# Histograms and Distributions

# So far…

**Scatterplots - pairs of attributes**

**How does a single attribute behave?**

**Example**

**Life expectancy in GapMinder data set**

**Example**

**Life expectancy in GapMinder data set**

**Varies between countries**

**1967 Average Life Expectancy in Morocco, Bangladesh = 50.34, 43.45**

**Example**

**Life expectancy in GapMinder data set**

**Varies between countries**

**1967 Average Life Expectancy in Morocco, Bangladesh = 50.34, 43.45**

**Varies over time**

**Average life expectancy in Namibia in 1952, 2007 = 41.73, 52.90**

**Aside - life expectancy for country itself an estimate**

**Varies according to**

**Year of birth**

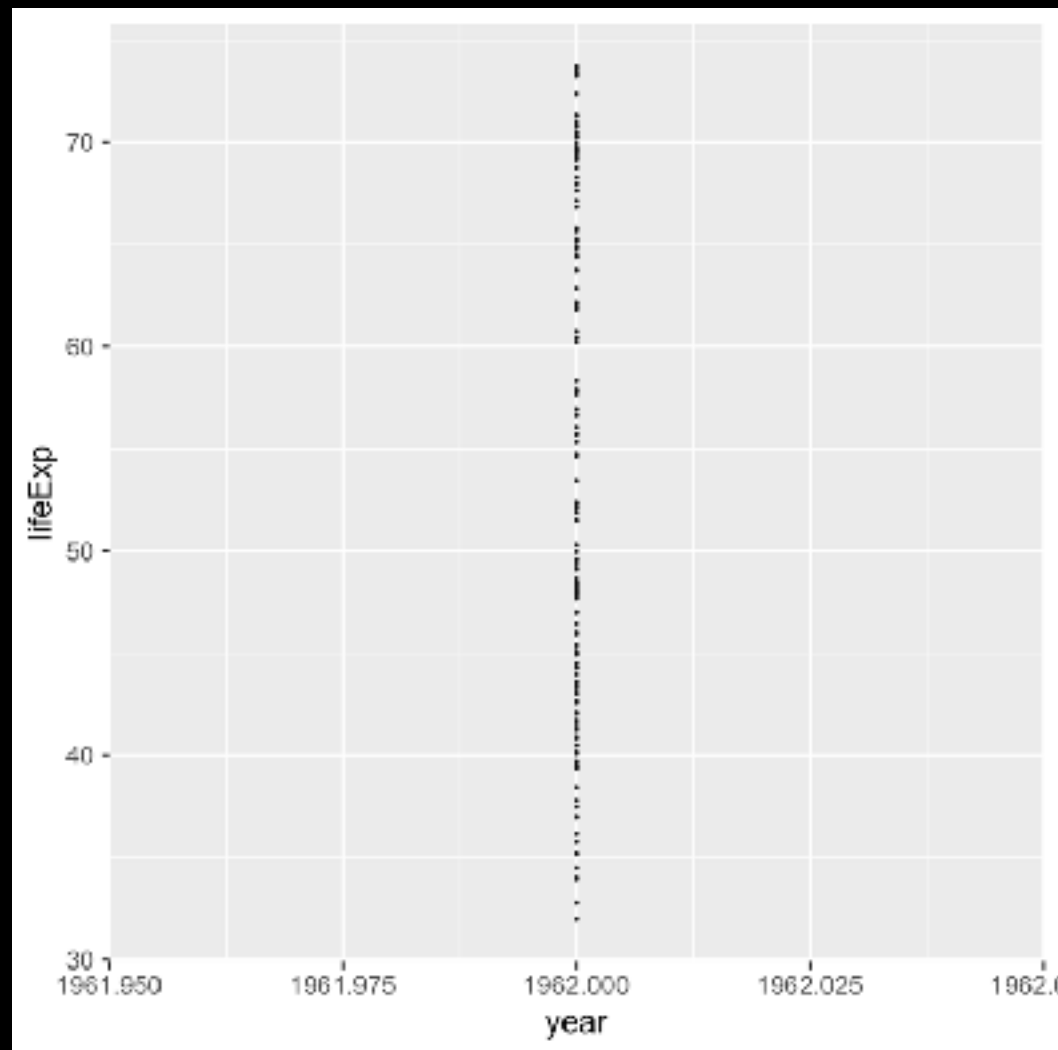**Socioeconomic group (wealthy/poor)**

**Geographical location**

**First attempt - just plot points**

**Example - Life Expectancy in 1962 for all countries**

**First attempt - just plot points**
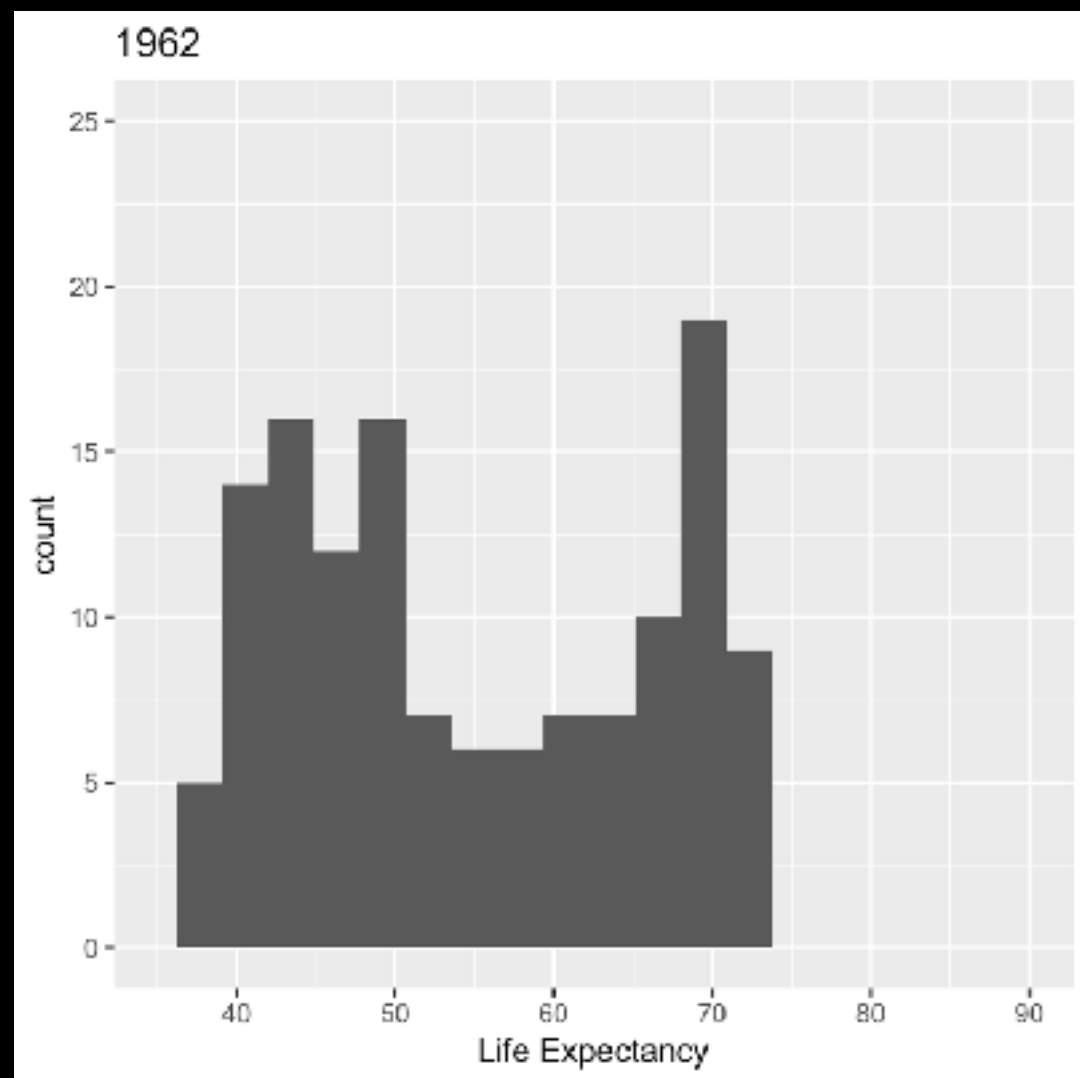
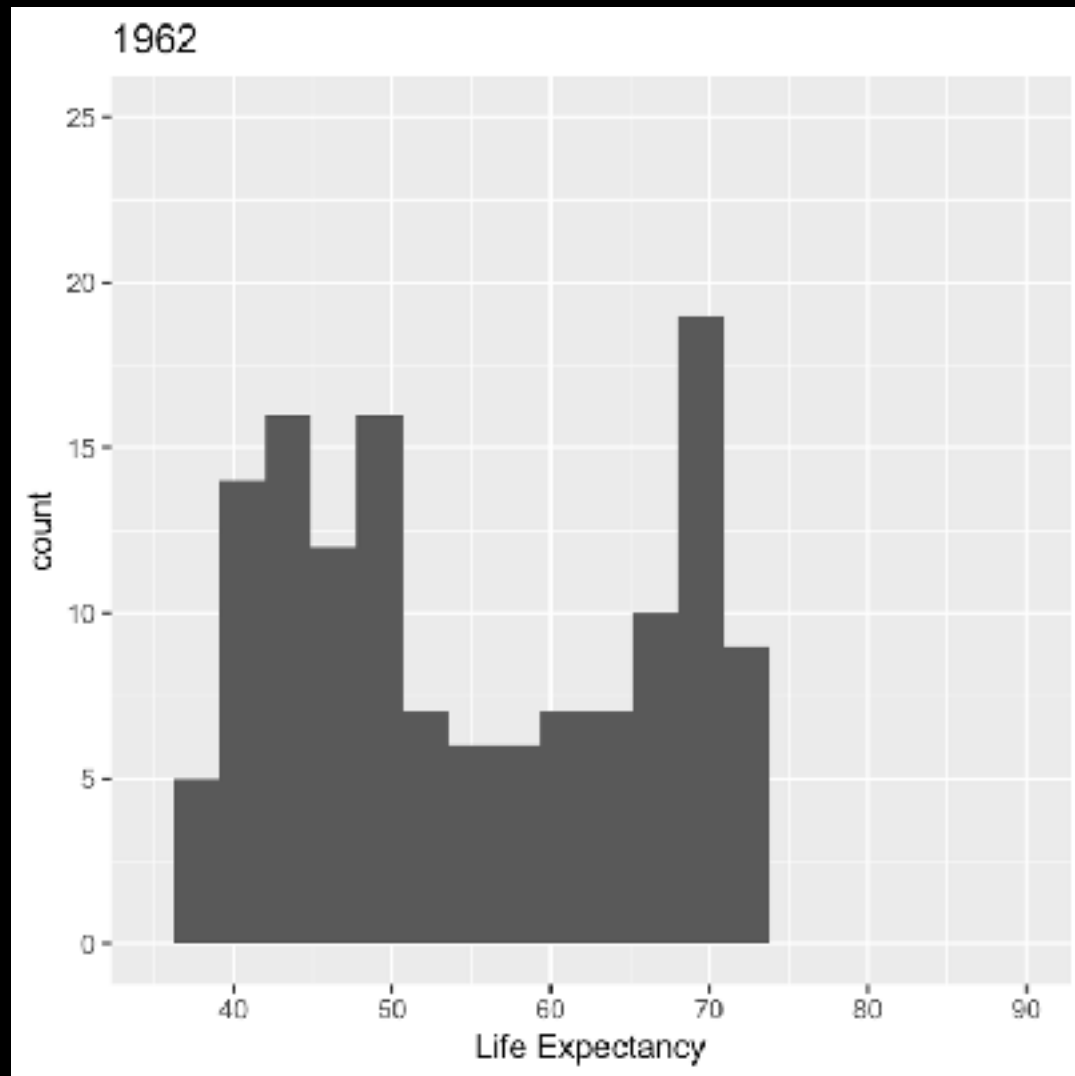**Example - Life Expectancy in 1962 for all countries**

# Histogram

Means of converting Ordinal to Categorical data

Create 'bins' -
how many countries have an Life Expectancy of 40-45 etc.

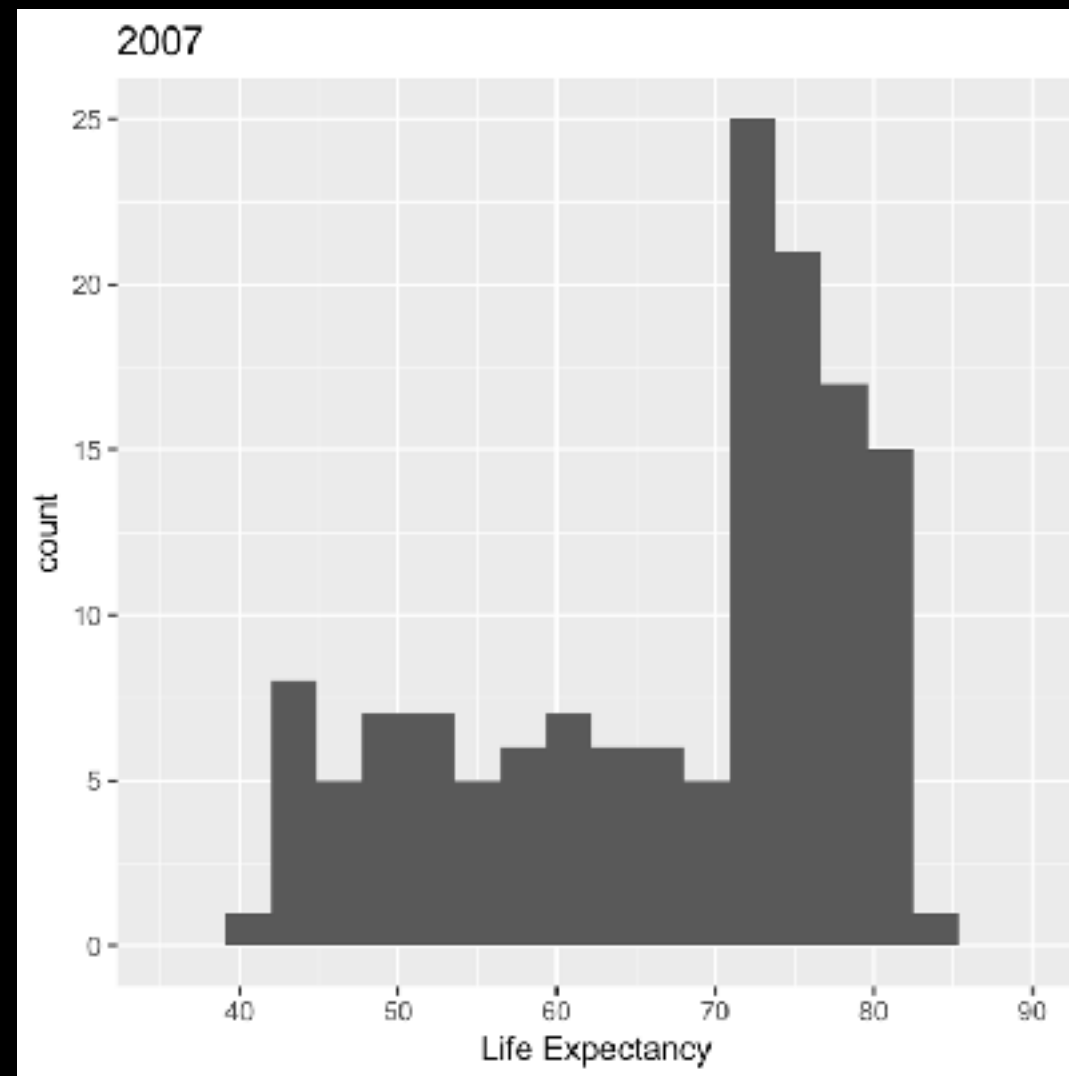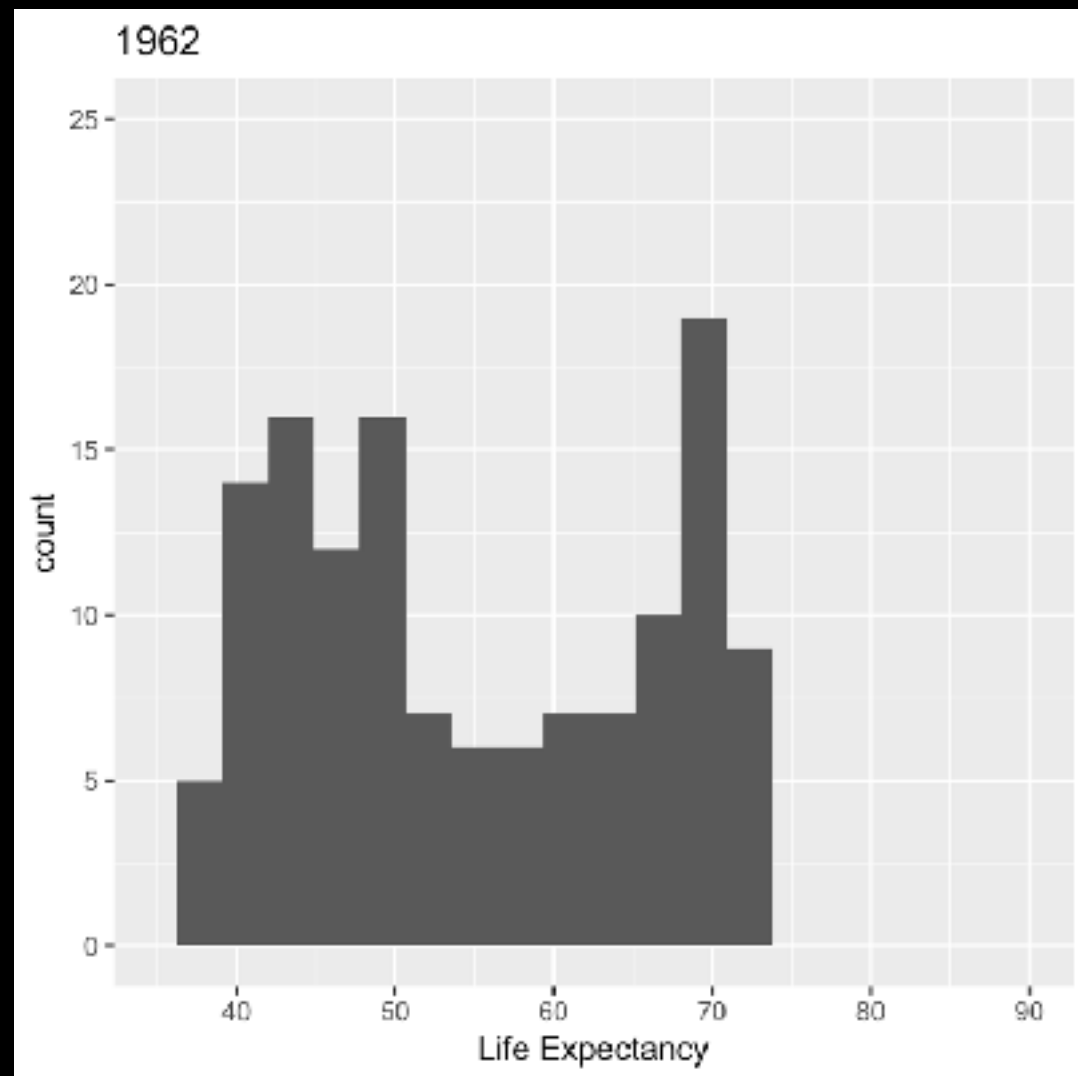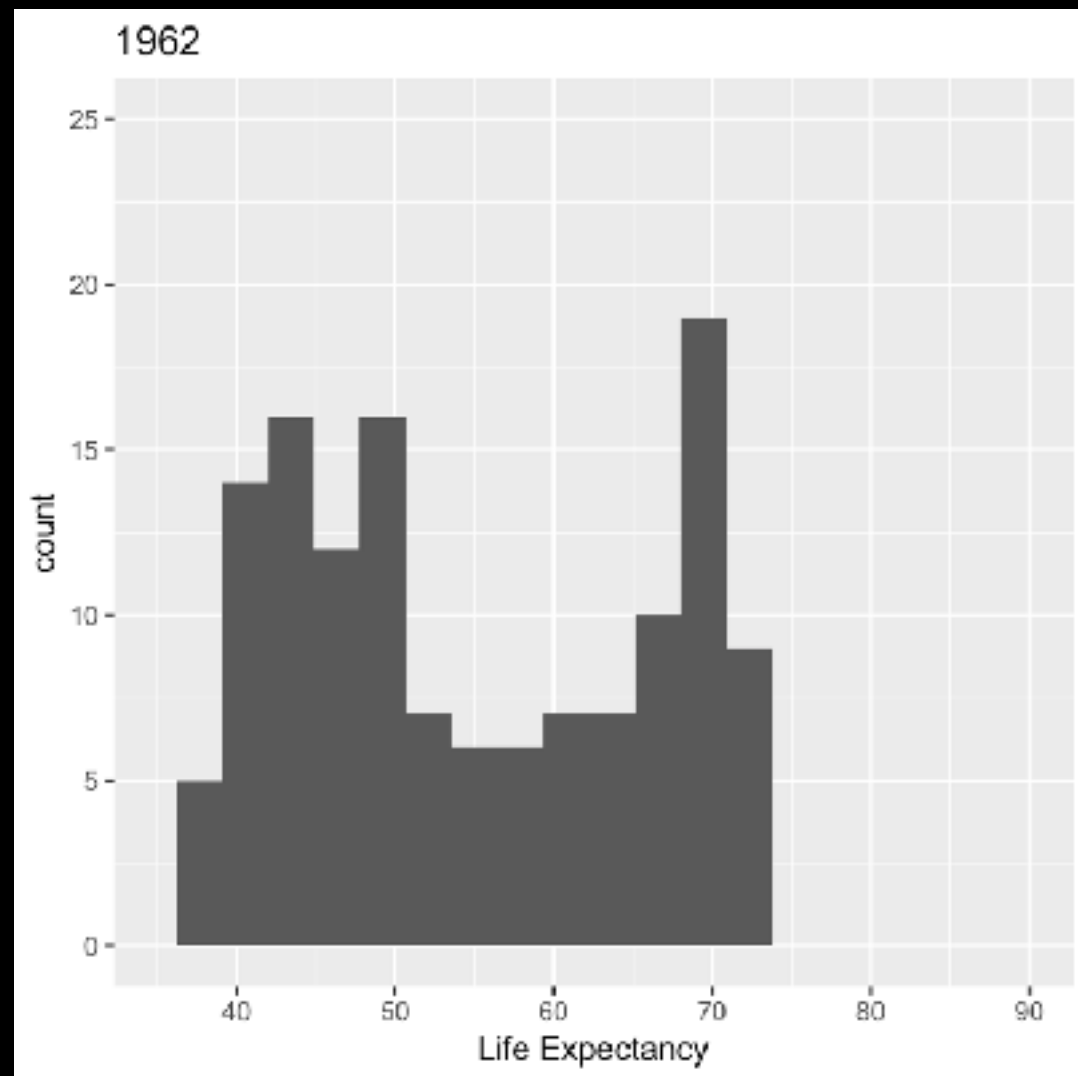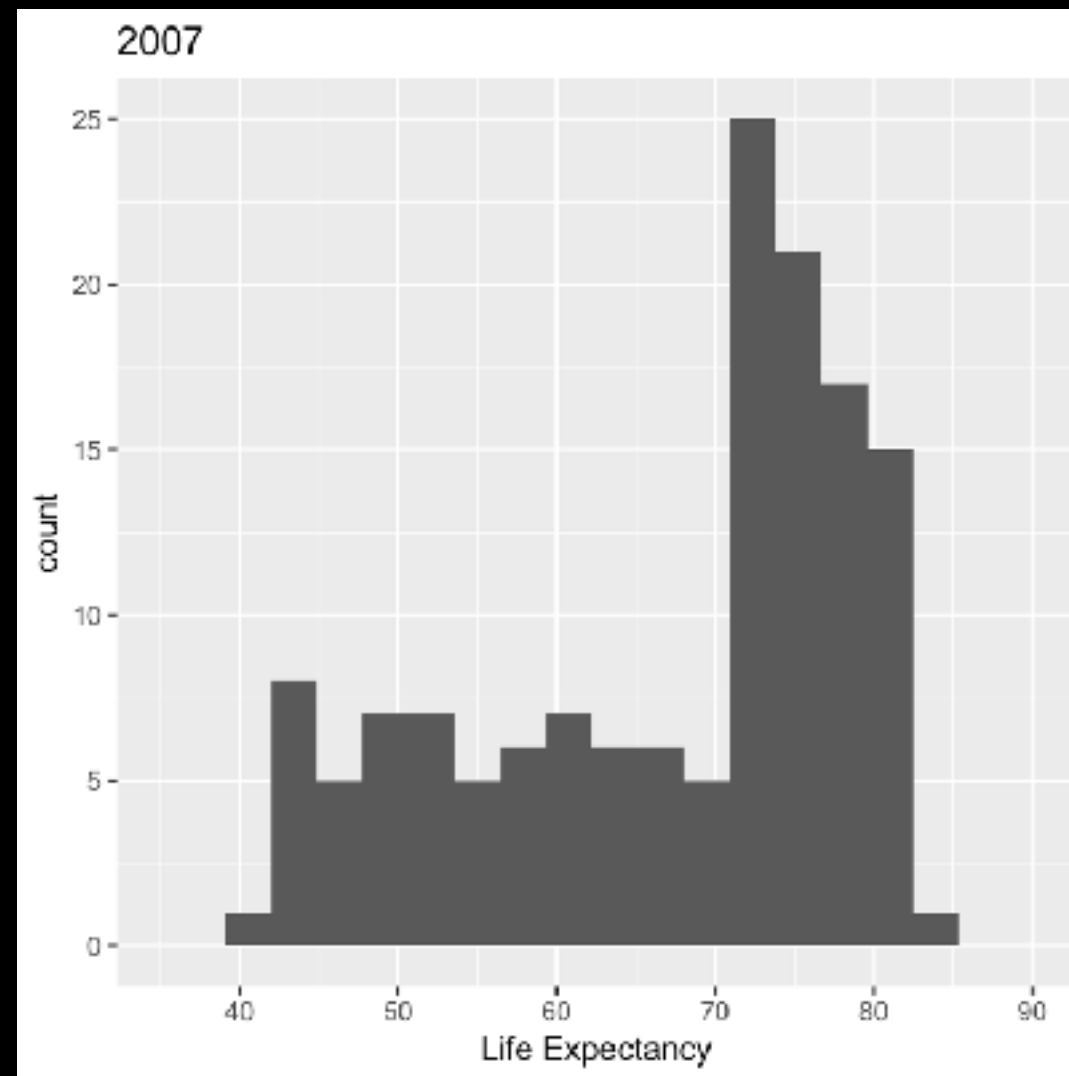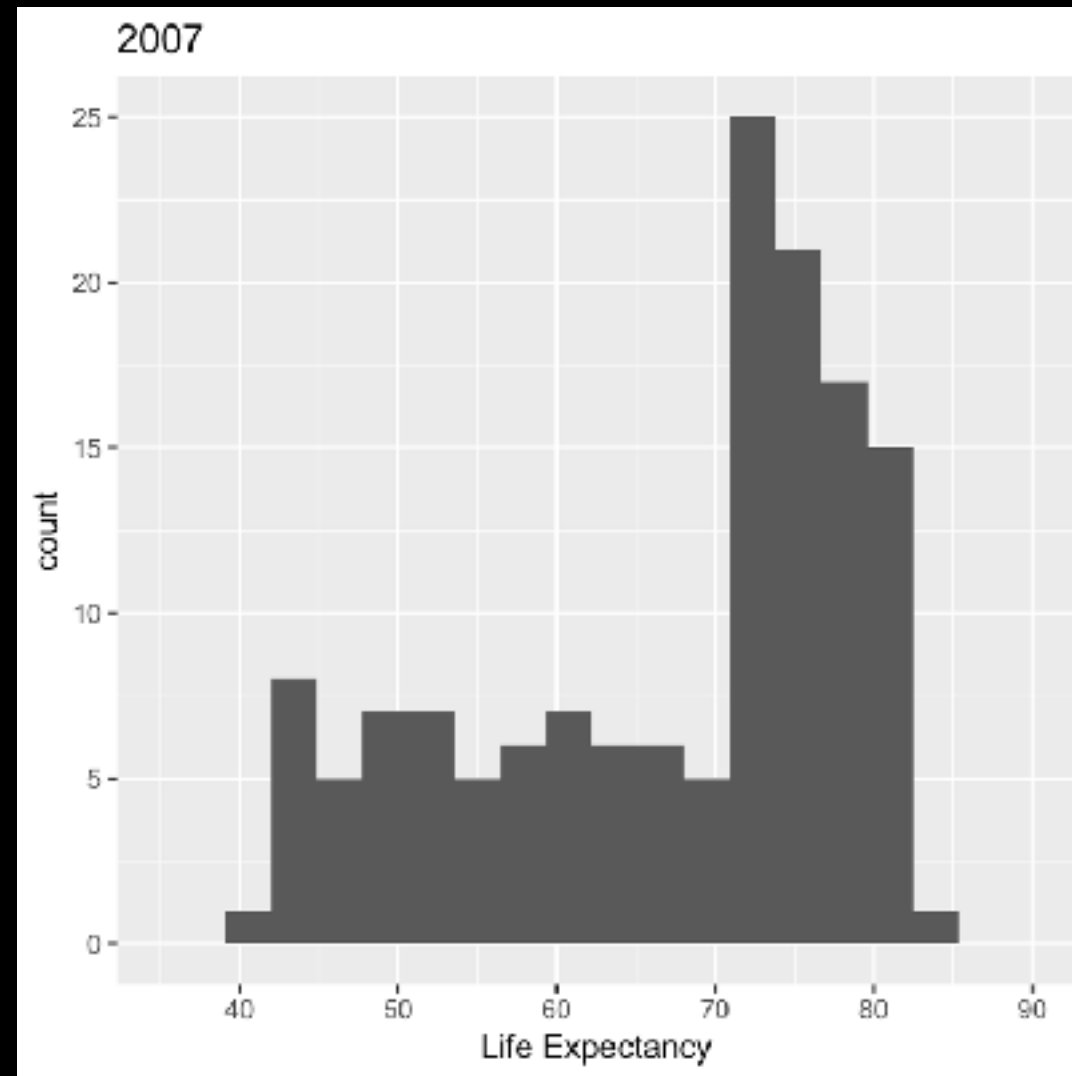Create bar plot of counts
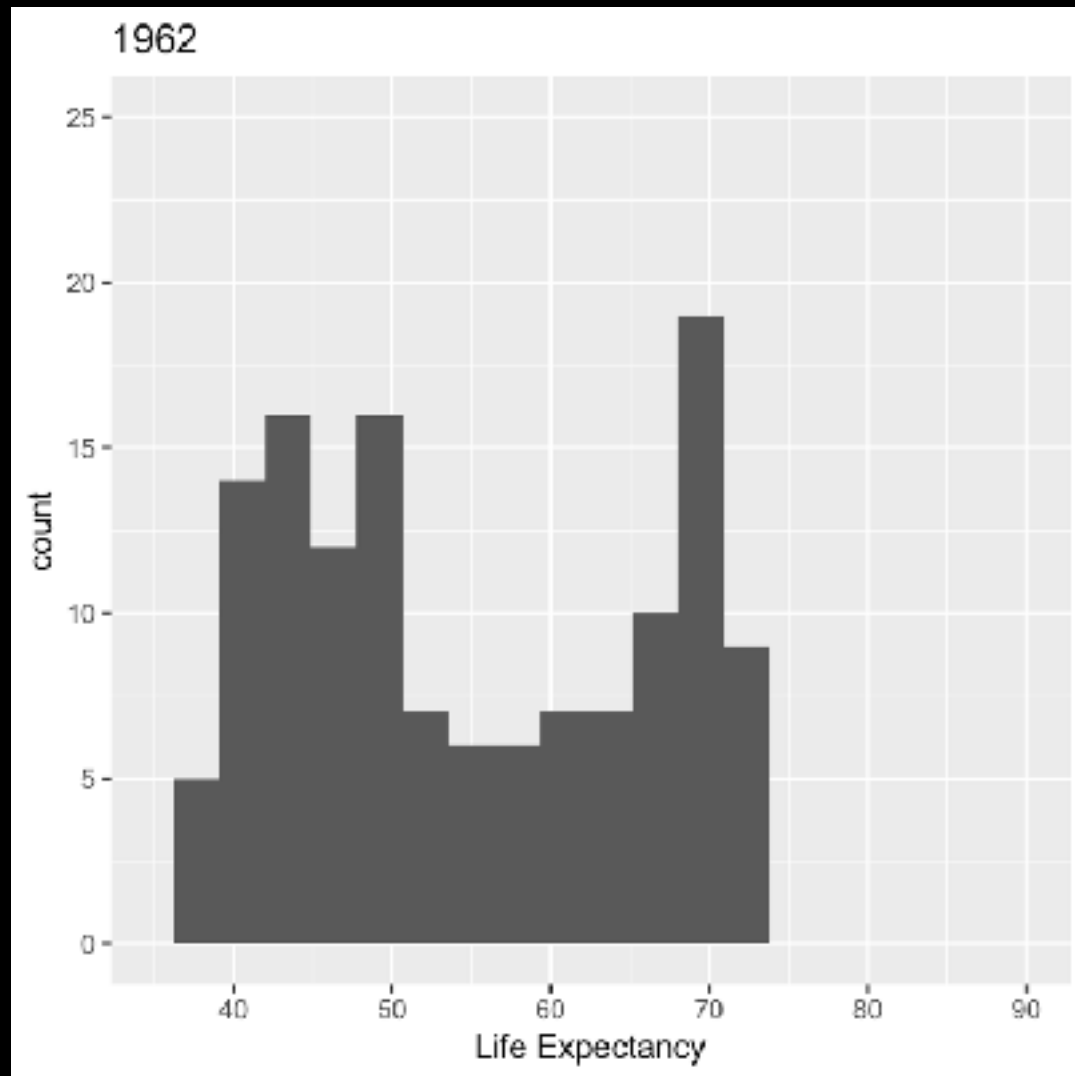
**What are seeing here?**

**What are seeing here?**

**What are seeing here?**

**What's different?**

**What are seeing here?**     **What's different?**
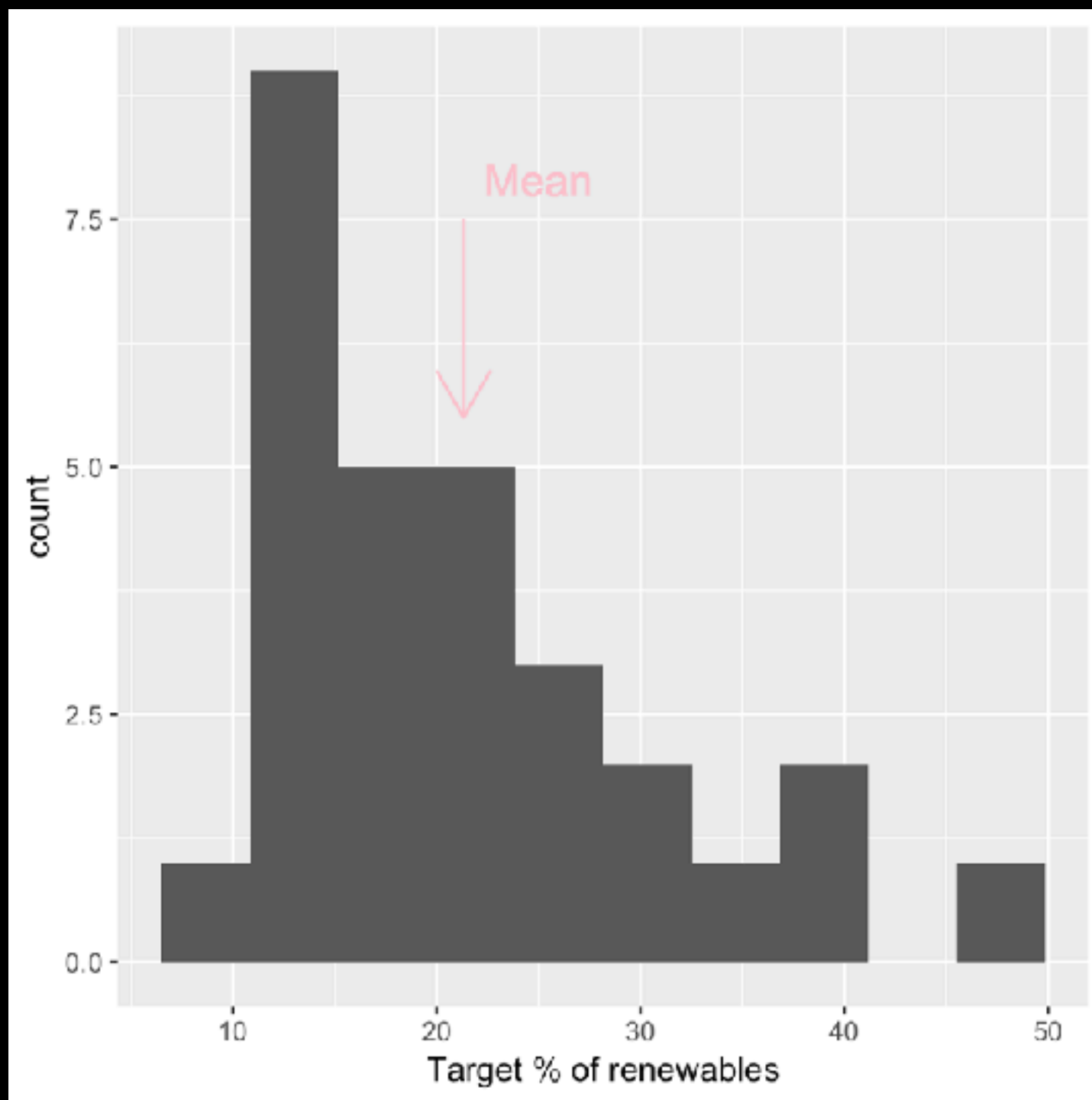
**What did I do to make this comparison?**

**Number of bins**

**Too many - just see noise**
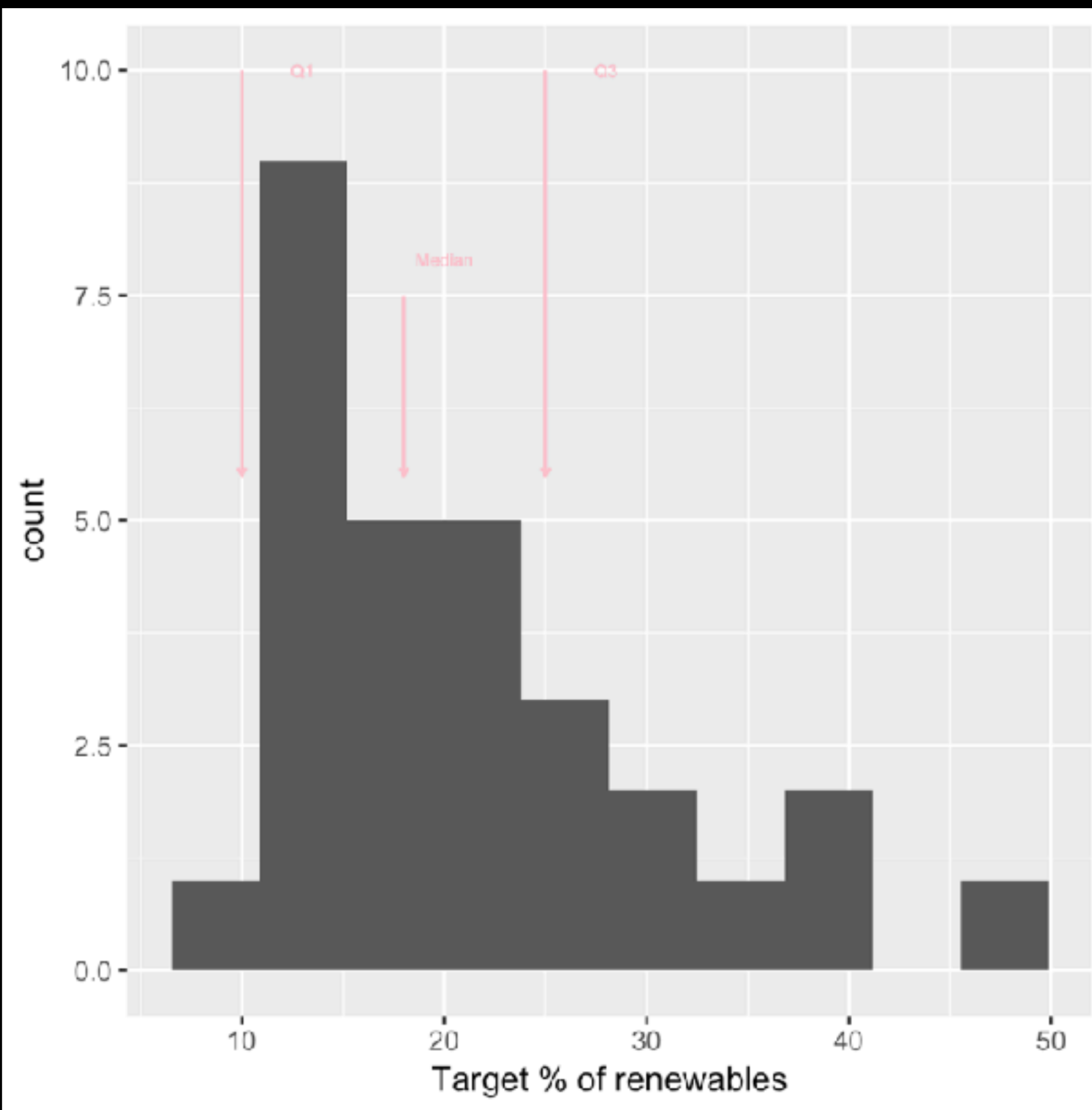
**Too few - don't see any features**

**What about this ?**

**Better to think of quartiles**

**1st quartile (Q1) - value of y which is greater than 25% of the $y_i$**

**Median - value of y which is greater than 50% of the $y_i$**

**3rd quartile (Q3) - value of y which is greater than 75% of the $y_i$**

**What if we want to compare distributions?**

**Life expectancy of countries between years?**

**What if we want to compare distributions?**

**Life expectancy of countries between years?**

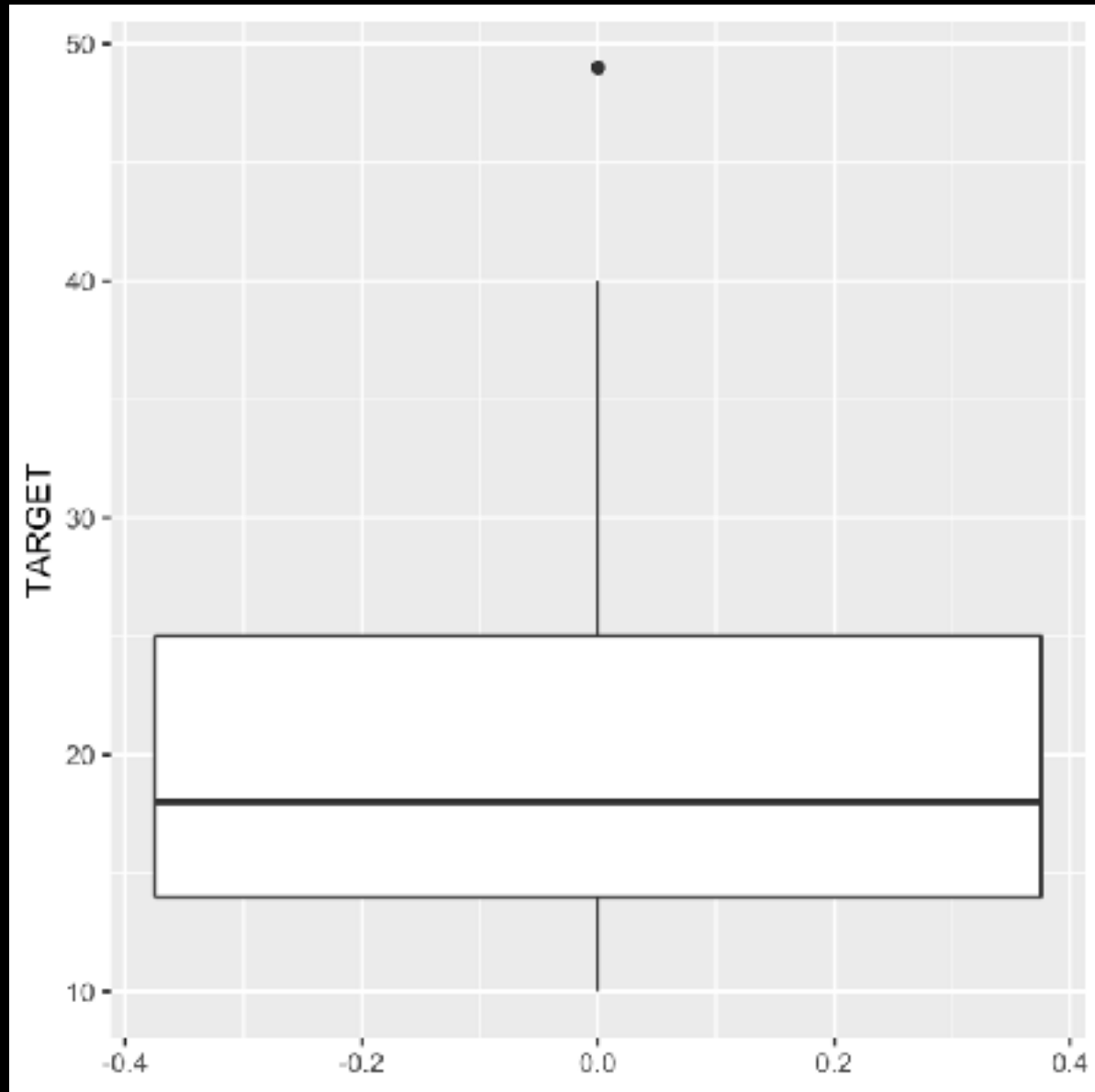**Could do a facet plot of histograms**

**What if we want to compare distributions?**

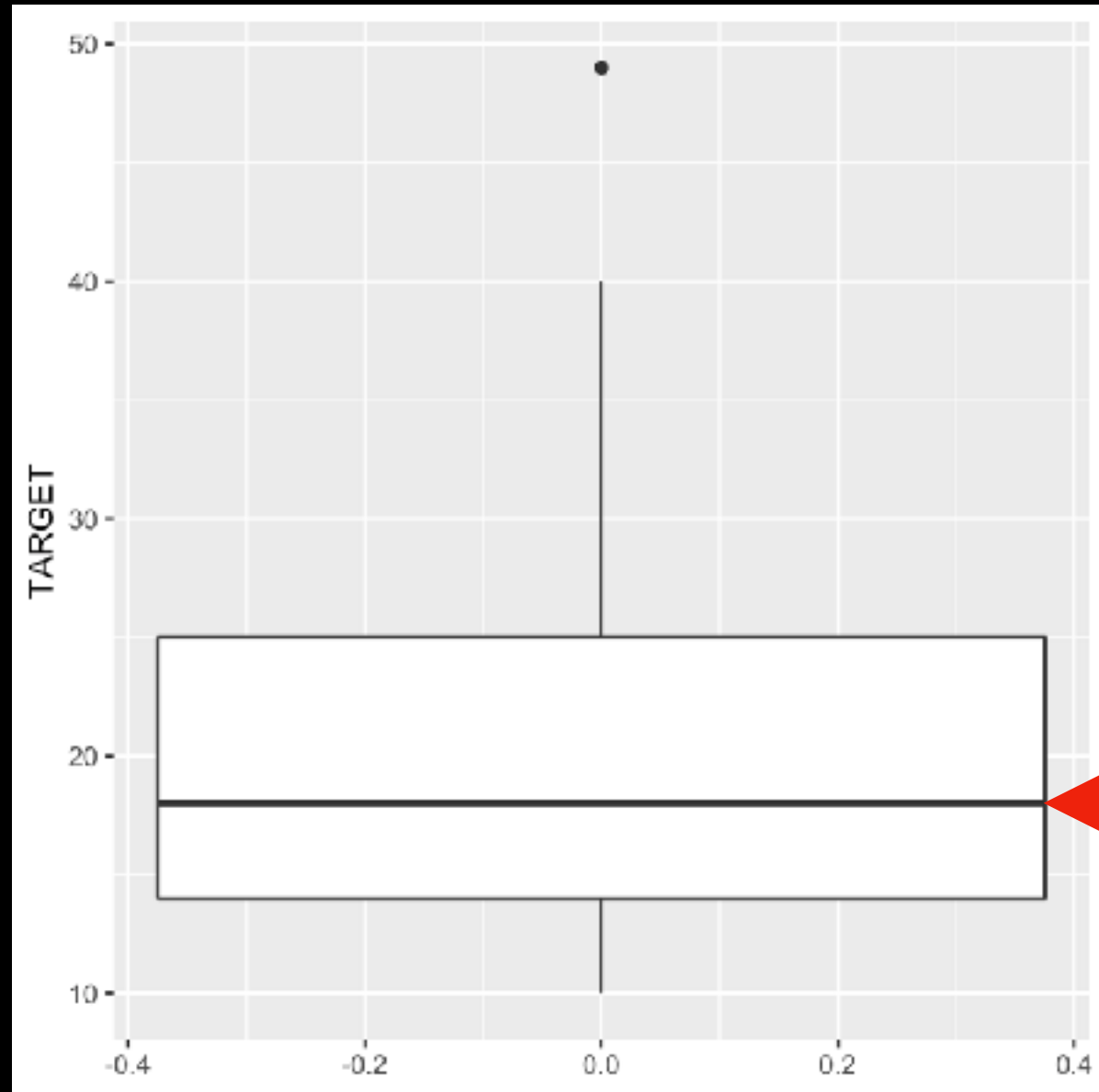**Life expectancy of countries between years?**

**Could do a facet plot of histograms**

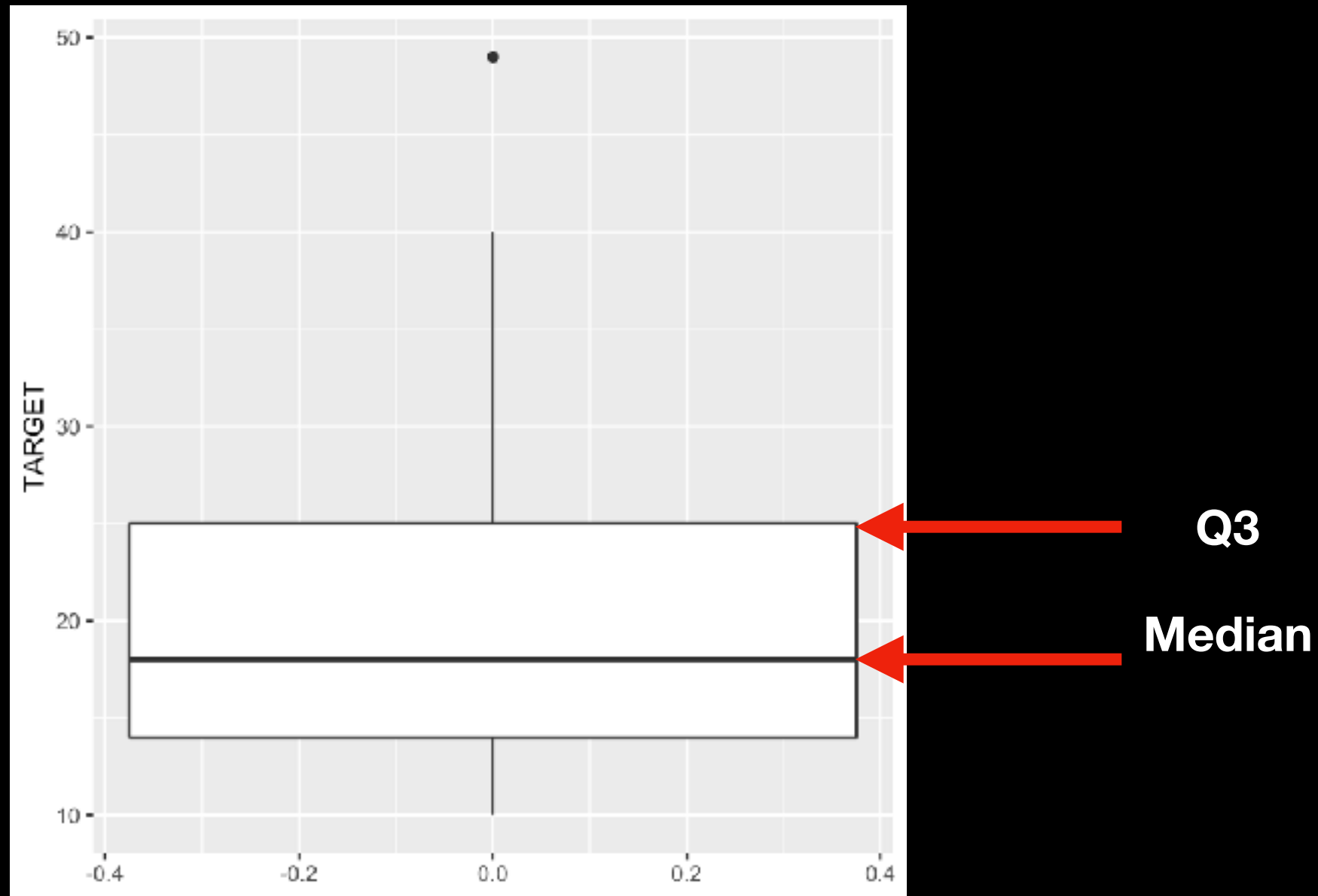**But can also other comparisons**

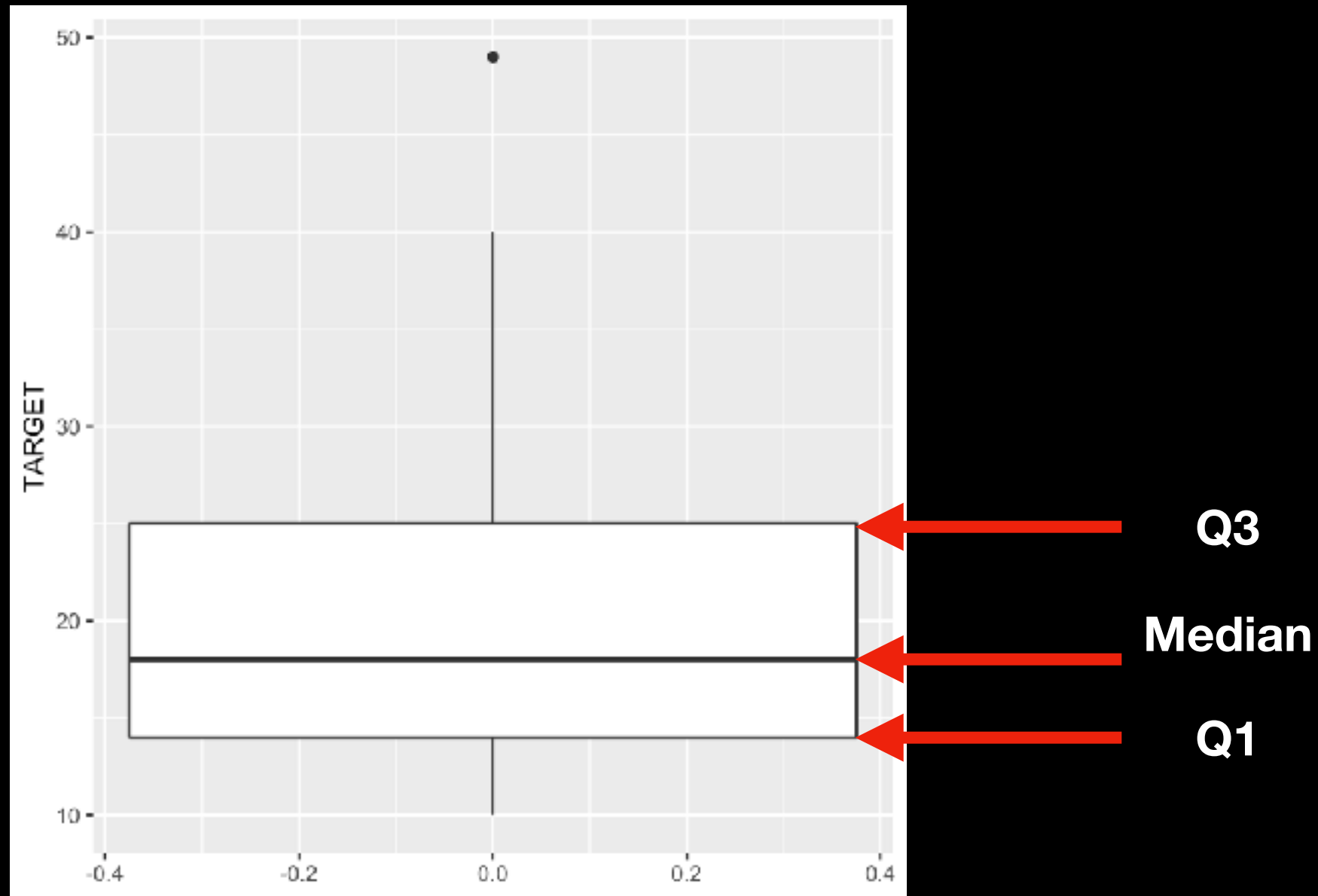# Box plot of renewables target

# Box plot of renewables target

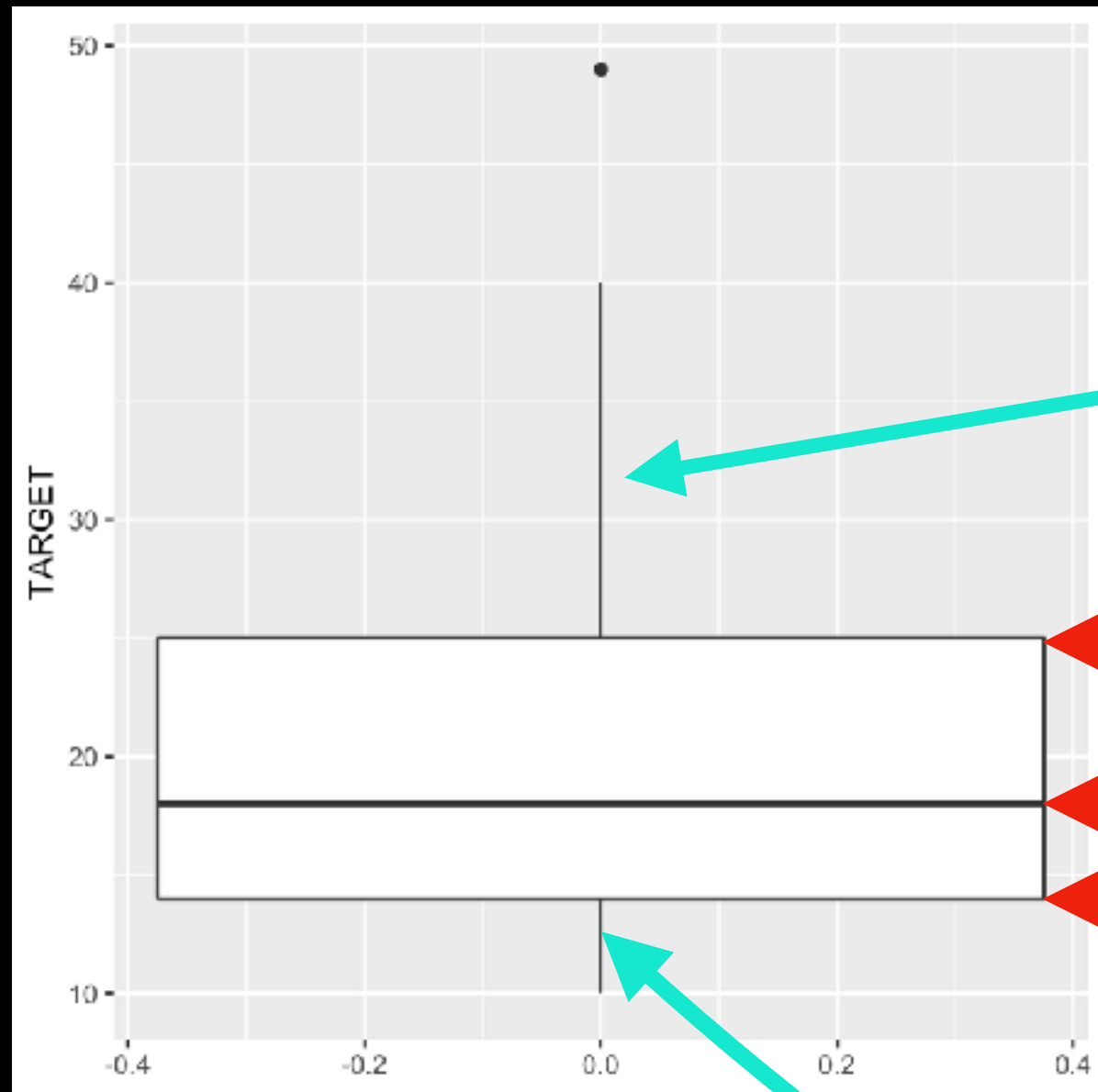# Box plot of renewables target

# Box plot of renewables target

# Box plot of renewables target

# Box plot of renewables target

**Whisker length**
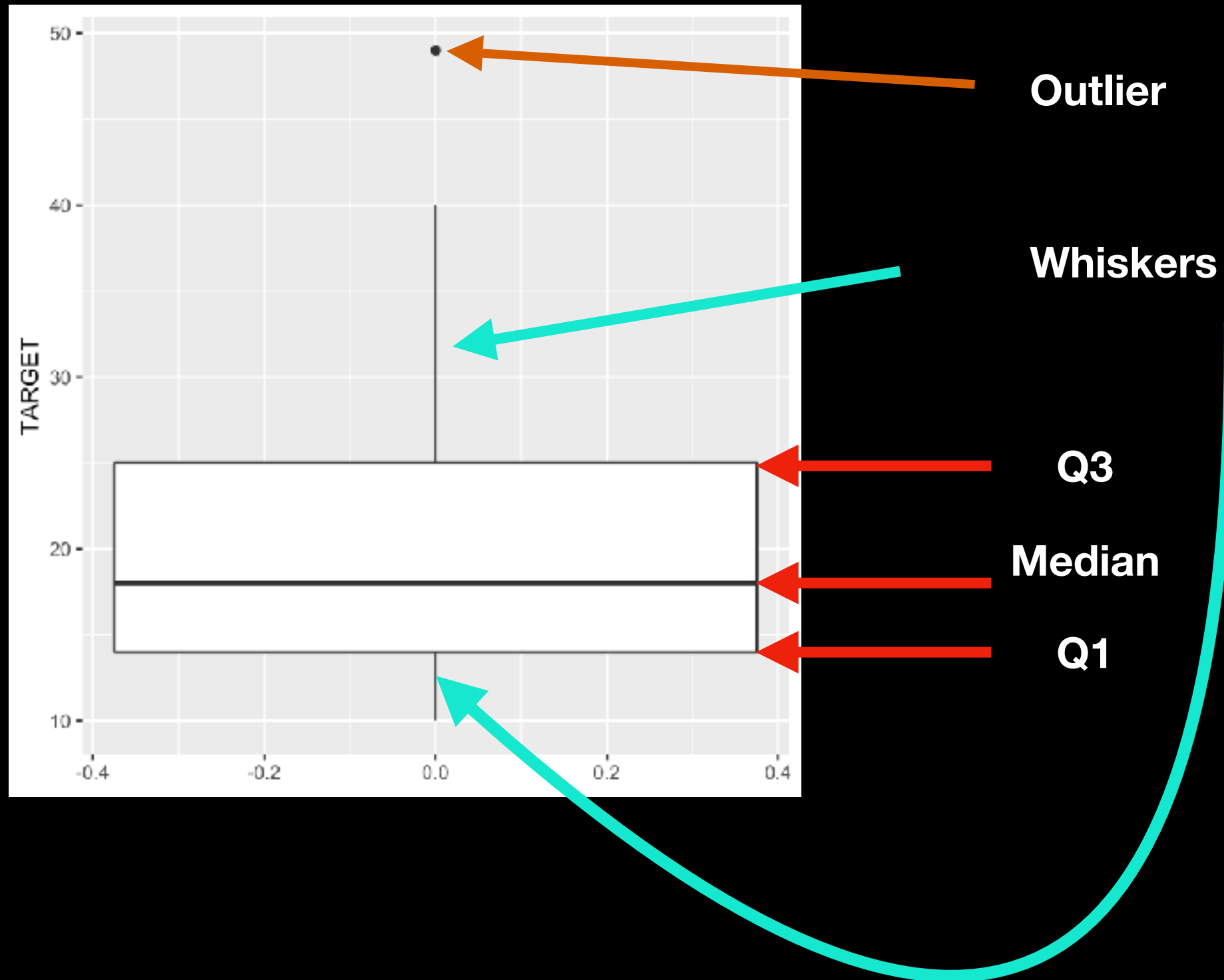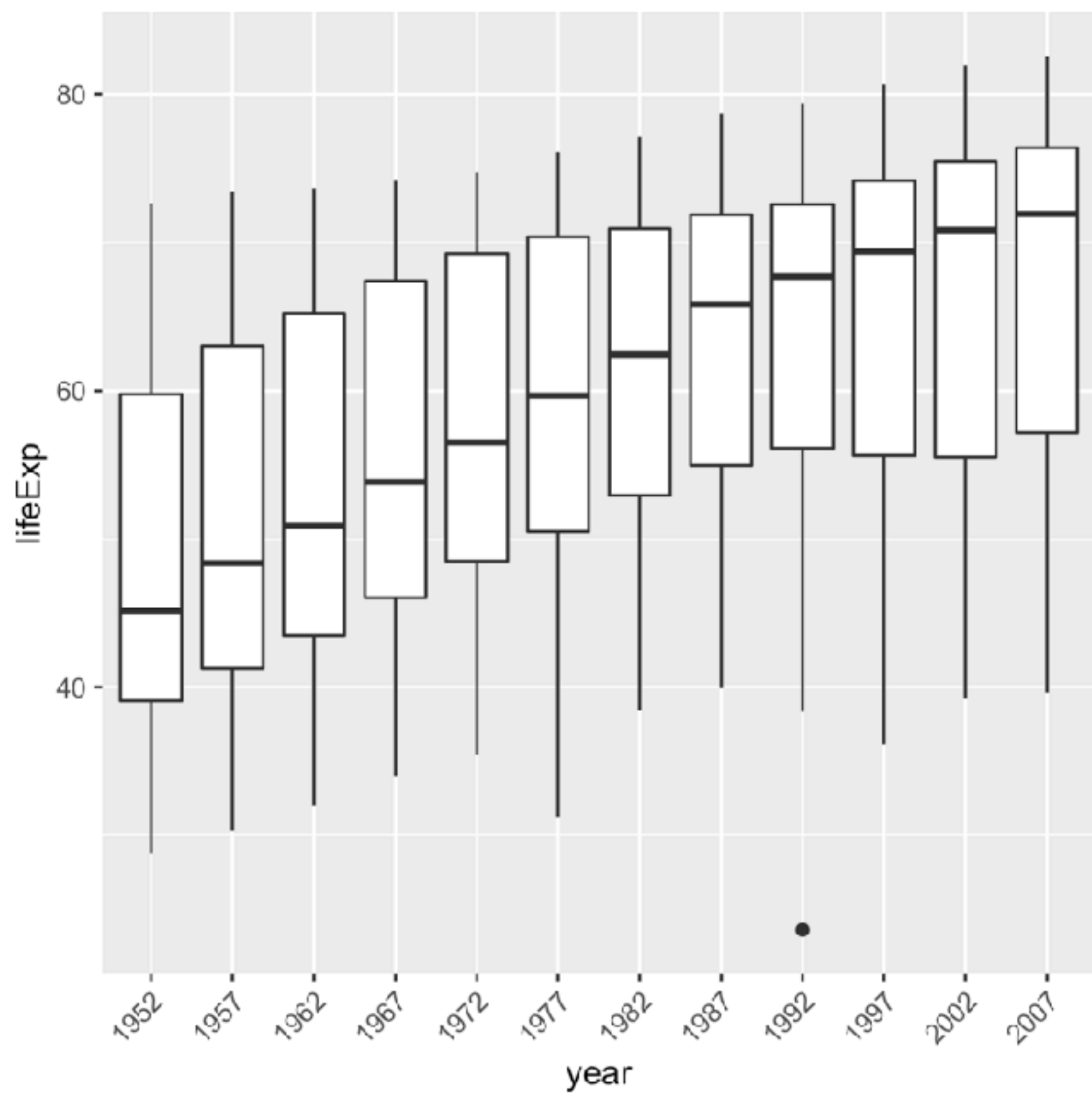
**No bigger than 1.5 x IQR**

**IQR = Q3 - Q1**

**If any data is greater than maximum whisker length
then plotted as point**

**But Life Expectancy isn't unimodal !**

**Try Violin plots**

**Estimate distribution - width of shape indicates height of distribution**