

Notes on basic probability and statistics

①

The sample space Ω is the set of all possible outcomes of an experiment

ex. 1: toss a coin twice $\Omega = \{HH, HT, TH, TT\}$

ex 2: measure a distance $\Omega = [0, \infty)$ \leftarrow "outcomes"

Subsets of Ω are called events

ex 1: first toss is heads: $A = \{HH, HT\}$

ex 2: distance is less than 10: $A = [0, 10)$

A function P that assigns a real number $P(A)$ to each event A is a probability measure if it satisfies:

1. $P(A) \geq 0$ for every A \leftarrow aka probability distribution

2. $P(\Omega) = 1$

3. If A_1, A_2, \dots are disjoint, then

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$$

Two events A and B are independent if and only if

$$P(A \text{ and } B) = P(A) \cdot P(B)$$

ex. $P(\text{H on first toss and H on second toss})$

$$= P(\text{H on first toss}) \times P(\text{H on second toss})$$

If $P(B) > 0$ then the conditional probability of A given B is:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \leftarrow \text{"A and B"}$$

NOTE:
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

This is a version of Bayes Theorem

A random variable X is a mapping $X: \Omega \rightarrow \mathbb{R}$ that assigns a real number $X(\omega)$ to each outcome ω .

ex 1. $X =$ number of heads

<u>ω</u>	<u>$P(\{\omega\})$</u>	<u>$X(\omega)$</u>	<u>x</u>	<u>$P(X=x)$</u>
HH	$\frac{1}{4}$	2	0	$\frac{1}{4}$
HT	$\frac{1}{4}$	1	1	$\frac{1}{2}$
TH	$\frac{1}{4}$	1	2	$\frac{1}{4}$
TT	$\frac{1}{4}$	0		

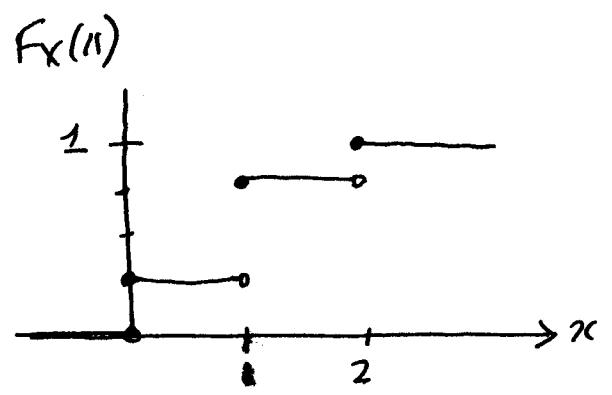
The cumulative distribution function (CDF) is the function

$F_X: \mathbb{R} \rightarrow [0, 1]$ defined by:

$$F_X(x) = P(X \leq x)$$

ex 1. $X =$ number of heads

$$F_X(x) = \begin{cases} 0 & x < 0 \\ 1/4 & 0 \leq x < 1 \\ 3/4 & 1 \leq x < 2 \\ 1 & x \geq 2 \end{cases}$$



X is discrete if it takes countably many values.

The probability function for X is $f_X(x) = P(X=x)$

X is continuous if there exists a function f_X such that $f_X(x) \geq 0$ for all x , $\int_{-\infty}^{\infty} f_X(x) dx = 1$ and for

every $a \leq b$ $P(a < X < b) = \int_a^b f_X(x) dx$

In this case, f_X is called a probability density function. (PDF)

An important discrete probability function: Poisson

X takes values in $\{0, 1, 2, \dots\}$

$$f_X(x) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad \lambda > 0 \text{ "parameter"}$$

An important continuous probability ^{density} function: Normal

X takes values in \mathbb{R} . $N(\mu, \sigma^2) \rightarrow f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}, \quad \sigma > 0, \mu \in \mathbb{R}$

The expected value $E(X)$ of a random variable X is:

$$E(X) = \int x dF_x(x) = \begin{cases} \sum x f_x(x) & \text{if } X \text{ is discrete} \\ \int x f_x(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

for Poisson, $E(X) = \lambda$

for Normal, $E(X) = \mu$

The variance $V(X)$ of a random variable is:

$$V(X) = E(X - E(X))^2 = E(X^2) - (E(X))^2$$

for Poisson, $V(X) = \lambda$

for Normal, $V(X) = \sigma^2$

Central Limit Theorem

Suppose X_1, X_2, \dots, X_n are independent and identically distributed (iid), then $\bar{X}_n = \frac{1}{n} \sum X_i$ has a distribution which is approximately normal with mean μ and variance σ^2/n .

Note if you have f , can figure out F or vice versa

Statistical Inference

Given a sample $X_1, \dots, X_n \sim F$, how do "infer" F ?
(or failing that, some features of F , such as the expectation)
iid draws

A parametric statistical model is a set of distributions that can be parameterized by a finite number of parameters:

EX: a normal model

$$\left\{ f_X(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right), \mu \in \mathbb{R}, \sigma^2 > 0 \right\}$$

more generally:

$$\left\{ f_X(x; \theta) : \theta \in \Theta \right\}$$

Statistical inference reduces to inferring parameters

Basic concepts:

An estimator $\hat{\theta}_n$ of a parameter θ is consistent if $\hat{\theta}_n$ converges to θ as $n \rightarrow \infty$

The quality of a point estimate can be assessed by mean square error: $MSE = E(\hat{\theta}_n - \theta)^2$

$$MSE = \underbrace{(\bar{\theta}_n - \theta)^2}_{\text{"bias"}} + E(\hat{\theta}_n - \bar{\theta}_n)^2_{\text{"variance"}}$$

EX. IN a normal model, can estimate μ by:

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n x_i$$

can show that $\hat{\mu}_n$ is consistent.

Since $\hat{\mu}_n$ is a function of x_1, \dots, x_n it is itself a random variable and has a PDF. Can show that

$$\hat{\mu}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

This is a "sampling distribution"

The standard deviation of $\hat{\mu}_n$ (more generally $\hat{\theta}_n$) is called the standard error (SE)

$$SE = SE(\hat{\theta}_n) = \sqrt{V(\hat{\theta}_n)}$$

A $1-\alpha$ confidence interval for a parameter θ is an interval $C_n = (a, b)$ where a and b depend on the data such that

$$P(\theta \in C_n) \geq 1-\alpha \text{ for all } \theta$$

Note: C_n is random.

EX. IN a normal model, a 95% confidence interval for μ is given by

$$\left(\hat{\mu}_n - 1.96 \frac{\hat{\sigma}_n}{\sqrt{n}}, \hat{\mu}_n + 1.96 \frac{\hat{\sigma}_n}{\sqrt{n}} \right)$$

an estimator of σ , e.g. the sample standard deviation

In the normal model we know that

$$\hat{\mu}_n \sim N(\mu, \sigma^2/n).$$

Can show that, approximately,

$$\hat{\mu}_n \sim N(\mu, \hat{\sigma}_n^2/n),$$

Where $\hat{\sigma}_n$ is the sample standard deviation.

Suppose your sample comprises:

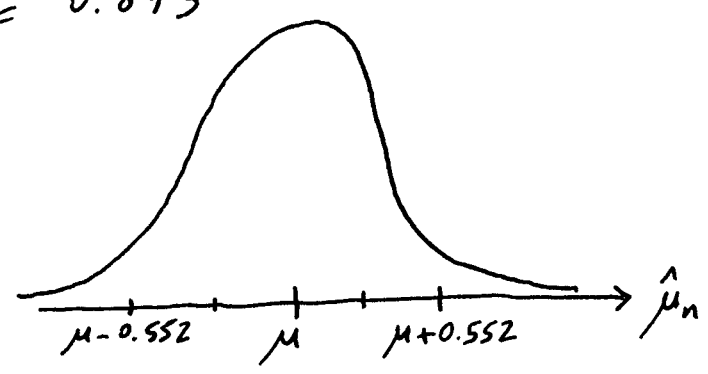
- 0.164, -0.051, 0.882, 0.283, -1.424
- 1.120, -0.342, -1.490, -0.394, -1.700

Then $\hat{\mu}_n = \frac{\sum x_i}{10} = -0.519$

$$\hat{\sigma}_n = \sqrt{\frac{1}{9} \sum (x_i - \hat{\mu}_n)^2} = 0.873$$

$$SE = \frac{\hat{\sigma}_n}{\sqrt{n}} = 0.276$$

So $\hat{\mu}_n \sim N(\mu, 0.276^2)$

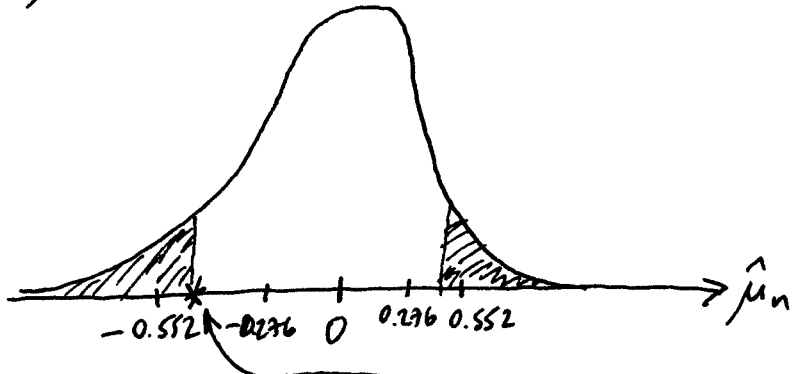


the distribution of the $\hat{\mu}_n$'s that you might have gotten. Actually have one $\hat{\mu}_n$ and it is -0.519 .

In hypothesis testing you play the following game:

Suppose μ really is 0 (or 1 or 21 or whatever)
(this is the "null hypothesis")

Then $\hat{\mu}_n \sim N(0, 0.276^2)$

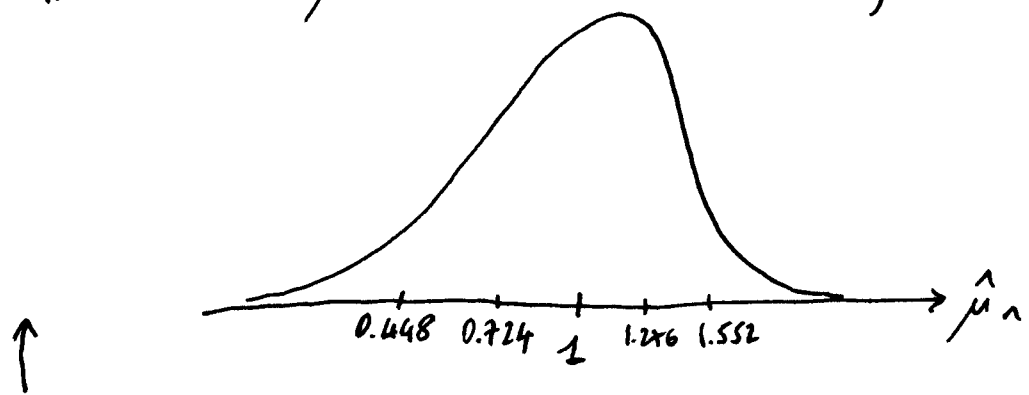


The actual observed $\hat{\mu}_n$ is here... not so close to zero!

How likely is it that you^{ie $\hat{\mu}_n$} would be that far from the true value of μ (i.e. zero)? Answer: 0.06

This quantity is called a p-value, two-sided in this case.

How about $\mu=1$? does that seem plausible?



The actual observed $\hat{\mu}_n$ is over here somewhere! p-value is essentially zero. You can essentially rule out 1 as a plausible value of μ .

Hypothesis testing is a very stylized and often quite unhelpful way to think about statistical problems.

back to estimation for parametric models...

For our normal example it seems very natural to estimate μ with $\frac{\sum x_i}{n}$. (and in fact it is a fine estimator) but could use, for example, a trimmed mean or the median.

More generally, it might not be obvious how to construct an estimator. Maximum Likelihood Estimation is a general purpose scheme for producing (usually good) estimators

let x_1, \dots, x_n be iid with PDF $f_x(x; \theta)$.

The likelihood function $L_n(\theta) = \prod_{i=1}^n f(x_i; \theta)$

Often work with the log likelihood function

$$l_n(\theta) = \log L(\theta) = \sum_{i=1}^n \log f(x_i; \theta)$$

The maximum likelihood estimator (MLE) denoted

by $\hat{\theta}_n$ is the value of θ that maximizes $L_n(\theta)$

(or equivalently maximizes $l_n(\theta)$)

Ex. $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ (think coin tossing with a possibly biased coin) (10)

the probability function is: $f_X(x; p) = p^x(1-p)^{1-x}$ for $x=0, 1$.

$$\text{Then } \mathcal{L}_n(p) = \prod_{i=1}^n f(X_i; p) = \prod_{i=1}^n p^{X_i} (1-p)^{1-X_i} = p^S (1-p)^{n-S}$$

$$\text{where } S = \sum_{i=1}^n X_i$$

$$\therefore \ell_n(p) = S \log p + (n-S) \log(1-p)$$

take derivative and set equal to zero to find MLE:

$$\hat{p}_n = \frac{S}{n}$$

Ex. $X_1, \dots, X_n \sim N(\mu, \sigma^2)$

$$\mathcal{L}_n(\mu, \sigma) = \prod_{i=1}^n \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right)$$

$$= \sigma^{-n} \exp\left(-\frac{nS^2}{2\sigma^2}\right) \exp\left(-\frac{n(\bar{X} - \mu)^2}{2\sigma^2}\right)$$

$$\text{where } \bar{X} = \frac{1}{n} \sum X_i \text{ and } S^2 = \frac{1}{n} \sum (X_i - \bar{X})^2$$

Take log, take derivatives with respect to μ and σ^2 , and set equal to zero:

$$\hat{\mu} = \bar{X} \quad \hat{\sigma} = S$$

Properties of MLEs

(11)

1. The MLE is consistent $\hat{\theta}_n \rightarrow \text{true } \theta_{\wedge}^{(\theta^*)}$ as $n \rightarrow \infty$

2. The MLE is asymptotically normal:

$$\frac{\hat{\theta}_n - \theta^*}{\hat{SE}} \rightarrow N(0, 1)$$

(so we know the sampling distribution and can thus get confidence intervals and p-values)

3. The MLE is the optimal estimator in a particular technical sense.

Important stuff that I am omitting for now...

- Likelihood ratio tests.
- AIC/BIC etc.
- The bootstrap
- Monte Carlo
- Bayesian Methods