

PassengerId คือ เลขประจำตัวผู้โดยสาร

Pclass คือ ระดับของผู้โดยสาร

1 = ระดับ 1 ระดับสูง

2 = ระดับ 2 ระดับปานกลาง

3 = ระดับ 3 ระดับล่าง

Name คือ ชื่อของผู้โดยสาร

Sex คือ เพศของผู้โดยสาร

Age คือ อายุของผู้โดยสาร

SibSp คือ จำนวนบุคคลที่เดินทางมาด้วย (ญาติพี่น้อง สามี ภรรยา)

Parch คือ จำนวนบุคคลที่เดินทางมาด้วย (พ่อแม่ ลูก)

Ticket คือ หมายเลขตั๋ว

Fare คือ ค่าโดยสาร

Cabin คือ หมายเลขห้องพัก

Embarked คือ สถานที่ขึ้นเรือ

C = Cherbourg

Q = Queenstown

S = Southampton

Survived คือ ผู้ที่รอดชีวิต

0=ไม่รอด

1=รอด

#write by 5605104046 นาย วิศิษฐ์ เลิศศักดิ์วิมาน

#write by 5605104043 นาย อานนท์ กันทา

Load the data sets

```
setwd("D:/Com/ICT_UTCC/DATA MINING/titanic-master/")
```

```
train <- read.csv("Data/CSV/train.csv", stringsAsFactors = FALSE) # 891 obj ชุด
```

```
test <- read.csv("Data/CSV/test.csv", stringsAsFactors = FALSE) # 418 obj ชุด
```

#Import library

```
library(rpart)
```

```
library(randomForest)
```

```
library(party)
```

```
library(rattle)
```

```
library(FileWriter)
```

Clean data sets (1)- (9)

#(1) ทำการใส่ค่า NA ลงไปในช่องว่างในชุดข้อมูลฝึกสอน

```
train[train == ""] <- NA
```

#(2) ทำการใส่ค่า NA ลงไปในช่องว่างของตัวแปรSurvivedของข้อมูล

```
test$Survived <- NA
```

#(3) ทำการรวมชุดข้อมูลฝึกสอนและข้อมูลทดสอบไว้ด้วยกันในตัวแปรอบเจ็ค combi

```
combi <- rbind(train, test)
```

#(4) ทำสร้างตัวแปร Title และตัดเอาหน้าหน้าชื่อจากตัวแปร Name ไปสร้างตัวแปร Title

```
#   และทำการใส่ค่า Mlle แทนลงไปในช่วงที่มีคำว่า Mme หรือ Mlle
```

```
#   และทำการใส่ค่า Lady แทนลงไปในช่วงที่มีคำว่า Dona หรือ Lady หรือ the Countess
```

```
#   และทำการใส่ค่า Sir แทนลงไปในช่วงที่มีคำว่า Capt หรือ Don หรือ Major หรือ Sir หรือ Jonkheer หรือ Dr
```

```
#   และทำการแปลง Title จากข้อมูลแบบvectorให้เป็นข้อมูลแบบfactor โดยฟังก์ชัน factor
```

```
combi$Title <- sapply(combi$Name, FUN=function(x) {strsplit(x, split='[,.')] [[1]] [2]})
```

```
combi$Title <- sub(' ', '', combi$Title)
```

```
combi$Title[combi$Title %in% c('Mme', 'Mlle')] <- 'Mlle'
```

```
combi$Title[combi$Title %in% c('Dona', 'Lady', 'the Countess')] <- 'Lady'
```

```
combi$Title[combi$Title %in% c('Capt', 'Don', 'Major', 'Sir', 'Jonkheer', 'Dr')] <- 'Sir'
```

```
combi$Title <- factor(combi$Title)
```

#(5) ทำสร้างตัวแปร FamilySize และทำการใส่ขนาดของครอบครัวลงไปโดย นำ

```
# จำนวนบุคคลที่เดินทางมาด้วย (ญาติพี่น้อง สามี ภรรยา) + จำนวนบุคคลที่เดินทางมาด้วย (พ่อแม่ ลูก) + 1(จำนวนของตัวเอง)
```

```
combi$FamilySize <- combi$SibSp + combi$Parch + 1
```

#(6) ทำสร้างตัวแปร Surname และทำการใส่สกุลที่ตัดได้จาก ตัวแปรName

```
combi$Surname <- sapply(combi$Name, FUN=function(x) {strsplit(x, split='[,.]')[[1]][1]})
```

#(7) ทำสร้างตัวแปร FamilyID เพื่อแยกประเภทขนาดความใหญ่ของครอบครัว

โดยทำการใส่ Small ลงไปใน FamilyID ที่ FamilySize <= 2

```
combi$FamilyID <- paste(as.character(combi$FamilySize), combi$Surname, sep="")
```

```
combi$FamilyID[combi$FamilySize <= 2] <- 'Small'
```

```
famIDs <- data.frame(table(combi$FamilyID))
```

```
famIDs <- famIDs[famIDs$Freq <= 2,]
```

```
combi$FamilyID[combi$FamilyID %in% famIDs$Var1] <- 'Small'
```

```
combi$FamilyID <- factor(combi$FamilyID)
```

#(8) ทำการใส่ค่าลงไปในตัวFare ที่มีIndex = 1044 โดยค่าที่ใส่เป็นค่าเฉลี่ยจาก Fareทั้งหมด

```
combi$Fare[1044] <- median(combi$Fare, na.rm=TRUE)
```

#(9) ทำการใส่ค่าลงไปในตัว Embarked colum ลำกับที่ 62และ 830 ด้วยค่า S เนื่องจากค่าSถูกแทนมากที่สุด

```
combi$Embarked[c(62,830)] = "S"
```

#ทำการสร้างตัวแปรอบเจ็ค Agefit มาเพื่อพยากรณ์หาค่า Ageที่ว่างอยู่โดยมีตัวแปรเป้าหมายคือ Age และมีตัวที่ใช้ในการพยากรณ์ คือ

Pclass , Sex , SibSp , Parch , Fare , Embarked , Title , FamilySize

#โดยใช้ฟังก์ชัน rpart และใช้ method="anova"

```
Agefit <- rpart(Age ~ Pclass + Sex + SibSp + Parch + Fare + Embarked + Title + FamilySize, data=combi[!is.na(combi$Age),], method="anova")
```

```
combi$Age[is.na(combi$Age)] <- predict(Agefit, combi[is.na(combi$Age),])
```

#ทำการแปลงหรือencode (Sex, Embarked ,Survived) จากข้อมูลแบบvectorให้เป็นข้อมูลแบบfactor โดยฟังก์ชัน factor

```
combi$Sex <- factor(combi$Sex)
```

```
combi$Embarked <- factor(combi$Embarked)
```

```
combi$Survived <- factor(combi$Survived)
```

ทำการแยกชุดข้อมูลฝึกสอนและข้อมูลทดสอบ (combi) ออกจากกัน ไว้ในตัวแปรอบเจ็ค train กับ test

```
train <- combi[1:891,]
```

```
test <- combi[892:1309,]
```

#ทำการนำข้อมูลจากตัวแปรอบเจ็ค train เขียนข้อมูลลงไปใน Path “ Data/CSV/ “ ไฟล์ชื่อ Clean_train .csv

#ทำการนำข้อมูลจากตัวแปรอบเจ็ค test เขียนข้อมูลลงไปใน Path “ Data/CSV/ “ ไฟล์ชื่อ Clean_test.csv

```
write.csv(train, file = "Data/CSV/Clean_train.csv", row.names = FALSE)
```

```
write.csv(test, file = "Data/CSV/Clean_test.csv", row.names = FALSE)
```

ทำการสร้างตัวแปรอบเจ็ค fit มาเพื่อพยากรณ์หาผู้รอดชีวิตโดยมีตัวแปรเป้าหมายคือ Survived และมีตัวที่ใช้ในการพยากรณ์ คือ

Pclass , Sex , Age , SibSp , Parch, Fare, Embarked , Title , FamilySize , FamilyID

โดยใช้ฟังก์ชัน cforest กำหนดให้ ntree=5000, mtry=2

set.seed(415)

fit <- cforest(Survived ~ Pclass + Sex + Age + SibSp + Parch + Fare + Embarked + Title + FamilySize + FamilyID, data = train,

controls=cforest_unbiased(ntree=5000, mtry=2))

สร้างตัวแปรอบเจ็ค Prediction ขึ้นมา และนำค่า fit ที่ได้จากการพยากรณ์ใส่ลงไปในข้อมูลทดสอบ test โดยใช้ฟังก์ชัน predict

Prediction <- predict(fit, test, OOB=TRUE, type = "response")

สร้างตัวแปรอบเจ็ค submit ขึ้นมา ทำการสร้าง frame ข้อมูลโดยใส่

column แรกชื่อ PassengerId และใส่ค่า PassengerId ของข้อมูลทดสอบลงไป

column สองชื่อ Survived และใส่ค่า Prediction ที่ได้จากการพยากรณ์ลงไป

submit <- data.frame(PassengerId = test\$PassengerId, Survived = Prediction)

#ทำการนำข้อมูลจากตัวแปรอบเจ็ค submit เขียนข้อมูลลงไปใน Path “ Data/CSV/ “ ไฟล์ชื่อ CF.csv

FileWriter()

write.csv(submit, file = "Data/CSV/CF.csv", row.names = FALSE)

Public Leaderboard

https://www.kaggle.com/c/titanic/leaderboard?submissionId=2762439

Public Leaderboard

48	13	JonathanSlomka	0.82775	1	Fri, 19 Feb 2016 15:07:04
49	13	Sonia B	0.82775	13	Thu, 25 Feb 2016 07:38:57 (-3.5d)
50	13	VeronicaD	0.82775	22	Fri, 26 Feb 2016 22:50:57
51	13	Bryan Balajadia	0.82775	6	Wed, 09 Mar 2016 03:54:18 (-8.9d)
52	13	Samuel 4	0.82775	20	Fri, 18 Mar 2016 15:19:52 (-13.6h)
53	new	Wisit Lertsakwimarn	0.82775	16	Sun, 27 Mar 2016 11:15:40

Your Best Entry ↑

You improved on your best score by 0.00957.

You just moved up 93 positions on the leaderboard.

Tweet this!

54	14	MVP24	0.82297	2	Thu, 28 Jan 2016 01:44:24
55	14	ferranarroyo	0.82297	1	Sat, 30 Jan 2016 17:51:17
56	14	akmnko	0.82297	1	Tue, 02 Feb 2016 03:37:51
57	14	zamp	0.82297	1	Tue, 02 Feb 2016 13:58:33
58	14	sravanthi P	0.82297	15	Sun, 21 Feb 2016 23:31:20 (-15.9d)
59	14	nkelly13	0.82297	9	Tue, 09 Feb 2016 19:51:27 (-20.6h)
60	14	mayhem!	0.82297	1	Fri, 12 Feb 2016 19:54:35

Activate Windows
Go to Settings to activate Windows.

18:16
27/3/2559