

IB Statistics

Ishan Nath, Lent 2023

Based on Lectures by Dr. Sergio Bacallado

January 24, 2023

Contents

1	Introduction	2
1.1	Probability Review	2
1.2	Moment Generating Functions	4
1.3	Limit Theorems	4
1.4	Conditioning	5
1.5	Change of Variables	6
1.6	Important Distributions	6
2	Estimation	7
2.1	Bias-Variance Decomposition	8
	Index	10

1 Introduction

Statistics is the science of making informed decisions. It can include:

- Design of experiments,
- Graphical exploration of data,
- Formal statistical inference (part of Decision theory),
- Communication of results.

Let X_1, X_2, \dots, X_n be independent observations from a distribution $f(x \mid \theta)$, with parameter θ . We wish to make inferences about the value of θ from X_1, X_2, \dots, X_n . Such inference can include:

- Estimating θ ,
- Quantifying uncertainty in estimates,
- Testing a hypothesis about θ .

1.1 Probability Review

Let Ω be the *sample space* of outcomes in an experiment. A measurable subset of Ω is called an *event*. We denote the set of events as \mathcal{F} .

A function $\mathbb{P} : \mathcal{F} \rightarrow [0, 1]$ is called a *probability measure* if:

- $\mathbb{P}(\emptyset) = 0$,
- $\mathbb{P}(\Omega) = 1$,
- $\mathbb{P}(\bigcup_i A_i) = \sum_i \mathbb{P}(A_i)$, if (A_i) are disjoint and countable.

A *random variable* is a (measurable) function $X : \Omega \rightarrow \mathbb{R}$.

The *distribution function* of X is

$$F_X(x) = \mathbb{P}(X \leq x).$$

A *discrete random variable* takes values in a countable subset $E \subset \mathbb{R}$, and its *probability mass function* or pmf is $p_X(x) = \mathbb{P}(X = x)$.

We say X has *continuous* distribution if it has a *probability density function* or pdf, satisfying

$$\mathbb{P}(X \in A) = \int_A f_X(x) \, dx,$$

for any measurable A . The *expectation* of X is defined

$$\mathbb{E}[X] = \begin{cases} \sum_{x \in X} x \cdot p_X(x) & X \text{ discrete,} \\ \int x \cdot f_X(x) dx & X \text{ continuous.} \end{cases}$$

If $g : \mathbb{R} \rightarrow \mathbb{R}$, then

$$\mathbb{E}[g(x)] = \int g(x) f_X(x) dx.$$

The *variance* of X is

$$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2].$$

We say that X_1, X_2, \dots, X_n are *independent* if for all x_1, x_2, \dots, x_n ,

$$\mathbb{P}(X_1 \leq x_1, \dots, X_n \leq x_n) = \mathbb{P}(X_1 \leq x_1) \cdots \mathbb{P}(X_n \leq x_n).$$

If the variables have probability density functions, then

$$f_X(x) = \prod_{i=1}^n f_{X_i}(x_i),$$

where X is the vector of variables (X_1, \dots, X_n) and x is the vector (x_1, \dots, x_n) .

Importantly, if $a_1, \dots, a_n \in \mathbb{R}$,

$$\mathbb{E}[a_1 X_1 + \cdots + a_n X_n] = a_1 \mathbb{E}[X_1] + \cdots + a_n \mathbb{E}[X_n].$$

Moreover,

$$\text{Var}(a_1 X_1 + \cdots + a_n X_n) = \sum_{i,j} a_i a_j \text{Cov}(X_i, X_j).$$

Here the *covariance* of X_i and X_j is

$$\text{Cov}(X_i, X_j) = \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])].$$

If $X = (X_1, \dots, X_n)^T$ and $\mathbb{E}[X] = (\mathbb{E}[X_1], \dots, \mathbb{E}[X_n])$, then the linearity of expectation can be rewritten as

$$\mathbb{E}[a^T X] = a^T \mathbb{E}[X],$$

and moreover

$$\text{Var}(a^T X) = a^T \text{Var}(X) a,$$

where $\text{Var}(X)$ is the *covariance matrix*: $(\text{Var}(X))_{ij} = \text{Cov}(X_i, X_j)$.

1.2 Moment Generating Functions

The *moment generating function* of a variable X is

$$M_X(t) = \mathbb{E}[e^{tx}].$$

This may only exist for t in some neighbourhood of 0. The important properties of MGFs is that

$$\mathbb{E}[X^n] = \frac{d^n}{dt^n} M_X(0),$$

and from this we obtain $M_X = M_Y \iff F_x = F_y$.

MGFs also make it easy to find the distribution function of sums of iid variables.

Example 1.1.

Let X_1, \dots, X_n be iid Poisson(μ). Then

$$\begin{aligned} M_{X_1}(t) &= \mathbb{E}[e^{tX_1}] = \sum_{x=0}^{\infty} e^{tx} \cdot \frac{e^{-\mu} \mu^x}{x!} \\ &= e^{-\mu} \sum_{x=0}^{\infty} \frac{(e^t \mu)^x}{x!} = e^{-\mu} e^{\mu \exp(t)} = e^{-\mu(1-e^t)}. \end{aligned}$$

If $S_n = X_1 + \dots + X_n$, then

$$\begin{aligned} M_{S_n}(t) &= \mathbb{E}[e^{t(X_1 + \dots + X_n)}] = \prod_{i=1}^n \mathbb{E}[e^{tX_i}] \\ &= e^{-\mu(1-e^t)n} \end{aligned}$$

This is the same as a Poisson(μn) MGF, so $S_n \sim \text{Poisson}(\mu \cdot n)$.

1.3 Limit Theorems

We list some important limit theorems, starting with the *weak law of large numbers* (WLLN). This says if X_1, \dots, X_n are iid with $\mathbb{E}[X_1] = \mu$, then let $\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ be the sample mean. WLLN says that for all $\varepsilon > 0$,

$$\mathbb{P}(|\overline{X}_n - \mu| > \varepsilon) \rightarrow 0,$$

as $n \rightarrow \infty$.

The *strong law of large numbers* (SLLN) says a stronger result, namely

$$\mathbb{P}(\overline{X}_n \rightarrow \mu) = 1,$$

i.e. $\overline{X_n}$ converges to μ almost surely.

The *central limit theorem* is another important limit theorem. If we take

$$Z_n = \frac{\sqrt{n}(\overline{X_n} - \mu)}{\sigma},$$

where $\sigma^2 = \text{Var}(X_i)$, then Z_n is “approximately” $N(0, 1)$ as $n \rightarrow \infty$.

What this means is that $\mathbb{P}(Z_n \leq z) \rightarrow \Phi(z)$ as $n \rightarrow \infty$ for all $z \in \mathbb{R}$, where Φ is the distribution function of a $N(0, 1)$ variable.

1.4 Conditioning

Let X and Y be discrete random variables. Their *joint pmf* is

$$p_{X,Y}(x, y) = \mathbb{P}(X = x, Y = y).$$

The *marginal pmf* is

$$p_X(x) = \mathbb{P}(X = x) = \sum_{y \in Y} p_{X,Y}(x, y).$$

The *conditional pmf* of X given $Y = y$ is

$$p_{X|Y}(x | y) = \mathbb{P}(X = x | Y = y) = \frac{\mathbb{P}(X = x, Y = y)}{\mathbb{P}(Y = y)} = \frac{p_{X,Y}(x, y)}{p_Y(y)}.$$

This is defined to be 0 if $p_Y(y) = 0$.

For continuous random variables X, Y , the *joint pdf* $f_{X,Y}$ has

$$\mathbb{P}(X \leq x', y \leq y') = \int_{-\infty}^{x'} \int_{-\infty}^{y'} f_{X,Y}(x, y) \, dy \, dx.$$

The *marginal pdf* of Y is

$$f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) \, dx.$$

The *conditional pdf* of X given Y is

$$f_{X|Y}(x | y) = \frac{f_{X,Y}(x, y)}{f_Y(y)}.$$

The *conditional expectation* is given by

$$\mathbb{E}[X | Y] = \begin{cases} \sum_x x \cdot p_{X|Y}(x | y) & X, Y \text{ discrete,} \\ \int_{-\infty}^{\infty} x \cdot f_{X|Y}(x | y) dx & X, Y \text{ continuous.} \end{cases}$$

This is a random variable, which is a function of Y . The *tower property* says that

$$\mathbb{E}[\mathbb{E}[X | Y]] = \mathbb{E}[X].$$

Hence we can write the variance of X as follows:

$$\begin{aligned} \text{Var}(X) &= \mathbb{E}[X^2] - (\mathbb{E}[X])^2 \\ &= \mathbb{E}[\mathbb{E}[X^2 | Y]] - (\mathbb{E}[\mathbb{E}[X | Y]])^2 \\ &= \mathbb{E}[\mathbb{E}[X^2 | Y] - (\mathbb{E}[X | Y])^2] + \mathbb{E}[\mathbb{E}[X | Y]^2] - \mathbb{E}[\mathbb{E}[X | Y]]^2 \\ &= \mathbb{E}[\text{Var}(X | Y)] + \text{Var}(\mathbb{E}[X | Y]). \end{aligned}$$

1.5 Change of Variables

The *change of variables* formula is as follows:

Let $(x, y) \mapsto (u, v)$ be a differentiable bijection. Then,

$$\begin{aligned} f_{U,V}(u, v) &= f_{X,Y}(x(u, v), y(u, v)) \cdot |\det J|, \\ J &= \frac{\partial(x, y)}{\partial(u, v)} = \begin{pmatrix} \partial x / \partial u & \partial x / \partial v \\ \partial y / \partial u & \partial y / \partial v \end{pmatrix}. \end{aligned}$$

1.6 Important Distributions

$X \sim \text{Negbin}(k, p)$ if X models the time in successive iid $\text{Ber}(p)$ trials to achieve k successes. If $k = 1$, this is the same as a geometric distribution.

$X \sim \text{Poisson}(\lambda)$ is the limit of $\text{Bin}(n, \lambda/n)$ random variables, as $n \rightarrow \infty$.

If $X_i \sim \Gamma(\alpha_i, \lambda)$ for $i = 1, \dots, n$ with X_1, \dots, X_n independent, then if $S_n = X_1 + \dots + X_n$,

$$M_{S_n}(t) = \prod_{i=1}^n M_{X_i}(t) = \left(\frac{\lambda}{\lambda - 1} \right)^{\alpha_1 + \dots + \alpha_n}$$

which is the mgf of a $\Gamma(\sum \alpha_i, \lambda)$ random variable. Hence $S_n \sim \Gamma(\sum \alpha_i, \lambda)$.

Also, if $X \sim \Gamma(a, \lambda)$, then for any $b \in (0, \infty)$, $bX \sim \Gamma(a, \lambda/b)$.

Special cases of the Gamma distribution include $\Gamma(1, \lambda) = \text{Exp}(\lambda)$, and $\Gamma(\frac{k}{2}, \frac{1}{2}) = \chi_k^2$, the Chi-squared distribution with k degrees of freedom. This can be thought of as the sum of k independent squared $N(0, 1)$ random variables.

2 Estimation

Suppose we observe data X_1, X_2, \dots, X_n , which are iid from some pdf (or pmf) $f_X(x | \theta)$, with θ unknown. We let $X = (X_1, \dots, X_n)$.

Definition 2.1. An *estimator* is a statistic or a function of the data $T(X) = \hat{\theta}$, which we use to approximate the true parameter θ . The distribution of $T(X)$ is called the *sampling distribution*.

Example 2.1.

If X_1, \dots, X_n are iid $N(\mu, 1)$, we can define an estimator for the mean as

$$\hat{\mu} = T(X) = \frac{1}{n} \sum_{i=1}^n X_i.$$

The sampling distribution of $\hat{\mu}$ is $N(\mu, \frac{1}{n})$.

Definition 2.2. The *bias* of $\hat{\theta} = T(X)$ is

$$\text{bias}(\hat{\theta}) = \mathbb{E}_\theta[\hat{\theta}] - \theta.$$

Remark. In general, the bias is a function of θ , even if the notation $\text{bias}(\hat{\theta})$ does not make that explicit.

Definition 2.3. We say that $\hat{\theta}$ is *unbiased* if $\text{bias}(\hat{\theta}) = 0$ for all $\theta \in \Theta$.

Example 2.2.

Our previous estimator

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

is unbiased because $\mathbb{E}_\mu[\hat{\mu}] = \mu$ for all $\mu \in \mathbb{R}$.

Definition 2.4. The *mean squared error* (mse) of $\hat{\theta}$ is

$$\text{mse}(\hat{\theta}) = \mathbb{E}_\theta[(\hat{\theta} - \theta)^2].$$

Like the bias, the mean squared error of $\hat{\theta}$ is a function of θ .

2.1 Bias-Variance Decomposition

We can write the mean squared error as

$$\begin{aligned} \text{mse}(\hat{\theta}) &= \mathbb{E}_\theta[(\hat{\theta} - \theta)^2] = \mathbb{E}_\theta[(\hat{\theta} - \mathbb{E}_\theta[\hat{\theta}] + \mathbb{E}_\theta[\hat{\theta}] - \theta)^2] \\ &= \text{Var}_\theta(\hat{\theta}) + \text{bias}^2(\hat{\theta}) + 2 \underbrace{[\mathbb{E}_\theta[(\hat{\theta} - \mathbb{E}_\theta[\hat{\theta}])]]}_{0} (\mathbb{E}_\theta[\hat{\theta}] - \theta). \end{aligned}$$

The two terms on the right hand side are non-negative, so there is a trade off between bias and variance.

Example 2.3.

Let $X \sim \text{Bin}(n, \theta)$, where n is known, and we wish to estimate θ . The standard estimator is

$$T_u = \frac{X}{n}, \quad \mathbb{E}_\theta[T_u] = \frac{\mathbb{E}_\theta[X]}{n} = \theta.$$

Hence T_u is unbiased. We can also calculate the mean squared error as

$$\text{mse}(T_u) = \text{Var}_\theta(T_u) = \frac{\text{Var}_\theta(X)}{n^2} = \frac{n\theta(1-\theta)}{n^2} = \frac{\theta(1-\theta)}{n}.$$

Consider a second estimator

$$T_B = \frac{X+1}{n+2} = w \frac{X}{n} + (1-w) \frac{1}{2},$$

for $w = \frac{n}{n+2}$. In this case T_B is interpolating between our unbiased estimator, and the constant estimator. The bias of T_B is

$$\text{bias}(T_B) = \mathbb{E}_\theta[T_B] - \theta = \mathbb{E}\left[\frac{X+1}{n+2}\right] - \theta = \frac{1}{n+2} - \frac{2}{n+2}\theta.$$

This is not equal to zero for all but one value of θ . Hence, T_B is biased. We can also calculate the variance

$$\begin{aligned} \text{Var}_\theta(T_B) &= \frac{1}{(n+2)^2} n\theta(1-\theta) - w^2 \frac{\theta(1-\theta)}{n}, \\ \text{mse}(T_B) &= \text{Var}_\theta(T_B) + \text{bias}^2(T_B) \\ &= w^2 \frac{\theta(1-\theta)}{n} + (1-w)^2 \left(\frac{1}{2} - \theta\right)^2. \end{aligned}$$

Hence the mse of the biased estimator is a weighted average of the mse of the unbiased estimator, and a parabola. For θ around $1/2$, the biased estimator has a lower mse than the unbiased estimator.

The message here is that our prior judgements about θ affect our choice of estimator, and unbiasedness is not always desirable.

Example 2.4.

Suppose $X \sim \text{Poisson}(\lambda)$. We wish the estimate $\theta = \mathbb{P}(X = 0)^2 = e^{-2\lambda}$. For an estimator $T(X)$ to be unbiased, we must have for all λ ,

$$\begin{aligned}\mathbb{E}_\lambda[\hat{\theta}] &= \sum_{x=0}^{\infty} T(x) \frac{e^{-\lambda} \lambda^x}{x!} = e^{-2\lambda} = \theta \\ \iff \sum_{x=0}^{\infty} T(x) \frac{\lambda^x}{x!} &= e^{-\lambda} = \sum_{x=0}^{\infty} (-1)^x \frac{\lambda^x}{x!}.\end{aligned}$$

For this to hold for all $\lambda \geq 0$, we should take $T(X) = (-1)^X$. But this estimator makes no sense.

Index

- bias, 7
- central limit theorem, 5
- change of variables, 6
- conditional expectation, 6
- conditional probability density function, 5
- conditional probability mass function, 5
- continuous random variable, 2
- covariance, 3
- discrete random variable, 2
- distribution function, 2
- estimator, 7
- event, 2
- expectation, 3
- independence, 3
- joint probability density function, 5
- joint probability mass function, 5
- marginal probability density function, 5
- marginal probability mass function, 5
- mean squared error, 7
- probability density function, 2
- probability mass function, 2
- probability measure, 2
- random variable, 2
- sample space, 2
- sampling distribution, 7
- strong law of large numbers, 4
- tower property, 6
- unbiased estimator, 7
- variance, 3
- weak law of large numbers, 4