

Further Probability & Statistics

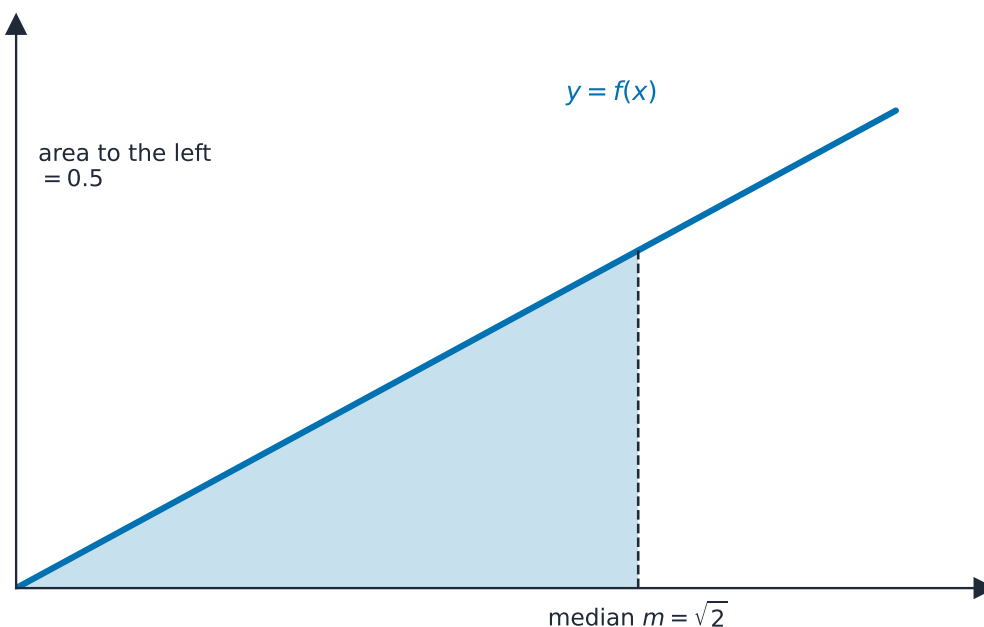
A-Level Further Mathematics

This handout covers Topic 4: **Further Probability & Statistics** 进阶概率统计. It adds continuous distributions, small-sample inference, the chi-squared and non-parametric tests, and probability generating functions.

Continuous random variables

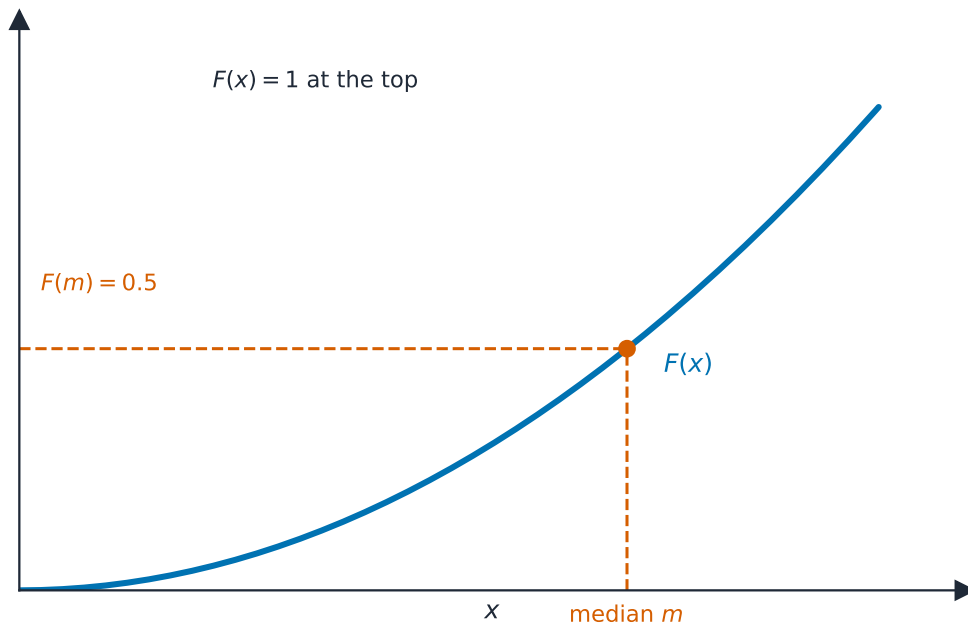
A continuous variable X is described by a **probability density function** 概率密度函数 $f(x)$, which may be defined piecewise. The probability over a range is the area under f , and the mean of any function of X is

$$E(g(X)) = \int g(x) f(x) dx.$$



The median sits where the area under $f(x)$ to its left is exactly 0.5.

The **cumulative distribution function** 累积分布函数 $F(x) = P(X \leq x)$ is the running total: $F(x) = \int_{-\infty}^x f(t) dt$, and $f(x) = F'(x)$. Use F to find probabilities and **percentiles** 百分位数 (for example, the **median** 中位数 solves $F(x) = 0.5$).



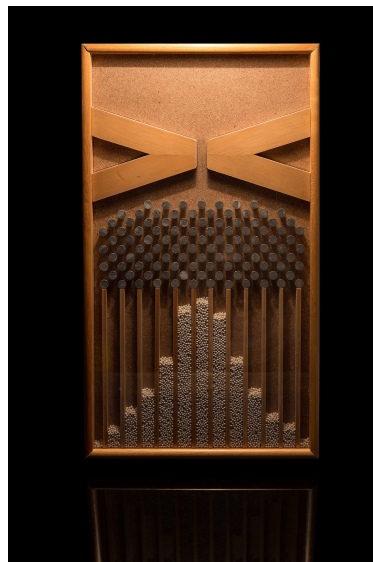
The cumulative graph $F(x)$ rises from 0 to 1; the height 0.5 is reached at the median.

Worked example. A variable has $f(x) = \frac{1}{2}x$ for $0 \leq x \leq 2$. Find the median.

The cumulative distribution function is $F(x) = \int_0^x \frac{1}{2}t \, dt = \frac{1}{4}x^2$. Set $F(m) = 0.5$:

$$\frac{1}{4}m^2 = 0.5 \Rightarrow m^2 = 2 \Rightarrow m = \sqrt{2} = 1.41.$$

Inference using normal and t-distributions



A Galton board: balls falling through pins pile up into the bell-shaped normal distribution.

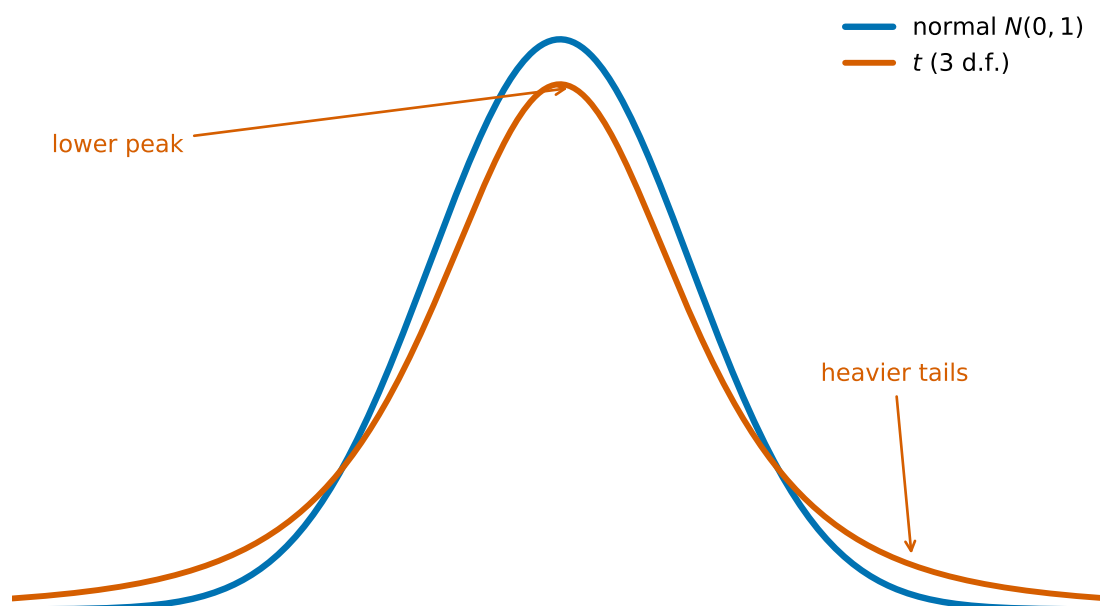
Image: Exhibit made by Estes Objethos Atelier, photo by Rodrigo.Argenton, CC BY-SA 4.0 (commons.wikimedia.org)

When a sample is small and the population variance is unknown, base your **hypothesis test** 假设检验 on the t -distribution instead of the normal. The same idea gives a

confidence interval 置信区间 for the mean:

$$\bar{x} \pm t \frac{s}{\sqrt{n}},$$

where t comes from the t -tables with $n - 1$ degrees of freedom. To compare two populations, use a 2-sample or paired-sample t -test, after finding a **pooled estimate** 合并估计 of the shared variance when appropriate.



For a small sample the t -distribution is flatter with heavier tails, so its critical values are larger.

Worked example. A sample of $n = 10$ has mean $\bar{x} = 50$ and standard deviation $s = 4$. Find a 95% confidence interval for the mean (use $t = 2.262$ for 9 degrees of freedom).

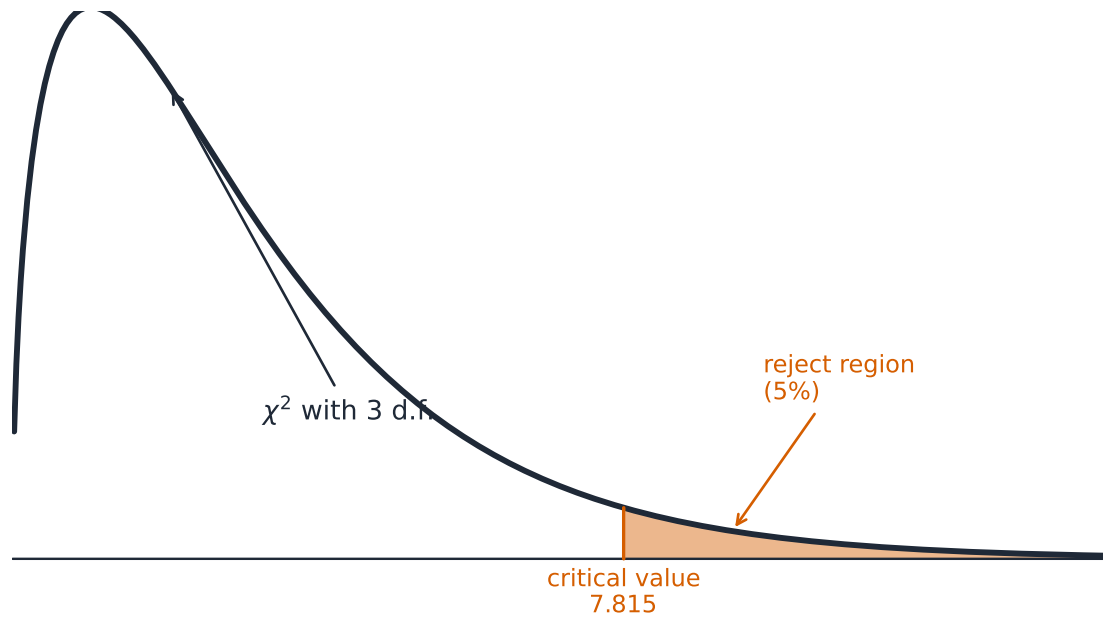
$$50 \pm 2.262 \times \frac{4}{\sqrt{10}} = 50 \pm 2.86 \Rightarrow (47.1, 52.9).$$

Chi-squared tests

A χ^2 -test (**chi-squared test** 卡方检验) compares observed counts O with expected counts E from a **theoretical distribution** 理论分布:

$$\chi^2 = \sum \frac{(O - E)^2}{E}.$$

Compare this with a table value for the right number of **degrees of freedom** 自由度. Two uses: a **goodness of fit** 拟合优度 test (does the data follow the proposed model?), and a test for **independence** 独立性 of two variables in a **contingency table** 列联表.



The test rejects the model when χ^2 exceeds the critical value, landing in the shaded 5% tail.

Worked example. Four equally likely categories give observed counts 20, 30, 25, 25 (so each expected count is 25). Test the fit at the 5% level.

$$\chi^2 = \frac{(20 - 25)^2 + (30 - 25)^2 + 0 + 0}{25} = \frac{25 + 25}{25} = 2.$$

With $4 - 1 = 3$ degrees of freedom the table value is 7.815. Since $2 < 7.815$, do not reject the model.

Non-parametric tests

A **non-parametric test** 非参数检验 makes no assumption that the data is normal, so it is useful when that assumption fails. The basic ones are:

- the **sign test** 符号检验: count how many values fall above and below a proposed median, and test those counts with a binomial model;
- the **Wilcoxon signed-rank test** 威尔科克森符号秩检验, which also uses the sizes of the differences, not just their signs;
- the **Wilcoxon rank-sum test** 威尔科克森秩和检验, for comparing two separate samples.

Probability generating functions



Dice: a starting point for probability and discrete random variables.

Image: Diacritica, CC BY-SA 3.0 (commons.wikimedia.org)

The **probability generating function** 概率母函数 of a discrete variable X is

$$G(t) = E(t^X) = \sum_x P(X = x) t^x.$$

It packs the whole distribution into one function. The mean and variance come from its derivatives at $t = 1$: $E(X) = G'(1)$ and $\text{Var}(X) = G''(1) + G'(1) - (G'(1))^2$. Also, the PGF of a sum of independent variables is the product of their PGFs.

Worked example. X has $P(X = 0) = 0.5$, $P(X = 1) = 0.3$, $P(X = 2) = 0.2$. Find $E(X)$ using the PGF.

Here $G(t) = 0.5 + 0.3t + 0.2t^2$, so $G'(t) = 0.3 + 0.4t$ and

$$E(X) = G'(1) = 0.3 + 0.4 = 0.7.$$