

• code for plots on last memo.

STATISTICAL ANALYSIS.

R - notes.

Statistical Analysis (TUT 3) (TUT 4).

- transform numerical to words

%>% mutate(.) function

1. basic stats analysis.

① Basic summary

② Find missing values

③ Remove rows, na.omit(.)

④ Summary after clean

⑤ Ensure no na values

+ / " " = p , " " = x , " " = g(fit) edsl

⑥ Impute missing values with column means

airquality - impute = airquality %>% mutate (across (everything), ~ ifelse (is.na(.), mean (., na.rm = TRUE), .))

+ ((g(x)=p, q(x)=x) edsl)

⑦ Replace actual data with imputed data

+ (g(x)=x, "na"=b(na) d(na)-na)

⑧ Create histogram. (each variable except month and day)

+ (l(b(na) g(x)) edsl)

⑨ Boxplot by month => Ozone

Boxplot by month => Temp

Boxplot by month => Solar.R

Boxplot by month each variable.

⑩ Scatterplot with smoothing

① Pairwise Scatterplot matrix

Linear regression model:

model = lm(Ozone ~ temp + wind + Solar.R, data = airquality)

summary(model)

② TimeSeries plot (Ozone over time (index))

airquality \$ DayIndex = 1:nrow(airquality)

ggplot(airquality, aes(x = DayIndex, y = Ozone)) +

geom_line(color = "darkgreen") +

geom_smooth(method = "loess") +

labs(title = "", x = "", y = "") +

theme_minimal()

3. Facet Scatter plot.

ggplot(airquality, aes(x = Temp, y = Ozone)) +

geom_point(aes(color = month)) +

geom_smooth(method = "lm", se = FALSE) +

facet_wrap(~month) +

theme_minimal() +

labs(title = "")

Statistical analysis:

- Load data
- Scatterplot of GDP vs unemployed
- ggplot
- Fit a linear regression model.

model = lm(unemployment_rate ~ GDP_per_capita, data = dc)
summary(model)

Correlation Analysis.

x <- ggplot(mtcars[, c(" ", " ", " ", " ", " ")],
title = " ")

Print(p5).

Linear regression:

- model <- lm(mpg ~ hp, data = mtcars)
- summary(model)
- Extract residuals and fitted values

residuals <- resid(model)

fitted <- fitted(model)

(TUT4)

- residuals vs fitted values plot
- normal Q-Q plot
- Scale - Location Plot (spread vs fitted).
- Histogram and density of residuals

Shapiro-Wilk Normality test (includes log transform)
code (TUT4) ~~memo~~ or ~~gaw~~ to ~~July~~ ~~2023~~ ~~2023~~

TUT3 (code for following plots):

- * Histogram
 - * Boxplots
 - * Scatterplot with smoothing
 - * Pairwise scatterplot matrix
 - * Time Series plot
 - * Faceted scatterplot
 - * Scatterplot of GDP vs Unemp
 - * ggpairs.
- } "airquality" dataset used in TUT3.
- } economics dataset from BB.

TUT4 (code for following plots):

- * Histogram (horsepower)
 - * Boxplot of MPG grouped by cylinder
 - * Density plot of MPG
 - * Pairwise plot.
 - * Residuals vs fitted model plot
 - * Scale location plot.
 - * Histogram, Density of residuals
 - * normal Q-Q Plot.
- } mtcars dataset from

SVM - models (TUT 5)

```
rm (list = ls())  
library (MASS)  
library (e1071)
```

- ① Load and preview data set:

```
data ("iris")  
head ( )
```

- ② Create BINARY RESPONSE variable:

Affairs\$had-Affairs <- as.factor (ifelse (Affairs\$Affairs > 0, 1, 0))

data = Affairs[, -1] # remove ID
str (data)

- ③ Manually train-split (70/30):
Set seed 100

```
n <- nrow (data)
```

train_index <- sample (1:n, size = 0.7 * n)

train_data <- data[train_index,]

test_data <- data[-train_index,]

- ④ Train SVM model with linear kernel.

```
svm_model = svm (had-affair ~ ., data = train_data,  
kernel = "linear", scale = TRUE,  
probability = TRUE)
```

⑤ Predict on test data

```
pred-test <- predict(svm-model, newdata = test-data)
```

⑥ Calculate training error

```
training-error <- mean(pred-test != test-data$hp) * 100  
print(training-error).
```

⑦ Evaluate and comment on the prediction accuracy of both models.

```
prediction <- predict(mtcars-test, mtcars-train)
```

```
svm-accuracy <- postResample(prediction, mtcars-test, mtcars-train)
```

Scaling the data:

```
mtcars-scaled = as.data.frame(scale(mtcars))
```

To fit one, all algos → kfold-algo

Crossval and kfold → cv-algo

Crossval-algo -> kfold -> kfold-test

With-algo - algo, ~ initial-best algo = kfold-algo

algo - algo, "algo" = kfold

"algo" = gridsearch

Plot filtering multiple line charts. (Assignment 4)

```
rm(list = TRUE)
```

```
library(tidyverse)
```

```
library(ggplot2)
```

```
data = read.csv(" ", header = TRUE)
```

- ① Melr the dataset to a long format for plotting

```
data_long = gather(data, key = "time", value = "value",  
-dummy-variable, -id)
```

```
data_long$time <- as.numeric(as.charactergsub("v",  
data_long$time)))
```

- ② Plot using ggplot with facet-wrap.

```
P = ggplot(data_long, aes(x = __, y = __, group = id, color  
= dummy-variable)) +  
geom_line(alpha = 0.7) +  
labs(title = " ", x = " ", y = " ")
```

```
scale_color_manual(values = c("T1" = "darkblue", "T2"  
, "T3" = "red")) +
```

```
theme_classic() +  
theme(legend.title = element_blank()) modify length
```

```
print(p)
```

ANOVA MODEL:

```
model <- aov(weight ~ group, data = Plant.growth)  
summary(model)
```

NORMALISING THE DATA.

```
maxs <- apply(rock, 2, max)
```

```
mins <- apply(lmin, 2, min)
```

```
Scaled.rock <- as.data.frame(scale(rock,  
centre = mins, scale = maxs - mins))
```

2 amount of variables.