

# Visualizing Data

Guy J. Abel

# R graphics Package

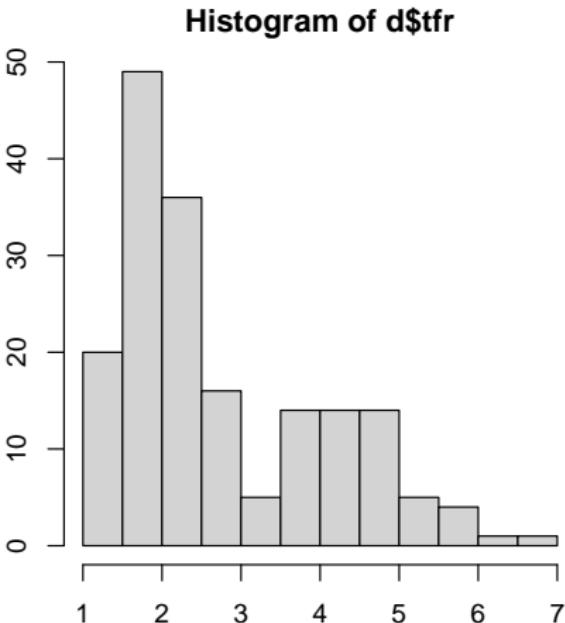
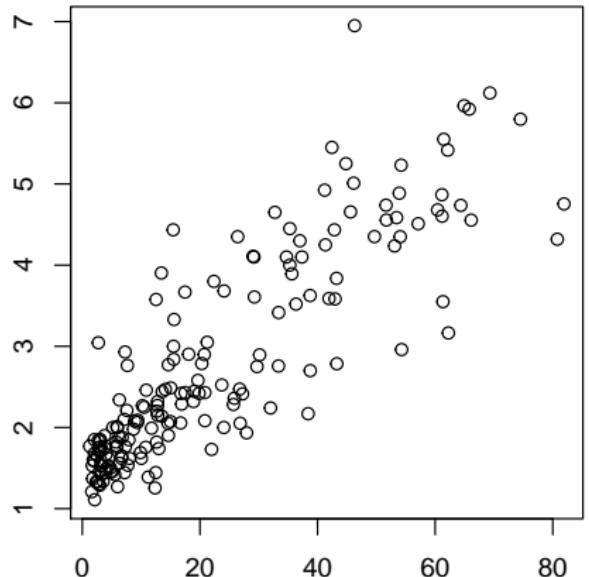
- R has fantastic graphic capabilities.
  - There are many approaches and packages designed for this specific task.
  - The `graphics` package (loaded when start R) has many plotting functions.
  - Simple look, may not be consistent across different types of plots

# R graphics Package

```
> d
# A tibble: 179 x 8
  name      region_name area_name   imr     tfr sex_ratio developed growth_policy
  <chr>    <chr>       <chr>    <dbl>   <dbl>   <dbl> <chr>    <chr>
1 Burundi  Eastern Afr~ Africa    42.4    5.45   1.03 Less    Lower
2 Comoros  Eastern Afr~ Africa   53.1    4.24   1.05 Less    Lower
3 Djibouti Eastern Afr~ Africa  33.4    2.76   1.04 Less    Lower
4 Eritrea  Eastern Afr~ Africa  34.7    4.1    1.05 Less    Lower
5 Ethiopia Eastern Afr~ Africa  37      4.3    1.04 Less    Lower
6 Kenya    Eastern Afr~ Africa  36.3    3.52   1.03 Less    Lower
7 Madagas~ Eastern Afr~ Africa  29.0    4.11   1.03 Less    Lower
8 Malawi   Eastern Afr~ Africa  41.3    4.25   1.03 Less    Lower
9 Mauriti~ Eastern Afr~ Africa 11.2    1.39   1.04 Less    Maintain
10 Mozambi~ Eastern Afr~ Africa 53.9    4.89   1.02 Less    Lower
# ... with 169 more rows
```

# R graphics Package

```
> # left  
> plot(x = d$imr, y = d$tfr)  
> # right  
> hist(x = d$tfr)
```

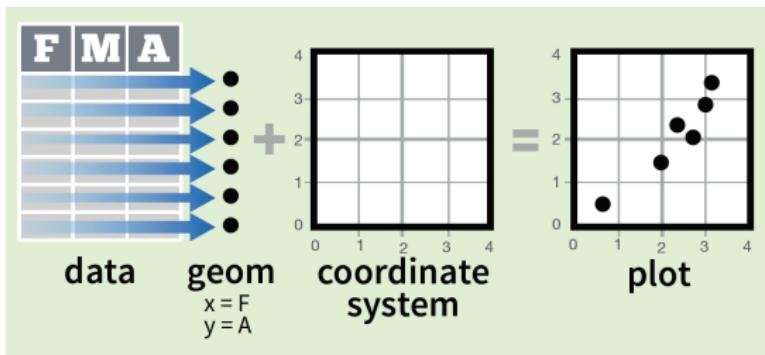


# Grammar of Graphics

- The `ggplot2` package by Hadley Wickham implements the grammar of graphics, a coherent system for describing and building graphs.
  - One of the most beautiful and most versatile.
  - Flexible and consistent approach across different plots
  - Do more faster as a single learning one system for all plots types that you can apply in many data situations.
- Included in the `tidyverse` package (more on `tidyverse` later)
  - `library(tidyverse)` loads `ggplot2` along with other useful packages

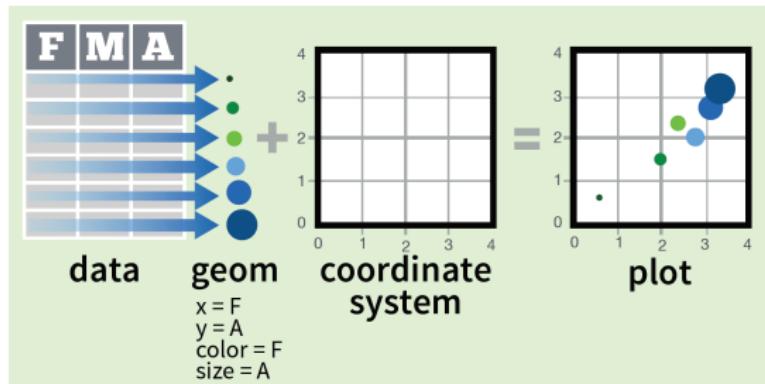
# Grammar of Graphics

- The grammar of graphics builds every graph from the same components:
  - data set
  - coordinate system
  - geoms (visual marks that represent data points)



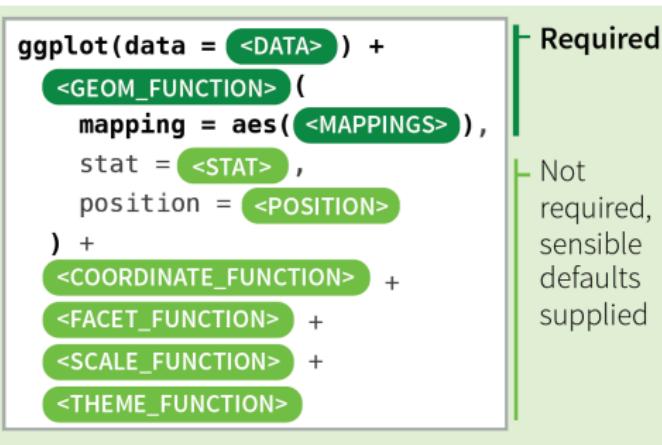
# Grammar of Graphics

- When variables are mapped to data they can take different visual properties of the geom (aesthetics) like size, color, and x and y locations.
- Aesthetics meaning something you can see.



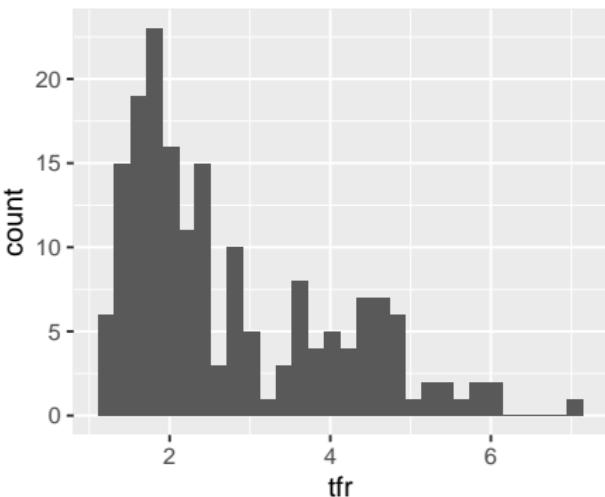
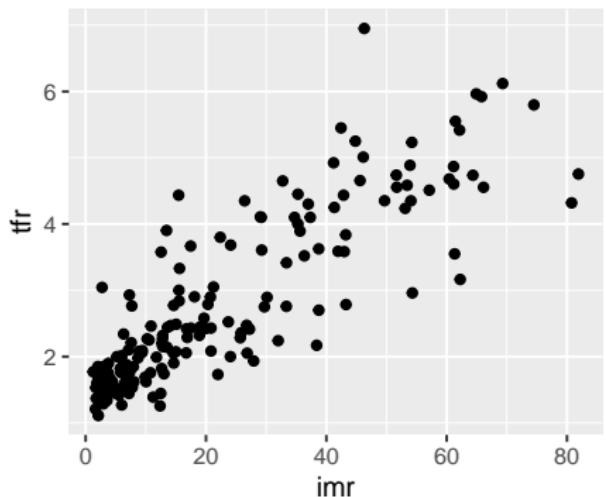
# Grammar of Graphics

- The R code follows a common template



# Grammar of Graphics

```
> library(tidyverse)
> # left
> ggplot(data = d, mapping = aes(x = imr, y = tfr)) +
+   geom_point()
> # right
> ggplot(data = d, mapping = aes(x = tfr)) +
+   geom_histogram()
```

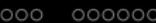
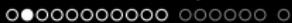


# Grammar of Graphics

- The function `ggplot()` creates a coordinate system to add layers onto.
- The first argument of `ggplot()` is the data set used to graph
  - For example, `ggplot(data = swiss)` creates an empty graph that will use the `swiss` data set.
  - Can then build on the empty graphs by adding one or more layers to `ggplot()`.
- The `geom` functions adds a layer of points to your plot.
  - The `ggplot2` library comes with many `geom` functions
  - Each add a different type of layer to a plot.
  - Each `geom` function in `ggplot2` takes a `mapping` argument.
- The `mapping` argument explains where your points should go
  - Set in the `mapping` argument with a call to the `aes()` function.
  - For example the `x` and `y` arguments of `aes()` explain which variables to map to the `x-` and `y-axes` of your plot.
  - The `ggplot()` function looks for those variables in your data set.
  - The `mapping` arguments are different for each `geom` function

# Geom

- The ggplot2 package provides over 30 geom functions
  - Further packages provide even more.
- Each geom function is suitable for visualizing a certain type of data or relationship.
  - Listed on the first page of the ggplot cheat sheet.
  - Next to geom is a visual representation of the geom.
  - Beneath geom is a list of aesthetics that apply to the geom.
  - Required aesthetics are in bold.



# Geoms

## GRAPHICAL PRIMITIVES

a <- ggplot(economics, aes(date, unemploy))  
 b <- ggplot(seals, aes(x = long, y = lat))

a + geom\_blank()  
 (Useful for expanding limits)

b + geom\_curve(aes(yend = lat + 1,  
 xend=long+1,curvature=z)) -> x, xend, y, yend,  
 alpha, angle, color, curvature, linetype, size

a + geom\_path(lineend="butt", linejoin="round",  
 linemitre=i)  
 x, y, alpha, color, group, linetype, size

a + geom\_polygon(aes(group = group))  
 x, y, alpha, color, fill, group, linetype, size

b + geom\_rect(aes(xmin = long, ymin=lat, xmax=  
 long + 1, ymax = lat + 1)) -> xmax, xmin, ymax,  
 ymin, alpha, color, fill, linetype, size

a + geom\_ribbon(aes(ymax=unemploy - 900,  
 ymáx=unemploy + 900)) -> x, ymax, ymin,  
 alpha, color, fill, group, linetype, size

## LINE SEGMENTS

common aesthetics: x, y, alpha, color, linetype, size

b + geom\_abline(aes(intercept=0, slope=1))  
 b + geom\_hline(aes(yintercept = lat))  
 b + geom\_vline(aes(xintercept = long))

b + geom\_segment(aes(yend=lat+1, xend=long+1))  
 b + geom\_spoke(aes(angle=1:1155, radius = 1))

## ONE VARIABLE continuous

c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)

## TWO VARIABLES

continuous x , continuous y  
 e <- ggplot(mpg, aes(cty, hwy))

e + geom\_label(aes(label = cty), nudge\_x = 1,  
 nudge\_y = 1, check\_overlap = TRUE) x, y, label,  
 alpha, angle, color, family, fontface, hjust,  
 lineheight, size, vjust

e + geom\_jitter(height = 2, width = 2)  
 x, y, alpha, color, fill, shape, size

e + geom\_point(), x, y, alpha, color, fill, shape,  
 size, stroke

e + geom\_quantile(), x, y, alpha, color, group,  
 linetype, size, weight

e + geom\_rug(sides = "bl") x, y, alpha, color,  
 linetype, size

e + geom\_smooth(method = lm), x, y, alpha,  
 color, fill, group, linetype, size, weight

e + geom\_text(aes(label = cty), nudge\_x = 1,  
 nudge\_y = 1, check\_overlap = TRUE), x, y, label,  
 alpha, angle, color, family, fontface, hjust,  
 lineheight, size, vjust

## discrete x , continuous y

f <- ggplot(mpg, aes(class, hwy))

f + geom\_col(), x, y, alpha, color, fill, group,  
 linetype, size

f + geom\_boxplot(), x, y, lower, middle, upper,  
 ymax, ymin, alpha, color, fill, group, linetype,  
 shape, size, weight

## continuous bivariate distribution

h <- ggplot(diamonds, aes(carat, price))

h + geom\_bin2d(binwidth = c(0.25, 500))  
 x, y, alpha, color, fill, linetype, size, weight

h + geom\_density2d()  
 x, y, alpha, colour, group, linetype, size

h + geom\_hex()  
 x, y, alpha, colour, fill, size

## continuous function

i <- ggplot(economics, aes(date, unemploy))

i + geom\_area()  
 x, y, alpha, color, fill, linetype, size

i + geom\_line()  
 x, y, alpha, color, group, linetype, size

i + geom\_step(direction = "hv")  
 x, y, alpha, color, group, linetype, size

## visualizing error

df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)  
 j <- ggplot(df, aes(grp, fit, ymin = fit-se, ymax = fit+se))

j + geom\_crossbar(fatten = 2)  
 x, y, ymax, ymin, alpha, color, fill, group, linetype,  
 size

j + geom\_errorbar(), x, y, max, min, alpha, color,  
 group, linetype, size, width (also  
 geom\_errorbarh())

i + geom\_linerange()



# Geoms

common aesthetics: x, y, alpha, color, linetype, size

```
b + geom_abline(aes(intercept=0, slope=1))
b + geom_hline(aes(yintercept = lat))
b + geom_vline(aes(xintercept = long))

b + geom_segment(aes(yend=lat+1, xend=long+1))
b + geom_spoke(aes(angle = 1:1155, radius = 1))
```

## ONE VARIABLE continuous

c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)



c + geom\_area(stat = "bin")  
x, y, alpha, color, fill, linetype, size



c + geom\_density(kernel = "gaussian")  
x, y, alpha, color, fill, group, linetype, size, weight



c + geom\_dotplot()  
x, y, alpha, color, fill



c + geom\_freqpoly()  
x, y, alpha, color, group, linetype, size



c + geom\_histogram(binwidth = 5)  
x, y, alpha, color, fill, linetype, size, weight

---

discrete

d <- ggplot(mpg, aes(fl))  
d + geom\_bar()  
x, alpha, color, fill, linetype, size, weight

discrete x , continuous y

```
f <- ggplot(mpg, aes(class, hwy))

f + geom_col(), x, y, alpha, color, fill, group,
linetype, size

f + geom_boxplot(), x, y, lower, middle, upper,
ymin, ymax, alpha, color, fill, group, linetype,
shape, size, weight

f + geom_dotplot(binaxis = "y", stackdir =
"center"), x, y, alpha, color, fill, group

f + geom_violin(scale = "area"), x, y, alpha, color,
fill, group, linetype, size, weight
```

discrete x , discrete y

g <- ggplot(diamonds, aes(cut, color))

g + geom\_count(), x, y, alpha, color, fill, shape,
size, stroke

## THREE VARIABLES

seals\$z <- with(seals, sqrt(delta\_long^2 + delta\_lat^2)); l <- ggplot(seals, aes(long, lat))

l + geom\_contour(aes(z = z))
x, y, z, alpha, colour, group, linetype,
size, weight

## visualizing error

df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)

j <- ggplot(df, aes(grp, fit, ymin = fit-se, ymax = fit+se))



j + geom\_crossbar(fatten = 2)
x, y, ymax, ymin, alpha, color, fill, group, linetype,
size



j + geom\_errorbar(), x, ymax, ymin, alpha, color,
group, linetype, size, width (also
geom\_errorbarh())



j + geom\_linerange()
x, ymin, ymax, alpha, color, group, linetype, size



j + geom\_pointrange()
x, y, ymin, ymax, alpha, color, fill, group, linetype,
shape, size

## maps

data <- data.frame(murder = USArrests\$Murder,
state = tolower(rownames(USArrests)))
map <- map\_data("state")
k <- ggplot(data, aes(fill = murder))



k + geom\_map(aes(map\_id = state), map = map)
+ expand\_limits(x = map\$long, y = map\$lat),
map\_id, alpha, color, fill, linetype, size

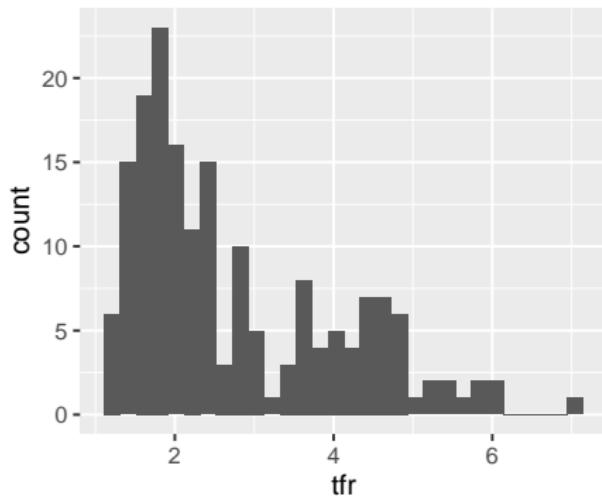
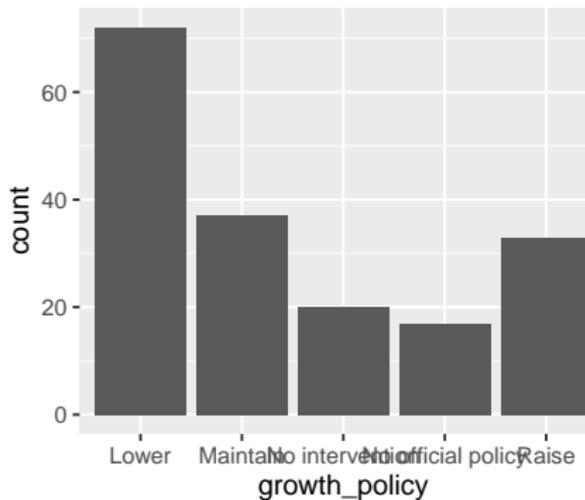
l + geom\_raster(aes(fill = z), hjust=0.5, vjust=0.5,
interpolate=FALSE)
x, y, alpha, fill



l + geom\_tile(aes(fill = z)), x, y, alpha, color, fill,
linetype, size, width

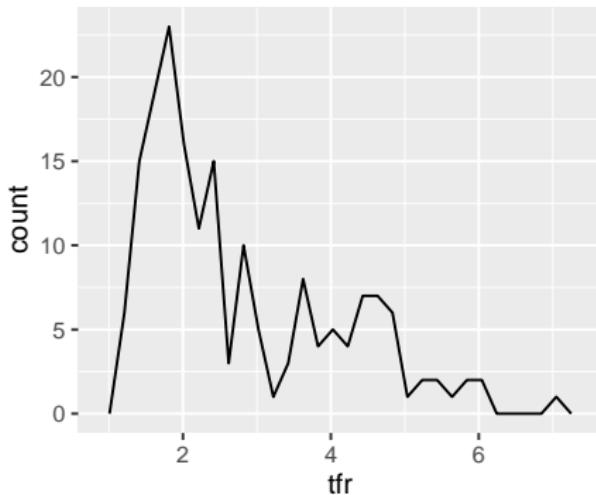
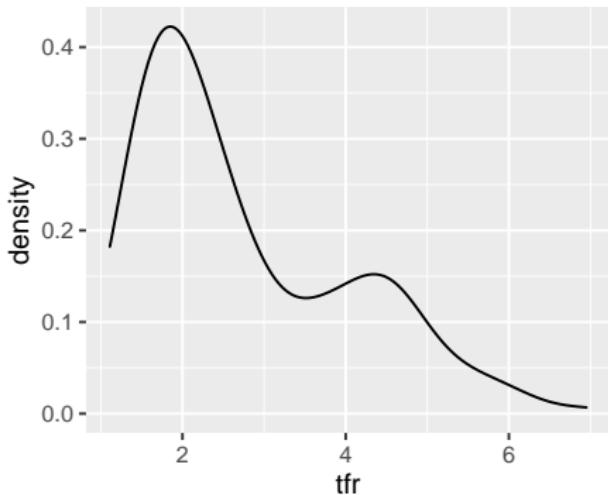
# Geom Examples: One Variable

```
> # left: one discrete variable
> ggplot(data = d, mapping = aes(x = growth_policy)) +
+   geom_bar()
> # right: one continuous variable
> ggplot(data = d, mapping = aes(x = tfr)) +
+   geom_histogram()
```



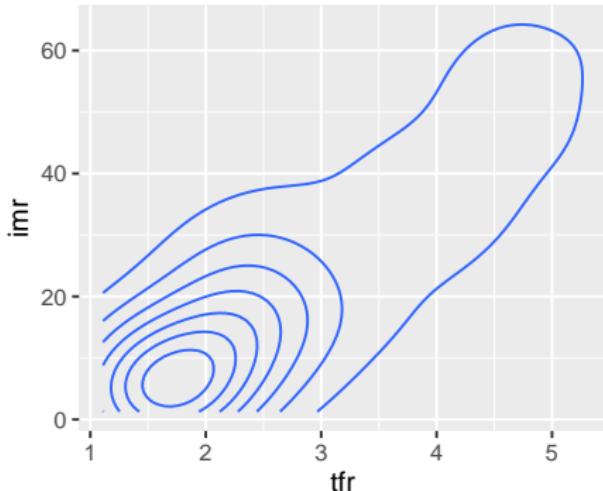
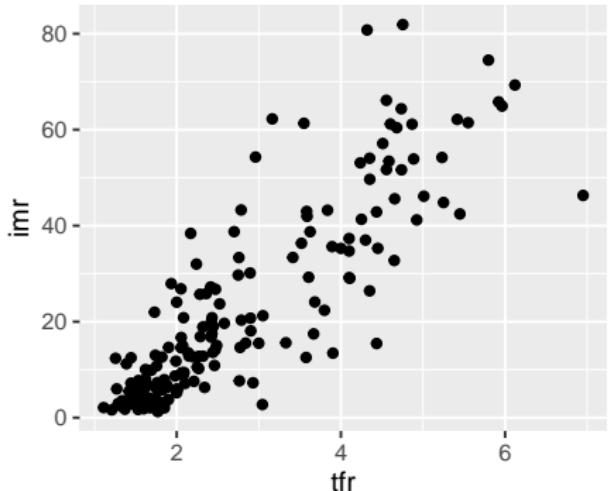
# Geom Examples: One Variable

```
> # left: one continuous variable  
> ggplot(data = d, mapping = aes(x = tfr)) +  
+   geom_density()  
> # right: one continuous variable  
> ggplot(data = d, mapping = aes(x = tfr)) +  
+   geom_freqpoly()
```



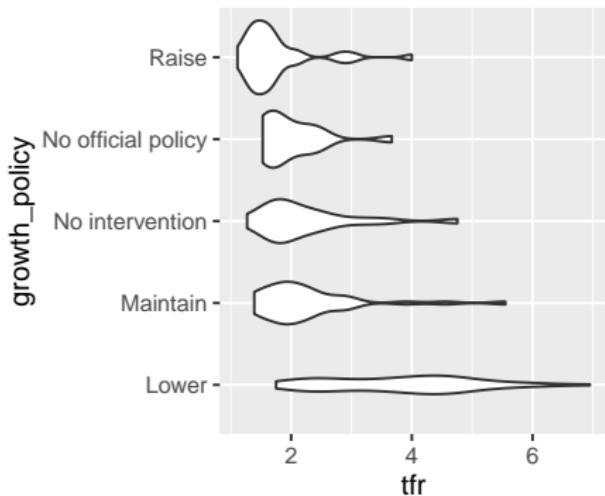
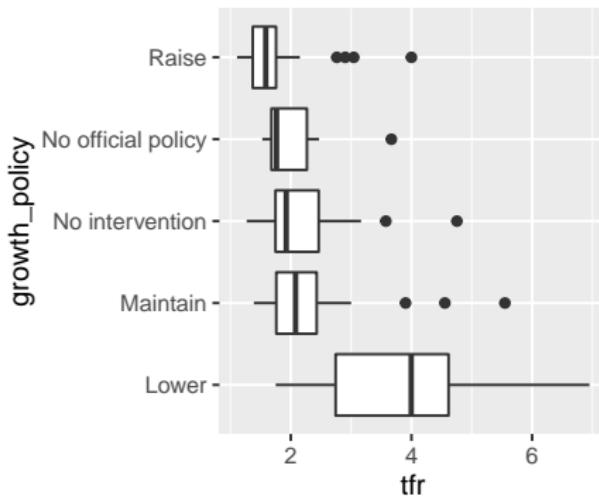
# Geom Examples: Two Variables

```
> # left: both continuous variable  
> ggplot(data = d, mapping = aes(x = tfr, y = imr)) +  
+   geom_point()  
> # right: both continuous variable  
> ggplot(data = d, mapping = aes(x = tfr, y = imr)) +  
+   geom_density2d()
```



# Geom Examples: Two Variables

```
> # left: one continuous, one discrete variable
> ggplot(data = d, mapping = aes(x = tfr, y = growth_policy)) +
+   geom_boxplot()
> # right: one continuous, one discrete variable
> ggplot(data = d, mapping = aes(x = tfr, y = growth_policy)) +
+   geom_violin()
```

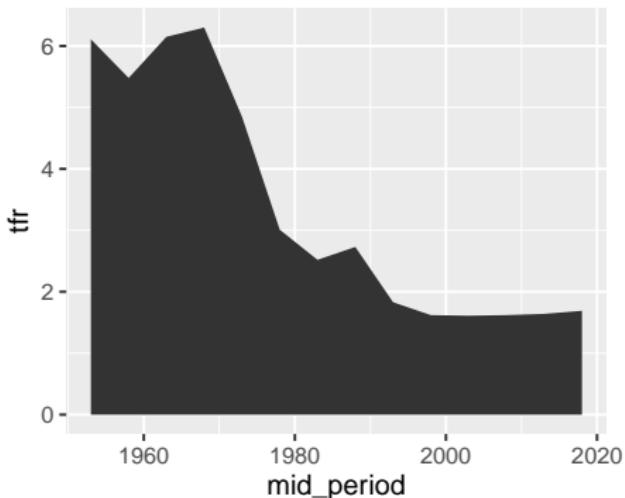
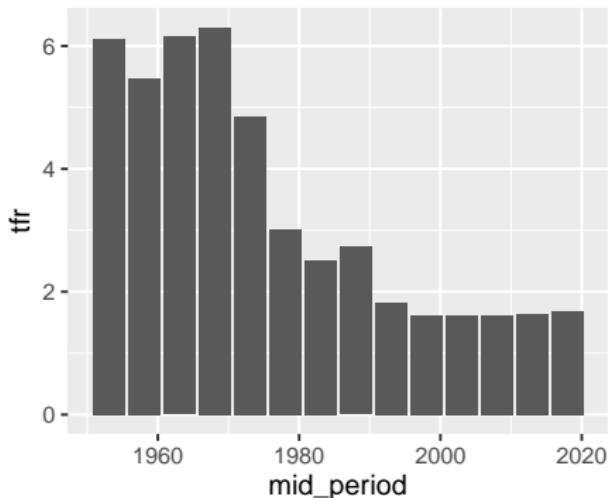


# Geom Examples: Time Variable

```
> china
# A tibble: 14 x 5
  name period mid_period imr    tfr
  <chr> <chr>     <dbl> <dbl> <dbl>
1 China 1950-1955     1953 129.   6.11
2 China 1955-1960     1958 131.   5.48
3 China 1960-1965     1963 135.   6.15
4 China 1965-1970     1968 93.8   6.3 
5 China 1970-1975     1973 71.9   4.85
6 China 1975-1980     1978 55.0   3.01
7 China 1980-1985     1983 44.6   2.52
8 China 1985-1990     1988 42.3   2.73
9 China 1990-1995     1993 41.3   1.83
10 China 1995-2000    1998 35.1   1.62
11 China 2000-2005    2003 27.1   1.61
12 China 2005-2010    2008 18.1   1.62
13 China 2010-2015    2013 12.3   1.64
14 China 2015-2020    2018 9.89   1.69
```

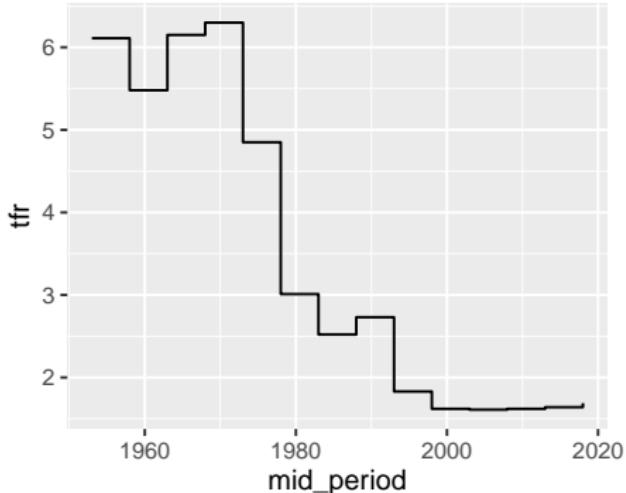
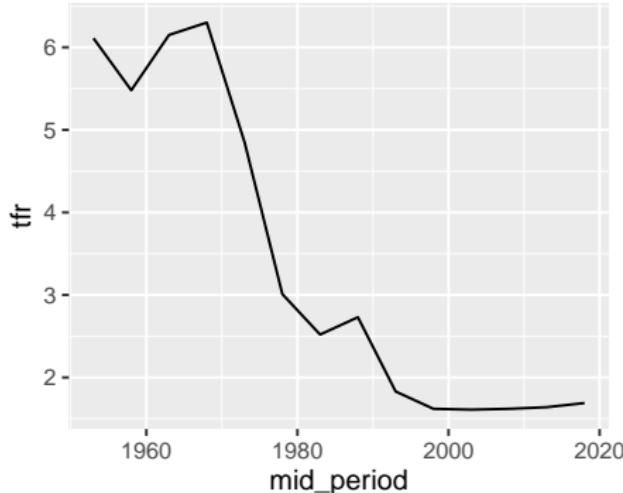
# Geom Examples: Time Variable

```
> # left  
> ggplot(data = china, mapping = aes(x = mid_period, y = tfr)) +  
+   geom_col()  
> # right  
> ggplot(data = china, mapping = aes(x = mid_period, y = tfr)) +  
+   geom_area()
```



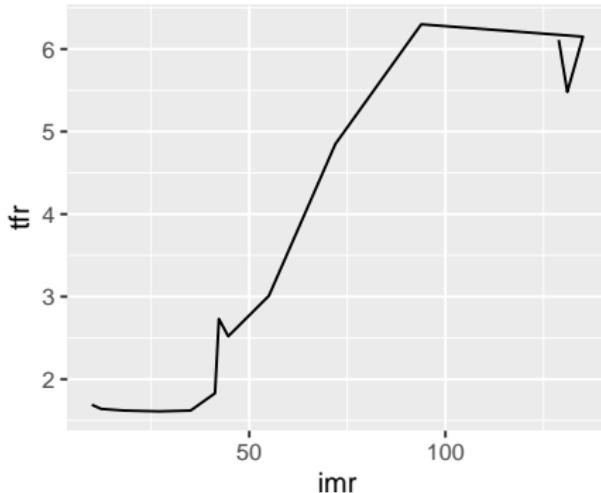
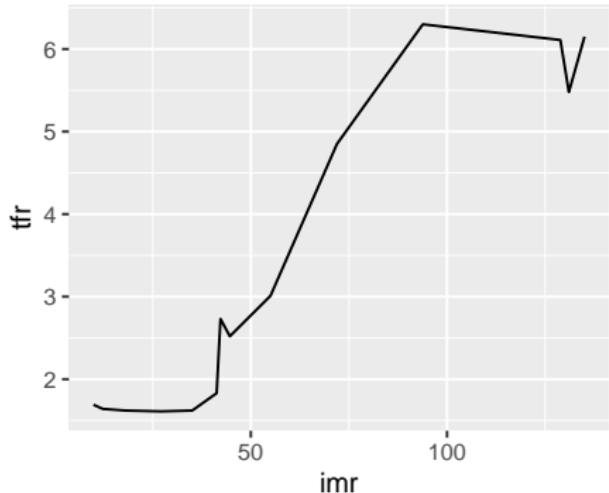
# Geom Examples: Time Variable

```
> # left  
> ggplot(data = china, mapping = aes(x = mid_period, y = tfr)) +  
+   geom_line()  
> # right  
> ggplot(data = china, mapping = aes(x = mid_period, y = tfr)) +  
+   geom_step()
```



# Geom Examples: Time Variable

```
> # left: follows order on x-axis (if x is not time can cause problem)
> ggplot(data = china, mapping = aes(x = imr, y = tfr)) +
+   geom_line()
> # right: follows order in data
> ggplot(data = china, mapping = aes(x = imr, y = tfr)) +
+   geom_path()
```

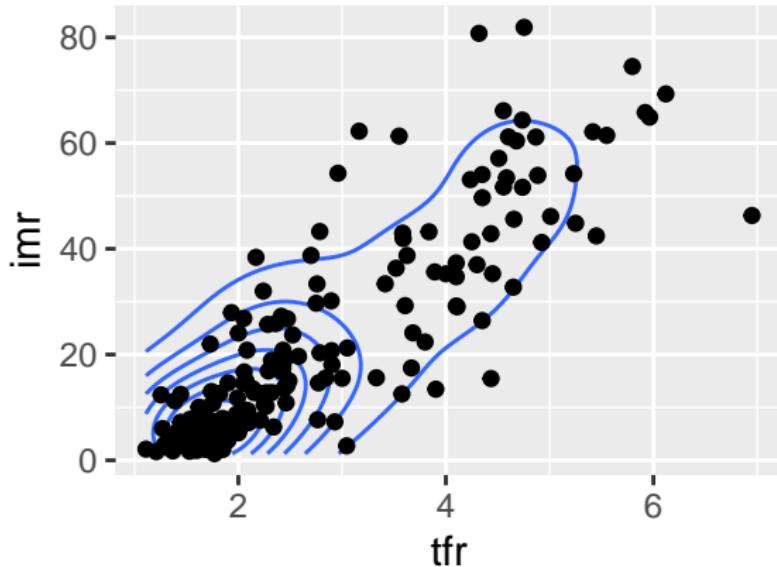


# Layers

- Geoms can be added in layers
  - Adds more detail the plots
- The data and `mapping` aesthetics in the `ggplot` function are shared by all the geom functions
- To overwrite for a individual geom you can place `mappings` or `data` arguments in the geom function
  - This will treat `mappings` or `data` local to the layer only.

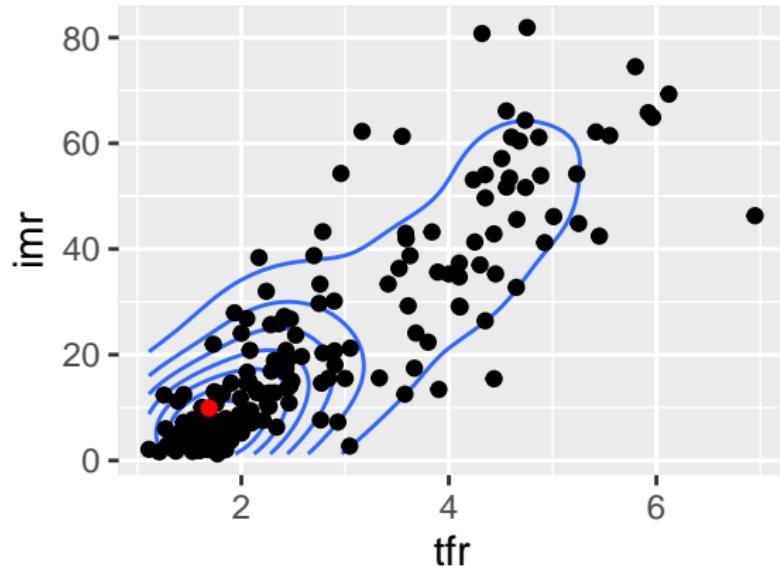
# Layers

```
> ggplot(data = d, mapping = aes(x = tfr, y = imr)) +  
+   geom_density2d() +  
+   geom_point()
```



# Layers

```
> # add a layer of red points for latest China data - more on aesthetics later!
> ggplot(data = d, mapping = aes(x = tfr, y = imr)) +
+   geom_density2d() +
+   geom_point() +
+   geom_point(data = china[14,], colour = "red")
```

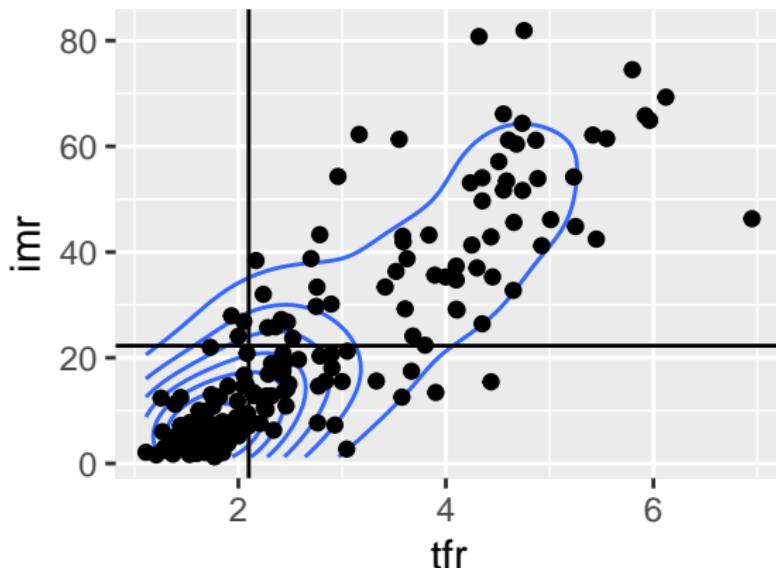


# Layers

- Line geoms can help indicate important thresholds
  - The `geom_hline()` function for horizontal line with `yintercept` arguments.
  - The `geom_vline()` function for vertical line with `xintercept` arguments.
  - The `geom_abline()` function for a line with `intercept` and `slope` arguments.
- We can apply a R function to a column before plotting
  - Use within `mapping` argument

# Layers

```
> ggplot(data = d, mapping = aes(x = tfr, y = imr)) +  
+   geom_density2d() +  
+   geom_point() +  
+   geom_vline(xintercept = 2.1) +  
+   # use mapping when applying a function to a column in the data  
+   geom_hline(mapping = aes(yintercept = mean(imr)))
```



# Exercise 1 (ex21.R)

```
# 0. a) Set the working directory to the course folder on your computer by loading
#      b) Check the working directory
getwd()
#      c) Load the tidyverse package (which loads the ggplot2 package amongst others)
library(#####)
#      d) Run the code in ex21_prelim.R to import the UN data for this exercise
source("./exercise./ex21_prelim.R")
#      e) Get familiar with the data by printing to console...
# all countries, 2010-2015 only
d
# UK data, all periods
uk
# South Eastern Asia data, all periods
se_asia
##
##
##
## 1. Create a scatter plot of infant mortality rates (x) against total fertility rate
#      in all countries in the 2010-2015 period
#      (Hint: a) see above to help select the correct data b) use imr and tfr columns)
#
# 2. Create a bar plot of the immigration policies in different countries from d
#      (Hint: policy data in the immigration_policy variable)
```

# Groups

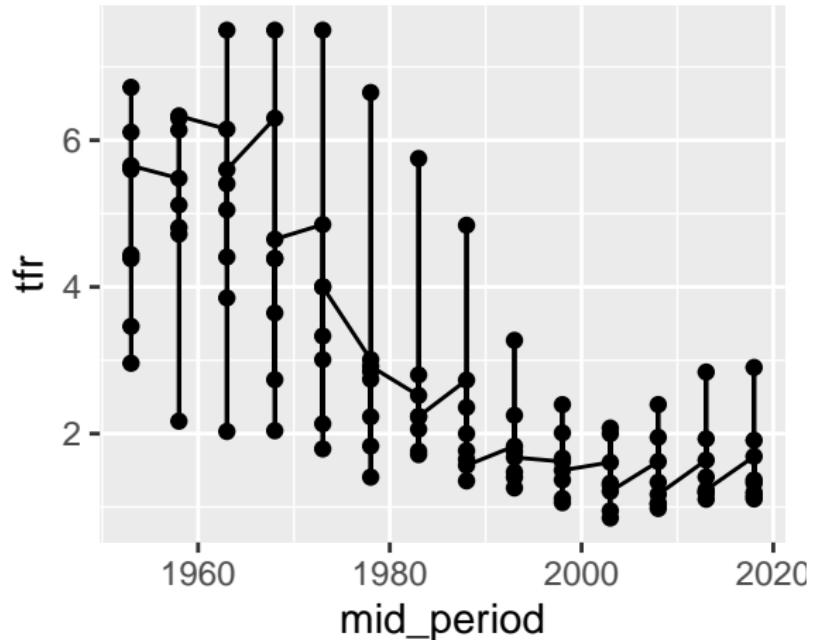
- Many `geoms_` functions use a single object to describe all of the data.
- You can set the `group` aesthetic in the `mappings` to a **discrete** variable to draw multiple objects.
  - This will tell `ggplot2` to draw separate object for each unique value of the grouping variable.
- Will not add a legend or distinguishing features to the plot.
  - These are added when we have some more aesthetics for the `mappings`, e.g. `colour` or `linetype` (will get to these soon).

# Groups

```
> east_asia
# A tibble: 112 x 7
  name              period   mid_period    imr     tfr sex_ratio     pop
  <chr>            <chr>      <dbl> <dbl> <dbl> <dbl> <dbl>
1 China           1950-1955    1953  129.   6.11  1.07  5.93e5
2 China, Hong Kong SAR 1950-1955    1953   61.9  4.44  1.09  2.25e3
3 China, Macao SAR 1950-1955    1953   65.9  4.39  1.05  1.99e2
4 China, Taiwan Province of ~ 1950-1955    1953   78.8  6.72  1.05  8.57e3
5 Dem. People's Republic of ~ 1950-1955    1953  121.   3.46  1.05  9.96e3
6 Japan            1950-1955    1953   47.7  2.96  1.06  8.69e4
7 Mongolia         1950-1955    1953  183.    5.6   1.03  8.23e2
8 Republic of Korea 1950-1955    1953  159.    5.65  1.06  2.03e4
9 China            1955-1960    1958  131.    5.48  1.07  6.40e5
10 China, Hong Kong SAR 1955-1960   1958   46.3  4.72  1.08  2.83e3
# ... with 102 more rows
```

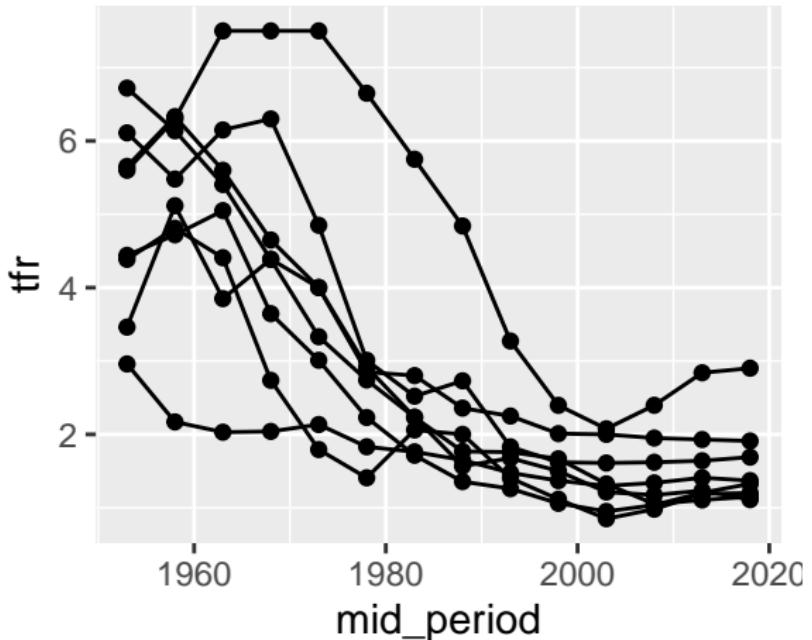
# Groups

```
> # eastern asia: line through all points  
> ggplot(data = east_asia, mapping = aes(x = mid_period, y = tfr)) +  
+   geom_point() +  
+   geom_line()
```



# Groups

```
> # eastern asia: line for each country (identified by name)
> ggplot(data = east_asia, mapping = aes(x = mid_period, y = tfr, group = name)) +
+   geom_point() +
+   geom_line()
```



# Optional Aesthetics

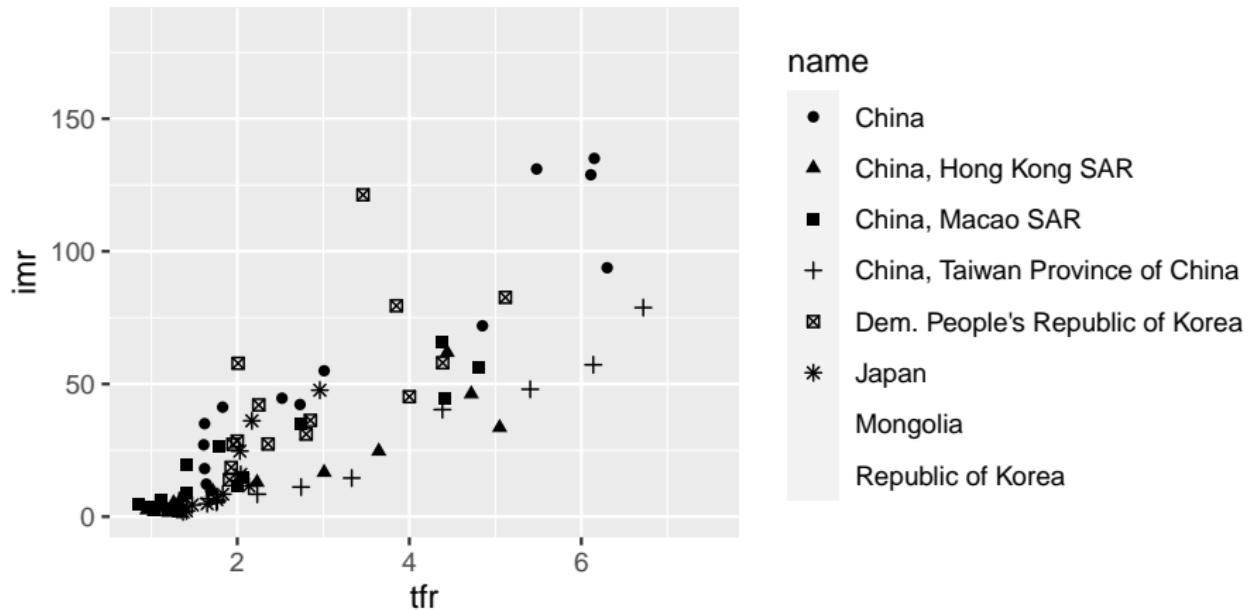
- Beyond group, there are other aesthetics (aes) that can be used depending on the geom, for example:
  - shape: shape of point
  - linetype: type of line used
  - size: size of object
  - colour: border colour (i.e. for dimensional object)
  - fill: internal colour (i.e. for two dimensional object)
  - alpha: transparency (between 0 and 1)
  - and many more ...
- Not all aesthetics work with all geom.
  - For example `geom_point()` has no linetype option.
  - Check the cheatsheet to see what is available.

# Optional Aesthetics

- You can set the aesthetic to a
  - Single value to change the appearance for all plotted objects
    - Use outside of `mapping = aes()`
  - Variables (columns) in the data
    - Use inside of `mapping = aes()`
- When set to variable names in the data
  - If discrete will operate by each group (category) of the variable.
  - If continuous will operate along the scale of the variable.
- Not all aesthetics can be set to continuous variables
  - For example, `linetype = mid_period`
- Can set different aesthetics to different variables, allowing you to show more many dimensions of your data.

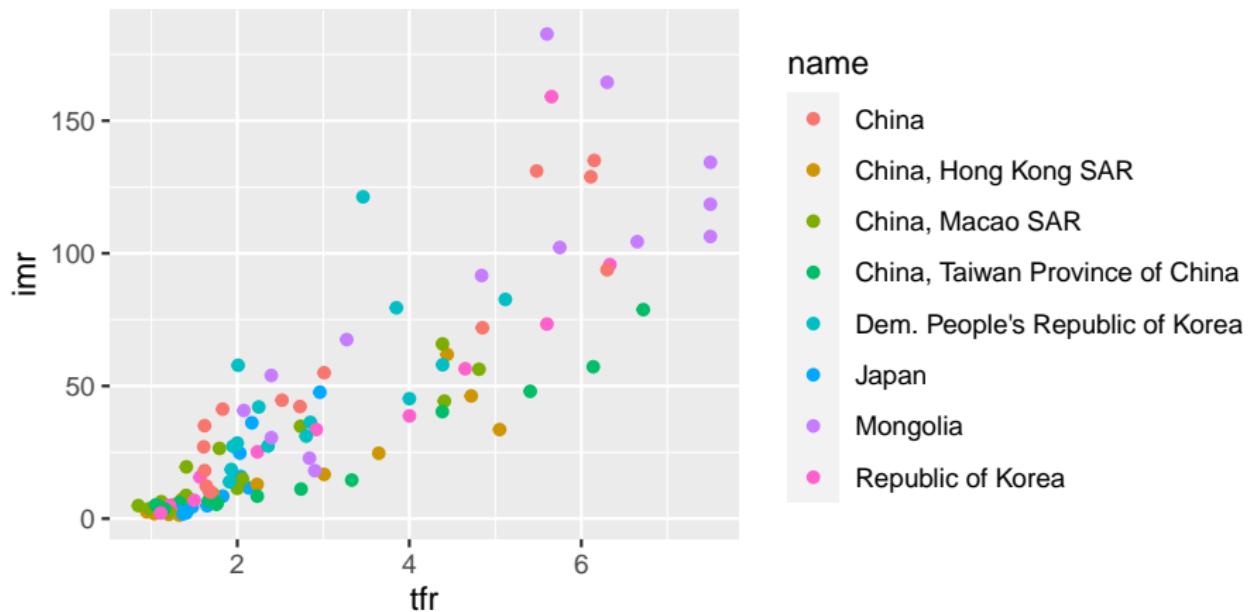
# Optional Aesthetics

```
> # eastern asia: set shape from discrete variable  
> ggplot(data = east_asia, mapping = aes(x = tfr, y = imr, shape = name)) +  
+   geom_point()
```



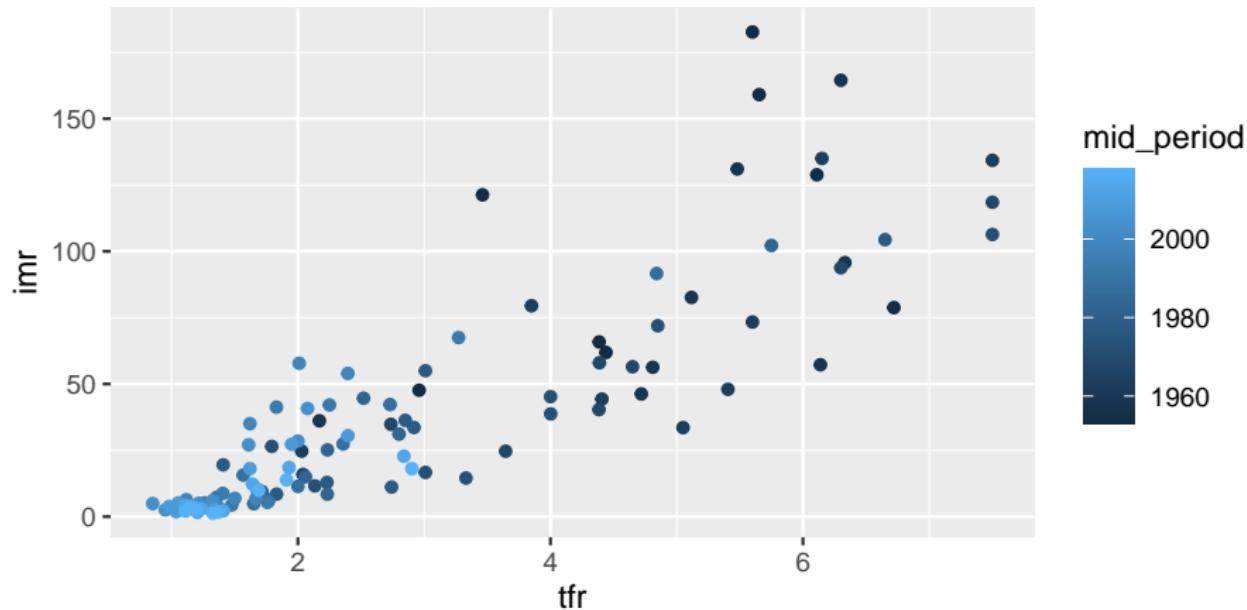
# Optional Aesthetics

```
> # eastern asia: set colour from continuous variable  
> ggplot(data = east_asia, mapping = aes(x = tfr, y = imr, colour = name)) +  
+   geom_point()
```



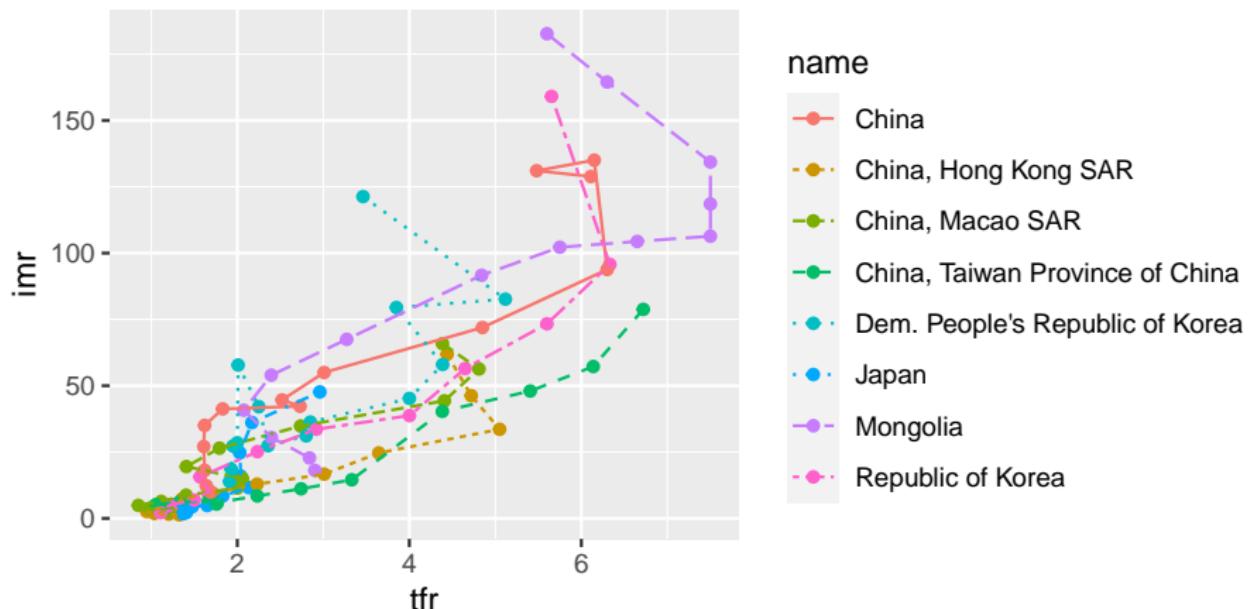
# Optional Aesthetics

```
> # eastern asia: set colour from continuous variable  
> ggplot(data = east_asia, mapping = aes(x = tfr, y = imr, colour = mid_period)) +  
+   geom_point()
```



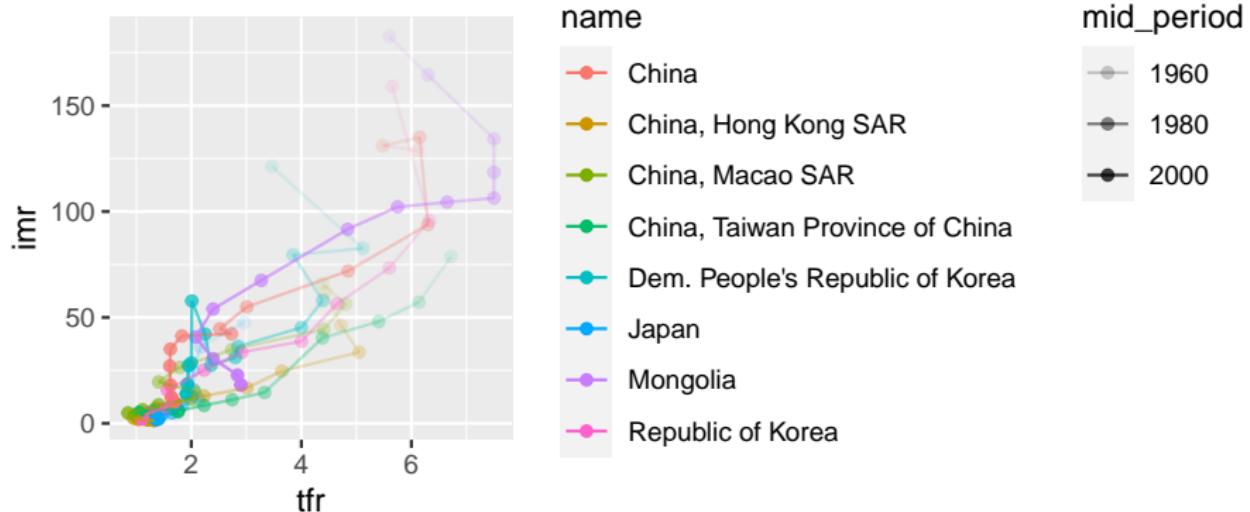
# Optional Aesthetics

```
> # eastern asia: set colour and linetype from discrete variable  
> ggplot(data = east_asia,  
+         mapping = aes(x = tfr, y = imr, colour = name, linetype = name)) +  
+     geom_point() +  
+     geom_path()
```



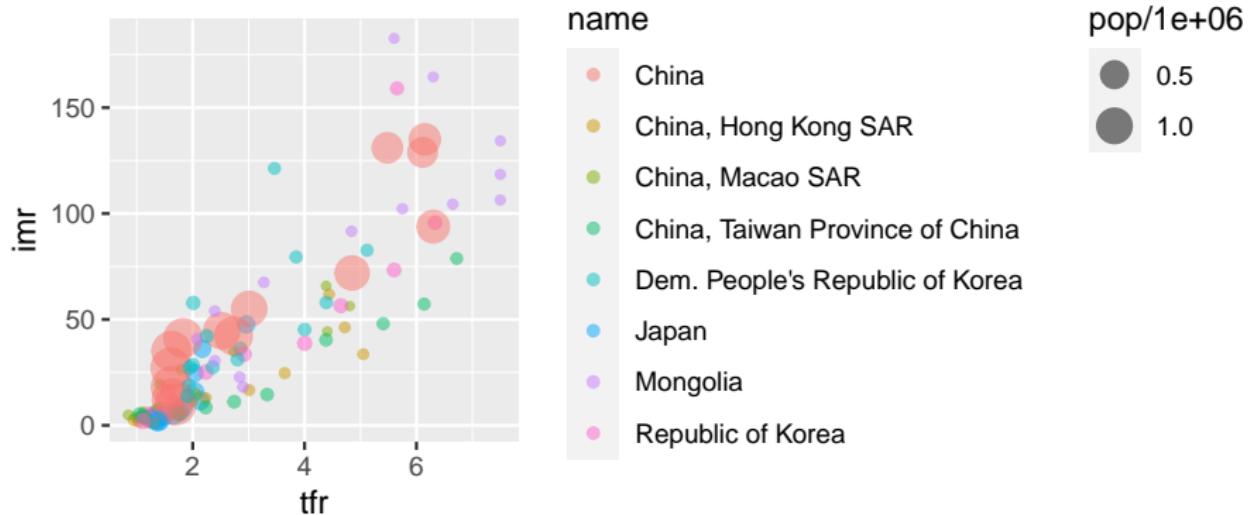
# Optional Aesthetics

```
> # eastern asia: set colour from discrete and alpha from continuous variable  
> ggplot(data = east_asia,  
+         mapping = aes(x = tfr, y = imr, colour = name, alpha = mid_period)) +  
+     geom_point() +  
+     geom_path() +  
+     # to fit legend on my slides  
+     theme(legend.box = "horizontal")
```



# Optional Aesthetics

```
> # eastern asia: set colour and size from data, set alpha to 0.5 for all.  
> ggplot(data = east_asia,  
+         mapping = aes(x = tfr, y = imr, colour = name, size = pop/1e6)) +  
+         # all points with same transparency  
+         geom_point(alpha = 0.5) +  
+         # to fit legend on my slides  
+         theme(legend.box = "horizontal")
```

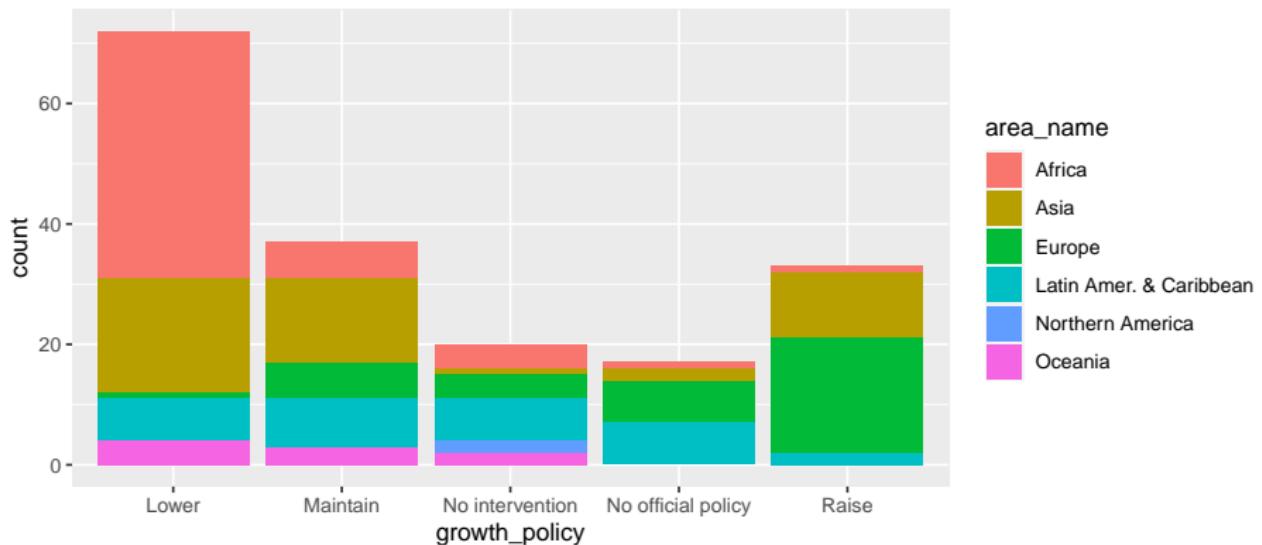


# Position

- Each geoms has three main arguments (`mapping`, `position`, `stat`)
- The position adjustments determine how to arrange geoms that would otherwise occupy the same space.
- For most geoms you are unlikely to want to change from the default.
- However, for `geom_bar()` it is useful to know
  - `stack`: stacks elements on top of one another
  - `dodge`: arrange elements side by side
  - `fill`: stack elements on top of one another, normalize height
- For other geoms there are other position adjustment possibilities
  - `jitter`: adds random noise to x and y position of each element to avoid over plotting
  - `nudge`: nudge labels (in `geom_label()`) away from points

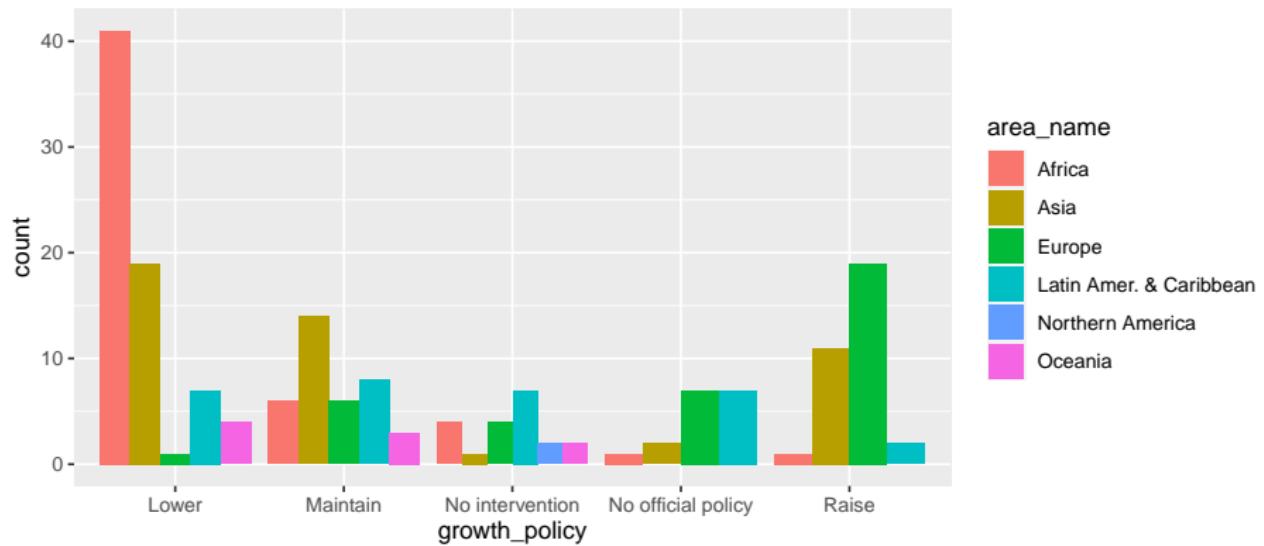
# Position

```
> # stack position (default for geom_bar())
> ggplot(data = d, mapping = aes(x = growth_policy, fill = area_name)) +
+   geom_bar()
```



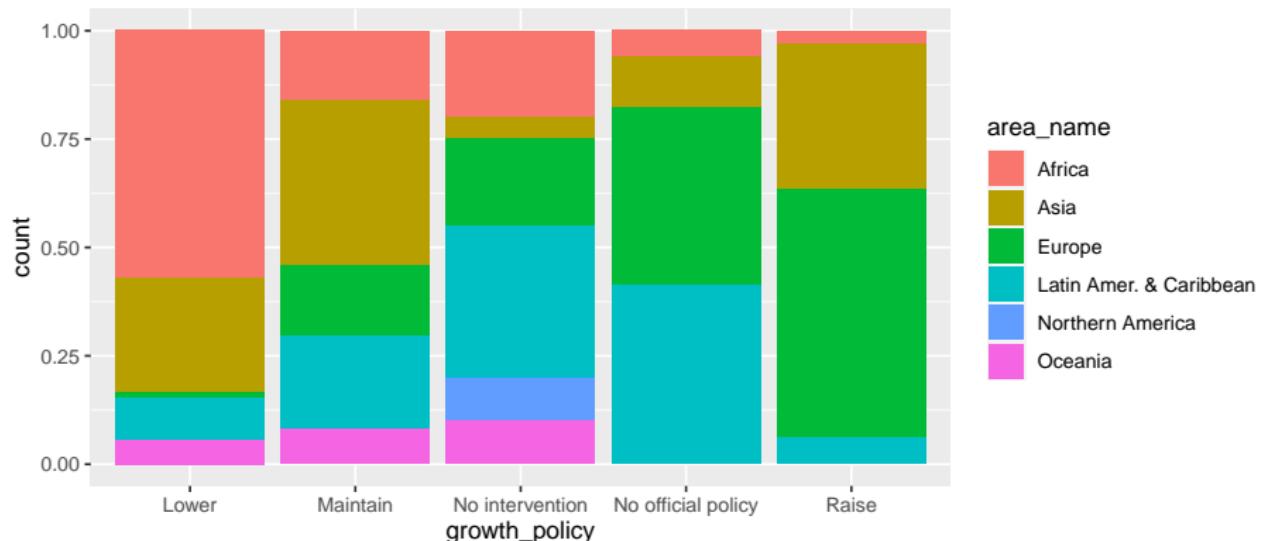
# Position

```
> # dodge position  
> ggplot(data = d, mapping = aes(x = growth_policy, fill = area_name)) +  
+   geom_bar(position = "dodge")
```



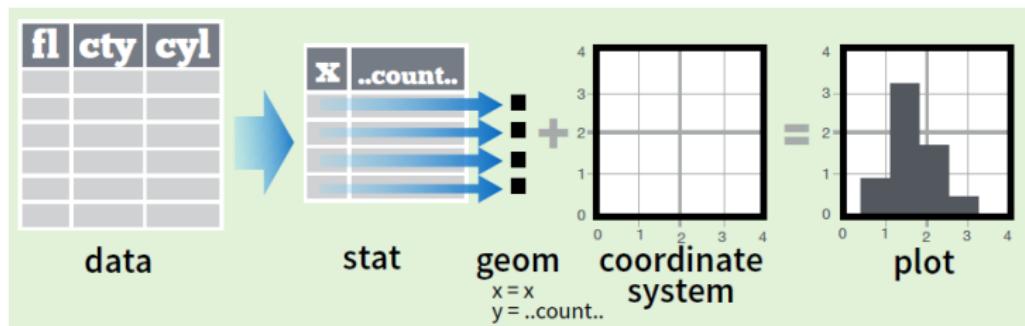
# Position

```
> # fill position  
> ggplot(data = d, mapping = aes(x = growth_policy, fill = area_name)) +  
+   geom_bar(position = "fill")
```



# Stats

- Each geom has three main arguments (`mapping`, `position`, `stat`)
- The `stat` argument builds new variables to plot (e.g., `count`, `prop`).
  - Short for statistical transformation.
- Each geom is associated with a default stat that it uses to calculate values to plot.
  - Applied the transformation and stores the results behind the scenes.
  - Uses an intuitive set of defaults.
  - Rarely need to adjust a geom stat.



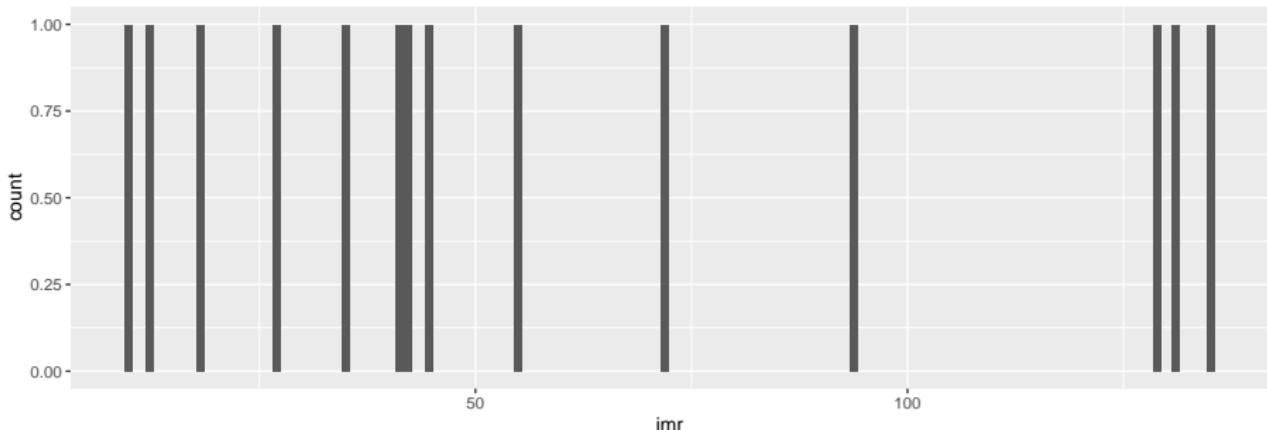
# Stats

- Some stats we have already seen in action:
  - `identity`: use raw values, e.g. `geom_point()`
  - `count`: computes two new variables, `count` and `proportion`, e.g. `geom_bar()`
  - `bin`: arrange data into bins and counts, e.g. `geom_histogram()`
  - `density`: computes kernel density estimates, e.g. `geom_density()`
  - `smooth`: fit a model to your data and then plot the model line,  
e.g. `geom_smooth()`
  - `boxplot`: computes box plots statistics (min, max, IQR, median),  
e.g. `geom_boxplot()`
- For most geoms you are unlikely to want to change from the default.
  - Show how to change with `geom_bar()` in to better understand...

# Stats

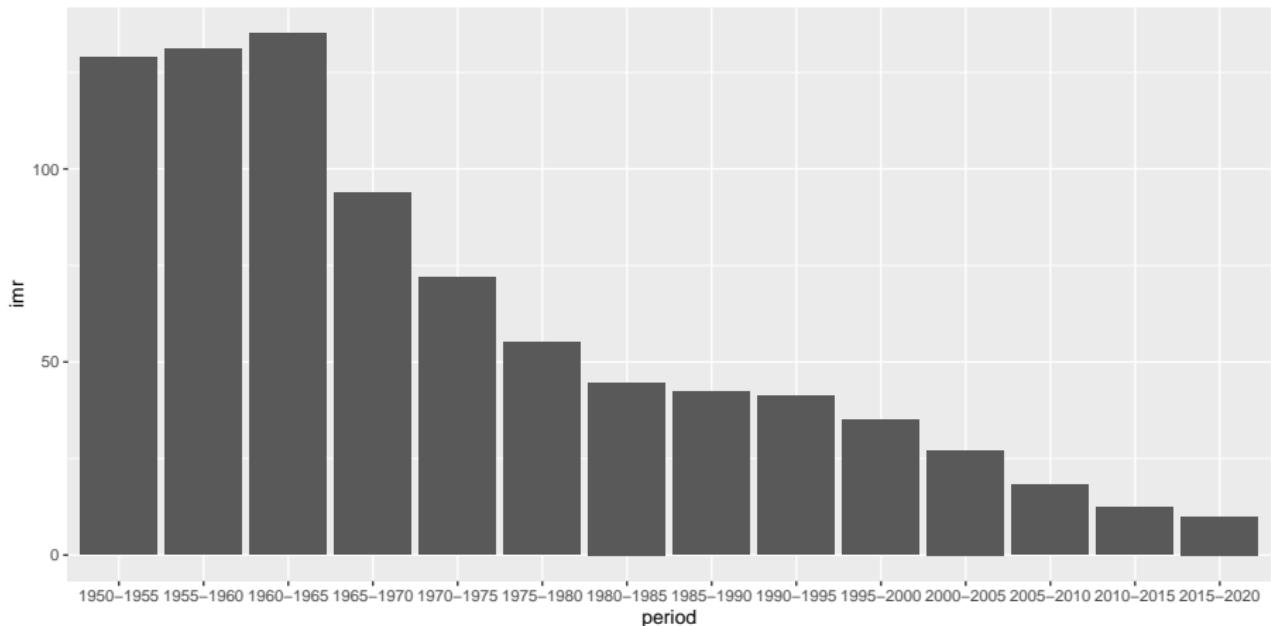
```
> # how can we plot a bar for imr in every period?... this gives an error...
> ggplot(data = china, mapping = aes(x = period, y = imr)) +
+   geom_bar()
Error: stat_count() can only have an x or y aesthetic.
```

```
> # this is not what we want...
> ggplot(data = china, mapping = aes(x = imr)) +
+   geom_bar()
```



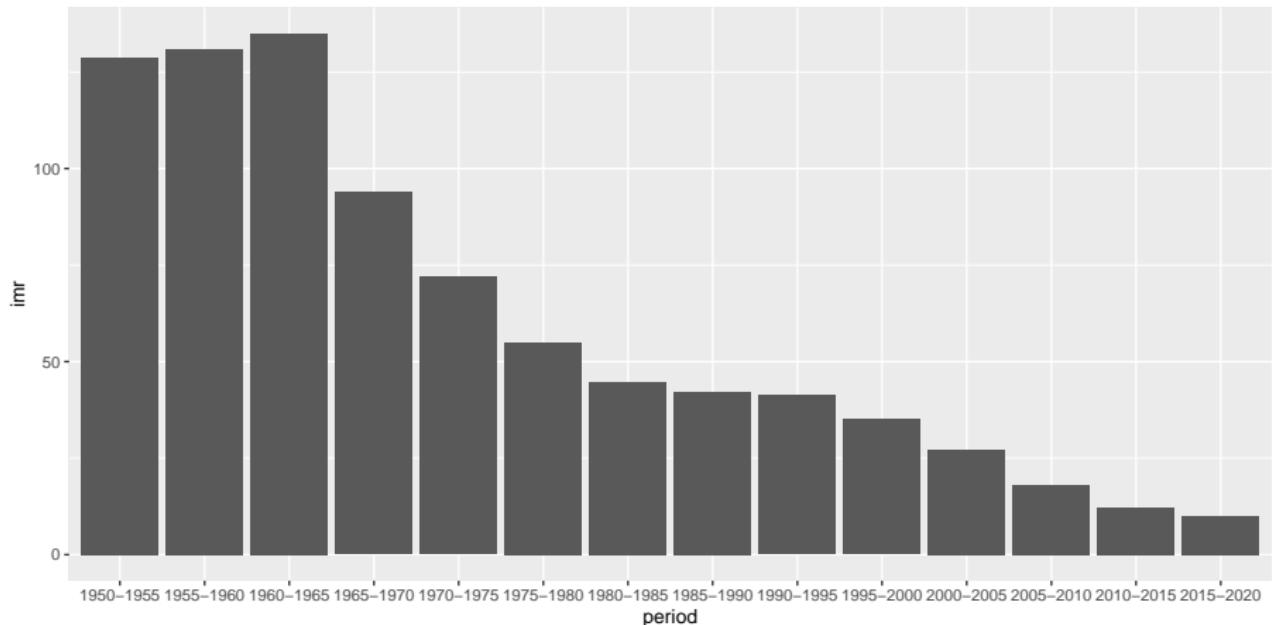
# Stats

```
> # china: identity stat  
> ggplot(data = china, mapping = aes(x = period, y = imr)) +  
+   geom_bar(stat = "identity")
```



# Stats

```
> # geom_bar() with identity stat is the same as geom_col()  
> ggplot(data = china, mapping = aes(x = period, y = imr)) +  
+   geom_col()
```

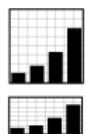


## Exercise 2 (ex22.R)

```
# 0. a) Check your working directory is in the course folder. Load .Rproj file if needed  
getwd()  
#     b) Load the tidyverse package (which loads the ggplot2 package amongst others)  
library(tidyverse)  
#     c) Run the code in ex22_prelim.R to import the UN data for this exercise  
source(#####)  
##  
##  
##  
# 1. Create a scatter plot of infant mortality rates (x) against total fertility rate (y)  
#     a) point colours matching area names  
#     b) point sizes matching the population size in millions (divide pop by 1000)  
  
# 2. Adapt the plot above to add  
#     a) horizontal line where tfr is 2.1  
#     b) all points to have a transparency of 0.5  
  
# 3. Create a plot of paths for the infant mortality rate (x) against total fertility rate (y)  
#     relationship for South East Asian countries in se_asia using  
#     a) separate paths for each country  
#     b) transparency to change by the years in the mid_period column
```

# Coordinates

- Plots are based on a coordinate system.
  - The `ggplot()`, like most software, uses the Cartesian coordinate system by default
- Can change the coordinates system of easily using a `coord_` function with `ggplot()`



`r <- d + geom_bar()`  
`r + coord_cartesian(xlim = c(0, 5))`

xlim, ylim  
The default cartesian coordinate system



`r + coord_fixed(ratio = 1/2)`

ratio, xlim, ylim  
Cartesian coordinates with fixed aspect ratio between x and y units



`r + coord_flip()`  
xlim, ylim  
Flipped Cartesian coordinates



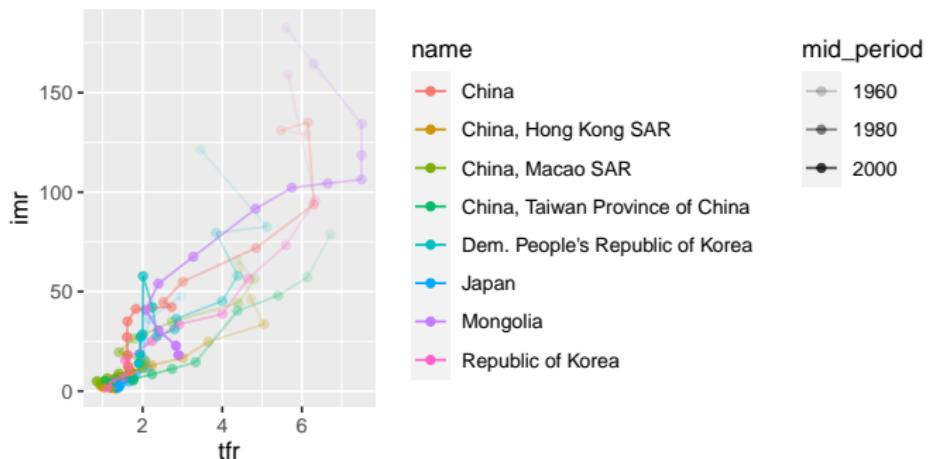
`r + coord_polar(theta = "x", direction=1)`  
theta, start, direction  
Polar coordinates

`r + coord_trans(ytrans = "sqrt")`

xtrans, ytrans, limx, limy  
Transformed cartesian coordinates. Set xtrans and ytrans to the name of a window function.

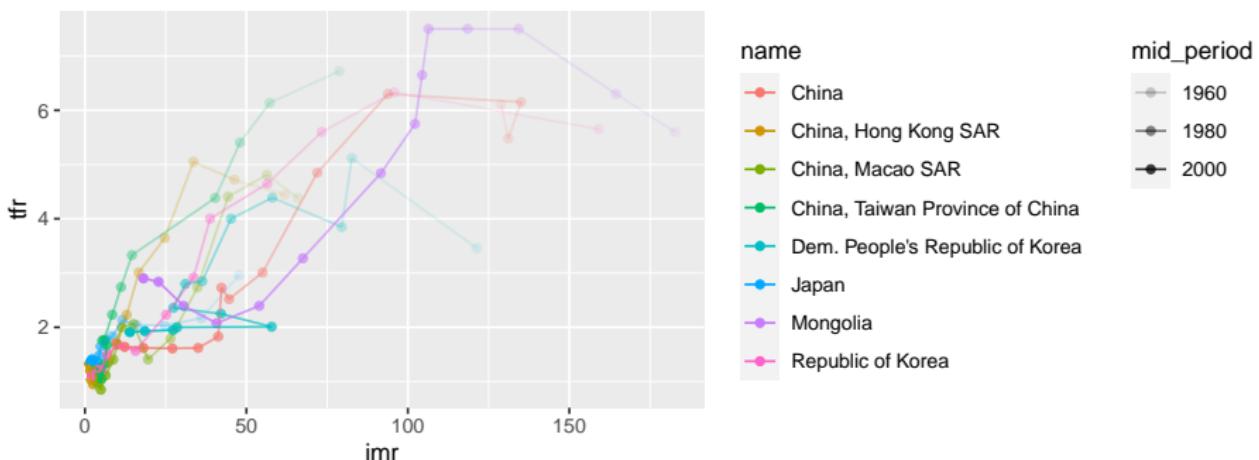
# Coordinates

```
> # eastern asia: fix coordinates: 1 coord unit of imr (y) same as 0.05 tfr (x)
> ggplot(data = east_asia,
+         mapping = aes(x = tfr, y = imr, colour = name, alpha = mid_period)) +
+         geom_point() +
+         geom_path() +
+         coord_fixed(ratio = 0.05) +
+         # to fit legend on my slides
+         theme(legend.box = "horizontal")
```



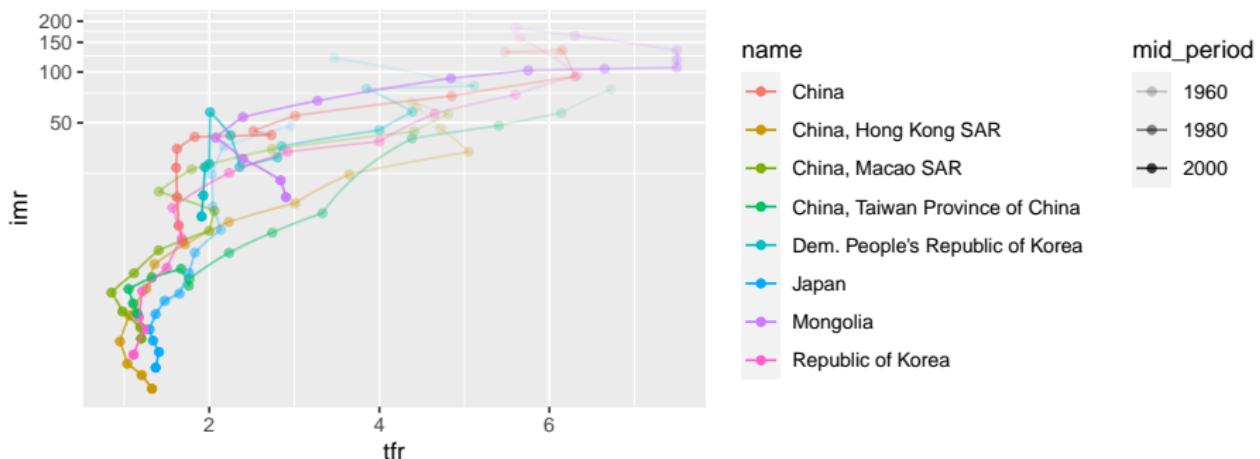
# Coordinates

```
> # eastern asia flip coordinates: x is y, y is x
> ggplot(data = east_asia,
+         mapping = aes(x = tfr, y = imr, colour = name, alpha = mid_period)) +
+         geom_point() +
+         geom_path() +
+         coord_flip() +
+         # to fit legend on my slides
+         theme(legend.box = "horizontal")
```



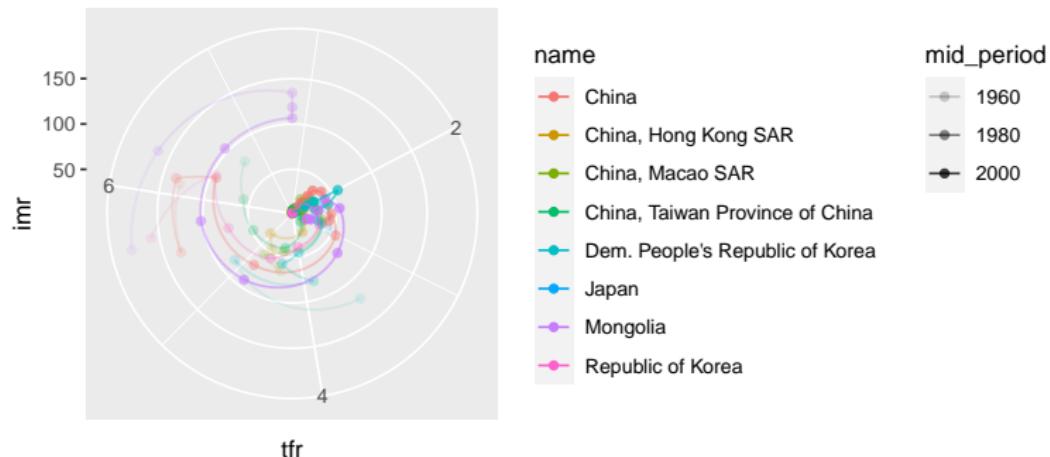
# Coordinates

```
> # log y axes
> ggplot(data = east_asia,
+         mapping = aes(x = tfr, y = imr, colour = name, alpha = mid_period)) +
+     geom_point() +
+     geom_path() +
+     coord_trans(y = "log") +
+     # to fit legend on my slides
+     theme(legend.box = "horizontal")
```



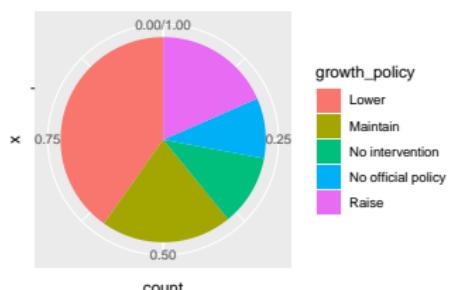
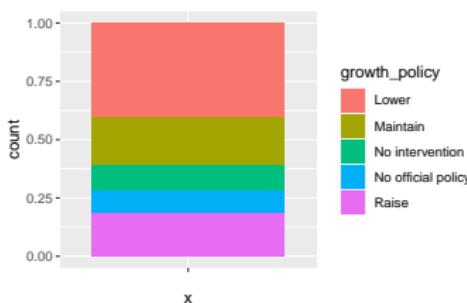
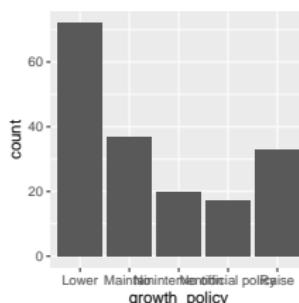
# Coordinates

```
> # polar coordinates: x is now the angle in circle and y is distance from centre
> ggplot(data = east_asia,
+         mapping = aes(x = tfr, y = imr, colour = name, alpha = mid_period)) +
+         geom_point() +
+         geom_path() +
+         coord_polar() +
+         # to fit legend on my slides
+         theme(legend.box = "horizontal")
```



# Pie Charts

```
> ggplot(data = d, mapping = aes(x = growth_policy)) +
+   geom_bar()
>
> ggplot(data = d, mapping = aes(x = "", fill = growth_policy)) +
+   geom_bar(position = "fill")
>
> ggplot(data = d, mapping = aes(x = "", fill = growth_policy)) +
+   geom_bar(position = "fill") +
+   coord_polar(theta = "y")
```



# Facets

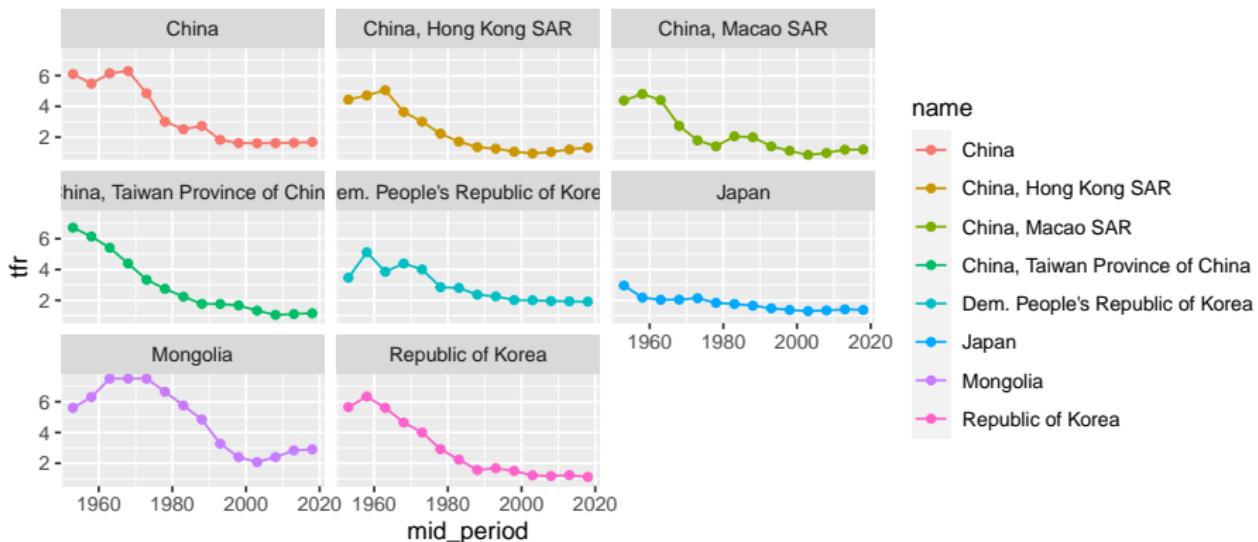
- Facets divide a plot into subplots based on the values of one or more discrete variables
  - Subplots can be a very effective way to compare across discrete (categorical) values
- Two functions:
  - `facet_wrap()` wraps a sequence of panels defined by a discrete variable(s) onto a page in roughly rectangular form.
  - `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.

# Facets

- Define facet layouts using arguments depending on function
  - `facet_wrap()` use `facet = "a"` or `facet = vars(a)`
  - `facet_grid()` use `row = vars(a), col = vars(b)`
- Other arguments are similar
  - `nrow, ncol`: Number of rows and columns.
  - `scales`: fixes or frees scales in facets
    - "fixed": both scales fixed (default)
    - "free": both scales free
    - "free\_x" or "free\_y": scales free in one dimension.
  - `labeller`: adjust facet strip labels,
    - for example `labeller = label_wrap_gen(10)` to make sure each facet strip label is only a maximum 10 characters wide.

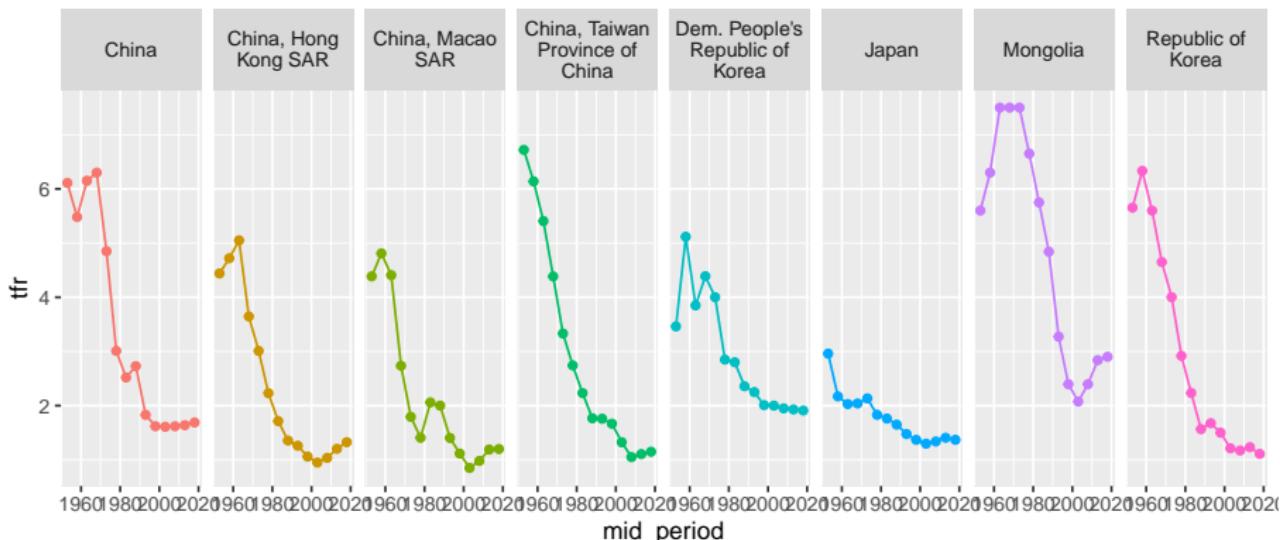
# Facets

```
> # eastern asia: facet_wrap() function for each country (name)
> ggplot(data = east_asia, mapping = aes(x = mid_period, y = tfr, colour = name)) +
+   geom_point() +
+   geom_line() +
+   facet_wrap(facets = "name") +
+   # to fit legend on my slides
+   theme(legend.box = "horizontal")
```



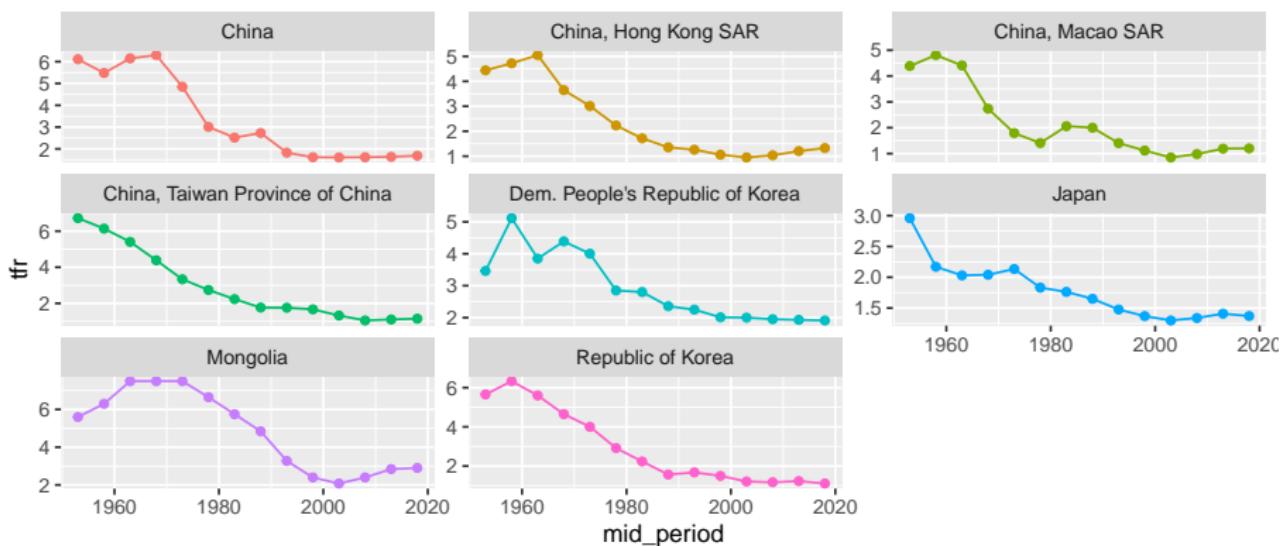
# Facets

```
> # facet_wrap() function using one row and restrict label width
> ggplot(data = east_asia, mapping = aes(x = mid_period, y = tfr, colour = name)) +
+   geom_point() +
+   geom_line() +
+   facet_wrap(facets = vars(name), nrow = 1, labeller = label_wrap_gen(15)) +
+   # drop country name guide
+   guides(colour = FALSE)
```



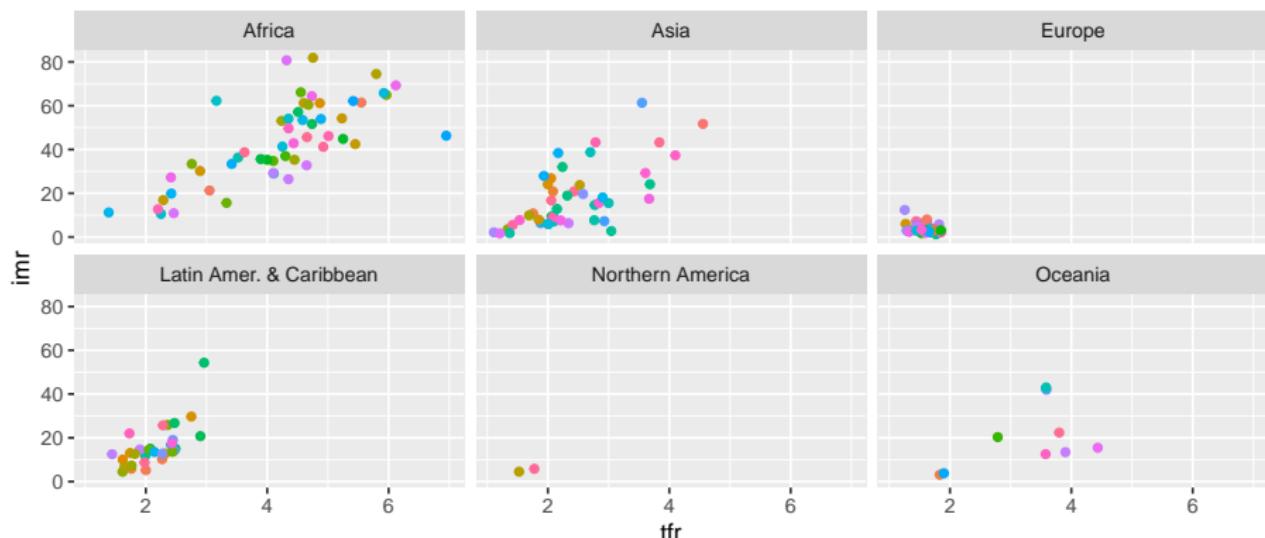
# Facets

```
> # facet_wrap() function with different y axis limits
> ggplot(data = east_asia, mapping = aes(x = mid_period, y = tfr, colour = name)) +
+   geom_point() +
+   geom_line() +
+   facet_wrap(facets = vars(name), scales = "free_y") +
+   guides(colour = FALSE)
```



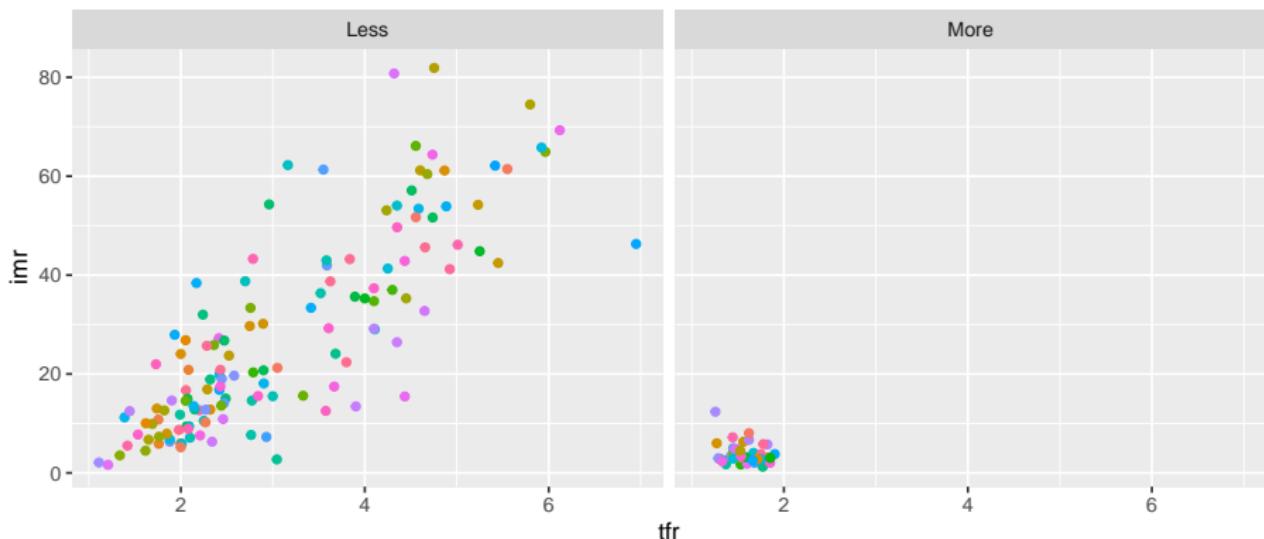
# Facets

```
> # facet_wrap() function scatter plot
> ggplot(data = d,
+         mapping = aes(x = tfr, y = imr, colour = name)) +
+     geom_point() +
+     facet_wrap(facets = vars(area_name)) +
+     guides(colour = FALSE)
```



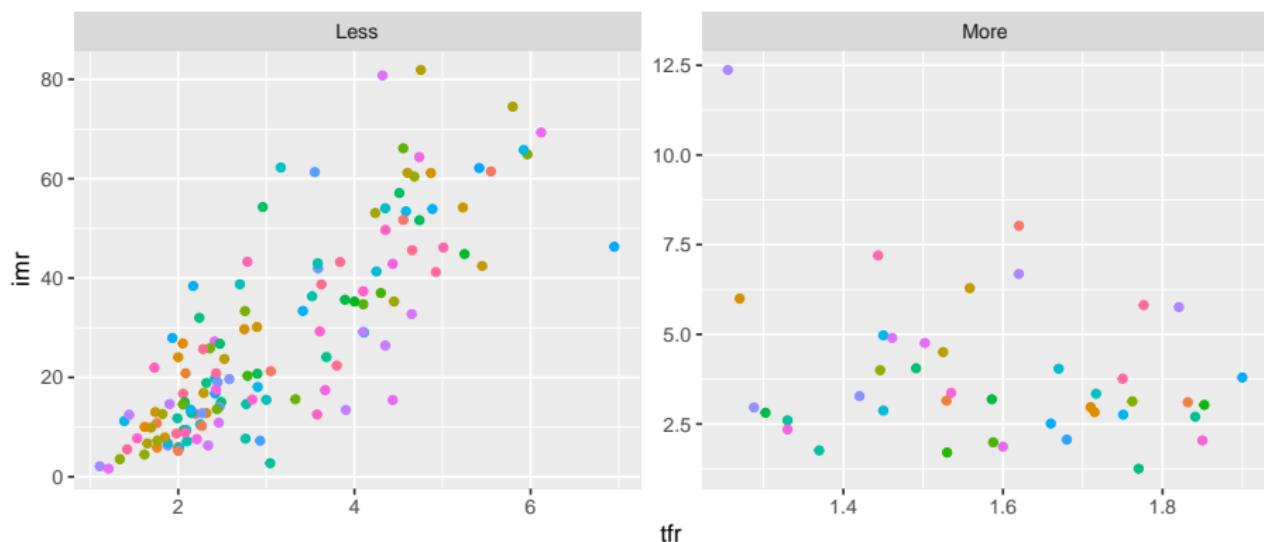
# Facets

```
> # facet_wrap() function scatter plot using development level
> ggplot(data = d,
+         mapping = aes(x = tfr, y = imr, colour = name)) +
+     geom_point() +
+     facet_wrap(facets = vars(developed)) +
+     guides(colour = FALSE)
```



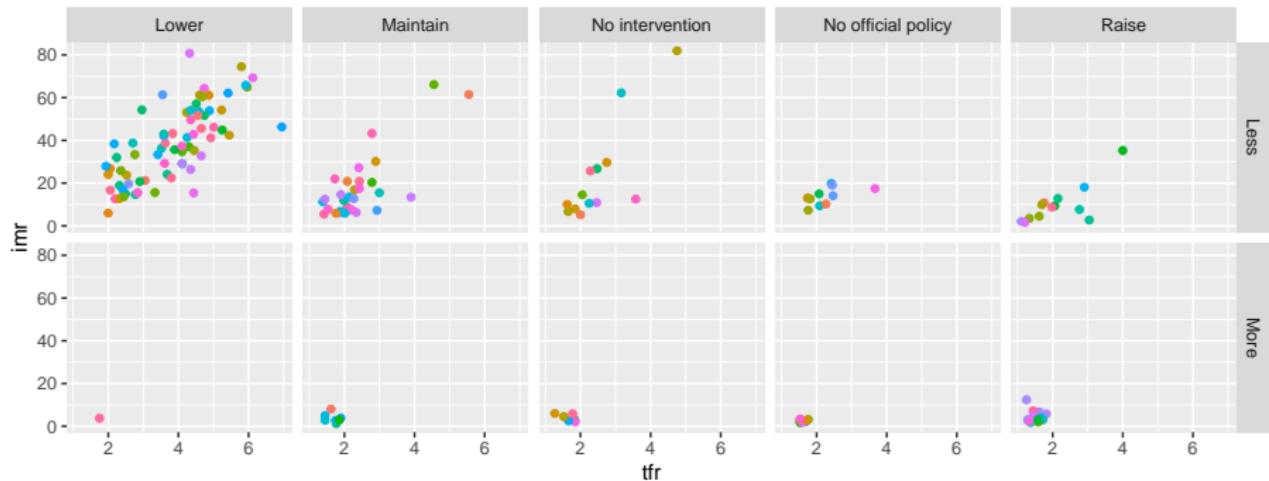
# Facets

```
> # all data: facet_wrap() function scatter plot with x and y scales free
> ggplot(data = d,
+         mapping = aes(x = tfr, y = imr, colour = name)) +
+         geom_point() +
+         facet_wrap(facets = "developed", scales = "free") +
+         guides(colour = FALSE)
```



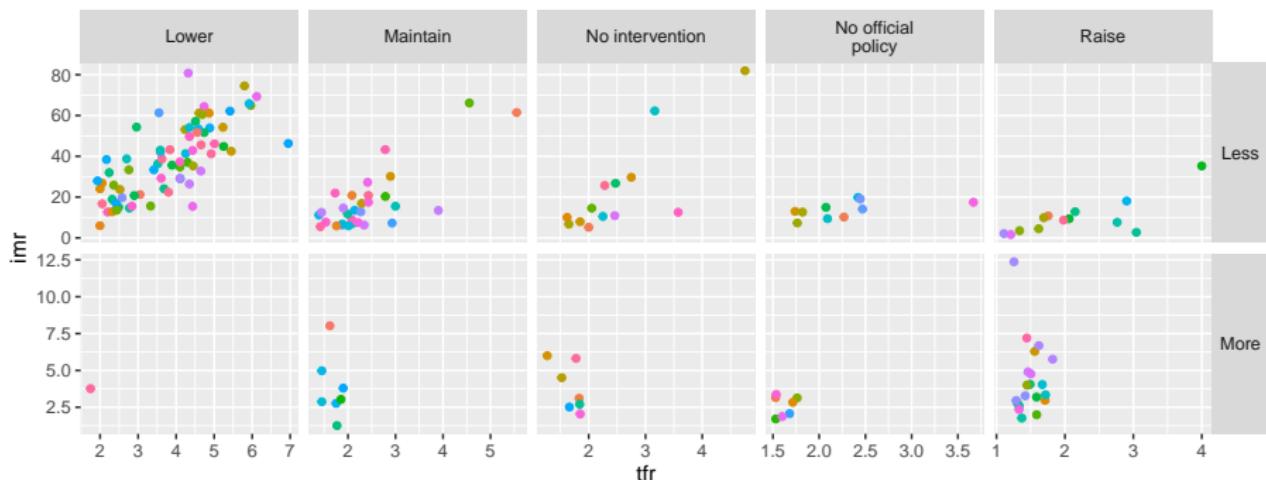
# Facets

```
> # facet_grid() function scatter plot using development status and population growth
> ggplot(data = d,
+         mapping = aes(x = tfr, y = imr, colour = name)) +
+         geom_point() +
+         facet_grid(row = vars(developed), col = vars(growth_policy)) +
+         guides(colour = FALSE)
```



# Facets

```
> # ... alter labels and free grid scales
> ggplot(data = d,
+         mapping = aes(x = tfr, y = imr, colour = name)) +
+     geom_point() +
+     # rotate growth policy labels to fit legend on my slides
+     facet_grid(row = vars(developed), col = vars(growth_policy),
+                scales = "free", labeller = label_wrap_gen(18)) +
+     guides(colour = FALSE) +
+     # rotate growth policy labels to fit legend on my slides
+     theme(strip.text.y = element_text(angle = 0))
```



## Exercise 3 (ex23.R)

```
# 0. a) Check your working directory is in the course folder. Load .Rproj file if needed  
getwd()  
# b) Load the tidyverse package (which loads the ggplot2 package amongst others)  
  
# c) Run the code in ex23_prelim.R to import the UN data for this exercise  
source("./exercise/ex23_prelim.R")  
##  
##  
##  
# 1. Create a pie chart of the immigration policies in different countries from different areas  
  
# 2. Adapt the code from the question above to create pie charts for the growth_policy column  
# a. separate facets for each area  
# b. no axis  
# (Hint: use theme_void() to drop axis)  
  
# 3. Uncomment the solution to Question 1 of ex22 and then adapt so that the infant mortality rate  
# transformed to the log10 scale  
# qplot(data = d, mapping = aes(x = imr, y = tfr, colour = area name, size = pop/100000000))
```

# Scales

- The scale functions change aesthetics from their default.
- Each aesthetic has its own functions. Take a general form:
  - `scale_*_continuous()`: alter aesthetics based on continuous variables
  - `scale_*_discrete()`: alter aesthetics based on discrete variables
  - `scale_*_manual(values = c())`: map discrete values to manually chosen visual ones

## Color and fill scales (Discrete)

```
n + scale_fill_brewer(palette = "Blues")
```

For palette choices: RColorBrewer::display.brewer.all()

```
n + scale_fill_grey(start = 0.2, end = 0.8, na.value  
= "red")
```

## Color and fill scales (Continuous)

```
o <- c + geom_dotplot(aes(fill = ..x..))
```

```
o + scale_fill_distiller(palette = "Blues")
```

```
o + scale_fill_gradient(low="red", high="yellow")
```

```
o + scale_fill_gradient2(low="red", high="blue",  
mid = "white", midpoint = 25)
```

```
o + scale_fill_gradientn(colours=topo.colors(6))  
Also: rainbow(), heat.colors(), terrain.colors(),  
cm.colors(), RColorBrewer::brewer.pal()
```

## Shape and size scales

```
p <- e + geom_point(aes(shape = fl, size = cyl))
```

```
+x p + scale_shape() + scale_size()
```

```
+p p + scale_shape_manual(values = c(3:7))
```

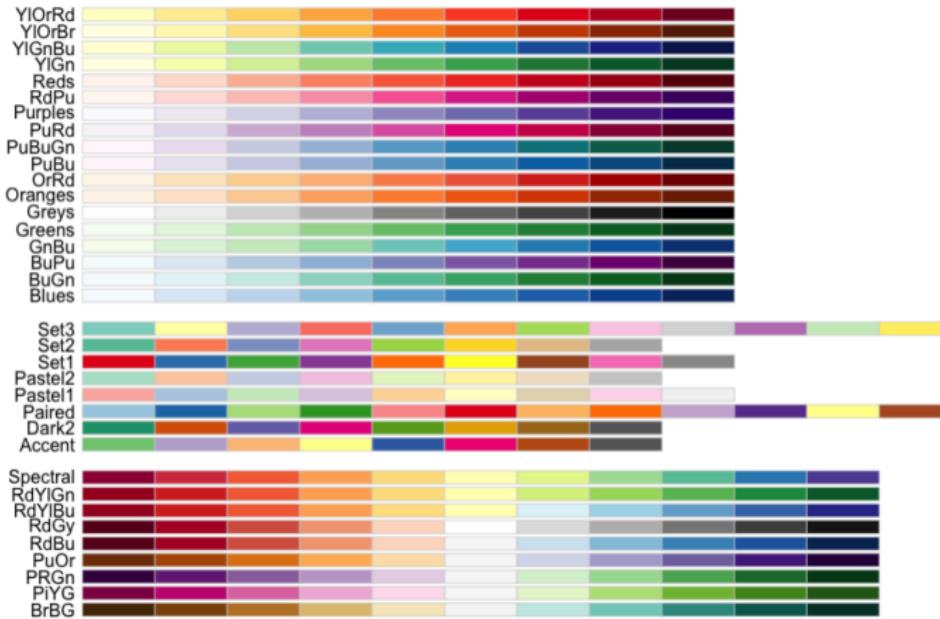
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25

□○△+×◊▽\*◆⊕☒☒☒☒●▲◆●○□◆△▽

```
p + scale_radius(range = c(1,6)) Maps to radius of  
p + scale_size_area(max_size = 6) circle, or area
```

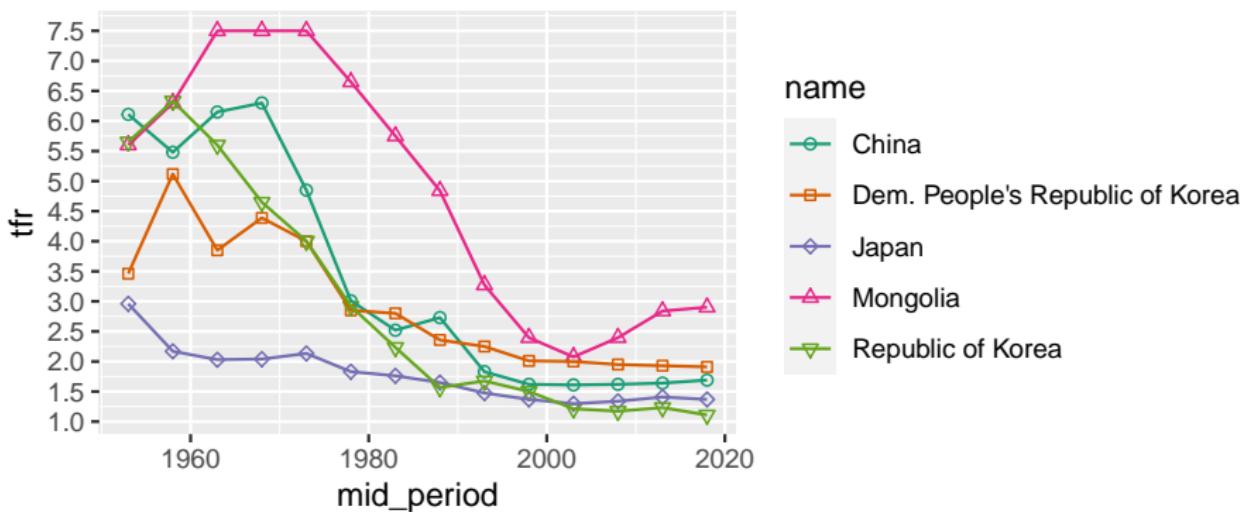
## Scales

- Colour schemes from ColorBrewer that are loaded with ggplot2
    - <http://colorbrewer2.org>



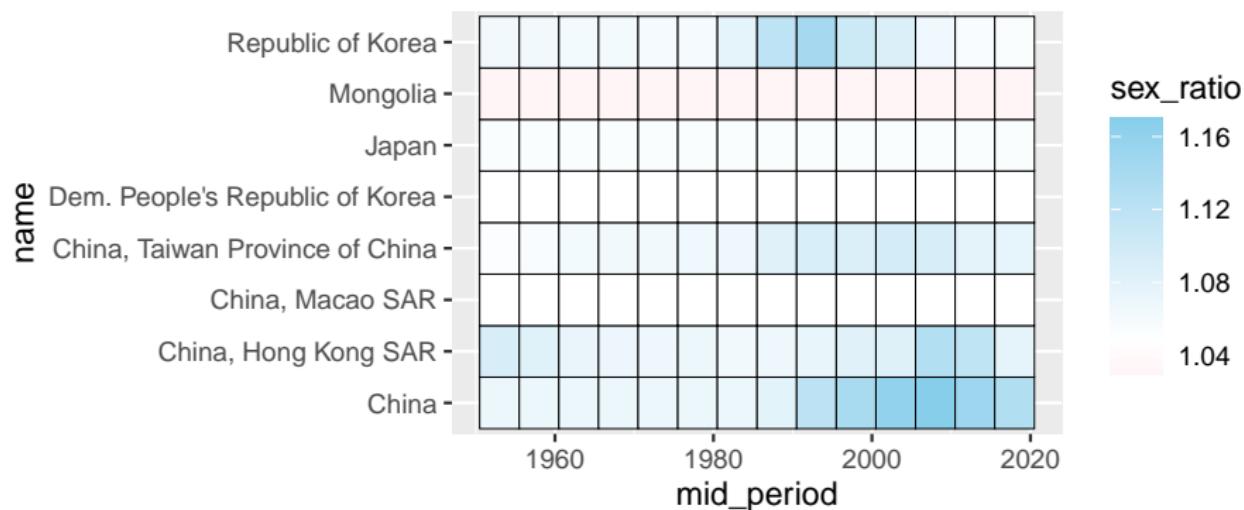
# Scales

```
> # eastern asia without Hong Kong, Macau and Taiwan
> ggplot(data = east_asia %>%
+         filter(!str_detect(string = name, pattern = ",,")),
+         mapping = aes(x = mid_period, y = tfr, colour = name, shape = name)) +
+     geom_point() +
+     geom_line() +
+     scale_color_brewer(palette = "Dark2") +
+     scale_shape_manual(values = 21:25) +
+     scale_y_continuous(breaks = seq(from = 1, to= 8, by = 0.5))
```



# Scales

```
> ggplot(data = east_asia, mapping = aes(x = mid_period, y = name, fill = sex_ratio)
+   geom_tile(colour = "black") +
+   scale_fill_gradient2(low="pink", high="skyblue", mid="white", midpoint=1.05)
```



# Labels

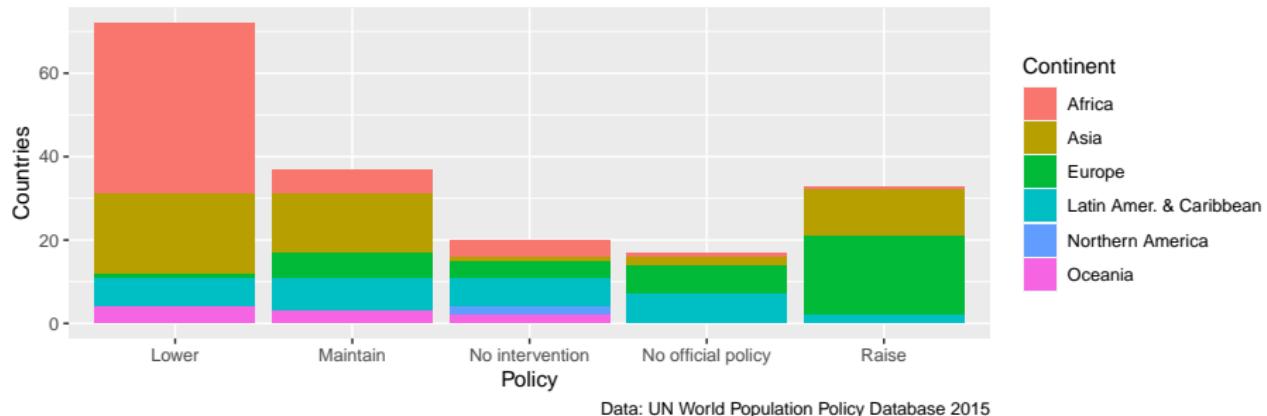
- Good labels are critical for making your plots accessible to a wider audience.
- The `labs()` function can add
  - `title`: main title
  - `subtitle`: highlight main message
  - `caption`: data source information
  - `x, y, fill, colour, alpha, ...` : full variable names for aesthetics

# Labels

```
> # all countries 2010-15 data
> ggplot(data = d, mapping = aes(x = growth_policy, fill = area_name)) +
+   geom_bar() +
+   labs(title = "Population Growth Policies",
+        subtitle = "Most African countries have lowering population growth policies",
+        caption = "Data: UN World Population Policy Database 2015",
+        x = "Policy", y = "Countries", fill = "Continent")
```

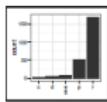
## Population Growth Policies

Most African countries have lowering population growth policies. Most European countries have rising population policies.



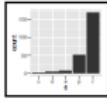
# Themes

- The theme functions customize the “look” of plots
  - Changes how the plot looks without changing the information that the plot displays.
- There are eight theme functions in ggplot2
- The ggthemes package has many themes to match publication styles e.g. Economist or FiveThirtyEight



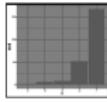
`r + theme_bw()`

White background  
with grid lines



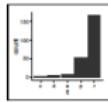
`r + theme_gray()`

Grey background  
(default theme)



`r + theme_dark()`

dark for contrast



`r + theme_classic()`

`r + theme_light()`

`r + theme_linedraw()`

`r + theme_minimal()`

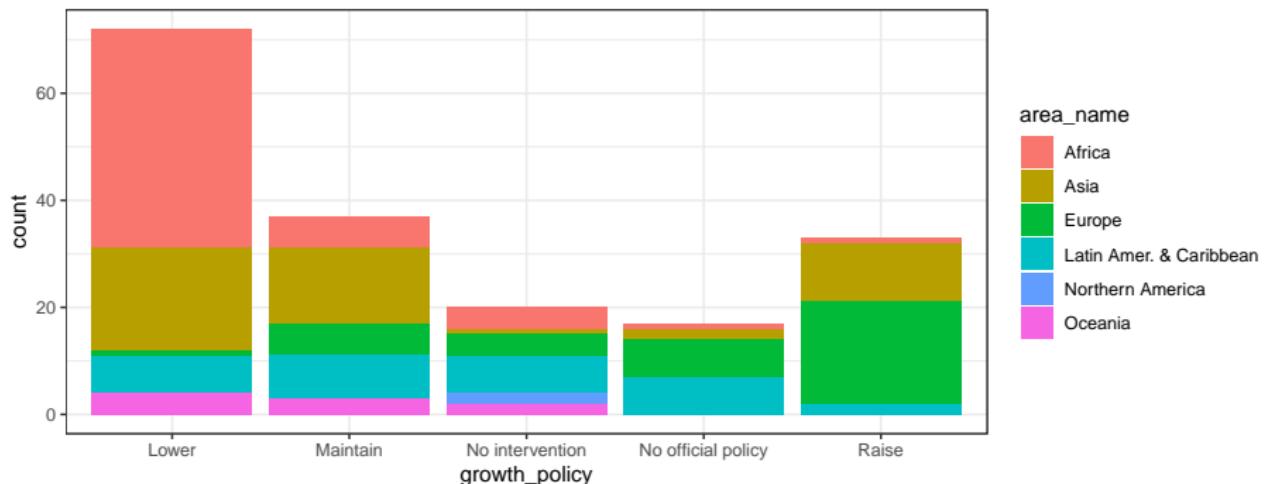
Minimal themes

`r + theme_void()`

Empty theme

# Themes

```
> # change theme  
> ggplot(data = d, mapping = aes(x = growth_policy, fill = area_name)) +  
+   geom_bar() +  
+   theme_bw()
```



# Saving

- Can assign a plot to an R object
  - Can then keep adding to the plot with +

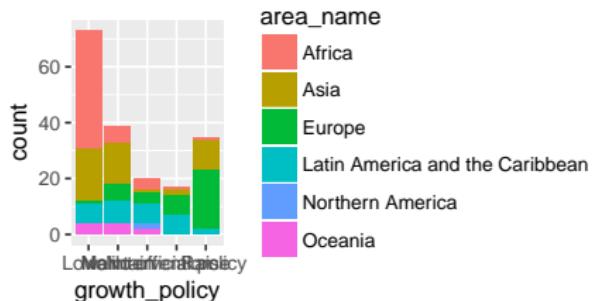
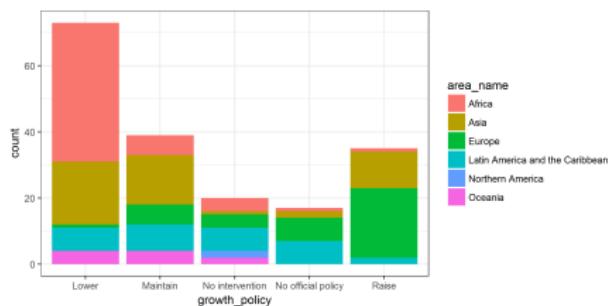
```
> # save a plot to an object
> g <- ggplot(data = d, mapping = aes(x = growth_policy, fill = area_name)) +
+   geom_bar()
> # plot g with the black and white theme
> g + theme_bw()
```

- The ggsave() function to save plot from
  - filename: matches file type to file extension, e.g. filename = myplot.pdf will create a PDF. If no path is given in the file name, the file will appear in the current working directory (getwd()).
  - width, height, units: plot size, uses RStudio window size by default
  - scale: scales the size of the plots objects

```
> # saves the last plot from RStudio as a PNG
> ggsave(filename = "myplot.png", width = 10, height = 5, unit = "cm", scale = 2)
> # saves the plot from object g as a PDF
> ggsave(filename = "myplot.pdf", width = 10, height = 5, unit = "cm", plot = g)
```

# Scale

- Scale alters the point size:
  - Left: PNG with scale = 2
  - Right: PDF with scale = 1 (default)



# Combining Plots

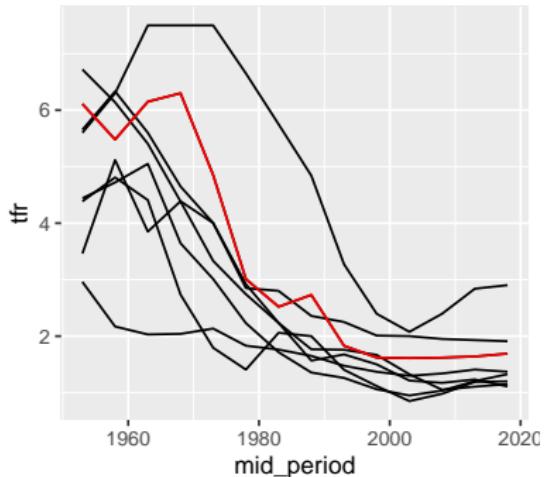
- Use the patchwork package to combine plots
- Simple syntax to arrange plots
  - + to place side by side
  - / to place on top of each other
  - () to combine plots on a row
  - plot\_layout(widths = c(2, 1)) to control layout options

```
> g1 <- ggplot(data = east_asia,  
+                 mapping = aes(x = mid_period, y = tfr, group = name)) +  
+     geom_line() +  
+     geom_line(data = china, colour = "red") +  
+     labs(title = "East Asia TFR (China, Red)")  
> g2 <- ggplot(data = china,  
+                 mapping = aes(x = imr, y = tfr, colour = mid_period)) +  
+     geom_point() +  
+     geom_path() +  
+     labs(title = "China TFR vs IMR Trace")
```

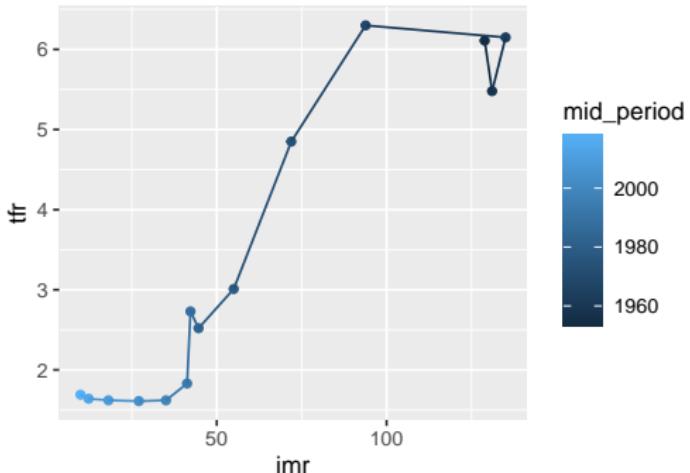
# Combining Plots

```
> library(patchwork)  
> g1 + g2
```

East Asia TFR (China, Red)



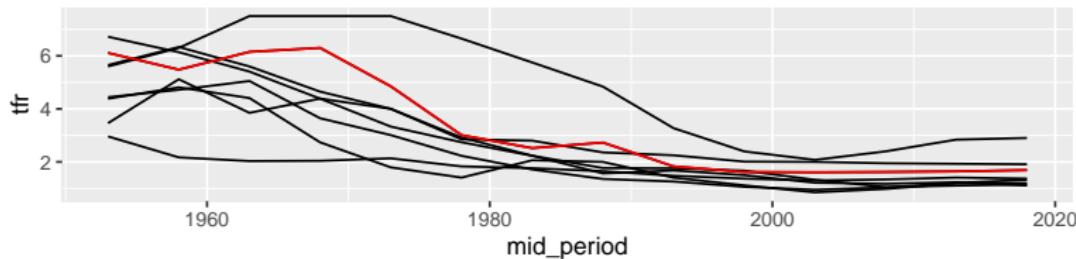
China TFR vs IMR Trace



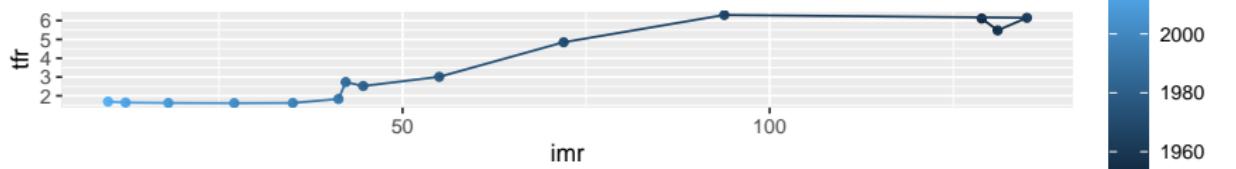
# Combining Plots

```
> g1 / g2 + plot_layout(heights = c(2, 1))
```

East Asia TFR (China, Red)

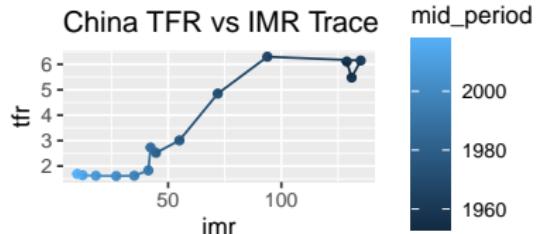
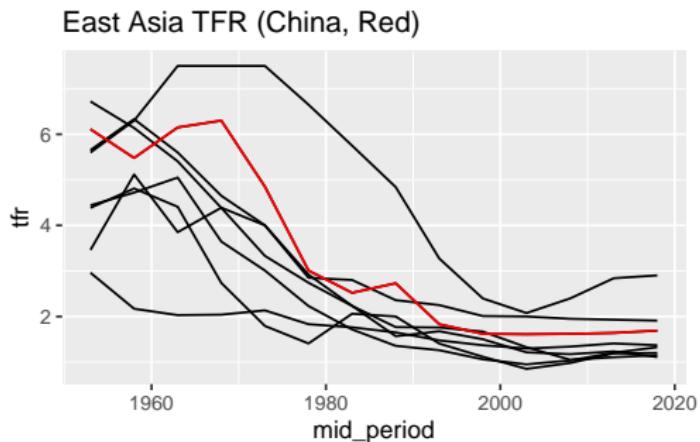


China TFR vs IMR Trace



# Combining Plots

```
> layout <- "
+ AAB
+ AA#
+ "
> g1 + g2 + plot_layout(design = layout)
```





# RStudio Cheatsheet

## Basics

**ggplot2** is based on the **grammar of graphics**, the idea that you can build every graph from the same components. It includes **data**, **geoms**, **aesthetics**, **stat**, **position** and **geom**s—visual marks that represent data points.



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.



Complete the template below to build a graph.

```
ggplot(data = <DATA> +
  <GEOM_FUNCTION>(<mapping> = aes(<MAPPINGS>),
  stat = <STAT>, position = <POSITION>) +
  <COORD_FUNCTION> +
  <FACET_FUNCTION> +
  <SCALE_FUNCTIONS> +
  <THEME_FUNCTION>)
```

**ggplot(data = mpg, aes = cty ~ hwy)}** Begins a plot that you finish by adding layers to. Add one geom function per layer.

```
+> geom_point()
```

**spike[=cty] ~ hwy**: data = mpg, geom = "point")  
Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

**last\_plot()** Returns the last plot

**ggsave("plot.png", width = 5, height = 5)** Saves last plot at 5x5 inches to file plot.png in working directory.  
Hatches file type to file extension.

## Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

### GRAPHICAL PRIMITIVES

```
a + geom_rect(mapping = aes(fill = cty), nudge_x = 1,
  nudge_y = 1, check_overlap = TRUE) x, y, label,
  alpha, angle, color, curvature, linetype, size
```

```
b + geom_curve(aes(x0 = lat + 1,
  xend = long + 1, curvature = -1), x, y, end,
  alpha, angle, color, curvearrow, group, size)
```

```
c + geom_path(inend = "butt", linejoin = "round",
  x, y, alpha, color, group, size)
```

```
d + geom_polygon(aes(group = group))
  x, y, alpha, color, fill, group, linetype, size
```

```
e + geom_rect(mapping = long ~ ymax, xmin =
  long - 1, ymin = long - 1 ~ ymax, xmn, ymax,
  ymn, fill, alpha, color, group, size)
```

```
a + geom_ribbon(aes(group = unemploy = 900,
  fill = unemploy * 900)) x, y, max, ymn,
  alpha, color, fill, group, linetype, size
```

### LINE SEGMENTS

```
common aesthetics: x, y, alpha, color, linetype, size
```

```
b + geom_abline(aes(slope = 0, intercept = 0), x,
  y, alpha, color, fill, linetype, size)
```

```
b + geom_hline(aes(intercept = 1), x,
  y, alpha, color, fill, linetype, size)
```

```
b + geom_segment(aes(endx = lat + 1, endy = long + 1))
  x, y, alpha, color, fill, group, radius = 1)
```

### ONE VARIABLE continuous

```
c + geom_area(aes(y = "bin"))
  x, y, alpha, color, fill, linetype, size
```

```
c + geom_density2d(mapping = "gaussian")
  x, y, alpha, color, fill, group, linetype, size
```

```
c + geom_dotplot()
  x, y, alpha, color, fill, group, linetype, size
```

```
c + geom_freqpoly(aes(x = x), x, y, alpha, color,
  group, linetype, size)
```

```
c + geom_histogram(mapping = 5) x, y, alpha,
  color, fill, linetype, size
```

```
c + geom_qq(aes(sample = hwy)) x, y, alpha,
  color, fill, linetype, size, weight
```

```
discrete
d + geom_bar(mapping = aes(x = hwy))
  d + geom_bar()
```

x, alpha, color, fill, group, linetype, size

### TWO VARIABLES

```
continuous x, continuous y
e + geom_rect(aes(xmin = x, xmax = x + 1,
  ymin = y, ymax = y + 1), x, y, label,
  alpha, angle, color, family, fontface, hjust =
  "center", vjust = "top")
```

```
e + geom_hex(x, y, alpha, color, fill, shape,
  size, stroke)
```

```
e + geom_quantile(x, y, alpha, color, group,
  linetype, size, weight)
```

```
e + geom_rug(sides = "tl", x, y, alpha, color,
  linetype, size)
```

```
e + geom_smooth(method = "loess", x, y, alpha,
  color, fill, group, linetype, size, weight)
```

```
e + geom_text(aes(label = cyl), x, y, alpha, color,
  family, fontface, hjust = "center", vjust =
  "middle")
```

```
discrete x, continuous y
f <- ggplot(mpg, aes(cty, hwy))
```

```
f + geom_col(x, y, alpha, color, fill, group,
  linetype, size)
```

```
f + geom_boxplot(x, y, lower, middle, upper,
  ymn, ymx, alpha, color, fill, group, linetype,
  shape, size, weight)
```

```
f + geom_dotsplot(binaxis = "y", stackdir =
  "center", x, y, alpha, color, fill, group,
  linetype, size)
```

```
f + geom_violin(mapping = "area", x, y, alpha,
  color, fill, group, linetype, size, weight)
```

```
discrete x, discrete y
g <- ggplot(mpg, aes(cty, hwy))
```

```
g + geom_count(x, y, alpha, color, fill, shape,
  size, stroke)
```

### THREE VARIABLES

```
ss <- ggplot(mtcars, aes(cyl, long + 2, delta + lat + 2)) l <- ggplotify(ss, as_layer, id)
```

```
+> geom_contour(aes(z = z))
  x, y, z, alpha, color, group, linetype,
  size, weight
```

### continuous bivariate distribution

```
h <- ggplot(diamonds, aes(carat, price))
  h + geom_bivariate(binswidth = c(0.25, 500),
  x, y, alpha, color, fill, linetype, size, weight)
```

```
h + geom_hex(x, y, alpha, colour, group, linetype, size)
```

```
h + geom_hex2(x, y, alpha, color, fill, size)
```

continuous function

```
i + geom_function(mapping = aes(x, y, alpha, color))
  x, y, alpha, color, fill, linetype, size
```

```
i + geom_line(mapping = aes(x, y, alpha, color))
  x, y, alpha, color, group, linetype, size
```

```
i + geom_step(direction = "hv")
  x, y, alpha, color, group, linetype, size
```

### visualizing error

```
d <- data.frame(grp = c("A", "B"), fit = 4.5, se = 1.2)
j <- ggplot(d, aes(grp, fit, ymin = fit - se, ymax = fit + se))
  j + geom_crazier(fatten = 2)
  x, y, max, ymn, alpha, color, fill, group, linetype,
  size
```

```
j + geom_errorbar(x, y, ymin, ymn, alpha, color,
  fill, group, linetype, size)
```

```
j + geom_linerange(x, y, ymin, ymx, alpha, color, group,
  linetype, size)
```

```
j + geom_pointrange(x, y, ymin, ymx, alpha, color, fill,
  group, linetype, size)
```

### maps

```
data <- data.frame(murder = USA$restr5Murder,
  state = abbreviate(USStateNames)[USStateIDs])
  map <- leaflet(data, zoom = 2)
```

```
k <- ggplot(data, aes(lab = murder))
```

```
k + geom_map(as.map = statel,
  expand = TRUE, limit = map$long + map$lat,
  map_id, alpha, color, fill, linetype, size)
```





# Exercise 4 (ex24.R)

```
# 0. a) Check your working directory is in the course folder. Load the .Rproj file  
getwd()  
#     b) Load the tidyverse package (which loads the ggplot2 package amongst others)  
library(####)  
#     c) Run the code in ex24_prelim.R to import the UN data for this exercise  
source("./exercise/ex24_prelim.R")  
##  
##  
##  
# 1. Uncomment the code below (from ex23.R) and adapt to  
#     a. change the colour of the points to come from the "Set1" palette  
#     b. add more breaks in the size legend  
#         (Hint: use scale_size_continuous() with breaks argument set to a vector of  
#         typical population values such as c(50, 250, 500, 1000) )  
#     c. x-axis label: "Infant Mortality Rate"  
#     d. y-axis label: "Total Fertility Rate"  
#     e. colour label: "Continent"  
#     f. size label: "Population (m)"  
# ggplot(data = d, mapping = aes(x = imr, y = tfr, colour = area_name, size = pop/10  
#     geom_point() +  
#     coord_trans(x = "log10")
```