

The Central Limit Theorem

In the past few lectures we have studied statistical inference and have often assumed we have normally distributed observations. Two reasons for this are

- Quite often real data is (approximately) normal (although this must be verified, for example, with a q-q-plot)
- The analysis is tractable.

This lecture discusses a third reason, which hinges on a key result called the central limit theorem.

The Central Limit Theorem

Let X_1, \dots, X_n be a random sample
(not necessarily normal) with mean μ
and variance $\sigma^2 > 0$. Denote

$$\text{We have } \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \text{ and } S_n = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}.$$

Then

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \xrightarrow{D} N(0, 1)$$

as $n \rightarrow \infty$.

(The precise nature of the convergence
is not covered in this course.)

Practically, we can use the CLT
to say, for n "sufficiently large,"

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \approx N(0, 1).$$

Recall, ^{when} we derived confidence intervals
and tests for a normal sample (σ^2
known), we only used the fact

that $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$.

By the CLT, even if the sample is
not normal, it is still true that

$$P\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} > z_\alpha\right) \approx \alpha.$$

Therefore, by the same analysis as before,
we can show that

$$\left[\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$$

is an approximate $(100(1-\alpha))\%$
confidence interval for μ .

Similarly, the test

$$\begin{cases} H_0 & \frac{|\bar{X} - \mu|}{\sigma/\sqrt{n}} < z_{\alpha/2} \\ H_1 & \frac{|\bar{X} - \mu|}{\sigma/\sqrt{n}} > z_{\alpha/2} \end{cases}$$

has level of significance α .
approximately

Note that the CLT does not assume the data is continuous. It also hold for discrete data.

Example $X_1, \dots, X_n \sim \text{Poisson}(\lambda)$.

Then $\mu = \lambda$, $\sigma^2 = \lambda$.

By the CLT,

$$\frac{\bar{X} - \lambda}{\sqrt{\lambda/n}} \approx N(0, 1).$$

Observe

$$\frac{\bar{X} - \lambda}{\sqrt{\lambda/n}} = \frac{\sum_{i=1}^n X_i - n\lambda}{\sqrt{n\lambda}}$$

Thus, $\sum_{i=1}^n X_i \approx N(n\lambda, \lambda)$

But $\sum_{i=1}^n X_i \sim \text{Poisson}(n\lambda)$. Making the substitution

$n\lambda \rightarrow \lambda$, we conclude

$\text{Poisson}(\lambda) \approx N(\lambda, \lambda)$ for
~~Poisson~~ large λ .

How large does n need to be
for the CLT approximation to
hold? There are no guarantees,
but here are some rules of thumb:

- If $n \approx 25$ or so, the approximation probably holds.
- If X is symmetric, unimodal (only one peak), and continuous, then $n \approx 5$ may be sufficient
(Example: t , χ^2 , Gamma).

See examples in book.