# Notes for video 8

Jonas Moss

23/3/2021

Let's take a look at the null-hypothesis

$$H_0 : \text{the data is normal with underlying mean } \mu \text{ and standard deviation } \sigma.$$

This is a simple null-hypothesis, usually denoted as $N(\mu, \sigma)$. We can easily simulate $n$ data points from this null-hypothesis in Python using `rng.normal(mu, sigma, n)`.

A famous test statistic for $H_0$ is the $Z$-test

$$Z = \sqrt{n}\frac{\overline{x} - \mu}{\sigma} \tag{1}$$

which you have learned is normally distributed with mean 0 and standard deviation 1 under the null-hypothesis $N(\mu, \sigma)$. Again, this is easy to calculate in Python. This test statistic makes most sense if the alterntive hypothesis has a different mean than $H_0$, something like

$$H_a : \text{the data does not have underlying mean } \mu, \text{ but has standard deviation } \sigma.$$

Recall the definition of a $p$-value, taken from Wikipedia, "the probability of obtaining test results at least as extreme as the results actually observed, under the assumption that the null hypothesis is correct."

In terms of probability, the definition of the observed two-sided $p$-value is

$$p = P(|Z| \geq |z|) = 1 - F(z), \tag{2}$$

where $Z$ is the random test statistic and $z$ is the observed value of the test statistic and $F(z) = P(|Z| \leq z)$ is the distribution function of $|Z|$. The variable $|Z|$ is well-known and called a folded normal, see Wikipedia for details. Since $\mu = 0$ and $\sigma = 1$ it also equals a standard half-normal, which we will use in programming. I repeat that this is the definition of the *observed* $p$-value. Since $z$ isn't random, $p$ is not random either.

We will need the distribution of the $p$-value under the null-hypothesis. As it's written now, $p$ is an observed quantity. To calculate its distribution, we must endow $z$ with a distribution. So let $Z \sim N(0, 1)$ and observe that the $p$-value, as a random variable, equals

$$p = 1 - F(|Z|).$$

(We would usually use capital letters for random variables, but $P$ is already taken, and we'll have to use $p$ even though it is random.) The distribution of this random variable equals

$$P(p \leq x) = P(1 - F(Z) \leq x) = P(1 - x \leq F(Z)).$$

Since $F$ is strictly increasing, its inverse $F^{-1}$ exists, and

$$P(1 - x \leq F(Z)) = P(F^{-1}(1 - x) \leq Z).$$
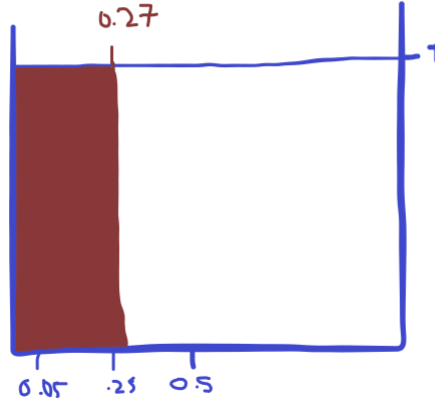
area of red: 0.27

0.27

1

0.05  .25  0.5

Figure 1: Example of a $p$-value

However, $P(z \le Z) = 1 - F(z)$, where $F(z)$ is the distribution of $Z$. Hence

$$P(F^{-1}(1 - x) \le Z) = 1 - F(F^{-1}(1 - x)),$$

and as $F^{-1}$ is the inverse of $F$, we get that $F(F^{-1}(1 - x)) = 1 - x$. Combining all of this we obtain

$$P(p \le x) = P(F^{-1}(1 - x) \le Z) = 1 - (1 - x) = x.$$

We recognize this as the distribution function of a uniform random variable.

**Proposition 1.** *p-values are uniformly distributed.*

This result has a curious implications for simulation, which we will explore.

**Example 2.** If the observed $p$-value is 0.27, the probability of observing a $p$-value less than or equal to 0.27 is equal to 0.27 under the nullhypothesis. This is the area of the red rectangle in the image of Figure 1.