

Generazione di modelli per predire l'interesse dei clienti per l'acquisto di una polizza

Francesco Gozzoli

29/10/2021

INTRODUZIONE

Il dataset contiene informazioni sull'interesse di clienti di una compagnia di assicurazioni per l'acquisto di polizze auto. Inizialmente le colonne che compongono il dataset si trovano nella situazione mostrata dalla funzione `summary`.

```
insurance <- read.csv("../insurance.csv")
summary(insurance)
```

```
##      id      Gender      Age      Driving_License
## Min.   :      1  Length:381109  Min.   :20.00  Min.   :0.0000
## 1st Qu.: 95278  Class :character  1st Qu.:25.00  1st Qu.:1.0000
## Median :190555  Mode  :character  Median :36.00  Median :1.0000
## Mean   :190555              Mean   :38.82  Mean   :0.9979
## 3rd Qu.:285832              3rd Qu.:49.00  3rd Qu.:1.0000
## Max.   :381109              Max.   :85.00  Max.   :1.0000
## Region_Code  Previously_Insured  Vehicle_Age  Vehicle_Damage
## Min.   : 0.00  Min.   :0.0000  Length:381109  Length:381109
## 1st Qu.:15.00  1st Qu.:0.0000  Class :character  Class :character
## Median :28.00  Median :0.0000  Mode  :character  Mode  :character
## Mean   :26.39  Mean   :0.4582
## 3rd Qu.:35.00  3rd Qu.:1.0000
## Max.   :52.00  Max.   :1.0000
## Annual_Premium  Policy_Sales_Channel  Vintage  Response
## Min.   : 2630  Min.   : 1      Min.   : 10.0  Min.   :0.0000
## 1st Qu.: 24405  1st Qu.: 29      1st Qu.: 82.0  1st Qu.:0.0000
## Median : 31669  Median :133      Median :154.0  Median :0.0000
## Mean   : 30564  Mean   :112      Mean   :154.3  Mean   :0.1226
## 3rd Qu.: 39400  3rd Qu.:152      3rd Qu.:227.0  3rd Qu.:0.0000
## Max.   :540165  Max.   :163      Max.   :299.0  Max.   :1.0000
```

Per poter essere utilizzate efficientemente, alcune colonne hanno bisogno di essere manipolate. In particolare:

- La colonna `Gender`, viene trasformata nella colonna `Male`, che ammette valori binari 0 e 1 dove 1 rappresenta il sesso maschile e 0 quello femminile.

```
insurance[which(insurance$Gender == "Male"),]$Gender = 1
insurance[which(insurance$Gender == "Female"),]$Gender = 0
colnames(insurance)[2] = "Male"
insurance$Male <- as.factor(insurance$Male)
```

- La colonna `Vehicle_Age` viene convertita da categorica ordinata a numerica. I valori utilizzati sono:
 - -1 per le auto immatricolate da meno di un anno

- 0 per le auto immatricolate tra gli 1 e i 2 anni precedenti
- 1 per le auto immatricolate da più di 2 anni

```
insurance[which(insurance$Vehicle_Age == "< 1 Year"),]$Vehicle_Age <- -1
insurance[which(insurance$Vehicle_Age == "1-2 Year"),]$Vehicle_Age <- 0
insurance[which(insurance$Vehicle_Age == "> 2 Years"),]$Vehicle_Age <- 1
insurance$Vehicle_Age <- as.numeric(insurance$Vehicle_Age)
```

- La colonna Vehicle_Damage viene convertita da categorica nominale a factor. I valori utilizzati sono:
 - 0 per le auto che non sono mai state coinvolte in incidenti
 - 1 per le auto che sono state coinvolte in incidenti

```
insurance[which(insurance$Vehicle_Damage == "Yes"),]$Vehicle_Damage <- 1
insurance[which(insurance$Vehicle_Damage == "No"),]$Vehicle_Damage <- 0
insurance$Vehicle_Damage <- as.factor(insurance$Vehicle_Damage)
```

Le colonne categoriche e binarie restanti vengono convertite in factor

```
insurance$Driving_License <- as.factor(insurance$Driving_License)
insurance$Region_Code <- as.factor(insurance$Region_Code)
insurance$Previously_Insured <- as.factor(insurance$Previously_Insured)
insurance$Policy_Sales_Channel <- as.factor(insurance$Policy_Sales_Channel)
insurance$Response <- as.factor(insurance$Response)
```

Attraverso la funzione ggplot si generano i grafici che mostrano la distribuzione di frequenza delle variabili numeriche Age, Vintage, Annual_Premium

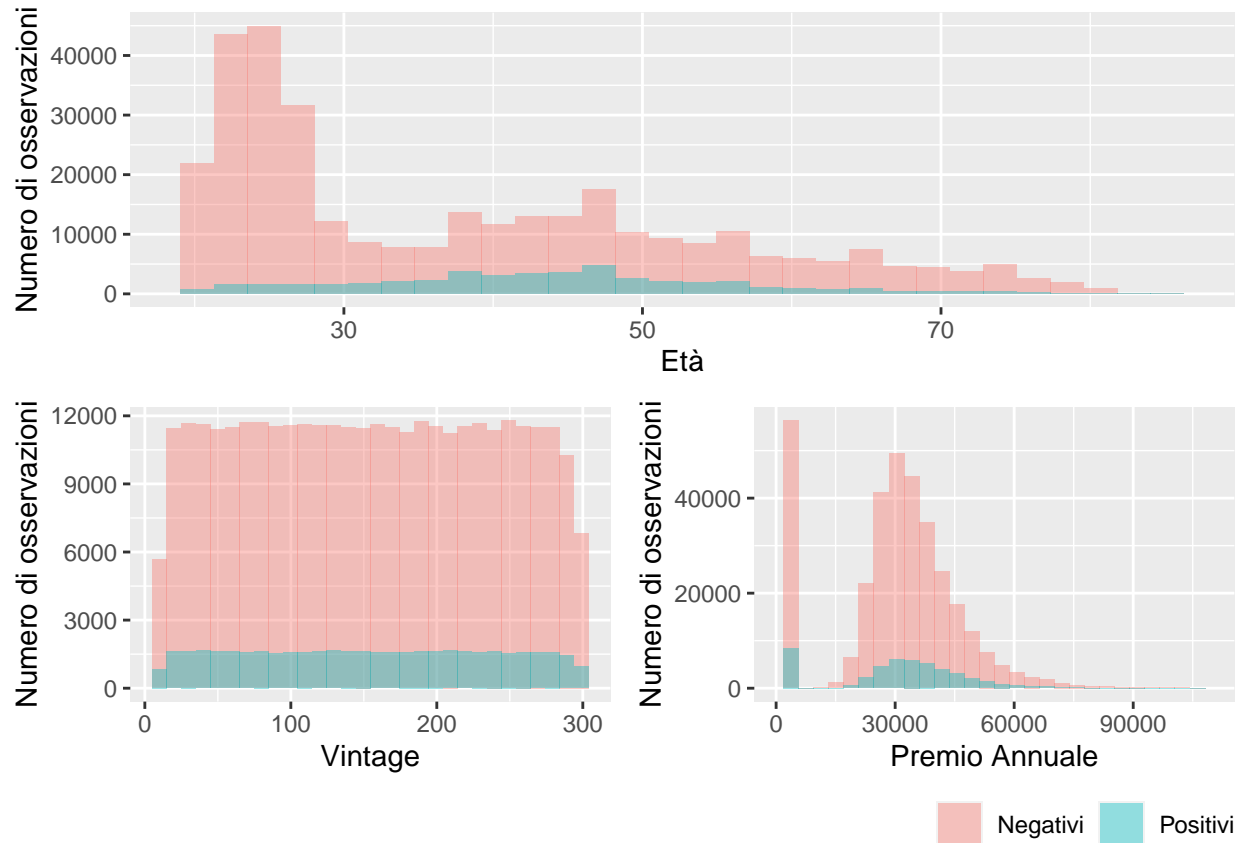
```
library(ggplot2)

age_plot <- ggplot(insurance,
  aes(x=Age, fill=Response))+
  geom_histogram(alpha=0.4,position="identity")+
  xlab("Età")+
  scale_y_continuous("Numero di osservazioni") + guides(fill=guide_legend(title=NULL)) +
  scale_fill_discrete(labels=c("Negativi","Positivi")) +
  theme(legend.position = c(1,1),legend.justification=c(1,1))

#Frequenze Vintage
vintage_plot <- ggplot(insurance,
  aes( x=Vintage, fill=Response))+
  geom_histogram(alpha=0.4,position="identity")+
  xlab("Vintage")+
  scale_y_continuous("Numero di osservazioni")

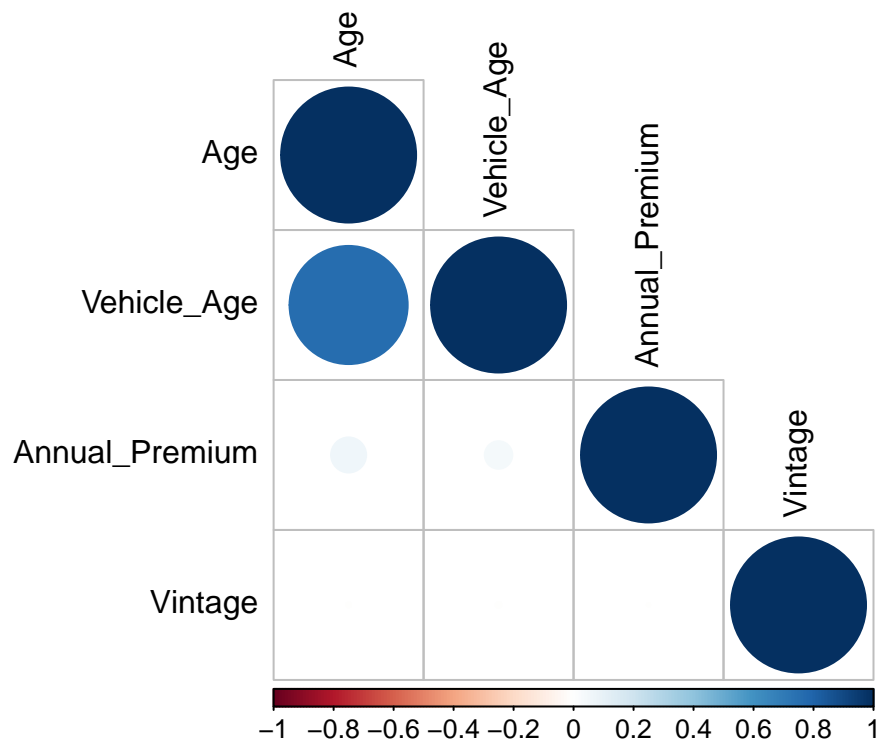
#Frequenze annual_premium
premium_plot <- ggplot( insurance,
  aes( x=Annual_Premium, fill=Response))+
  geom_histogram(alpha=0.4,position="identity")+
  xlab("Premio Annuale")+
  scale_y_continuous("Numero di osservazioni") +
  scale_x_continuous(limits = c(0, 110000))

library(ggpubr)
ggarrange(age_plot, ggarrange(vintage_plot, premium_plot, ncol = 2, legend = "none"),
  nrow = 2, common.legend = T, legend = "bottom")
```



Nonostante il numero non elevato di features presenti nel dataset possiamo cercare la correlazione presente tra esse. Dalla correlation matrix si può notare che le features **Vehicle_Age** e **Age** sono piuttosto correlate tra loro. Nonostante questo ho deciso di mantenerle entrambe poiché il vantaggio dal punto di vista computazionale è minimo, evitando così perdita di informazioni

```
library(corrplot)
correlationMatrix <- stats::cor(insurance[c(3, 7, 9, 11)])
corrplot(correlationMatrix, method="circle", type="lower", tl.col="black")
```



Uno degli aspetti più critici del dataset in analisi è lo sbilanciamento tra le classi, come si evince dal seguente piechart

```
freq <- as.data.frame(table(insurance$Response))
colnames(freq)[1] = "Response"
freq$perc <- prop.table(freq$Freq)
print(table(insurance$Response))
```

```
##
##      0      1
## 334399 46710
```

```
print(freq)
```

```
##  Response  Freq      perc
## 1         0 334399 0.8774366
## 2         1  46710 0.1225634
```

```
ggplot(freq, aes(x = "", y = perc, fill = Response)) +
  geom_col(color = "black") +
  coord_polar(theta = "y") +
  xlab("") +
  ylab("") +
  ggtitle("Frequenza delle classi") +
  theme(plot.title = element_text(hjust = 0.5))
```

Frequenza delle classi

