

Lesson 7: Programming Data Visualisation with R

Dr. Kam Tin Seong

Assoc. Professor of Information Systems

**School of Computing and Information Systems,
Singapore Management University**

2020-2-15 (updated: 2021-06-27)

Content

- Introducing Tidyverse
- *ggplot2*, The Layered Grammar of Graphics
 - The Essential Grammatical Elements in *ggplot2*
 - Designing Analytical Graphics with *ggplot2*
- Interactive Data Visualisation with R
 - *ggiraph* methods
 - *plotly* R methods
- Coordinated Link Views with R
 - *crosstalk* methods

Introducing Tidyverse

tidyverse is an opinionated collection of R packages designed for data science. All packages share an underlying design philosophy, grammar, and data structures.

The tidyverse

Components

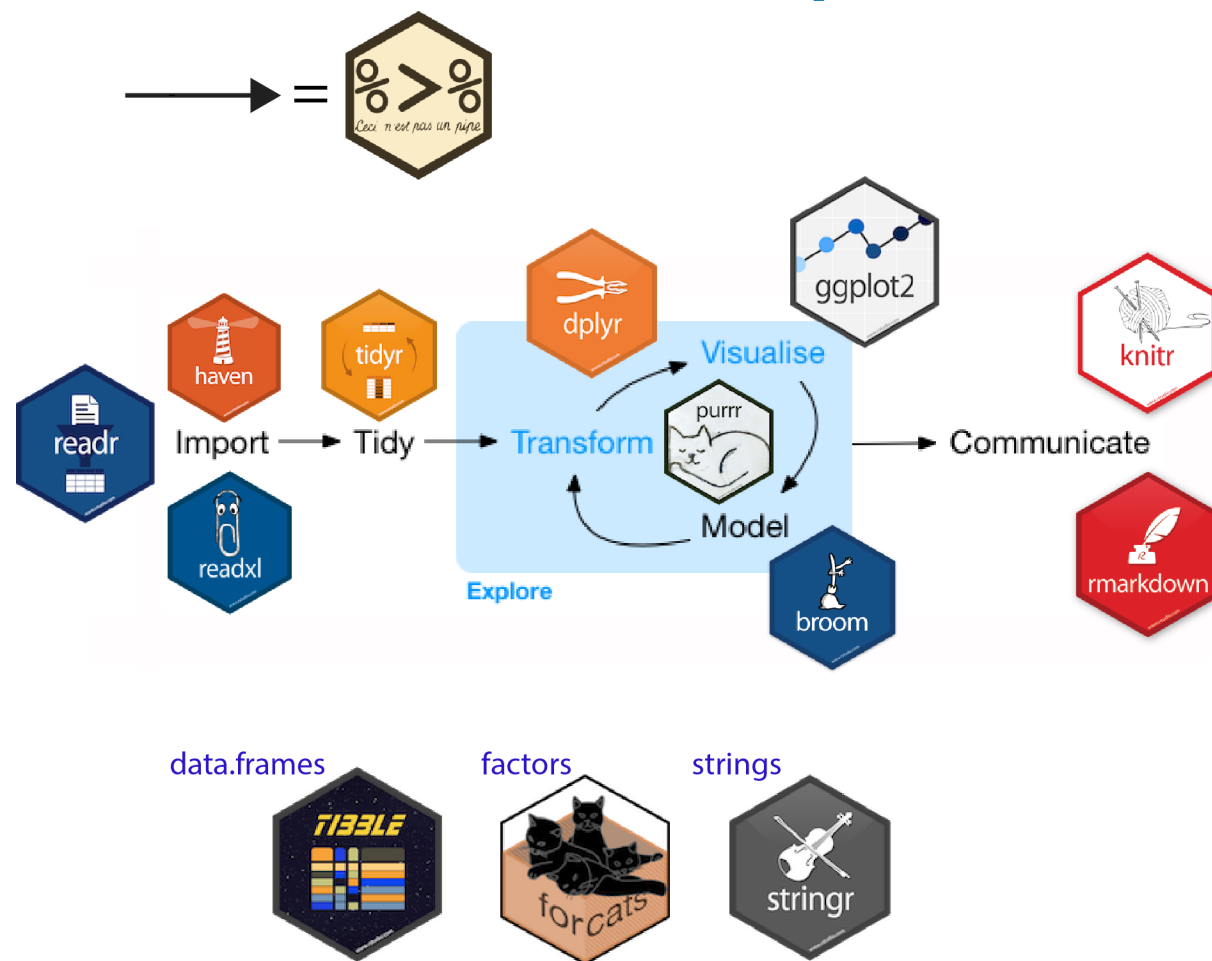


Reference: [tidyverse](#)

Core Tidyverse packages

- **dplyr** is a grammar of data manipulation, providing a consistent set of verbs that help you solve the most common data manipulation challenges.
- **tidyr** helps R users to create tidy data.
- **stringr** provides a cohesive set of functions designed to make working with strings as easy as possible.
- **forcats** provides a suite of tools that solve common problems with factors, including changing the order of levels or the values.
- **readr** provides a fast and friendly way to read rectangular data (like csv, tsv, and fwf).
- **tibble** is a modern reimagining of the data.frame, keeping what time has proven to be effective, and throwing out what is not.
- **ggplot2** is a system for declaratively creating graphics, based on The Grammar of Graphics.
- **purrr** enhances R's functional programming (FP) toolkit by providing a complete and consistent set of tools for working with functions and vectors.

Data Science Workflow with Tidyverse



Reference: [Introduction to the Tidyverse: How to be a tidy data scientist](#)

Getting started

Installing and loading the required libraries

- Before we get started, it is important for us to ensure that the required R packages have been installed. If yes, we will load the R packages. If they have yet to be installed, we will install the R packages and load them onto R environment.
- The chunk code on the right will do the trick.

```
packages = c('DT', 'ggiraph', 'patchwork',  
             'plotly', 'tidyverse')  
  
for(p in packages){  
  if(!require(p, character.only = T)){  
    install.packages(p)  
  }  
  library(p, character.only = T)  
}
```

Importing data

- The code chunk below imports *exam_data.csv* into R environment using *read_csv()* function of **readr** package.
- **readr** is one of the tidyverse package.

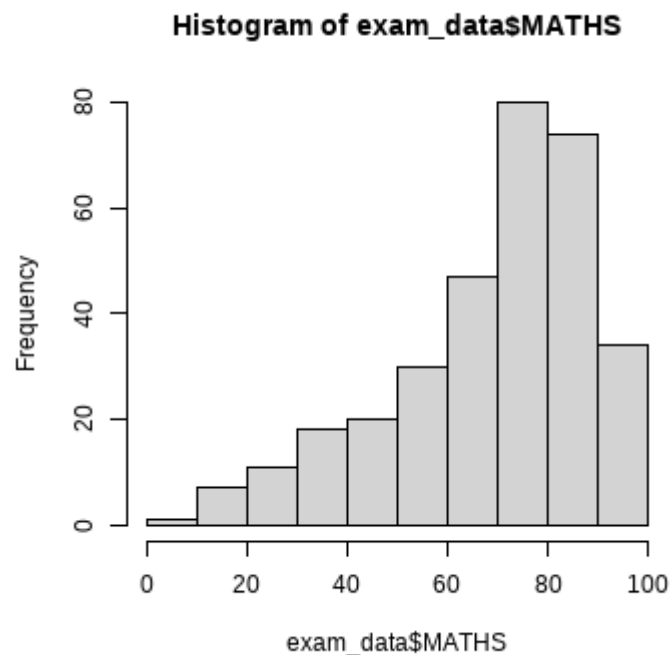
```
exam_data <- read_csv("data/Exam_data.csv")
```

- Year end examination grades of a cohort of primary 3 students from a local school.
- There are a total of seven attributes. Four of them are categorical data type and the other three are in continuous data type.
 - The categorical attributes are: ID, CLASS, GENDER and RACE.
 - The continuous attributes are: MATHS, ENGLISH and SCIENCE.

Introducing ggplot2

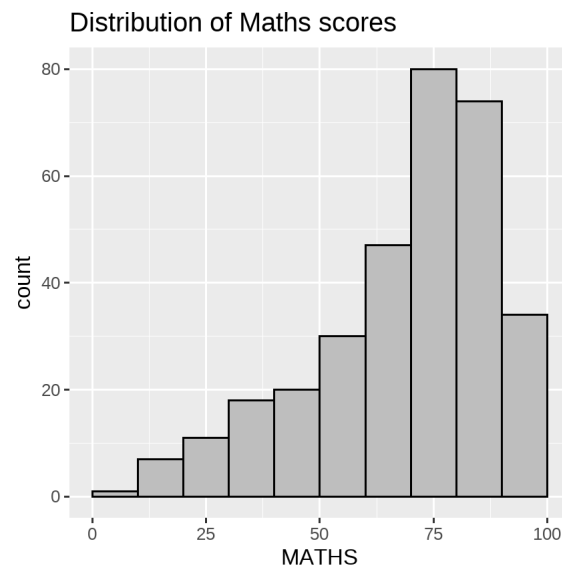
R Graphics

```
hist(exam_data$MATHS)
```



ggplot2

```
ggplot(data=exam_data, aes(x = MATHS)) +  
  geom_histogram(bins=10,  
                 boundary = 100,  
                 color="black",  
                 fill="grey") +  
  ggtitle("Distribution of Maths scores")
```



Then, why ggplot2

The transferable skills from ggplot2 are not the idiosyncrasies of plotting syntax, but a powerful way of thinking about visualisation, as a way of mapping between variables and the visual properties of geometric objects that you can perceive.

Hadley Wickham

```
ggplot(data=exam_data, aes(x = MATHS)) +  
  geom_histogram(bins=10,  
                 boundary = 100,  
                 color="black",  
                 fill="grey") +  
  ggtitle("Distribution of Maths scores")
```

Grammar of Graphics

- Wilkinson, L. (1999) **Grammar of Graphics**, Springer.
- The grammar of graphics is an answer to a question:

What is a statistical graphic?

- Grammar of graphics defines the rules of structuring mathematical and aesthetic elements into a meaningful graph.
- Two principles
 - Graphics = distinct layers of grammatical elements
 - Meaningful plots through aesthetic mapping

- A good grammar will allow us to gain insight into the composition of complicated graphics, and reveal unexpected connections between seemingly different graphics (Cox 1978).
- A grammar provides a strong foundation for understanding a diverse range of graphics.
- A grammar may also help guide us on what a well-formed or correct graphic looks like, but there will still be many grammatically correct but nonsensical graphics.

Essential Grammatical Elements in ggplot2

A Layered Grammer of Graphics

- **Data**: The dataset being plotted.
- **Aesthetics** take attributes of the data and use them to influence visual characteristics, such as position, colours, size, shape, or transparency.
- **Geometrics**: The visual elements used for our data, such as point, bar or line.
- **Facets** split the data into subsets to create multiple variations of the same graph (paneling, multiple plots).
- **Statistics**, statistical transformations that summarise data (e.g. mean, confidence intervals).
- **Coordinate systems** define the plane on which data are mapped on the graphic.
- **Themes** modify all non-data components of a plot, such as main title, sub-title, y-axis title, or legend background.



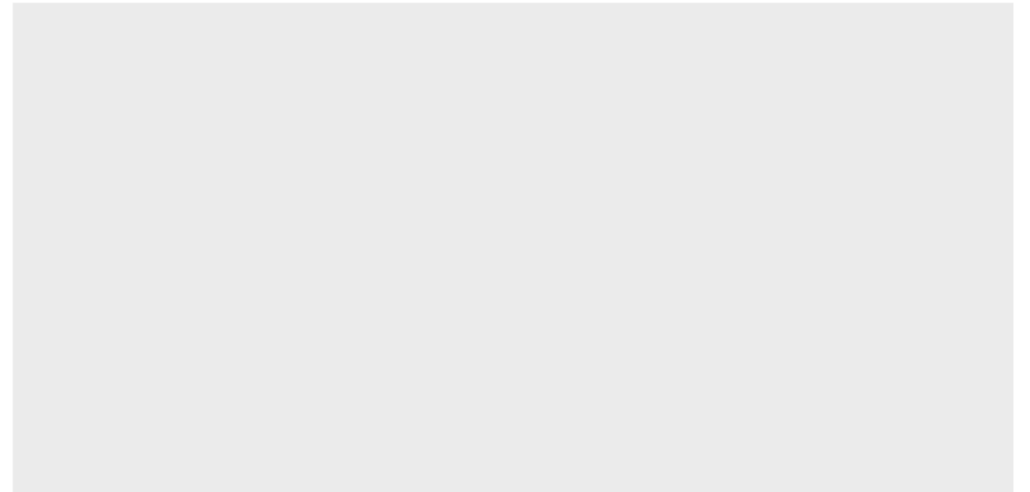
Reference: Hadley Wickham (2010) "A layered grammar of graphics."
Journal of Computational and Graphical Statistics, vol. 19, no. 1, pp. 3–28.
<https://vita.had.co.nz/papers/layered-grammar.html>

Essential Grammatical Elements in ggplot2

The *ggplot()* function and *data* argument

- Let us call the *ggplot()* function using the code chunk on the right.
- Notice that a blank canvas appears.
- *ggplot()* initializes a ggplot object.
- The *data* argument defines the dataset to be used for plotting.
- If the dataset is not already a `data.frame`, it will be converted to one by *fortify()*.

```
ggplot(data=exam_data)
```



Essential Grammatical Elements in ggplot2

The Aesthetic mappings

- The aesthetic mappings take attributes of the data and use them to influence visual characteristics, such as position, colour, size, shape, or transparency.
- Each visual characteristic can thus encode an aspect of the data and be used to convey information.
- All aesthetics of a plot are specified in the `aes()` function call (in later part of this lesson, you will see that each *geom* layer can have its own aes specification)

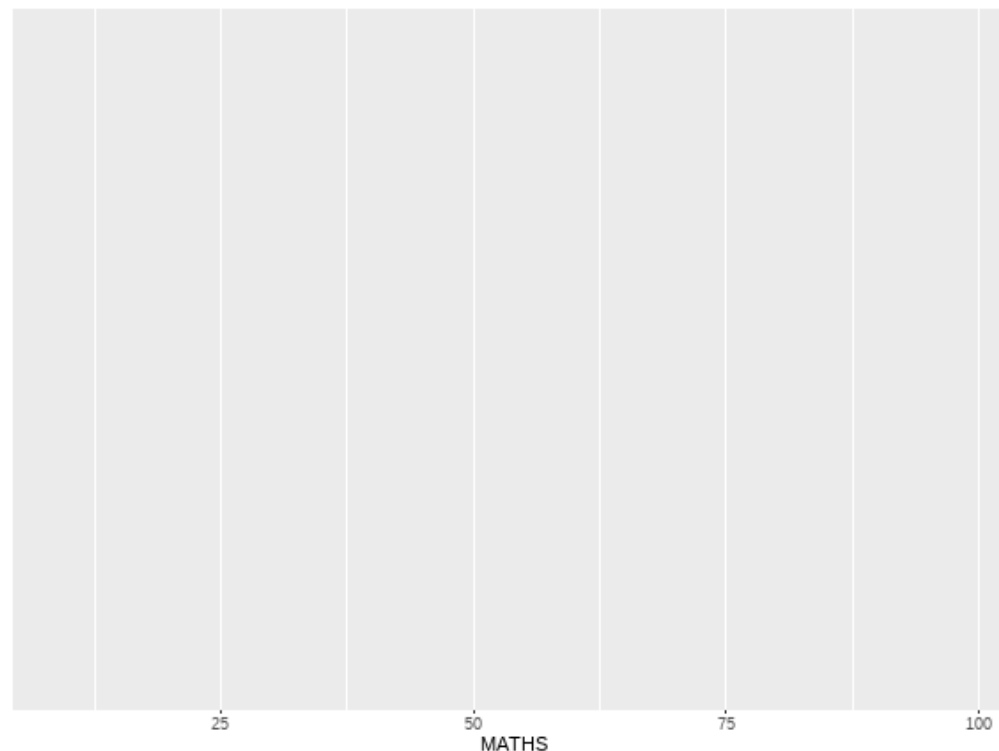
Essential Grammatical Elements in ggplot2

Working with *aes()*

- The code chunk on the right add the aesthetic element into the plot.

```
ggplot(data=exam_data,  
       aes(x= MATHS))
```

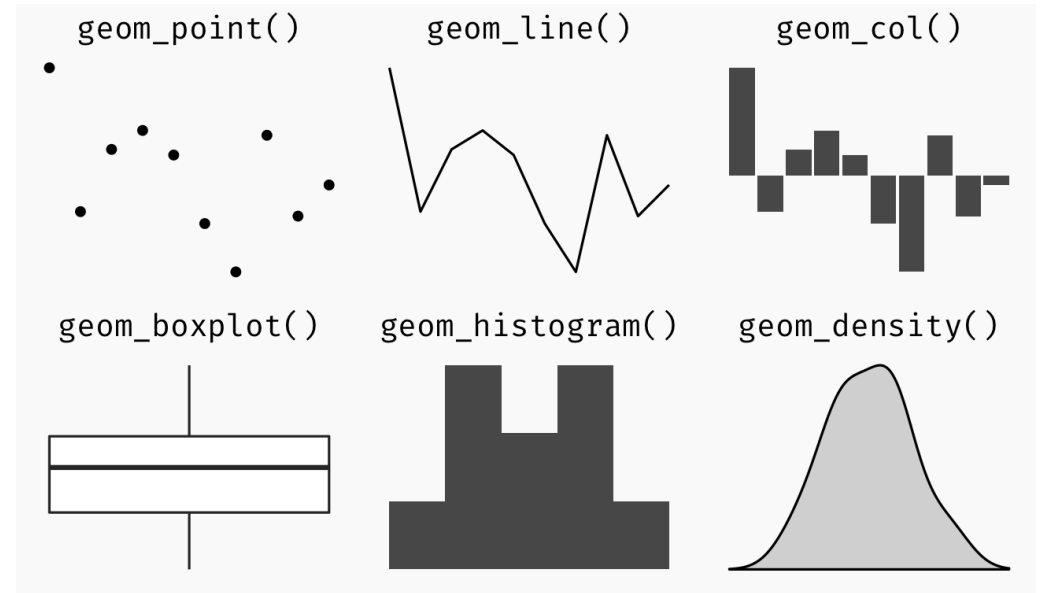
- Notice that ggplot includes the x-axis and the axis's label.



Essential Grammatical Elements in ggplot2

Geometric Objects: *geom*

- Geometric objects are the actual marks we put on a plot. Examples include:
 - *geom_point* for drawing individual points (e.g., a scatter plot)
 - *geom_line* for drawing lines (e.g., for a line charts)
 - *geom_smooth* for drawing smoothed lines (e.g., for simple trends or approximations)
 - *geom_bar* for drawing bars (e.g., for bar charts)
 - *geom_histogram* for drawing binned values (e.g. a histogram)
 - *geom_polygon* for drawing arbitrary shapes
 - *geom_map* for drawing polygons in the shape of a map! (You can access the data to use for these maps by using the `map_data()` function).



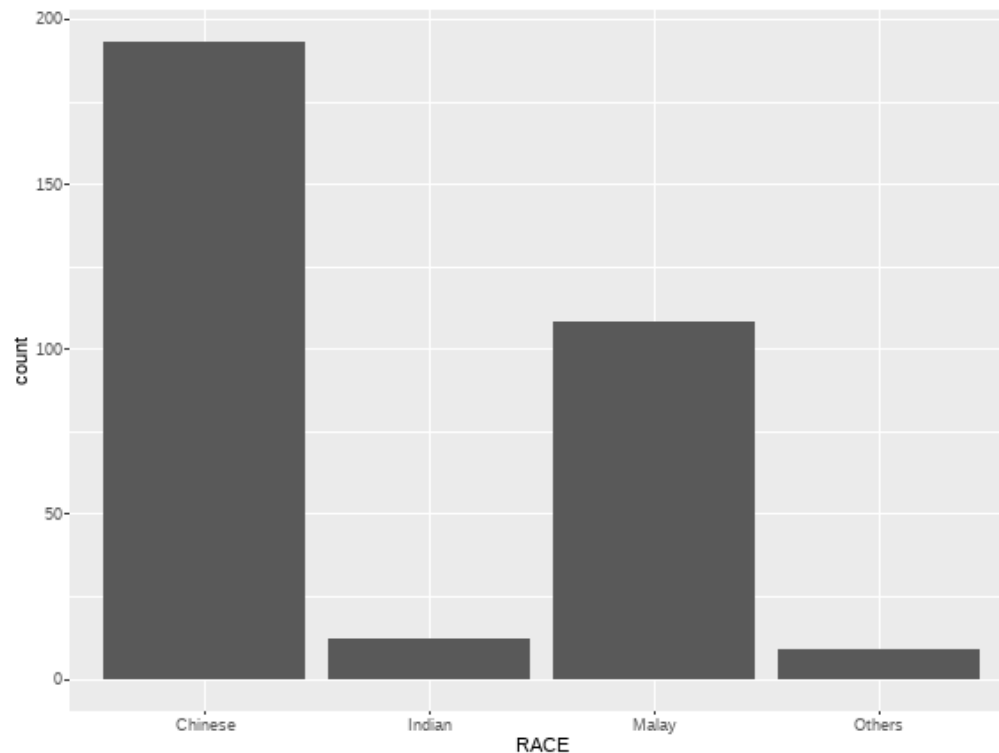
- A plot must have at least one geom; there is no upper limit. You can add a geom to a plot using the `+` operator.
- For complete list, please refer to [here](#).

Essential Grammatical Elements in ggplot2

Geometric Objects: *geom_bar*

- The code chunk below plots a bar chart.

```
ggplot(data=exam_data,  
       aes(x=RACE)) +  
  geom_bar()
```



Essential Grammatical Elements in ggplot2

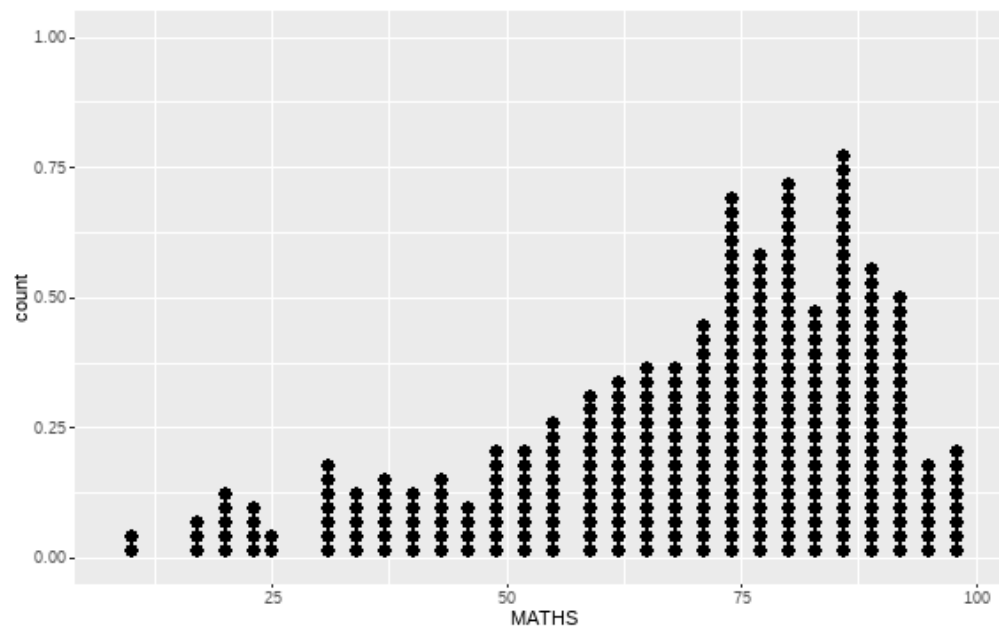
Geometric Objects: *geom_dotplot*

In a dot plot, the width of a dot corresponds to the bin width (or maximum width, depending on the binning algorithm), and dots are stacked, with each dot representing one observation.

```
ggplot(data=exam_data,  
       aes(x = MATHS)) +  
  geom_dotplot(dotsize = 0.5)
```

Note that y scale is not very useful, in fact it is very misleading.

`stat_bindot()` using `bins = 30`. Pick better val



Reference: Dot plot (statistics) [https://en.wikipedia.org/wiki/Dot_plot_\(statistics\)](https://en.wikipedia.org/wiki/Dot_plot_(statistics))

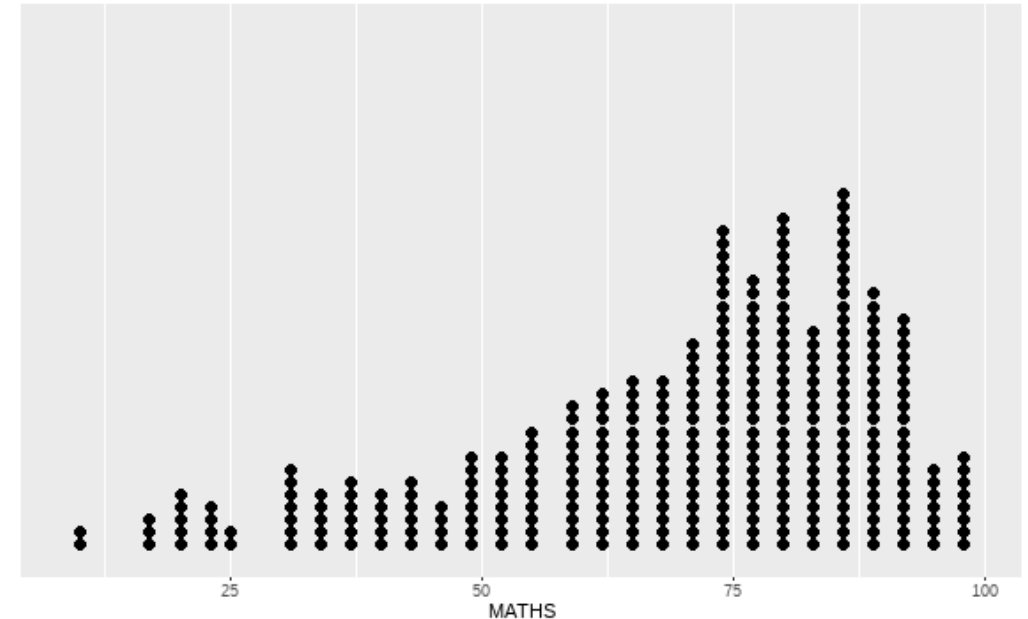
Essential Grammatical Elements in ggplot2

Geometric Objects: *geom_dotplot*

The code chunk below performs the following two steps:

- *scale_y_continuous()* is used to turn off the y-axis, and
- *binwidth* argument is used to change the binwidth to 2.5.

```
ggplot(data=exam_data,  
       aes(x = MATHS)) +  
  geom_dotplot(binwidth=2.5,  
              dotsize = 0.5) +  
  scale_y_continuous(NULL,  
                    breaks = NULL)
```



Essential Grammatical Elements in ggplot2

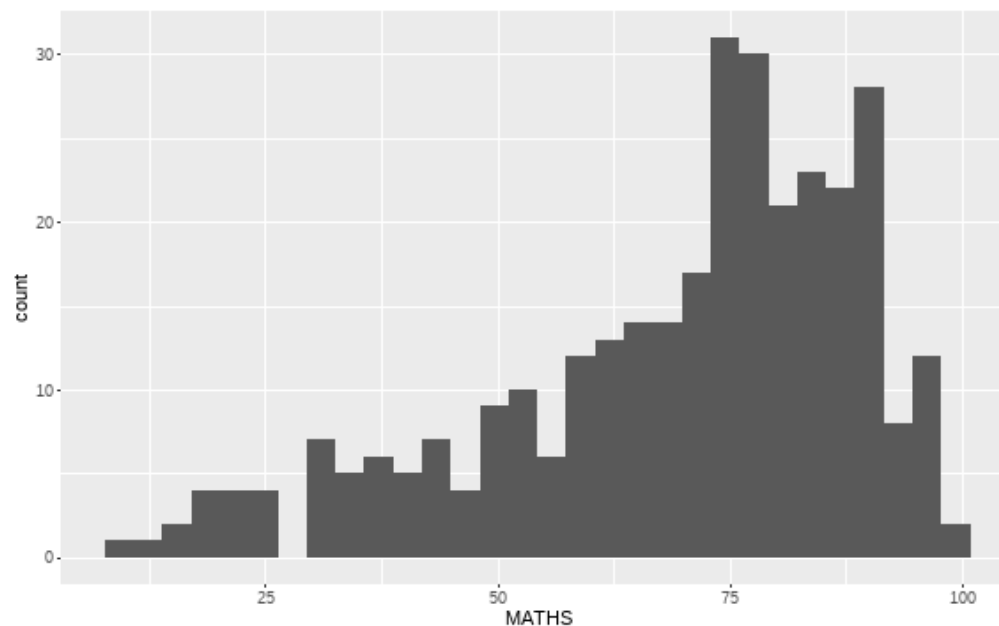
Geometric Objects: *geom_histogram*

In the code chunk below, *geom_histogram()* is used to a simple histogram by using values in *MATHS* field of *exam_data*.

```
ggplot(data=exam_data,  
       aes(x = MATHS)) +  
  geom_histogram()
```

- Note that the default bin is 30.

`stat_bin()` using `bins = 30`. Pick better value



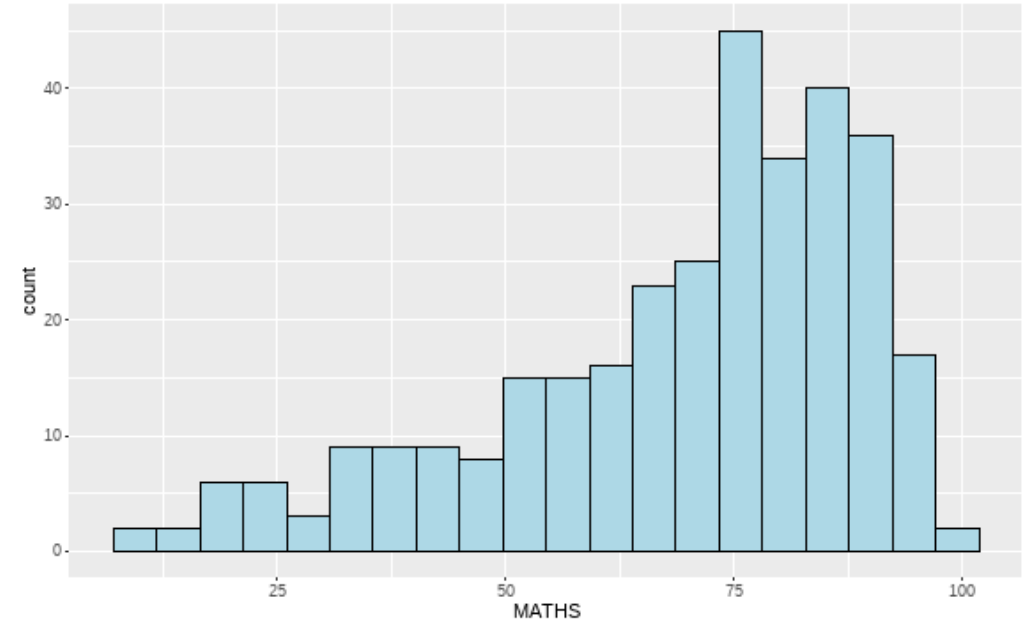
Essential Grammatical Elements in ggplot2

Modifying a geometric object by changing *geom()*

In the code chunk below,

- *bins* argument is used to change the number of bins to 20,
- *fill* argument is used to shade the histogram with light blue color, and
- *color* argument is used to change the outline colour of the bars in black

```
ggplot(data=exam_data,  
       aes(x= MATHS)) +  
  geom_histogram(bins=20,  
                color="black",  
                fill="light blue")
```



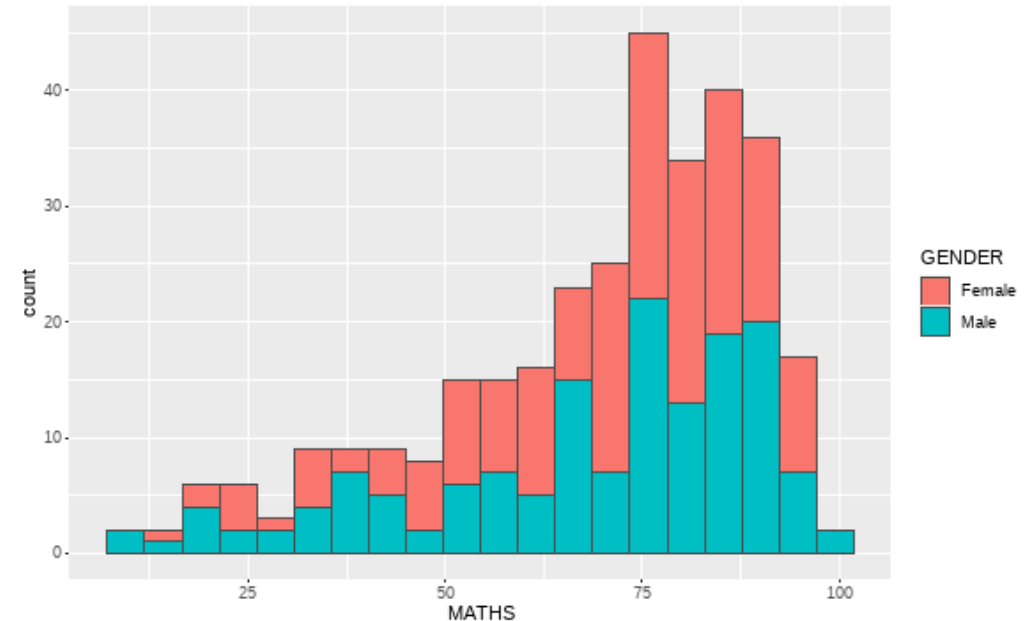
Essential Grammatical Elements in ggplot2

Modifying a geometric object by changing *aes()*

- The code chunk below changes the interior colour of the histogram (i.e. *fill*) by using subgroup of *aesthetic()*.

```
ggplot(data=exam_data,  
       aes(x= MATHS,  
           fill = GENDER)) +  
  geom_histogram(bins=20,  
                color="grey30")
```

- Note that this approach can be used to colour, fill and alpha of the geometric.



Essential Grammatical Elements in ggplot2

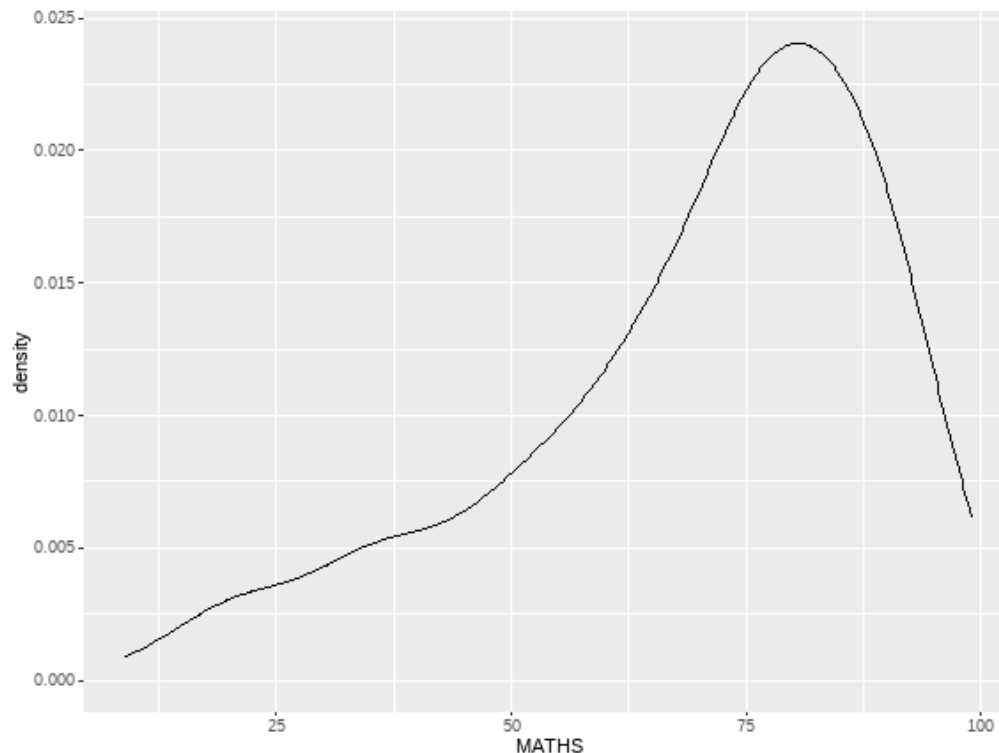
Geometric Objects: *geom-density*

geom-density() computes and plots kernel density estimate, which is a smoothed version of the histogram.

It is a useful alternative to the histogram for continuous data that comes from an underlying smooth distribution.

The code below plots the distribution of Maths scores in a kernel density estimate plot.

```
ggplot(data=exam_data,  
       aes(x = MATHS)) +  
  geom_density()
```



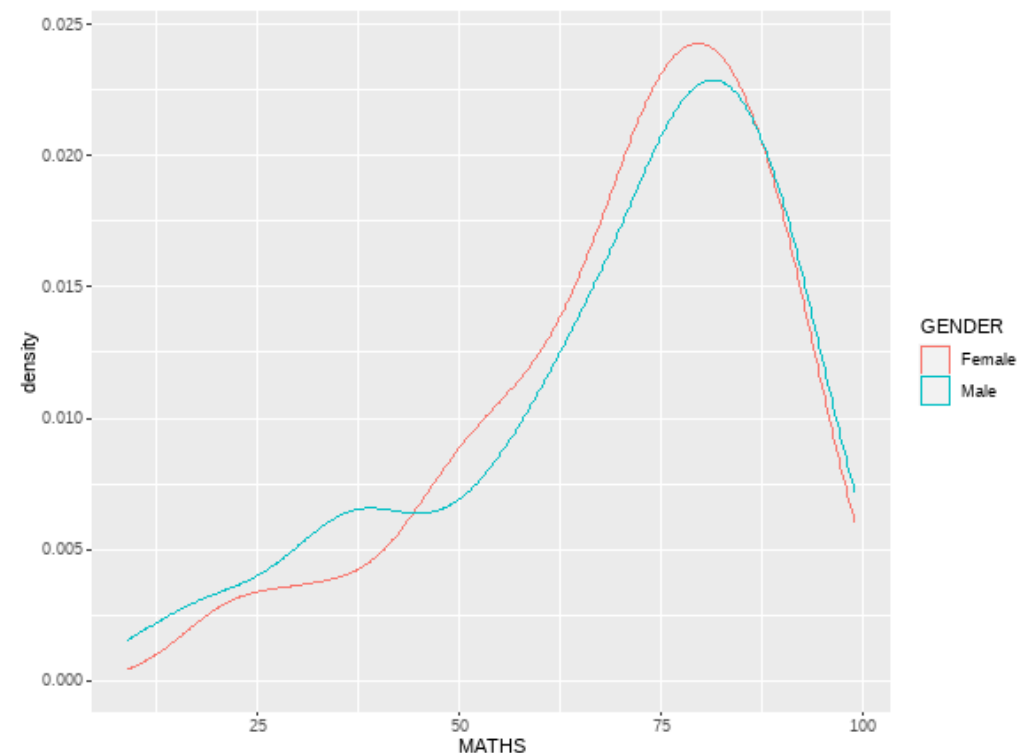
Reference: Kernel density estimation https://en.wikipedia.org/wiki/Kernel_density_estimation

Essential Grammatical Elements in ggplot2

Geometric Objects: *geom-density*

The code chunk below plots two kernel density lines by using *colour* or *fill* arguments of *aes()*

```
ggplot(data=exam_data,  
       aes(x = MATHS,  
           colour = GENDER)) +  
  geom_density()
```

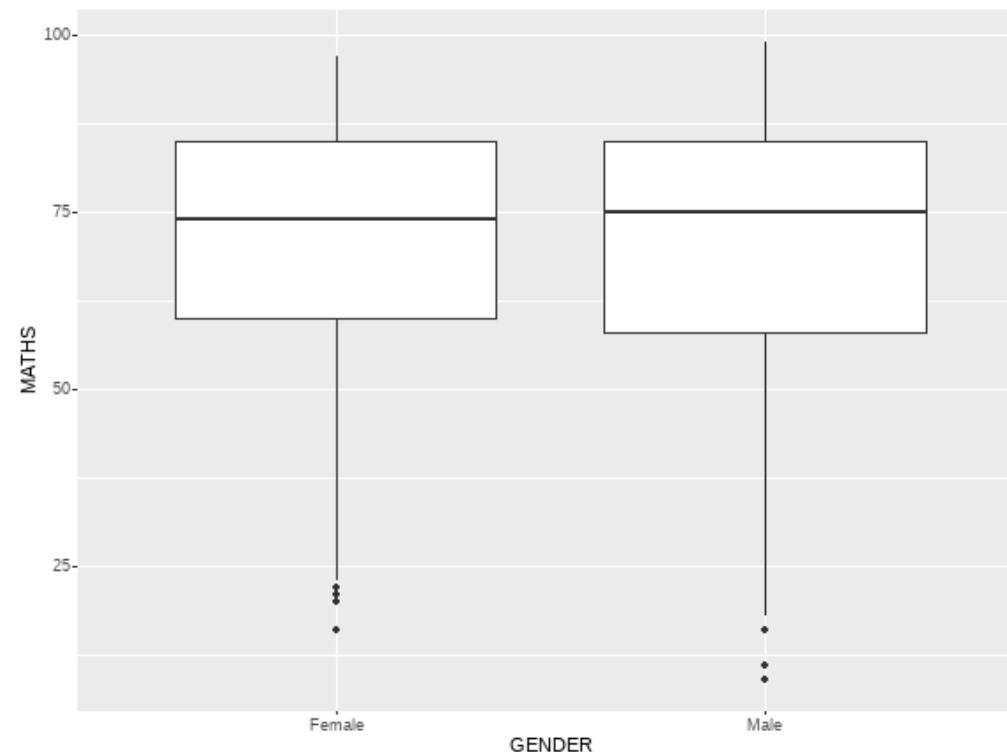


Essential Grammatical Elements in ggplot2

Geometric Objects: *geom_boxplot*

The code chunk below plots boxplots by using *geom_boxplot()*. Notice that two aesthetic properties are used, namely x and y.

```
ggplot(data=exam_data,  
      aes(y = MATHS,  
          x= GENDER)) +  
  geom_boxplot()
```



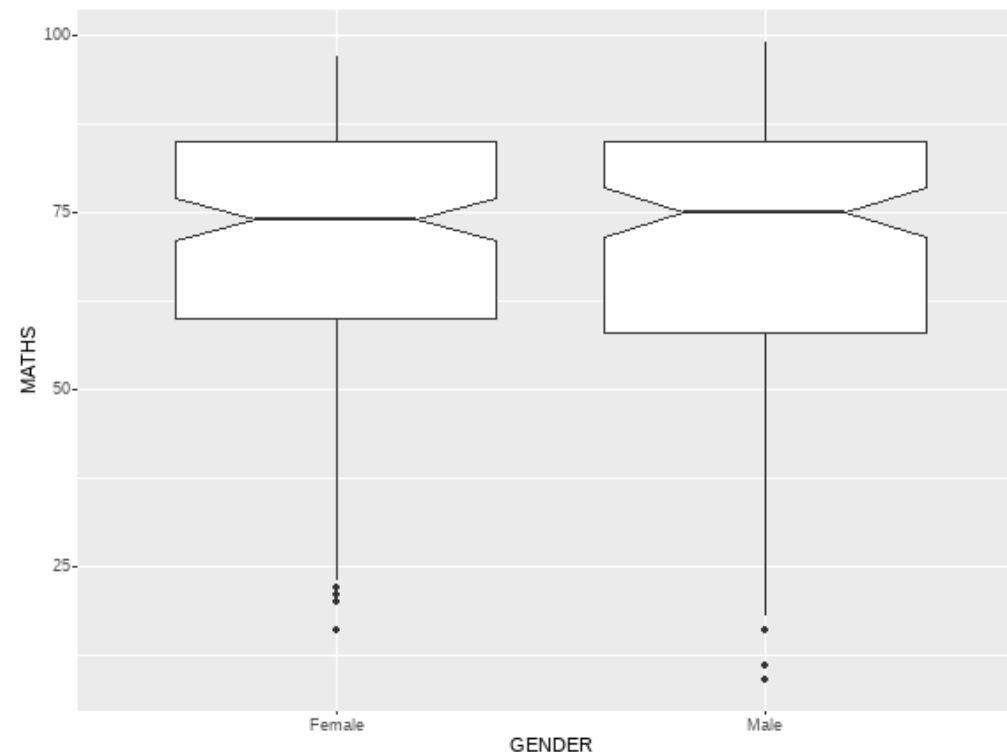
Essential Grammatical Elements in ggplot2

Geometric Objects: *geom_boxplot*

Notches are used in box plots to help visually assess whether the medians of distributions differ. If the notches do not overlap, this is evidence that the medians are different.

The code chunk below plots the distribution of Maths scores by gender in notched plot instead of boxplot.

```
ggplot(data=exam_data,  
       aes(y = MATHS,  
           x= GENDER)) +  
  geom_boxplot(notch=TRUE)
```



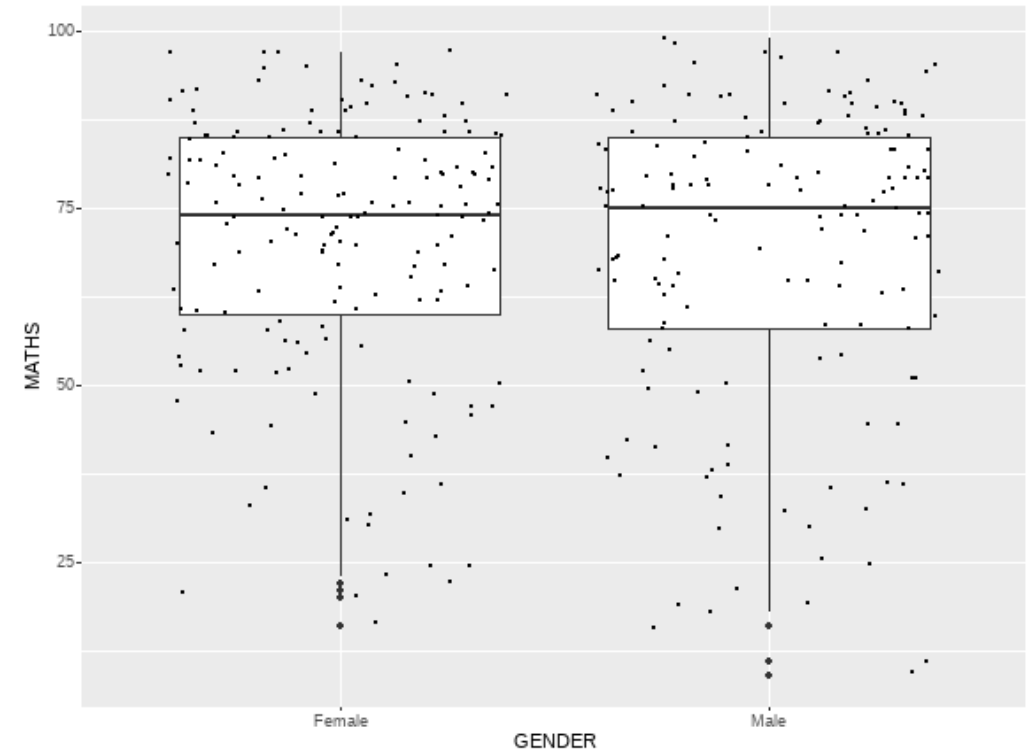
Reference: Notched Box Plots <https://sites.google.com/site/davidsstatistics/home/notched-box-plots>

Essential Grammatical Elements in ggplot2

geom objects can be combined

The code chunk below plots the data points on the boxplots by using both *geom_boxplot()* and *geom_point()*.

```
ggplot(data=exam_data,  
      aes(y = MATHS,  
          x= GENDER)) +  
  geom_boxplot() +  
  geom_point(position="jitter",  
            size = 0.5)
```



Essential Grammatical Elements in ggplot2

Geometric Objects: *geom_violin*

Violin plots are a way of comparing multiple data distributions. With ordinary density curves, it is difficult to compare more than just a few distributions because the lines visually interfere with each other. With a violin plot, it's easier to compare several distributions since they're placed side by side.

The code below plot the distribution of Maths score by gender in violin plot.

```
ggplot(data=exam_data,  
       aes(y = MATHS,  
           x= GENDER)) +  
  geom_violin()
```



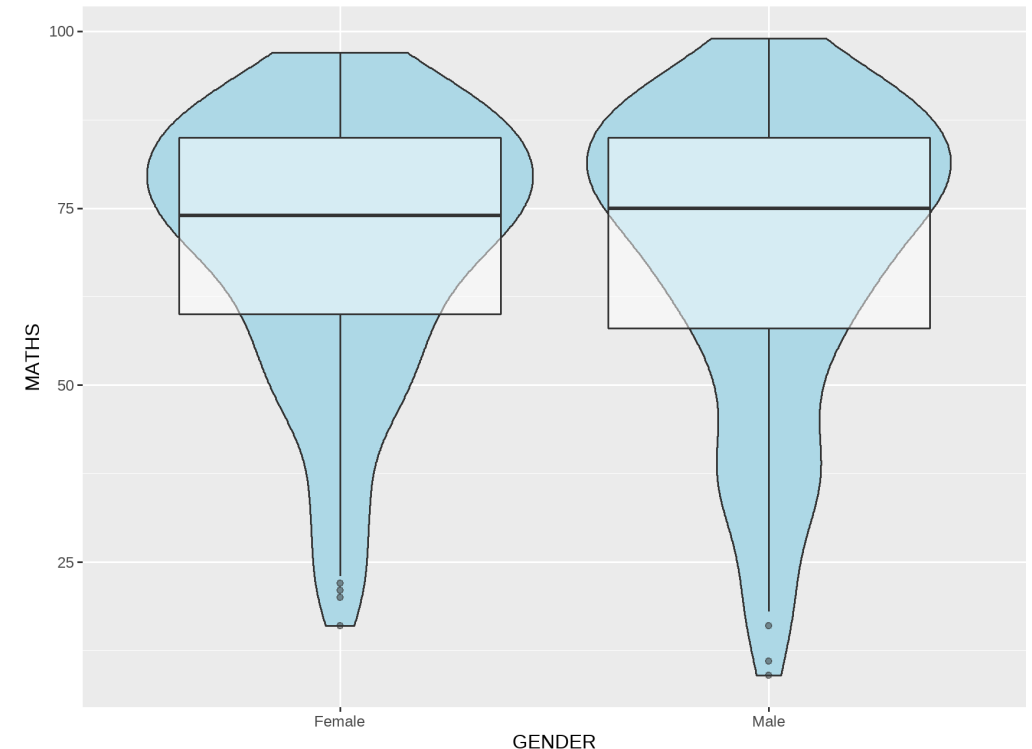
Essential Grammatical Elements in ggplot2

Geometric Objects: *geom_violin()* and *geom_boxplot()*

The code chunk below combined a violin plot and a boxplot to show the distribution of Maths scores by gender.

```
ggplot(data=exam_data,  
       aes(y = MATHS,  
           x= GENDER)) +  
  geom_violin(fill="light blue") +  
  geom_boxplot(alpha=0.5)
```

Notice that the violin plot is plotted first, then follows by the boxplot.

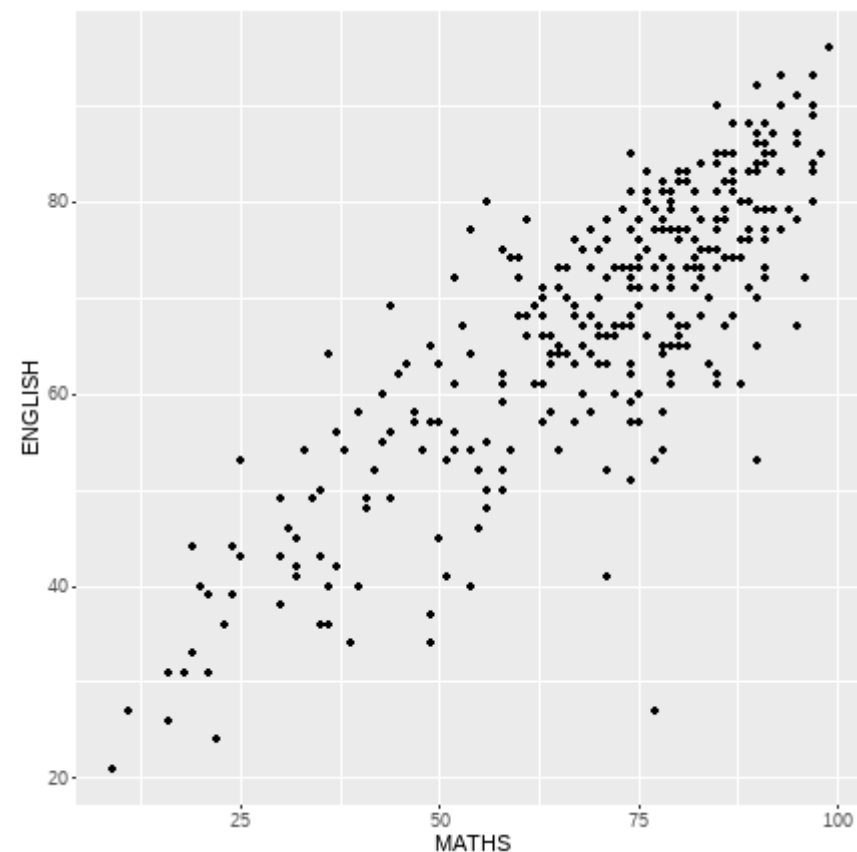


Essential Grammatical Elements in ggplot2

Geometric Objects: *geom_point()*

The code chunk below plots a scatterplot showing the Maths and English grades of pupils by using *geom_point()*.

```
ggplot(data=exam_data,  
       aes(x= MATHS,  
           y=ENGLISH)) +  
  geom_point()
```



Essential Grammatical Elements in ggplot2

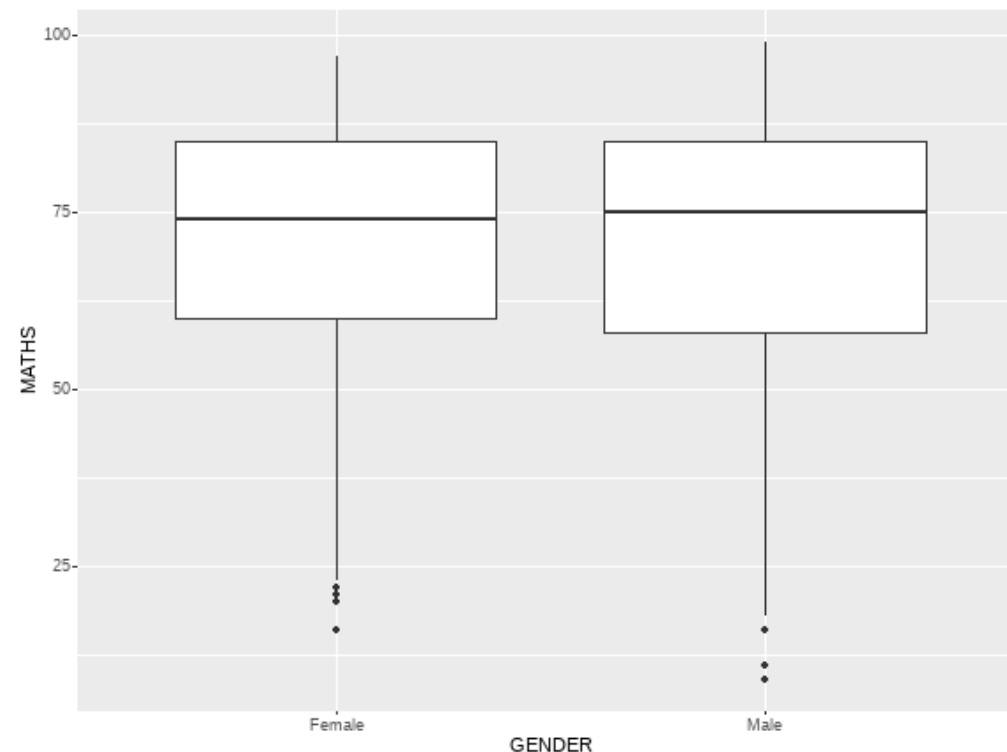
Statistics, *stat*

- The Statistics functions statistically transform data, usually as some form of summary. For example:
 - frequency of values of a variable (bar graph)
 - a mean
 - a confidence limit
- There are two ways to use these functions:
 - add a *stat_()* function and override the default geom, or
 - add a *geom_()* function and override the default stat.

Essential Grammatical Elements in ggplot2

Working with stat

- The boxplots on the right are incomplete because the positions of the means were not shown.
- Next two slides will show you how to add the mean values on the boxplots

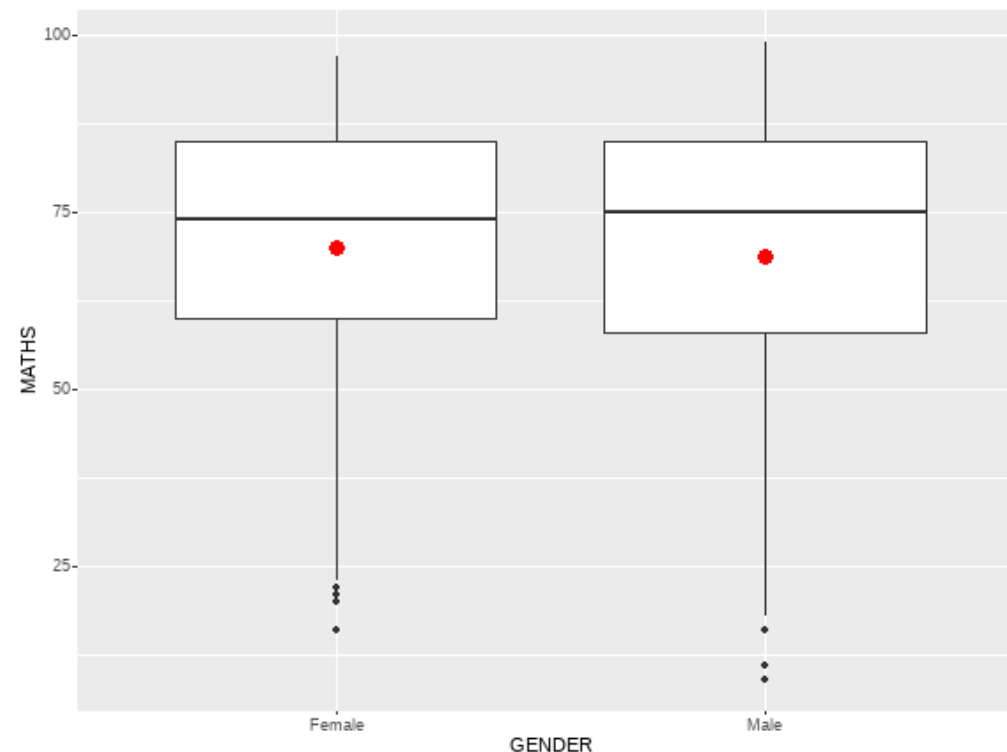


Essential Grammatical Elements in ggplot2

Working with stat - the *stat_summary()* method

The code chunk below adds mean values by using *stat_summary()* function and overriding the default geom.

```
ggplot(data=exam_data,  
      aes(y = MATHS, x= GENDER)) +  
  geom_boxplot() +  
  stat_summary(geom = "point",  
              fun.y="mean",  
              colour = "red",  
              size=4)
```



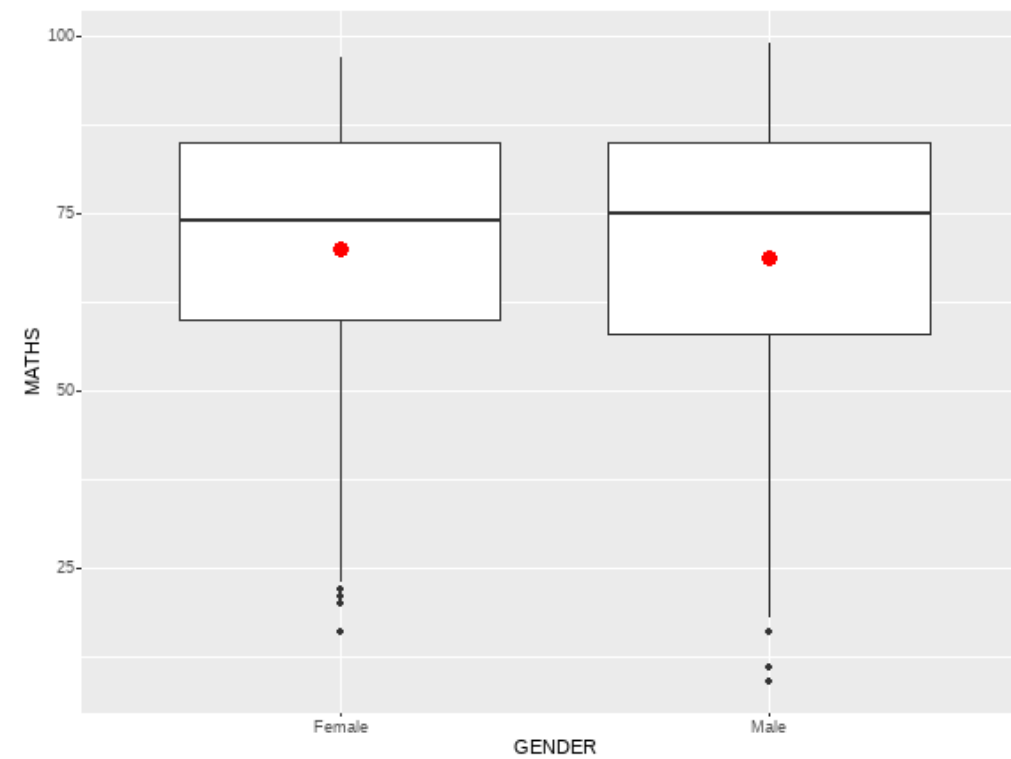
Essential Grammatical Elements in ggplot2

Working with stat - the *geom()* method

The code chunk below adding mean values by using *geom_()* function and overriding the default stat.

```
ggplot(data=exam_data,  
      aes(y = MATHS, x= GENDER)) +  
  geom_boxplot() +  
  geom_point(stat="summary",  
            fun.y="mean",  
            colour ="red",  
            size=4)
```

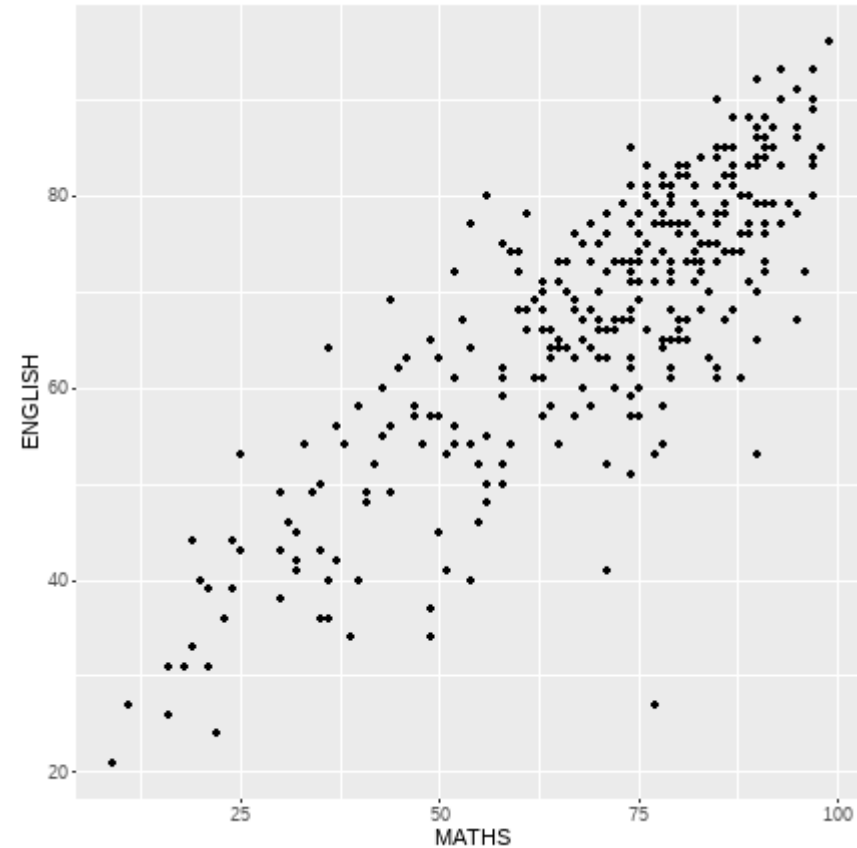
No summary function supplied, defaulting to `mean_`



Essential Grammatical Elements in ggplot2

How to add a best fit curve on a scatterplot?

- The scatterplot on the right shows the relationship of Maths and English grades of pupils.
- The interpretability of this graph can be improved by adding a best fit curve.



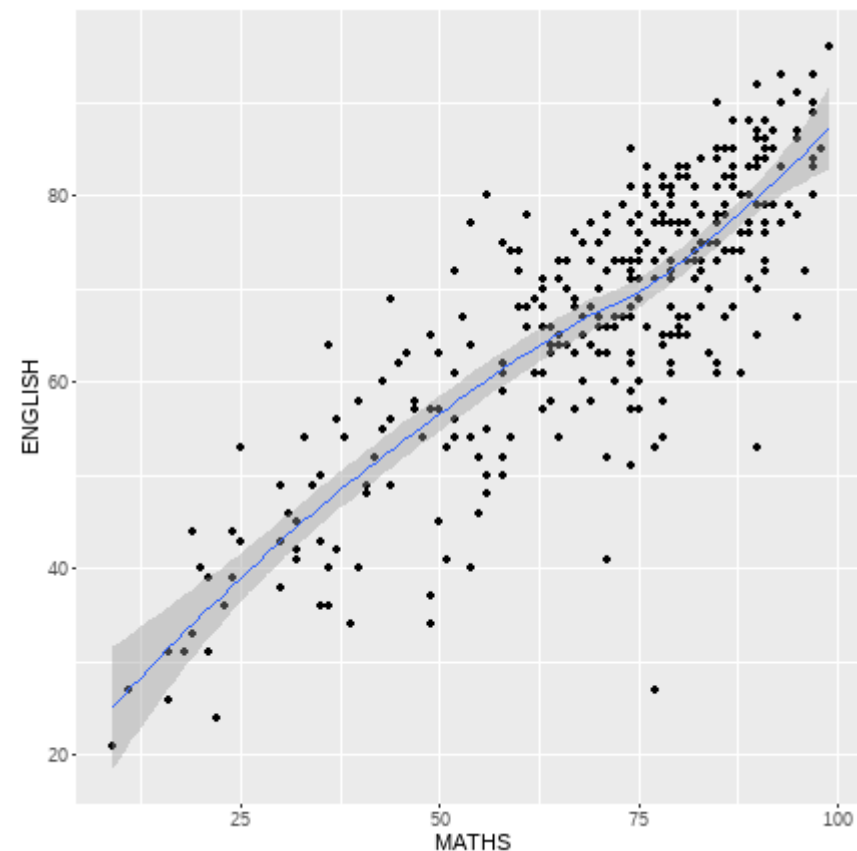
Essential Grammatical Elements in ggplot2

How to add a best fit curve on a scatterplot?

In the code chunk below, *geom_smooth()* is used to plot a best fit curve on the scatterplot.

- The default method used is *loess*.

```
ggplot(data=exam_data,  
       aes(x= MATHS, y=ENGLISH)) +  
  geom_point() +  
  geom_smooth(size=0.5)
```

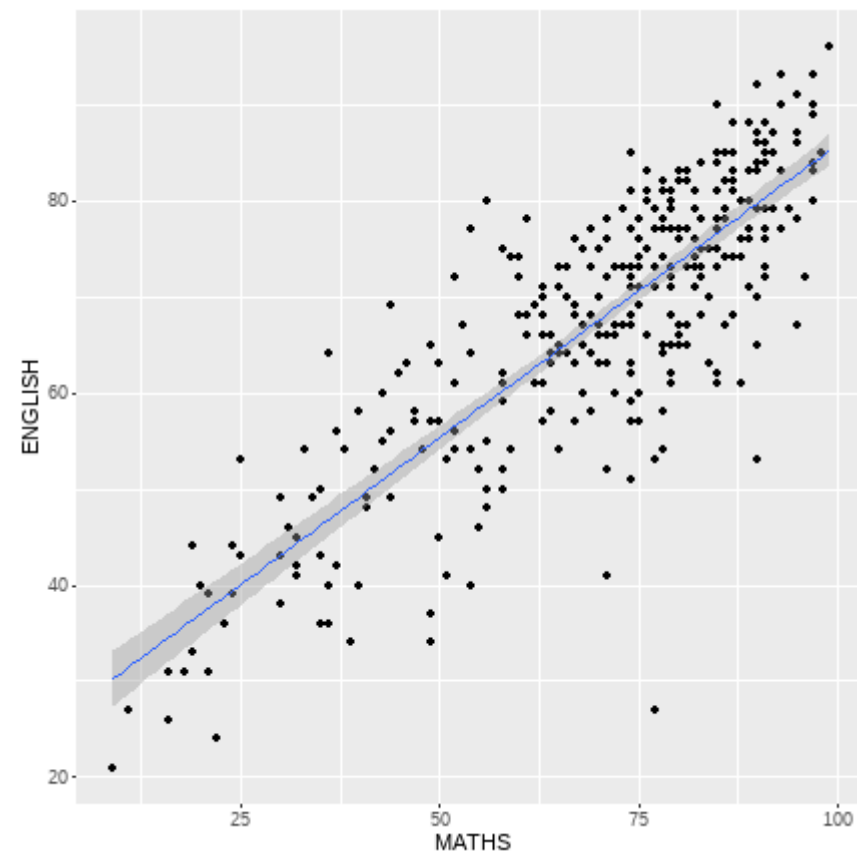


Essential Grammatical Elements in ggplot2

How to add a best fit curve on a scatterplot?

The default smoothing method can be overridden as shown below.

```
ggplot(data=exam_data,  
       aes(x= MATHS,  
           y=ENGLISH)) +  
  geom_point() +  
  geom_smooth(method=lm,  
             size=0.5)
```



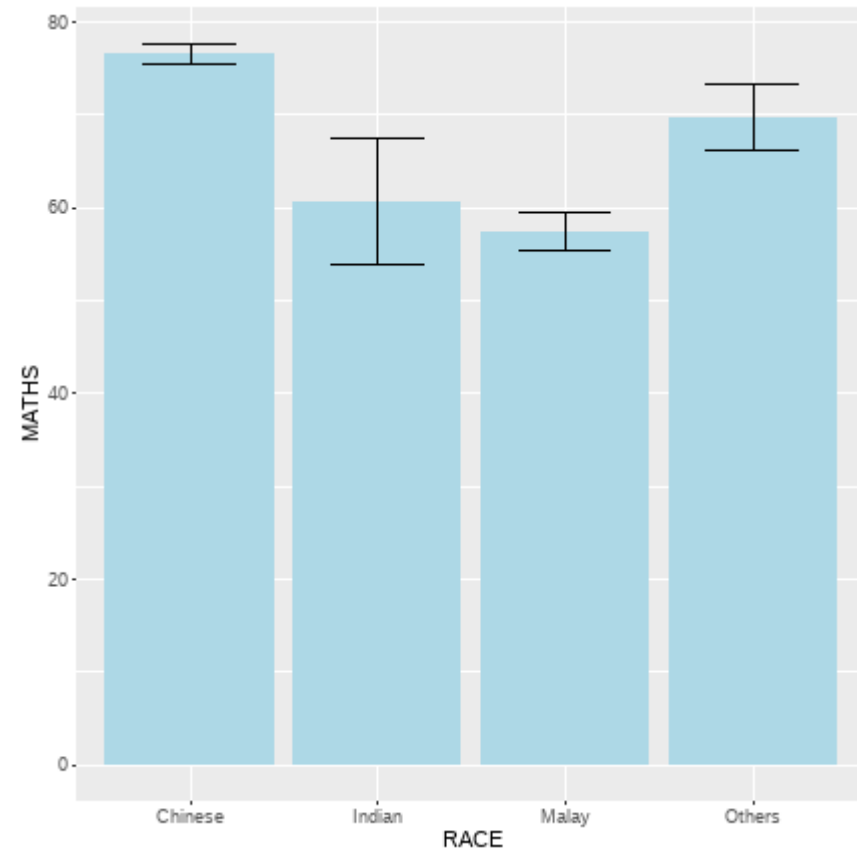
Essential Grammatical Elements in ggplot2

Exercise: Plotting bar chart with error bar using ggplot2

The code chunk below plots the average maths scores by race as bars and their corresponding standard deviation as error bar plot.

```
ggplot(data=exam_data,  
       aes(RACE, MATHS)) +  
  stat_summary(geom = "bar",  
              fun.y = mean,  
              position = "dodge",  
              fill = "light blue") +  
  stat_summary(geom = "errorbar",  
              fun.data = mean_se,  
              position = "dodge",  
              width = 0.50)
```

For more examples, please refer to [Plotting means and error bars \(ggplot2\)](#) and [ggplot2 error bars](#)



Essential Grammatical Elements in ggplot2

Facets

- Facetting generates small multiples (sometimes also called trellis plot), each displaying a different subset of the data.
- Facets are an alternative to aesthetics for displaying additional discrete variables.
- ggplot2 supports two types of faceting, namely: *facet_grid* and *facet_wrap*.

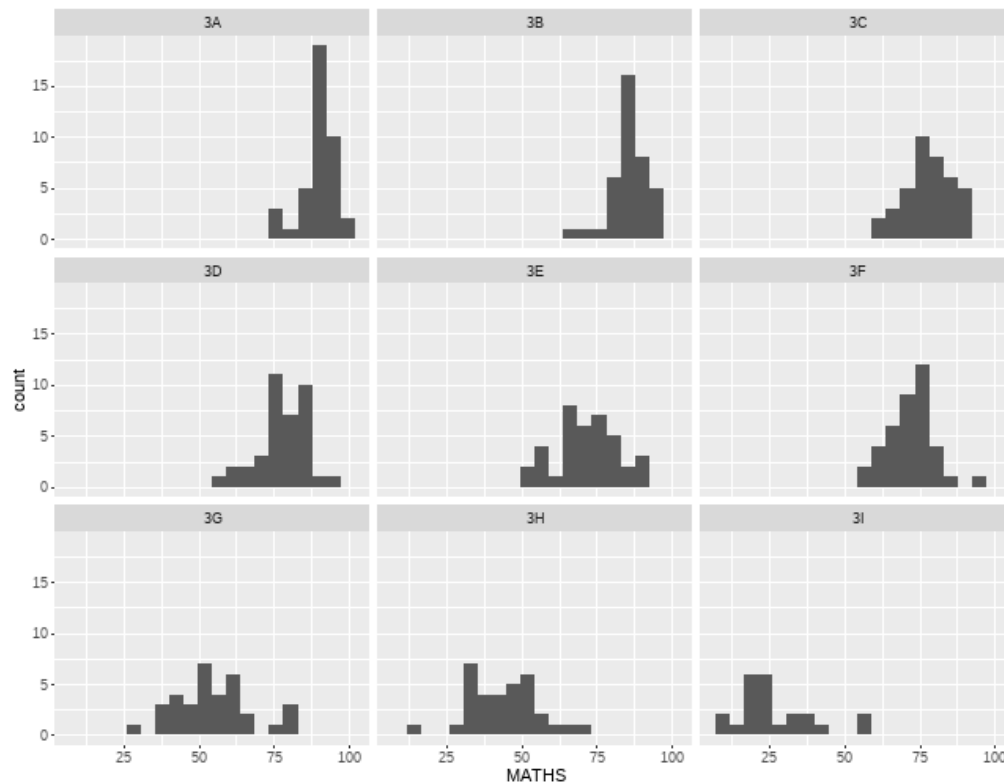
Essential Grammatical Elements in ggplot2

facet_wrap()

facet_wrap wraps a 1d sequence of panels into 2d. This is generally a better use of screen space than *facet_grid* because most displays are roughly rectangular.

The code chunk below plots a trellis plot using *facet_wrap()*.

```
ggplot(data=exam_data,  
       aes(x= MATHS)) +  
  geom_histogram(bins=20) +  
  facet_wrap(~ CLASS)
```



Essential Grammatical Elements in ggplot2

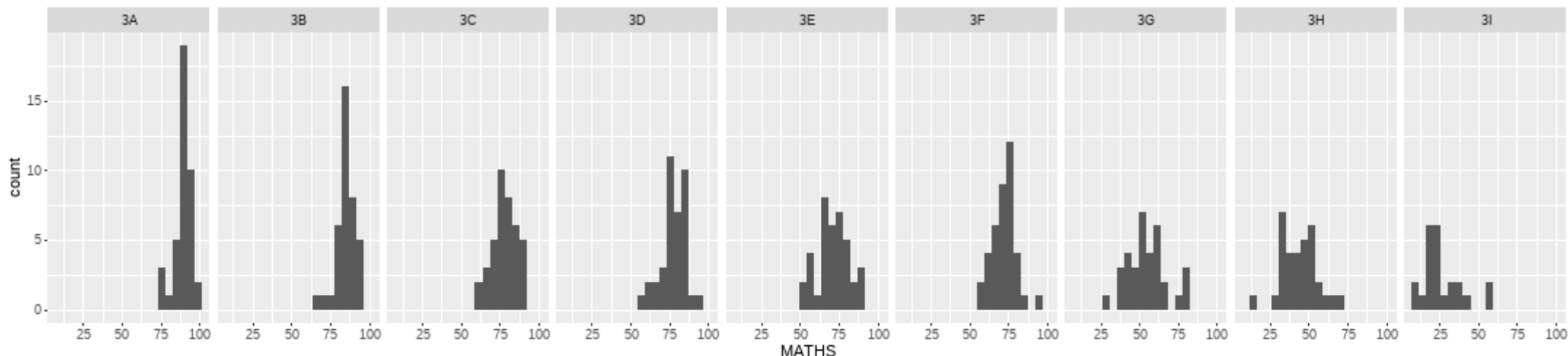
facet_grid() function

facet_grid() forms a matrix of panels defined by row and column facetting variables.

It is most useful when you have two discrete variables, and all combinations of the variables exist in the data.

The code chunk below plots a trellis plot using *facet_grid()*.

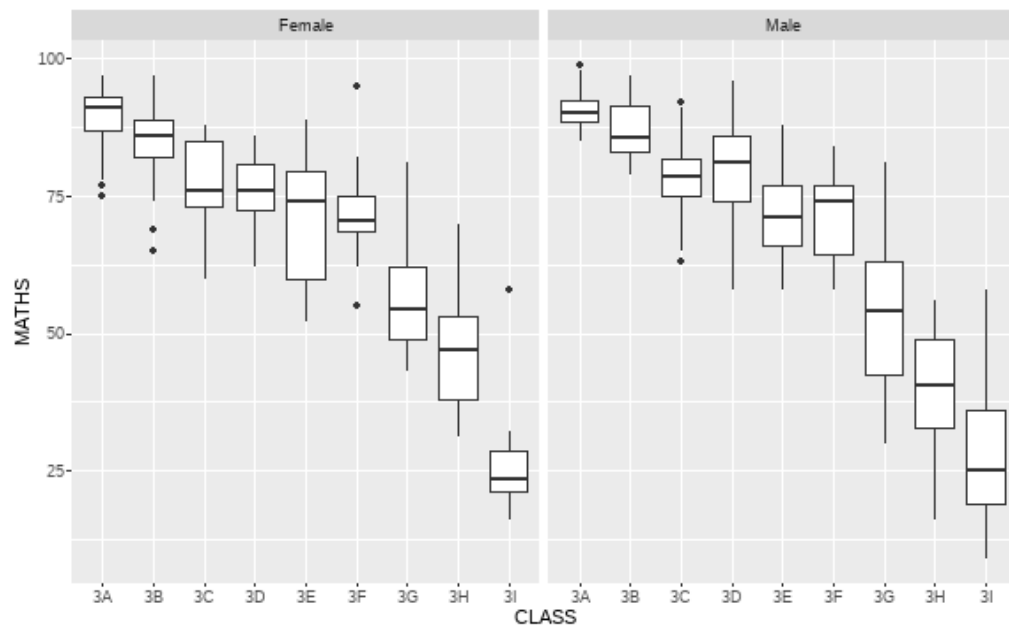
```
ggplot(data=exam_data,  
       aes(x= MATHS)) +  
  geom_histogram(bins=20) +  
  facet_grid(~ CLASS)
```



Facetting

Exercise:

Plot a trellis boxplot looks similar to the figure below.



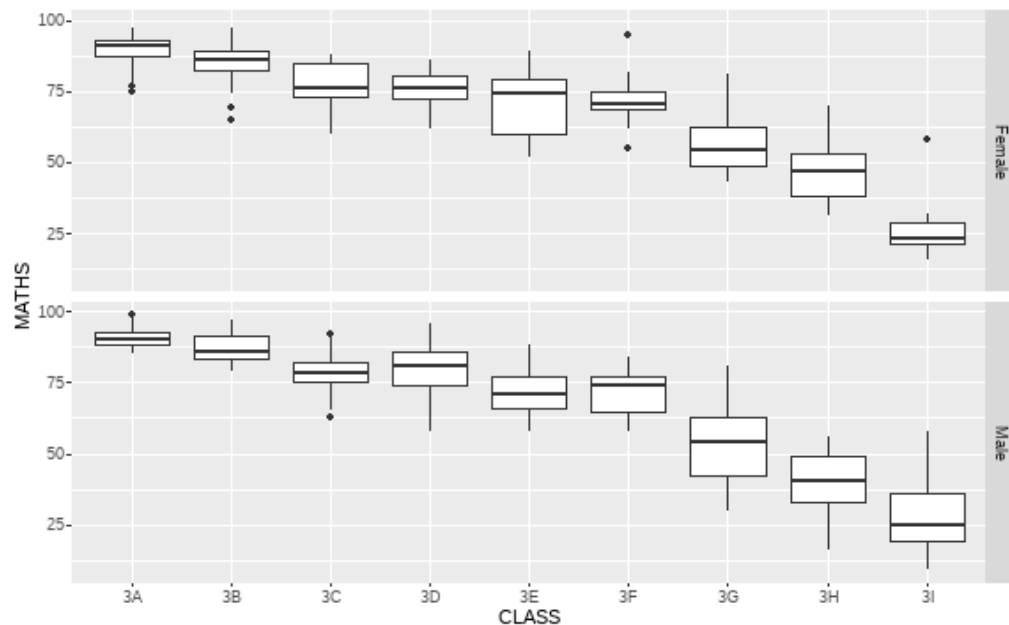
The solution:

```
ggplot(data=exam_data,  
       aes(y = MATHS, x= CLASS)) +  
  geom_boxplot() +  
  facet_grid(~ GENDER)
```

Facetting

Exercise:

Plot a trellis boxplot looks similar to the figure below.



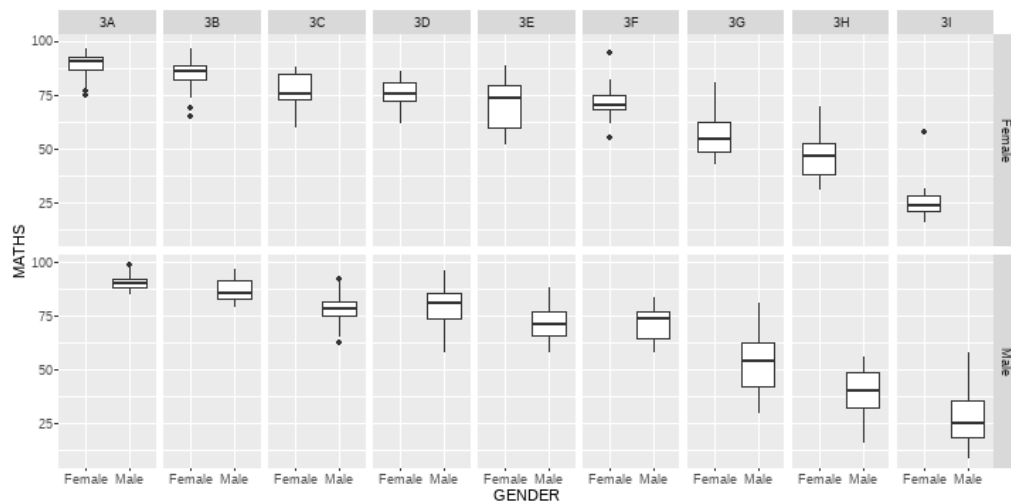
The solution:

```
ggplot(data=exam_data,  
       aes(y = MATHS, x= CLASS)) +  
  geom_boxplot() +  
  facet_grid(GENDER ~.)
```

Facetting

Exercise:

Plot a trellis boxplot looks similar to the figure below.



The solution:

```
ggplot(data=exam_data,  
       aes(y = MATHS, x= GENDER)) +  
  geom_boxplot() +  
  facet_grid(GENDER ~ CLASS)
```

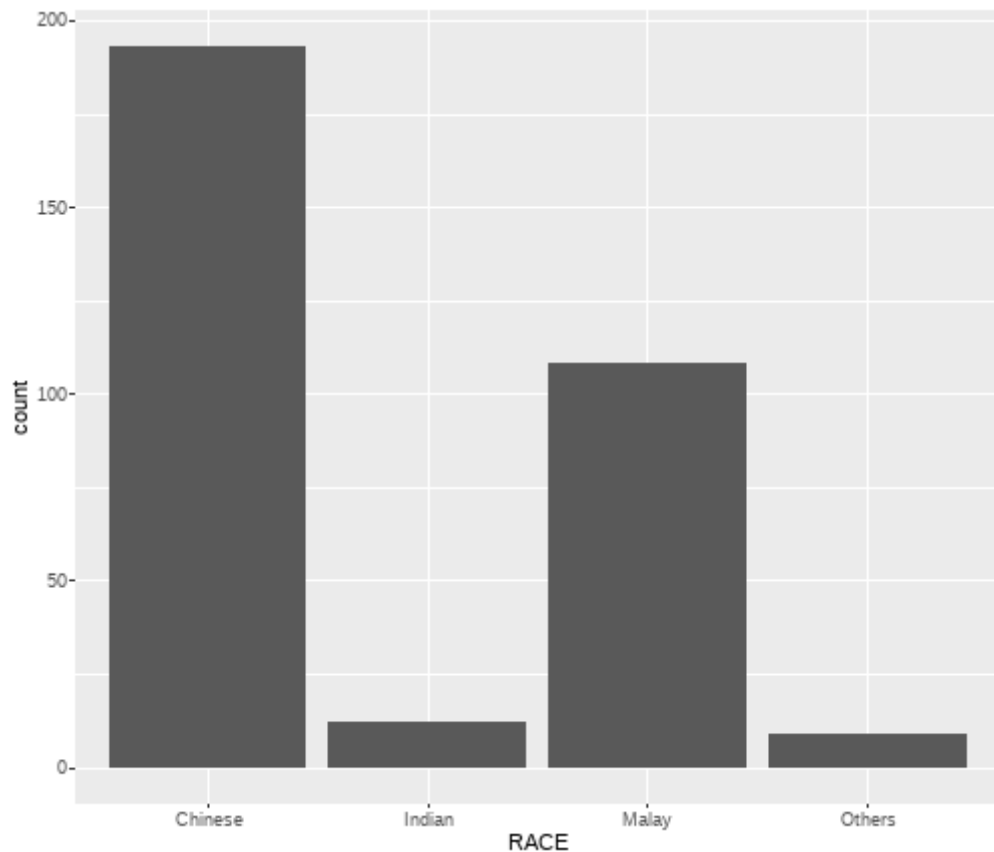
Essential Grammatical Elements in ggplot2

Coordinates

- The *Coordinates* functions map the position of objects onto the plane of the plot.
- There are a number of different possible coordinate systems to use, they are:
 - *coord_cartesian()*: the default cartesian coordinate systems, where you specify x and y values (e.g. allows you to zoom in or out).
 - *coord_flip()*: a cartesian system with the x and y flipped.
 - *coord_fixed()*: a cartesian system with a "fixed" aspect ratio (e.g. 1.78 for a "widescreen" plot).
 - *coord_quickmap()*: a coordinate system that approximates a good aspect ratio for maps.

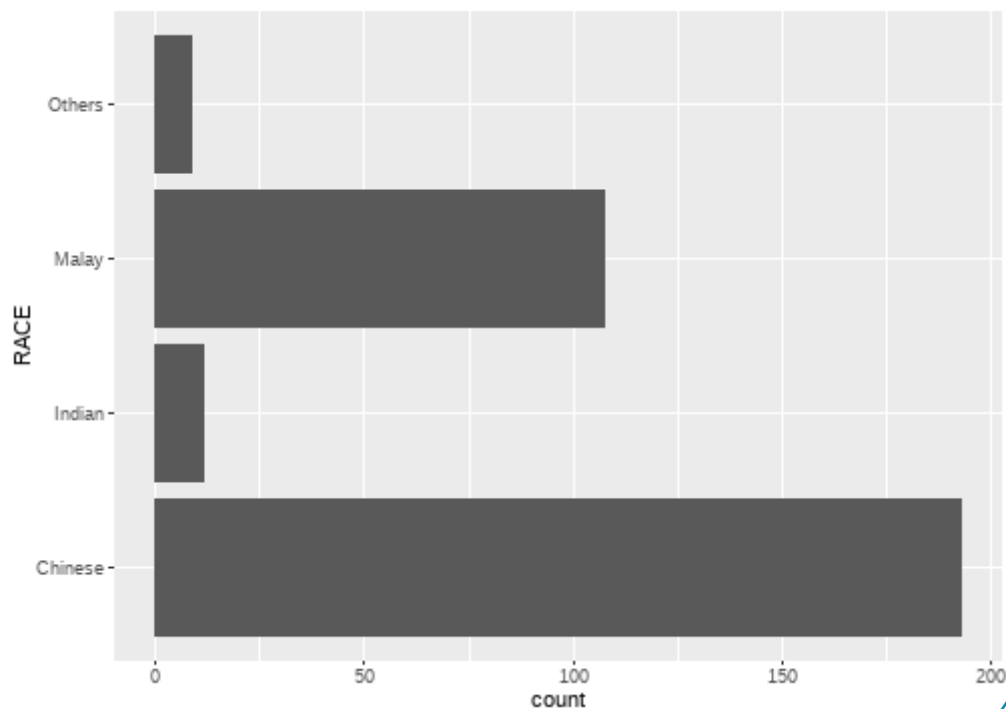
Working with Coordinate

By the default, the bar chart of ggplot2 is in vertical form.



The code chunk below flips the horizontal bar chart into vertical bar chart by using `coord_flip()`.

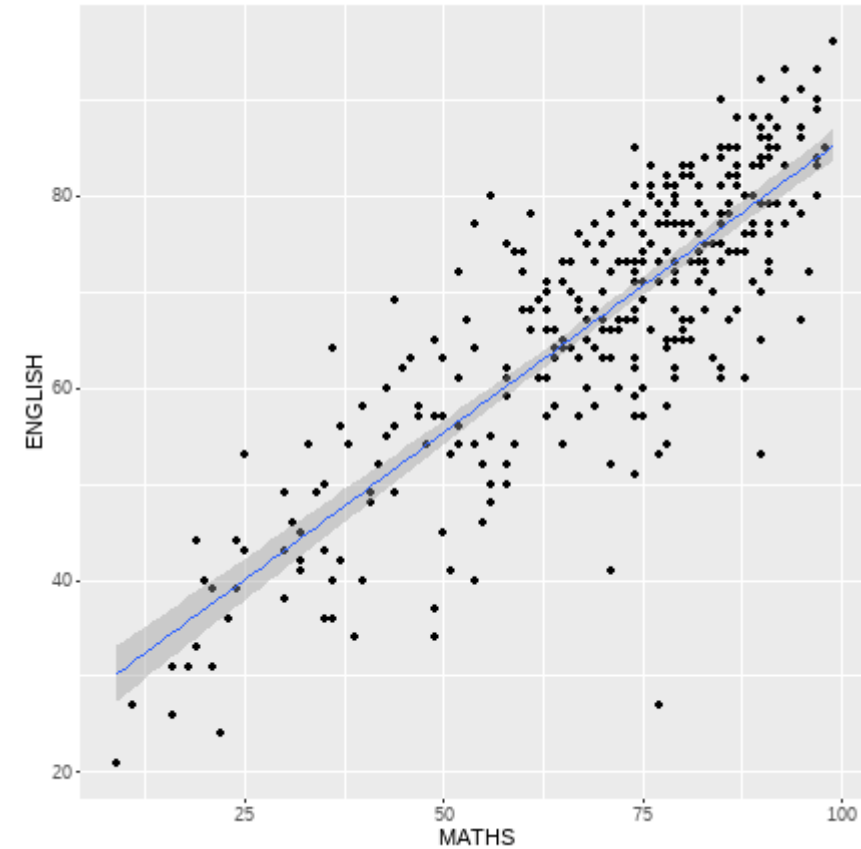
```
ggplot(data=exam_data,  
       aes(x=RACE)) +  
  geom_bar() +  
  coord_flip()
```



Working with Coordinate

Changing the range of y- and x-axis

The scatterplot on the right is slightly misleading because the y-axis and x-axis range are not equal.

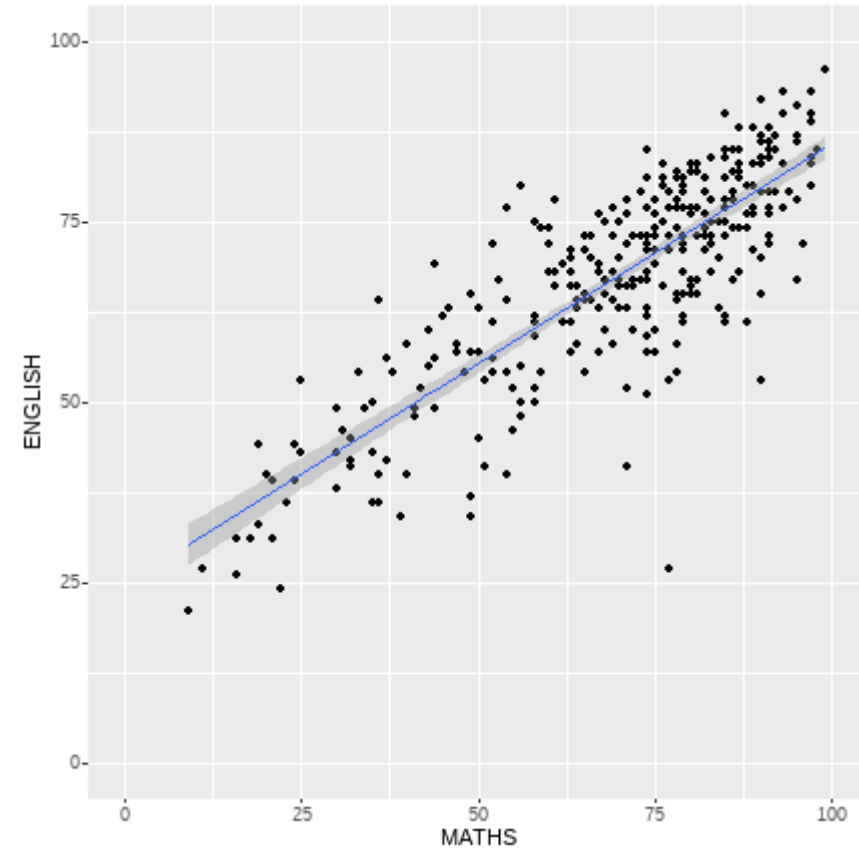


Working with Coordinate

Changing the range of y- and x-axis

The code chunk below fixed both the y-axis and x-axis range from 0-100.

```
ggplot(data=exam_data,  
       aes(x= MATHS, y=ENGLISH)) +  
  geom_point() +  
  geom_smooth(method=lm,  
             size=0.5) +  
  coord_cartesian(xlim=c(0,100),  
                 ylim=c(0,100))
```



Essential Grammatical Elements in ggplot2

Themes

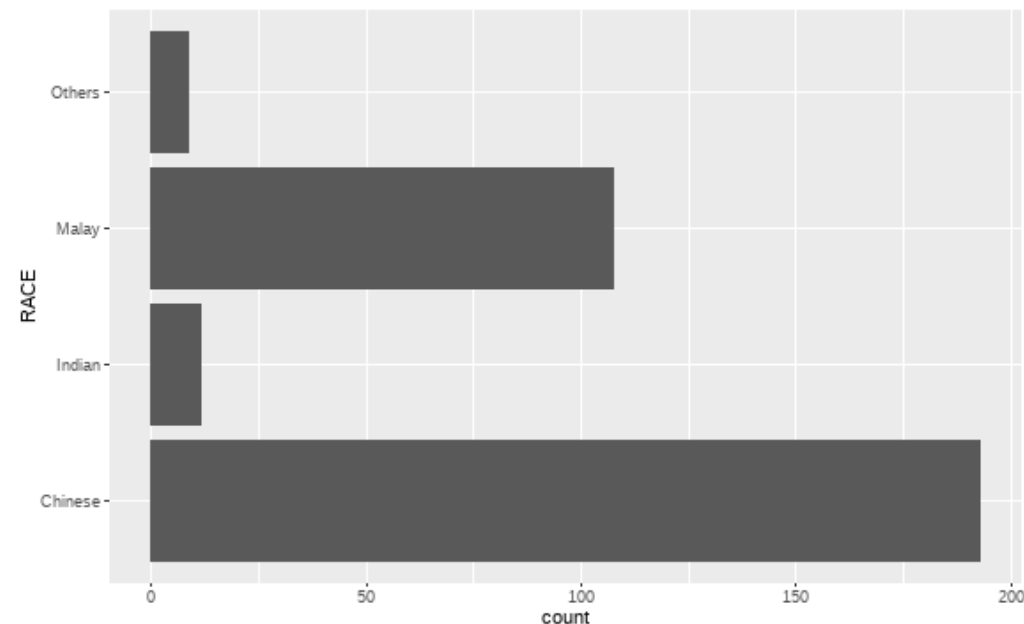
- Themes control elements of the graph not related to the data. For example:
 - background colour
 - size of fonts
 - gridlines
 - colour of labels
- Built-in themes include:
 - *theme_gray()* (default)
 - *theme_bw()*
 - *theme_classic()*
- A list of theme can be found at <http://ggplot2.tidyverse.org/reference/theme.html>.
- Each theme element can be conceived of as either a line (e.g. x-axis), a rectangle (e.g. graph background), or text (e.g. axis title).

Essential Grammatical Elements in ggplot2

Working with theme

The code chunk below plot a horizontal bar chart using *theme_gray()*.

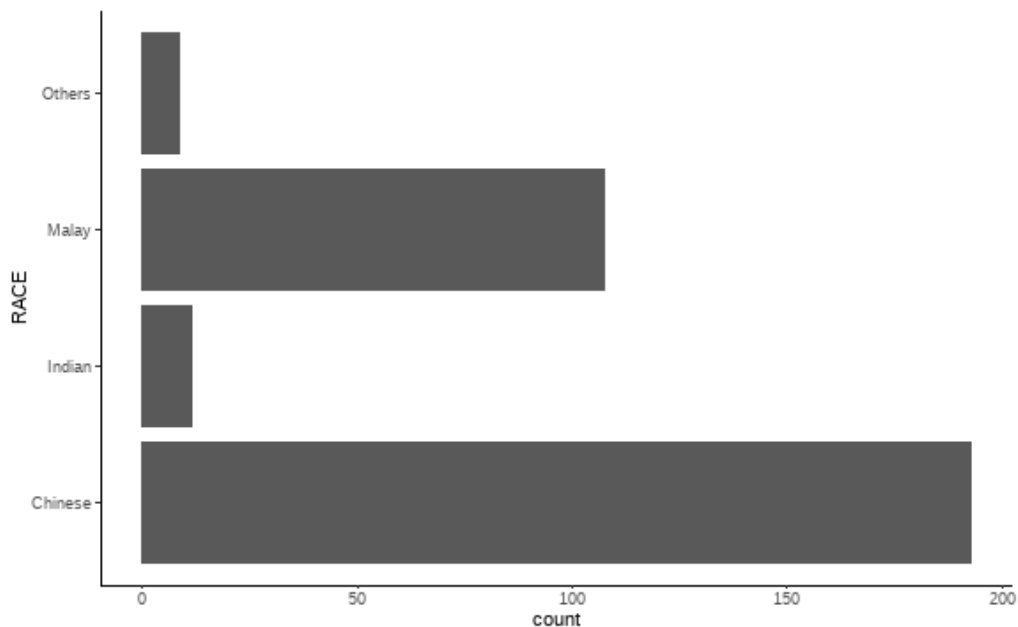
```
ggplot(data=exam_data,  
       aes(x=RACE)) +  
  geom_bar() +  
  coord_flip() +  
  theme_gray()
```



Working with theme

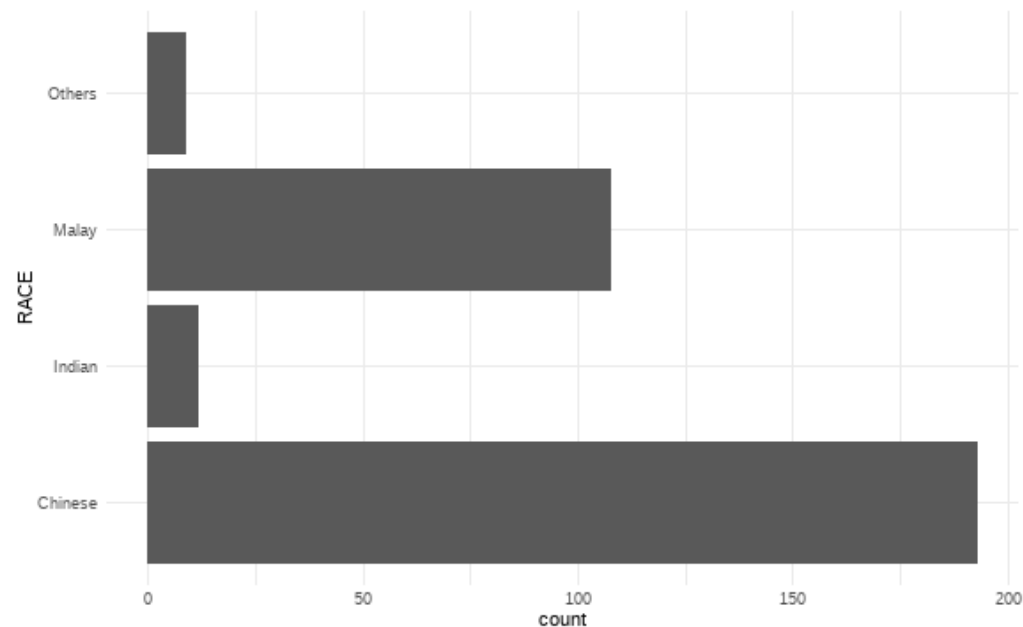
A horizontal bar chart plotted using *theme_classic()*.

```
ggplot(data=exam_data, aes(x=RACE)) +  
  geom_bar() +  
  coord_flip() +  
  theme_classic()
```



A horizontal bar chart plotted using *theme_minimal()*.

```
ggplot(data=exam_data, aes(x=RACE)) +  
  geom_bar() +  
  coord_flip() +  
  theme_minimal()
```

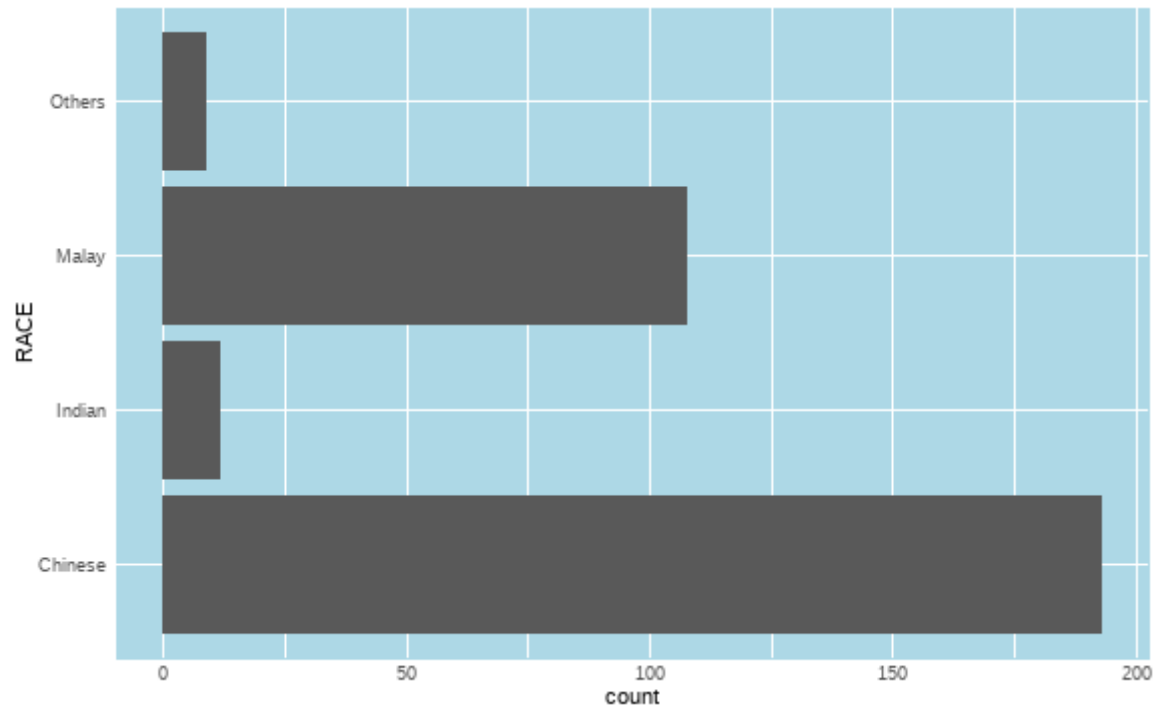


Working with theme

Exercise

Plot a horizontal bar chart looks similar to the figure below.

- Changing the colors of plot panel background of *theme_minimal* to lightblue and the color of grid lines to white.



Working with theme

The solution

```
ggplot(data=exam_data, aes(x=RACE)) +  
  geom_bar() +  
  coord_flip() +  
  theme_minimal() +  
  theme(panel.background = element_rect(fill = "lightblue",  
                                         colour = "lightblue",  
                                         size = 0.5,  
                                         linetype = "solid"),  
        panel.grid.major = element_line(size = 0.5,  
                                         linetype = 'solid',  
                                         colour = "white"),  
        panel.grid.minor = element_line(size = 0.25,  
                                         linetype = 'solid',  
                                         colour = "white"))
```

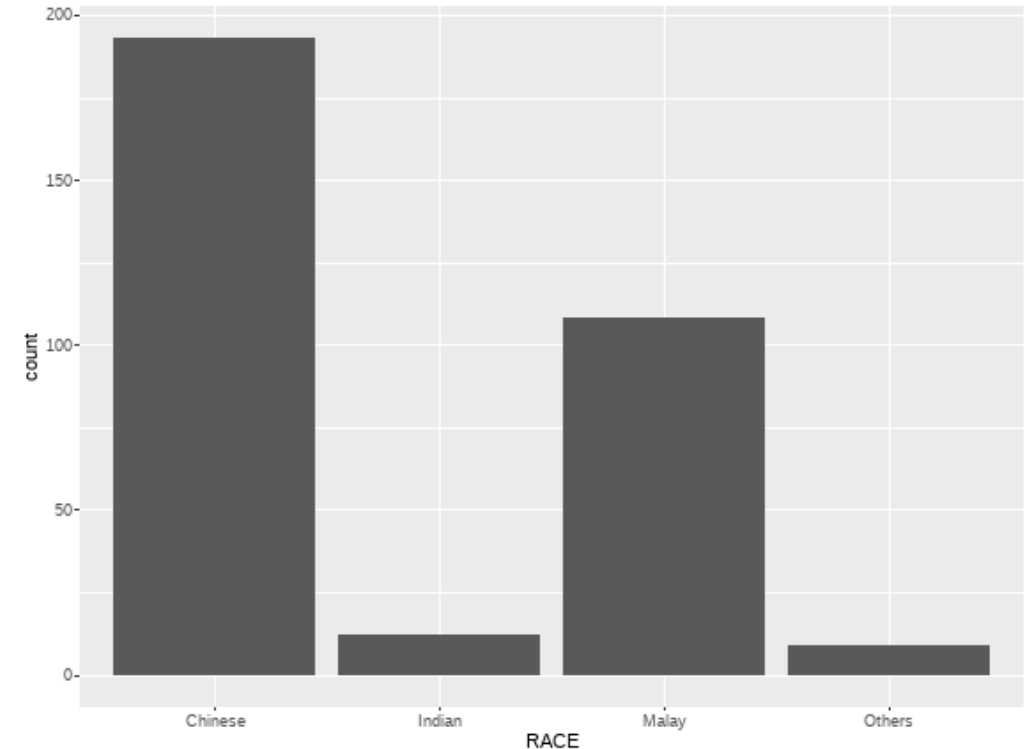
Designing Data-drive Graphics for Analysis I

The original design

A simple vertical bar chart for frequency analysis.

Critics:

- y-axis label is not clear (i.e. count)
- To support effective comparison, the bars should be sorted by their respective frequencies.
- For static graph, frequency values should be added to provide additional information.

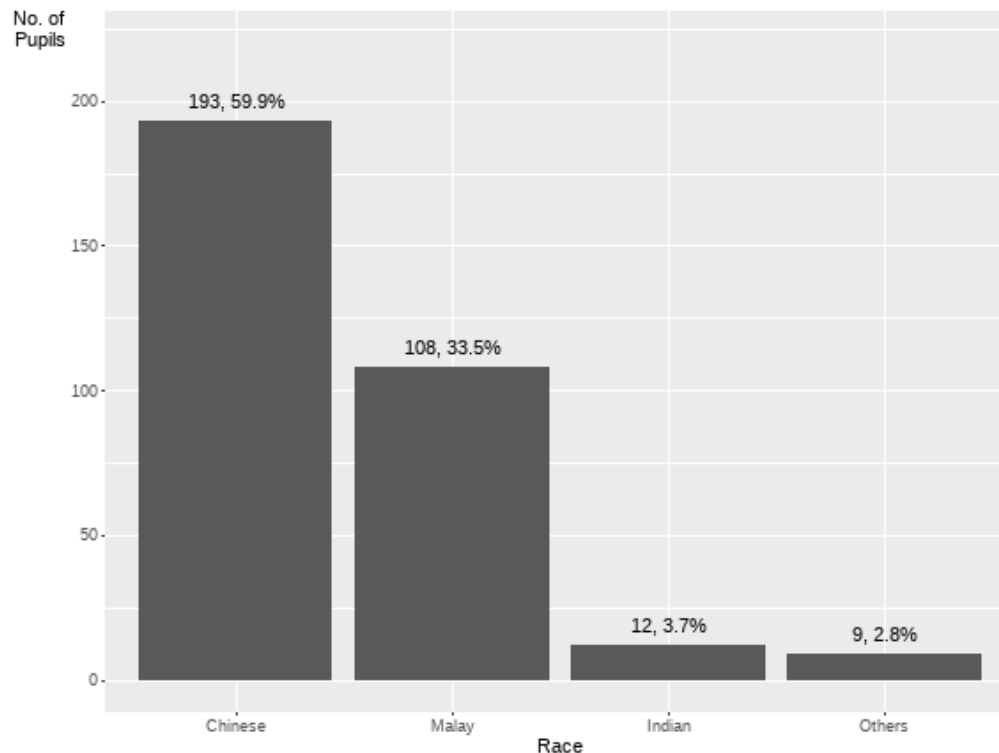


Designing Data-drive Graphics for Analysis I

The makeover design

The code chunk.

```
ggplot(data=exam_data,  
      aes(x=reorder(RACE,RACE,  
                    function(x)-length(x))))+  
  geom_bar() +  
  ylim(0,220) +  
  geom_text(stat="count",  
    aes(label=paste0(..count.., ", ",  
    round(..count../sum(..count..)*100,  
          1), "%")),  
    vjust=-1) +  
  xlab("Race") +  
  ylab("No. of\nPupils") +  
  theme(axis.title.y=element_text(angle = 0))
```

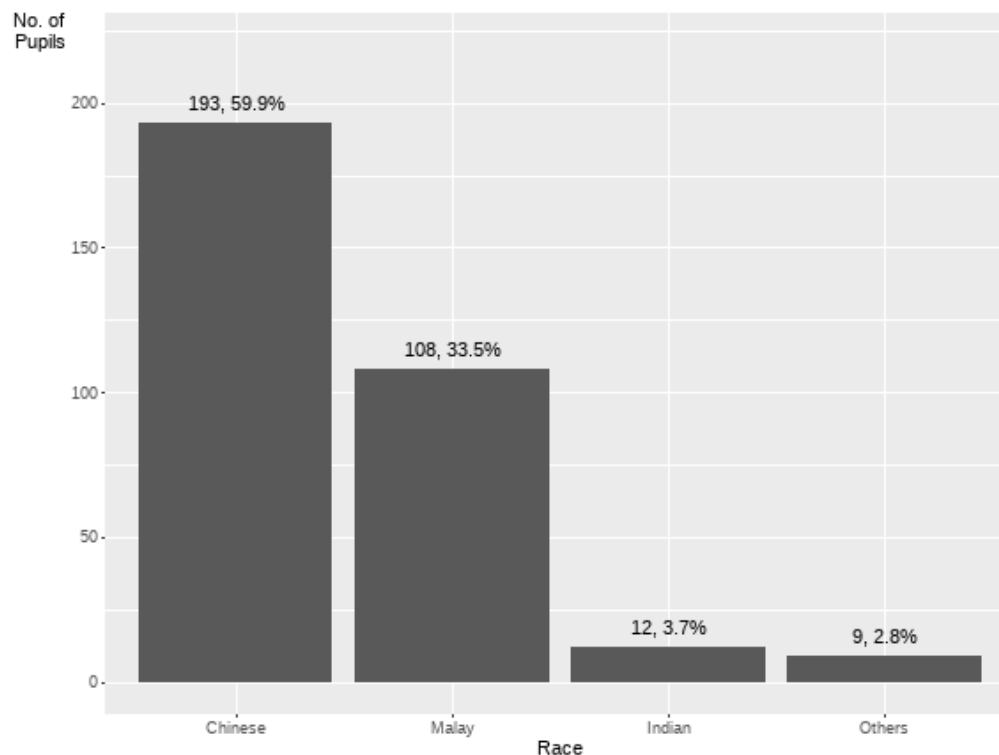


Designing Data-drive Graphics for Analysis I

The makeover design

This code chunk uses `fct_infreq()` of **forcats** package.

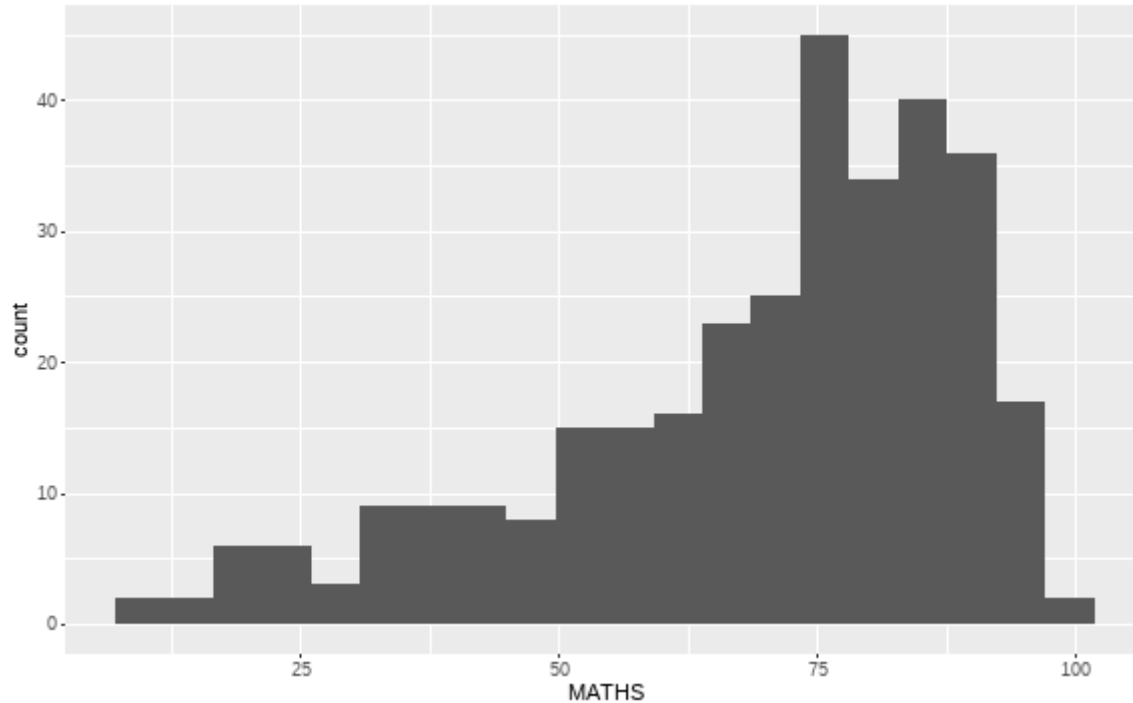
```
exam_data %>%  
  mutate(RACE = fct_infreq(RACE)) %>%  
  ggplot(aes(x = RACE)) +  
  geom_bar()+  
  ylim(0,220) +  
  geom_text(stat="count",  
    aes(label=paste0(..count.., ", ",  
      round(..count../sum(..count..)*100,  
        1), "%")),  
    vjust=-1) +  
  xlab("Race") +  
  ylab("No. of\nPupils") +  
  theme(axis.title.y=element_text(angle = 0))
```



Credit: I learned this trick from [Getting things into the right order](#) of Prof. Claus O. Wilke, the author of **Fundamentals of Data Visualization**.

Designing Data-drive Graphics for Analysis II

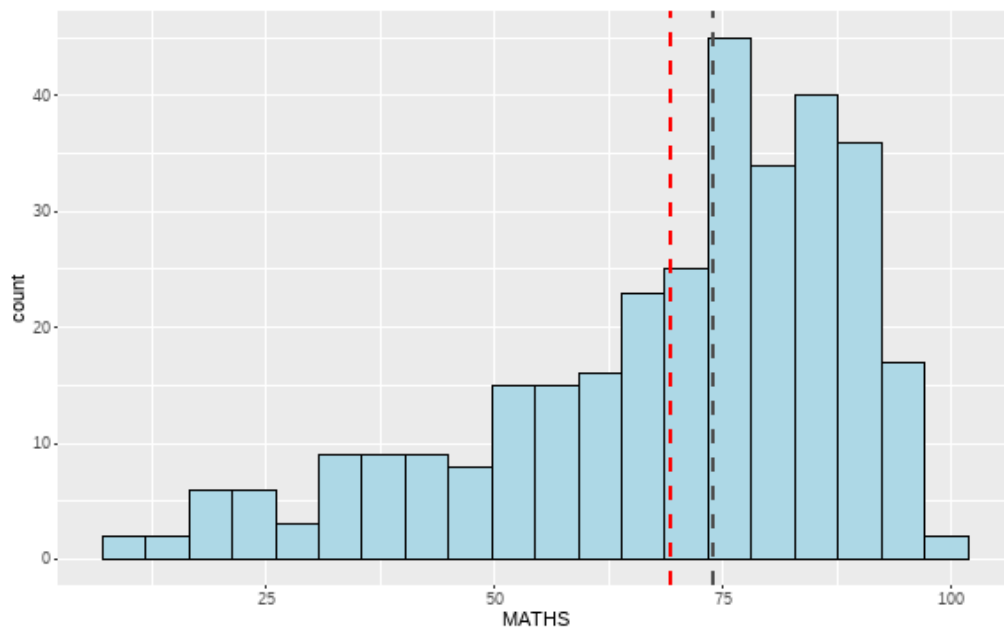
The original design



Designing Data-drive Graphics for Analysis II

The makeover design

- Adding mean and median lines on the histogram plot.
- Change fill color and line color

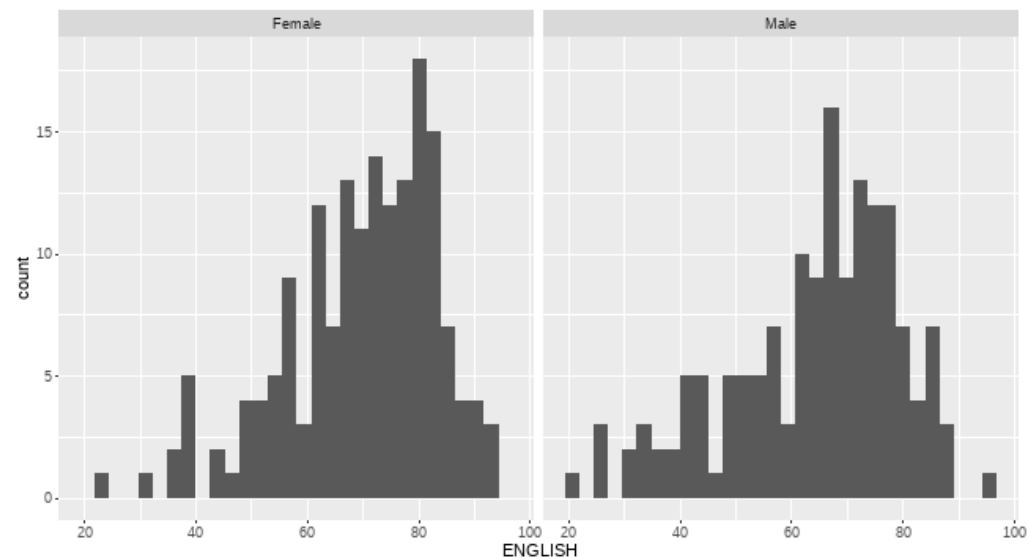


The code chunk

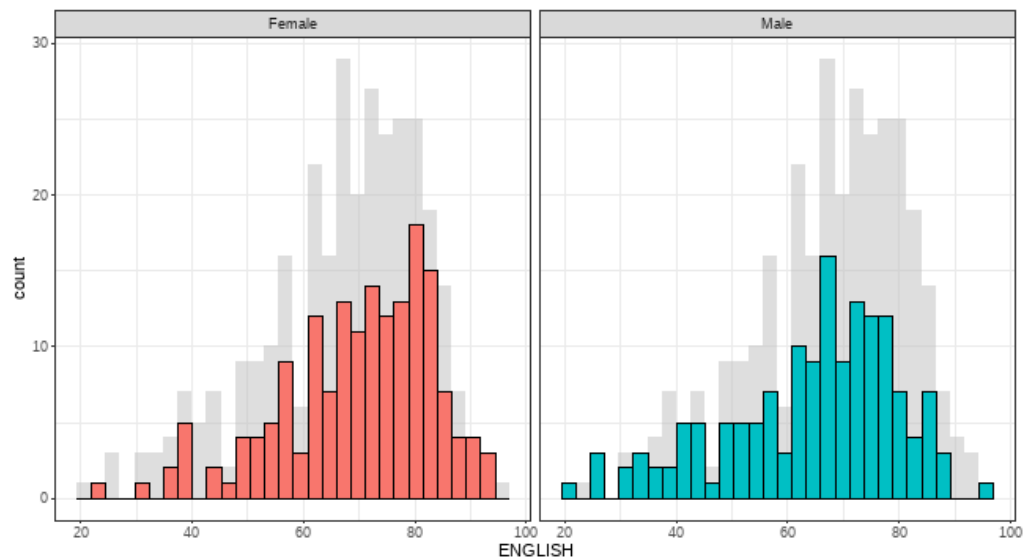
```
ggplot(data=exam_data,  
       aes(x= MATHS)) +  
  geom_histogram(bins=20,  
                color="black",  
                fill="light blue") +  
  geom_vline(aes(xintercept=mean(MATHS,  
                                na.rm=T)),  
            color="red",  
            linetype="dashed",  
            size=1) +  
  geom_vline(aes(xintercept=median(MATHS,  
                                  na.rm=T)),  
            color="grey30",  
            linetype="dashed",  
            size=1)
```

Designing Data-drive Graphics for Analysis III

The histograms on the left are elegantly designed but not informative. This is because they only reveal the distribution of English scores by gender but without context such as all pupils.



The makeover histograms on the right are not only elegantly designed but also informative. This is because they reveal the distribution of English scores by gender with reference to all pupils.



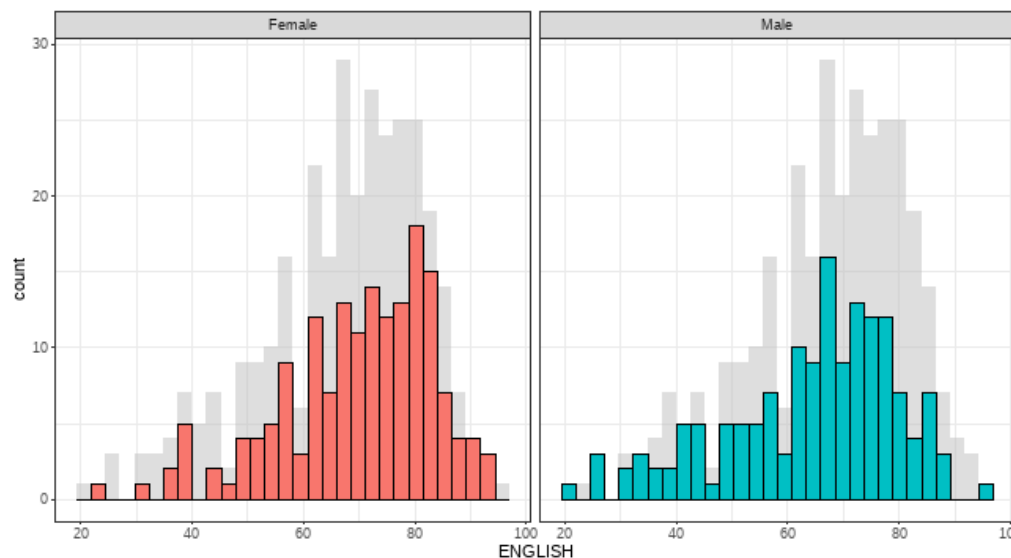
Designing Data-drive Graphics for Analysis III

The code chunk below is used to create the makeover design on the right. Note that the second line is used to create the so called **Background Data** - full without the 3th column (GENDER).

```
d <- exam_data
d_bg <- d[, -3]

ggplot(d, aes(x = ENGLISH,
              fill = GENDER)) +
  geom_histogram(data = d_bg,
                fill = "grey",
                alpha = .5) +
  geom_histogram(colour = "black") +
  facet_wrap(~ GENDER) +
  guides(fill = FALSE) +
  theme_bw()
```

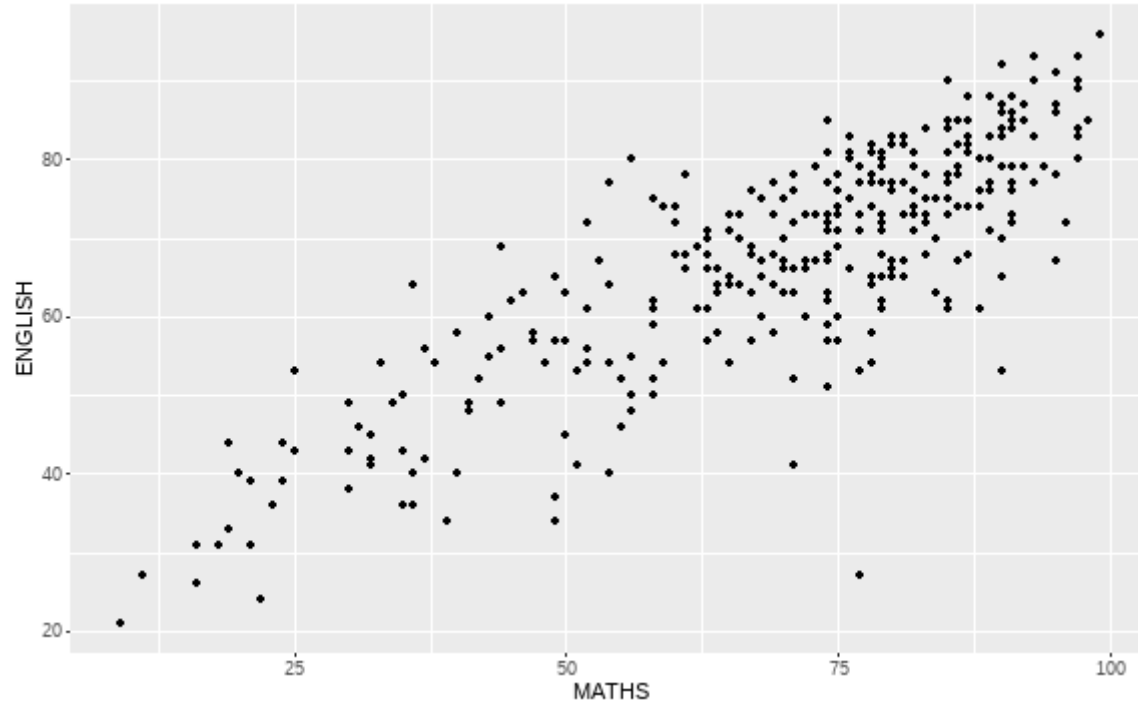
The makeover histograms on the right are not only elegantly designed but also informative. This is because they reveal the distribution of English scores by gender with reference to all pupils.



Credit: I learned this trick from an article entitled [Plotting background data for groups with ggplot2](#) from RBLOG maintain Simon Jackson.

Designing Data-drive Graphics for Analysis IV

The original design.



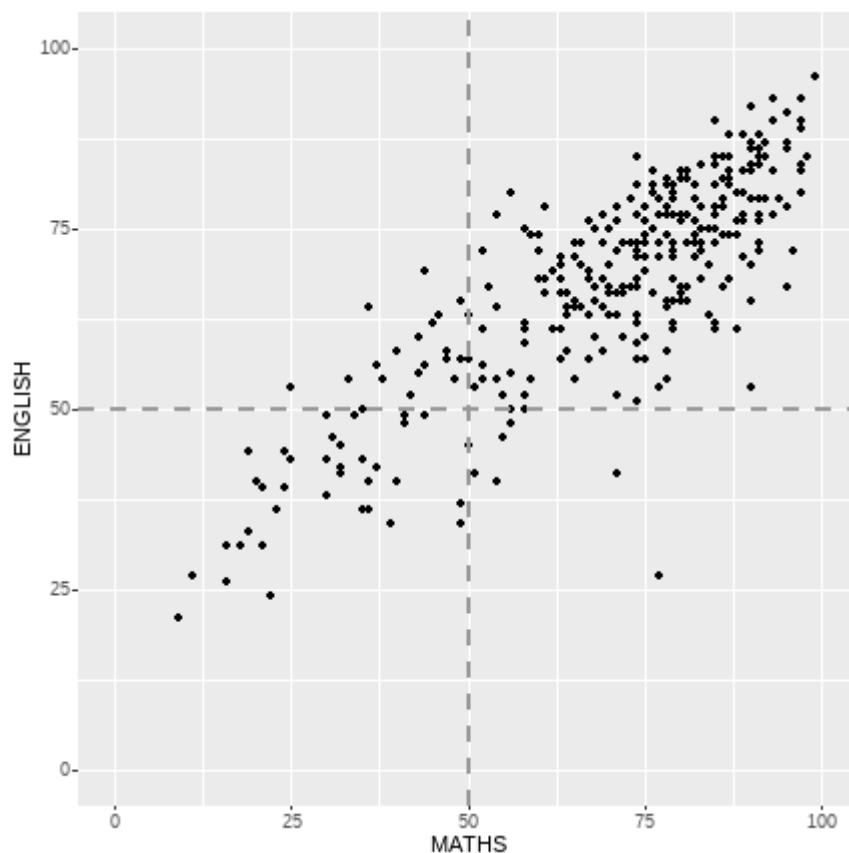
Designing Data-drive Graphics for Analysis IV

The makeover design

The code chunk used to create the makeover.

```
ggplot(data=exam_data,  
       aes(x=MATHS, y=ENGLISH)) +  
  geom_point() +  
  coord_cartesian(xlim=c(0,100),  
                 ylim=c(0,100)) +  
  geom_hline(yintercept=50,  
            linetype="dashed",  
            color="grey60",  
            size=1) +  
  geom_vline(xintercept=50,  
            linetype="dashed",  
            color="grey60",  
            size=1)
```

A within group scatterplot with reference lines.



Interactive with R - A Tale of Two R Packages

- ggiraph
- plotlyr



Interactive Data Visualisation - ggiraph methods

- An htmlwidget and a ggplot2 extension. It allows ggplot graphics to be interactive by exporting them as SVG documents and using special attributes on the various elements.
- Interactive is made with ggplot geometries that can understand three arguments:
 - **Tooltip**: a column of data-sets that contain tooltips to be displayed when the mouse is over elements.
 - **OnClick**: a column of data-sets that contain a JavaScript function to be executed when elements are clicked.
 - **Data_id**: a column of data-sets that contain an id to be associated with elements.
- If it used within a shiny application, elements associated with an id (data_id) can be selected and manipulated on client and server sides.

Reference: [ggiraph](#) package

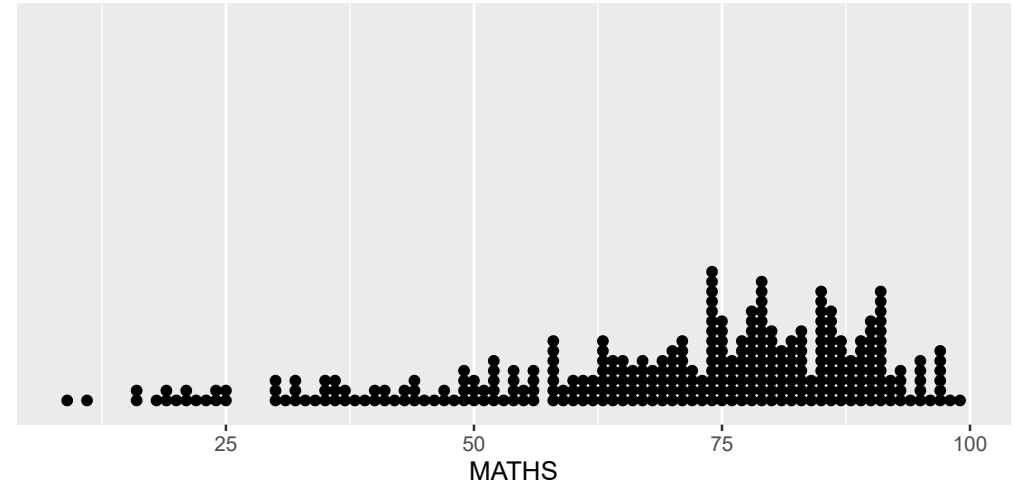
ggiraph Method: Interactive tooltip

The code chunk below uses `geom_dotplot_interactive()` of **ggiraph** package to implement an interactive tooltip on the dot plot.

```
p <- ggplot(data=exam_data,
  aes(x = MATHS)) +
  geom_dotplot_interactive(
    aes(tooltip = ID),
    stackgroups = TRUE,
    binwidth = 1,
    method = "histodot") +
  scale_y_continuous(NULL,
    breaks = NULL)

girafe(
  ggobj = p,
  width_svg = 6,
  height_svg = 6*0.5
)
```

Interactivity: hovering displays student's ID



Things to learn from the code chunk on the left:

- Instead of using `geom_dotplot()`, `geom_dotplot_interactive()` is used,
- An visual attribute (i.e. ID) is assigned to aesthetics tooltip to create an interactive element.
- `girafe()` is used with the ggplot object to translate the graphic as a web interactive graphics.

ggiraph Methods

What are the differences?

The original ggplot2 code chunk.

```
ggplot(data=exam_data,  
       aes(x = MATHS)) +  
  geom_dotplot(binwidth=2.5,  
              dotsize = 0.5) +  
  scale_y_continuous(NULL,  
                    breaks = NULL)
```

The ggiraph code chunk.

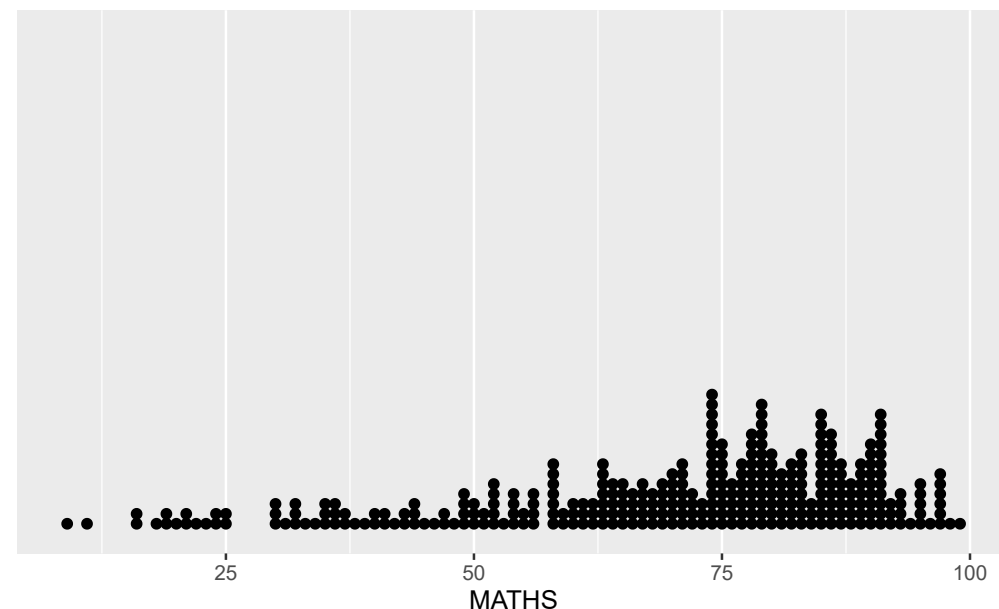
```
p <- ggplot(data=exam_data,  
           aes(x = MATHS)) +  
  geom_dotplot_interactive(  
    aes(tooltip = ID),  
    stackgroups = TRUE,  
    binwidth = 1,  
    method = "histodot") +  
  scale_y_continuous(NULL,  
                    breaks = NULL)  
  
girafe(  
  ggobj = p,  
  width_svg = 6,  
  height_svg = 6*0.618  
)
```

A complete list of geometries supported by ggiraph and their corresponding command syntax can be found [here](#).

ggiraph Methods

Hover effect with *data_id* aesthetic

```
p <- ggplot(data=exam_data,  
  aes(x = MATHS)) +  
  geom_dotplot_interactive(  
    aes(data_id = CLASS),  
    stackgroups = TRUE,  
    binwidth = 1,  
    method = "histodot") +  
    scale_y_continuous(NULL,  
      breaks = NULL)  
  
girafe(  
  ggobj = p,  
  width_svg = 6,  
  height_svg = 6*0.618  
)
```



Interactivity: Elements associated with a *data_id* (i.e CLASS) will be highlighted upon mouse over.

Note that the default value of the hover css is *hover_css = "fill:orange;"* .

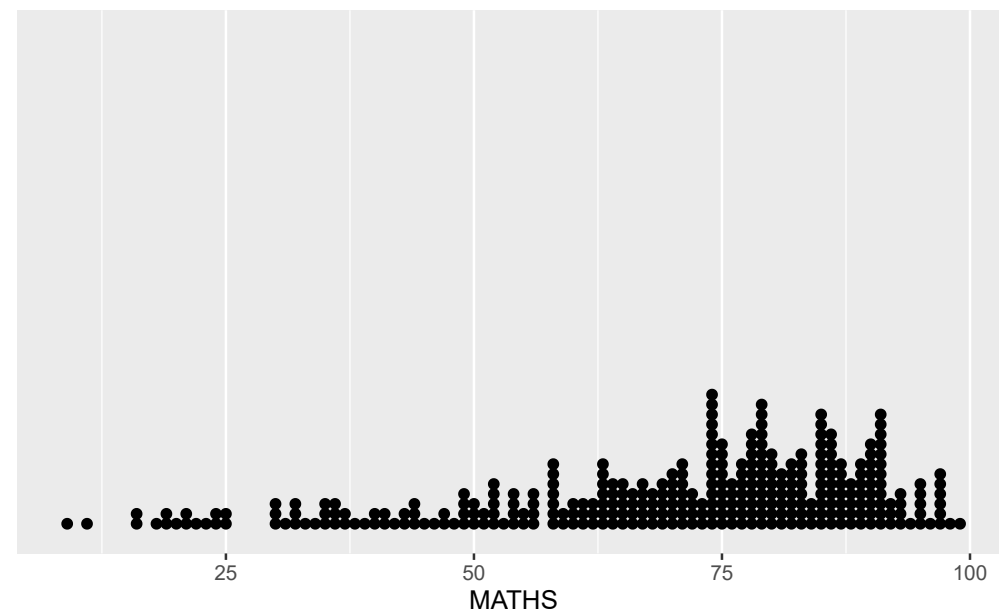
ggiraph Methods

Styling hover effect

In the code chunk below, css codes are used to change the highlighting effect.

```
p <- ggplot(data=exam_data,
  aes(x = MATHS)) +
  geom_dotplot_interactive(
    aes(data_id = CLASS),
    stackgroups = TRUE,
    binwidth = 1,
    method = "histodot") +
  scale_y_continuous(NULL,
    breaks = NULL)

girafe(
  ggobj = p,
  width_svg = 6,
  height_svg = 6*0.618,
  options = list(
    opts_hover(css = "fill: #202020;"),
    opts_hover_inv(css = "opacity:0.2;")
  )
)
```



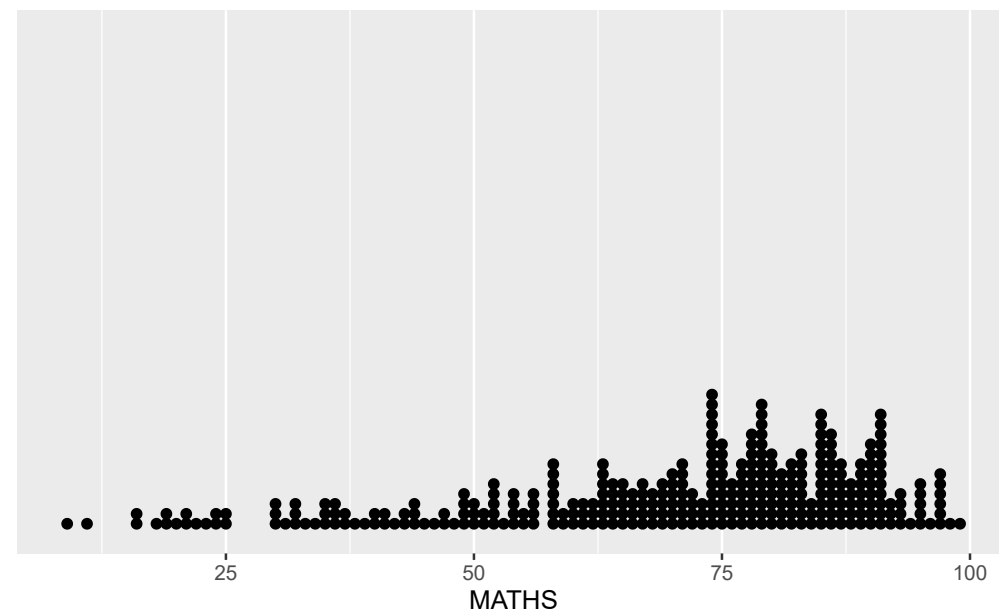
Interactivity: Elements associated with a *data_id* (i.e CLASS) will be highlighted upon mouse over.

ggiraph Methods

Combining tooltip and hover effects

```
p <- ggplot(data=exam_data,
  aes(x = MATHS)) +
  geom_dotplot_interactive(
    aes(data_id = CLASS,
      tooltip = CLASS),
    stackgroups = TRUE,
    binwidth = 1,
    method = "histodot") +
  scale_y_continuous(NULL,
    breaks = NULL)

girafe(
  ggobj = p,
  width_svg = 6,
  height_svg = 6*0.618,
  options = list(
    opts_hover(css = "fill: #202020;"),
    opts_hover_inv(css = "opacity:0.2;")
  )
)
```

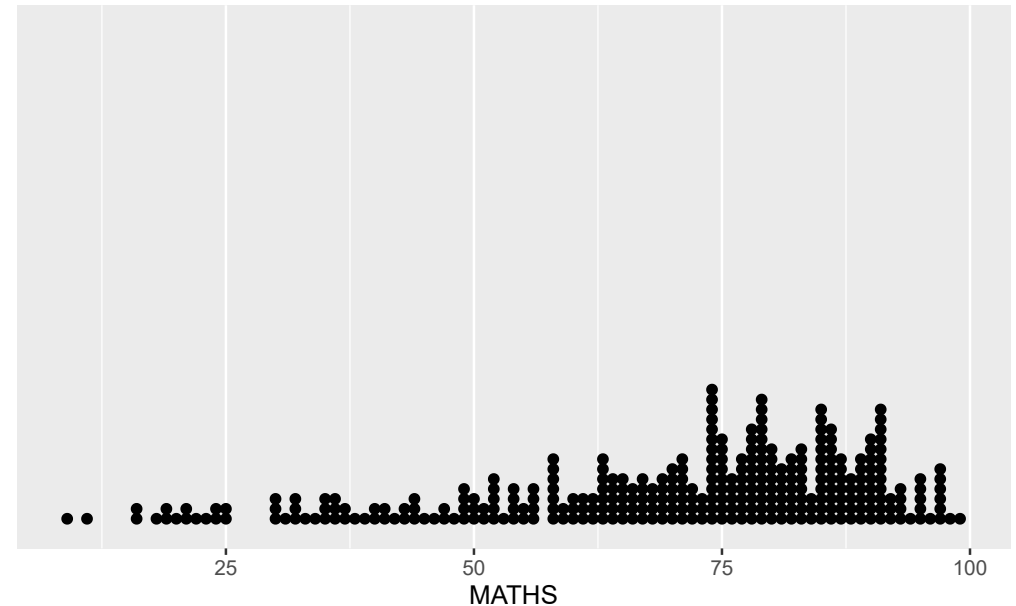


Interactivity: Elements associated with a *data_id* (i.e CLASS) will be highlighted upon mouse over and tooltip appears too.

ggiraph Methods

Click effect with onclick

```
exam_data$onclick <- sprintf("window.open(\"%s\",  
"https://www.moe.gov.sg/schoolfinder?journey=Pi  
  
p <- ggplot(data=exam_data,  
  aes(x = MATHS)) +  
  geom_dotplot_interactive(  
    aes(onclick = onclick),  
    stackgroups = TRUE,  
    binwidth = 1,  
    method = "histodot") +  
  scale_y_continuous(NULL,  
    breaks = NULL)  
  
girafe(  
  ggobj = p,  
  width_svg = 6,  
  height_svg = 6*0.618)
```

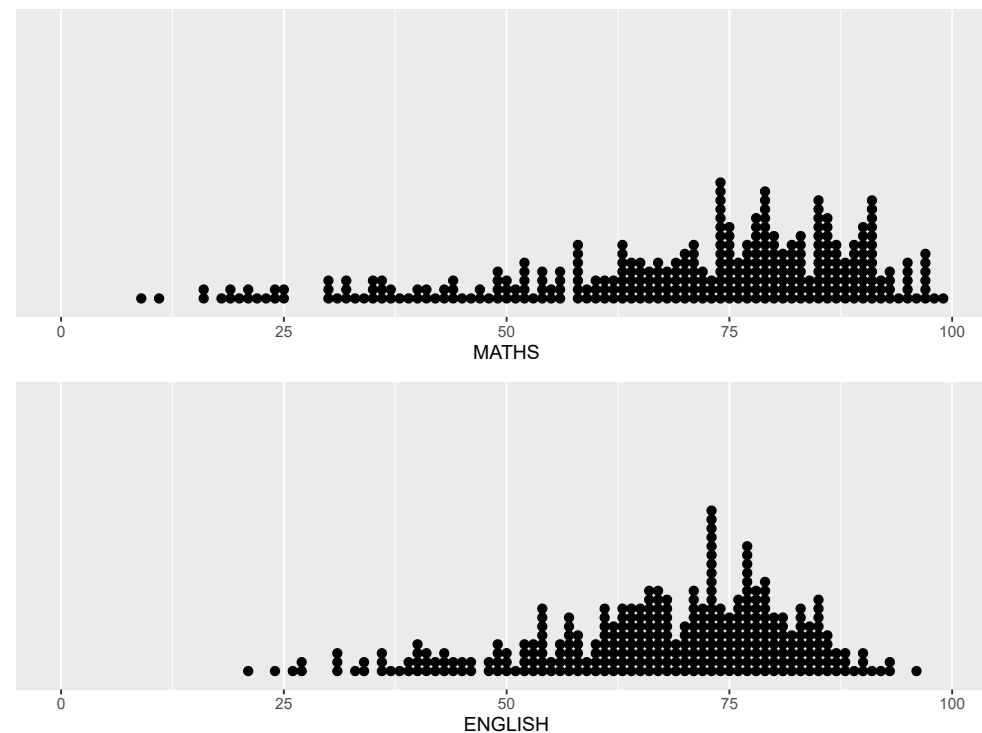


Interactivity: Web document link with a data object will be displayed on the web browser upon mouse click.

ggiraph Methods

Coordinated Multiple Views

The two dotplots on the right are linked views. In fact, they are coordinated. In this case, when a data point is selected on one of the dotplot, the corresponding data point with the similar ID will be highlighted too.



ggiraph Methods

Coordinated Multiple View

```
p1 <- ggplot(data=exam_data,  
  aes(x = MATHS)) +  
  geom_dotplot_interactive(  
    aes(data_id = ID,  
      tooltip = ID),  
    stackgroups = TRUE,  
    binwidth = 1,  
    method = "histodot") +  
  scale_y_continuous(NULL,  
    breaks = NULL) +  
  coord_cartesian(xlim=c(0,100))
```

```
p2 <- ggplot(data=exam_data,  
  aes(x = ENGLISH)) +  
  geom_dotplot_interactive(  
    aes(data_id = ID,  
      tooltip = ID),  
    stackgroups = TRUE,  
    binwidth = 1,  
    method = "histodot") +
```

```
  scale_y_continuous(NULL,  
    breaks = NULL) +  
  coord_cartesian(xlim=c(0,100))  
  
girafe(  
  code = print(p1 / p2),  
  width_svg = 8,  
  height_svg = 6,  
  options = list(  
    opts_hover(css = "fill: #202020;"),  
    opts_hover_inv(css = "opacity:0.2;")  
  )  
)
```

Things to learn from the code chunk above:

- As ggiraph is mainly only new geoms, **patchwork** package can be used seamlessly.

Reference

Textbook

Hadley Wickham (2011) **ggplot2: Elegant Graphics for Data Analysis**, Springer. [Online version](#).

Winston Chang (2013) **R Graphics Cookbook 2nd edition**. [Online version](#).

R packages

ggplot2

ggiraph

patchwork