same sex

$$= 1$$

$\hookrightarrow g(I, S) \text{ & } g(F, S) \text{ are equally likely.}$

\Rightarrow Fraternal & Identical are equally likely

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

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

MMD If a random variable is distributed binomially

$$x \sim \text{Bin}(n, p) \Rightarrow P(X = k | p, n)$$

\hookrightarrow number of successes being k

$p = \text{probability of success}$

$n = \text{number of observations (trials)}$

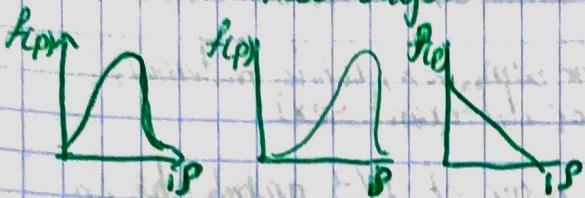
$$\binom{n}{k} = \frac{n!}{k!(n-k)!} = \frac{T(n+1)}{T(k+1) \cdot T(n-k+1)}$$

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

Probability that the number of successes is exactly k

↳ observing $n-k$ failures
↳ observing k successes

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



⇒ all of these can be modeled using the Beta distribution family of

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

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

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

↳ 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:

$$E(p|\alpha, \beta) = \frac{\alpha}{\alpha + \beta}$$

The mean

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

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

↳ under this parametrization the expected value E is:

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

↳ expected value of ρ

↳ variance of ρ

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

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

↳ posterior dist.

↳ likelihood function

↳ prior

our belief about ρ
after observing k successes

* → we are applying the Bayes rule, without the normalizing constant $/B(\alpha, \beta)$
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 *

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

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

▷ 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

$$x = [0, 0, 1, 0, 1, \dots, 0, 1] \quad \sum^n_i$$

$$k = \sum x_i$$

K = number of successes

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

M = is equivalent to having had a sample size n in the past.
(aka. the # 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:

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

\downarrow

marginal distribution of
the # of successes in
n future trials, given
parameters $\mu, M & n$

$$= \frac{T(M)}{T(M\mu) \cdot T(M(1-\mu))} \binom{n}{k} \int_0^1 p^{k+M\mu-1} (1-p)^{n-k+M(1-\mu)-1} dp$$

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

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

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

failures

= is the posterior predicted distribution

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

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