

Machine Learning Project

Customer Personality Analysis

Mario Avolio
880995

Rocco Gianni Rapisarda
845197

11 gennaio 2022

Sommario

L'apprendimento automatico e la statistica sono discipline strettamente collegate. Secondo [Michael I. Jordan](#), le idee dell'apprendimento automatico, dai principi metodologici agli strumenti teorici, sono stati sviluppati prima in statistica. In questo elaborato si riporta l'attività di sviluppo e sperimentazione di diversi modelli di **Machine Learning** per l'analisi approfondita dei clienti ideali per una generica azienda. La concentrazione è stata focalizzata soprattutto sul **Clustering** mediante l'algoritmo **K-Means**, sebbene nel corso della trattazione si esporranno anche altre metodologie utilizzate per l'analisi dei dati.

1 Descrizione del dominio di riferimento e obiettivi dell'elaborato

Molto spesso lo sviluppo di nuovi prodotti o servizi è attivato dall'imitazione dei concorrenti e da analisi di mercato generiche, mentre al cliente si dedica poca attenzione. L'acquisto è prima di tutto un'esperienza ed è necessario comprendere quali bisogni la guidano: solo così ogni segmento di mercato individuato sarà connesso con la capacità dell'azienda di soddisfare le aspettative dei clienti, comprese quelle inesprese. Progettare, sviluppare e vendere prodotti non connessi con il proprio target rappresenta un costo insostenibile, mentre è necessario progettare uno sviluppo in linea con la *customer satisfaction*. Per questo motivo [Customer Personality Analysis](#) riguarda un'analisi dettagliata dei clienti ideali per una generica azienda. Il compito fondamentale è quello di aiutare un'attività commerciale a comprendere meglio i propri compratori al fine di rendere più semplice la modifica e la scelta dei propri prodotti, in relazione alle esigenze richieste dagli acquirenti. L'obiettivo che ha spinto ad analizzare questo insieme di dati è inerente alle **diverse** personalità e comportamenti che gli acquirenti assumono durante il ruolo di potenziali clienti aziendali. Per questo motivo le aziende non possono adottare lo stesso approccio per ogni tipologia di plausibile compratore.

2 Scelte di design, ipotesi e assunzioni

TODO

3 Descrizione del training set

3.1 Attributi

Si fornisce la descrizione originale degli attributi analizzati.

People

- ID: Customer's unique identifier
- Year_Birth: Customer's birth year
- Education: Customer's education level
- Marital_Status: Customer's marital status
- Income: Customer's yearly household income
- Kidhome: Number of children in customer's household
- Teenhome: Number of teenagers in customer's household
- Dt_Customer: Date of customer's enrollment with the company
- Recency: Number of days since customer's last purchase
- Complain: 1 if the customer complained in the last 2 years, 0 otherwise

Products

- MntWines: Amount spent on wine in last 2 years
- MntFruits: Amount spent on fruits in last 2 years
- MntMeatProducts: Amount spent on meat in last 2 years
- MntFishProducts: Amount spent on fish in last 2 years
- MntSweetProducts: Amount spent on sweets in last 2 years
- MntGoldProds: Amount spent on gold in last 2 years

Promotion

- NumDealsPurchases: Number of purchases made with a discount
- AcceptedCmp1: 1 if customer accepted the offer in the 1st campaign, 0 otherwise
- AcceptedCmp2: 1 if customer accepted the offer in the 2nd campaign, 0 otherwise
- AcceptedCmp3: 1 if customer accepted the offer in the 3rd campaign, 0 otherwise
- AcceptedCmp4: 1 if customer accepted the offer in the 4th campaign, 0 otherwise
- AcceptedCmp5: 1 if customer accepted the offer in the 5th campaign, 0 otherwise
- Response: 1 if customer accepted the offer in the last campaign, 0 otherwise

Place

- NumWebPurchases: Number of purchases made through the company's website
- NumCatalogPurchases: Number of purchases made using a catalogue
- NumStorePurchases: Number of purchases made directly in stores
- NumWebVisitsMonth: Number of visits to company's website in the last month

La tabella 1 fornisce un'iniziale descrizione della tipologia di variabili presenti nel dataset.

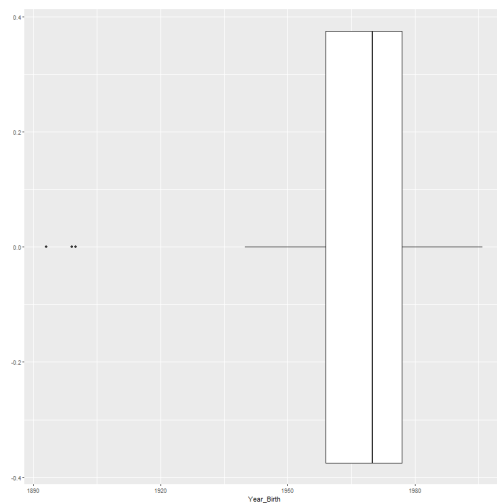
	supply(customers, class)
ID	integer
Year_Birth	integer
Education	character
Marital_Status	character
Income	integer
Kidhome	integer
Teenhome	integer
Dt_Customer	character
Recency	integer
MntWines	integer
MntFruits	integer
MntMeatProducts	integer
MntFishProducts	integer
MntSweetProducts	integer
MntGoldProds	integer
NumDealsPurchases	integer
NumWebPurchases	integer
NumCatalogPurchases	integer
NumStorePurchases	integer
NumWebVisitsMonth	integer
AcceptedCmp3	integer
AcceptedCmp4	integer
AcceptedCmp5	integer
AcceptedCmp1	integer
AcceptedCmp2	integer
Complain	integer
Z_CostContact	integer
Z_Revenue	integer
Response	integer

Tabella 1: Output funzione *supply(customers, class)*

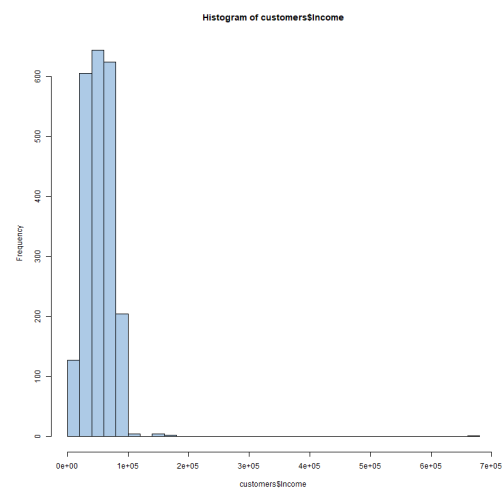
3.2 Prime analisi

Prima di fornire un'analisi dettagliata degli elementi del dataset si è ritenuto necessario effettuare una prima ispezione di alto livello, senza entrare nel dettaglio di ciascun attributo. Il dataset viene importato mediante la funzione:

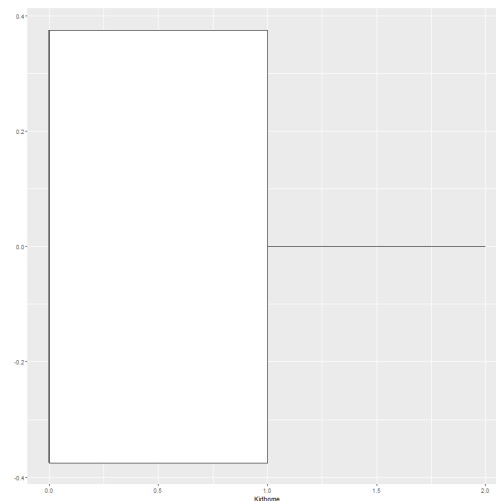
```
customers <- read.csv(paste(getwd(), "/Data/marketing_campaign.csv", sep = ""),  
  header=TRUE, sep="\t", stringsAsFactors=F) # use TAB as separator!
```



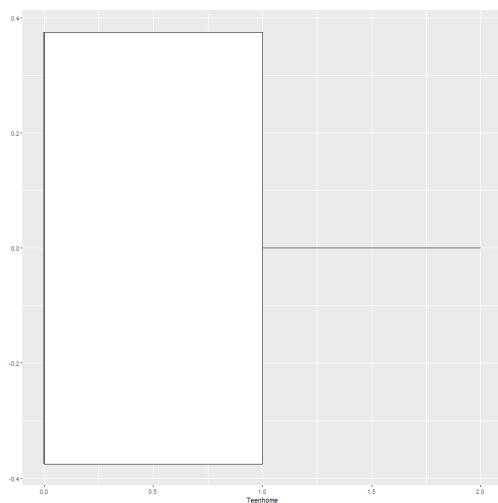
(a) BoxPlot Year_Birth



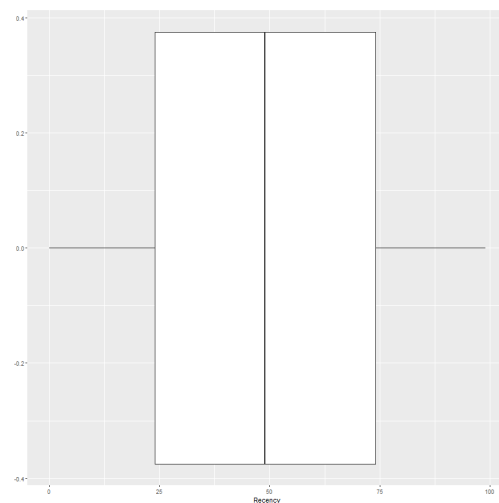
(b) BoxPlot Income



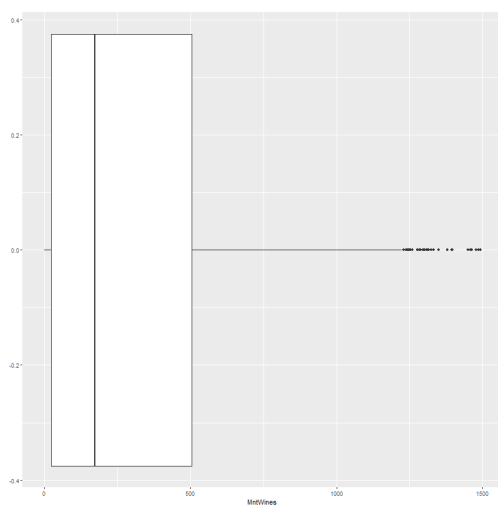
(c) BoxPlot KidHome



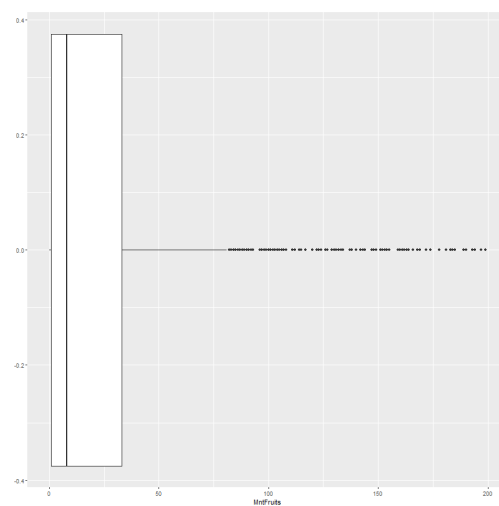
(d) BoxPlot TeenHome



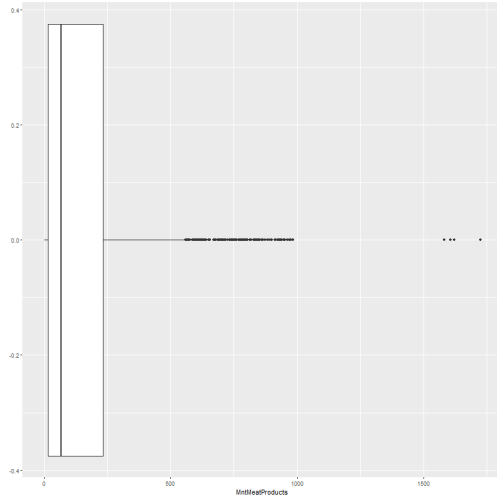
(e) BoxPlot Recency



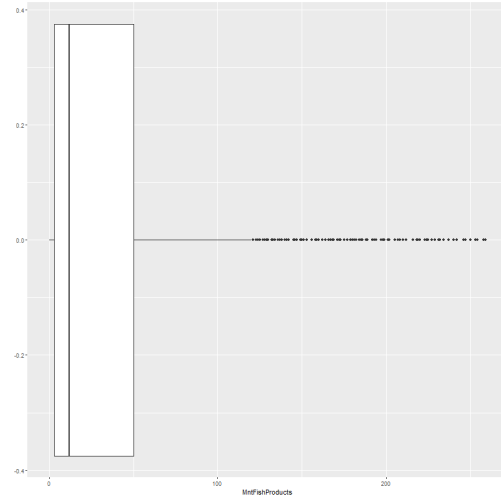
(f) BoxPlot MntWines



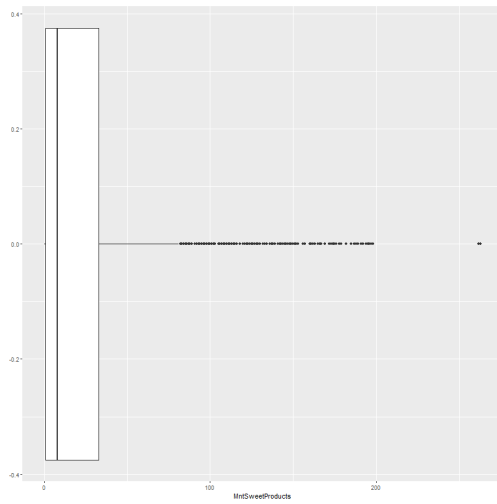
(g) BoxPlot MntFruits



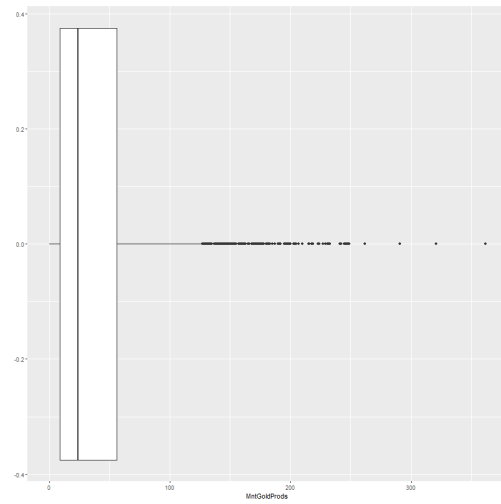
(h) BoxPlot MntMeatProducts



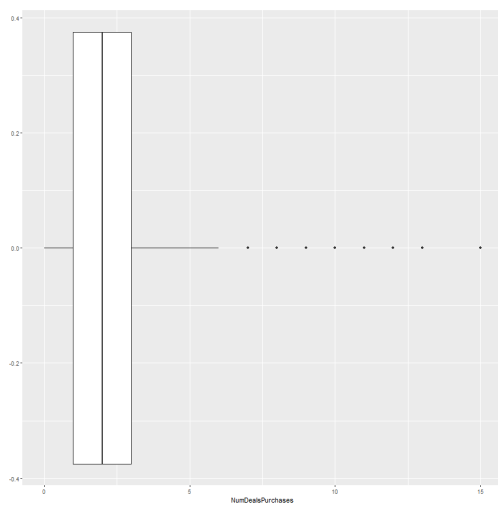
(i) BoxPlot MntFishProducts



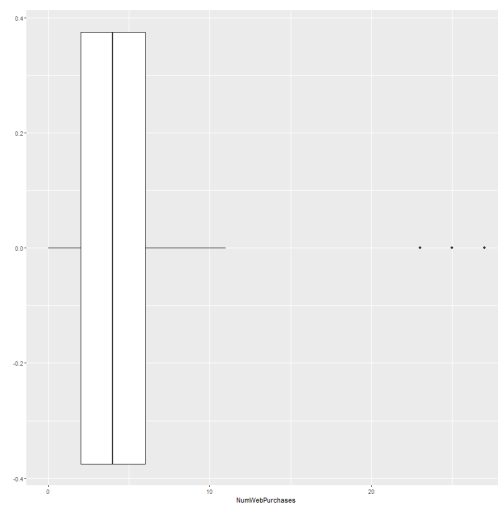
(j) BoxPlot MntSweetProducts



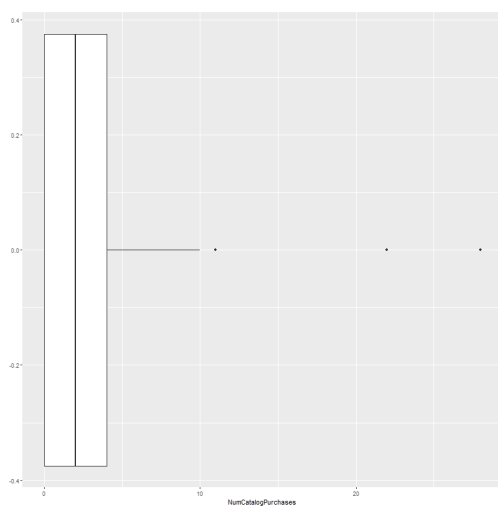
(k) BoxPlot MntGoldProds



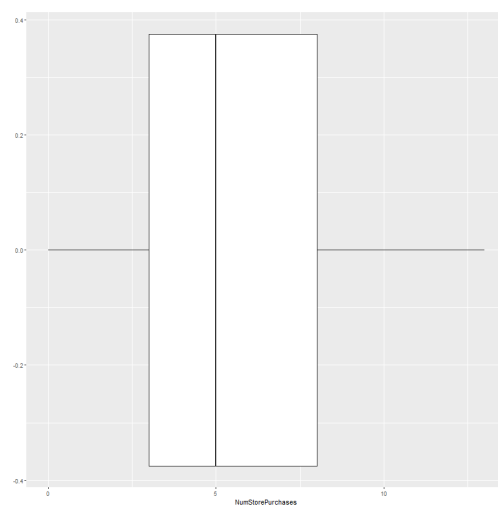
(l) BoxPlot numDealsPurchases



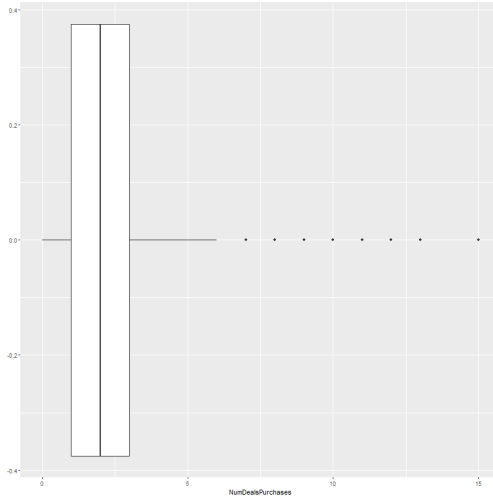
(m) BoxPlot NumWebPurchases



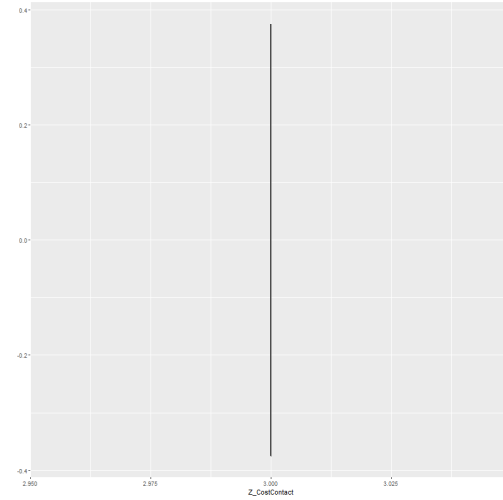
(n) BoxPlot NumCatalogPurchases



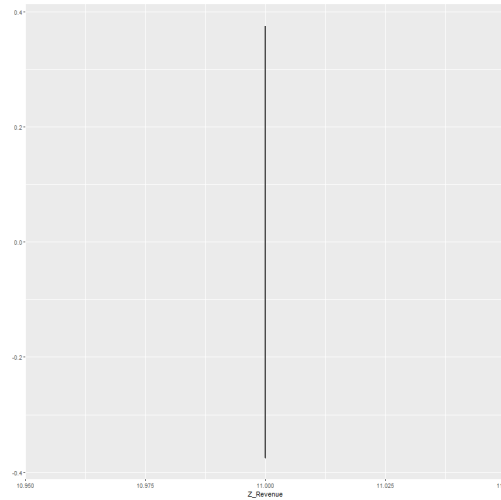
(o) BoxPlot NumStorePurchases



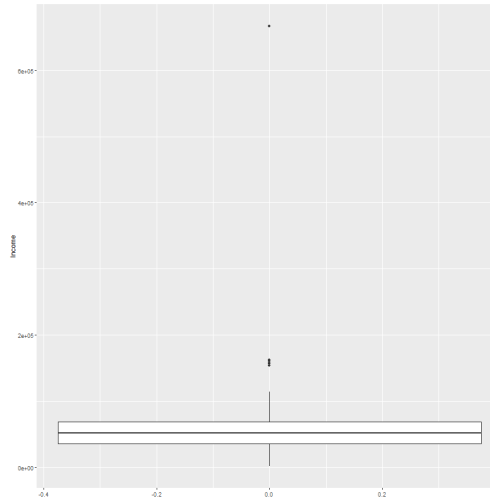
(p) BoxPlot NumDealsPurchases



(q) BoxPlot Z_CostContact



(r) BoxPlot Z_Revenue



(s) BoxPlot Income

Durante questa prima indagine sono stati effettuati controlli di base sulle variabili. In particolare, dalle analisi dei boxplot riportati rispettivamente nelle figure 1(r) e 1(q) è doveroso notare la mancanza di **varianza** di tali variabili, questo aspetto sarà successivamente preso in considerazione durante la fase di preprocessing.

Inoltre durante l'analisi delle figure 1(s) e 1(b), mediante un apposito *warning* durante l'esecuzione del codice R, si è notata la presenza di alcuni valori mancanti.

```
hist(customers$Income,40,col="#adcae6")
ggplot(customers, aes(y = Income)) + geom_boxplot()
n_miss(customers) # counting the total number of missing values in the data
miss_var_summary(customers) # Summarizing missingness in each variable
```

La tabella 2 mostra l'output del comando `miss_var_summary(customers)`, dove si possono notare 24 valori mancanti nella variabile *Income*.

	variable	n_miss	pct_miss
1	Income	24	1.07
2	ID	0	0.00
3	Year_Birth	0	0.00
4	Education	0	0.00
5	Marital_Status	0	0.00
6	Kidhome	0	0.00
7	Teenhome	0	0.00
8	Dt_Customer	0	0.00
9	Recency	0	0.00
10	MntWines	0	0.00
11	MntFruits	0	0.00
12	MntMeatProducts	0	0.00
13	MntFishProducts	0	0.00
14	MntSweetProducts	0	0.00
15	MntGoldProds	0	0.00
16	NumDealsPurchases	0	0.00
17	NumWebPurchases	0	0.00
18	NumCatalogPurchases	0	0.00
19	NumStorePurchases	0	0.00
20	NumWebVisitsMonth	0	0.00
21	AcceptedCmp3	0	0.00
22	AcceptedCmp4	0	0.00
23	AcceptedCmp5	0	0.00
24	AcceptedCmp1	0	0.00
25	AcceptedCmp2	0	0.00
26	Complain	0	0.00
27	Z_CostContact	0	0.00
28	Z_Revenue	0	0.00
29	Response	0	0.00

Tabella 2: Output funzione *miss_var_summary(customers)*

3.3 Data Preprocessing

La fase di *data preprocessing* è risultata fondamentale per la buona riuscita dell'analisi dei dati preposti. Essa si basa su quattro principali stadi:

- Refactor del dataset
- Risoluzione dei valori mancanti nella variabile *income*
- Splitting del dataset in *trainingSet* e *testSet*
- Feature Scaling

3.3.1 Refactor del Dataset

In questa fase ci si è concentrati su diversi aspetti migliorativi. In primo luogo si è voluto esettuare un'azione di incorporamento di dati. La motivazione è dovuta principalmente alla presenza di dati

ridondati all'interno di molte variabili, in particolare si vuole porre l'attenzione su: *Marital_Status* e *Education*. Difatti utilizzando la funzione *unique*, come riportato nelle tabelle 4 e 3, si è potuto rilevare la presenza di troppi elementi superflui.

	unique(customers\$Marital_Status)
1	Single
2	Together
3	Married
4	Divorced
5	Widow
6	Alone
7	Absurd
8	YOLO

Tabella 3: Output *unique(customers\$Marital_Status)*

	unique(customers\$Education)
1	Graduation
2	PhD
3	Master
4	Basic
5	2n Cycle

Tabella 4: Output *unique(customers\$Education)*

Come mostrato dal codice sottostante, l'obiettivo è stato quello di diminuire, per favorire l'analisi mediante i diversi algoritmi utilizzati, in numero delle categorie di valori per le relative variabili.

```
# ----- COLLAPSING
#Collapsing marital Status into two categories: Single & Couple
unique(customers$Marital_Status)
customers <- mutate(customers, Marital_Status = replace(Marital_Status, Marital_Status
  == "Divorced" | Marital_Status == "Widow" | Marital_Status == "Alone" |
  Marital_Status == "Absurd" | Marital_Status == "YOLO", "Single"))
customers <- mutate(customers, Marital_Status = replace(Marital_Status, Marital_Status
  == "Together" | Marital_Status == "Married", "Couple"))

#Collapsing the Education into two Categories: graduate and non-graduate
unique(customers$Education)
customers <- mutate(customers, Education = replace(Education, Education == "Graduation" |
  Education == "PhD" | Education == "Master", "graduate"))
customers <- mutate(customers, Education = replace(Education, Education == "Basic" |
  Education == "2n Cycle", "non-graduate"))
# -----
```

In particolar modo si è deciso di fornire due possibili valori riassuntivi per ciascuna variabile, come indicato nelle tabelle 5 e 6.

La fase di *refactoring* si è anche occupata della conversione in *factor* degli elementi *character* all'interno dell'insieme di dati. Il codice sottostante mostra la procedura seguita. Le tabelle 8 e 7 mostrano il risultato di tale procedura.

```
# ----- CONVERSION
```

unique(customers\$Education)	
1	graduate
2	non-graduate

Tabella 5: Output `unique(customers$Education)` dopo la procedura di *collapsing* dei dati

unique(customers\$Marital_Status)	
1	Single
2	Couple

Tabella 6: Output `unique(customers$Marital_Status)` dopo la procedura di *collapsing* dei dati

```
#Converting them to factors
customers <- mutate(customers, Marital_Status = as.factor(Marital_Status), Education =
  as.factor(Education))

# Encoding the categorical features to numeric
customers <- mutate(customers, Education = case_when(Education == "graduate" ~ 1,
  Education == "non-graduate" ~ 0))
customers <- mutate(customers, Marital_Status = case_when(Marital_Status == "Couple" ~ 1,
  Marital_Status == "Single" ~ 0))
# -----
```

head(customers\$Marital_Status)	
1	0.00
2	0.00
3	1.00
4	1.00
5	1.00
6	1.00

Tabella 7: Output `head(customers$Marital_Status)`

head(customers\$Education)	
1	1.00
2	1.00
3	1.00
4	1.00
5	1.00
6	1.00

Tabella 8: Output `head(customers$Education)`

Questa fase si è anche occupata della creazione di nuove variabili partendo da quelle già presenti all'interno da quelle già esistenti. Si ponga l'attenzione in particolar modo alle **categorie** di variabili *Mnt*, *Accepted*, *KidHome*, *TeenHome*. Tali categorie possono essere sommate per creare nuove variabili riassuntive. Il codice seguente e la tabella 9 ne riportano un esempio:

```
# ----- TOTAL
#Creating a new variable:Total_spent
```

```

customers <- mutate(customers, Total_spent = MntWines + MntFruits + MntMeatProducts +
  MntFishProducts + MntSweetProducts + MntGoldProds)

# Details about previous campains also combined. Creating a new variable:Total_Campains
customers <- mutate(customers, Total_Campains = AcceptedCmp1 + AcceptedCmp2 +
  AcceptedCmp3 + AcceptedCmp4 + AcceptedCmp5)

# These variables can be combined and we can get the no of children for the customers.
  Creating a new variable:Total_Childs
customers <- mutate(customers, Total_Childs = Kidhome + Teenhome)
# -----

```

	Total_spent	Total_Campains	Total_Childs
1	1617	0	0
2	27	0	2
3	776	0	0
4	53	0	1
5	422	0	1
6	716	0	1

Tabella 9: Primi valori delle variabili Total_spent & Total_Campains & Total_Childs

Per finire è doveroso sottolineare che in questa sezioni ci si è anche occupati dell'eliminazione delle variabili superflue, come *Z_CostContact* e *Z_Revenue* che non hanno varianza, e della sostituzione della varibile *Year_Birth* con *Age*.

```

# we can calculate customer age from the birth year. It will be more usefull to our
  analysis.
# creating a new variable Age from Year of Birth
thisYear <- as.numeric(format(as.Date(Sys.Date(), format="%d-%m-%Y"), "%Y"))
thisYear
customers <- mutate(customers, Age = thisYear - Year_Birth)

#Dropping some redundant features
customers <- select(customers, - ID, - Year_Birth, - Z_CostContact, - Z_Revenue,
  -Dt_Customer)

```

3.3.2 Risoluzione dei valori mancanti

Come mostrato nel codice sottostante, non si è deciso di eliminare completamente i valori mancanti all'interno della variabile *income*, bensì si è adottata un'altra strategia: la sostituzione con il valore medio della variabile stessa.

```

customers$Income <- ifelse(is.na(customers$Income), # is.na check is a value is not
  available
  ave(customers$Income, FUN = function(x) mean(x, na.rm = TRUE)), #
  if is not available change with average
  customers$Income # else
)

```

3.3.3 Splitting in TrainingSet e TestSet

Al fine di una buona analisi dei dati, si è deciso di suddividere il dataset iniziale in due più ridotti: l'insieme di allenamento e quello di testing. Il codice sottostante mostra la procedura effettuata.

```
set.seed(17538)
split <- sample.split(customers$Response, SplitRatio = 0.8)
trainingSet <- subset(customers, split == TRUE)
testSet <- subset(customers, split == FALSE)
```

3.3.4 Feature Scaling

E' opportuno sottolineare che si è deciso di normalizzare i valori delle variabili al fine di poter utilizzare al meglio gli algoritmi di machine learning che verranno descritti nel corso di questa trattazione. Il codice seguente mostra un esempio:

```
trainingSet_scaled <- as.data.frame(scale(trainingSet[, getIndipendentNumbersOfCol()]))
testSet_scaled <- as.data.frame(scale(testSet[, getIndipendentNumbersOfCol()])))
```

3.4 EDA

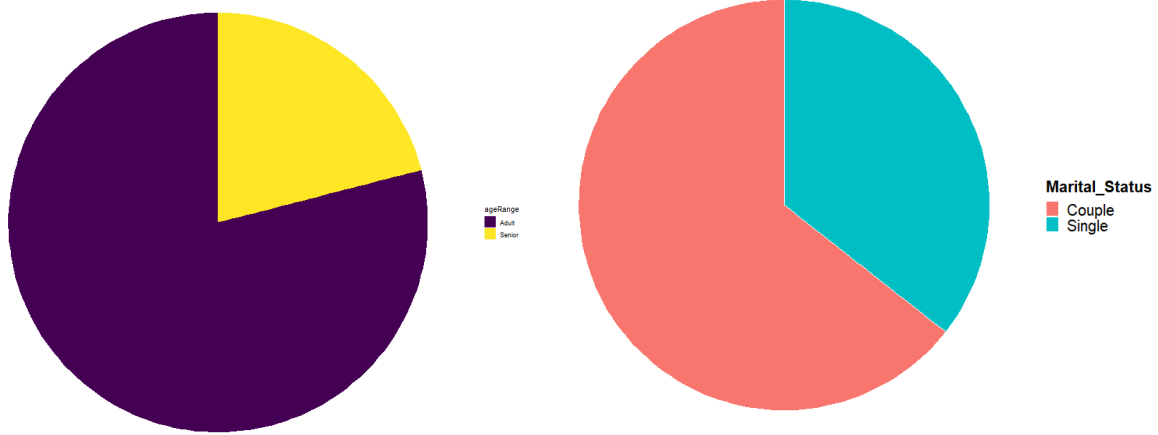
Dopo aver eseguito il preprocessing dei dati si è passati ad un'analisi esplorativa dei dati mutando i valori che può assumere una certa variabile tramite la funzione *cut*. Questa operazione è stata eseguita per *Age*, *Income* e *Total spent*.

```
# Age Range
ageRange <- cut(trainingSet$Age, breaks = c(24, 64, Inf), include.lowest = T,
               ordered_result = T, labels = c("Adult", "Senior"))
trainingSet <- mutate(trainingSet, Age_range = ageRange)

# Income Range
incomeRange <- cut(trainingSet$Income,
                  calculateBreaksFromSummary(trainingSet$Income),
                  labels = c("low", "low medium", "medium high", "high"))
trainingSet <- mutate(trainingSet, Income_range = incomeRange)

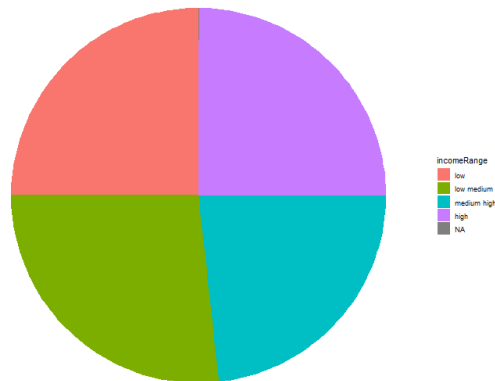
# Spent Range
spentRange <- cut(trainingSet$Total_spent,
                 calculateBreaksFromSummary(trainingSet$Total_spent),
                 labels = c("low", "low medium", "medium high", "high"))
trainingSet <- mutate(trainingSet, Spent_range = spentRange)
```

Pie chart of marital status



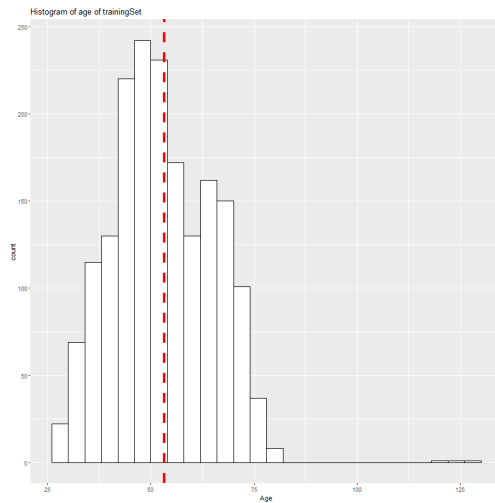
(t) PiePlot Year_Birth

(u) PiePlot Income

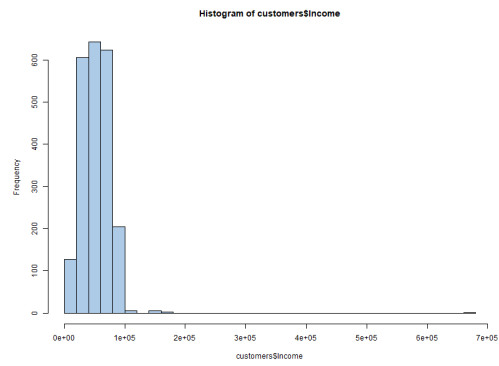


(v) PiePlot di Marital_Status

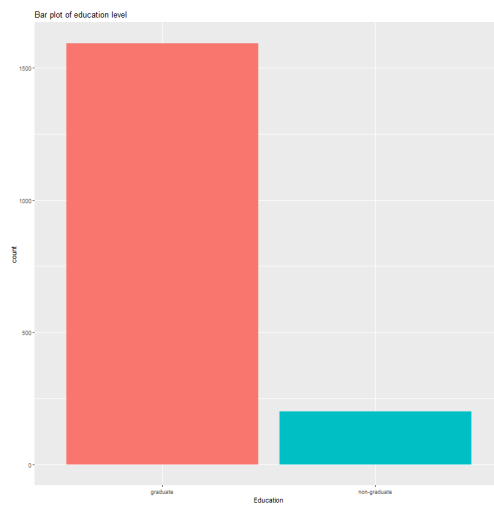
Il primo grafico a torta (t) è relativo alla variabile Age, il colore viola rappresenta il valore *Adult* mentre il colore giallo *Senior*. Dal grafico si ricava che la maggior parte degli individui è *Adult*. Il secondo rappresenta la variabile *income*, il dataset è equi-distribuito in questo caso. La distribuzione dei valori che assume la variabile *Total spent* è raffigurata nella Figura 3. Da essa si ricava che low e high sono simili mentre quello più frequente è *low medium*. Inoltre si può notare facendo l'istogramma della variabile *Age* che l'età media degli individui presenti nel dataset è superiore ai 50 anni e la maggioranza ha un titolo di studio maggiore o uguale alla laurea.



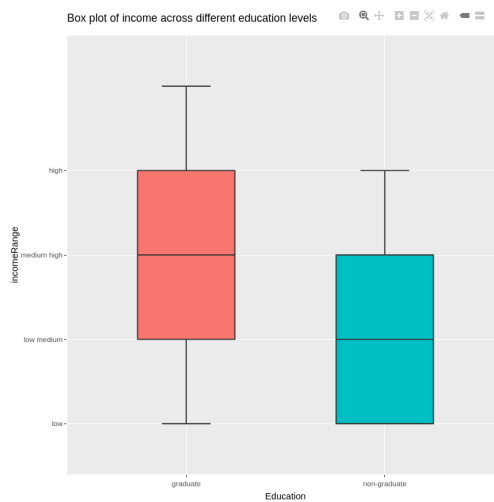
(w) Istogramma di Age.



(x) Istogramma di Income.



(y) Istogramma di Education.



(z) BoxPlot di Education.

3.4.1 Total Children

La maggior parte delle istanze del dataset ha 1 figlio, la variabile `Total_Children` è la somma tra la variabile `KidHome` e `TeenHome`. Questa informazione si è ricavata eseguendo il codice riportato di seguito.

```
ggplot(trainingSet, aes(x=Total_Childs)) + geom_histogram(binwidth = 0.5, colour =  
  "Black")
```

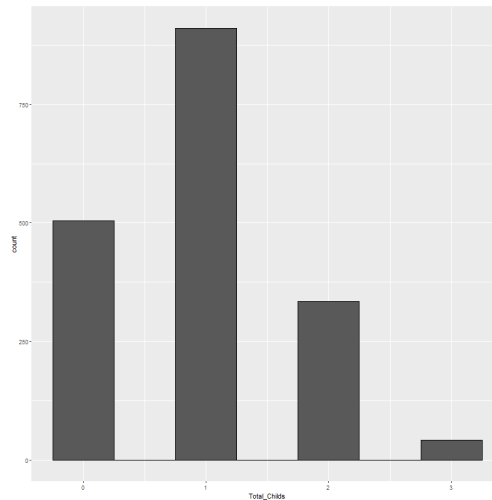


Figura 1: hist plot di total-children

3.4.1.1 Total_Childen e Age

Si è analizzata anche la relazione tra il numero totale di figli, *Total.children* ed *Age* ed è risultato che tra gli *Adult* il numero di figli più frequente è 1 e che rispetto ai *Senior* hanno più figli.

```
age_children_histogram <- ggplot(trainingSet, aes(x=Total_Spent)) +  
  geom_histogram(aes(fill=Age_range), binwidth = 0.5, colour = "Black")  
age_children_histogram + facet_grid(Age_range~.)
```

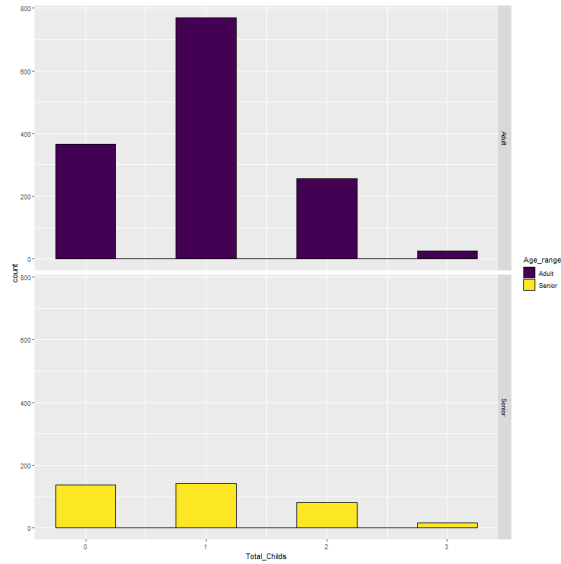



Figura 2: Istogramma di Total_Children e Age

3.4.1.2 Total_Children e Marital_Status

Nel pre processing i valori che può assumere la variabile *Marital_Status* sono stati collassati in due possibili valori. *Marital_Status* è 0 se l'istanza ha come valore dell'attributo *Marital_Status* è *single*, *widow* oppure *divorced*. E' uguale a 1 se il valore assunto da *Marital_Status* è *couple* oppure *together*. Rappresentando il diagramma a barre considerando *Total_Children* e *Marital_Status* si ricava che la maggior parte delle istanze ha un figlio. Il numero di istanze con *Marital_Status* pari a 0 è minore rispetto a quelle con valore uguale a 1 per i casi di zero e due figli, mentre per il caso di tre figli sono molto simili.

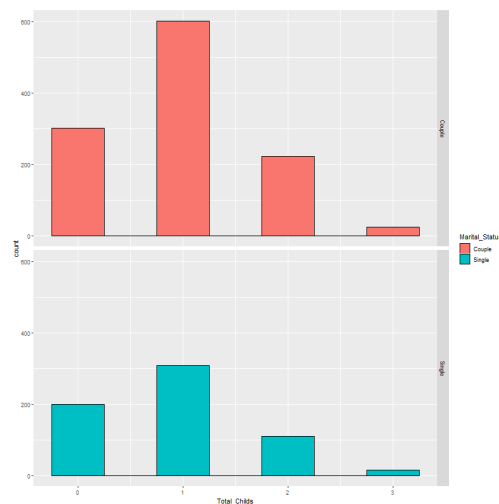


Figura 3: Istogramma di Total_Children e Marital_Status

Il grafico prendendo in considerazione la variabile *Education* è simile.

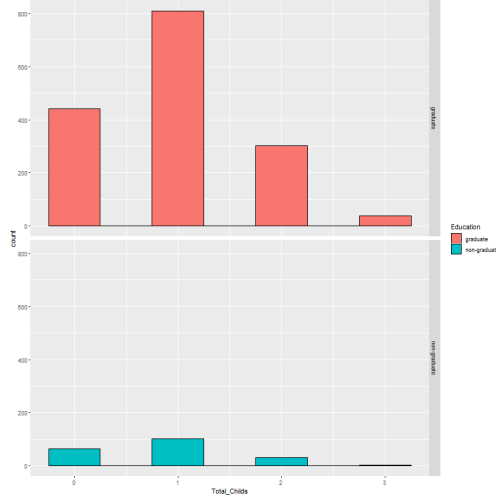


Figura 4: hist plot di Total_Children e Education

3.4.1.3 Total_Children e Income

Sia dall' hist plot che dal jitted rplot si nota che più è basso il guadagno più si hanno figli. Nel caso di due figli ci sono più casi di persone con stipendio *low medium* rispetto a chi ha uno stipendio *low*. Nel caso di istanze con stipendio *high* è comune che si abbiano figli.

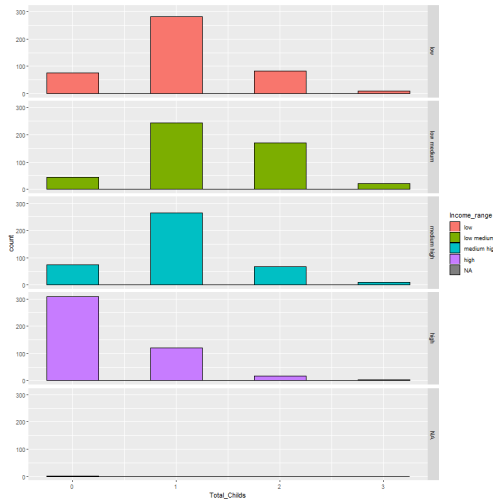


Figura 5: Istogramma Income e Total_Children

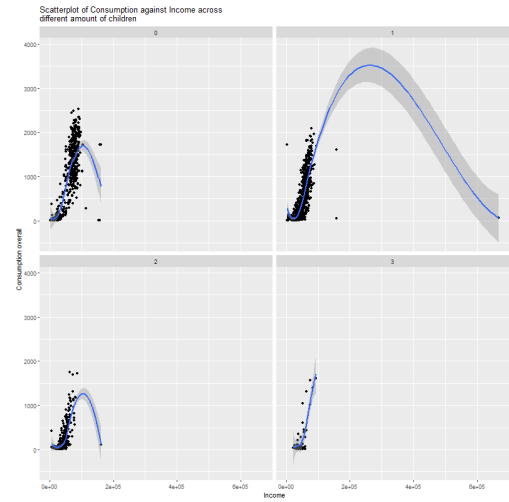


Figura 6: ScatterPlot , Income e Total_Children

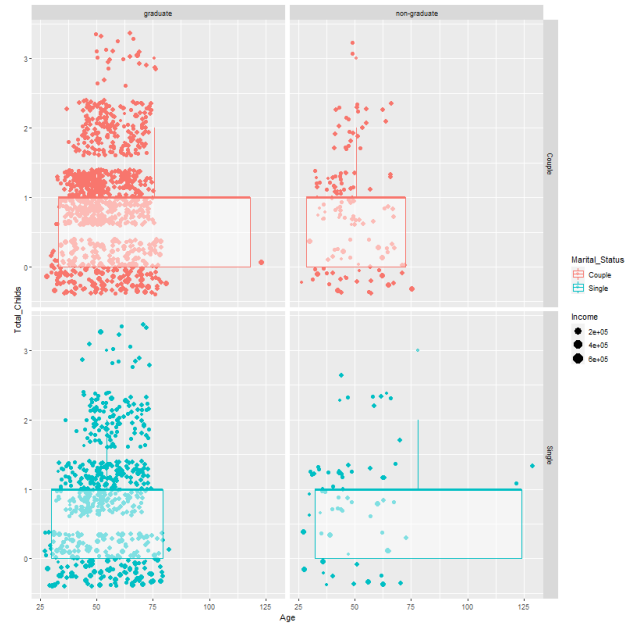


Figura 7: JitterPlot di Total_Children e Income

3.4.2 Total Spent

La maggior parte degli individui spende meno di 500\$. Facendo il bar-plot della variabile *Total_Spent* si nota che le istanze presenti nel dataset spendono più per il vino e per la carne. Per poter produrre il grafico della Figura 8 sono stati sommati i totali che ogni istanza ha speso per un certo prodotto.

```
# Plot histogram of total_spent
ggplot(trainingSet, aes(x=Total_spent)) + geom_histogram(binwidth = 15, colour = "Black")

# Create dataAcceptedCmp
dataTotalSpent <- data.frame(
  name = c("MntWines", "MntFruits", "MntMeatProducts", "MntFishProducts",
    "MntSweetProducts", "MntGoldProds"),
  value = c(sum(trainingSet$MntWines), sum(trainingSet$MntFruits),
    sum(trainingSet$MntMeatProducts),
    sum(trainingSet$MntFishProducts), sum(trainingSet$MntSweetProducts),
    sum(trainingSet$MntGoldProds)))

# Barplot of total_spent for each type
ggplot(dataTotalSpent, aes(x=name, y=value)) +
  geom_bar(stat = "identity", color = "Black") + xlab("Total Spent for each type")
```

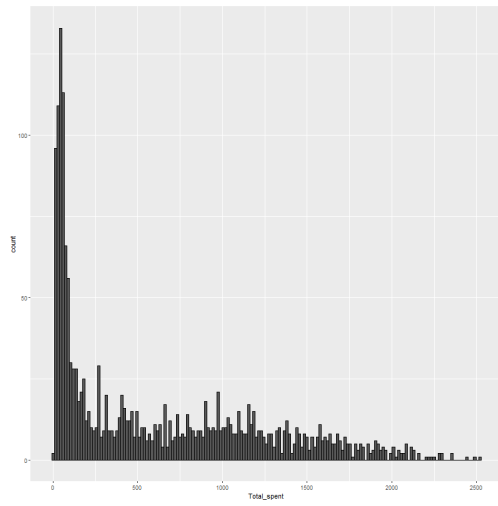


Figura 8: Istogramma Total_Spent.

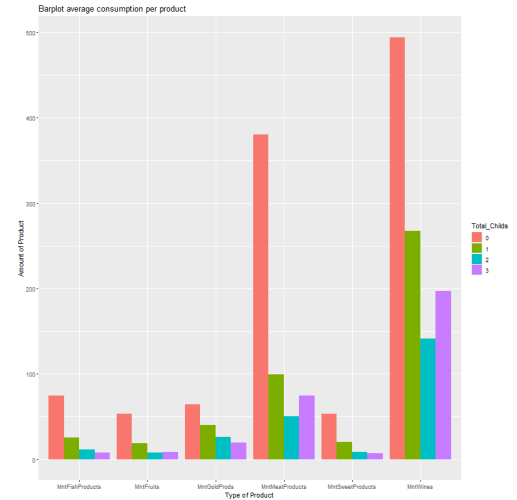


Figura 9: Istogramma di Total_Spent distinguendo i tipi di prodotto.

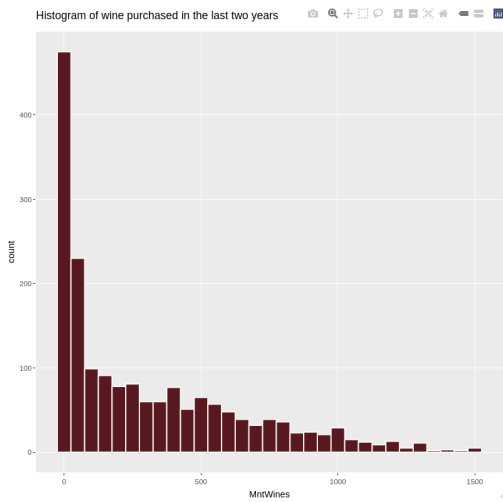


Figura 10: Istogramma Wine

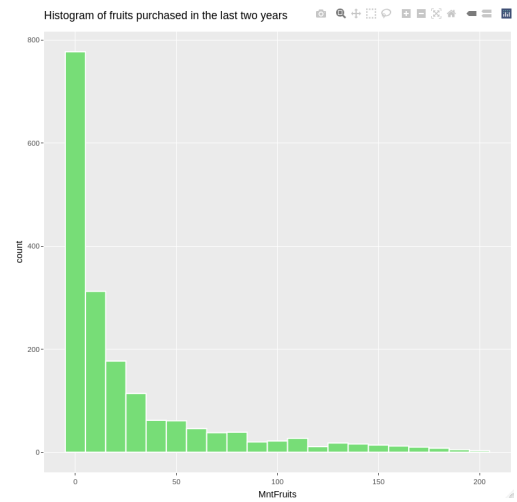


Figura 11: Istogramma Fruit

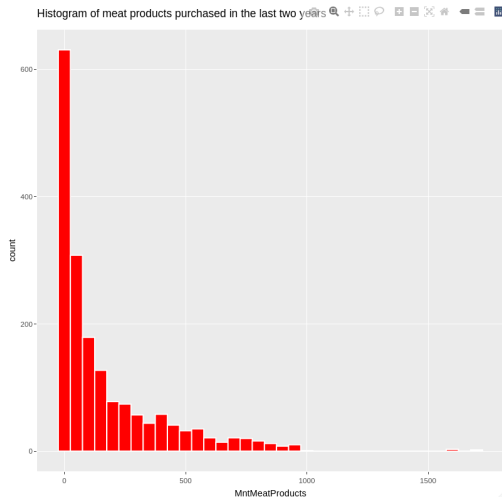


Figura 12: Istogramma Meat

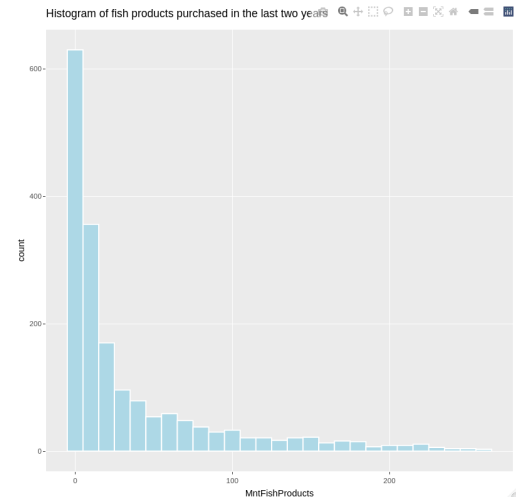


Figura 13: Istogramma Fish

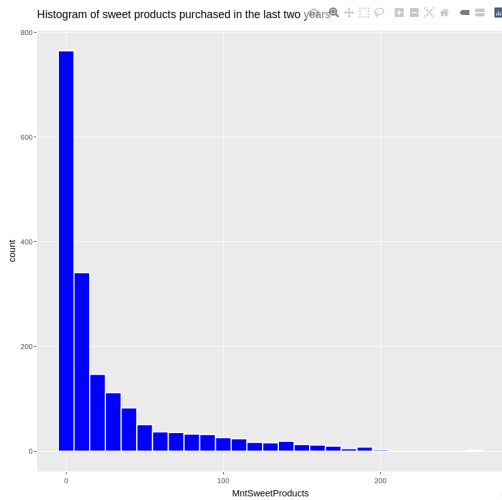


Figura 14: Istogramma Sweet

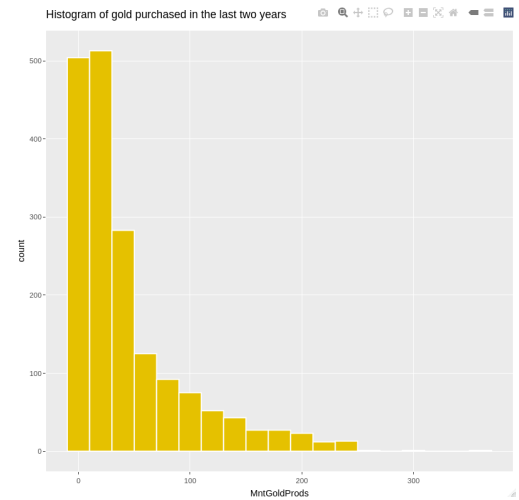


Figura 15: Istogramma Gold

3.4.2.1 Total Spent e Age

Osservando la variabile *Total_Spent* considerando *Age* si ha che la maggior parte degli individui spende meno di 500\$.

```
age_total_spent_histogram <- ggplot(trainingSet, aes(x=Total_spent)) +
  facet_grid(Age_range~.)
age_total_spent_histogram + geom_histogram(aes(fill=Age_range), binwidth = 15,
  color="Black")
age_total_spent_histogram + geom_density(aes(fill=Age_range), position = "Stack")
```

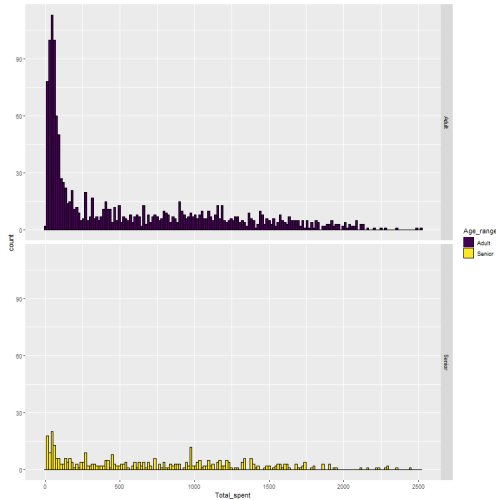


Figura 16: Istogramma di Total_Spent e Age

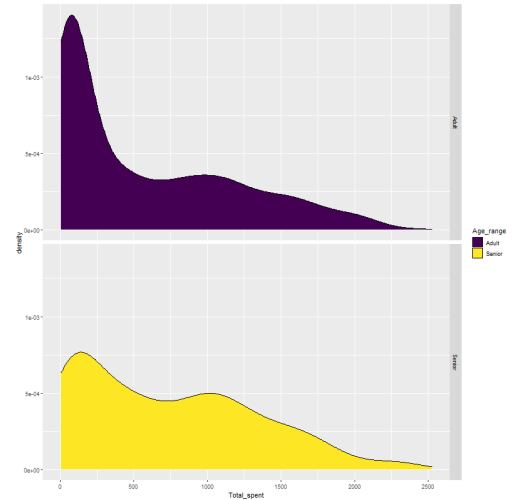


Figura 17: Diagramma di Densità di Total_Spent e Age

3.4.2.2 Total Spent e Marital_Status

In questo caso sembrano simili. la maggior parte delle coppie spende un totale inferiore a 500. La situazione in single è un po' più rilassata che può essere dovuto al fatto che la maggior parte degli individui nel dataset sono coppie.

```
marital_status_total_spent_histogram <- ggplot(trainingSet, aes(x=Total_spent)) +
  facet_grid(Marital_Status~.)
marital_status_total_spent_histogram + geom_histogram(aes(fill=Marital_Status), binwidth =
  15, colour = "Black")
marital_status_total_spent_histogram + geom_density(aes(fill=Marital_Status), position =
  "Stack")
```

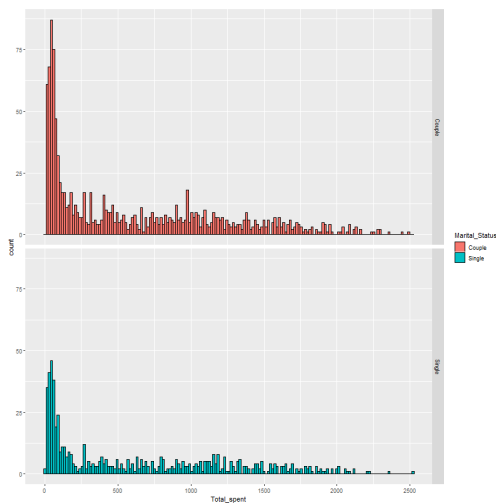


Figura 18: Istogramma di Total_Spent e Marital_Status

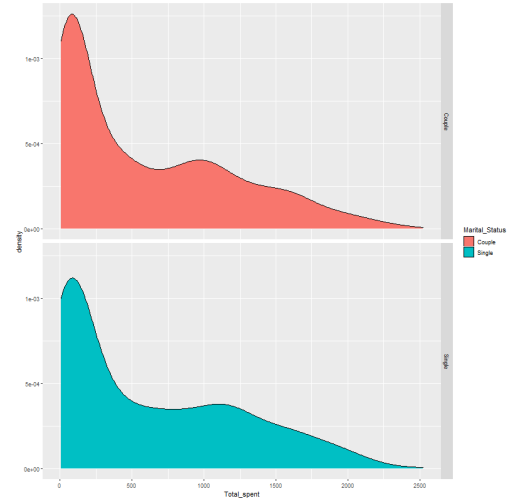


Figura 19: Diagramma di Densità di Total_Spent e Marital_Status

3.4.2.3 Total Spent ed Education

I laureati in genere spendono più dei non laureati. la maggior parte dei non laureati spende da 0 a 1500. da 1500 ci sono più casi di laureati che non laureati. Il grafico della densità è molto simile a quello della variabile *Marital_Status*.

```
education_total_spent_histogram <- ggplot(trainingSet, aes(x=Total_spent)) +  
  facet_grid(Education~.)  
education_total_spent_histogram + geom_histogram(aes(fill=Education), binwidth = 15,  
  colour = "Black")  
education_total_spent_histogram + geom_density(aes(fill=Education), position = "Stack")
```

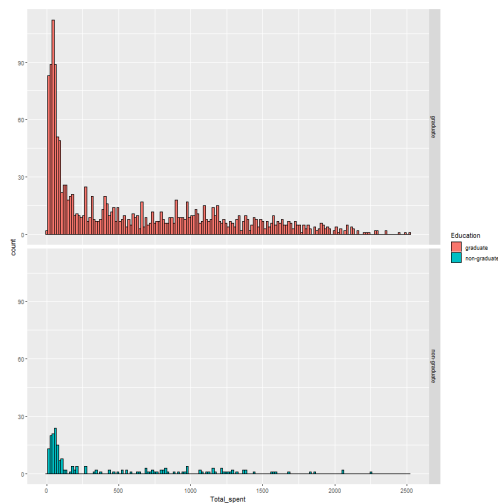


Figura 20: Istogramma di Total_Spent e Education

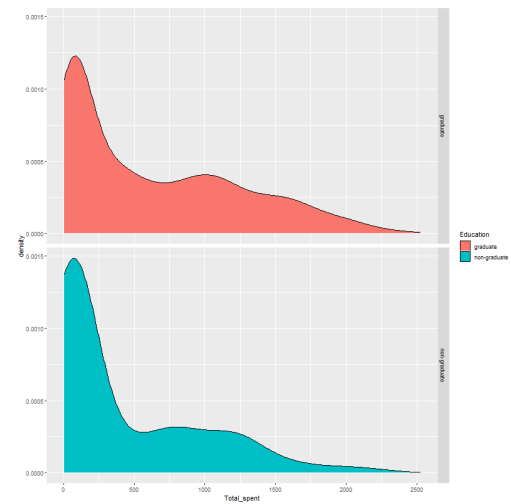


Figura 21: Diagramma di Densità di Total_Spent e Education

3.4.2.4 Total Spent e Total_Children

```
children_total_spent_histogram <- ggplot(trainingSet, aes(x=Total_spent)) +  
  facet_grid(Total_Childs~.)  
children_total_spent_histogram + geom_histogram(aes(fill=Total_Childs), binwidth = 15,  
  colour = "Black")  
children_total_spent_histogram + geom_density(aes(fill=Total_Childs), position = "Stack")
```

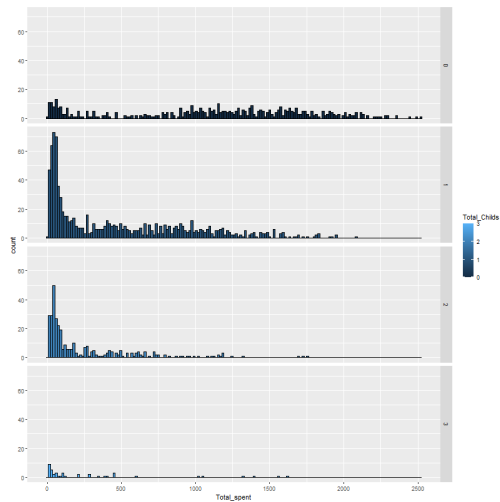


Figura 22: Istogramma di Total_Spent e Total_Children

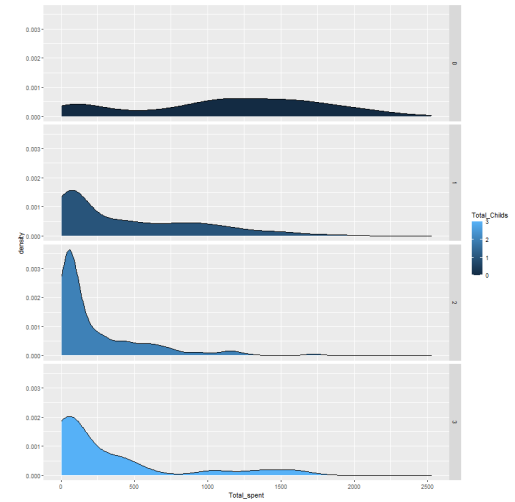


Figura 23: Diagramma di Densità di Total_Spent e Total_Children

3.4.2.5 Total_Spent e Income

```
income_total_spent_histogram <- ggplot(trainingSet, aes(x=Total_spent)) +
  geom_histogram(aes(fill=Income_range), binwidth = 15, colour = "Black")
income_total_spent_histogram + facet_grid(Income_range~.)
ggplot(trainingSet, aes(y=Income, x=Total_spent)) + geom_jitter() + geom_smooth()

#-----

ggplot(trainingSet, aes(x=Total_spent, y=Income, colour=Marital_Status, size=Income)) +
  facet_grid(Marital_Status~Education) +
  geom_jitter() + ylim(0,100000) + geom_smooth() + geom_boxplot(size=0.7, alpha=0.5)

ggplot(trainingSet, aes(x=Total_spent, y=Income, colour=Total_Childs, size=Income)) +
  facet_grid(Total_Childs~Education) +
  geom_jitter() + ylim(0,100000) + geom_smooth() + geom_boxplot(size=0.7, alpha=0.5)

ggplot(trainingSet, aes(x=Total_spent, y=Income, colour=Age_range, size=Income)) +
  facet_grid(Age_range~Education) +
  geom_jitter() + ylim(0,100000) + geom_smooth() + geom_boxplot(size=0.7, alpha=0.5)
```

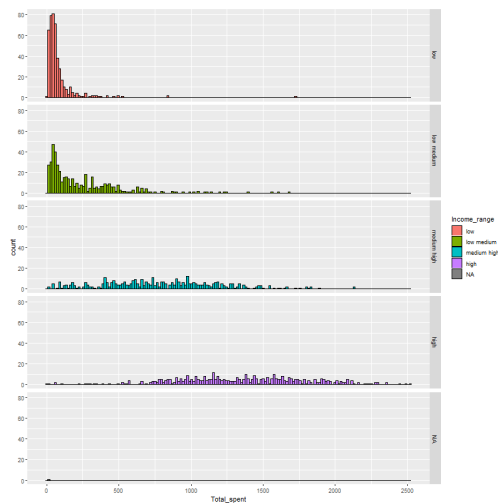


Figura 24: Istogramma di Total_Spent e Income

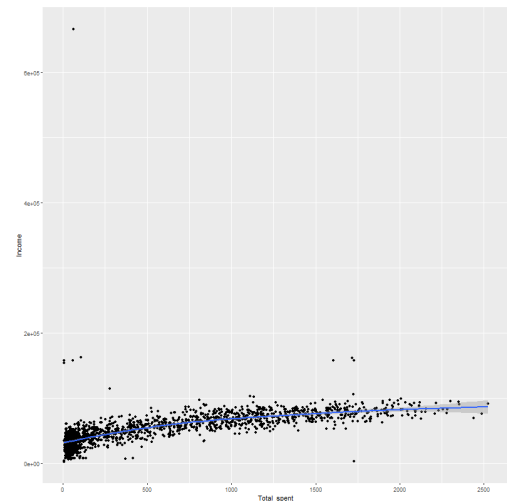


Figura 25: Diagramma di Densità di Total_Spent e Income

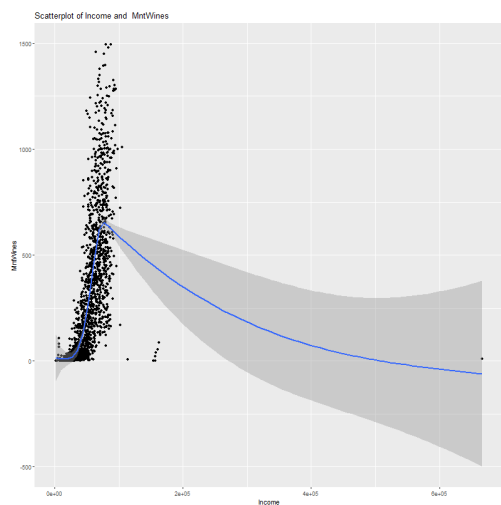


Figura 26: ScatterPlot Income e Mnt-Wines

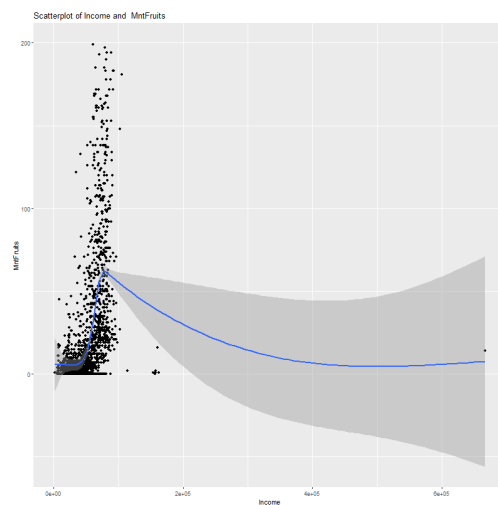


Figura 27: ScatterPlot Income e Mnt-Fruits

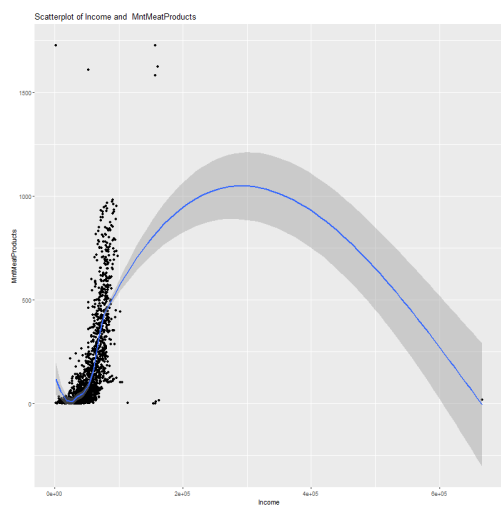


Figura 28: ScatterPlot Income e Mnt-MeatProducts

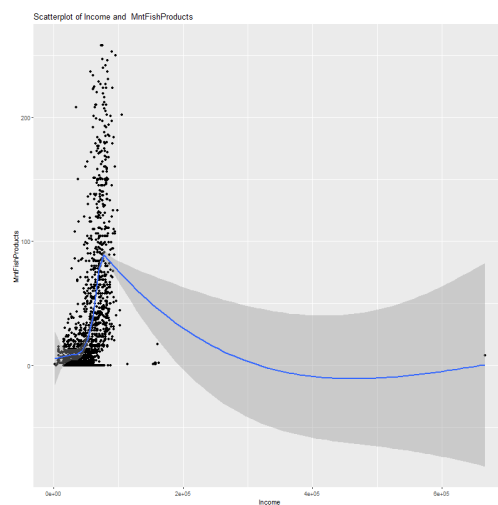
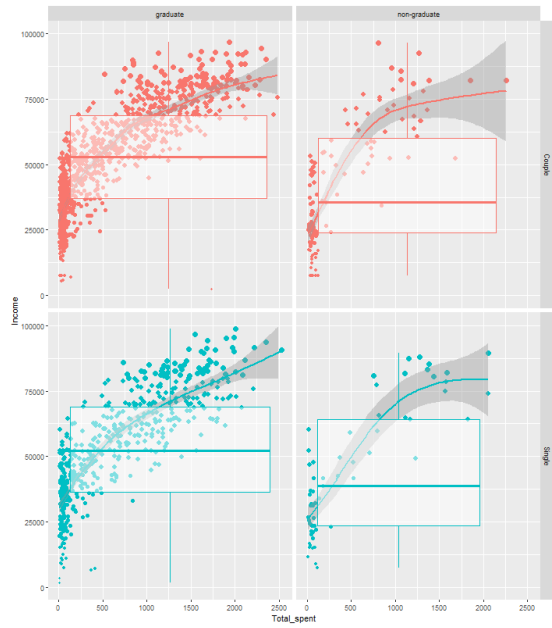
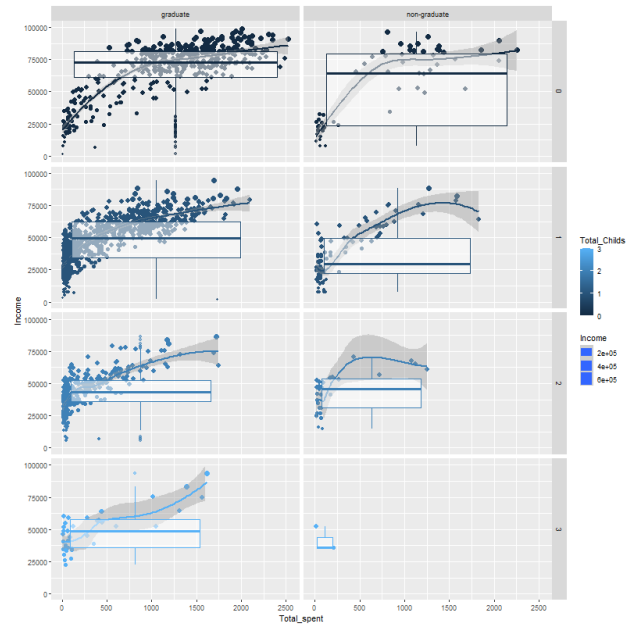


Figura 29: ScatterPlot Income e Mnt-FishProducts

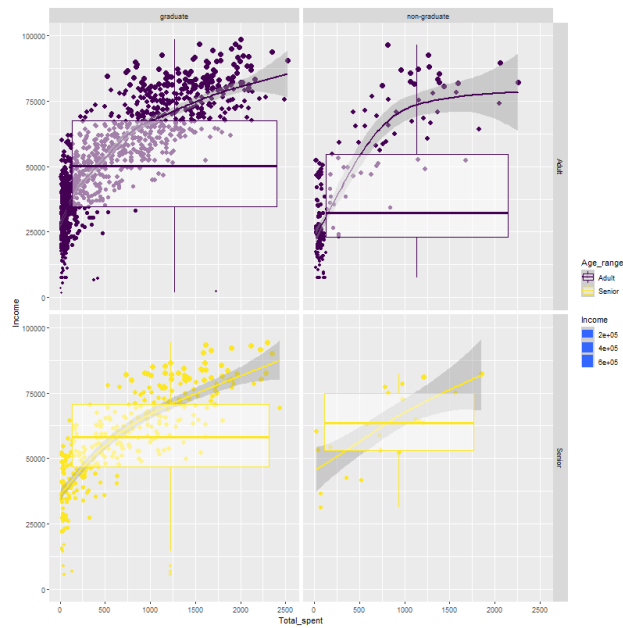
3.4.2.6 Total_Spent JitterPlot+BoxPlot



(a) JitterPlot+BoxPlot di Total_Spent, Income e Mari-



(b) JitterPlot+BoxPlot di Total_Spent, Income e Total_Children.



(c) JitterPlot+BoxPlot di Total_Spent, Income ed Age.

3.4.3 Campaign Analysis

Si è eseguita un'analisi anche sulle campagne accettate da ogni individuo. Per ogni istanza del dataset si sa se ha accettato o meno una campagna. In questo caso le campagne considerate sono cinque e come è possibile osservare dai barplot più di mille istanze non hanno accettato alcuna campagna, mentre chi ha accettato una o più campagne ha scelto la campagna 4.

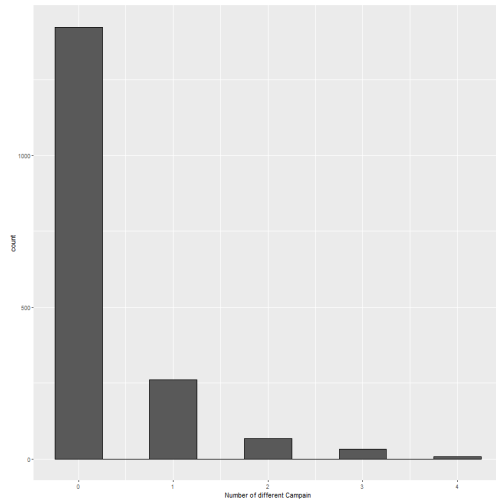


Figura 30: Istogramma del numero totale di campagne accettate da un'istanza. Total_Campaign

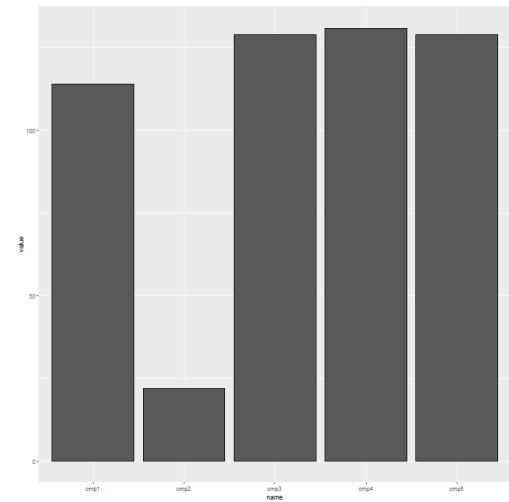


Figura 31: Istogramma del numero di istanze che ha accettato la campagna i-esima.

3.5 PCA

La PCA è stata prevalentemente sfruttata al fine di ridurre il numero elevato di variabili che descrivono l'insieme di dati a un numero minore di variabili latenti, limitando il più possibile la perdita di informazioni. Il codice seguente mostra parte del codice riportato durante l'analisi dei dati:

```
pca <- PCA(trainingSet_scaled, graph = FALSE)

#Getting the variance of the first 9 new dimensions
pca$eig[,2][1:9]

#Getting the cummulative variance
pca$eig[,3][1:5]

#Getting the most correlated variables
dimdesc(pca, axes = 1:2)

# get eigenvalue
get_eigenvalue(pca)

# visualize pca
fviz_eig(pca, addlabels = TRUE, ylim = c(0, 50))
fviz_contrib(pca, choice = "var", axes = 1, top = 5)
fviz_pca_biplot(pca)

#Creating a factor map for the variable contributions
fviz_pca_var(pca, col.var = "contrib", repel = TRUE)

fviz_pca_var(pca, select.var = list(contrib = 5), col.var = "contrib", repel = TRUE)

# Extract the principal components
trainingSet_input <- data.frame(get_pca_ind(pca)$coord)
```

Da esso si vuole dare particolare attenzione alla funzione *get_eigenvalue(pca)* che fornire le informazioni rappresentate nella tabella 10. Si vuole anche fornire un riferimento grafico a quest'ultima tramite l'output della funzione *fviz_eig(pca, addlabels = TRUE, ylim = c(0, 50))* descritto dalla figura 32. Da essa si può notare che le prime 5 dimensioni fornite dalla PCA forniscono il 70% della varianza cumulativa, per questo motivo si è deciso di prendere in considerazione tali dimensioni. Inoltre si vuole sottolineare l'importanza della prima dimensione che riesce a spiegare più del 40% della varianza dei dati.

	eigenvalue	variance.percent	cumulative.variance.percent
Dim.1	6.98	41.06	41.06
Dim.2	1.75	10.32	51.38
Dim.3	1.15	6.78	58.16
Dim.4	1.05	6.19	64.35
Dim.5	1.00	5.87	70.22
Dim.6	0.79	4.66	74.88
Dim.7	0.66	3.91	78.79
Dim.8	0.63	3.70	82.49
Dim.9	0.57	3.33	85.82
Dim.10	0.47	2.77	88.59
Dim.11	0.42	2.49	91.08
Dim.12	0.39	2.30	93.38
Dim.13	0.35	2.05	95.43
Dim.14	0.31	1.81	97.24
Dim.15	0.25	1.46	98.69
Dim.16	0.22	1.31	100.00
Dim.17	0.00	0.00	100.00

Tabella 10: Output funzione *get_eigenvalue(pca)*

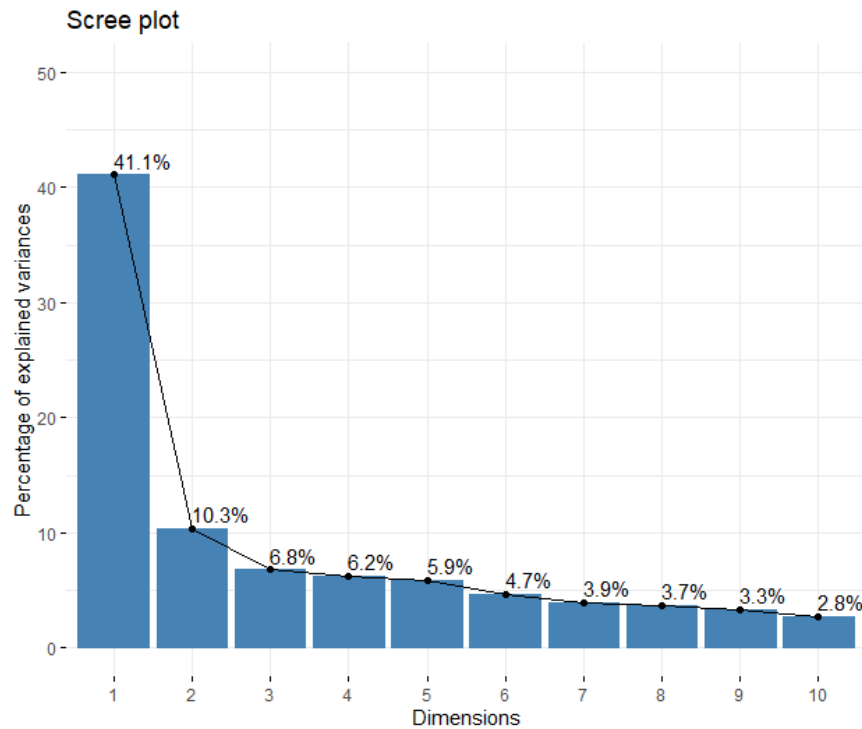


Figura 32: Output funzione *fviz_eig(pca, addlabels = TRUE, ylim = c(0, 50))*

In particolare la tabella 11 vuole far notare il contributo di ciascun attributo del dataset nella creazione delle dimensioni della pca. Da essa possiamo notare le cinque principali variabili che hanno contribuito maggiormente nella creazione della prima dimensione della *principal component*

analysis: Total_spent, MntMeatProducts, NumCatalogPurchases, MntWines e MntFishProducts. La figura 33 ne mostra un grafico più esplicativo.

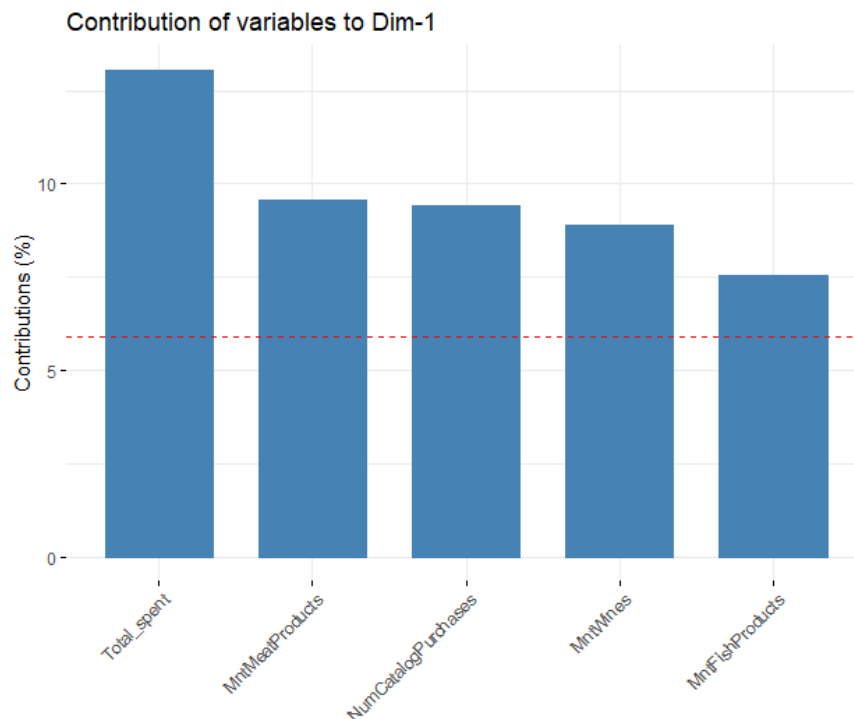


Figura 33: Output funzione *fviz_contrib(pca, choice = "var", axes = 1, top = 5)*

	Dim.1	Dim.2	Dim.3	Dim.4	Dim.5
Income	7.31	0.16	2.59	3.75	0.92
Recency	0.00	0.01	1.24	9.79	87.66
MntWines	8.89	4.77	9.07	1.61	0.88
MntFruits	6.99	1.23	11.85	0.00	0.35
MntMeatProducts	9.55	1.20	0.11	0.00	0.12
MntFishProducts	7.56	1.32	8.51	0.16	0.33
MntSweetProducts	6.94	0.87	9.33	0.06	0.15
MntGoldProds	4.69	2.38	6.52	1.19	0.31
NumDealsPurchases	0.21	36.53	3.87	0.16	0.01
NumWebPurchases	4.30	17.88	0.97	2.21	0.00
NumCatalogPurchases	9.43	0.13	0.70	0.20	0.23
NumStorePurchases	7.53	4.02	0.50	0.32	0.37
NumWebVisitsMonth	5.78	8.73	0.35	12.79	1.36
Total_spent	13.06	0.49	0.82	0.58	0.34
Total_Campains	2.68	0.01	35.64	12.53	3.01
Total_Childs	4.79	14.87	0.20	2.76	0.15
Age	0.28	5.40	7.71	51.88	3.81

Tabella 11: Output *pca\$var\$contrib*

Le funzioni *fviz_pca_var* hanno permesso di analizzare graficamente le dimensioni che spiegano il contributo delle variabili considerate nelle prime due dimensioni (fig. 34 e fig. 35).

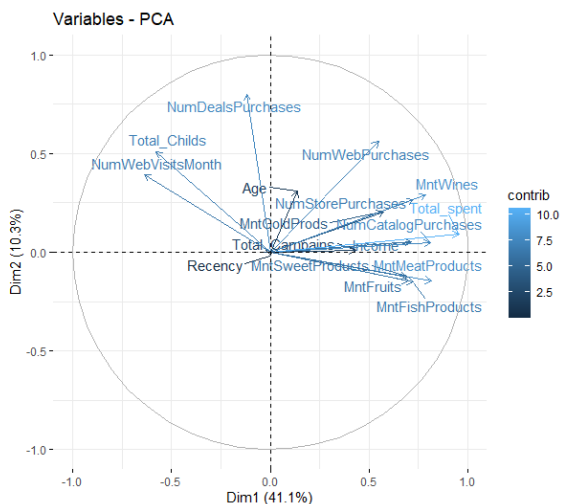


Figura 34: Contributo di tutte la variabili sulle prime due dimensioni.

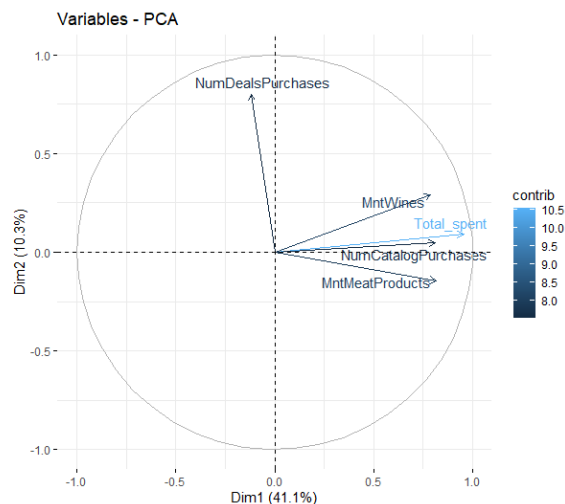


Figura 35: Contributo delle cinque variabili più significative sulle prime due dimensioni.

In fine si è preso in considerazione un nuovo dataset, nato dalle dimensioni fornite dalla PCA, descritto dalla tabella 12. Esso verrà sfruttato durante l'analisi del dataset tramite i diversi algoritmi utilizzati.

	Dim.1	Dim.2	Dim.3	Dim.4	Dim.5
1	4.15	0.41	1.40	0.12	0.27
2	-2.56	-0.14	-0.62	1.48	-0.78
3	1.81	-0.22	0.42	0.13	-1.15
4	-2.46	-0.91	0.29	-0.99	-0.58
5	-0.24	0.50	1.19	0.05	1.42
6	0.65	0.73	-0.16	-0.19	-1.22

Tabella 12: *Head* del Dataset fornito dalla PCA

4 Modelli utilizzati

Il dataset non presenta una variabile target specifica, per questo motivo di è ritenuto opportuno analizzare il tutto mediante algoritmi non supervisionati come quelli di **clustering**. In particolare si è deciso di utilizzare principalmente l'algoritmo **K-Means**. Oltre a ciò si è anche cercato di trovare un target su cui poter fare predizioni mediante alberi decisionali.

4.1 K-Means

La natura stessa dei dati ha comportato l'obbligo di analizzare il tutto mediante un algoritmo come **K-Means**.

4.1.1 Silhouette e Elbow Method

Si è ritenuto opportuno utilizzare metodi come **Elbow Method** e **Silhouette** al fine di determinare il numero ottimale di **clusters** da utilizzare.

```
set.seed(6)
wcss <- vector()
for (i in 1:10) {
  wcss[i] <- sum(kmeans(trainingSet_input, i)$withinss)
}
plot(1:10, wcss, type="b", main = paste('Clusters'), xlab='Number of clusters',
     ylab="WCSS")
```

OR

```
fviz_nbclust(trainingSet_input, kmeans, method="wss")+geom_vline(xintercept=2, linetype=2)
```

Il codice mostra una prima analisi mediante *elbow method*, il cui output si può notare dalle figure 36 e 37. Esse mostrano un'inclinazione tra un numero di clusters compreso tra 2 e 3. Per maggior

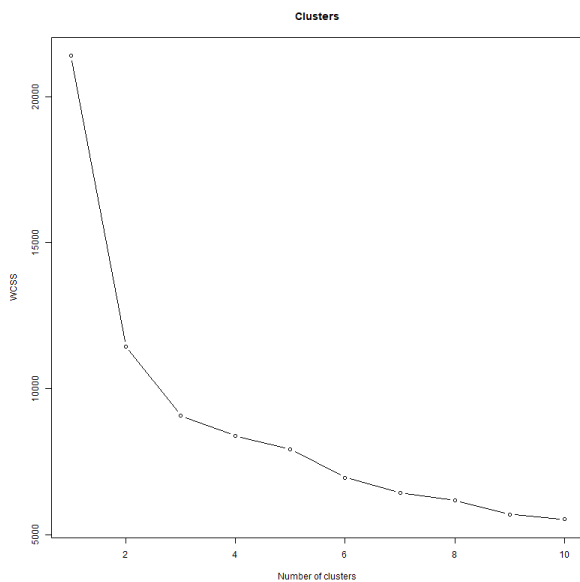


Figura 36: Elbow Method effettuato manualmente

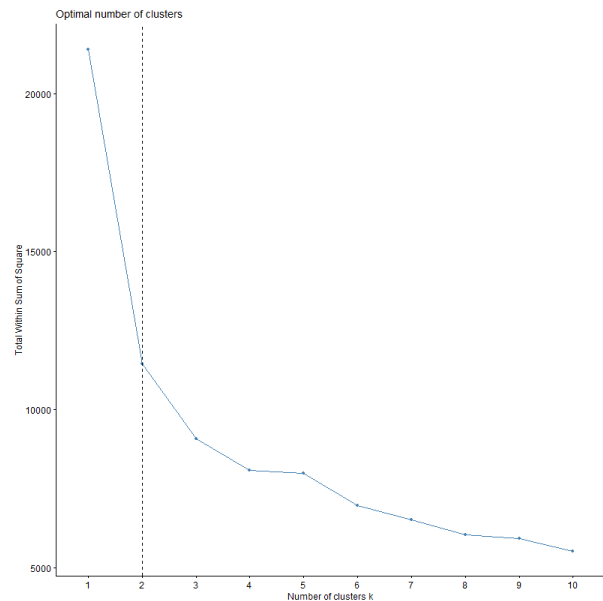


Figura 37: Elbow Method effettuato automaticamente dal metodo *fviz_nbclust*

sicurezza si è quindi deciso di sfruttare anche il metodo *Silhouette* mediante il codice sottostante.

```
k <- 2:10
avg_sil <- sapply(k, silhouette_score)
plot(k, type='b', avg_sil, xlab='Number of clusters', ylab='Average Silhouette Scores',
     frame=FALSE)
avg_sil # <<<- important

# OR
fviz_nbclust(trainingSet_input, kmeans, method="silhouette")
```

Il codice presenta in output i grafici espressi delle figure 38 e 39, da cui si è potuto constatare con maggior sicurezza l'utilizzo di un numero di clusters pari a 2.

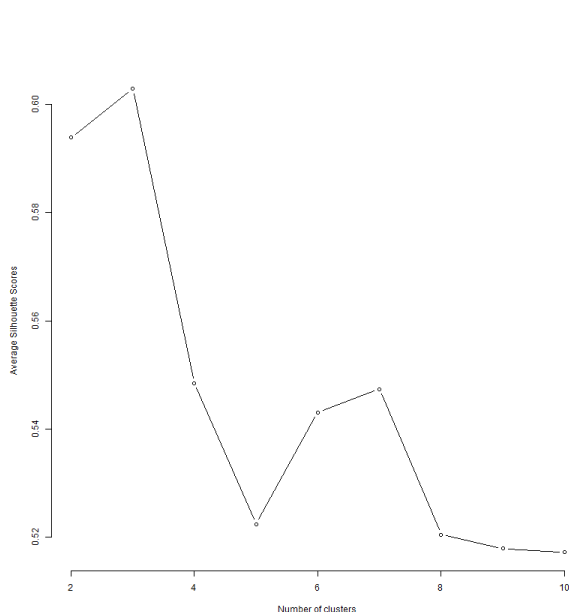


Figura 38: Silhouette effettuata manualmente

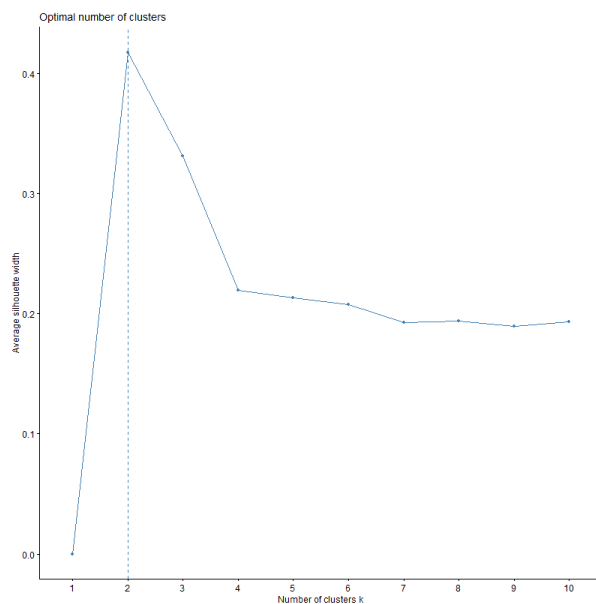


Figura 39: Silhouette effettuata automaticamente dal metodo *fviz_nbclust*

4.1.2 Algoritmo e Analisi

Il codice sottostante viene espresso mediante le figure 40 e 41 che rappresentano rispettivamente il **partizionamento per ogni cluster** e la **Matrice di dissimilarità**. E' doveroso notare che entrambe le figure mostrano un buon partizionamento degli elementi del cluster, sebbene ci sia qualche elemento di intersezione tra i due.

```
set.seed(29)
km <- kmeans(trainingSet_input, 2, nstart = 10)
print(km$centers)
fviz_cluster(km, trainingSet_input, geom = "point", ellipse.type = "norm", repel = TRUE)
cl <- km$cluster
dissplot(dist(trainingSet_input), labels=cl, options=list(main="Kmeans Clustering With
k=2"))
```

In particolare la figura 40 presenta gli elementi all'interno di ogni cluster. La tabella 13 mostra il numero preciso di elementi per ogni cluster. Si è ritenuto necessario analizzare i risultati ottenuti

	cluster	n
1	1	1081
2	2	711

Tabella 13: Numero di elementi per ogni cluster

mediante le variabili più significative, sorte durante la *Principal Component Analysis*.

```
wines <- ggplot(trainingSet, aes(MntWines)) +
```



Figura 40: Partizionamento in clusters dei dati

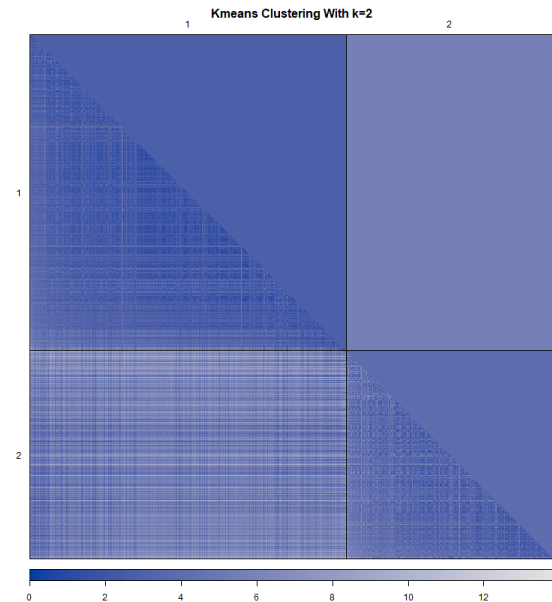


Figura 41: Dissimilarity matrix

```
facet_grid(cluster~.)

wines + geom_histogram(color = "black", fill = "red")
wines + geom_density(fill="red", position = "Stack")
ggplot(trainingSet,
  aes(x=cluster,y=MntWines,fill=cluster))+geom_boxplot(outlier.colour="black")
```

Il codice soprastante è descritto dalle figure 44, 43 e 42. Dall'analisi dei grafici si nota una spesa maggiore di vini per i *customers* all'interno del secondo cluster. In particolare la maggior parte dei clienti all'interno del primo cluster non ha acquistato vini negli ultimi due anni.

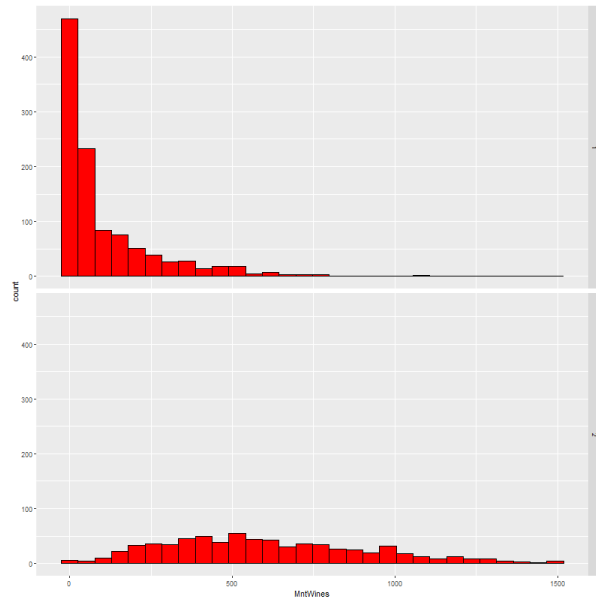


Figura 42: Istogramma della variabile Wines in relazione al numero di cluster

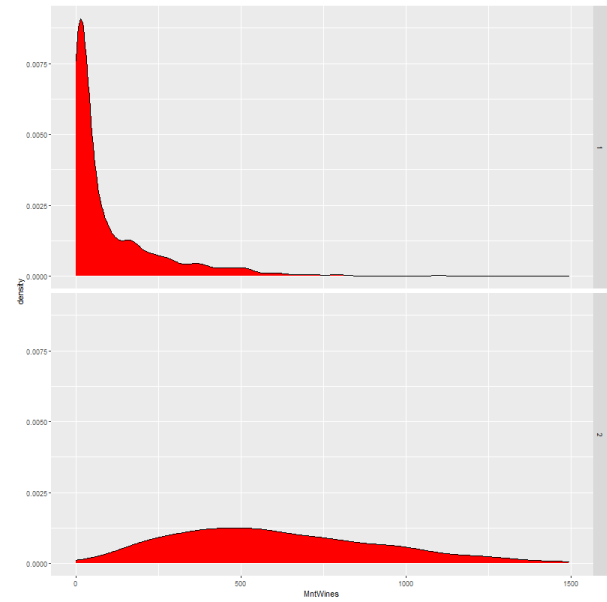


Figura 43: Diagramma di densità della variabile Wines in relazione al numero di cluster

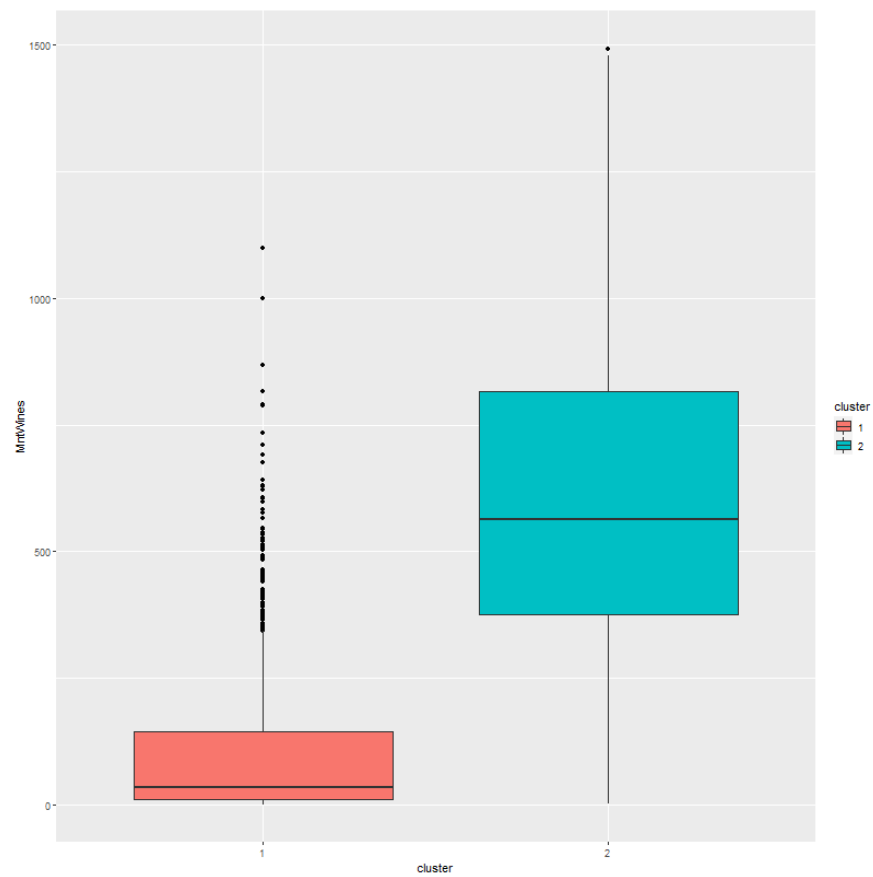


Figura 44: BoxPlot della variabile Wines in relazione al numero di cluster

```
income <- ggplot(trainingSet, aes(Income))+
  facet_grid(cluster~.) +
  xlim(0,200000)

income + geom_histogram(color = "black", fill = "green")
income + geom_density(fill="green", position = "Stack")
ggplot(trainingSet,
  aes(x=cluster,y=Income,fill=cluster))+geom_boxplot(outlier.colour="black") +
  ylim(0,200000)
```

Il codice appena descritto ha il compito di analizzare gli elementi dei cluster in relazione al reddito di ciascun utente, ovvero in relazione alla variabile *income*. Le figure 46, 45 e 47 ne mostrano la rappresentazione dei grafici. Da esse è doveroso notare che nel primo cluster la maggior parte dei clienti possiede un reddito generalmente più basso rispetto ai *customers* facenti parte del secondo. Calcolando il reddito medio dei compratori, pari a 52247 dollari, si può anche notare che la maggior parte degli elementi nel primo cluster possiedono un il reddito inferiore alla media. E' opportuno notare anche che, la maggior parte degli utenti del secondo cluster, contrariamente ai primi, possiedono un reddito superiore alla media.

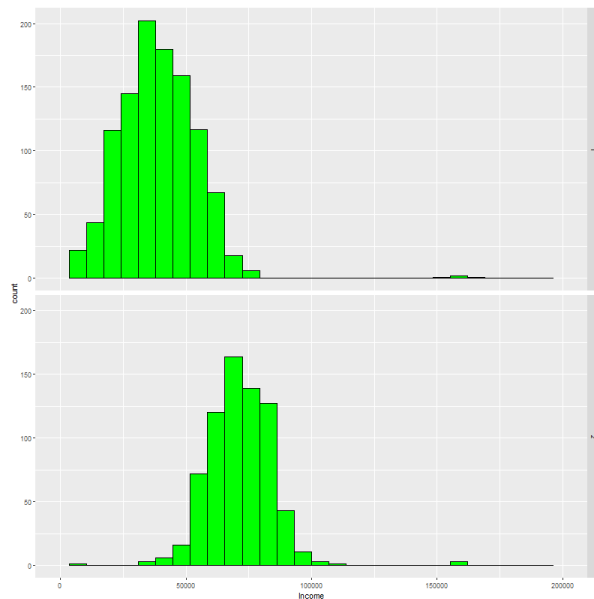


Figura 45: Istogramma della variabile Income in relazione al numero di cluster

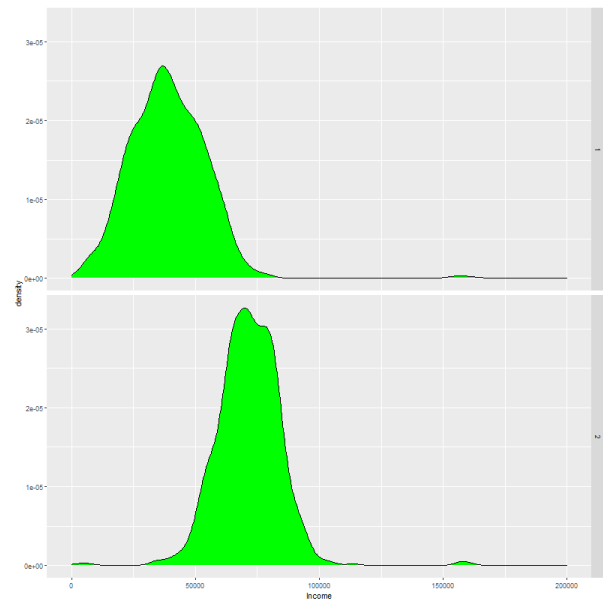


Figura 46: Diagramma di densità della variabile Income in relazione al numero di cluster

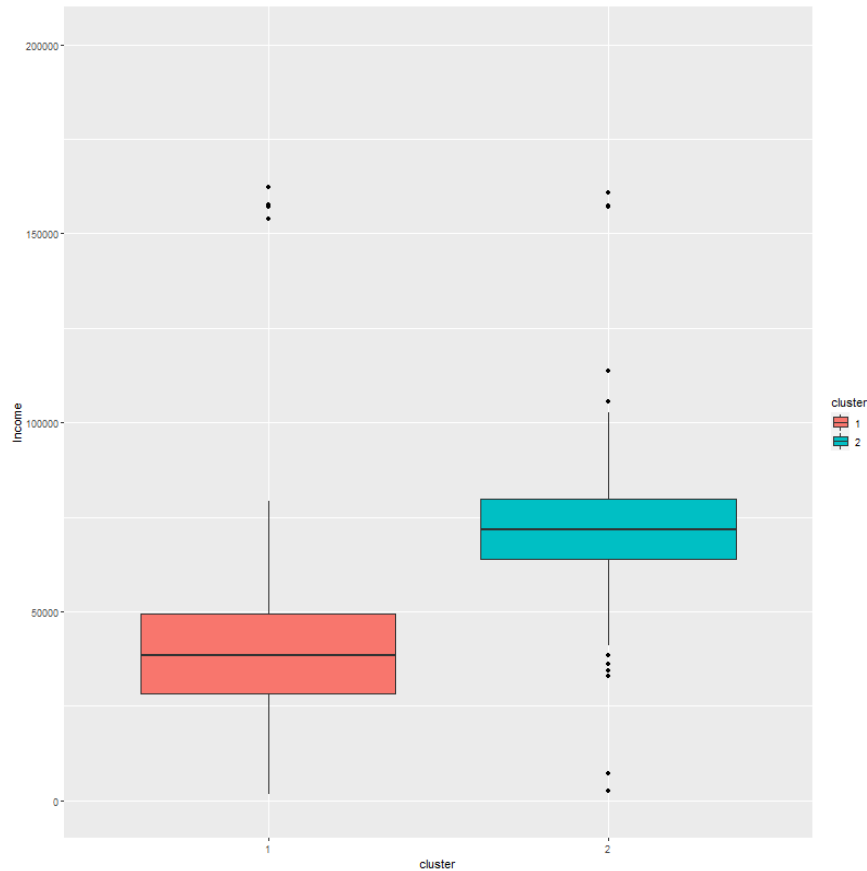


Figura 47: BoxPlot della variabile income in relazione al numero di cluster

Un'altra variabile analizzata è Total_spent, che spiega la gran parte della varianza della prima dimensione della PCA. Il codice sottostante spiega le figure 50, 49 e 48. In particolare, dall'analisi di esse è emerso che i compratori del primo cluster generalmente molto meno denaro rispetto a quelli del secondo. Oltre a ciò si può anche notare che quest'ultimi spendono mediamente più di mille dollari ogni due anni per prodotti come: vino, frutta, pesce e carne.

```
ts <- ggplot(trainingSet, aes(Total_spent), colour=cluster) + facet_grid(cluster~.)
ts + geom_histogram(color = "black", fill = "purple")
ts + geom_density(fill="purple", position = "Stack")
ggplot(trainingSet,
  aes(x=cluster,y=Total_spent,fill=cluster))+geom_boxplot(outlier.colour="black")
```

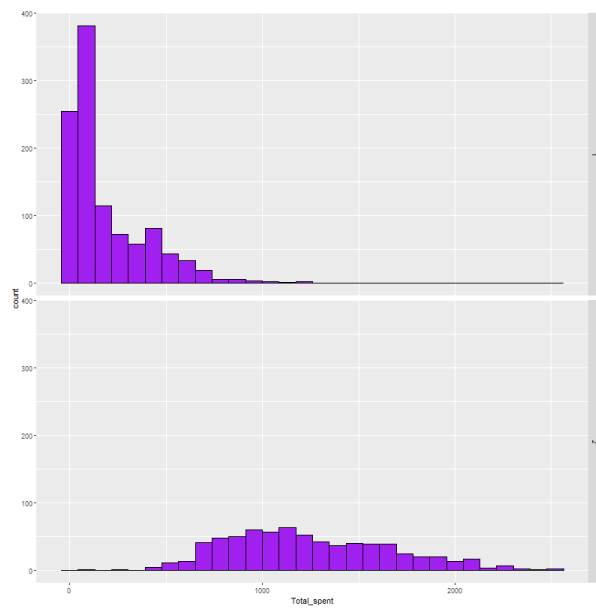


Figura 48: Istogramma della variabile Total_spent in relazione al numero di cluster

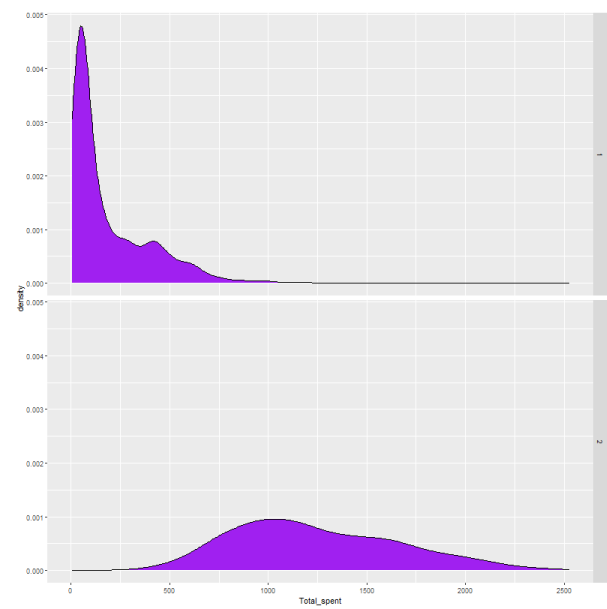


Figura 49: Diagramma di Densità della variabile Total_spent in relazione al numero di cluster

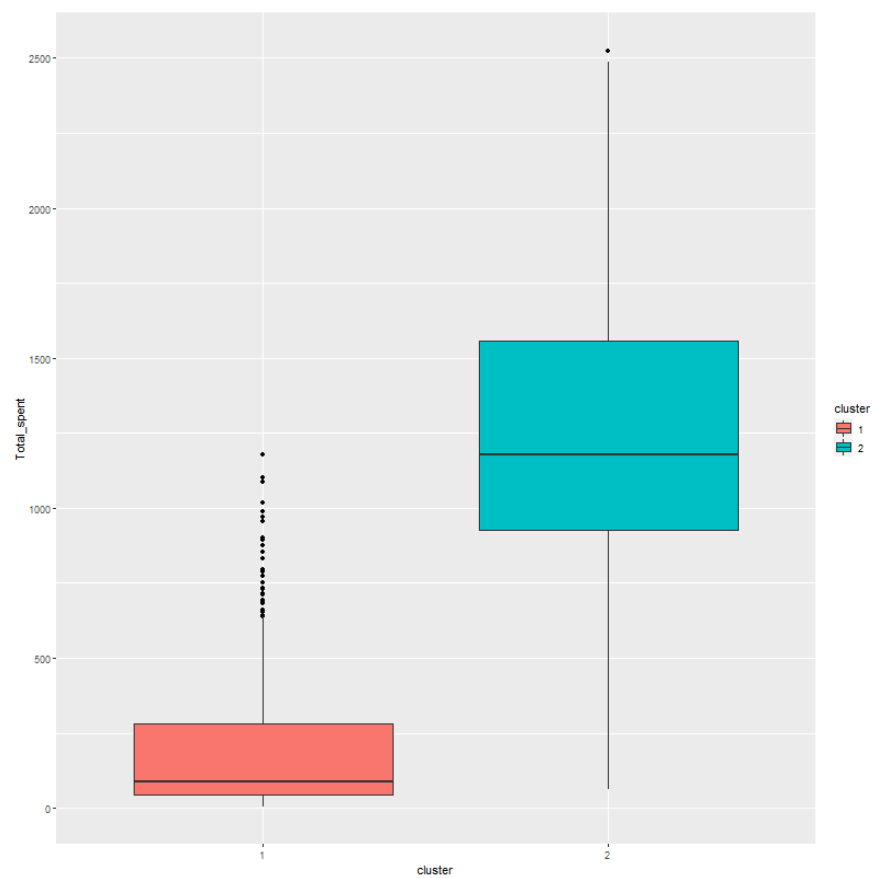


Figura 50: BoxPlot della variabile Total_spent in relazione al numero di cluster

Le figure 53, 52 e 51 mostrano che i clienti del primo cluster effettuino compere sul catalogo generalmente in quantità minore rispetto a quelli del secondo. Difatti ques'ultimi acquistano mediamente 5 prodotti dal catalogo.

```
numCatalogPurchases <- ggplot(trainingSet, aes(NumCatalogPurchases)) +
  facet_grid(cluster~.)
numCatalogPurchases + geom_histogram(color = "black", fill = "blue")
numCatalogPurchases + geom_density(fill="blue", position = "Stack")
ggplot(trainingSet,
  aes(x=cluster,y=NumCatalogPurchases,fill=cluster))+geom_boxplot(outlier.colour="black")
+ ylim(0,10)
```

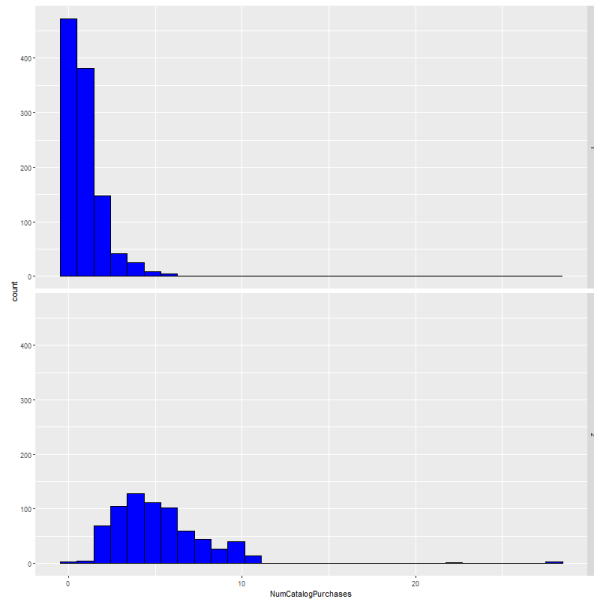


Figura 51: Istogramma della variabile NumCatalogPurchases in relazione al numero di cluster

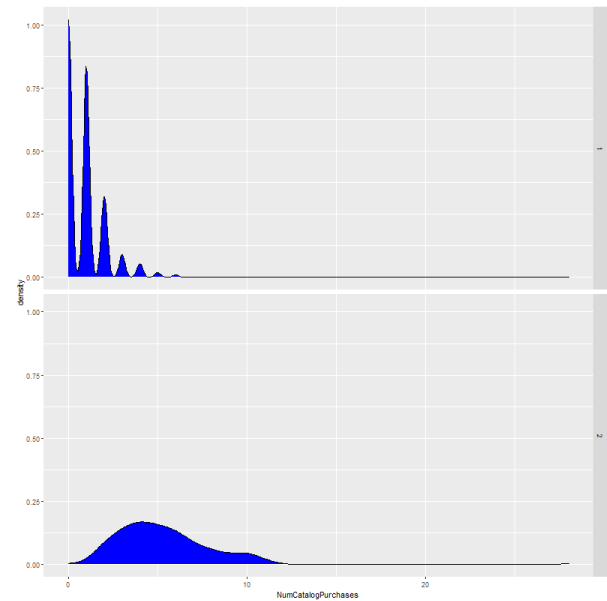


Figura 52: Diagramma di Densità della variabile NumCatalogPurchases in relazione al numero di cluster

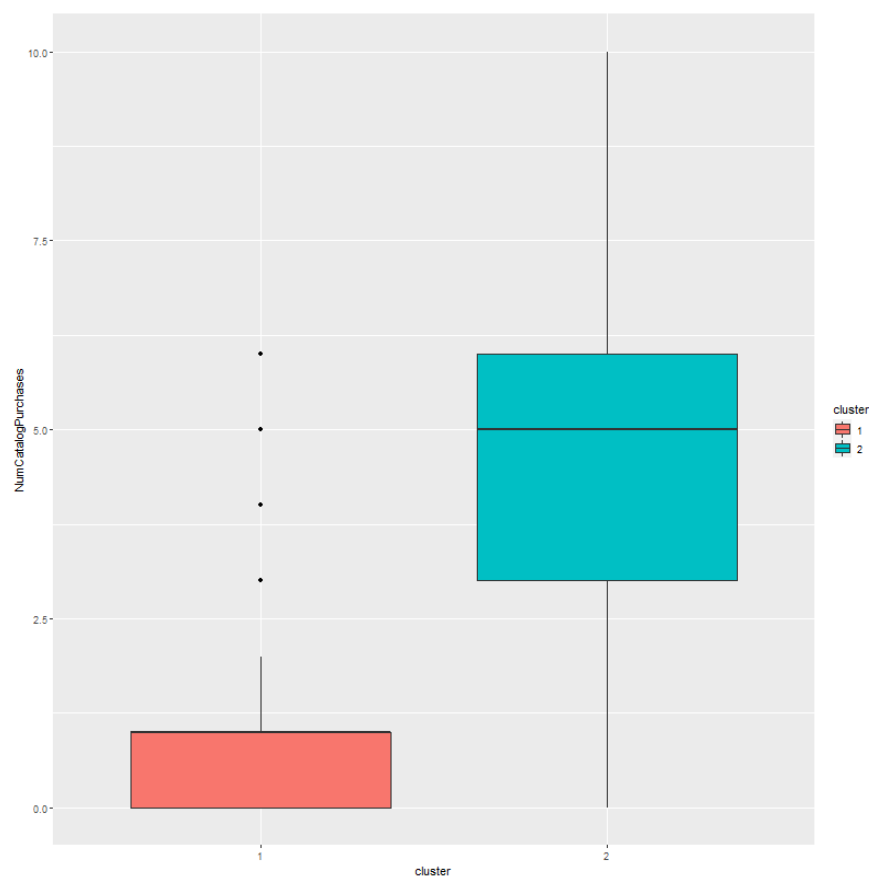


Figura 53: BoxPlot della variabile NumCatalogPurchases in relazione al numero di cluster

Dall'analisi della variabile *MntMeatProducts*, inerente all'acquisto di prodotti di carne negli ultimi due anni, si possono evincere alcune importanti informazioni correlate anche alla precedente analisi di *Total_spent*. In particolare le figure 56, 55 e 54 mostrano come gli acquirenti del primo cluster tendano a spendere generalmente di meno rispetto a quelli del secondo. Difatti i primi acquistano sicuramente molto meno rispetto alla media di circa 170 dollari ogni due anni spesi per prodotti di carne.

```
meat <- ggplot(trainingSet, aes(MntMeatProducts)) + facet_grid(cluster~.)
meat + geom_histogram(color = "black", fill = "brown")
meat + geom_density(fill="brown", position = "Stack")
ggplot(trainingSet,
  aes(x=cluster,y=MntMeatProducts,fill=cluster))+geom_boxplot(outlier.colour="black")
```

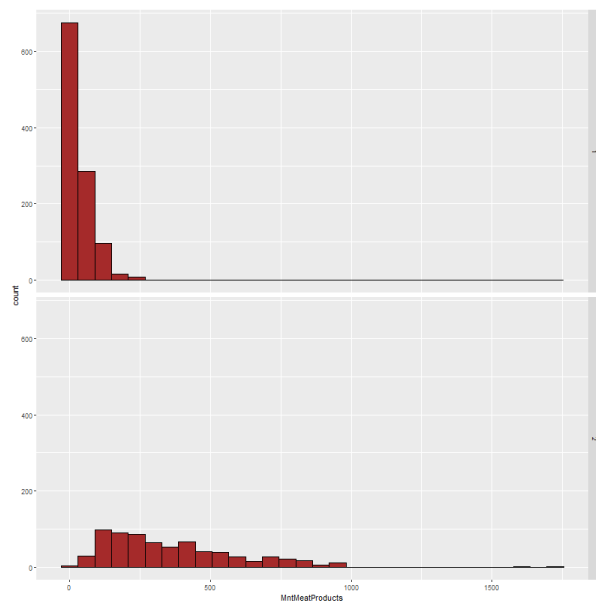


Figura 54: Istogramma della variabile MntMeatProducts in relazione al numero di cluster

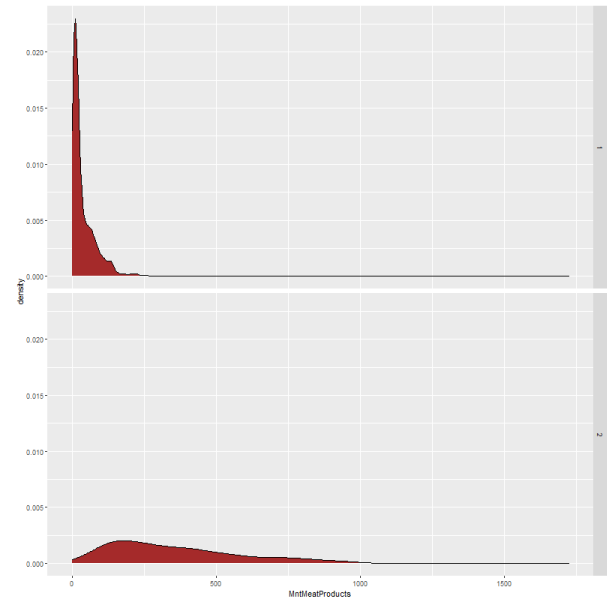


Figura 55: Diagramma di Densità della variabile MntMeatProducts in relazione al numero di cluster

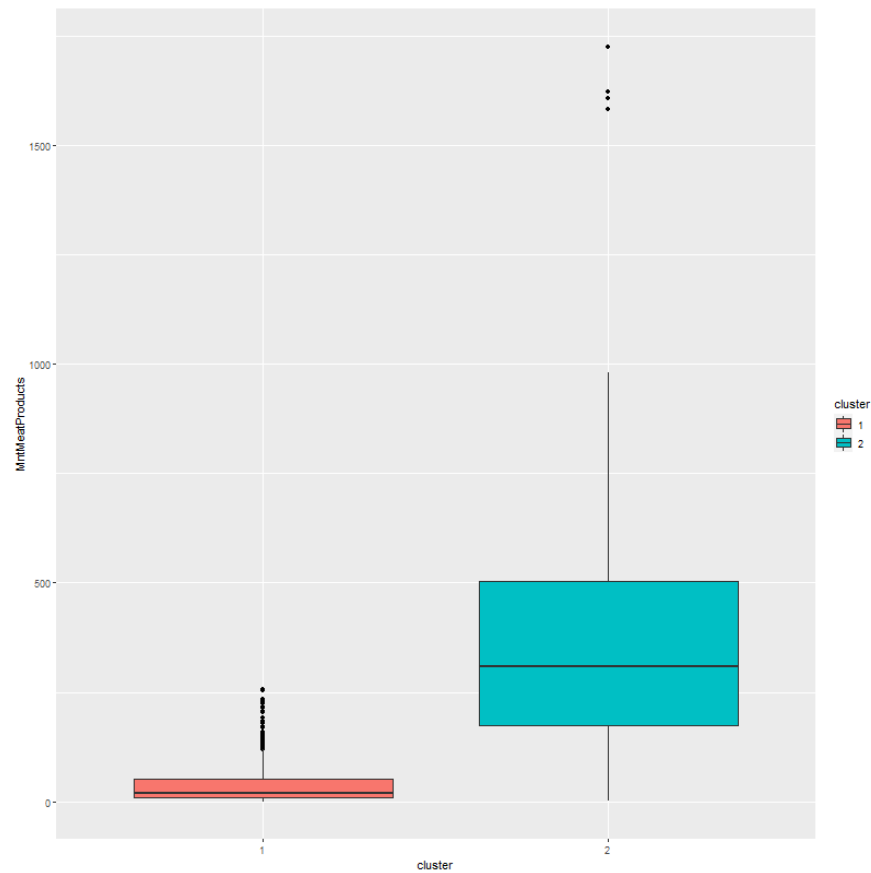


Figura 56: BoxPlot della variabile MntMeatProducts in relazione al numero di cluster

4.2 D-Tree

Usando il dataset ricavato dalla PCA ovvero *trainingSet_input* si è cercato di fare una previsione sul valore che assume la variabile *Response*. Per farlo è stata aggiunta una variabile al dataset.

```
trainingSet_input$Response<-trainingSet$Response
```

Successivamente è stata usata la funzione *rpart* per costruire un albero di decisione indicando come formula la variabile *Response*.

```
response.default.tree <- rpart(Response ~ ., data = trainingSet_input, method = "class")
```

Per visualizzarlo si è eseguita la funzione *prp* del *package*, *rpart.plot* indicando per i campi richiesti un parametro come per esempio *type* che assegna a tutti i nodi un etichetta, non solo le foglie. *extra*, che se uguale ad 1 visualizza il numero di osservazioni che cadono nel nodo.

```
prp(response.default.tree,  
  type = 1, extra = 1, varlen = -10,  
  box.col = ifelse(customer.default.tree$frame$var == "<leaf>", 'gray', 'white'))
```

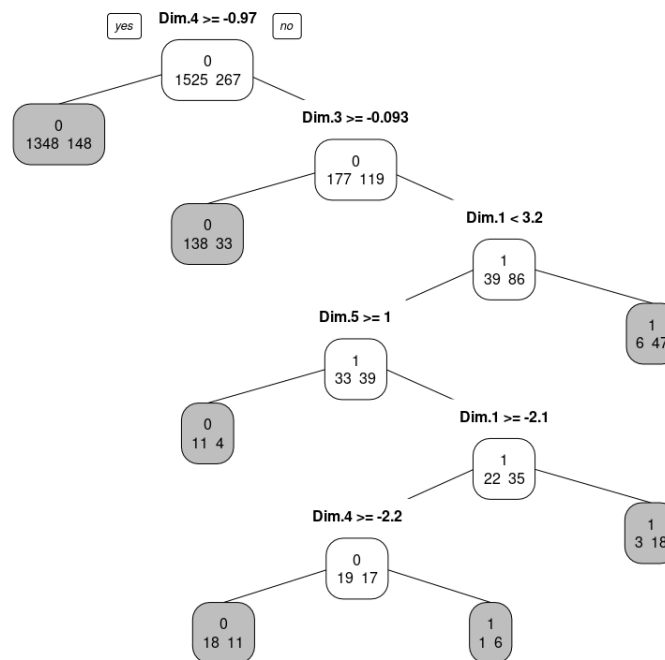


Figura 57: DecisionTree response.default.tree

Successivamente si è eseguita una previsione sulla variabile *Response* tramite la funzione *predict* passando come parametro anche l'albero *response.default.tree*.

```
response.default.tree.pred <- predict(response.default.tree, trainingSet_input, type =  
  "class")
```

Per confrontare il valore effettivo della variabile *Response* e della previsione si è fatta la matrice di confusione.

```
confusionMatrix.default<-confusionMatrix(response.default.tree.pred,
  as.factor(trainingSet_input$Response), positive = "1")
```

	Reference		
	-	0	1
Prediction	0	1515	196
	1	10	71

Dalla matrice di confusione sommando i valori presenti nella diagonale e facendo il rapporto con la somma di tutti i valori della matrice si ricava che l'accuratezza è del 88.5%. Si può ricavare anche la precisione che è pari al rapporto tra il numero di istanze classificate con il valore 1 e che corrisponde al valore che ha *Response* e il numero di istanze classificate con il valore 1 al 87.6%. Di seguito sono riportati altri dati statistici che è possibile ricavare.

Positive Class: 1
Sensitivity: 0.26592 Pos Pred Value: 0.87654
Specificity: 0.99344 Neg Pred Value: 0.88545

Per capire meglio l'albero decisionale si usa la funzione printcp().

```
cv.ct <- rpart( Response ~ ., data = trainingSet_input, method ="class", cp = 0,
  minsplit = 2, xval = 10)
printcp(cv.ct)
```

Variables actually used in tree construction: Dim.1 Dim.2 Dim.3 Dim.4 Dim.5
Root node error: 267/1792 = 0.149
n=1792

	CP	nsplit	rel error	xerror	xstd
1	0.0880150	0	0.00000	1.00000	0.056456
2	0.0131086	2	0.823970	0.84270	0.052535
3	0.0074906	6	0.771536	0.85393	0.052833
4	0.0056180	9	0.749064	0.91386	0.054375
5	0.0049938	32	0.554307	0.91386	0.054375
6	0.0037453	38	0.524345	0.91386	0.054375
7	0.0028090	98	0.299625	0.99625	0.056369
8	0.0024969	114	0.250936	1.04494	0.057483
9	0.0022472	123	0.228464	1.04494	0.057483
10	0.0018727	128	0.217228	1.11236	0.058955
11	0.0012484	202	0.052434	1.11236	0.058955
12	0.000000	205	0.048689	1.11985	0.059113

Succesivamente si è cercato il miglior albero di decisione prendendo prima il minimo tra una matrice di informazioni sulle strutture ottimali in base ad un parametro di complessità, cp dalla *lactable* che è pari a 0.8426966. Poi si è cercato l'errore standard corrispondente al minimo errore che è 0.05253461.

```
# min error
minerror <- min(cv.ct$cptable[ , 4])
minerror
```

```

minerrorstd <- cv.ct$cptable[cv.ct$cptable[,4] == minerror, 5]
minerrorstd

# set of trees where xerror is less than minerror + minerrorstd
simplertrees <- cv.ct$cptable[cv.ct$cptable[,4] < minerror + minerrorstd, ]
simplertrees

# cp of the simplest of those trees
bestcp <- simplertrees[1, 1]
bestcp

```

	CP	nsplit	rel error	xerror	xstd
2	0.013108614	2	0.8239700	0.8576779	0.05293184
3	0.007490637	6	0.7715356	0.8726592	0.05332376

Si ricava l'insieme di alberi in cui xerror è minore della somma tra minerror e minerrorstd e da esso si trova che il parametro di complessità, cp è 0.01310861.

```

bestcp <- simplertrees[1, 1]
bestcp

```

L'albero più semplice si può visualizzare eseguendo le istruzioni successive.

```

response.best.tree <- prune( cv.ct, cp = bestcp )
prp(response.best.tree, type = 1, extra = 1, varlen = -15, cex = 0.5,
     box.col = ifelse(response.best.tree$frame$var == "<leaf>", 'gray', 'white' ))

```

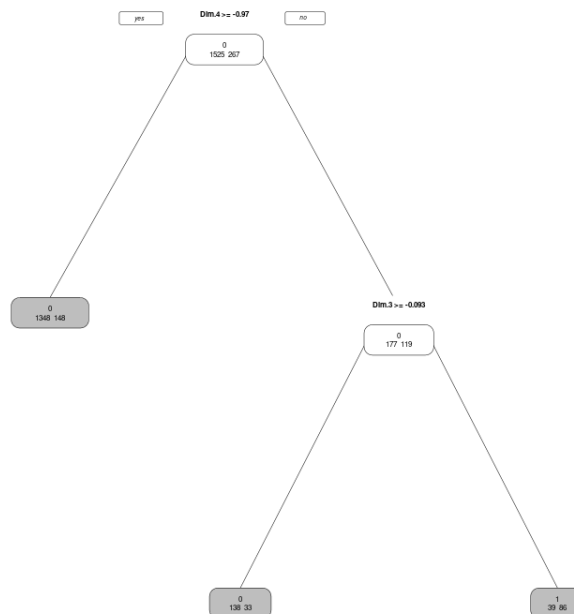


Figura 58: DecisionTree response.best.tree

Come per l'albero precedente prima di visualizzare la matrice di confusione ed altri dati statistici si è fatta una previsione.

```
response.best.tree.pred <- predict(response.best.tree, trainingSet_input, type = "class")
```

Ed in fine per confrontare il valore effettivo della variabile *Response* e della previsione si è fatta la matrice di confusione.

```
confusionMatrix2<-confusionMatrix(response.best.tree.pred,  
  as.factor(trainingSet_input$Response), positive = "1")
```

	Reference		
	-	0	1
Prediction	0	1486	181
	1	39	86

5 Esperimenti

TODO

5.1 Analisi dei risultati ottenuti

Conclusioni

Questa sperimentazione ha avuto l'obiettivo di analizzare un insieme di dati mediante tecniche di machine learning differenti, in particolare mediante l'algoritmo **K-Means**. Dall'analisi dei dati riscontrati si può giungere alla conclusione che una buona suddivisione dei dati riportati può avvenire mediante l'utilizzo di due cluster. In particolare il primo cluster presenta clienti con un reddito generalmente al di sotto della media e sicuramente minore rispetto alla maggior parte dei compratori facenti parte della seconda divisione. Secondo i dati analizzati ciò ha comportato sicuramente una riduzione delle spese totali da parte dei primi. La riduzione del numero di acquisti di prodotti ha generalmente toccato elementi che variano dai vini fino alla carne.