



6.12 Exercise: Inference with iNZight

This exercise will enable you to make comparisons between sub-groups allowing for sampling error.

Background understanding: see [Step 6.9](#)

The skills addressed in this Exercise are:

- Use iNZight to get inferential mark-ups of plots so that you can make visual comparisons between sub-groups allowing for sampling error.
- To obtain numerical confidence limits for true between-group differences.

[*iNZight Lite version [linked here](#)*]

INSTRUCTIONS

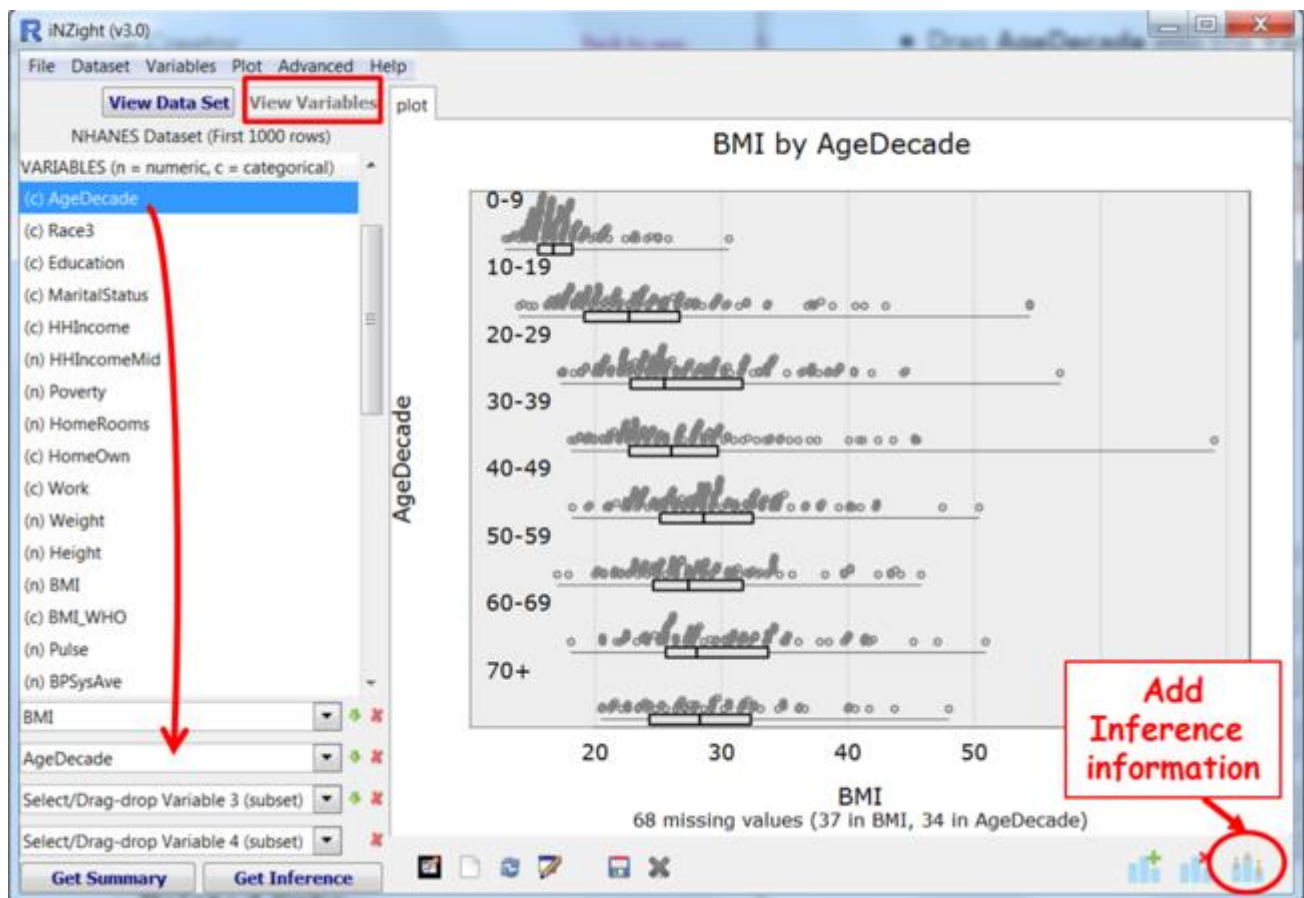
Follow the instructions below to generate the graphs. Or you may prefer to [print the instructions](#). If you have a problem doing the exercise, scroll down to **Common questions**.

Start iNZight (not VIT) and load the `nhanes_1000` dataset into iNZight using **File > Example data ...**. You will find the data set in **Module (package) FutureLearn**.

Construct confidence intervals for sub-groups of a numeric outcome

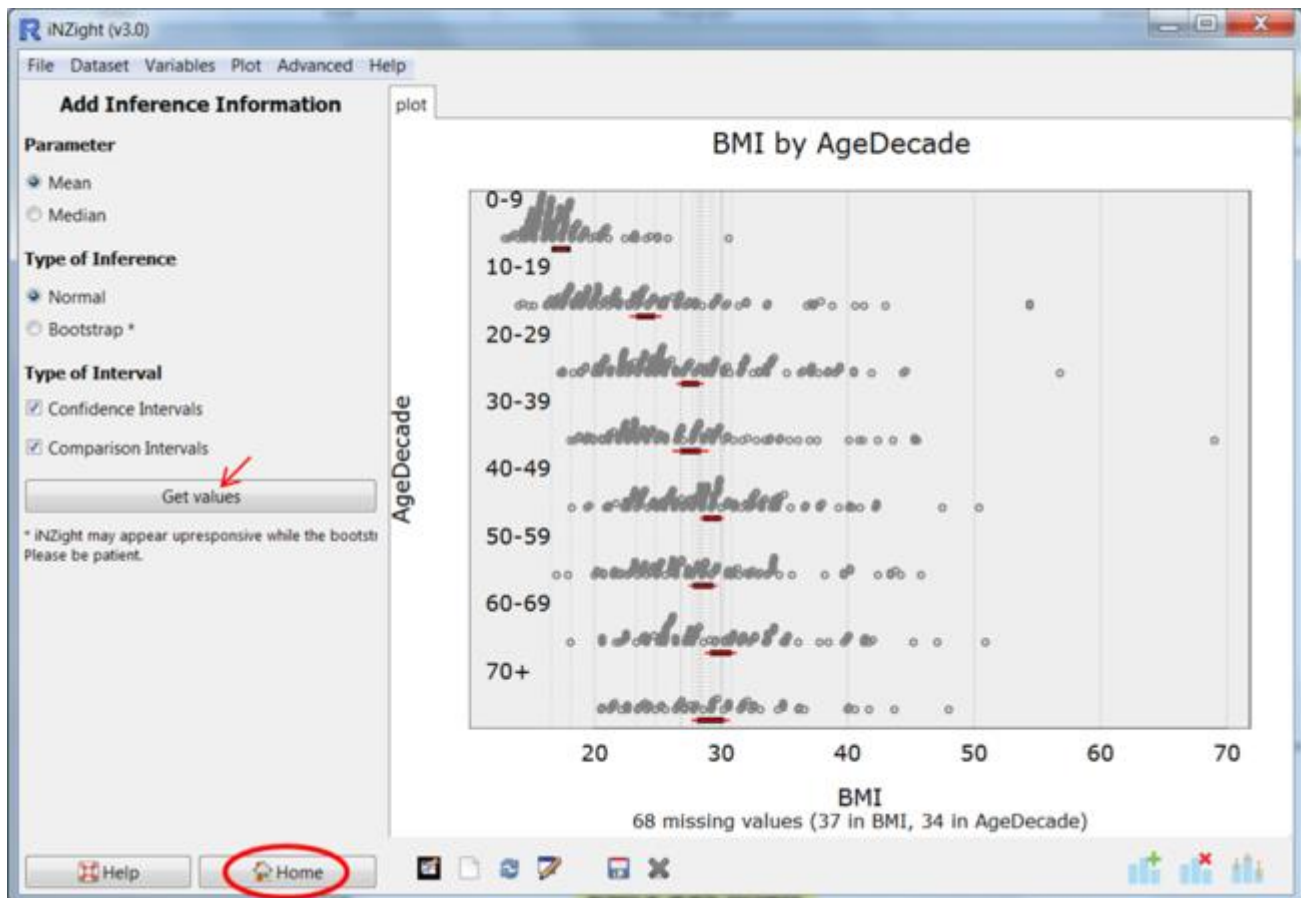
We are revisiting BMI (Body Mass Index) for different age groups and using the nhanes-1000 data, a sample of the American population. We will use our sample to estimate the mean BMI for different age groups. To do this we need to construct intervals around our estimates in order to allow for sampling error.

- Drag BMI into the Variable 1 slot
- Drag AgeDecade into the Variable 2 slot.



You should have a series of dot plots of BMI for different age groups in your plot window.

- Click **Add Inference Information** icon (or go *Plot > Add Inference ...*)



You can squash your plot window vertically so that it is easier to see how much overlap there is between each age group.

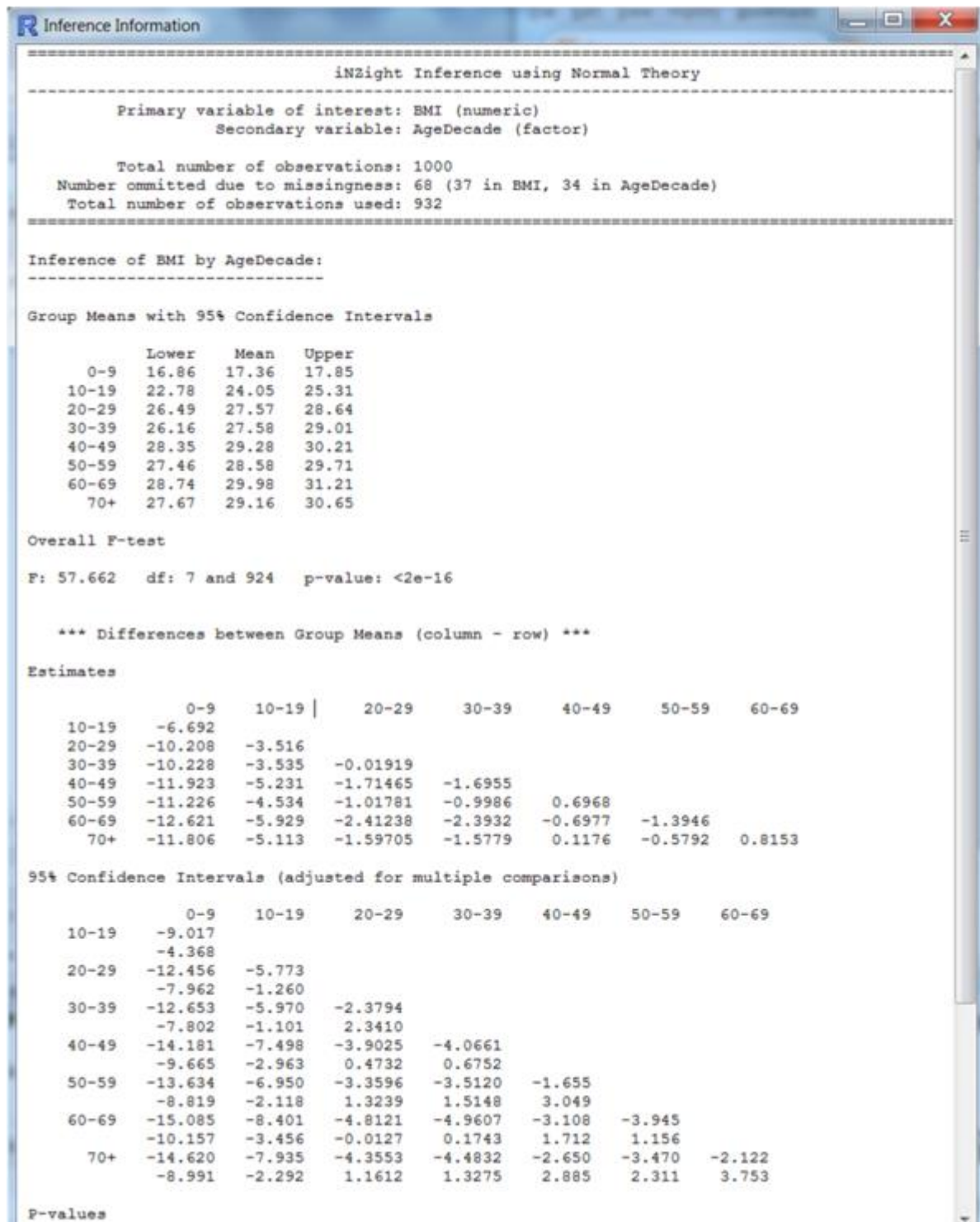
What do you see here? The thick black lines are called 'comparison intervals' and are the lines that we look at when observing any overlap. The thin red lines are the individual confidence intervals for each mean/median. [The endpoints of the lines are available (if you really need them) by clicking *Get values*.]

Post a comment if you have any interesting observations about the estimation of the population means for each age group and the differences between them.

For the actual **confidence intervals for the true differences** in means/medians:

- Return to the **Home** command panel
- Click **Get Inference**
- Select **Normal** inference
- Click **OK**.

The **Summary window** will pop up with the group means estimate and the 95% Confidence Interval around each estimate. You will also see differences between each group mean and their respective confidence intervals as well. There is a lot of extra information in the Summary window. Further study of statistics will help you understand this content.



EXPLORE (~5 min)

Find another numeric variable in the NHANES-1000 dataset that you might like to explore across age groups, e.g. **Height** and move it to the **Variable 1** slot. Use iNZight to explore the estimates for the American population and build intervals around your estimates.

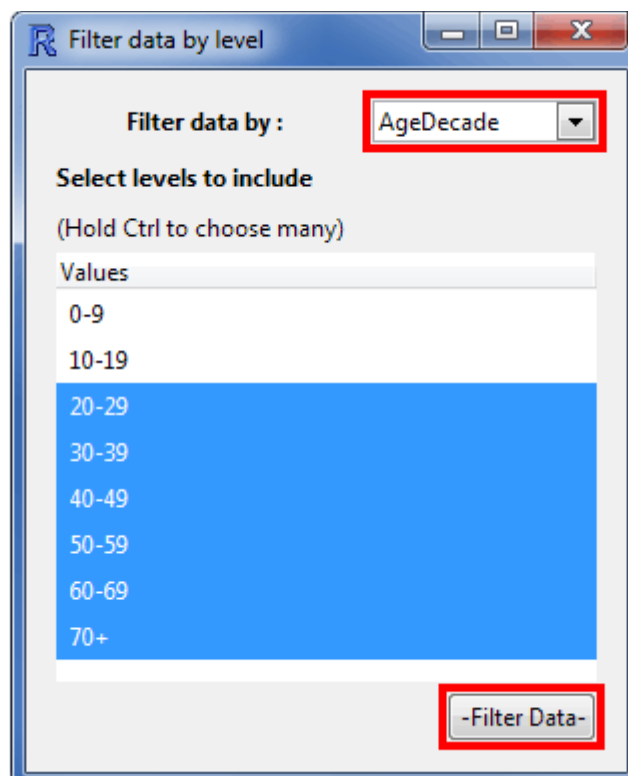
Post a comment if you see anything interesting.

Construct confidence intervals for sub categories of a categorical outcome

Now we'll use the NHANES-1000 dataset to form confidence intervals around estimates of how people in the American population rate their general health at the time of the survey. Is there any difference between age groups?

First, filter out people below the ages of 20 (because "general health" was not recorded for them):

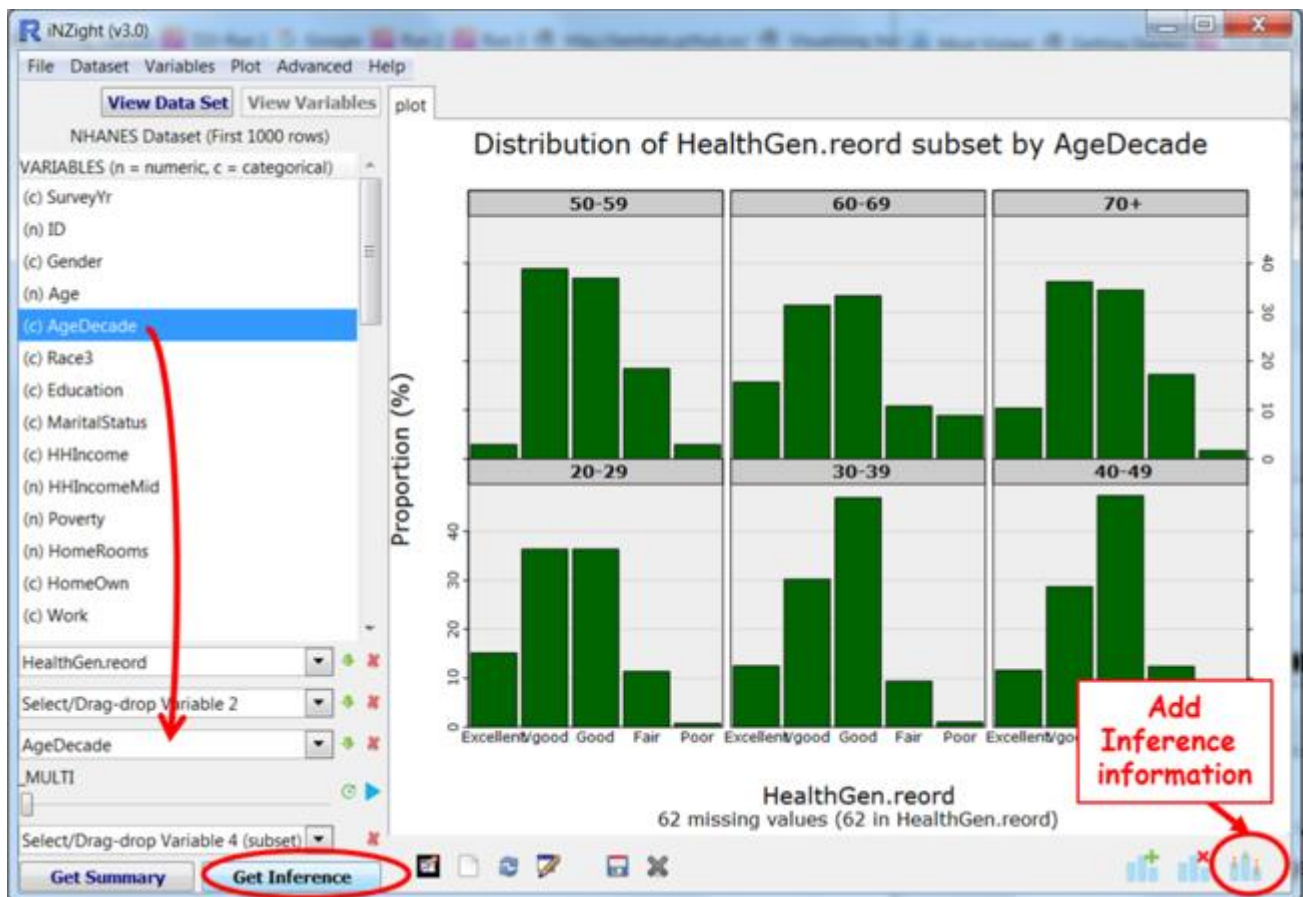
- Click **Dataset > Filter Dataset** and select levels of a categorical variable. Then click **Proceed**.
- **Filter data by:** Select **AgeDecade**
- Using the **Ctrl** or **Shift** keys select all of the age groups you want to include (ages 20 and above)
- Click **-Filter Data-**.



Now re-order the values for **HealthGen** into a sensible order (it is currently alphabetic). We will use the order from 1=Excellent to 5=Poor using **Variables > Categorical Variables > Reorder Levels** to produce the new variable **HealthGen.reord**.

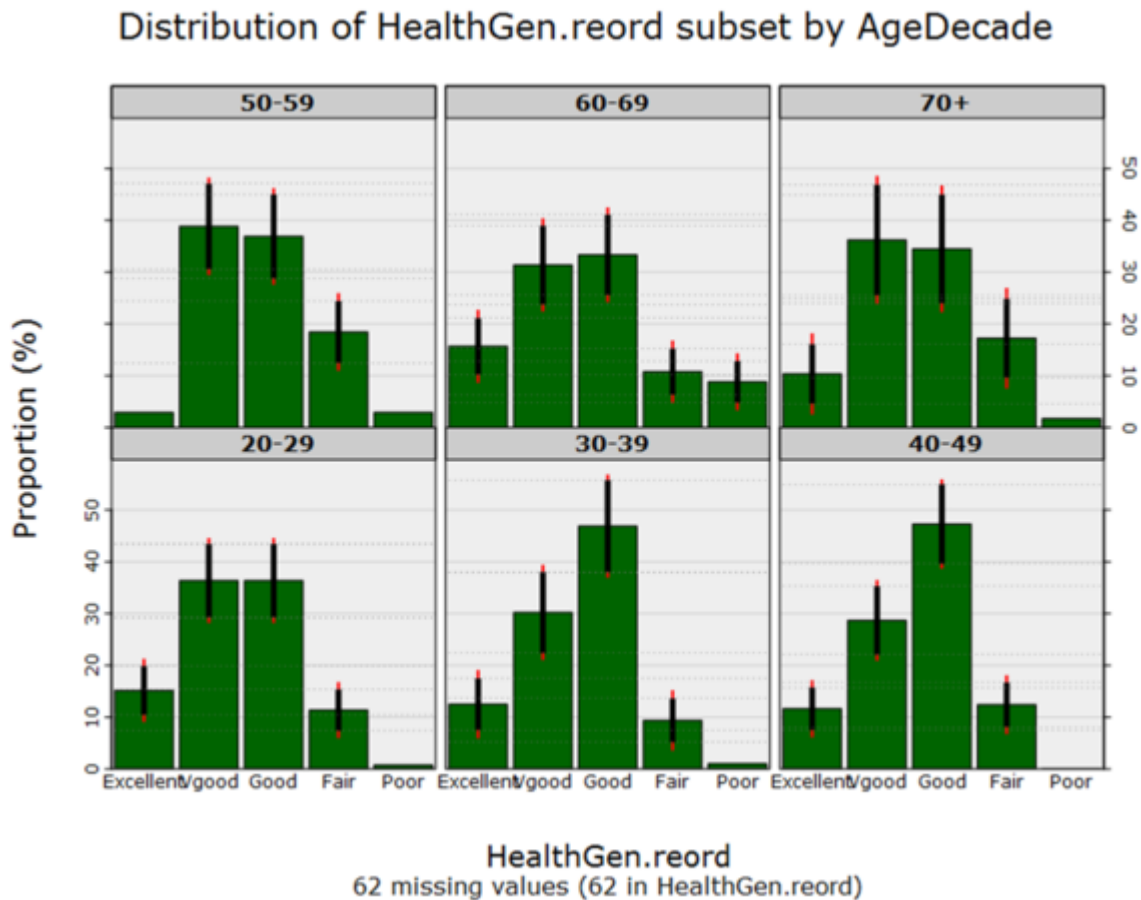
Now construct the series of plots for **HealthGen.reord** and **AgeDecade**:

- Drag **HealthGen.reord** into the **Variable 1** slot
- Drag **AgeDecade** into the **Variable 3 (subset)** slot to create separate bar charts for each age group.



A series of graphs will appear in the plot window. It will help you see the detail if you enlarge your plot window.

- Click **Add Inference Information** icon (or go *Plot > Add Inference ...*)



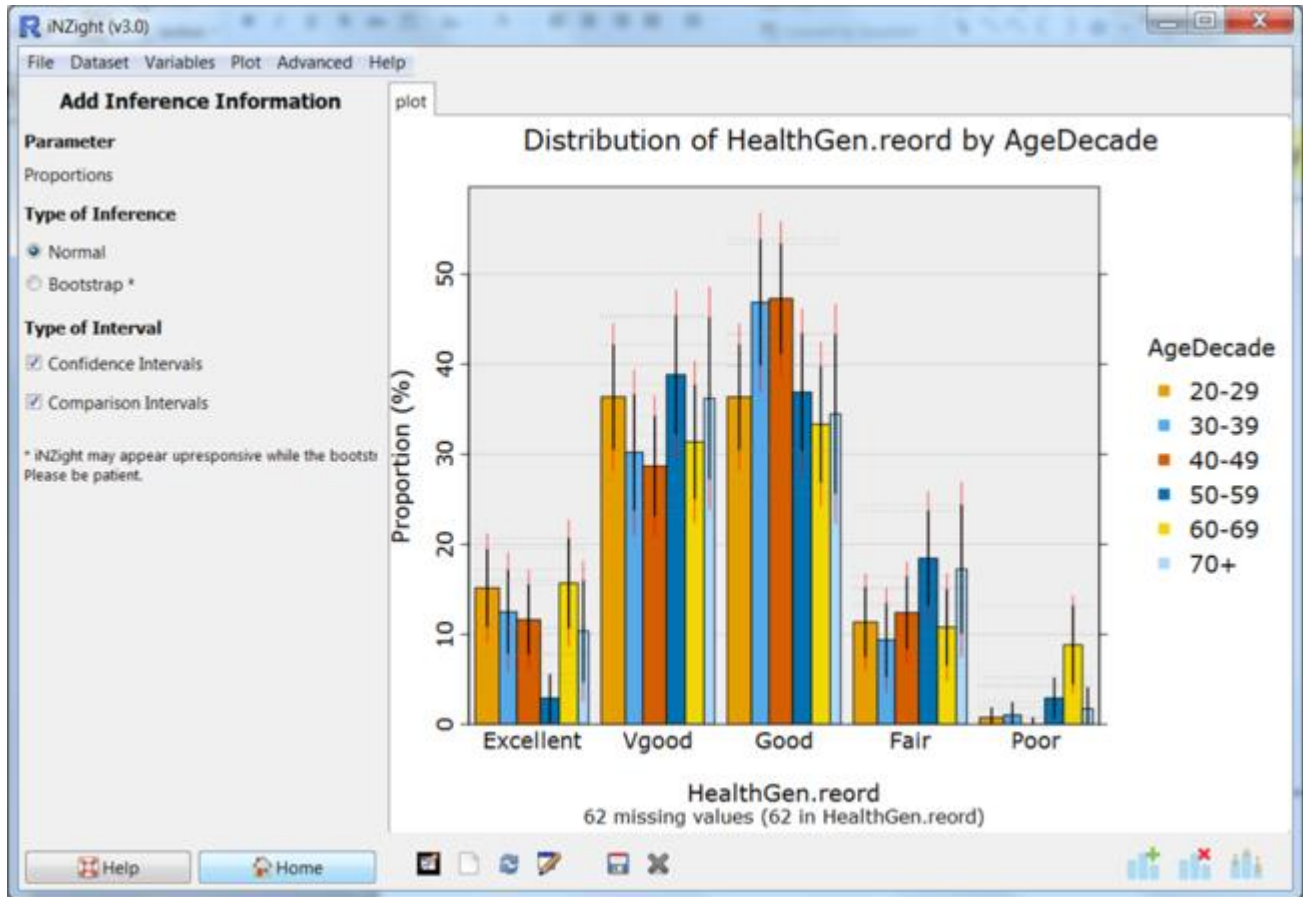
[Again, the endpoints of the lines are available (if you really need them) by clicking *Get values*.]

What does this graph tell you? What can you infer about the true differences between the percentages in each health category?

For the actual **confidence intervals for the true differences** in the proportions, go back to the **Home** command panel and click the **Get Inference** button.

Now construct side-by-side bar charts for HealthGen.reord by AgeDecade

- use the little green "switch" arrow by slot 2 so that HealthGen.reord remains in the Variable 1 slot and AgeDecade is switched into the Variable 2 slot.
- Click the Add Inference Information icon to add the intervals.



What does this graph tell you about age-differences in the percentages for each health category (e.g. differences between age groups in the percentages saying they are in very good health) ?

Again, if you want the actual **confidence intervals** for the **true differences** in the proportions, go back to the **Home** command panel and click the **Get Inference** button.