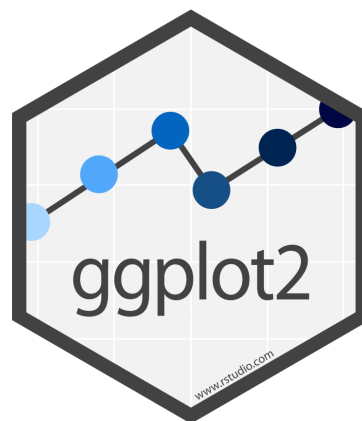


# Data Visualization

# esquisse and ggplot2



# Why learn ggplot2?

More customization:

- branding
- making plots interactive
- combining plots

Easier plot automation (creating plots in scripts)

Faster (eventually)

# ggplot2

- A package for producing graphics - gg = *Grammar of Graphics*
- Created by Hadley Wickham in 2005
- Belongs to “Tidyverse” family of packages
- “*Make a ggplot*” = Make a plot with the use of ggplot2 package

## Resources:

- <https://ggplot2-book.org/>
- <https://www.opencasestudies.org/>

# ggplot2

Based on the idea of:

layering

plot objects are placed on top of each other with +

□

□

□

# ggplot2

- Pros: extremely powerful/flexible – allows combining multiple plot elements together, allows high customization of a look, many resources online
- Cons: ggplot2-specific “grammar of graphic” of constructing a plot
- [ggplot2 gallery](#)

# Tidy data

To make graphics using `ggplot2`, our data needs to be in a **tidy** format

## Tidy data:

1. Each variable forms a column.
2. Each observation forms a row.

## Messy data:

- Column headers are values, not variable names.
- Multiple variables are stored in one column.
- Variables are stored in both rows and columns.

## Tidy data: example

Each variable forms a column. Each observation forms a row.

religion	income	freq
Agnostic	<\$10k	27
Agnostic	\$10-20k	34
Agnostic	\$20-30k	60
Agnostic	\$30-40k	81
Agnostic	\$40-50k	76
Agnostic	\$50-75k	137
Agnostic	\$75-100k	122
Agnostic	\$100-150k	109
Agnostic	>150k	84
Agnostic	Don't know/refused	96



# Messy data: example

Column headers are values, not variable names

religion	<\$10k	\$10-20k	\$20-30k	\$30-40k	\$40-50k	\$50-75k
Agnostic	27	34	60	81	76	137
Atheist	12	27	37	52	35	70
Buddhist	27	21	30	34	33	58
Catholic	418	617	732	670	638	1116
Don't know/refused	15	14	15	11	10	35
Evangelical Prot	575	869	1064	982	881	1486
Hindu	1	9	7	9	11	34
Historically Black Prot	228	244	236	238	197	223
Jehovah's Witness	20	27	24	24	21	30
Jewish	19	19	25	25	30	95

Table 4: The first ten rows of data on income and religion from the Pew Forum. Three columns, \$75-100k, \$100-150k and >150k, have been omitted

Read more about tidy data and see other examples: [Tidy Data](#) tutorial by Hadley Wickham

It's also helpful to have data in long format!!!

# Making data to plot

```
set.seed(3)
var_1 <- seq(from = 1, to = 30)
var_2 <- rnorm(30)
my_data = tibble(var_1, var_2)
my_data
```

```
# A tibble: 30 × 2
  var_1    var_2
  <int>    <dbl>
1     1 -0.962
2     2 -0.293
3     3  0.259
4     4 -1.15
5     5  0.196
6     6  0.0301
7     7  0.0854
8     8  1.12
9     9 -1.22
10    10  1.27
# ... with 20 more rows
```

First plot with **ggplot2** package

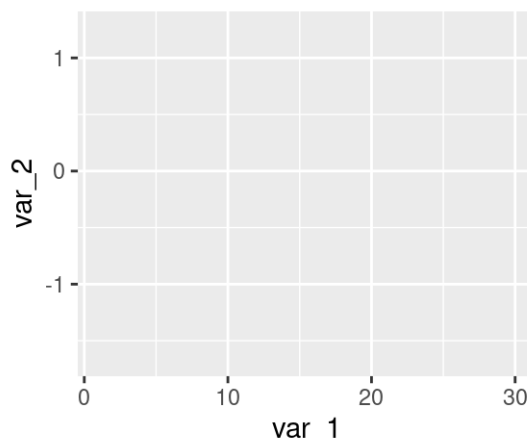
## First layer of code with **ggplot2** package

Will set up the plot - it will be empty!

- **Aesthetic mapping** (`mapping= aes(x= , y =)`) describes how variables in our data are mapped to elements of the plot

```
library(ggplot2) # don't forget to load ggplot2  
# This is not code but shows the general format  
ggplot({data_to_plot}, mapping = aes(x = {var in data to plot},  
                                     y = {var in data to plot}))
```

```
ggplot(my_data, mapping = aes(x = var_1, y = var_2))
```



## Next layer code with **ggplot2** package

There are many to choose from, to list just a few:

- `geom_point()` – points (we have seen)
- `geom_line()` – lines to connect observations
- `geom_boxplot()`
- `geom_histogram()`
- `geom_bar()`
- `geom_col()`
- `geom_errorbar()`
- `geom_density()`
- `geom_tile()` – blocks filled with color

## Next layer code with **ggplot2** package

Need the + sign to add the next layer to specify the type of plot

```
ggplot({data_to_plot}, mapping = aes(x = {var in data to plot},  
                                     y = {var in data to plot})) +  
  geom_{type of plot}</div>
```

```
ggplot(my_data, mapping = aes(x = var_1, y = var_2)) +  
  geom_point()
```



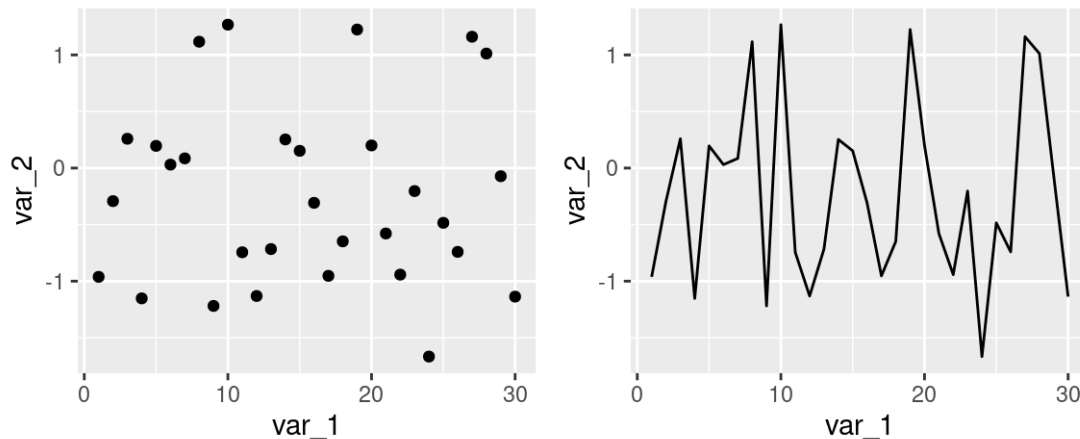
Read as: *add points to the plot (use data as provided by the aesthetic mapping)*

# Specifying plot layers: examples

```
plt1 <-  
  ggplot(my_data, aes(x = var_1, y = var_2)) +  
  geom_point()
```

```
plt2 <-  
  ggplot(my_data, aes(x = var_1, y = var_2)) +  
  geom_line()
```

`plt1; plt2` *# to have 2 plots printed next to each other on a slide*

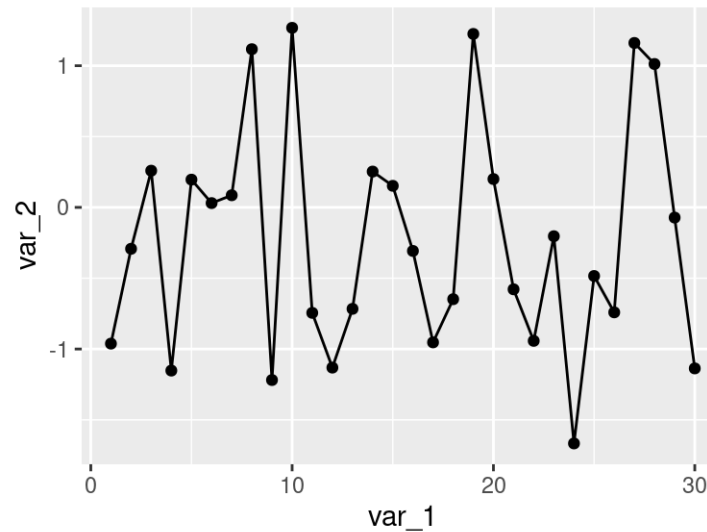


Also check out the [patchwork package](#)

# Specifying plot layers: combining multiple layers

Layer a plot on top of another plot with +

```
ggplot(my_data, aes(x = var_1, y = var_2)) +  
  geom_point() +  
  geom_line()
```

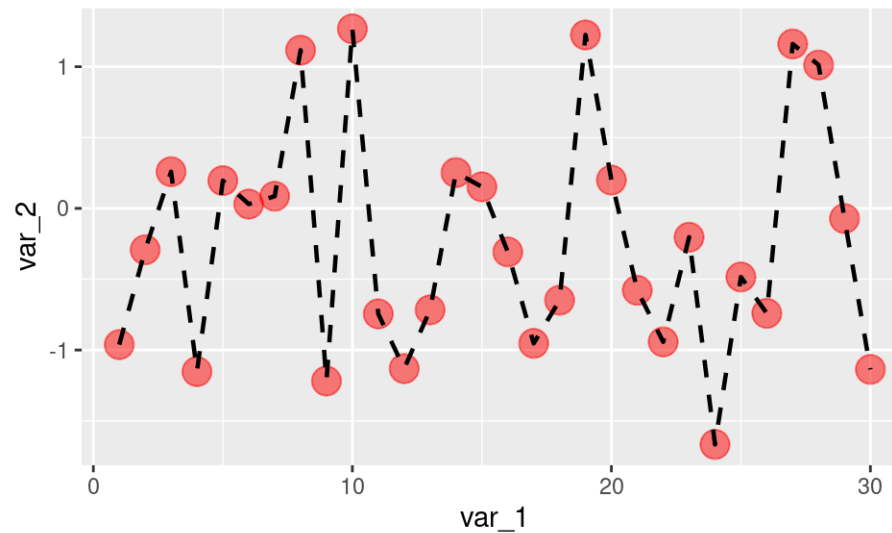




# Customize the look of the plot

You can change look of each layer separately.

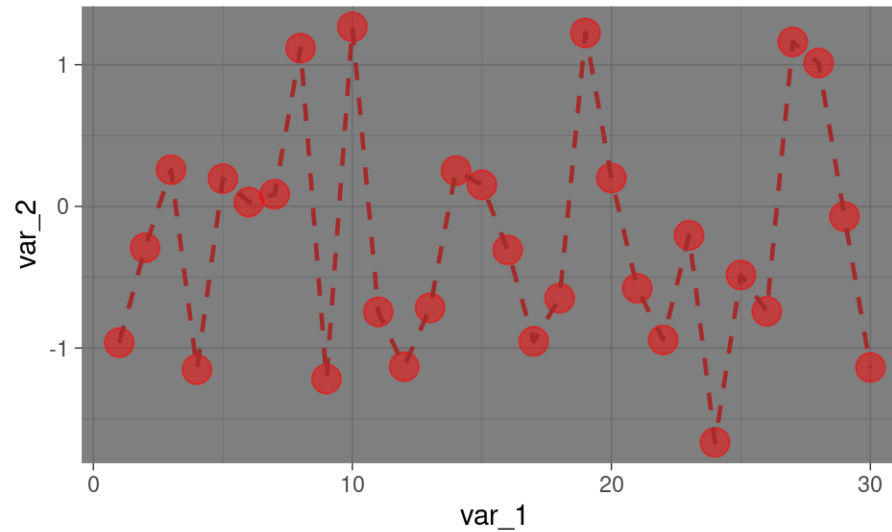
```
ggplot(my_data, aes(x = var_1, y = var_2)) +  
  geom_point(size = 5, color = "red", alpha = 0.5) +  
  geom_line(size = 0.8, color = "black", linetype = 2)
```



# Customize the look of the plot

You can change the look of whole plot using `theme_*()` functions.

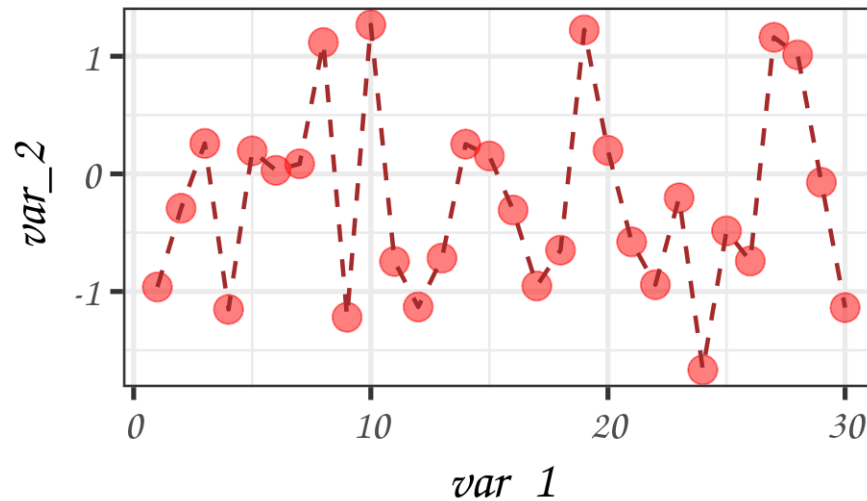
```
ggplot(my_data, aes(x = var_1, y = var_2)) +  
  geom_point(size = 5, color = "red", alpha = 0.5) +  
  geom_line(size = 0.8, color = "brown", linetype = 2) +  
  theme_dark()
```



# Customize the look of the plot

You can change the look of whole plot - **specific elements, too** - like changing [font](#) and font size - or even more [fonts](#)

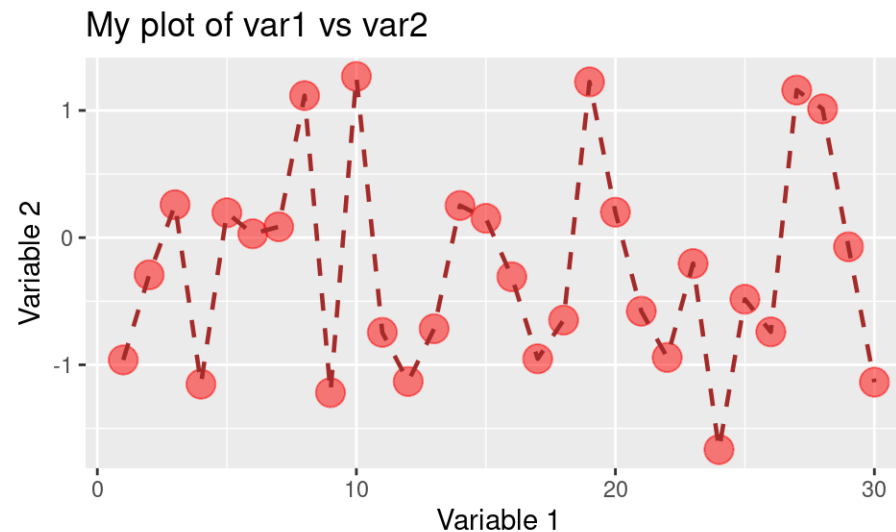
```
ggplot(my_data, aes(x = var_1, y = var_2)) +  
  geom_point(size = 5, color = "red", alpha = 0.5) +  
  geom_line(size = 0.8, color = "brown", linetype = 2) +  
  theme_bw(base_size = 20, base_family = "Comic Sans MS")
```



# Adding labels

The `labs()` function can help you add or modify titles on your plot.

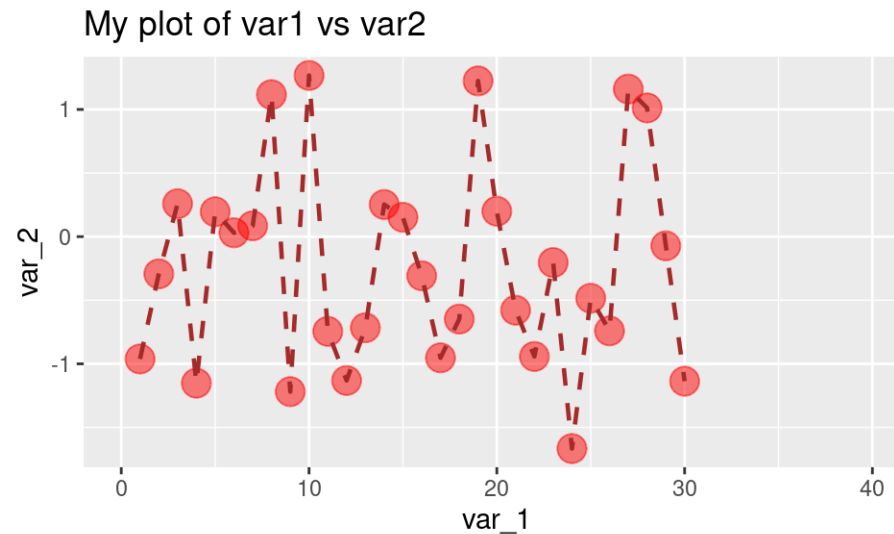
```
ggplot(my_data, aes(x = var_1, y = var_2)) +  
  geom_point(size = 5, color = "red", alpha = 0.5) +  
  geom_line(size = 0.8, color = "brown", linetype = 2) +  
  labs(title = "My plot of var1 vs var2",  
       x = "Variable 1",  
       y = "Variable 2")
```



# Changing axis

`xlim()` and `ylim()` can specify the limits for each axis

```
ggplot(my_data, aes(x = var_1, y = var_2)) +  
  geom_point(size = 5, color = "red", alpha = 0.5) +  
  geom_line(size = 0.8, color = "brown", linetype = 2) +  
  labs(title = "My plot of var1 vs var2") +  
  xlim(0, 40)
```



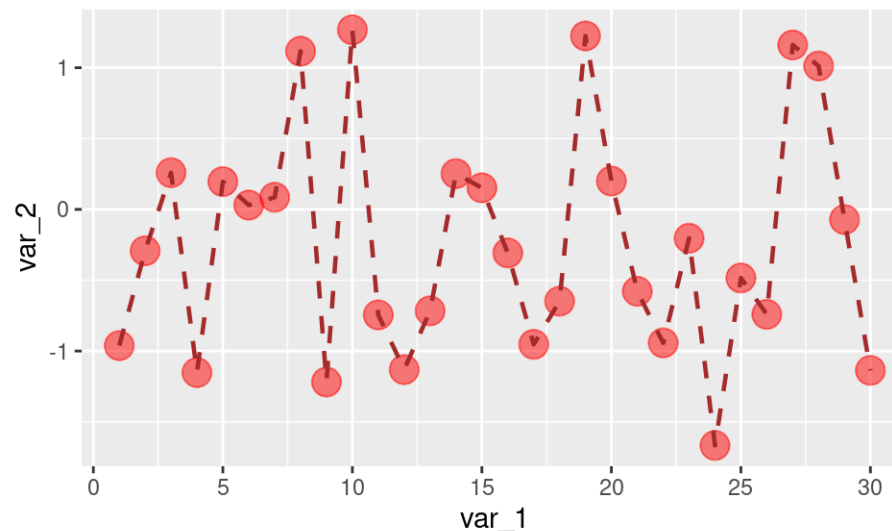
# Changing axis

`scale_x_continuous()` and `scale_y_continuous()` can change how the axis is plotted. Can use the `breaks` argument to specify how you want the axis ticks to be.

```
seq(from = 0, to = 30, by = 5)
```

```
[1] 0  5 10 15 20 25 30
```

```
ggplot(my_data, aes(x = var_1, y = var_2)) +  
  geom_point(size = 5, color = "red", alpha = 0.5) +  
  geom_line(size = 0.8, color = "brown", linetype = 2) +  
  scale_x_continuous(breaks = seq(from = 0, to = 30, by = 5))
```



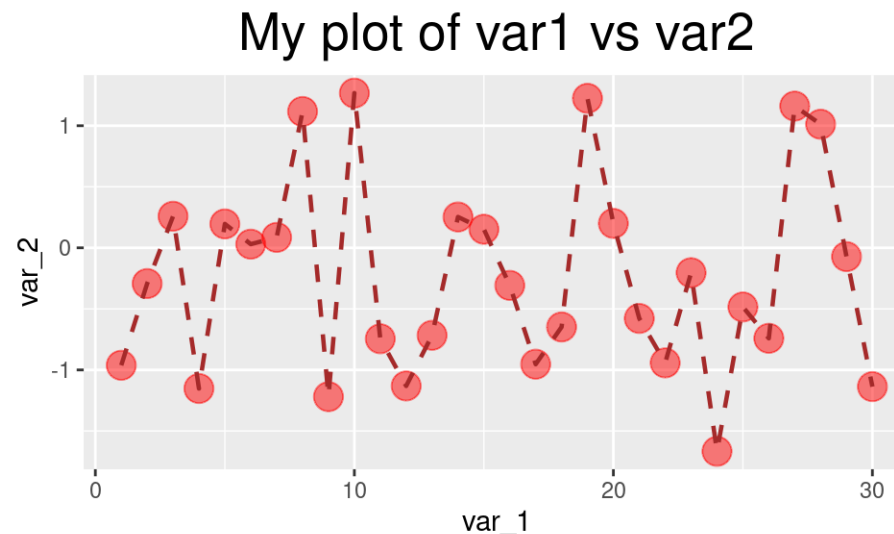
# Lab 1

[Lab document](#)

## theme() function

The `theme()` function can help you modify various elements of your plot. Here we will adjust the horizontal justification (`hjust`) of the plot title.

```
ggplot(my_data, aes(x = var_1, y = var_2)) +  
  geom_point(size = 5, color = "red", alpha = 0.5) +  
  geom_line(size = 0.8, color = "brown", linetype = 2) +  
  labs(title = "My plot of var1 vs var2") +  
  theme(plot.title = element_text(hjust = 0.5, size = 20))
```





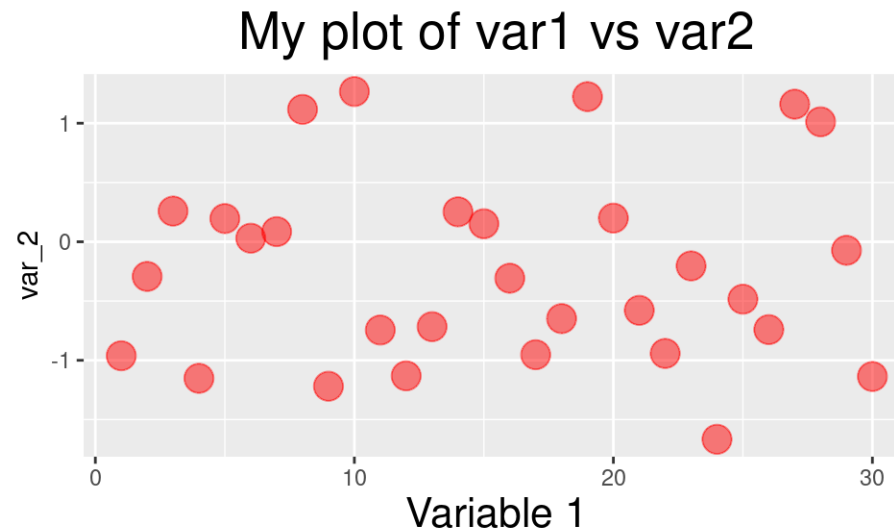
## theme() function

The `theme()` function always takes:

1. an object to change (use `?theme()` to see - `plot.title`, `axis.title`, `axis.ticks` etc.)
2. the aspect you are changing about this: `element_text()`, `element_line()`, `element_rect()`, `element_blank()`
3. what you are changing:
  - text: size, color, fill, face, alpha, angle
  - position: "top", "bottom", "right", "left", "none"
  - rectangle: size, color, fill, linetype
  - line: size, color, linetype

## theme() function

```
ggplot(my_data, aes(x = var_1, y = var_2)) +  
  geom_point(size = 5, color = "red", alpha = 0.5) +  
  labs(title = "My plot of var1 vs var2", x = "Variable 1") +  
  theme(plot.title = element_text(hjust = 0.5, size = 20),  
        axis.title.x = element_text(size = 16))
```



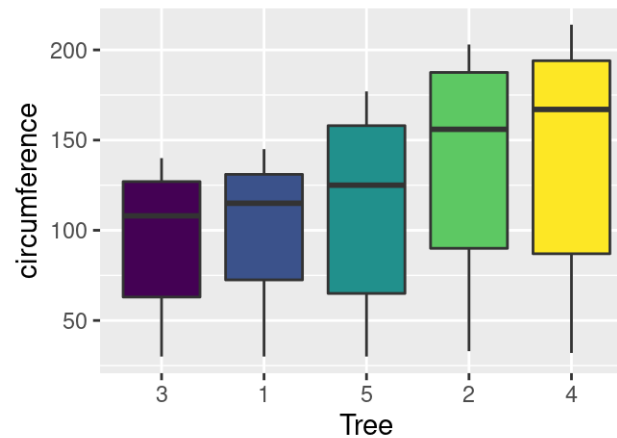
# theme() function

```
head(Orange, 3)
```

	Tree	age	circumference
1	1	118	30
2	1	484	58
3	1	664	87

If specifying position - use: "top", "bottom", "right", "left", "none"

```
ggplot(Orange, aes(x = Tree, y = circumference, fill = Tree)) +  
  geom_boxplot() +  
  theme(legend.position = "none")
```



## Can make your own theme to use on plots!

Guide on how to: <https://rpubs.com/mclaire19/ggplot2-custom-themes>

## Group and/or color by variable's values

First, we will generate some data frame for the purpose of demonstration.

- 2 different categories (e.g. pasta, rice)
- 4 different items (e.g. 2 of each category)
- 10 price values changes collected over time for each item

```
# create 4 vectors: 2x character class and 2x numeric class
item_categ <- rep(c("pasta", "rice"), each = 20)
item_ID <- rep(seq(from = 1, to = 4), each = 10)
item_ID <- paste0("ID_", item_ID)
observation_time <- rep(seq(from = 1, to = 10), times = 4)
item_price_change <- c(sample(0.5:2.5, size = 10, replace = TRUE),
                      sample(0:1, size = 10, replace = TRUE),
                      sample(2:5, size = 10, replace = TRUE),
                      sample(6:9, size = 10, replace = TRUE))
# use 4 vectors to create data frame with 4 columns
food <- tibble(item_ID, item_categ, observation_time, item_price_change)
```

## Group and/or color by variable's values

food

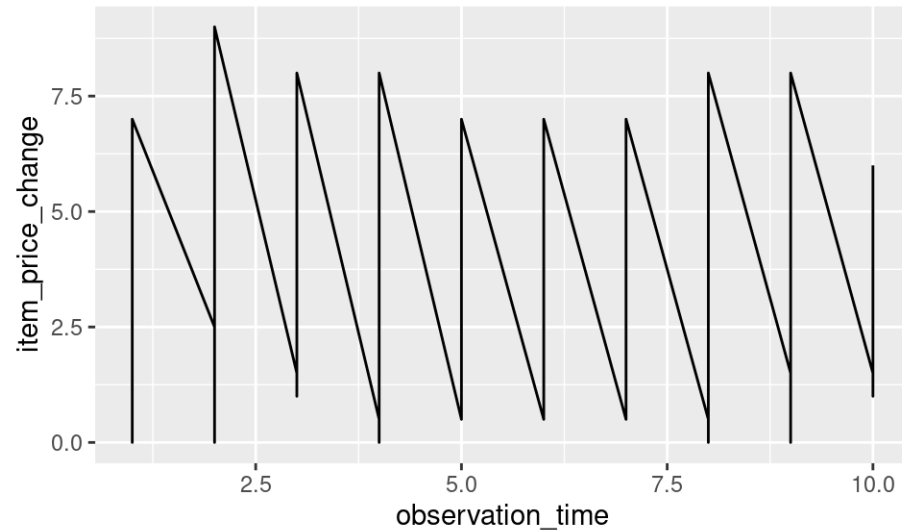
```
# A tibble: 40 × 4
```

	item_ID	item_categ	observation_time	item_price_change
	<chr>	<chr>	<int>	<dbl>
1	ID_1	pasta	1	0.5
2	ID_1	pasta	2	2.5
3	ID_1	pasta	3	1.5
4	ID_1	pasta	4	0.5
5	ID_1	pasta	5	0.5
6	ID_1	pasta	6	0.5
7	ID_1	pasta	7	0.5
8	ID_1	pasta	8	0.5
9	ID_1	pasta	9	1.5
10	ID_1	pasta	10	1.5

```
# ... with 30 more rows
```

## Starting a plot

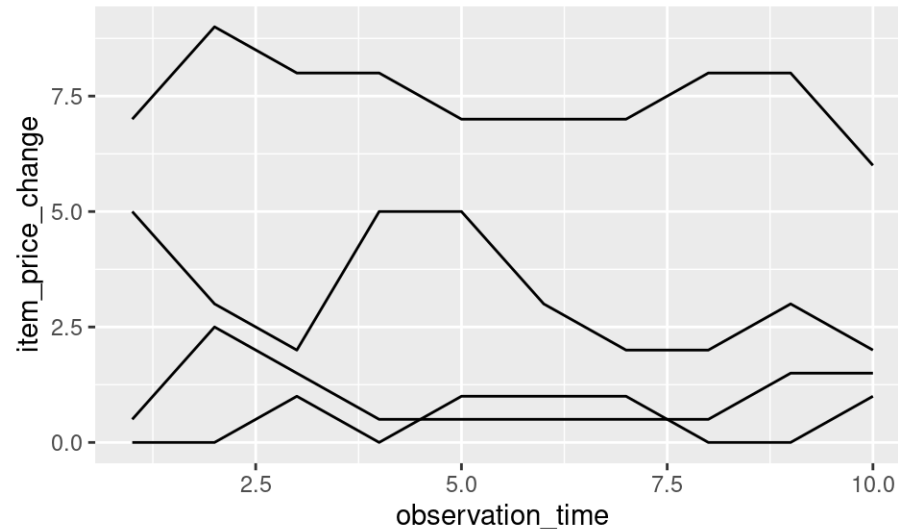
```
ggplot(food, aes(x = observation_time,  
                  y = item_price_change)) +  
  geom_line()
```



## Using **group** in plots

You can use **group** element in a mapping to indicate that each **item\_ID** will have a separate price line.

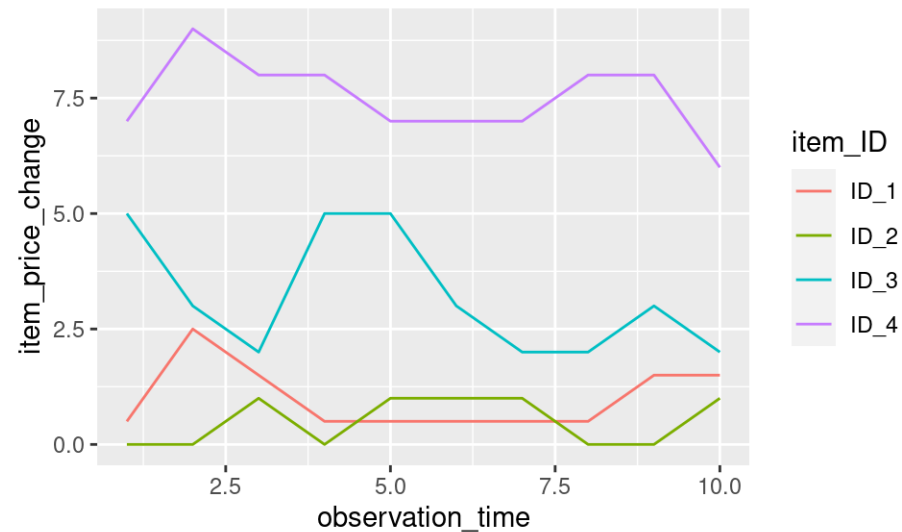
```
ggplot(food, aes(x = observation_time,  
                  y = item_price_change,  
                  group = item_ID)) +  
  geom_line()
```





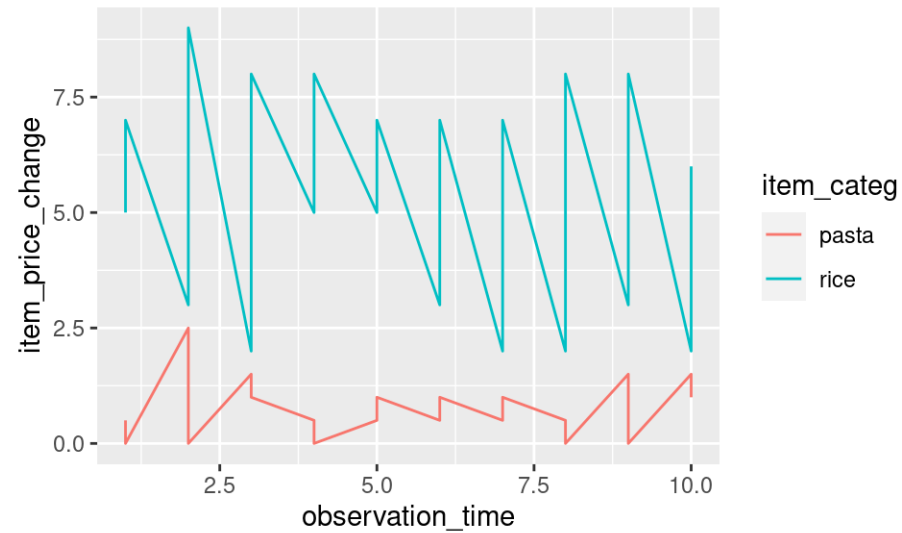
## Adding color will automatically group the data

```
ggplot(food, aes(x = observation_time,  
                 y = item_price_change,  
                 color = item_ID)) +  
  geom_line()
```



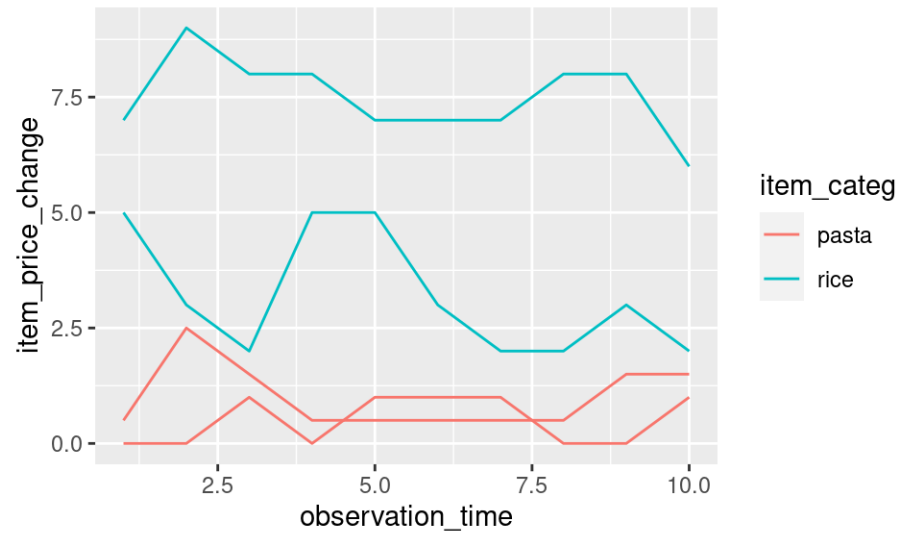
# Adding color will automatically group the data

```
ggplot(food, aes(x = observation_time,  
                 y = item_price_change,  
                 color = item_categ)) +  
  geom_line()
```



## Sometimes you need group and color

```
ggplot(food, aes(x = observation_time,  
                 y = item_price_change,  
                 group = item_ID,  
                 color = item_categ)) +  
  geom_line()
```

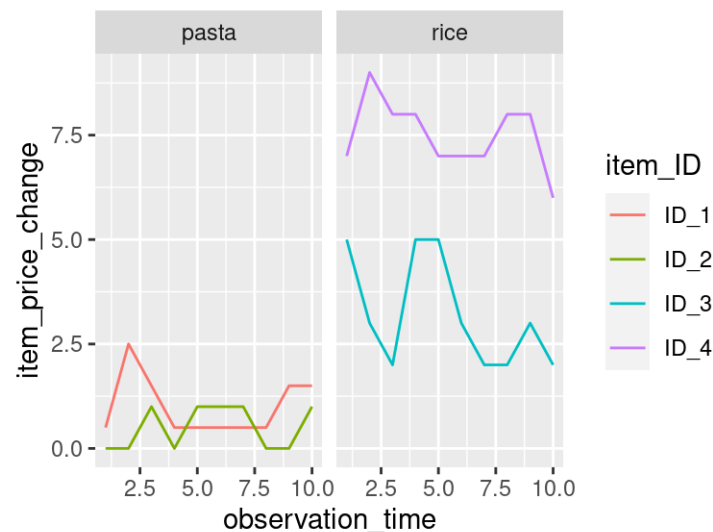


## Adding a facet can help make it easier to see what is happening

Two options: `facet_grid()` - creates a grid shape `facet_wrap()` - more flexible

Need to specify how you are faceting with the `~` sign.

```
ggplot(food, aes(x = observation_time,  
                 y = item_price_change,  
                 color = item_ID)) +  
  geom_line() +  
  facet_grid( ~ item_categ)
```

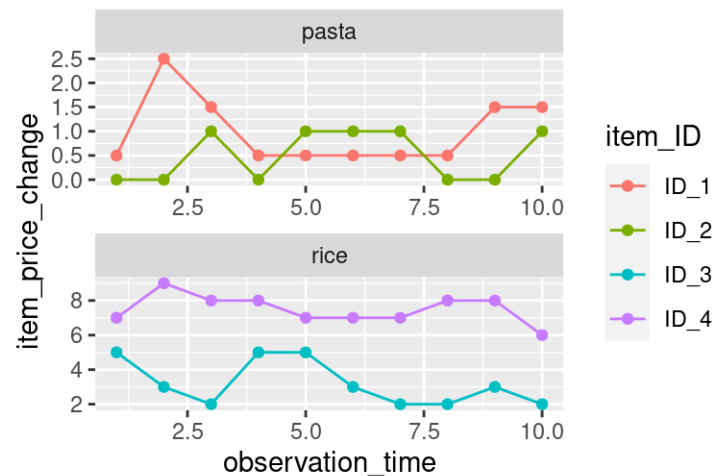


## facet\_wrap()

- more flexible - arguments `ncol` and `nrow` can specify layout
- can have different scales for axes using `scales = "free_x"`, `scales = "free_y"`, or `scales = "free"`

```
rp_fac_plot <- ggplot(food, aes(x = observation_time,  
                                y = item_price_change,  
                                color = item_ID)) +  
  geom_line() +  
  geom_point() +  
  facet_wrap(~ item_categ, ncol = 1, scales = "free")
```

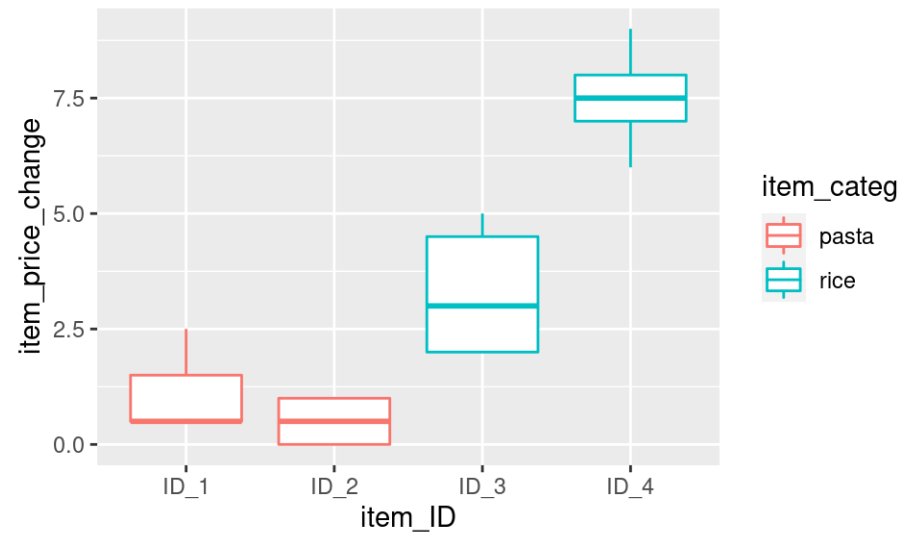
rp\_fac\_plot



## Tips - Color vs Fill

NOTE: color is needed for points and lines, fill generally needed for boxes and bars

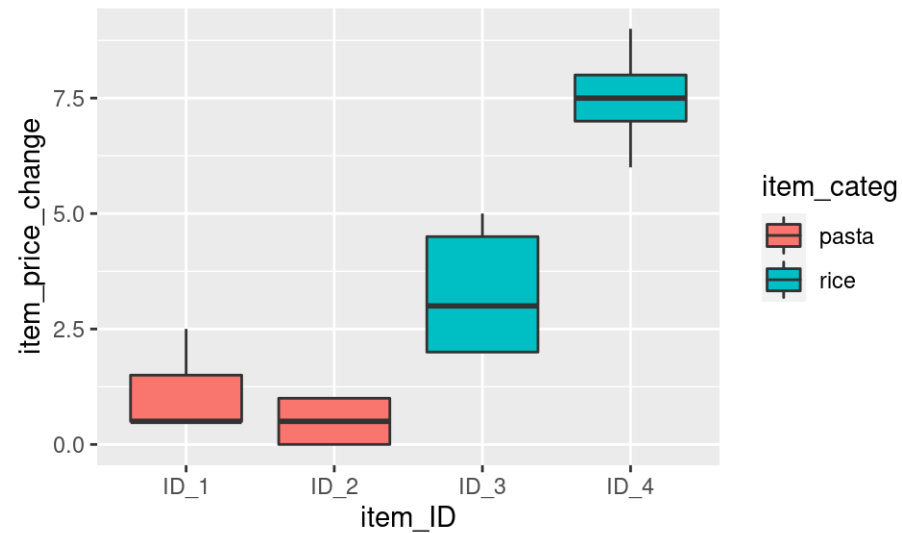
```
ggplot(food, aes(x = item_ID,  
                  y = item_price_change,  
                  color = item_categ)) +  
  geom_boxplot()
```



## Tips - Color vs Fill

NOTE: color is needed for points and lines, fill generally needed for boxes and bars

```
ggplot(food, aes(x = item_ID,  
                  y = item_price_change,  
                  fill = item_categ)) +  
  geom_boxplot()
```



## Tips - plus sign + can't come at start of a new line

This will not work! Also don't use pipes instead of +!

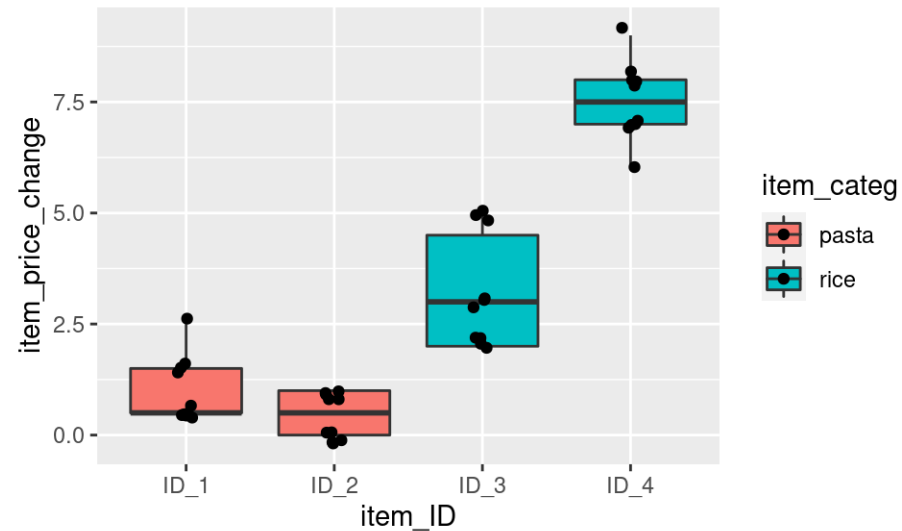
```
ggplot(food, aes(x = item_ID,  
                 y = item_price_change,  
                 fill = item_categ))  
+ geom_boxplot()
```



## Tip - Good idea to add jitter to top of box plots

Can add `width` argument to make the jitter more narrow.

```
ggplot(food, aes(x = item_ID,  
                 y = item_price_change,  
                 fill = item_categ)) +  
  geom_boxplot() +  
  geom_jitter(width = .06)
```

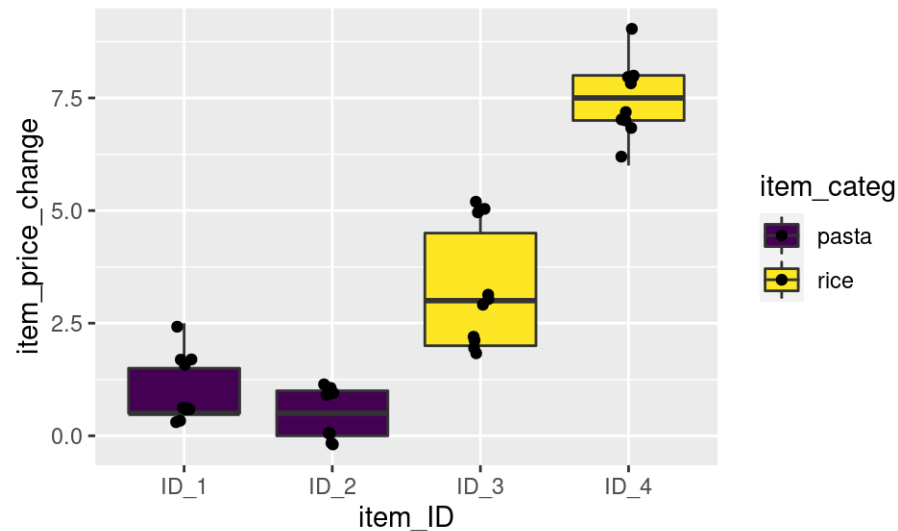


## Tip - be careful about colors for color blindness

`scale_fill_viridis_d()` for discrete /categorical data

`scale_fill_viridis_c()` for continuous data

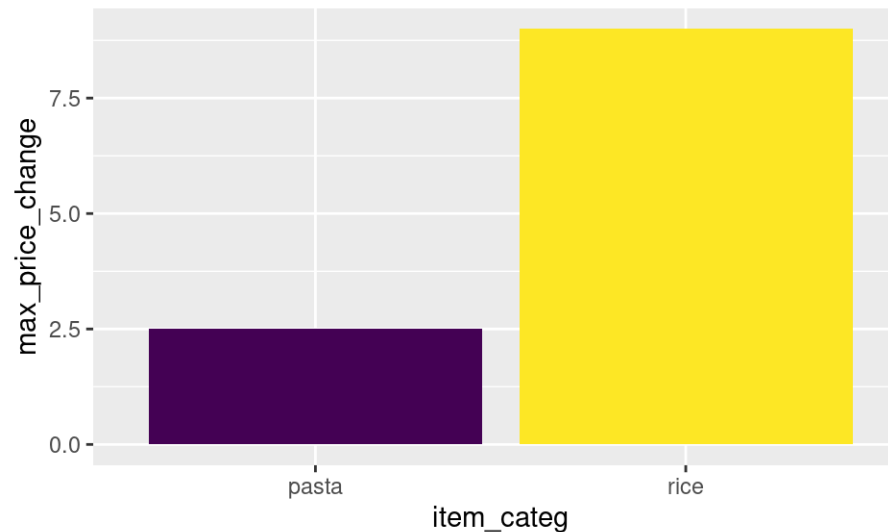
```
ggplot(food, aes(x = item_ID,  
                 y = item_price_change,  
                 fill = item_categ)) +  
  geom_boxplot() +  
  geom_jitter(width = .06) +  
  scale_fill_viridis_d()
```



## Tip - can pipe data after wrangling into ggplot()

```
food_bar <- food %>%  
  group_by(item_categ) %>%  
  summarize("max_price_change" = max(item_price_change)) %>%  
  ggplot(aes(x = item_categ,  
            y = max_price_change,  
            fill = item_categ)) +  
  scale_fill_viridis_d()+  
  geom_col() +  
  theme(legend.position = "none")
```

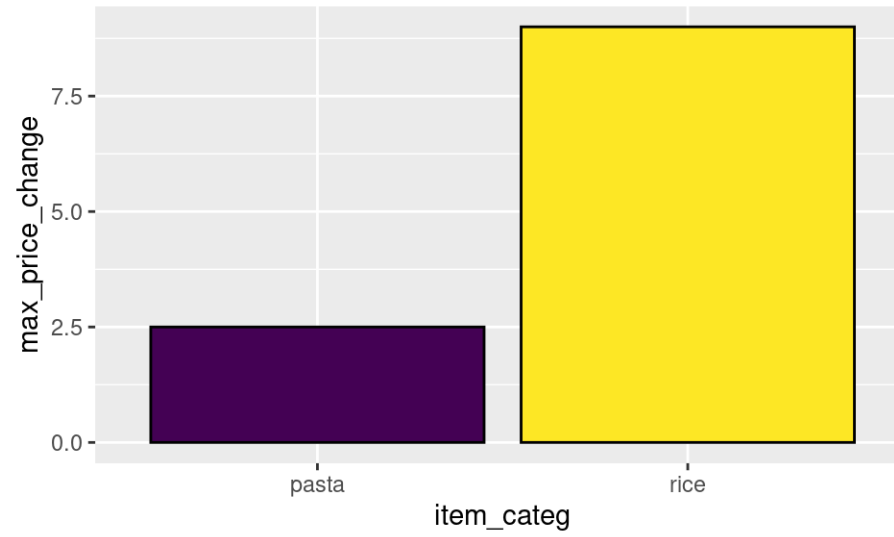
food\_bar



## Tip - color outside of `aes()`

Can be used to add an outline around column/bar plots.

```
food_bar +  
  geom_col(color = "black")
```

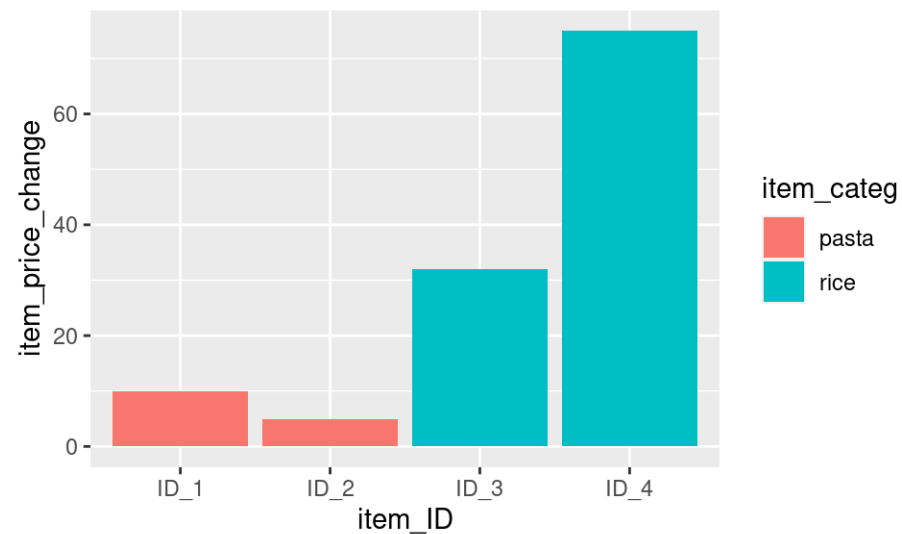


## Tip - col vs bar

`geom_bar()` can only one aes mapping & `geom_col()` can have two

□ May not be plotting what you think you are...

```
ggplot(food, aes(x = item_ID,  
                  y = item_price_change,  
                  fill = item_categ)) +  
  geom_col()
```



## What did we plot?

```
head(food)
```

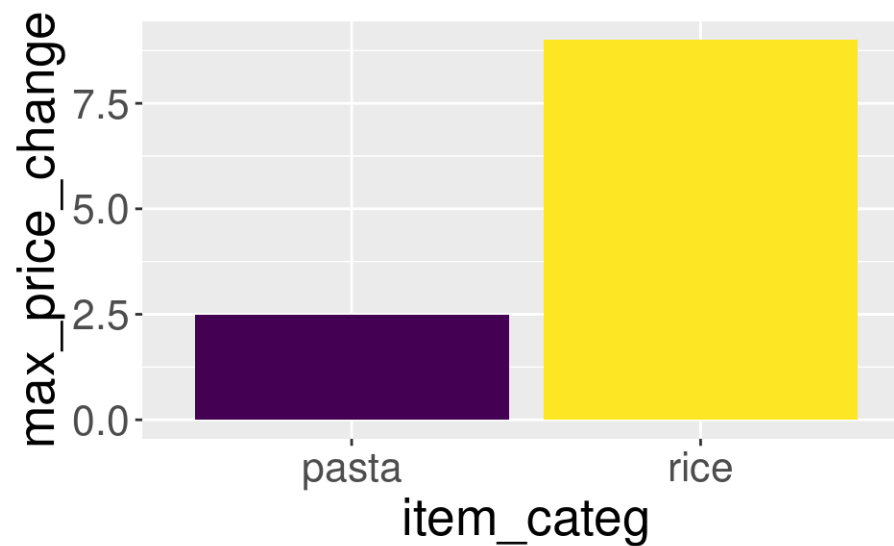
```
# A tibble: 6 × 4
  item_ID item_categ observation_time item_price_change
  <chr>    <chr>             <int>             <dbl>
1 ID_1    pasta                1                0.5
2 ID_1    pasta                2                2.5
3 ID_1    pasta                3                1.5
4 ID_1    pasta                4                0.5
5 ID_1    pasta                5                0.5
6 ID_1    pasta                6                0.5
```

```
food %>% group_by(item_ID) %>%
  summarize(sum = sum(item_price_change))
```

```
# A tibble: 4 × 2
  item_ID sum
  <chr>    <dbl>
1 ID_1    10
2 ID_2     5
3 ID_3    32
4 ID_4    75
```

## Tip - make sure labels aren't too small

```
food_bar +  
  theme(text = element_text(size = 20))
```



**Extensions**

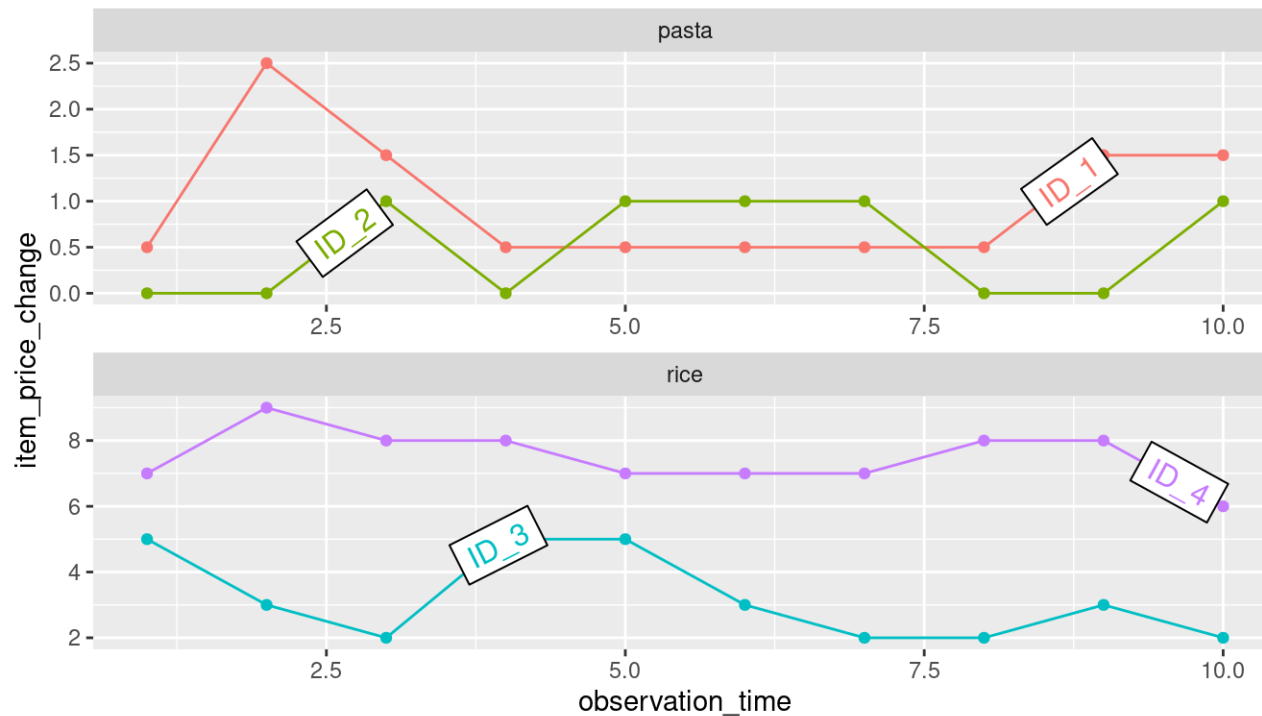


# directlabels package

Great for adding labels directly onto plots

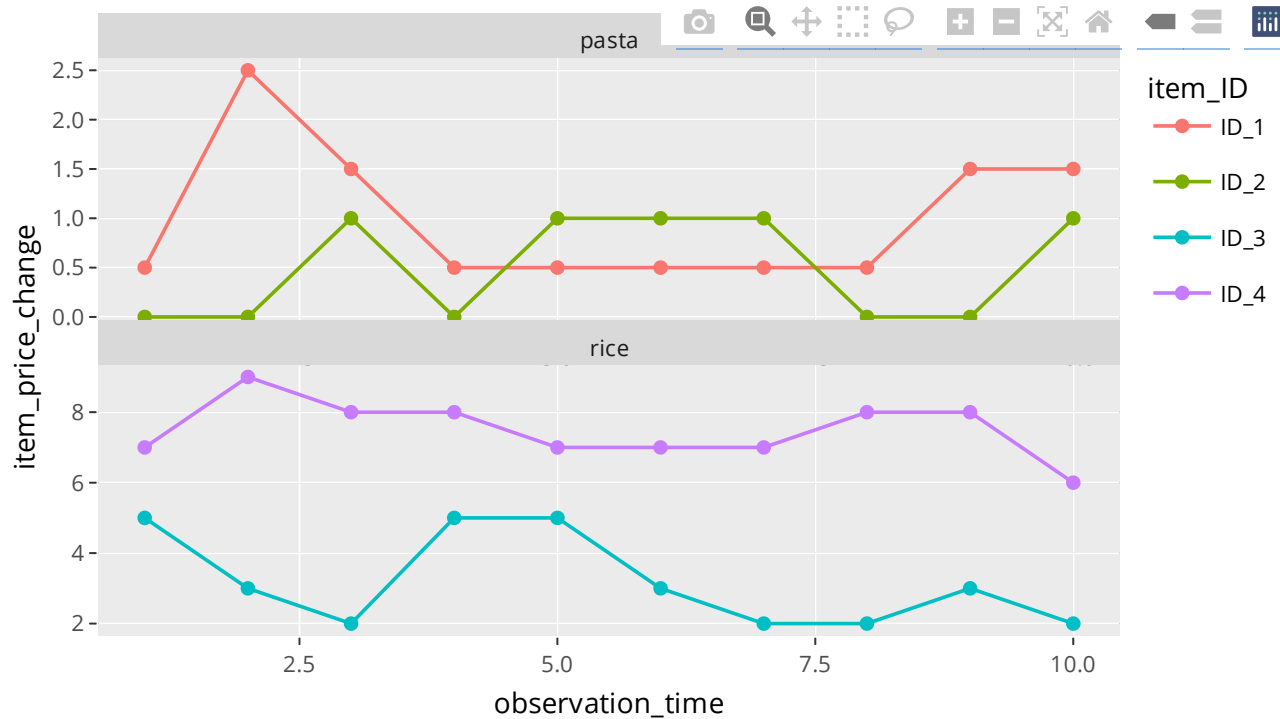
<https://www.opencasestudies.org/ocs-bp-co2-emissions/>

```
#install.packages("directlabels")  
library(directlabels)  
direct.label(rp_fac_plot, method = list("angled_boxes"))
```



# plotly

```
#install.packages("plotly")  
library("plotly")  
ggplotly(rp_fac_plot)
```



Also check out the [ggiraph package](#)

# Saving a ggplot to file

A few options:

- RStudio > Plots > Export > Save as image / Save as PDF
- RStudio > Plots > Zoom > [right mouse click on the plot] > Save image as
- In the code

```
ggsave(filename = "saved_plot.png", # will save in working directory
        plot = rp_fac_plot,
        width = 6, height = 3.5)      # by default in inch
```

# Lab 2

[Lab document](#)