

Lab Assignment Ten

Noah Gallego

2024-10-11

Lab Assignment 10: PCA Analysis on the MTCars Dataset

Import Libraries

```
library(dplyr)      # For data manipulation

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

library(ggplot2)    # For plotting
library(ggrepel)     # For Last Plot
library(factoextra) # For Scree Plot

## Welcome! Want to learn more? See two factoextra-related books at
https://goo.gl/ve3WBa
```

Exploring the Dataset

```
data("mtcars")
df = mtcars
# Remove Unessecary Columns
df = df %>% select(c("mpg", "disp", "hp", "drat", "wt", "qsec"))

# Display First Few Columns
head(df)
```

	mpg	disp	hp	drat	wt	qsec
## Mazda RX4	21.0	160	110	3.90	2.620	16.46
## Mazda RX4 Wag	21.0	160	110	3.90	2.875	17.02
## Datsun 710	22.8	108	93	3.85	2.320	18.61
## Hornet 4 Drive	21.4	258	110	3.08	3.215	19.44
## Hornet Sportabout	18.7	360	175	3.15	3.440	17.02
## Valiant	18.1	225	105	2.76	3.460	20.22

```
# Scale Dataset for PCA
```

```
df_scaled = scale(df)
```

```
apply(df_scaled, 2, sd)
```

```
## mpg disp hp drat wt qsec
```

```
## 1 1 1 1 1 1
```

Perform PCA

```
# Perform PCA
```

```
pca = prcomp(df_scaled)
```

```
summary(pca)
```

```
## Importance of components:
```

```
## PC1 PC2 PC3 PC4 PC5 PC6
```

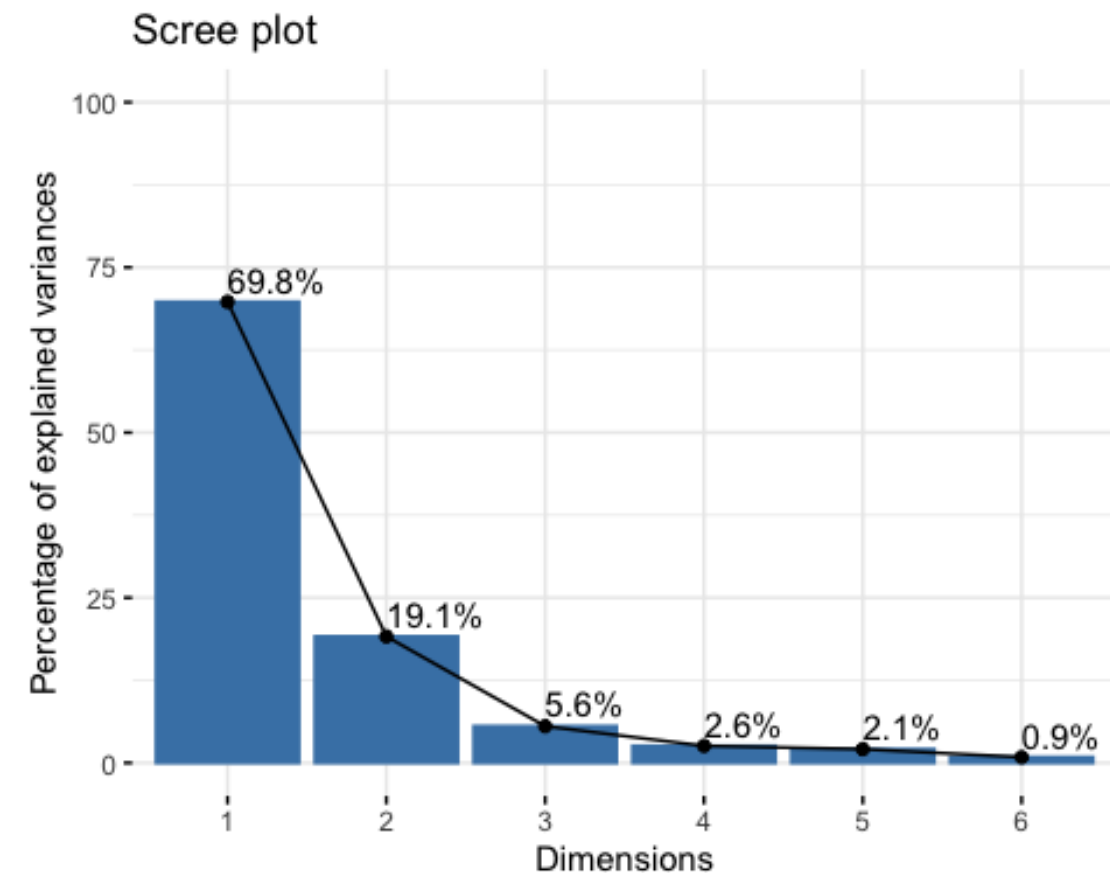
```
## Standard deviation 2.0463 1.0715 0.57737 0.39289 0.3533 0.22799
```

```
## Proportion of Variance 0.6979 0.1913 0.05556 0.02573 0.0208 0.00866
```

```
## Cumulative Proportion 0.6979 0.8892 0.94481 0.97054 0.9913 1.00000
```

```
# Display Proportion of Variance
```

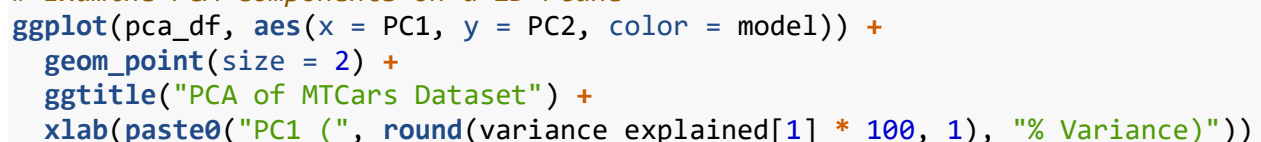
```
fviz_eig(pca, addlabels = TRUE, ylim = c(0, 100)) # Scree plot showing  
variance explained by components
```



```
# Examine Loadings
```

```
pca$rotation
```

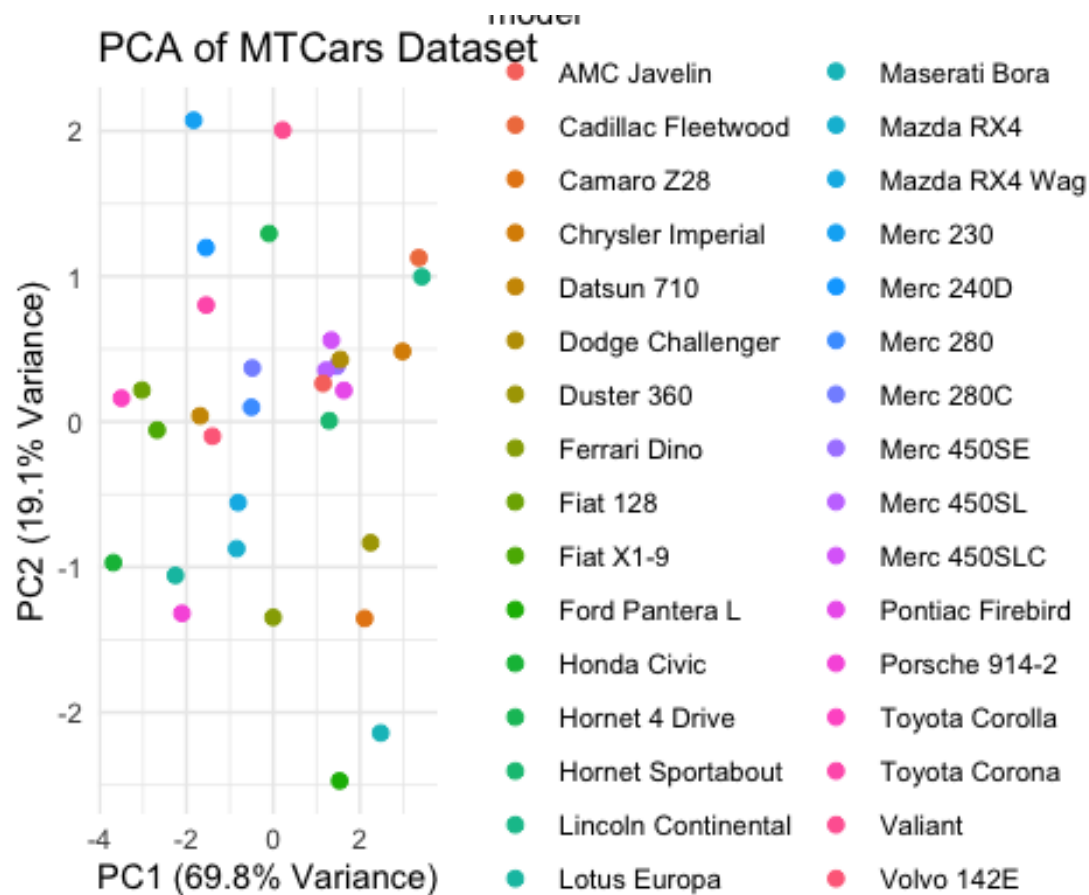
```
biplot(pca, scale = 0, cex = c(0.3, 1))
```



```

+ ylab(paste0("PC2 (", round(variance_explained[2] * 100, 1), "% Variance)"))
+ theme_minimal()

```



Outlier Detection

```

pca_distances = sqrt(rowSums(pca$x[, 1:2]^2))
outlier_threshold = quantile(pca_distances, 0.95)
outliers = pca_distances > outlier_threshold

pca_df$outliers = outliers

# Plot the PCA with outliers labeled
ggplot(pca_df, aes(x = PC1, y = PC2, color = outliers)) +
  geom_point(size = 2) +
  geom_text_repel(data = subset(pca_df, outliers), aes(label = model),
    size = 3, box.padding = 0.5, max.overlaps = Inf) +
  labs(title = "Outliers Detected Using PCA")

```

Outliers Detected Using PCA

