

Probability Theory

Inclusion-Exclusion

$$|A \cap B| = |A| + |B| - |A \cup B|$$

Permutations: $k!$, order is imp.

Combinations: $\binom{n}{k} = \frac{n!}{k! \cdot (n-k)!}$

$$P(A \cup B) = P(A) + P(B) \text{ for } A, B \text{ disjoint}$$

$$P(\bar{A}) = 1 - P(A) \quad P(\emptyset) = 0$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Conditional Probability

probability mass function = a function giving the prob. of each event

$$\sum p_i = 1 \leftarrow P(A)$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, P(B) > 0$$

dependent:

$$P(A \cap B) = P(B|A) \cdot P(A), P(A) > 0$$

$$P(A \cap B^c) = P(A) \cdot P(B^c)$$

$$P(A^c \cap B^c) = P(A^c) \cdot P(B^c)$$

independent:

$$P(A \cap B) = P(A) \cdot P(B)$$

$$P(A \cap C) = P(A) \cdot P(C)$$

$$P(B \cap C) = P(B) \cdot P(C)$$

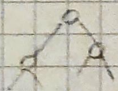
$$P(A \cap B \cap C) = P(A) \cdot P(B) \cdot P(C)$$

Law of Total Prob.

$$P(A) = \sum_{i=1}^n P(A \cap C_i) = \sum_{i=1}^n P(A|C_i) \cdot P(C_i)$$

Bayes Theorem

$$P(R_1|R_2) = \frac{P(R_1 \cap R_2)}{P(R_2)} = \frac{P(R_2|R_1) \cdot P(R_1)}{P(R_2)}$$



$$E(ax+by) = E(x)a + E(y)b$$

$$\text{Var}(ax+by) = E((x-E(x))^2) = \sigma^2$$

Random Variables & Distr.

continuous Random Variable:

$$\text{pdf: } P(c=x=d) = \int_c^d f(x) dx$$

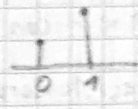
$$\text{cdf: } F(x) = P(X \leq x)$$

$$X \sim \text{ber}(p)$$

$$E(x) = p$$

$$\text{Var}(x) = p \cdot (1-p)$$

	x	0	1
pmf: P(x)	1-p	p	
cdf: F(x)	1-p	1	

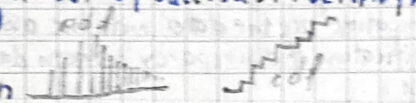


$X \sim B(n, p)$ Number of successes in n trials

$$B(n, p) = \text{ber}(p)$$

$$E(x) = np$$

$$\text{Var} = p \cdot (1-p) \cdot n$$



$$P(x) = \binom{n}{x} \cdot p^x \cdot (1-p)^{n-x}$$

pmf: binom(x, n, p)

$$\text{cdf: } \text{pbinom}(x, n, p, \text{lower.tail} = \text{false}) = 0.5^x$$

$X \sim \text{Geom. Distr.}$ = total number of attempts before success

$$E(x) = \frac{1}{p}$$

$$\text{Var}(x) = \frac{1-p}{p^2}$$

$$\text{pmf: } p(x) = (1-p)^{x-1} \cdot p = \text{dgeom}(n, p)$$

$$\text{cdf: } \text{pgeom}(n, p)$$

$X \sim P(\lambda)$ λ = intensity param.

$$P(x) = \frac{\lambda^x}{x!} \cdot e^{-\lambda}, x \in \{0, 1, 2, \dots\} = \text{dpois}(n, \lambda)$$

$$E(x) = \text{Var}(x) = \lambda$$

$$X \sim U(0, b)$$

$$\text{pdf: } f(x) = \begin{cases} \frac{1}{b-a}, & x \in (0, b) \\ 0, & \text{else} \end{cases}$$

$$E(x) = \frac{a+b}{2}$$

$$\text{Var}(x) = \frac{(b-a)^2}{12}$$

$$\text{cdf: } F(x) = \begin{cases} \frac{x-a}{b-a}, & a \leq x < b \\ 0, & x < a \\ 1, & x \geq b \end{cases}$$

$X \sim \text{exp}(\lambda)$ = models waiting time (continuous analog of geometric distrib.)

$$\text{pdf: } \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

$$E(x) = \frac{1}{\lambda}$$

$$\text{Var}(x) = \frac{1}{\lambda^2}$$

$X \sim N(\mu, \sigma^2)$ Normal distribution

$$\text{pdf: } f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}, x \in \mathbb{R}$$

$$E(x) = \mu$$

$$\text{Var}(x) = \sigma^2$$

Standard Normal distr. $Z \sim N(0, 1)$

$$\text{cdf: } \Phi(t) = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} \cdot e^{-\frac{t^2}{2}} dt$$

$$Z = \frac{X-\mu}{\sigma} \Rightarrow P(X > 200) = P\left(\frac{X-\mu}{\sigma} > \frac{200-\mu}{\sigma}\right)$$

Quantiles:

16: ~69%, 26: ~95%, 36: ~99.7%

0.k Quantile: $Z_k = \sigma^{-1}(k) + \mu, 0 < k < 1$

$$P(Z > 2) = 1 - \Phi(2.01) = 0.9778\%$$

Descriptive Statistics

mean = center of mass of data $\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i$

empirical variance of the data $s^2 := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

empirical standard deviation $s = \sqrt{s^2}$

Hypothesis Testing

How "compatible" are the data with the assertion μ_0 ?
 > quantification of discrepancy between data & assertion $\bar{x} - \mu_0$
 if large \leftrightarrow data hardly compatible with assertion

Tests

test statistics: $Z = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \approx Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} \sim N(0,1)$

Just measure the discrepancy between data & assertion.
 p-value: quantifies the discrepancy
 $p = P_{H_0}(|Z| \geq |z|) \approx 9 \cdot 10^{-10}$ example
 if p is small, reject hypothesis
 \hookrightarrow the observed discrepancy was significant

significance level: α

Rejection Area R

2-sided: $H_0: \mu = \mu_0$ ($H_1: \mu \neq \mu_0$)

$$R = (-\infty, q_{\alpha/2}] \cup [q_{1-\alpha/2}, \infty)$$

$$P = P_{H_0}(|Z| \geq |z|)$$

left-sided: $H_0: \mu \geq \mu_0$ ($H_1: \mu < \mu_0$)

$$R = (-\infty, q_{\alpha/2}]$$

$$P = P_{H_0}(Z \leq z)$$

right-sided: $H_0: \mu \leq \mu_0$ ($H_1: \mu > \mu_0$)

$$R = [q_{1-\alpha/2}, \infty)$$

$$P = P_{H_0}(Z \geq z)$$

t-Test (one-sample)

$$T := \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \sim t(n-1)$$

$n-1$ degrees of freedom
 $SEM_n := \frac{s}{\sqrt{n}}$
 standard error of mean

measures the discrepancy $\bar{x} - \mu_0$ in the units SEM

$$T = \frac{\bar{x} - \mu_0}{SEM}$$

\Rightarrow if $H_0: \mu = \mu_0$ holds true
 $|T| = 1 \leftrightarrow |\bar{x} - \mu_0| = 1 \cdot SEM$
 unlikely \leftarrow not very
 $|T| = 5 \leftrightarrow |\bar{x} - \mu_0| = 5 \cdot SEM$

Confidence Interval

$$I := (\bar{x} - q_{1-\alpha/2} \cdot SEM, \bar{x} + q_{1-\alpha/2} \cdot SEM)$$

overlaps the parameter μ_0 with prob. $1-\alpha$

t-Test (two sample)

\Rightarrow discrepancy large if SEM small

$$T = \frac{\bar{y} - \bar{x}}{SEM_{\bar{y} - \bar{x}}}$$

$T := \dots \sim t(r)$ deg. of freedom

$H_0: \mu_x = \mu_y$

$$I := (\bar{y} - \bar{x}) - q_{1-\alpha/2} \sqrt{SEM_y^2 + SEM_x^2}$$

$$(\bar{y} - \bar{x}) + q_{1-\alpha/2} \sqrt{SEM_y^2 + SEM_x^2}$$

Overlaps the parameter μ_0 with prob $1-\alpha$

Proportions

absolute frequency
 relative $= \frac{h}{n} = \text{proportion of } n = \dots \approx 0$

$$h_i = \frac{1}{n} \sum_{i=1}^n Y_i$$

$$E(H) = p \quad \text{Var}(H) = \frac{p(1-p)}{n}$$

$$n \rightarrow \infty \quad \frac{H-p}{\sqrt{p(1-p)/n}} \rightarrow N(0,1)$$

$$\sigma_H = \sqrt{p \cdot (1-p)/n}$$

$$SEM_H := \frac{\sqrt{H(1-H)}}{n} \quad \text{when } p \text{ is unknown}$$

discrepancy = $|h - p| \approx \dots \cdot SEM$

$$Z := \frac{H - p_0}{SEM_H} \approx N(0,1)$$

$$I := (H - q_{1-\alpha/2} \cdot SEM_H, H + q_{1-\alpha/2} \cdot SEM_H)$$

overlaps the parameter p_0 with $1-\alpha$

$$\text{two sample: } Z := \frac{(H_2 - H_1) - 0}{\sqrt{SEM^2 + SEM^2}}$$

$I := (H_2 - H_1) - \dots$
 with prob. $1-\alpha$, 0 is overlapped

χ^2 -test (goodness of fit)

$$\chi = (x_1, \dots, x_d) \sim \text{mult}_d(n, p)$$

with $p \in (0,1)^d$ & $\sum_{k=1}^d p_k = 1$

$$H_0: p = (p_1, \dots, p_d)$$

χ^2 -stat: 'observed' \rightarrow 'expected'

$$\chi^2 := \sum_{k=1}^d \frac{(x_k - E_{H_0}[x_k])^2}{E_{H_0}[x_k]}$$

k	1	2	3	4	5	6	Σ
x_k	21	22	46	17	19	25	170
$E_{H_0}[x_k]$	20	20	20	20	20	20	120

\leftarrow Bsp.