

Lab 1 Group 3, Legislator age

[Code ▾](#)

Team03: Savita Chari, Tymon Silva, Denny Lehman

Study question 2: Wrong Test, Right Data

Our Analysis:

The response categories in Likert scales creates rank order/ordinal data, but the intervals between values cannot be presumed equal. Therefore, standard statistical analysis such as the mean, standard deviation etc. are inappropriate for ordinal data. The data generated by this survey is not continuous so its distribution cannot be measured, So A-B is not a normal distribution. It also violates the assumption that each pair (A1,B1) is an i.i.d. Hence, this data is not appropriate for a paired t-test.

Study question 3: Legislator age

You would like to test whether Democratic or Republican senators are older. List all assumptions for a Wilcoxon rank-sum test and evaluate them.

Hypothesis

H_0 : Median age of Democrat is same as median age of Republican senators.

Our null hypothesis is that there is no difference in the mean age between all the Democrat legislators and the Republican legislators. The expectation of X equals the expectation of Y . In other words, $\Delta = 0$.

What is Wilcoxon rank-sum test?

It is a way of examining the relationship between a numeric outcome variable (Y) and a categorical explanatory variable (X, with 2 levels) when the groups are Independent. It is a nonparametric test for unpaired data. It is used to compare one random variable X against another random variable Y .

The Wilcoxon test for comparing two population means makes the following assumptions:

1. The two samples are independent of one another and are Metric scale.
2. The two populations have equal variance or spread
3. The two populations are normally distributed (iid)

The following rules must be observed

1. The first assumption must be satisfied for a two-sample t-test.
2. But when assumptions #2 and #3 (equal variance and normality) are not satisfied, the samples size must be larger than 30 for the results are approximately correct
3. But when our samples are small and our data skew or not normal, we probably should not perform two-sample t-test

[Code](#)

Data Cleaning

The Data set has 538 rows and 34 columns. But we need a Subset of this data to answer the research question

```
[1] "The Dataset had 538 Rows and 34 Columns"
```

In order to perform Wilcoxon test we need 2 types of columns from our dataset

1. $[X]$ A categorical column with 2 distinct groups : We chose the 'party' column for this requirement but the data set includes more than two categories. We filtered only the Republican and Democrat data in order to perform Wilcoxon rank-sum test as we are interested in only those categories. Thus we fulfill the requirement of having only 2 categories.
2. $[Y]$ A column with numeric outcome: We choose the 'birthday' field for this requirement. Though this field is of class 'Date', some manipulation will be needed to make it a numeric outcome

Data Transformation

The birthday field is of 'Date' class and is in YYYY/MM/DD format In order to find out the age of the legislator we calculate the current age of the legislator by doing a datediff between current date and the birthday

Hide

```
currentDate <- Sys.Date()
legislator_DemRep_Age_data <- legislator_subset_data %>%
  mutate(AgeInYears = round((as.numeric(difftime(currentDate , birthday,units = "weeks")))/52))
```

Exploring Assumption 1 : Are the data sets independent?

Well, yes, a legislator cannot belong to 2 parties simultaneously, so the data is independent but still we can very easily check the count of rows for each party in the data sets

Code

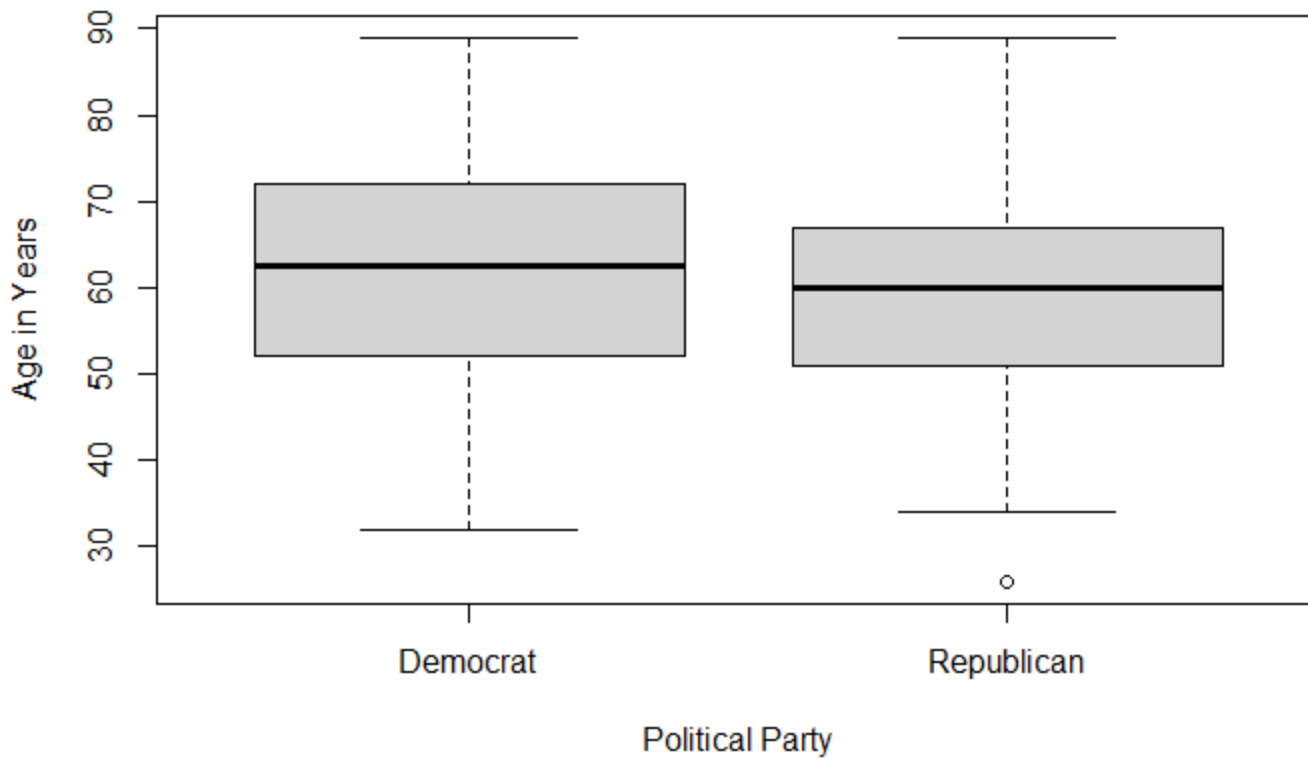
party	n
<chr>	<int>
Democrat	272
Republican	264
2 rows	

Hence, Assumption 1 is True, the sample-size for Democrats and Republicans is not same. Also, we have 272 rows for Democrats and 264 rows for Republicans. This sample size is larger than the minimum required 30 rows in order to avoid data skews and to make this a valid Wilcoxon rank-sum test.

Exploring Assumption 2 : Does the population have equal variance?

With our large data set we do not need to fulfill this assumption but from the box plot below, it is very clear that the Democrates have a bigger variance than the Republicans wrt. their age. In the Republican data set there is an outlier too.

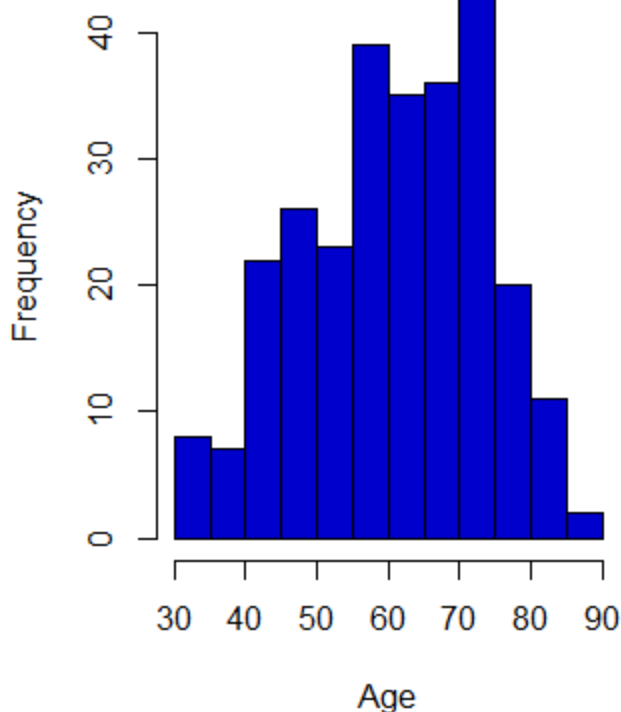
Variance in the Age



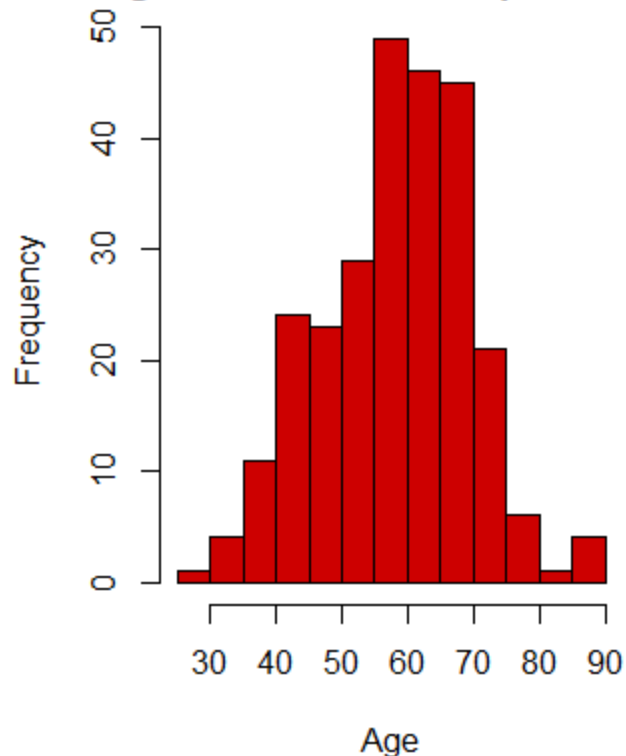
Exploring Assumption 3 : Is the populations normally distributed(iid) ?

Our sample size is larger than 30 so we do not have to fulfill this assumption, but from the graph below we can tell that the age distribution between republicans and democrats is not normal, but not too skewed too.

Age distribution of Democrats



Age distribution of Republicans



Performing two sided Wilcoxon rank-sum test

In order to prove the null hypothesis we perform the Wilcoxon rank-sum test. our data satisfies the requirements for conducting this test.

```
Wilcoxon rank sum test with continuity correction
```

```
data: legislator_DemRep_Age_data$AgeInYears by legislator_DemRep_Age_data$party
W = 40787, p-value = 0.006443
alternative hypothesis: true location shift is not equal to 0
95 percent confidence interval:
 0.9999557 5.0000343
sample estimates:
difference in location
                 3
```

Test Outcome and it's correctness

1. We used mean to create our test statistics which is well accepted distributions for statistical testing.
2. we followed the decision rules of a hypothesis testing which gives us the guarantee that our false positive rejection rate (the type 1 error rate) is bounded.
3. Our test is significant as our p-value is greater than 0.0025

With 95% confidence interval, We are rejecting the null hypothesis