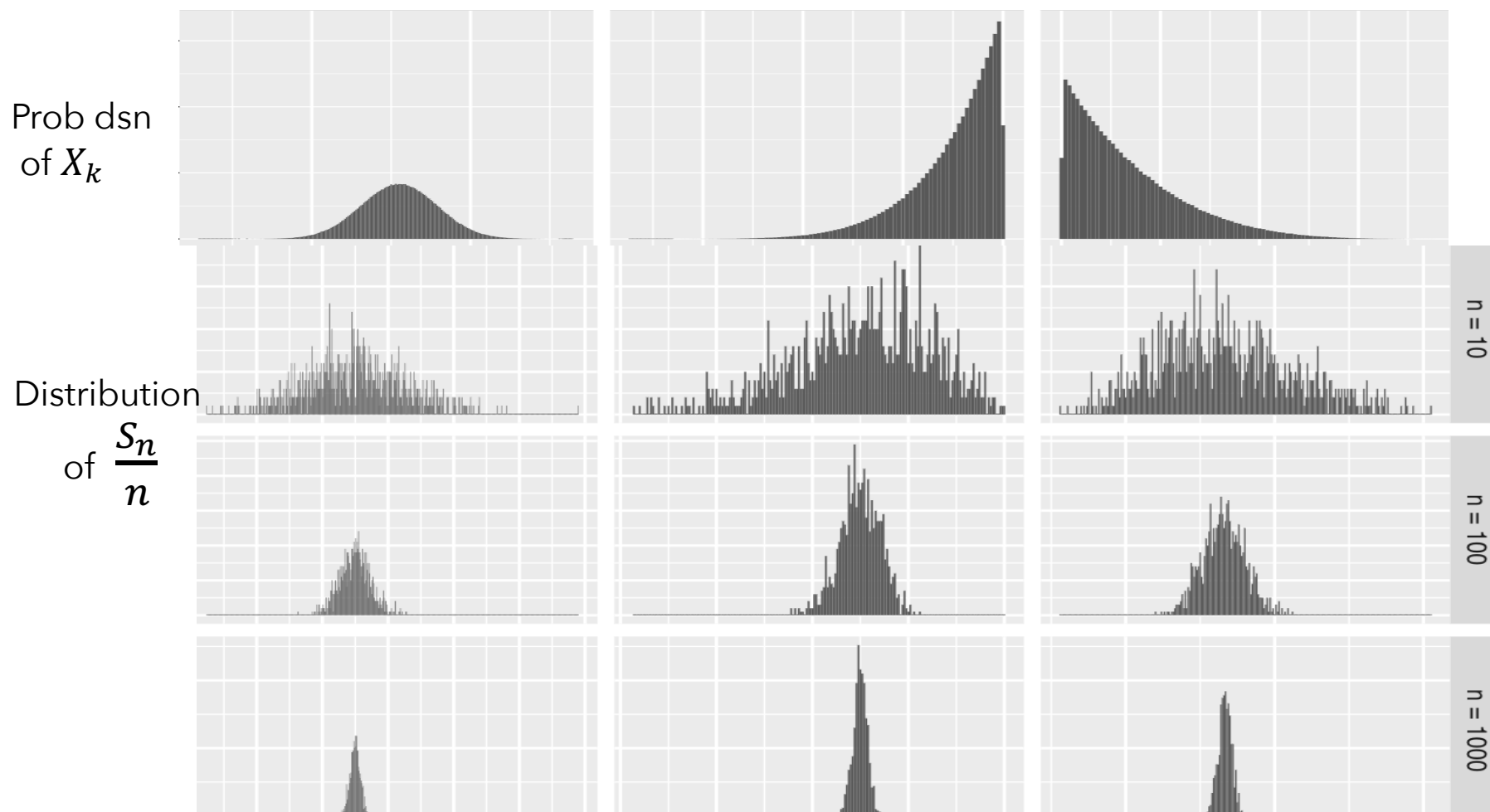


Stat 88: Prob. & Mathematical Statistics in Data Science



Lecture 30 PART 1: 4/8/2024

Finishing up chapter 8 and the Central Limit Theorem

Standard normal cdf, symmetries, percentiles

If you think of the standard normal curve as a probability histogram, then it is natural to think of areas under the curve as probabilities. In particular, the function Φ defined by:

(Go to Jupyter)

(a) Verify empirical rule

(b) Find the probability that Z is between -0.3 and 0.9

(c) Find the probability that Z is outside -1.5 and 1.5 .

(d) Find z such that $\Phi(z) = 0.95$

(e) Find z so that the area in the **middle** is 0.95 ($\Phi(z) - \Phi(-z) = 0.95$)

Example

- Suppose that each of 300 patients has a one in three chance of being helped by a treatment, independently of its effect on other patients. Find the probability that at least half the patients are helped by the treatment.

The Central Limit Theorem

- Suppose that X_1, X_2, \dots, X_n are iid with mean μ and SD σ
- Let $S_n = X_1 + \dots + X_n$ be the sample sum, and $A_n = \frac{S_n}{n}$ be the sample mean
- Then the distribution of S_n (and A_n) is *approximately normal* for large enough n .
- For S_n , the distribution is approximately normal (bell-shaped), centered at $\mathbf{E}(S_n) = \mathbf{n}\mu$ and with spread given by $\mathbf{SD}(S_n) = \sqrt{\mathbf{n}} \sigma$.
- For A_n , the distribution is centered at $\mathbf{E}(A_n) = \mu$ with spread $\mathbf{SD}(A_n) = \sigma/\sqrt{\mathbf{n}}$

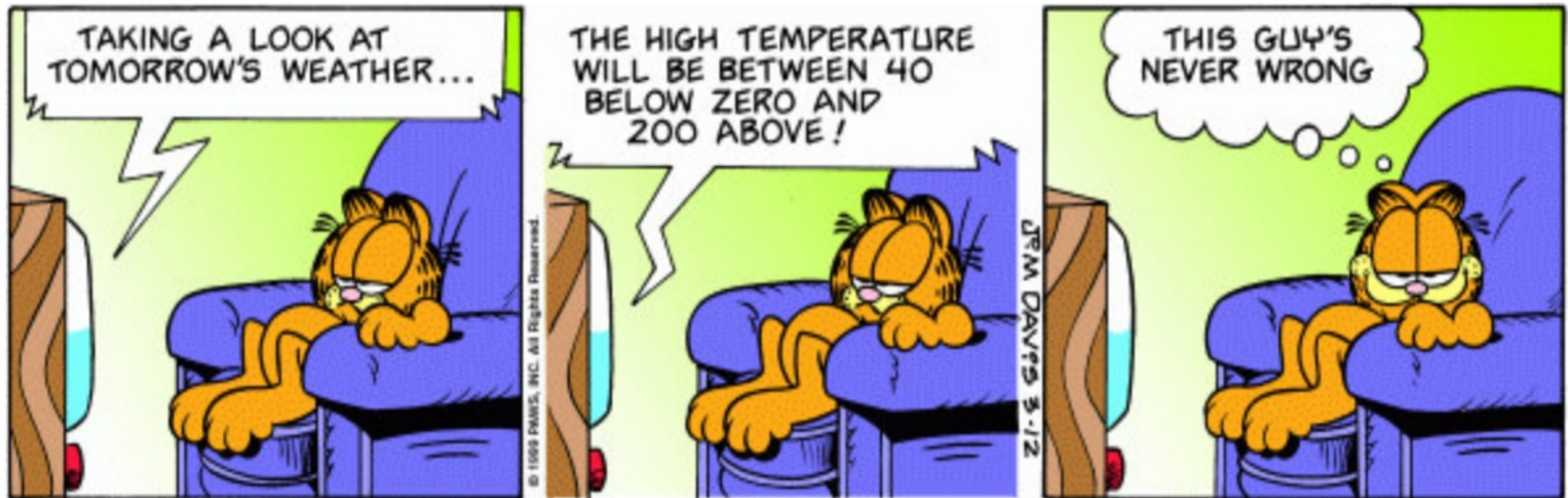
Normal approximations : standard units

- Let X be any random variable, with expectation μ and SD σ , consider a new random variable that is a linear function of X , created by shifting X to be centered at 0, and dividing by the SD. If we call this new rv X^* , then X^* has expectation _____ and SD _____.
- $E(X^*) =$
- $SD(X^*) =$
- This new rv does not have units since it measures how far above or below the average a value is, in SD's. Now we can compare things that we may not have been able to compare.
- Because we can convert anything to standard units, ***every normal curve is the same.***

How to decide if a distribution could be normal

- Need enough SDs on both sides of the mean.
- In 2005, the murder rates (per 100,000 residents) for 50 states and D.C. had a mean of 8.7 and an SD of 10.7. Do these data follow a normal curve?
- If you have indicators, then you are approximating binomial probabilities. In this case, if n is very large, but p is small, so that np is close to 0, then you can't have many sds on the left of the mean. So need to increase n , stretching out the distribution and then the normal curve begins to appear.
- If you are not dealing with indicators, then might bootstrap the distribution of the sample mean and see if it looks approximately normal.

Stat 88: Prob. & Mathematical Statistics in Data Science



Lecture 22 PART 2 : 4/8/2024

Section 9.1, 9.2

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%
- Based on data from a *random sample*, we will construct an interval of estimates for some unknown (but fixed) population parameter.

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?

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.

Example

- In a simple random sample of 400 voters in a state, 23% are undecided about which way they will vote. Find a 95% CI for the proportion of undecided voters in the state.
- In the above problem, find 99.7% confidence interval.

Interpretating a Confidence Interval

- 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?

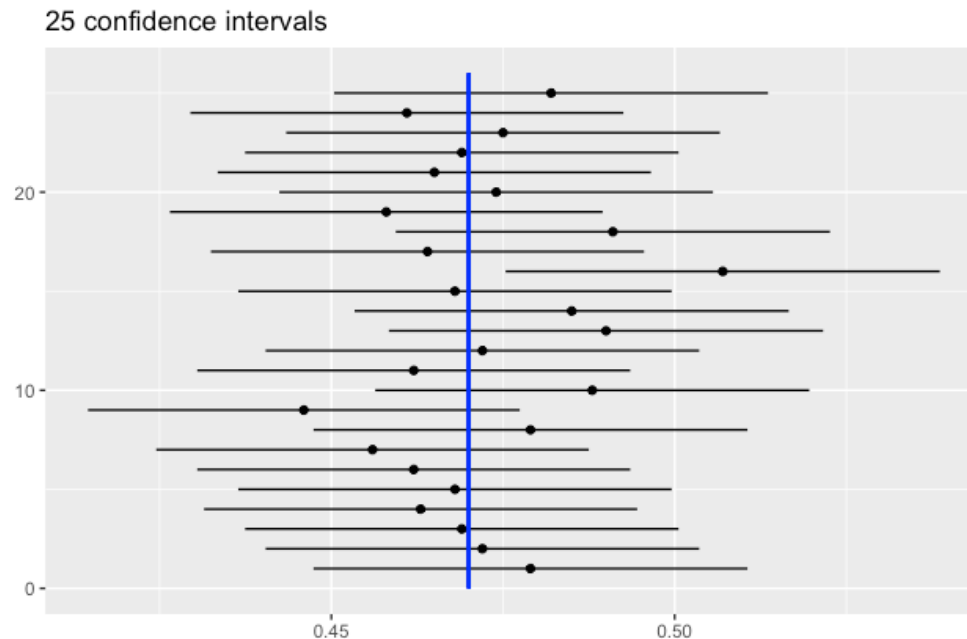
- Note that if we use the sample mean and extend one or two SDs in either direction, we *may* or *may not* cover the true population percentage.
- The *interval* is random, since we use a realization of the random variable (\bar{X}) to compute it.
- What fraction of such intervals (each interval computed 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?

