

Chapter 1 continued

Tests and confidence intervals for the median

The two procedures introduced in the previous lecture were a test and a confidence interval for the population median $\theta_{.5}$ of a continuous random variable, both based on using the binomial distribution, since $P(X_i < \theta_{.5}) = P(X_i > \theta_{.5}) = .5$. When testing, for example, $H_0 : \theta_{.5} = \theta_H$ versus $H_a : \theta_{.5} > \theta_H$, we can consider the number of observations X_i that fall above the hypothesized value of $\theta_{.5}$. If we denote this number by B then our test for the median is equivalent to a test of the binomial distribution probability p , and our p value is the probability of getting B or more successes in n trials assuming that $p = .5$. This p value can be calculated directly as:

$$\sum_{x=B}^n \binom{n}{x} (.5)^n, \text{ or approximated by using } Z_B = \frac{B - n/2}{\sqrt{n/4}} \text{ or } Z_B = \frac{(B - 1/2) - n/2}{\sqrt{n/4}},$$

and the standard normal distribution. The second Z_B value uses a continuity correction.

To construct a confidence interval for the median, we use similar ideas about the binomial distribution. We use the order statistics $X_{(1)} < X_{(2)} < \dots < X_{(n)}$, and find values a and b so that $P(X_{(a)} < \theta_{.5} < X_{(b)}) = 1 - \alpha$ for some specified value α . This probability can be expressed as a sum of binomial probabilities, because the event $X_i < \theta_{.5}$ follows a binomial distribution with parameters n and $p = .5$. If $X_{(a)} < \theta_{.5} < X_{(b)}$ is true then at least a X_i 's satisfy $X_i < \theta_{.5}$, and at most $b - 1$ X_i 's satisfy $X_i < \theta_{.5}$. Thus we want the probability that between a and $b - 1$ values of X_i are less than $\theta_{.5}$, which can be directly calculated by

$$\sum_{x=a}^{b-1} \binom{n}{x} (.5)^n.$$

We can either find a and $b - 1$ directly so that the above probability equals $1 - \alpha$ (approximately), or we can use a normal approximation to find a and b so that:

$$\frac{a - n/2}{\sqrt{n/4}} = -z_{(1-\alpha/2)} \text{ and } \frac{(b - 1) - n/2}{\sqrt{n/4}} = z_{(1-\alpha/2)}$$

and use an appropriate integer near a and b to identify the desired order statistics. Note that if a continuity correction is used in the normal approximation, then we find a and b according to:

$$\frac{(a - 1/2) - n/2}{\sqrt{n/4}} = -z_{(1-\alpha/2)} \text{ and } \frac{(b + 1/2 - 1) - n/2}{\sqrt{n/4}} = z_{(1-\alpha/2)}.$$

The empirical cumulative distribution function $\widehat{F}(x)$ and a confidence interval for $\widehat{F}(x)$ at a single point

The empirical cumulative distribution function (ecdf) $\widehat{F}(x)$ is defined by

$$\widehat{F}(x) = \text{the proportion of observations } \leq x .$$

Since the number of observations less than x is a binomial random variable with parameter $p = F(x)$, we can calculate the sample standard deviation of the ecdf as:

$$\widehat{SD}(\widehat{F}(x)) = \sqrt{\frac{\widehat{F}(x)(1 - \widehat{F}(x))}{n}},$$

and an approximate 100 $(1 - \alpha)\%$ confidence interval as:

$$\widehat{F}(x) \pm z_{(1-\alpha/2)} \sqrt{\frac{\widehat{F}(x)(1 - \widehat{F}(x))}{n}}.$$

Inference for other percentiles besides the median

The procedures considered for the median $\theta_{.5}$ are just a special case of procedures that can be used for any other percentile θ_p . In particular, hypothesis tests for θ_p can be performed by calculating p values using binomial probabilities $\binom{n}{x} p^x (1 - p)^{n-x}$. A confidence interval for θ_p also finds values a and b so that $P(X_{(a)} < \theta_p < X_{(b)}) = 1 - \alpha$ for some specified value α . The same reasoning as used above leads to finding a and b so that:

$$\sum_{x=a}^{b-1} \binom{n}{x} p^x (1 - p)^{n-x} \approx 1 - \alpha,$$

or by approximation in large samples with

$$\frac{a - np}{\sqrt{np(1-p)}} = -z_{(1-\alpha/2)} \text{ and } \frac{(b-1) - np}{\sqrt{np(1-p)}} = z_{(1-\alpha/2)}.$$

Note that the normal approximation above (or the continuity-corrected version) can be considered if n is large enough, for example, so that np and $n(1-p)$ are both greater than 5.

Comparison of tests