Quantitative Genomics and Genetics BioCB 4830/6830; PBSB.5201.03

Lecture 11: Hypothesis Testing

Jason Mezey Feb 29, 2024 (Th) 8:40-9:55

Announcements

 Homework #3 (!!) will be available tomorrow (Fri, March 1) and will be due Fri, March 15 by 11:59PM

Summary of lecture 11: Intro to Hypothesis Testing

- Last lecture, we began completed our discussion of Estimators
- ...and briefly introduced Confidence Intervals (CI)
- Today, we are going to introduce Hypothesis Testing (!!)

Conceptual Overview



Review of essential concepts

- **Inference** the process of reaching a conclusion about the true probability distribution (from an assumed family of probability distributions indexed by parameters) on the basis of a sample
- System, Experiment, Experimental Trial, Sample Space, Sigma Algebra, Probability Measure, Random Variable, Probability Distribution (pmf, pdf), Parameterized Probability Model, Sample, Random Vector, Sampling Distribution, Statistic, Statistic Sampling Distribution, Estimator, Estimator Sampling distribution

Review: Estimators



Review: Introduction to maximum likelihood estimators (MLE)

- We will generally consider *maximum likelihood estimators* (MLE) in this course
- Now, MLE's are very confusing when initially encountered...
- However, the critical point to remember is that an MLE is just an estimator (a function on a sample!!),
- i.e. it takes a sample in, and produces a number as an output that is our estimate of the true parameter value
- These estimators also have sampling distributions just like any other statistic!
- The structure of this particular estimator / statistic is complicated but just keep this big picture in mind

Review: Confidence intervals I

- For the estimation framework we have considered thus far, our goal was to define an estimator that provides a "reasonable guess given the sample" of the true value of the parameter
- This is called "point" estimation since the true parameter has a single value (i.e. it is a point)
- We could also estimate an interval, where our goal is to say something about the chances that the true parameter (the point) would fall in the interval
- **confidence interval** (CI) an estimate of an interval defined such that if it were estimated individually for an infinite number of samples, a specific percentage of the estimated intervals would contain the true parameter value
- Don't worry if this concept seems confusing (it is!) let's first consider an example and then discuss some basics

Estimation and Hypothesis Testing

- Thus far we have been considering a "type" of inference, estimation, where we are interested in determining the actual value of a parameter
- We could ask another question, and consider whether the parameter is NOT a particular value
- This is another "type" of inference called hypothesis testing
- We will use hypothesis testing extensively in this course

Statistics



Estimators



Hypothesis Tests



Hypothesis testing I

- To build this framework, we need to start with a definition of hypothesis
- Hypothesis an assumption about a parameter
- More specifically, we are going to start our discussion with a null hypothesis, which states that a parameter takes a specific value, i.e. a constant

$$H_0: \theta = c$$

• For example, for our height experiment / identity random variable, we have $Pr(X|\theta) \sim N(\mu, \sigma^2)$ and we could consider the following null hypothesis:

$$H_0: \mu = 0$$

Hypothesis testing II

- As example, consider our height experiment (reals as sample space) / identity random variable X / normal probability model $\theta = [\mu, \sigma^2]$ / sample n=1 (of one height measurement) / identity statistic T(x) = x (takes the height measured height)
- Let's assume that $\sigma^2 = 1$ and say we are interested in testing the following null hypothesis $H_0: \mu = 5.5$ such that we have the following probability distribution of the statistic under the null hypothesis:



Hypothesis testing III

- Our goal in hypothesis testing is to use a sample to reach a conclusion about the null hypothesis
- To do this, just as in estimation, we will make use of a statistic (a function on the sample), where recall we know the sampling distribution (the probability distribution) of this statistic
- More specifically, we will consider the probability distribution of this statistic, assuming that the null hypothesis is true:

$$Pr(T(\mathbf{X} = \mathbf{x}|\theta = c))$$

- Note that this means we have a probability distribution of the statistic given the null hypothesis!!
- We will use this distribution to construct a *p*-value

p-value l

- We quantify our intuition as to whether we would have observed the value of our statistics given the null is true with a *p*-value
- **p-value** the probability of obtaining a value of a statistic T(**x**), or more extreme, conditional on H0 being true
- Formally, we can express this as follows:

$$pval = Pr(|T(\mathbf{x})| \ge t|H_0 : \theta = c)$$

 Note that a p-value is a function on a statistic (!!) that takes the value of a statistic as input and produces a p-value as output in the range [0, 1]:

$$pval(T(x)): T(x) \to [0,1]$$

p-value II

- As an intuitive example, let's consider a continuous sample space experiment / identify r.v. / normal family / n=1 sample / identity statistic, i.e. T(x) = x
- Assume we know $\sigma^2 = 1$ (is this realistic?), let's say we are interested in testing the null hypothesis $H_0: \mu = 0$ and let's say that we assume that if we are wrong the value of μ will be greater than zero (why?)





p-value III

• Same example: let's consider a continuous sample space experiment / identify r.v. / normal family / n=1 sample / identity statistic, i.e. T(X) = X / assume we know $\sigma^2 = 1$ / we test the null hypothesis $H_0: \mu = 0$ and let's assume that if we are wrong the value of μ could be in either direction (again, why?)



p-value IV

- More technically a p-value is determined not just by the probability of the statistic given the null hypothesis is true, but also whether we are considering a "one-sided" or "two-sided" test
- For a one-sided test (towards positive values), the p-value is:

$$pval(T(\mathbf{x})) = \int_{T(\mathbf{x})}^{\infty} Pr(T(\mathbf{x})|\theta = c)dT(\mathbf{x})$$

$$pval(T(\mathbf{x})) = \sum_{T(\mathbf{x})}^{max(T(\mathbf{X}))} Pr(T(\mathbf{x})|\theta = c)$$

• For a two-sided test, the p-value is:

$$pval(T(\mathbf{x})) = \int_{-\infty}^{-|T(\mathbf{x}) - median(T(\mathbf{X}))|} Pr(T(\mathbf{x})|\theta = c)dT(\mathbf{x}) + \int_{|T(\mathbf{x})| - median(T(\mathbf{X}))|}^{\infty} Pr(T(\mathbf{x})|\theta = c)dT(\mathbf{x})$$

$$pval(T(\mathbf{x})) = \sum_{min(T(\mathbf{X}))}^{-|T(\mathbf{x})-median(T(\mathbf{X}))|} Pr(T(\mathbf{x})|\theta = c) + \sum_{|T(\mathbf{x})-median(T(\mathbf{X}))|}^{max(T(\mathbf{X}))} Pr(T(\mathbf{x})|\theta = c)$$

Hypothesis Testing IV

- To build a framework to answer a question about a parameter, we need to start with a definition of hypothesis
- **Hypothesis** an assumption about a parameter
- More specifically, we are going to start our discussion with a *null hypothesis*, which states that a parameter takes a specific value, i.e. a constant

$$H_0: \theta = c$$

• Once we have assumed a null hypothesis, we know the probability distribution of the statistic, assuming the null hypothesis is true:

$$Pr(T(\mathbf{X} = \mathbf{x}|\theta = c))$$

• **p-value** - the probability of obtaining a value of a statistic $T(\mathbf{x})$, or more extreme, conditional on H0 being true:

$$pval = Pr(|T(\mathbf{x})| \ge t | H_0 : \theta = c)$$
$$pval(T(x)) : T(x) \to [0, 1]$$

• Note that a p-value is a function of a statistic (!!)

Non-Intuitive Hypothesis Testing Concepts I

- We do not know what the true model is (=parameter values are) in a real case!
- We assess a null hypothesis that we define!
- We assess this null hypothesis by calculating a p-value which assumes that the null hypothesis is true!
- We assess this null hypothesis by calculating a p-value from a single sample!
- We make one of two decisions: cannot reject or reject!
 - We decide on the value p-value that allows us to decide
 - If we reject, we interpret this as strong evidence against the null hypothesis being correct but we do not know for sure!
 - If we cannot reject, we cannot say anything (i.e., we have no evidence that the null is wrong and we cannot say that the null is right)!

Hypothesis decisions I

- We use the p-value to make a decision about the null hypothesis
- Specifically, we use the p-value for our sample to decide whether we "accept" (or better stated: "cannot reject") the null hypothesis or "reject" the null hypothesis
- To do this, we use a value α such that if the p-value is below this value we "reject", if it is above we "cannot reject"
- Note that this value of α corresponds to a critical value ("threshold") of the test statistic c_{α}
- For example for a value $\alpha = 0.05$ we have the following for our previous examples:



Hypothesis decisions II

- Note that there are two possible outcomes of a hypothesis test: we reject or we cannot reject
- We never know for sure whether we are right (!!)
- If we cannot reject, this does not mean H0 is true (why? What if our p-value is 0.99?)
- The value α is called the type I error, the probability of incorrectly rejecting H0 when it is true
- The value $1-\alpha$ is the probability of making a correct decision not to reject H0
- Note that we can control the level of type I error because we decide on the value of $\boldsymbol{\alpha}$

Assume H0 is correct (!): $\mu = 0$







two-sided test

Results of hypothesis decisions I: when H0 is correct (!!)

	H_0 is true
cannot reject H_0	1- α , (correct)
reject H_0	α , type I error



Results of hypothesis decisions I: when H0 is correct (!!)



Results of hypothesis decisions I: when H0 is correct (!!)



Assume H0 is wrong (!): $\mu = 3$







Sample 11: T(x) = 2.8



Ó

4

6

two-sided test

Results of hypothesis decisions II: when H0 is wrong (!!)

	H_0 is true	H_0 is false
cannot reject H_0	1- α , (correct)	β , type II error
reject H_0	α , type I error	$1 - \beta$, power (correct)



Results of hypothesis decisions II: when H0 is wrong (!!)



Results of hypothesis decisions II: when H0 is wrong (!!)



Technical definitions

 Technically, correct decision given H0 is true is (for one-sided, similar for two-sided):

$$1 - \alpha = \int_{-\infty}^{c_{\alpha}} Pr(T(\mathbf{x})|\theta = c) dT(\mathbf{x})$$

• Type I error (H0 is true) is (for one-sided):

$$\alpha = \int_{c_{\alpha}}^{\infty} Pr(T(\mathbf{x})|\theta = c) dT(\mathbf{x})$$

• Type II error given H0 is false is (for one-sided):

$$\beta = \int_{-\infty}^{c_{\alpha}} Pr(T(\mathbf{x})|\theta) dT(\mathbf{x})$$

• Power is (for one-sided):

$$1 - \beta = \int_{c_{\alpha}}^{\infty} Pr(T(\mathbf{x})|\theta) dT(\mathbf{x})$$

That's it for today

• Next lecture, we will complete our discussion of hypothesis testing AND begin our discussion of Genetic Models (!!)