Evaluating Significance of Predictors Hypothesis Testing

CS109A Introduction to Data Science Pavlos Protopapas, Kevin Rader and Chris Tanner



Suppose our model for advertising is:

y = 1.01x + 120

Where y is the sales in \$1000, x is the TV budget.

Interpretation: for every dollar invested in advertising gets you 1.01 back in sales, which is 1% net increase.

But how certain are we in our estimation of the coefficient 1.01?

Now you know how certain you are in your estimates, will you want to change your answer?



Now we know how to generate these distributions we are ready to answer **two** *important questions:*

A. Which predictors are most important?B. And which of them really affect the outcome?





CS109A, PROTOPAPAS, RADER, TANNER

The example below is from Boston housing data. This dataset contains information collected by the U.S Census Service concerning housing in the area of Boston. The coefficients below are from a model that predicts prices given house size, age, crime, pupil-teacher ratio, etc.



Feature importance based on the absolute value of the coefficients.



Feature importance based on the absolute mean value of the coefficients over multiple bootstraps and includes the uncertainty of the coefficients.





Feature Importance

To incorporate the coefficients' uncertainty, we need to determine whether the estimates of $\beta's$ are sufficiently far from zero.

To do so, we define a new metric, which we call t-test statistic:

$$t = \frac{\mu_{\widehat{\beta}_1}}{\sigma_{\widehat{\beta}_1}}$$

which measures the distance from zero in units of standard deviation.











Feature importance base on the absolute value of the coefficients.

Feature importance base on the absolute value of the coefficients over multiple bootstraps and includes the uncertainty of the coefficients.

Feature importance base on t-test. Notice the rank of the importance has changed.



Because a predictor is ranked as the most important, it does not necessarily mean that the outcome depends on that predictor.

How do we assess if there is a true relationship between outcome and predictors?

As with R-squared, we should compare its significance (t-test) to the equivalent measure from a dataset where we know that there is no relationship between predictors and outcome.

<u>We are sure</u> that there will be no such relationship in data that are randomly generated. Therefore, we want to compare the t-test of the predictors from our model with t-test values calculated using random data.



- 1. For *n* random datasets fit *n* models.
- 2. Generate distributions for all predictors and calculate the means and standard errors $(\mu_{\hat{\beta}}, \sigma_{\hat{\beta}})$.
- 3. Calculate the t-tests.

Repeat and create a probability density function (pdf) for all the t-tests.

It turns out we do not have to do this, because this is a known distribution called student-t distribution.



Student-t distribution, where ν is the degrees of freedom (number of data points minus number of predictors).



To learn more about why student-t, what are degrees of freedom and more details see https://en.wikipedia.org/wiki/Student%27s t-test

To compare the t-test values of the predictors from our model, $|t^*|$, with the t-tests, calculated using random data, $|t^R|$, we estimate the probability of observing $|t^R| \ge |t^*|$.

We call this probability the p-value.

p-value = $P(|t^R| \ge |t^*|)$

small p-value indicates that it is unlikely to observe such a substantial association between the predictor and the response due to chance.

It is common to use p-value<0.05 as the threshold for significance.

To calculate the p-value we use the cumulative distribution function (CDF) of the student-t.

stats model a python library has a build-in function stats.t.cdf()
which can be used to calculate this.







Feature importance based on the absolute value of the coefficients over multiple bootstraps and includes the coefficients' uncertainty. Feature importance based on t-test. Notice the rank of the importance has changed.



05 06

0.4

0.1

0.8 0.9

lstař

Pinatio

Dis

Nor

Black

Indus

Tax

Crim

Age

01 02 03

Predictors



Hypothesis Testing

Hypothesis testing is a formal process through which we evaluate the validity of a statistical hypothesis by considering evidence **for** or **against** the hypothesis gathered by **random sampling** of the data.

- 1. State the hypotheses, typically a **null hypothesis**, H_0 and an **alternative hypothesis**, H_1 , that is the negation of the former.
- 2. Choose a type of analysis, i.e. how to use sample data to evaluate the null hypothesis. Typically this involves choosing a single test statistic.
- **3.** Sample data and compute the test statistic.
- 4. Use the value of the test statistic to either **reject** or not reject the null hypothesis.



1. State Hypothesis:

Null hypothesis:

 H_0 : There is no relation between X and Y

The alternative:

 H_a : There is some relation between X and Y

2. Choose test statistics

t-test

3. Sample:

Using bootstrap we can estimate $\hat{\beta}'_1$ s, and $\mu_{\hat{\beta}_1}$ and $\sigma_{\hat{\beta}_1}$ and the t-test.



Hypothesis testing

4. Reject or not reject the hypothesis:

We compute *p***-value**, the probability of observing any value equal to |t| or larger, from random data.

p-value < p-value-threshold we reject the null.









What to do? 🧐

Today's lucky student: The student whose country of origin is the furthest from Boston.

Instructions:

Listen to your peers' opinions and suggestions. Ask questions of each other ("What do you think"). Do not just lead others in the room without including everyone.

Make sure you do not cut-off or ignore what other students are trying to contribute.

If you have questions, please reach out to the teaching staff. You can buzz us to come help, or if all else fails, come to the main room.

