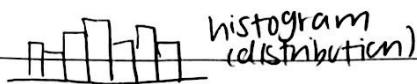
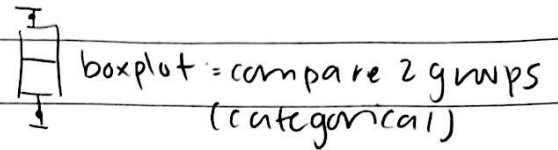
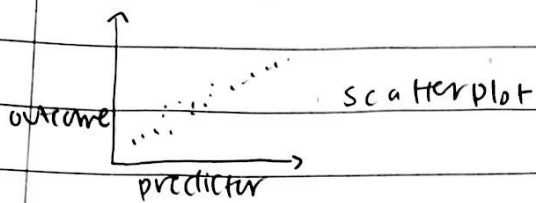


Week 3 Lecture Notes

Exploratory Data Analysis (EDA)



Basics of Probability Distribution

Bernoulli = single trial w/ 2 outcomes $P(\text{success})$ $P(\text{failure})$ $1 - P(\text{success})$

Binomial = multiple Bernoulli samples strung together

Geometric = # of failures before success when strung together Bernoulli trials

Normal Dist = μ σ $1 - \sigma$ $1 + \sigma$

Poisson = (# of events over a time period (exp))

Overview of Hypothesis Testing

1. come up w/ research question
2. Turn into hypothesis
3. collect samples / data
4. EDA
5. decide signif. level
6. Run your stats test
7. examine output make sure seems reasonable
8. Report your key stats / results

Statistical Tests (pt 1) - z + t-tests

population σ ?

yes no

$n > 30$ $n < 30$

z-test t-test

$H_0 \rightarrow$ means not statistically different

$H_A \rightarrow$ yes difference

• one-sample z / t test = you are comparing a sample's mean to a predetermined mean (hypothesized mean)

• paired = comparing means of 2 samples from the same pop.

• 2 sample = comparing means of 2 samples \rightarrow 2 indep. samples from 2 diff pop

Same indiv before
after treatment

indep = linked

t-test (assumption)

one-sample = 1) continuous

2) data follow a normal dist

3) Random sampling

paired = 1) continuous vars.

2) difference b/w

the matched sample pairs
is normally dist

3) Random sampling.

2 sample =

1) continuous data.

2) normally distributed

3) samples r independent

4) variances of 2 sample are the same

5) random samp

z test assumptions

1 sample = 1) continuous data

2) data \approx normal dist

3) sample is taken randomly

4) pop SD known

paired z = 1) continuous data

2) diff b/w the pairs of data points
follows a normal dist.

3) random sampling

4) pop SD of the difference is known.

2 sample 1) con. var.

2) random sample in both
pops

3) independence of samples.

4) variance of dis of
both pop. known

one tail = know if μ is either $<$, $>$
sample mean

2-tail = difference (not directionality)

typical outputs.

1) t-or z-stat 2) df (n-1) 3) p-value,
confidence interval

ANNOVA TEST

t-test, z-test \rightarrow 2 samples

ANNOVA \rightarrow 3 or more groups

One way @ two way

\rightarrow 1 variable differs
 \rightarrow 2 categorical factors differ and
how they effect each other

one way assumptions

1) normality (sample from approx normal)

2) indep. samples.

3) variance across group same.

4) outcome is continuous

H_0 = no difference, means equal

H_1 = means differ

two way assumptions.

1) outcome variable continuous/independent.

2) groups/vars should be
categorical + independent

3) sample are independent.

4) variance across groups the same

5) normality

* difference = # of groups.

p-value & confidence interval

significance level (α) = 0.05
p-value = prob that your result is
due to random chance.

$p < 50\%$ (reject null)

$p \geq 50\%$ (fail reject null)

that the pop mean is w/in
"95% of the CI will cover the true
pop mean"

$CI = (1 - \alpha) \rightarrow \alpha = 0.05 \rightarrow 95\% CI$
provides a range of means = (1.02, 3.5)

z or t with df (n-1) \rightarrow critical value
interpret test stat