

STEP Support Programme

STEP II Statistics Topic Notes

A lot of these formulae can be found in the STEP formulae book!

Probability A good introduction to basic probability can be found here.

- $P(A \cap B)$ means the probability that both A and B happen. If A and B are *mutually exclusive* then $P(A \cap B) = 0$. If A and B are *independent* then $P(A \cap B) = P(A) \times P(B)$.
- $P(A \cup B)$ means the probability that A or B (or both) happen.
- $P(A \cup B) = P(A) + P(B) P(A \cap B)$. To show this, draw a Venn diagram with two overlapping circles, the area inside one representing P(A) and the other representing P(B).
- P(A|B) means the probability that event A happens **given** that we know event B happens. This is a *conditional* probability. A lot of conditional probability questions can be done informally using tree diagrams, or by considering a population (e.g. 100,000 people).
- **Bayes' Theorem** this is not strictly on the specification, as questions can be answered using "informal methods", but it can be useful.

$$\mathbf{P}(A|B) = \frac{\mathbf{P}(A \cap B)}{P(B)}.$$

It can also be useful to consider the equivalent statement $P(A \cap B) = P(B) \times P(A|B)$.

Writing $P(A \cap B)$ in two different ways and equating gives us $P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$.

Combinations, permutations and arrangements

- The number of ways of choosing¹ r objects from n objects is ${}^{n}C_{r} = {n \choose r} = \frac{n!}{r!(n-r)!}$.
- The number of permutations² of r objects taken from a selection of n different objects is ${}^{n}P_{r} = \frac{n!}{(n-r)!}.$
- The number of different arrangements³ of r objects where a of them are all the same, another b are all the same (but different to the first lot) etc. is $\frac{r!}{a! \times b! \times \cdots}$.



 $^{^1}$ "Choosing" implies that the order doesn't matter.

² "Permutations" means that order does matter.

³ Order matters.



Discrete Probability Distributions

"Discrete" means that only certain values can be taken (such as the numbers on a dice — we cannot get a value between 2 and 3)⁴.

Let X be a discrete random variable.

• The *expectation* (or mean) is given by:

$$\mathbf{E}(X) = \sum i \times \mathbf{P}(X = i).$$

So if the possible values of X are 1, 2, ..., k then $E(X) = 1 \times P(X = 1) + 2 \times P(X = 2) + ... + k \times P(X = k).$

• The *variance* is given by:

$$Var(X) = E(X^2) - [E(X)]^2$$

where $E(X^2) = \sum i^2 \times P(X = i)$. Var(X) is never negative, and can be thought of as "the mean of the squares – the square of the mean"⁵.

• The *Mode* or *Modal Value* is the value of x for which P(X = x) is the greatest (there may be more than one mode).

The *binomial* distribution, $X \sim B(n, p)$, is the distribution of the number of "successes" in a sequence of n independent yes/no trials each of which has a probability p of success. An example would be the number of sixes you get when you roll a dice 10 times.

• $P(X=r) = \binom{n}{r} p^r (1-p)^{n-r}$

•
$$E(X) = np$$

•
$$\operatorname{Var}(X) = np(1-p)$$

The *Poisson* distribution, $X \sim Poi(\lambda)$, is the distribution of the number of occurrences of an event in a given "interval" (which can be time, length, etc.). An example could be the number of meteors greater than 1 meter diameter that strike earth in a year.

•
$$P(X=r) = \frac{e^{-\lambda}\lambda^n}{n!}$$

- $E(x) = Var(X) = \lambda$
- If the number of occurrences in an interval of length T follows a Poisson distribution with mean λ , then the number of occurrences in an interval of length kT follows a Poisson distribution with mean $k\lambda$.
- if X and Y are two *independent* Poisson random variables with means λ and μ respectively then X + Y has a Poisson distribution with mean $\lambda + \mu$.
- If n is "large" and p is "very small" then a Poisson distribution with mean np can be used to approximate a Binomial distribution, $X \sim B(n, p)$.

⁵The above definition of variance is usually the easiest to work with, but variance is really the mean squared distance of values from the mean. This gives $Var(X) = E((X - E(X))^2) = \sum (i - E(X))^2 \times P(X = i)$ which can be expanded to give the above result.



⁴ There can be infinitely many values, such as the number of coin tosses until you get a head, or non-integer values, such as the value of one dice roll divided by another dice roll — we can still only get certain values, such as $\frac{4}{3}$, but not others, such as $\frac{2}{7}$.



Continuous Probability Distributions

"Continuous" means that all the values in a certain range are possible. Examples include height of a person, or the half life of a radioactive element. Continuous random variables are usually defined by a probability distribution function, f(x).

- $P(a \leqslant X \leqslant b) = \int_a^b f(x) dx$
- $\int_{\text{all } x} f(x) \, \mathrm{d}x = 1$ (as the total probability must be 1)
- The *expectation* (or mean) is given by:

$$\mu = \int_{\text{all } x} x \mathbf{f}(x) \, \mathrm{d}x \,.$$

• The *variance* is given by:

$$\int_{\text{all }x} x^2 \mathbf{f}(x) \, \mathrm{d}x \, -\mu^2.$$

Note that the formulae for expectation and variance of a continuous distribution are very similar to the ones for a discrete distribution, all that has happened is that the sum has been replaced by an integral. In fact Leibniz considered integration to be an infinite sum of infinitesimal "bits". and so he based the integral symbol \int on the "long s" character (for "summation").

• The *cumulative distribution function* is defined by:

$$\mathbf{F}(x) = \mathbf{P}(X \leqslant x) = \int_{-\infty}^{x} \mathbf{f}(t) \, \mathrm{d}t$$

Here we have taken the lower limit as $-\infty$. It may be that f(x) = 0 for $-\infty < x < a$ say, in which case we could write the lower limit as a. Note the use of "dummy variable" t inside the integral — we cannot use x inside the integral as it is used as a limit.

- The median, m satisfies $P(X \le m) = P(X \ge m) = \frac{1}{2}$, i.e. $\int_{-\infty}^{m} f(x) dx = \int_{m}^{\infty} f(x) dx = \frac{1}{2}$.
- The *mode* is where the probability distribution function has a maximum (there may be more than one!).

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

- If $X \sim B(n, p)$ and n is "large" and/or p is "close to" $\frac{1}{2}$ then X can be approximated by a normal distribution, $X \sim N(np, np(1-p))$.
- If $X \sim Poi(\lambda)$ and λ is "large" then X can by approximated by a normal distribution, $X \sim N(\lambda, \lambda^2)$.

