

Stat 88: Prob. & Mathematical Statistics in Data Science



<https://xkcd.com/892/>

Lecture 21 : 4/7/2022

Section 9.2, 9.3, 9.4

A/B testing & confidence intervals

Hypotheses tests: Review of steps

1. State the *null* hypothesis (H_0) – that is, what is the assumption we are going to make. This will determine how we compute probabilities
2. State an *alternative* hypothesis (H_A) Note that this should not overlap with the null hypothesis, and it may or may not define probabilities (example: there was gender bias etc)
3. Decide on a test statistic to use that will help you decide which of the two hypotheses is supported by the evidence (data). Usually there is a natural choice. Use the null hypothesis to specify probabilities for the test statistic.
4. Find the observed value of the test statistic, and see if it is **consistent** with the null hypothesis. That is, compute the chance that we would see such an observed value, or more extreme values of the statistic (*p-value*).
5. State your conclusion: whether you reject the null hypothesis or not. This is based on your chosen cutoff (“level of the test”). Reject if the *p-value* is less than the cutoff.

Example: Woburn

In the early 1990s, a leukemia cluster was identified in the Massachusetts town of Woburn. Many more cases of leukemia, a malignant cancer that originates in a bone marrow cell, appeared in this small town than would be predicted. Was this evidence of a problem in the town or just chance?

Observed significance levels (a.k.a p -values)

- The p -value decides if observed values are *consistent* with the null hypothesis. It is a *tail* probability (also called *observed significance level*), and is the chance, **assuming that the null hypothesis is true**, of getting a test statistic equal to the one that was observed or even more in the direction of the alternative.
- If this probability is too small, then something is wrong, perhaps with your assumption (null hypothesis). That is, the data are unlikely if the null is true and therefore, your data are ***inconsistent*** with the null hypothesis.
- p -value is **not** the chance of null being true. The null is either true or not.
- The p -value is a *conditional probability* since it is computed *assuming* that the null hypothesis is true.
- The smaller the p -value, the stronger the evidence ***against*** the null and ***towards*** the alternative (in the direction of the alternative).
- Traditionally, below 5% ("result is statistically significant") and 1% ("result is highly significant") are what have been used. Significant means the p -value is small, not that the result is important.

Ex. 9.5.1

- All the patients at a doctor's office come in annually for a check-up when they are not ill. The temperatures of the patients at these check-ups are independent and identically distributed with unknown mean μ .
- The temperatures recorded in 100 check-ups have an average of 98.2 degrees and an SD of 1.5 degrees. Do these data support the hypothesis that the unknown mean μ is 98.6 degrees, commonly known as "normal" body temperature? Or do they indicate that μ is less than 98.6 degrees?

A/B testing: comparing *two* distributions

- Data 8, section 12.3, randomized controlled trial to see if botulinum toxin could help manage chronic pain.
- 31 patients → 15 in treatment group, 16 in control group. 2 patients in the control group reported pain relief and 9 in the treatment group.
- A/B testing is a term used to describe hypothesis tests which involve comparing the distributions of *two* random samples. (Earlier we had *one* sample and made a hypothesis about its distribution.)
- In particular, we can conduct an A/B test for hypothesis tests involving results of randomized controlled trials, A is the control group and B the treatment group.

Fisher's exact test

- Control group: 16 patients, 2 reported relief
- Treatment group: 15 patients, 9 reported relief
- H_0 : The treatment has no effect (there would have been 11 patients reporting pain relief no matter what, and it just so happens that 9 of them were in the treatment group)
- H_A : The treatment has an effect

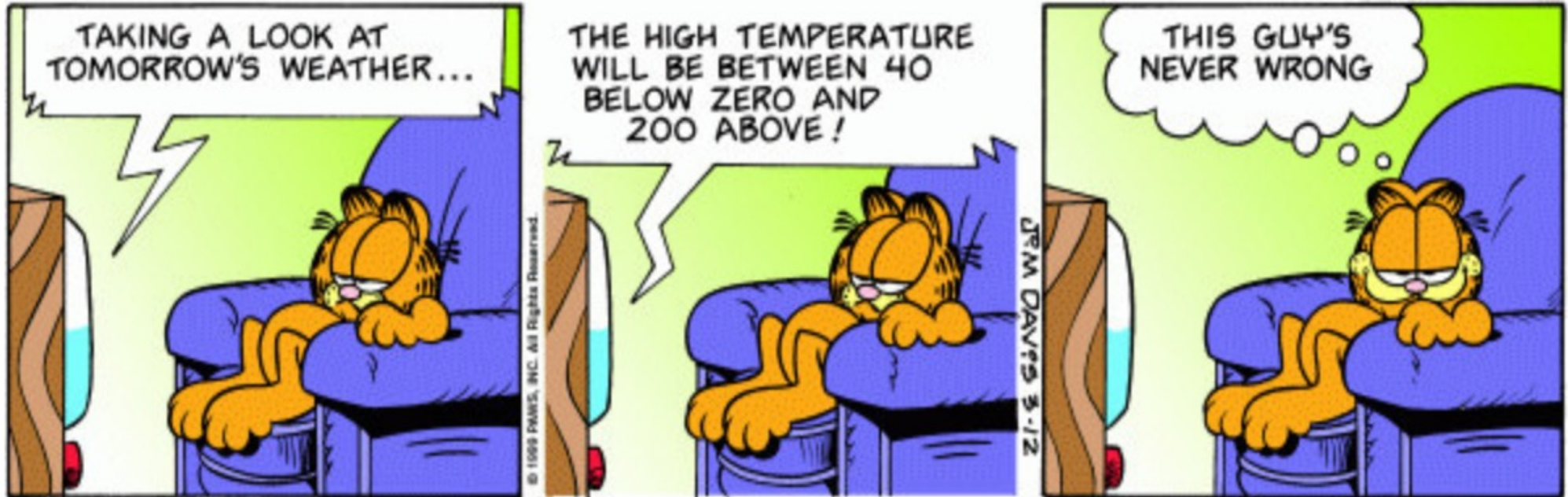
Example: The Lady Tasting Tea

- The first person to describe this sort of hypothesis test was the famous British statistician Ronald Fisher. In his book *The Design of Experiments*, he describes a tea party in which a lady of his acquaintance claimed that she could tell from tasting a cup of tea if the milk had been poured first or the tea.
- Fisher immediately set up an experiment in which she was given multiple cups of tea and asked to identify which of them had had the tea poured first. She tasted 8 cups of tea, of which 4 had the tea poured first, and identified 3 of them correctly. Does this data support her claim?

Example: Gender bias ?

- Rosen and Jerdee conducted several experiments using male bank supervisors (this was in 1974) who were given a personnel file and asked to decide whether to promote or hold the file. 24 were randomly assigned to a file labeled as that of a male employee and 24 to a female.
- 21 of the 24 males were promoted, and 14 of the females. Is there evidence of gender bias?

Section 9.3: Confidence intervals



Goal: Estimating a parameter

- Say we have a population whose average, μ , we want to estimate
- How would we do it? We could draw one data point X_1 and use it to estimate μ . Do you think this is a good method of estimation? If not, why not?
- What about if we draw a sample of size 2: X_1, X_2 where each of the X_i have expectation μ ? Is this better? Can we use the average of these two?
- We generally use a larger sample, say n is a large number and we draw an iid sample X_1, X_2, \dots, X_n . Why is this a better idea? The expectation of each of the X_i is μ , so the expectation of the sample mean is also μ . But this was true even for $n = 2$. Why use larger n ?

Using \bar{X} to estimate μ

- \bar{X} is an unbiased estimator of μ (what does that mean?)
- If we also know that each of the X_k had SD σ , what can we say about $SD(\bar{X})$?
- What does the Central Limit theorem say about the sample mean?
- We will use the CLT and the sample mean to define a random interval (why is it random?) that will cover the true mean with a specified probability, say 95%

Confidence intervals

- In the previous slide, we derived an ***approximate 95% Confidence Interval for the population mean μ***
- Why is the interval random?
- A *confidence interval* is an interval on the real line, that is, a collection of values, that are plausible estimates for the true mean μ .
- Using the CLT, we can estimate the chance that this interval contains the true mean. If we want the chance to be higher, we make the interval bigger. The interval is like a net. We are trying to catch the true mean in our net.
- The CLT takes the form: $\bar{X} \pm \text{margin of error}$, where the margin of error tells us how big our interval is, and depends on the SD of the sample mean.
- The margin of error = $z_{\alpha/2} \times SD(\bar{X})$, where $z_{\alpha/2}$ is the quantile we need to have an area of $1 - \alpha$ in the middle, that is, a ***coverage probability*** of $1 - \alpha$

Example

- A population distribution is known to have an SD of 20. The average of an iid sample of 64 observations is 55. What is your 95% confidence interval for the population mean?

Confidence levels

- The probability with which our *random* interval will cover the mean is called the confidence level.
- In reality (vs theory), we will have just one *realization* (observed value) of the sample mean (from our data sample), and we use that value to write down the **realization** of our random interval.
- What would we do differently if we wanted a 68% CI? 99.7% CI?
- What about an 80% CI? 99% CI?

Confidence intervals for the population mean: recap

- A *confidence interval* is an interval on the real line, that is, a collection of values, that are plausible estimates for the true mean μ .
- Using the CLT, we can estimate the chance that this interval contains the true mean. If we want the chance to be higher, we make the interval bigger. The interval is like a net. We are trying to catch the true mean in our net.
- The CLT takes the form: $\bar{X} \pm \text{margin of error}$, where the margin of error tells us how big our interval is, and depends on the SD of the sample mean.
- The margin of error = $z_{\alpha/2} \times SD(\bar{X})$, where $z_{\alpha/2}$ is the quantile we need to have an area of $1 - \alpha$ in the middle, that is, a **coverage probability** of $1 - \alpha$
- The probability with which our *random* interval will cover the mean is called the confidence level.
- In reality (vs theory), we will have just one *realization* (observed value) of the sample mean (from our data sample), and we use that value to write down the **realization** of our random interval.

Dealing with proportions

- A sample proportion is just the sample mean of a special population of 0's and 1's.
- This kind of population is so common since many of our problems deal with *classifying* and *counting*.
- We have a population of 1 million in a town. We take a SRS of size 400 and find that 22% of the sample is unemployed. Estimate the percentage of unemployed people in the town.

Section 9.4: Interpretation

- Chance that sample mean is less than 2 SDs away from population mean is about 0.95
- Therefore the chance that population mean is less than 2 SDs away from sample mean is about 0.95
- Which object is random in each of these sentences?
- Does it make sense to say "The probability that the number 2 is between 3 and 5 is 0.95" ?
- Does it make sense to say "The probability that the population mean is between 18 and 26 is 0.95"?

Interpretation

- Let's think about tossing coins. *Before* we toss a coin some number of times, we can say that the number of heads is random, since we *don't know* how many heads we will get.
- Suppose we have tossed the coin (say 100 times) and we see 53 heads, can we say 53 is a random number and the chance that 53 lies between 40 and 50 is 95%?
- 53 is our **realization** of the random "number of heads" in this *particular* instance of 100 tosses.

Confidence intervals: What is random?

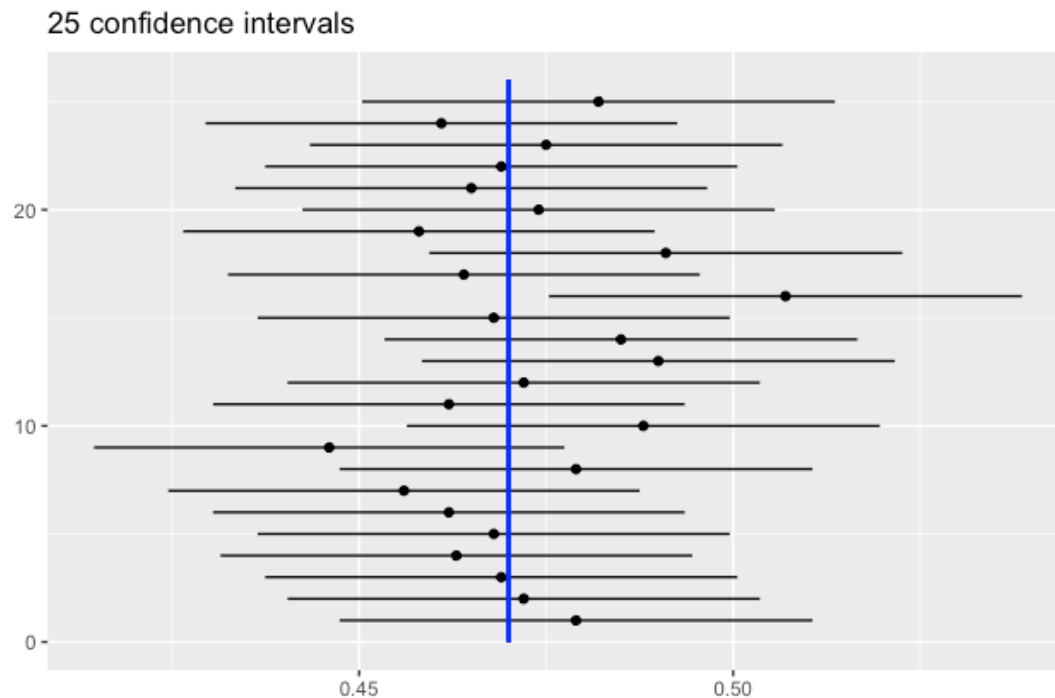
- Note that if we use the sample mean and extend one or two SDs in either direction, we *may* or *may not* cover the true population percentage.
- The *interval* is random, since we use a realization of the random variable (\bar{X}) to compute it.
- What fraction of such intervals (each interval computed from a random sample of data) will cover the true value μ ?
- This *coverage probability* (**before we actually collect the data**) is called the ***confidence level*** of the confidence interval.

Confidence Intervals

1. Which would be wider : a 99% CI or a 95% CI?
2. What about a 90% CI? 68%?
3. The _____ the confidence level, the _____ the interval
4. This does not make sense! Why are we using a normal distribution when the sample consists of Bernoulli random variables?
5. What is the chance that the population %, p , is in the interval (18%, 26%)?

Probability of coverage

- We draw 25 samples (sample size 100) from a Bernoulli distribution with $p=0.47$.
- Construct a 95% CI from each sample.
- How many intervals covered the blue line? How many did you expect?
- What is the *chance* that each CI will cover the true p (before you plug in #s)?
- If X =number of successful intervals, what is the distribution of X ?
- Why are the centers different? Are the widths the same?



Margin of error

- We have a confidence interval. Now we want to keep the **same confidence level**, but want to improve our accuracy. For example, say our *margin of error* is 4 percentage points, and we want it to be 1 percentage point. What should we do?
 - A. increase width of CI 4 times by increasing SD
 - B. Decrease width of CI by increasing n by 4 times
 - C. Decrease width of CI by increasing n by 16 times

Comparison with bootstrap CI

- How do you create a bootstrap CI for the population mean?

