USN					



InternalAssessmentTest1 Sep 2025

Sub:	Statistical Ma	ichine Lear	rning for Data	a Science (SM	L)	SubCod e:	BAD702	Brance:	h AIN	DS/CS	S(DS)
Date:	29/09/25	Duratio n:	90minute s	MaxMarks :	50	Sem	V	II	•	0	BE
				vFIVEOue ons				N	IARK S	СО	RBT
1	A startup has collected data on the daily time (in minutes) that 50 users spend on their new mobile learning app. The company wants to estimate the average daily usage in the population.					rs	10	CO1	L3		
		d, the da	•	large and the		• 0		ı			
	Questions:										
		-		e bootstrap n n daily usage?		od to estin	nate a 95%				
	•	hy might bootstrap be preferred here over using a standard for the confidence interval?									
		minutes	and the 97.3	esamples, the 5th percentile on mean?		-		ıld			
	1) How to mean	apply th	ne bootstr	ap to estim	ate	a 95% C	I for the				
	Procedure (step-by-s	tep):								
	Call 2. Repe 3. After {x-1: 4. The 1 97.5th CI _{95%}	it x. cat the foll Draw a Compu B resamp *,,x B* percentile th percent = (per	owing B tind bootstrap so the the sample oles you have bootstrap iles of that because the centile 2.5	ample of daily nes (common ample of size le mean of the re a bootstrap 95% CI is singular pootstrap-mean (\bar{x}^*) , percent potstrap standa	cho n w e boo dist mply ns d	ices: 1,000 ith replace otstrap same ribution of x the empiristribution: x x x y	— 10,000) ement from aple; store it the mean:	х.			

 $SE_{boot}=\mathrm{sd}(\bar{x}^*)$. You can also build a normal-approx CI using , $\bar{x}\pm z_{0.975}\cdot SE_{boot}$,but for skewed data the percentile or BCa intervals are preferred.

Quick checklist / assumptions: sample observations should be roughly independent and the sample should be representative of the population you want to infer about.

2) Why bootstrap is preferred here

- **Skewed underlying distribution:** Classical formula for the mean uses the Central Limit Theorem (CLT) and normal-based CIs. For small-to-moderate n and **skewed** data the CLT approximation can be poor; bootstrap captures the actual shape of the sampling distribution empirically.
- Small-ish sample size (n = 50): 50 is borderline; with heavy skew, normal approximations may be biased. Bootstrap does not require normality.
- No need to assume parametric form: If you're unsure about the population distribution (which you are), bootstrap is nonparametric and more robust.
- **Flexible to statistics beyond the mean:** Bootstrap easily extends to medians, percentiles, differences, etc.
- **Practical caveats:** bootstrap relies on the sample being representative and having independent observations. If those are violated (e.g., strong time dependence), bootstrap results will be unreliable.

If you need higher accuracy in skewed cases, consider the **BCa** (biascorrected and accelerated) bootstrap interval rather than plain percentile.

3) Interpretation of the given bootstrap percentiles

We performed 1,000 bootstrap resamples. The 2.5th and 97.5th percentiles of the bootstrap means are **32** and **48** minutes, respectively.

Conclusion:

A 95% bootstrap percentile confidence interval for the population mean is [32, 48] minutes. Practically, this means that based on the observed sample and the bootstrap procedure, the best empirical estimate is that the true **average daily usage** in the population lies between **32 and 48 minutes** with *approximate* 95% confidence.

Short interpretation we can report:

"We estimate the population mean daily usage to be between 32 and 48 minutes (95% CI, bootstrap percentile)."

Extra notes / caveats you should state:

• This interval is conditional on the sample being representative and

 observations independent. With only B = 1,000 resamples the interval is usually fine, but for more precise endpoints you can increase to B = 5,000–10,000 (computationally cheap). If the bootstrap distribution is strongly skewed, consider the BCa interval which corrects bias and skewness. 			
Bonus — Jupyter / Python code import numpy as np import pandas as pd def bootstrap_mean_ci(data, B=10000, alpha=0.05, seed=None): rng = np.random.default_rng(seed) n = len(data) means = np.empty(B) for b in range(B): sample = rng.choice(data, size=n, replace=True) means[b] = sample.mean() lower = np.percentile(means, 100*(alpha/2)) upper = np.percentile(means, 100*(1-alpha/2)) se_boot = means.std(ddof=1) return { "mean_observed": np.mean(data), "se_boot": se_boot, "ci_percentile": (lower, upper), "bootstrap_means": means }			
 What do you mean by Bootstrap samples? How does it help in building a Confidence Interval? Suppose you have a dataset of size n: \$x = {x_1, x_2,, x_n}\$ A bootstrap sample is a new sample of size n drawn with replacement from the original dataset. "With replacement" means the same observation can appear multiple times in the bootstrap sample, or not at all. For example, if your dataset is [10,20,30,40], a bootstrap sample could be [20,20,40,10] or [30,30,10,40]. We can generate many such bootstrap samples (say B=1000, 5000, or more). For each bootstrap sample, we compute the statistic of interest (mean, median, variance, regression coefficient, etc.). 	5+5	CO1	L3
This gives a bootstrap distribution of that statistic. How does it help in building a Confidence Interval?			

The key idea:

- In classical statistics, confidence intervals often rely on theoretical formulas (like assuming normality of errors, using t- or z-distributions, etc.).
- But what if the data are skewed, sample size is small, or you don't want to assume a specific distribution?

Bootstrap helps because:

- 1. It **mimics repeated sampling** from the population by resampling from the observed data.
- 2. The variability in the bootstrap distribution reflects the sampling variability of the statistic.
- 3. To build a **95% confidence interval**, we take the middle 95% of the bootstrap distribution (usually the 2.5th and 97.5th percentiles).
- 2(b) A pharmaceutical trial compares a new drug with a placebo. Suppose you reject the null hypothesis and claim the drug works, but later it turns out the drug has no effect.
 - What type of statistical error have you made?
 - How could you reduce the chances of this error in future experiments?

Question 1: What type of statistical error is this?

- The null hypothesis H0: "The drug has no effect."
- You rejected H0 and concluded the drug works.
- Later, it turned out the drug really has no effect → your rejection of H0 was wrong.

☐ This is a **Type I Error** (false positive).

Definition: Type I Error occurs when you reject a true null hypothesis.

Question 2: How to reduce the chances of this error?

The probability of making a Type I Error is the **significance level** (α) you choose for your test.

• Common choices: $\alpha = 0.05$ (5%) or $\alpha = 0.01$ (1%).

To reduce Type I error:

- 1. Lower the significance level (α) :
 - o For example, test at 1% instead of 5%.
 - o Makes it harder to reject H0.
- 2. **Use corrections for multiple testing** (e.g., Bonferroni correction) if many comparisons are being made.
- 3. Improve experimental design:

	0	Increase sample size for more reliable estimates.		
	0	Reduce bias and noise (randomization, blinding, proper		
		controls).		
4.	Replic	cation: Repeating the experiment independently reduces the		
	chance	e that one false positive drives the conclusion.		

3	Explain the practical relevance of Binomial and Poisson Distribution. Give the formula for their Pdf.	10	CO2	L2
	What is output of the following Python code:			
	from scipy import stats print(stats.binom.pmf(2, n=5, p=0.5))			
	1. Practical Relevance			
	Binomial Distribution			
	 Used when there are a fixed number of independent trials with two possible outcomes (Success/Failure). Examples: Tossing a coin 10 times and counting the number of heads. Quality control: number of defective items in a batch of 50. Medical trials: number of patients who respond positively to a new drug out of a fixed sample. 			
	Poisson Distribution			
	 Used when we count the number of events happening in a fixed interval of time/space, assuming events occur independently at a constant average rate. Examples: Number of customer arrivals at a bank per hour. Number of phone calls received at a call center per minute. Number of emergency cases arriving at a hospital per day. 			
	2. Probability Density Function (PMF)			
	• Binomial PMF: $P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}, k=0,1,2,\ldots,n$ where $\bullet \text{n = number of trials,}$ $\bullet \text{p = probability of success}$			
	where			

		Poisson PMF	?:				
		$P(\mathcal{I}$	$(X=k)=rac{\lambda^k e^{-\lambda}}{k!}, k=0,1$.,2,			
		where					
			verage rate of occurrence, umber of events.				
		3. Python	Code Output				
		from scipy	import stats				
		print(stats.	binom.pmf(2, n=5, p=0.5))				
		*	tes the probability of getting exprobability of success p=0.5.	eactly 2 successes in 5			
		P(X = k	$k = rac{\lambda^k e^{-\lambda}}{k!}, k = 0, 1, 2, \dots$	$\cdot (0.125) = 0.3125$			
4	a	Convergence	Always the same shape A	approaches the normal distribution as df $ ightarrow \infty$	5	CO3	L1
		Difference distribution	ce between Student's t-distri on	bution and Normal			
		Feature	Normal Distribution	Student's t-Dis			
		Shape	Symmetrical, bell-shaped	Symmetrical, bell-shaped			
		Parameters	Mean (μ), Standard deviation (σ)	Degrees of freedom ($df = n$			
		Tails	Thinner tails	Heavier tails (more probabi			
		Usage	When population σ is known, or sample size is very large	When σ is unknown and we sample (small/medium sam			
		Convergence	e Always the same shape	Approaches the normal dist			
		What is Sta	andard Error (SE)?				
		samp	lition: The standard error is the st ling distribution of a statistic (like ne sample mean:				

$SE = \frac{s}{\sqrt{n}}$			
 where: s = sample standard deviation, n = sample size. Interpretation: SE measures how much the sample mean is expected to vary from one sample to another. Key role: Smaller SE → the sample mean is a more precise estimate of the population mean. 			
b State the Central Limit Theorem? How is it relevant to Statistical Inference.	5	CO3	L
Central Limit Theorem (CLT)			
Statement: If we take many random samples of size n from any population with mean μ and finite variance σ2, then as nnn becomes large, the sampling distribution of the sample mean X ⁻ approaches a Normal distribution, regardless of the population's original distribution.			
Formally:			
$rac{ar{X} - \mu}{\sigma/\sqrt{n}} \;\; \stackrel{d}{\longrightarrow} \;\; N(0,1) ext{as } n o \infty$			
Relevance to Statistical Inference			
The CLT is the backbone of modern statistics because it justifies why we can use normal-based methods (like z-tests, t-tests, confidence intervals) even when the underlying population is not normal :			
 Approximate Normality of Sample Mean: Even if data is skewed or irregular, the mean of a reasonably large sample (usually n≥30) will be approximately normal. 			
 Confidence Intervals: CLT allows us to construct confidence intervals for population means using the normal (or t) distribution. 			
3. Hypothesis Testing: Test statistics (like the standardized sample mean) rely on the CLT to follow approximately normal distributions under H0.			
4. Practical Applications:Polling: estimating population proportions.			

5	a	State the difference between Violin plot and Box plot with diagram. When do we use Hexagonal Binning.				CO3	L1
		Feature	Box Plot	Violin P			
		Shows	Summary statistics (median, quartiles, whiskers, outliers)	Summary statistics and ful			
		Shape	Rectangular box with whiskers	Symmetric, violin-shaped			
		Data distribution	Hides the shape of the distribution (only gives 5-number summary)	Displays the kernel densit showing whether the data i multimodal, etc.			
		Use case	Good for comparing medians and spread across groups	Good for seeing detailed di especially multimodal data			
		Diagram (c	conceptual):				
		Box Plot:					
		whisker be	ox whisker				
							
		L	I				
		Media	an				
		Violin Plot:					
		Density shap	e				
		()					
		()					
		()					
		()					
		()					
		()					
		median sho	wn inside				
		So violin plot	ds = box plot + distribution shape.				
		When do w	ve use Hexagonal Binning?				

, ,			T		
and ove • Inst hex	skbin plots are used when you have two cont a large dataset , where scatter plots would s rplotting (points overlapping, making it have ead of plotting each point, the data space is a gonal bins , and the color of each hexagon equency) of points in that bin.	uffer from d to see density). divided into			
Use cases:					
• Exa	ualizing correlations in large datasets. Imple: plotting 100,000 values of height vs. Insity-based patterns become clear compared it.				
b What do you Population p	n mean by Test Statistic? Differentiate between Sparameter.	Sample Parameter and	5	CO3	L1
What do	you mean by a Test Statistic?				
is u • It n pop • The	est statistic is a numerical value calculated for sed in hypothesis testing. The sample statistic is from solulation parameter, relative to the variation in the test statistic is then compared to a theoretical formal, t, F, Chi-square, etc.) to decide whether othesis.	the hypothesized the data. al distribution			
Examples:					
• z-te	est statistic:				
	$z=rac{ar{x}-}{\sigma/\sqrt{s}}$ t-test statistic: $t=rac{ar{x}-}{s/\sqrt{s}}$ Chi-square test statistic (goodness- $\chi^2=\sumrac{(O)}{s}$	$\frac{\mu_0}{\sqrt{n}}$ of-fit):			
Difference b	petween Sample Parameter and Population Par	rameter			
Aspect	-	Sample			
Definition	A numerical summary that describes the entire population	A numerical summa sample (subset of p			
Examples	Population mean (μ), population variance (σ^2), population proportion (P)	Sample mean (x ⁻ \ba variance (s ²), sample			
Known / Unknown	Usually unknown (we can't measure the whole population)	e Known (we comput collected)			
Role	Fixed value (does not change)	Varies from sample estimate the popula			
Use in Inference	True value we want to infer about	Provides the basis for testing about the popular			

a Explain the intuition behind Hypothesis Testing, Level of Significance and p-value	5	CO3	1.3
Intuition behind Hypothesis Testing	<i>J</i>		נב
 Imagine you're a detective. You start with a presumption of innocence (null hypothesis, H₀). You collect evidence (sample data). You ask: "Is this evidence strong enough to reject H₀ and believe the alternative hypothesis, H₁?" 			
So, hypothesis testing is a structured decision-making process:			
 Assume H₀ is true (status quo). Calculate a test statistic from your sample. Compare it to a reference distribution (normal, t, chi-square, etc.). If the evidence is too unlikely under H₀, you reject H₀ in favor of H₁. Level of Significance (α)			
 Denoted by α, it is the threshold for risk you are willing to take of making a Type I Error (rejecting a true H₀). Common values: α = 0.05 → 5% chance of wrongly rejecting H₀. α = 0.01 → 1% chance. 			
Example: If $\alpha = 0.05$, you're saying "I'm okay with being wrong 5 times out of 100 in rejecting H_0 ."			
p-value			
 The p-value is the probability of observing a test statistic as extreme as (or more extreme than) the one you got, assuming H₀ is true. It answers: "If the null hypothesis were true, how surprising is my sample?" 			
Interpretation:			
 Small p-value (≤ α) → evidence is strong against H₀ → reject H₀. Large p-value (> α) → not enough evidence → fail to reject H₀. 			
Example:			
 p = 0.03, α = 0.05 → reject H₀ (evidence suggests effect exists). p = 0.40, α = 0.05 → fail to reject H₀ (data consistent with no effect). 			
b A sports scientist is tracking the performance of professional runners. In the first race of the season, one athlete unexpectedly runs much faster than their usual average time (an unusually good performance).	5	CO3	L3

entist predicts that in the next race, the athlete's performance will be and closer to their long-term average. Ons: What statistical concept explains why the athlete's performance is likely to decline toward their usual average in the next race? Does this mean the athlete is "getting worse"? Explain why or why not.	
What statistical concept explains why the athlete's performance is likely to decline toward their usual average in the next race? Does this mean the athlete is "getting worse"? Explain why or why	
likely to decline toward their usual average in the next race? Does this mean the athlete is "getting worse"? Explain why or why	
Give another real-world example (outside sports) where regression to the mean commonly occurs.	
t statistical concept explains this?	
ncept is Regression to the Mean.	
cans that if a random variable shows an extreme value (very high or w) in one measurement, the next measurement is likely to be closer verage, simply because of natural variability.	
case, the athlete's unusually fast race is partly due to chance factors weather, perfect mindset, competition, etc.). These lucky factors ll align again, so the next performance is expected to move closer to rage.	
s this mean the athlete is "getting worse"?	
t doesn't.	
lete is not suddenly performing worse. The unusually good nance was an outlier influenced by random variation. Their "true" nance ability is reflected by their long-term average.	
n performance drops back toward the average, it's not a decline in it's just the natural balancing out of random fluctuations.	
ther real-world example of Regression to the Mean	
s who score extremely high or extremely low on a test often score to the class average on the next test. This doesn't mean smart students ember" or struggling students got "smarter" — it's just that random (luck, question fit, mood, etc.) don't repeat in the same way.	
	ans that if a random variable shows an extreme value (very high or w) in one measurement, the next measurement is likely to be closer terage, simply because of natural variability. Tase, the athlete's unusually fast race is partly due to chance factors weather, perfect mindset, competition, etc.). These lucky factors ll align again, so the next performance is expected to move closer to rage. This mean the athlete is "getting worse"? It doesn't. There is not suddenly performing worse. The unusually good hance was an outlier influenced by random variation. Their "true" hance ability is reflected by their long-term average. In performance drops back toward the average, it's not a decline in it's just the natural balancing out of random fluctuations. There real-world example of Regression to the Mean exation: The world example of Regression to the Mean exation: The world example of Regression to the Mean exation: The real-world example of Regression to the Mean exation: The real-world example of Regression to the Mean exation: The real-world example of Regression to the Mean exation: The real-world example of Regression to the Mean exation: The real-world example of Regression to the Mean exation: The real-world example of Regression to the Mean exation: The real-world example of Regression to the Mean exation: The real-world example of Regression to the Mean exation: The real-world example of Regression to the Mean exation: The real-world example of Regression to the Mean exation: