

Probability And Statistics

Mathematical Methods

Log and Exponential

$\log = \log_e = \ln$; $\log(x^*) = \log(x) + \log(y)$; $\log(x^*) = y \log(x)$; $\log(e^x) = x$; $\lim_{x \rightarrow \infty} \log(x) = -\infty$
 $e^x = \exp(x)$; $\exp(x+y) = \exp(x)\exp(y)$; $\exp(x)^y = \exp(y)$; $\exp(\log(x)) = x$; $\exp(0) = 1$
Arithmetic Progressions n^{th} term = $a + (n-1)d$; $S_n = n/2(2a + (n-1)d)$; $S_{2n} = +\infty$ (unless $a = d = 0$).
Geometric Progressions n^{th} term = ar^{n-1} ; $S_n = (a(1-r^n))/(1-r)$ where $r \neq 1$; $S_{\infty} = a/(1-r)$ where $|r| < 1$.
Calculus (Differentiation opposite of Integration)

$$\frac{df}{dx} \equiv f'(x) \equiv \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

Chain Rule: $\frac{d}{dx} f(g(x)) = f'(g(x)) g'(x)$

Product Rule: $\frac{d}{dx} f(x)g(x) = f'(x)g(x) + f(x)g'(x)$

Quotient Rule: for $g(x) \neq 0$, $\frac{d}{dx} \frac{f(x)}{g(x)} = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2}$

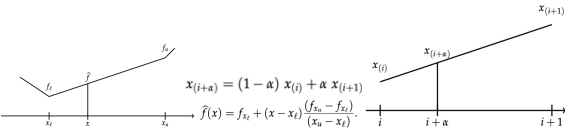
By parts: $\int f(x)g'(x)dx = [f(x)g(x)] - \int f'(x)g(x)dx$

Change of variable: if $y = g(x)$, $\int_a^b f(x)dx = \int_{g(a)}^{g(b)} f(g^{-1}(y))g^{-1}'(y)dy$

Images and inverses f: X → Y, A is subset of X (Image, Inverse, Inverse Image (B subset of Y)).

$f(A) = \{y \in Y | f(x) = y \text{ for some } x \in A\}$, $f^{-1}(f(A)) = X \cap f^{-1}(B) = \{x \in X | f(x) \in B\}$.

Interpolation



Numerical Summaries

Measures of Location Mean, Order statistic $x_{(i)}$ i^{th} smallest value, median, mode most frequent – multimodal.

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_{(i)} \quad x_{(n)} = \max(x_1, \dots, x_n) \quad \text{median} = x_{((n+1)/2)} = \begin{cases} x_{((n+1)/2)} & \text{if } n \text{ is odd} \\ \frac{x_{(n/2)} + x_{(n/2+1)}}{2} & \text{if } n \text{ is even} \end{cases}$$
$$x_{(1)} = \min(x_1, \dots, x_n)$$

Others Geometric mean (less affected by large values), Harmonic mean (useful when averaging rates).

$$x_G = \sqrt[n]{x_1 \cdots x_n}, x_H = \exp\left\{\frac{1}{n} \sum_{i=1}^n \log x_i\right\}, x_H = \left(\frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}\right)^{-1}$$

Measures of Dispersion Range, 1st and 3rd quartile,

first quartile = $x_{((n+1)/4)}$
range = $x_{(n)} - x_{(1)}$ third quartile = $x_{(3(n+1)/4)}$ interquartile range = third quartile – first quartile
Five-point summary of set of data lists, in order min, lower quartile, sample median, upper quartile, max.
Sample variance (mean square)/Sample standard deviation (root mean square), Skewness (asymmetry)

$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad s^2 = \left(\frac{1}{n} \sum_{i=1}^n x_i^2\right) - \left(\frac{1}{n} \sum_{i=1}^n x_i\right)^2 \quad s = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{skewness} = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s}\right)^3$$

	Least Robust	More Robust	Most Robust
Location	$\frac{x_{(1)} + x_{(n)}}{2}$	\bar{x}	$x_{((n+1)/2)}$
Dispersion	$x_{(n)} - x_{(1)}$	s^2	$x_{(3(n+1)/4)} - x_{((n+1)/4)}$

Covariance and correlation

$$s_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}), \quad s_{xy} = \frac{\sum_{i=1}^n x_i y_i}{n} - \bar{x}\bar{y}, \quad r_{xy} = \frac{s_{xy}}{s_x s_y} = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{n s_x s_y}$$

Box-and-Whisker Plots Median, 3rd, and 1st quartiles lines; whiskers extend to points within 3/2 IQR, extreme points plotted individually with x or o.

Empirical CDF Returns proportion of data having values which do not exceed x . $F_n(x) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq x)$

Elementary Set Theory

Sets and notation (below), complement (of A is set of elems not in A), subset (some elems), singleton (1).

	COMMUTATIVITY	ASSOCIATIVITY	DISTRIBUTIVITY	DE MORGAN'S LAWS
\in	"is an element of" (set membership)	$A \cup (B \cap C) = (A \cup B) \cap C$ $A \cap (B \cup C) = (A \cap B) \cup C$	$A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$ $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$	$\overline{A \cap B} = \overline{A} \cup \overline{B}$ $\overline{A \cup B} = \overline{A} \cap \overline{B}$
\Leftrightarrow	"if and only if" (equivalence)			
\Rightarrow	"implies"			
\exists	"there exists"			
\forall	"for all"			
s.t. or	"such that"			
wrt	"with respect to"			

Disjoint $(A \cap B) = \emptyset$ < Null set; **Partition of B** Disjoint sets whose union form B; **Difference** $A \setminus B = A \cap \overline{B}$.

Cartesian Product Set of all ordered pairs of elements of two sets. **Cardinality** No. of elems in set.

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

Ω universal set
 \emptyset empty set
 $A \subseteq \Omega$ Subset of Ω
 \overline{A} Complement of A
 $|A|$ Cardinality (or size) of A
 $A \cup B$ union (A or B)
 $A \cap B$ intersection (A and B)
 $A \subset B$ set inclusion (elements of A are also in B)
 $A \setminus B$ set difference (elements in A that are not in B)

Probability

Sample spaces (S) Set of all possible outcomes of experiment.
Event Subset of sample space (extreme events \emptyset (the null event) or S). E occurs if outcome $s \in E$.
Elementary events Singleton subsets of S which contain exactly one element from S.
Mutually Exclusive events if they are disjoint, i.e., at most one event can occur.
Sigma Algebra \mathcal{F} is subcollection of sets of all subsets of S with the following properties:
Definition 5.2.1. A collection \mathcal{F} of subsets of S is called a σ -field or σ -algebra if it satisfies the following conditions:
a) if $E, F \in \mathcal{F}$ then $E \cup F \in \mathcal{F}$ and $E \cap F \in \mathcal{F}$;
b) if $E \in \mathcal{F}$ then $\overline{E} \in \mathcal{F}$;
c) $\emptyset \in \mathcal{F}$.
Probability Measure
Definition 5.3.1. A probability measure P on (S, \mathcal{F}) is a mapping $P: \mathcal{F} \rightarrow [0, 1]$ satisfying
a) $P(S) = 1$;
b) if E_1, E_2, \dots is a collection of disjoint members of \mathcal{F} , so that $E_i \cap E_j = \emptyset$ for all pairs i, j with $i \neq j$, then
$$P\left(\bigcup_{i=1}^{\infty} E_i\right) = \sum_{i=1}^{\infty} P(E_i)$$

Interpretations of Probability
Classical Equally likely (uniform) probabilities.
Frequentist Repeated observations and find limiting value.
Subjective Degree of belief held by individual.
Independent Events
E, F independent iff $P(E \cap F) = P(E)P(F)$.

Theorem of Total Probability Let E_1, \dots, E_k be partition on S, then for any event $F \subseteq S$, we have $P(F) = \sum_{i=1}^k P(F|E_i)P(E_i)$. **Bayes' Theorem** Let E_1, \dots, E_k be partition on S $P(E_i|F) = (P(F|E_i)P(E_i))/P(F)$.

Conditional Probability $P(E|F)$; **Joint Probability** $P(E \cap F)$; **Marginal Probability** $P(E)$.

Summary of Conditional Probability
1. If $P(F) > 0$ then
$$P(E|F) = \frac{P(E \cap F)}{P(F)}$$

2. $P(\cdot|F)$ satisfies the axioms of probability, for fixed F. However, in general, $P(E|\cdot)$ does not satisfy the axioms of probability, for fixed E.

3. In general, $P(E|F) \neq P(F|E)$.

4. E and F are independent if and only if $P(E|F) = P(E)$.

Discrete Random Variables

Random Variables Measurable mapping $X: S \rightarrow \mathbb{R}$ with property that $\{s \in S : X(s) \leq x\} \in \mathcal{F}$ for each $x \in \mathbb{R}$.
Probability Distribution Function $P_X(x) \leq x\} = P(S_x)$; **Range** of random variable is image of S under X: $X(S) = \{x \in \mathbb{R} | \exists s \in S \text{ s.t. } X(s) = x\}$.
Cumulative Distributive Function $F_X(x) = P_X(X \leq x)$; Conditions of valid CDF (right side below). $P_X(a < X \leq b) = F_X(b) - F_X(a)$.

Definition 6.1.3. The cumulative distribution function (CDF) of a random variable X is the function $F_X: \mathbb{R} \rightarrow [0, 1]$, defined by
i) $0 \leq F_X(x) \leq 1, \forall x \in \mathbb{R}$;
ii) Monotonicity: $\forall x_1, x_2 \in \mathbb{R}, x_1 < x_2 \Rightarrow F_X(x_1) \leq F_X(x_2)$;
iii) $F_X(-\infty) = 0, F_X(\infty) = 1$.

Discrete Random Variable X is discrete if range of X, is countable (finite or infinite).

Probability Mass Function (Probability Function) – For discrete random variable X, $p_X(x) = P(X=x), x \in X$. If $x \notin X, p_X(x) = 0$. Properties of pmf (left), connection between F_X and p_X (right)

$$i) p_X(x_i) \geq 0;$$
$$ii) \sum_{x \in X} p_X(x) = 1. \quad F_X(x) = \sum_{x_i \leq x} p_X(x_i) \iff p_X(x_i) = F_X(x_i) - F_X(x_{i-1}), \quad i = 2, 3, \dots$$

Mean and Variance
Expectation or mean of a discrete random variable is defined to be $E_X(X)$ = sum of all $x(p_X(x))$.
$$E(g(X)) = \sum_x g(x)p_X(x) \quad E(ag(X) + bh(X)) = aE(g(X)) + bE(h(X))$$

(i) For a linear function, $g(X) = ax + b$ for constants, we have (from Theorem 6.4) that
$$E(g(X)) = \sum_x (ax + b)p_X(x) = a \sum_x xp_X(x) + b \sum_x p_X(x)$$

ii) F_X is continuous from the right on \mathbb{R} , that is, for $x \in \mathbb{R}$,
$$\lim_{h \rightarrow 0^+} F_X(x+h) = F_X(x)$$

iii) F_X is non-decreasing, that is,
$$a < b \implies F_X(a) \leq F_X(b).$$

iv) For $a < b$
$$P(a < X \leq b) = F_X(b) - F_X(a).$$

Definition 6.4.1. Let X be a random variable. The variance of X, denoted by σ^2 or σ_X^2 or $\text{Var}_X(X)$ is defined by
$$\text{Var}_X(X) = E\{[X - E_X(X)]^2\}.$$

Definition 6.4.3. The skewness (γ_1) of a discrete random variable X is given by
$$\gamma_1 = \frac{E_X\{[X - E_X(X)]^3\}}{s_X^3(X)}$$

6.4.1 Sums of Random Variables
Let X_1, X_2, \dots, X_n be n random variables, perhaps with different distributions and not necessarily independent.
Let $S_n = \sum_{i=1}^n X_i$ be the sum of those variables, and $\frac{S_n}{n}$ be their average.

Then the mean of S_n is given by
$$E(S_n) = \sum_{i=1}^n E(X_i), \quad E\left(\frac{S_n}{n}\right) = \frac{E(S_n)}{n} = \frac{\sum_{i=1}^n E(X_i)}{n}$$

However, for the variance of S_n , only if X_1, X_2, \dots, X_n are independent, we have
$$\text{Var}(S_n) = \sum_{i=1}^n \text{Var}(X_i), \quad \text{Var}\left(\frac{S_n}{n}\right) = \frac{\sum_{i=1}^n \text{Var}(X_i)}{n^2}$$

So if X_1, X_2, \dots, X_n are independent and identically distributed with $E(X_i) = \mu_X$ and $\text{Var}(X_i) = \sigma_X^2$ we get
$$E\left(\frac{S_n}{n}\right) = \mu_X, \quad \text{Var}\left(\frac{S_n}{n}\right) = \frac{\sigma_X^2}{n}$$

Geometric Distribution Potentially infinite sequence of independent Bernoulli(p) random variables
$$X = \min\{i | I_i = 1\}$$

to be the index of the first Bernoulli trial to result in a 1.
Example Tossing a coin, X is the number of tosses until the first head is obtained, $p = \frac{1}{2}$.
Then X is a random variable taking values in $\mathbb{Z}^+ = \{1, 2, \dots\}$, and we say $X \sim \text{Geometric}(p)$.
Clearly the pmf is given by
$$p_X(x) = p(1-p)^{x-1}, \quad x \in X = \{1, 2, \dots\}, \quad 0 \leq p \leq 1.$$

Notes
• The mean and variance are
$$\mu = E(X) = \frac{1}{p}, \quad \sigma^2 = \text{Var}(X) = \frac{1-p}{p^2}.$$

• The skewness is given by
$$\gamma_1 = \frac{2-p}{\sqrt{1-p}}.$$

and so is always positive.
Alternative Formulation
If $X \sim \text{Geometric}(p)$, let us consider $Y = X - 1$.
Then Y is a random variable taking values in $N = \{0, 1, 2, \dots\}$, and corresponds to the number of independent Bernoulli(p) trials before we obtain our first 1. (Some texts refer to this as the Geometric distribution).
Note we have pmf
$$p_Y(y) = p(1-p)^y, \quad y = 0, 1, 2, \dots,$$

and the mean becomes
$$\mu_Y = E_Y(Y) = \frac{1-p}{p}.$$

while the variance and skewness are unaffected by the shift.
6.5.5 Discrete Uniform Distribution
Let X be a random variable on $\{1, 2, \dots, n\}$ with pmf
$$p_X(x) = \frac{1}{n}, \quad x \in X = \{1, 2, \dots, n\}.$$

Then X is said to follow a discrete uniform distribution and we write $X \sim U(\{1, 2, \dots, n\})$.
Note The mean and variance are
$$\mu = E(X) = \frac{n+1}{2}, \quad \sigma^2 = \text{Var}(X) = \frac{n^2-1}{12}.$$

and the skewness is clearly zero.

Continuous Random Variables

Definition 7.0.1. A random variable X is (absolutely) continuous if $\exists f_X: \mathbb{R} \rightarrow \mathbb{R}$ (measurable) such that
$$P_X(B) = \int_B f_X(x) dx, \quad B \subseteq \mathbb{R},$$

in which case f_X is referred to as the probability density function, or pdf, of X.
7.0.2 Properties of Continuous F_X and f_X
By analogy with the discrete case, let X be the range of X, so that $X = \{x : f_X(x) > 0\}$.

i) For the cdf of a continuous random variable,
$$\lim_{x \rightarrow -\infty} F_X(x) = 0, \quad \lim_{x \rightarrow \infty} F_X(x) = 1.$$

ii) At values of x where F_X is differentiable
$$f_X(x) = \frac{d}{dx} F_X(x) \Big|_{x=x} = F_X'(x).$$

iii) If X is continuous,
$$f_X(x) \neq P(X=x) = \lim_{h \rightarrow 0^+} [P(X \leq x) - P(X \leq x-h)] = \lim_{h \rightarrow 0^+} [F_X(x) - F_X(x-h)] = 0$$

Warning! People usually forget, that $P(X=x) = 0$ for all x, when X is a continuous random variable.
iv) The pdf $f_X(x)$ is not itself a probability, then unlike the pmf of a discrete random variable we do not require $f_X(x) \leq 1$.

v) For $a < b$,
$$P(a < X \leq b) = P(a \leq X < b) = P(a \leq X \leq b) = P(a < X < b) = F_X(b) - F_X(a).$$

vi) From Definition 7.0.1 it is clear that the pdf of a continuous random variable X completely characterises its distribution, so we often just specify f_X .

1 follows that a function f_X is a pdf for a continuous random variable X if and only if
i) $f_X(x) \geq 0$,
ii) $\int_{-\infty}^{\infty} f_X(x) dx = 1$

This result follows direct from the definitions and properties of F_X .
Suppose we are interested in whether a continuous random variable X lies in an interval $[a, b]$. Well, $P_X(a < X \leq b) = P_X(X \leq b) - P_X(X \leq a)$, which in terms of the cdf and pdf gives
$$P_X(a < X \leq b) = F_X(b) - F_X(a) = \int_a^b f_X(x) dx.$$

Definition 7.0.2. The cumulative distribution function of CDF, F_X of a continuous random variable X is defined as
$$F_X(x) = P(X \leq x), \quad x \in \mathbb{R}.$$

Suppose that X is a continuous random variable X with pdf f_X and cdf F_X . Let $Y = g(X)$ be a function of X for some (measurable) function $g: \mathbb{R} \rightarrow \mathbb{R}$ s.t. g is continuous and strictly monotonic (so g^{-1} exists). We call $Y = g(X)$ a transformation of X.

Suppose g is monotonic increasing. We can compute the pdf and cdf of $Y = g(X)$ as follows:
The cdf of Y is given by
$$F_Y(y) = P_Y(Y \leq y) = P_Y(g(X) \leq y) = P_X(X \leq g^{-1}(y)) = F_X(g^{-1}(y))$$

The pdf of Y is given by using the chain rule of differentiation:
$$f_Y(y) = F_Y'(y) = f_X(g^{-1}(y))g^{-1}'(y)$$

Note $g^{-1}'(y) = \frac{d}{dy} g^{-1}(y)$ is positive since we assumed g was increasing.
If g monotonic decreasing, we have that
$$F_Y(y) = P_Y(Y \leq y) = P_Y(g(X) \leq y) = P_X(X \geq g^{-1}(y)) = 1 - F_X(g^{-1}(y))$$

By comparison with before, we would have
$$f_Y(y) = F_Y'(y) = -f_X(g^{-1}(y))g^{-1}'(y)$$

with $g^{-1}'(y)$ always negative.
Therefore, for $Y = g(X)$ we have
$$f_Y(y) = f_X(g^{-1}(y))|g^{-1}'(y)|. \quad (7.1)$$

Continuous Random Variables

Mean, Variance and Quantiles
Mean/Expectation

μ_X or E_X(X) = ∫_{-∞}^∞ x f_X(x) dx. E_X{g(X)} = ∫_{-∞}^∞ g(x) f_X(x) dx. E(aX + b) = aE(X) + b, E{g(X) + h(X)} = E{g(X)} + E{h(X)}.

Variance

σ²_X or Var_X(X) = E{(X - μ_X)²} = ∫_{-∞}^∞ (x - μ_X)² f_X(x) dx. Var_X(X) = ∫_{-∞}^∞ x² f_X(x) dx - μ²_X = E(X²) - {E(X)}².

Var(aX + b) = a² Var(X),

Quantiles

Definition 7.1.3. For a (continuous) random variable X we define the a-quantile Q_X(a), 0 ≤ a ≤ 1 to satisfy P(X ≤ Q_X(a)) = a,

Q_X(a) = F_X^{-1}(a).

Continuous Uniform Distribution

Pdf, cdf, a = 0 and b = 1 is referred to as Standard Uniform

f_X(x) = { 0 x ≤ a, (x-a)/(b-a) a < x < b, 1 x ≥ b }

- Suppose X ~ U(0,1), so F_X(x) = x, 0 ≤ x ≤ 1. We wish to map the interval (0,1) to the general interval (a,b), where a < b ∈ ℝ. So we define a new random variable Y = a + (b - a)X, so a < Y < b.

We first observe that for any y ∈ (a,b),

Y ≤ y ⇔ a + (b - a)X ≤ y ⇔ X ≤ (y - a) / (b - a).

From this we find Y ~ U(a,b), since

F_Y(y) = P(Y ≤ y) = P(X ≤ (y - a) / (b - a)) = F_X((y - a) / (b - a)) = (y - a) / (b - a).

- To find the mean of X ~ U(a,b),

E(X) = ∫_{-∞}^∞ x f(x) dx = ∫_a^b x * 1/(b-a) dx = [x²/(2(b-a))]_a^b = (b² - a²) / (2(b-a)) = (b-a)(b+a) / (2(b-a)) = (a+b)/2.

Similarly we get Var(X) = E(X²) - E(X)² = ((b-a)²)/12, so

μ = (a+b)/2, σ² = (b-a)²/12.

Exponential Distribution

Pdf, cdf, θ = 1/λ. parameter of distribution – expectation,

f_X(x) = λ e^{-λx}, x ≥ 0, F_X(x) = 1 - e^{-λx}, x > 0. E(X) = 1/λ = θ, Var(X) = 1/λ².

If X ~ Exp(λ), then, for all x, t > 0,

P(X > x + t | X > t) = P(X > x + t ∩ X > t) / P(X > t) = P(X > x + t) / P(X > t) = e^{-λ(x+t)} / e^{-λt} = e^{-λx} = P(X > x).

Thus, for all x, t > 0, P(X > x + t | X > t) = P(X > x) — this is known as the Lack of Memory Property, and is unique to the exponential distribution amongst continuous distributions.

Normal Distribution

Pdf, cdf,

f_X(x) = 1/(σ√2π) exp{-((x-μ)²)/(2σ²)} F_X(x) = 1/(σ√2π) ∫_{-∞}^x exp{-((t-μ)²)/(2σ²)} dt.

Special Case: If μ = 0 and σ² = 1, then X has a standard or unit normal distribution. The pdf of the standard normal distribution is written as φ(x) and simplifies to

φ(x) = 1/√2π exp{-1/2 x²}.

Also, the cdf of the standard normal distribution is written as Φ(x). Again, for the cdf, we can only write

Φ(x) = 1/√2π ∫_{-∞}^x e^{-t²/2} dt.

F_X(x) = P(X ≤ x) = P(Y ≤ (x-μ)/σ)

= Φ((x-μ)/σ). Φ(z) = 1 - Φ(-z).

Central Limit Theorem

lim_{n→∞} (Σ_{i=1}^n X_i - nμ) / (√nσ) ~ Φ. lim_{n→∞} (X̄ - μ) / (σ/√n) ~ Φ, where X̄ = (Σ_{i=1}^n X_i) / n.

Or for large n

X̄ ~ N(μ, σ²/n) Σ_{i=1}^n X_i ~ N(nμ, nσ²)

Example Consider the most simple example, that X_1, X_2, ... are i.i.d. Bernoulli(p) discrete random variables taking value 0 or 1.

Then the {X_i} have mean μ = p and variance σ² = p(1 - p). Then, by definition, we know that for any n we have

Σ_{i=1}^n X_i ~ Binomial(n, p).

which has mean np and variance np(1 - p).

But now, by the Central Limit Theorem (CLT), we also have for large n that approximately:

Σ_{i=1}^n X_i ~ N(nμ, nσ²) ≡ N(np, np(1 - p)).

So for large n

Binomial(n, p) ≈ N(np, np(1 - p)).

Jointly Distributed Random Variables

Marginal Probability distributions P_X(B) = P(X ∈ B). Joint probability density P_{X,Y}(B_X, B_Y) = P(X ∈ B_X | Y ∈ B_Y). Joint CDF; Recovering Marginal CDFs for X and Y; Joint PMF; Recovering Marginal PMFs p_X and p_Y

F_X(x) = P(X ≤ x, ∞) F_Y(y) = P_Y(∞, y) P_X(x, y) = P_X(x) P_Y(y) = Σ_{y'} P_{X,Y}(x, y')

Properties of valid CDF; Properties of Joint PMF: Conditional Probability

- 1. 0 ≤ F_{X,Y}(x, y) ≤ 1, ∀ x, y ∈ ℝ; E_{X,Y}(Y | X = x) = ∫_{-∞}^∞ y f(y | x) dy.
- 2. Monotonicity: ∀ x_1, x_2, y_1, y_2 ∈ ℝ; x_1 < x_2 ⇒ F_{X,Y}(x_1, y_1) ≤ F_{X,Y}(x_2, y_1) and y_1 < y_2 ⇒ F_{X,Y}(x_1, y_1) ≤ F_{X,Y}(x_1, y_2); 1. 0 ≤ P_{X,Y}(x, y) ≤ 1, ∀ x, y ∈ ℝ;
- 3. ∀ x, y ∈ ℝ, F_{X,Y}(x, -∞) = 0, F_{X,Y}(-∞, y) = 0 and F_{X,Y}(∞, ∞) = 1.

Joint PDFs

P_{X,Y}(B_{X,Y}) = ∫_{(x,y) ∈ B_{X,Y}} f_{X,Y}(x, y) dx dy, F_{X,Y}(x, y) = ∫_{-∞}^x ∫_{-∞}^y f_{X,Y}(s, t) ds dt, f_{X,Y}(x, y) = ∂² / ∂x ∂y F_{X,Y}(x, y)

Recovering marginal densities

Covariance

f_X(x) = d/dx F_X(x) = d/dx F_{X,Y}(x, ∞) = ∫_{-∞}^∞ f_{X,Y}(x, y) dy, f_Y(y) = ∫_{-∞}^∞ f_{X,Y}(x, y) dx.

Independence; Independence; Conditional probability; conditional pdf

P_{X,Y}(B_X, B_Y) = P_X(B_X) P_Y(B_Y) f_{X,Y}(x, y) = f_X(x) f_Y(y) P_{Y|X}(y | B_X) = P_{X,Y}(B_X, B_Y) / P_X(B_X) f_{Y|X}(y | x) = f_{X,Y}(x, y) / f_X(x)

Expectation (x4); Conditional expectation

E_{X,Y}{g(X, Y)} = Σ_{x,y} g(x, y) P_{X,Y}(x, y) E_{X,Y}{g(X, Y)} = ∫_{-∞}^∞ ∫_{-∞}^∞ g(x, y) f_{X,Y}(x, y) dx dy E_{X,Y}(XY) = E_X(X) E_Y(Y)

E_{X,Y}{g_1(X) + g_2(Y)} = E_X{g_1(X)} + E_Y{g_2(Y)} E_{Y|X}(Y | X = x) = ∫_{-∞}^∞ y p(y | x) dy

Estimation

Statistic Function T = T(X_1, ..., X_n) = T(X) and is itself a random variable.

Finding the MLE

In general, we have the following procedure to find MLEs.

- 1. Write down the likelihood function, L(θ) where

L(θ) = ∏_{i=1}^n f(x_i | θ)

that is, the product of the n mass/density functions viewed as a function of θ.

- 2. Take the natural log of the likelihood, and collect terms involving θ.
- 3. Find the value of θ for which log-likelihood is maximised. This is typically done by finding θ that solves

∂/∂θ ℓ(θ) = ∂/∂θ log(L(θ)) = 0

- 4. Check that the estimate θ obtained in step 3 corresponds to a maximum in the (log) likelihood function by inspecting the second derivative of ℓ(θ) wrt θ. If

∂²/∂θ² ℓ(θ) < 0

at θ = θ̂, then θ̂ is confirmed as the MLE of θ.

Confidence Interval

[x̄ - z_{1-α/2} * σ / √n, x̄ + z_{1-α/2} * σ / √n]

(x̄ - μ) / (s_{n-1} / √n) ~ t_{n-1}

where s_{n-1} = √(Σ_{i=1}^n (X_i - X̄)² / (n - 1)) is the bias-corrected sample standard deviation, and t_c is the Student's t-distribution with ν degrees of freedom.

Then it follows that an exact 100(1 - α)% confidence interval for μ is

[x̄ - t_{n-1, 1-α/2} * s_{n-1} / √n, x̄ + t_{n-1, 1-α/2} * s_{n-1} / √n]

where t_{c, ν} is the a-quantile of t_ν.

Estimation

Null Hypothesis (H_0) – Hypothesis that does not change current belief; Alternative Hypothesis (H_1) – Supports your claim.

10.1.1 Normal Distribution with Known Variance

Suppose X_1, ..., X_n are i.i.d. N(μ, σ²) with σ² known and μ unknown. We may wish to test if μ = μ_0 for some specific value μ_0 (e.g. μ_0 = 0, μ_0 = 9.8). Then we can state our null and alternative hypotheses as

H_0: μ = μ_0 versus H_1: μ ≠ μ_0.

Under H_0: μ = μ_0, we then know both μ and σ². So for the sample mean X̄ we have a known distribution for the test statistic

Z = (X̄ - μ_0) / (σ / √n) ~ Φ.

So if we define our rejection region R to be the 100α% tails of the standard normal distribution,

R = (-∞, -z_{1-α/2}) ∪ (z_{1-α/2}, ∞) = {z | |z| > z_{1-α/2}}

we have P(Z ∈ R | H_0) = α.

We thus reject H_0 at the 100α% significance level ⇔ our observed test statistic

z = (X̄ - μ_0) / (σ / √n) ∈ R.

10.2.2 Normal Distributions with Known Variances

Suppose

- X = (X_1, ..., X_n) are i.i.d. N(μ_X, σ_X²) with μ_X unknown;
- Y = (Y_1, ..., Y_n) are i.i.d. N(μ_Y, σ_Y²) with μ_Y unknown;
- the two samples X and Y are independent.

Then we still have that, independently,

X̄ ~ N(μ_X, σ_X²/n), Ȳ ~ N(μ_Y, σ_Y²/n)

From this it follows that the difference in sample means,

X̄ - Ȳ ~ N(μ_X - μ_Y, σ_X²/n + σ_Y²/n),

and hence

(X̄ - Ȳ) - (μ_X - μ_Y) / √(σ_X²/n + σ_Y²/n) ~ Φ.

So under the null hypothesis H_0: μ_X = μ_Y, we have

Z = (X̄ - Ȳ - (μ_X - μ_Y)) / √(σ_X²/n + σ_Y²/n) ~ Φ.

So if σ_X² and σ_Y² are known, we immediately have a test statistic

z = (X̄ - Ȳ) / √(σ_X²/n + σ_Y²/n)

which we can compare against the quantiles of a standard normal.

That is,

R = {z | |z| > z_{1-α/2}}

gives a rejection region for a hypothesis test of H_0: μ_X = μ_Y vs. H_1: μ_X ≠ μ_Y at the 100α% level.

10.2.3 Normal Distributions with Unknown Variances

On the other hand, suppose σ_X² and σ_Y² are unknown. Then if we know σ_X² = σ_Y² = σ² but σ² is unknown, we can still proceed.

We have

(X̄ - Ȳ) - (μ_X - μ_Y) / √(σ²/n + σ²/n) ~ Φ.

and so, under H_0: μ_X = μ_Y,

(X̄ - Ȳ) / √(s_{n1}²/n + s_{n2}²/n) ~ Φ.

Handwritten notes on probability distributions, including derivations for the binomial distribution and the normal distribution, and a table of values for the standard normal distribution.

Example Continuing the Binomial question... each of our Binomial(10, p) samples X_i have pmf

p_X(x_i) = C(10, x_i) p^{x_i} (1 - p)^{10 - x_i}, i = 1, 2, ..., 100.

Since the n = 100 data samples are assumed independent, the likelihood function for p for all of the data is

L(p | x) = L(p) = ∏_{i=1}^n p_X(x_i) = ∏_{i=1}^n C(10, x_i) p^{x_i} (1 - p)^{10 - x_i} = (∏_{i=1}^n C(10, x_i)) p^{Σ_{i=1}^n x_i} (1 - p)^{10n - Σ_{i=1}^n x_i}.

So the log-likelihood is given by

ℓ(p) = log {∏_{i=1}^n C(10, x_i)} + log(p) Σ_{i=1}^n x_i + log(1 - p) (10n - Σ_{i=1}^n x_i).

Next, we differentiate ℓ(p)

∂/∂p ℓ(p) = 0 + Σ_{i=1}^n x_i / p - 10n - Σ_{i=1}^n x_i / (1 - p).

Setting this derivative equal to zero, we get

Σ_{i=1}^n x_i / p - 10n - Σ_{i=1}^n x_i / (1 - p) = 0 ⇒ (1 - p) Σ_{i=1}^n x_i = p̂ (10n - Σ_{i=1}^n x_i) ⇒ Σ_{i=1}^n x_i = p̂ (10n - Σ_{i=1}^n x_i + Σ_{i=1}^n x_i) ⇒ p̂ = Σ_{i=1}^n x_i / 10n = x̄ / 10.

To check this point is a maximum of ℓ, we find the second derivative

∂²/∂p² ℓ(p) = -Σ_{i=1}^n x_i / p² - 10n - Σ_{i=1}^n x_i / (1 - p)² = -n x̄ / p² - 10n - n x̄ / (1 - p)² = -n (x̄ / p² + (10 - x̄) / (1 - p)²)

(which is in fact < 0 ∀ p, the likelihood is log concave).

Substituting p̂ = x̄ / 10, this gives

-100n (x̄ / 10 + (10 - x̄) / 10) = -1000n x̄ / 10 = -100n x̄.

which is clearly < 0. So the MLE for p is p̂ = x̄ / 10 = 0.257.

■

10.1.2 Normal Distribution with Unknown Variance

Similarly, if σ² in the previous example were unknown, we still have that

T = (X̄ - μ_0) / (s_{n-1} / √n) ~ t_{n-1}.

So for a test of H_0: μ = μ_0 vs. H_1: μ ≠ μ_0 at the α level, the rejection region of our observed test statistic t = (X̄ - μ_0) / (s_{n-1} / √n) is

R = (-∞, -t_{n-1, 1-α/2}) ∪ (t_{n-1, 1-α/2}, ∞) = {t | |t| > t_{n-1, 1-α/2}}.

Again, we have that P(T ∈ R | H_0) = α.

but with σ unknown.

We need an estimator for the variance using samples from two populations with different means. Just combining the samples together into one big sample would over-estimate the variance, since some of the variability in the samples would be due to the difference in μ_X and μ_Y.

So we define the bias-corrected pooled sample variance

S²_{n_1+n_2-2} = (Σ_{i=1}^{n_1} (X_i - X̄)² + Σ_{i=1}^{n_2} (Y_i - Ȳ)²) / (n_1 + n_2 - 2),

which is an unbiased estimator for σ².

We can immediately see that s²_{n_1+n_2-2} is indeed an unbiased estimate of σ² by noting

S²_{n_1+n_2-2} = (n_1 - 1) / (n_1 + n_2 - 2) S²_{n_1-1} + (n_2 - 1) / (n_1 + n_2 - 2) S²_{n_2-1}.

That is, s²_{n_1+n_2-2} is a weighted average of the bias-corrected sample variances for the individual samples x and y, which are both unbiased estimates for σ².

Then substituting S_{n_1+n_2-2} in for σ we get

(X̄ - Ȳ) - (μ_X - μ_Y) / (S_{n_1+n_2-2} √(1/n_1 + 1/n_2)) ~ t_{n_1+n_2-2}.

and so, under H_0: μ_X = μ_Y,

T = (X̄ - Ȳ) / (S_{n_1+n_2-2} √(1/n_1 + 1/n_2)) ~ t_{n_1+n_2-2}.

So we have a rejection region for a hypothesis test of H_0: μ_X = μ_Y vs. H_1: μ_X ≠ μ_Y at the 100α% level given by

R = {t | |t| > t_{n_1+n_2-2, 1-α/2}},

for the statistic

t = (X̄ - Ȳ) / (S_{n_1+n_2-2} √(1/n_1 + 1/n_2)).

Suppose we have a null hypothesis H_0: θ = θ_0 for the value of the unknown parameter(s). Then under H_0 we know the pmf {p_j}, and so we are able to calculate the expected frequency counts E = (E_1, ..., E_k) by E_j = np_j. (Note again we have Σ_{j=1}^k E_j = n.)

We then seek to compare the observed frequencies with the expected frequencies to test for goodness of fit.

To test H_0: θ = θ_0 vs. H_1: θ ≠ θ_0 we use the chi-square statistic

χ² = Σ_{i=1}^k (O_i - E_i)² / E_i.

If H_0 were true, then the statistic χ² would approximately follow a chi-square distribution with ν = k - 1 degrees of freedom.

• k is the number of values (categories) the simple random variable X can take.

• p is the number of parameters being estimated (dim(θ)).

• For the approximation to be valid, we should have ν_j E_j ≥ 5. This may require some merging of categories.

Handwritten calculations for the chi-square statistic, including the formula for χ² and the calculation of expected frequencies E_j.