

Homework 2 Complete

Problem 1 - Mercedes S-Class

```
library(class)
library(FNN)
```

```
##
## Attaching package: 'FNN'

## The following objects are masked from 'package:class':
##
##   knn, knn.cv
```

```
library(tidyverse)
```

```
## -- Attaching packages ----- tidyverse 1.3.0 --
```

```
## v ggplot2 3.2.1    v purrr   0.3.3
## v tibble  2.1.3    v dplyr  0.8.3
## v tidyr   1.0.0    v stringr 1.4.0
## v readr   1.3.1    v forcats 0.4.0
```

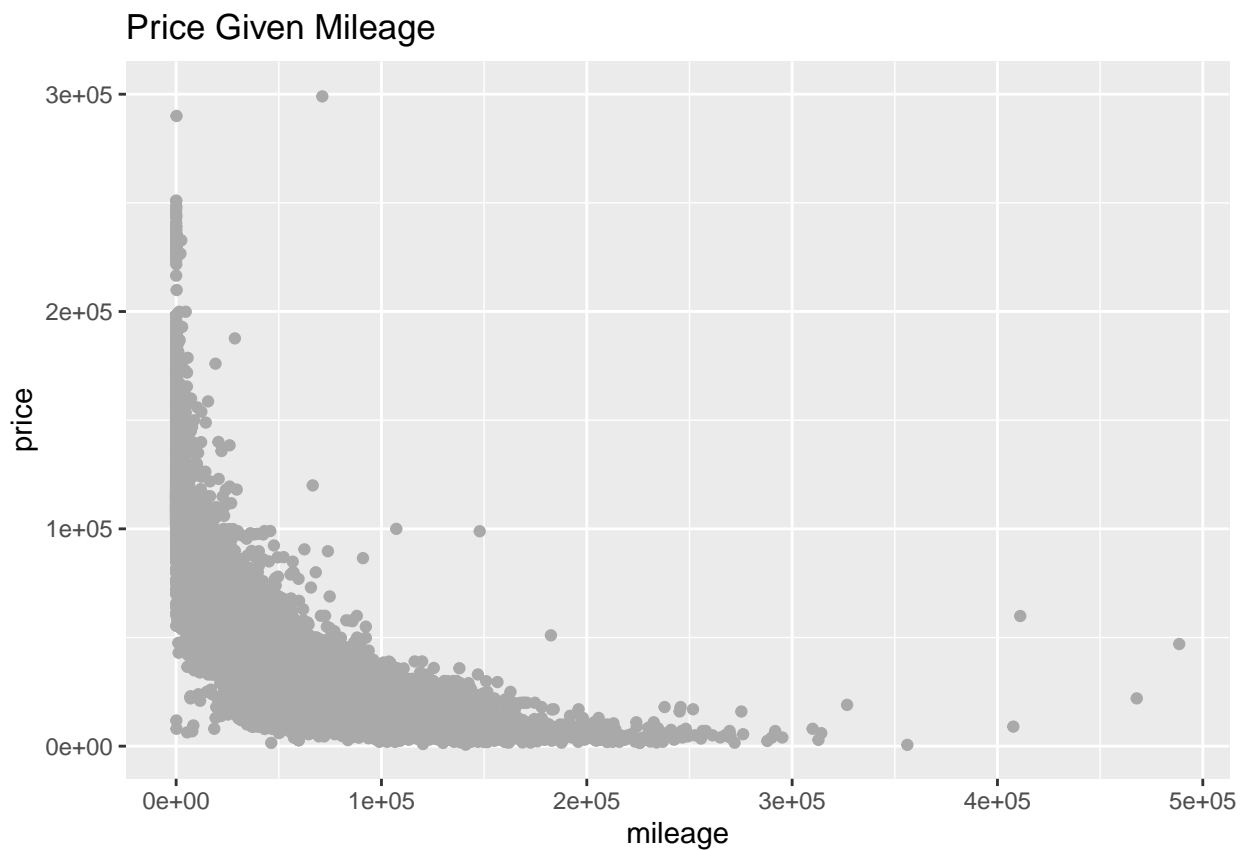
```
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
```

```
#Plot original data from sclass data to show relationship between the two variables before separating b
sclass <- read_csv("SDS323/data/sclass.csv")
```

```
## Parsed with column specification:
## cols(
##   id = col_double(),
##   trim = col_character(),
##   subTrim = col_character(),
##   condition = col_character(),
##   isOneOwner = col_logical(),
##   mileage = col_double(),
##   year = col_double(),
##   color = col_character(),
##   displacement = col_character(),
##   fuel = col_character(),
##   state = col_character(),
##   region = col_character(),
```

```
##   soundSystem = col_character(),
##   wheelType = col_character(),
##   wheelSize = col_character(),
##   featureCount = col_double(),
##   price = col_double()
## )
```

```
ggplot(data = sclass) +
  geom_point(mapping = aes(x = mileage, y = price), color='darkgrey') +
  ylim(0, 300000)+
  ggtitle("Price Given Mileage")
```



```
summary(sclass)
```

```
##           id           trim           subTrim           condition
##  Min.      :    2   Length:29466   Length:29466   Length:29466
##  1st Qu.:13231   Class :character   Class :character   Class :character
##  Median :26254   Mode  :character   Mode  :character   Mode  :character
##  Mean    :26269
##  3rd Qu.:39293
##  Max.    :52572
##  isOneOwner      mileage           year           color
##  Mode :logical   Min.      :    1   Min.      :1988   Length:29466
##  FALSE:25340     1st Qu.:   14   1st Qu.:2007   Class :character
##  TRUE :4126       Median : 26120   Median :2012   Mode  :character
```

```
##           Mean    : 40387    Mean    :2010
##           3rd Qu.: 68234    3rd Qu.:2015
##           Max.    :488525    Max.    :2015
## displacement      fuel          state          region
## Length:29466      Length:29466      Length:29466      Length:29466
## Class :character  Class :character  Class :character  Class :character
## Mode  :character  Mode  :character  Mode  :character  Mode  :character
##
##
##
## soundSystem      wheelType      wheelSize      featureCount
## Length:29466      Length:29466      Length:29466      Min.    : 0.00
## Class :character  Class :character  Class :character  1st Qu.: 18.00
## Mode  :character  Mode  :character  Mode  :character  Median : 53.00
##                                     Mean    : 46.48
##                                     3rd Qu.: 70.00
##                                     Max.    :132.00
##
## price
## Min.    : 599
## 1st Qu.: 28995
## Median : 56991
## Mean    : 67001
## 3rd Qu.:108815
## Max.    :299000
```

#The three variables we will be focusing on are the trim, mileage, and price. Mileage ranges from 1 to 100000.

*#In order to build a predictive model for price given mileage, we will separate the data based on trim.
#Below is the separation of 350 AMG:*

```
sclass350 = subset(sclass, trim == '350')
dim(sclass350)
```

```
## [1] 416 17
```

```
summary(sclass350)
```

```
##           id           trim           subTrim           condition
## Min.    : 282    Length:416    Length:416    Length:416
## 1st Qu.:14290    Class :character  Class :character  Class :character
## Median :26658    Mode  :character  Mode  :character  Mode  :character
## Mean    :26520
## 3rd Qu.:39599
## Max.    :52220
## isOneOwner      mileage          year          color
## Mode :logical    Min.    : 6    Min.    :1994    Length:416
## FALSE:310        1st Qu.: 19264  1st Qu.:2006    Class :character
## TRUE :106        Median : 29998  Median :2012    Mode  :character
##                                     Mean    : 42926  Mean    :2010
##                                     3rd Qu.: 63479  3rd Qu.:2012
##                                     Max.    :173000  Max.    :2013
## displacement      fuel          state          region
## Length:416        Length:416      Length:416      Length:416
```

```
## Class :character Class :character Class :character Class :character
## Mode :character Mode :character Mode :character Mode :character
##
##
##
## soundSystem wheelType wheelSize featureCount
## Length:416 Length:416 Length:416 Min. : 0.00
## Class :character Class :character Class :character 1st Qu.: 31.75
## Mode :character Mode :character Mode :character Median : 54.00
## Mean : 49.22
## 3rd Qu.: 70.00
## Max. :112.00
##
## price
## Min. : 6600
## 1st Qu.: 19401
## Median : 52900
## Mean : 46854
## 3rd Qu.: 61991
## Max. :106010
```

```
na.omit(sclass350)
```

```
## # A tibble: 416 x 17
## id trim subTrim condition isOneOwner mileage year color displacement
## <dbl> <chr> <chr> <chr> <lgl> <dbl> <dbl> <chr> <chr>
## 1 282 350 unsp CP0 FALSE 21929 2012 Black 3.0 L
## 2 284 350 unsp CP0 FALSE 17770 2012 Silv~ 3.0 L
## 3 285 350 unsp Used FALSE 29108 2012 Black 3.0 L
## 4 288 350 unsp CP0 FALSE 35004 2013 White 3.0 L
## 5 289 350 unsp Used TRUE 66689 2012 Black 3.0 L
## 6 290 350 unsp CP0 FALSE 19567 2012 Black 3.0 L
## 7 949 350 unsp CP0 TRUE 10616 2012 Black 3.0 L
## 8 950 350 unsp CP0 FALSE 2578 2013 Black 3.0 L
## 9 951 350 unsp CP0 FALSE 23677 2012 Black 3.0 L
## 10 952 350 unsp CP0 TRUE 28384 2012 Black 3.0 L
## # ... with 406 more rows, and 8 more variables: fuel <chr>, state <chr>,
## # region <chr>, soundSystem <chr>, wheelType <chr>, wheelSize <chr>,
## # featureCount <dbl>, price <dbl>
```

```
#Below is the separation of 65 AMG:
sclass65AMG = subset(sclass, trim == '65 AMG')
summary(sclass65AMG)
```

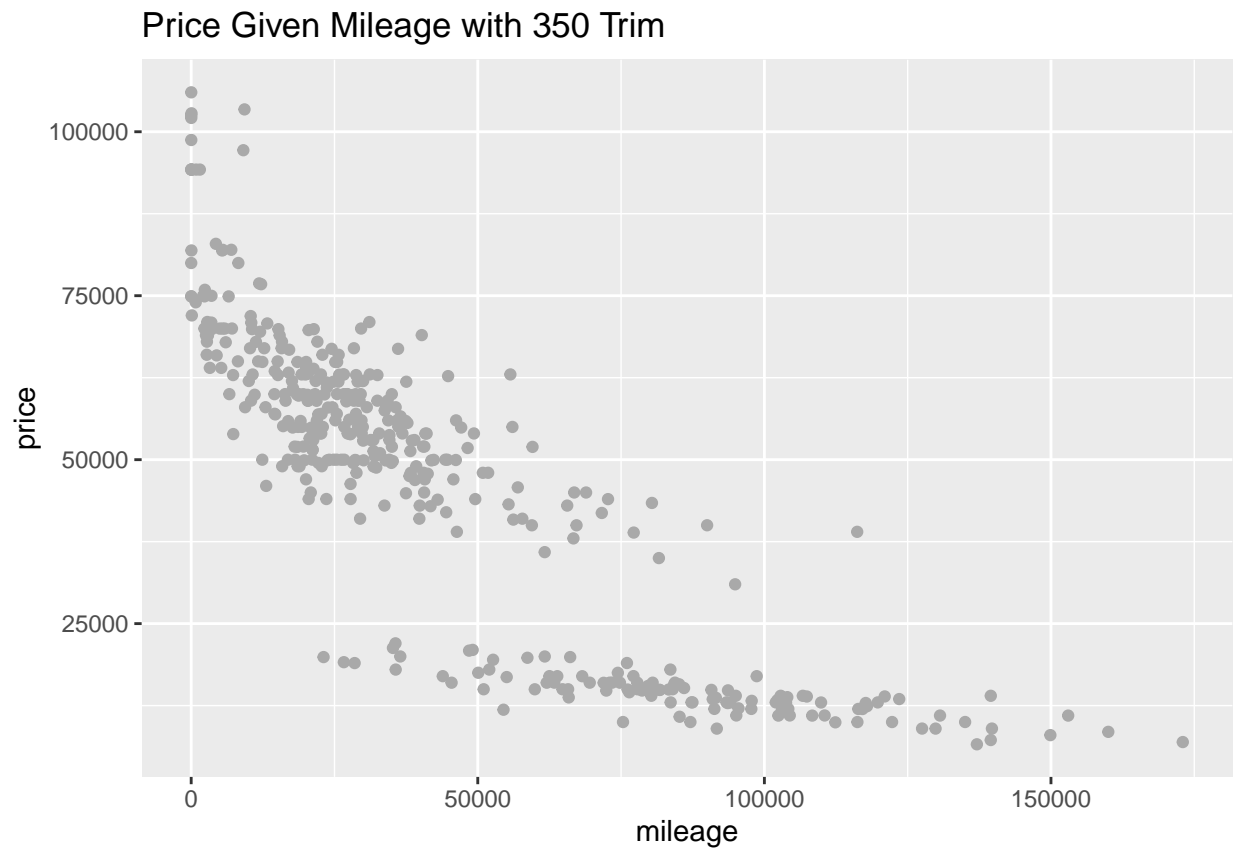
```
## id trim subTrim condition
## Min. : 1060 Length:292 Length:292 Length:292
## 1st Qu.:13977 Class :character Class :character Class :character
## Median :26557 Mode :character Mode :character Mode :character
## Mean :26444
## 3rd Qu.:38687
## Max. :52326
## isOneOwner mileage year color
## Mode :logical Min. : 1 Min. :2006 Length:292
## FALSE:254 1st Qu.: 20 1st Qu.:2007 Class :character
```

```
## TRUE :38          Median : 28803   Median :2010   Mode  :character
##                  Mean   : 33700   Mean   :2010
##                  3rd Qu.: 58496   3rd Qu.:2015
##                  Max.    :146975   Max.    :2015
## displacement      fuel            state            region
## Length:292        Length:292        Length:292        Length:292
## Class :character   Class :character   Class :character   Class :character
## Mode  :character   Mode  :character   Mode  :character   Mode  :character
##
##
## soundSystem        wheelType        wheelSize        featureCount
## Length:292        Length:292        Length:292        Min.   : 0.00
## Class :character   Class :character   Class :character   1st Qu.: 17.00
## Mode  :character   Mode  :character   Mode  :character   Median : 58.00
##                                     Mean   : 48.09
##                                     3rd Qu.: 72.00
##                                     Max.   :112.00
## price
## Min.   : 18990
## 1st Qu.: 48711
## Median : 79995
## Mean   :117121
## 3rd Qu.:225975
## Max.   :247075
```

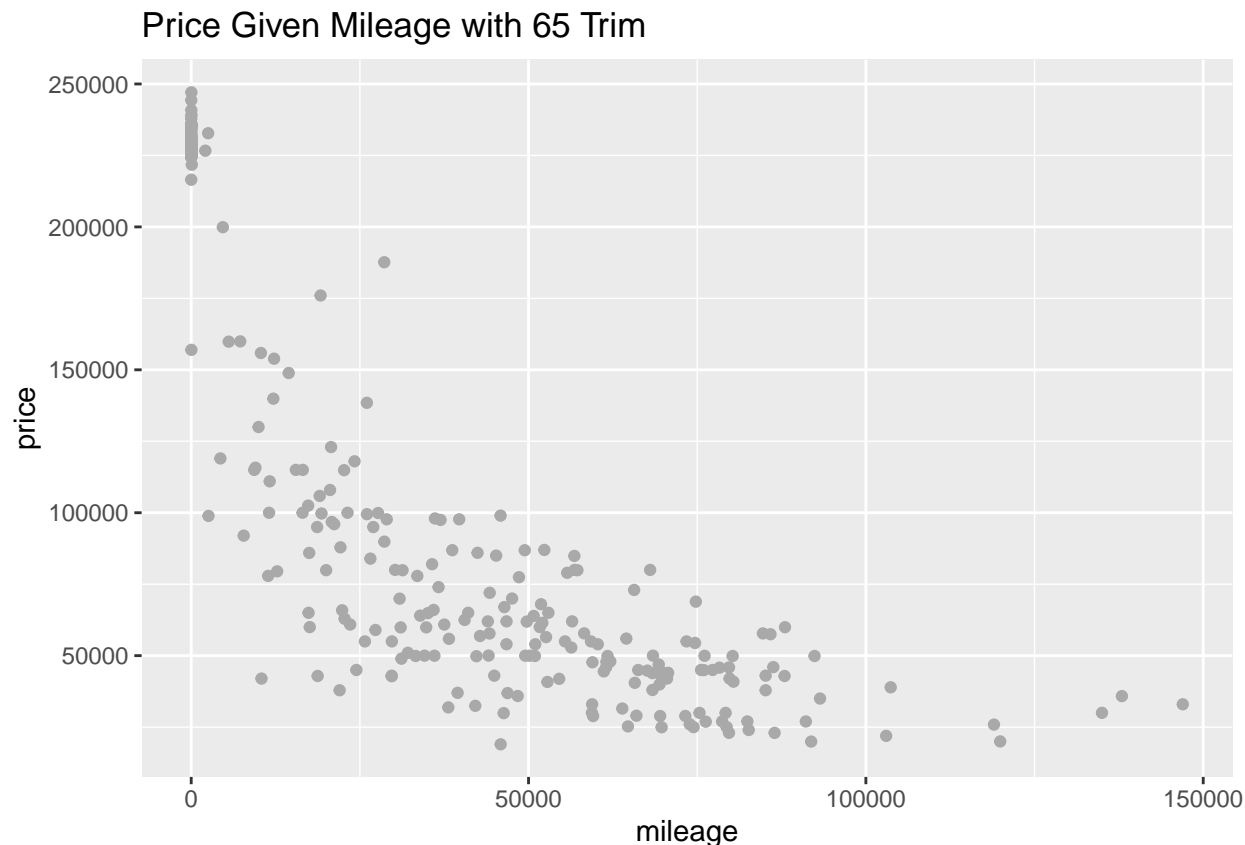
```
dim(sclass65AMG)
```

```
## [1] 292 17
```

```
#Before doing the train-test split, here are two plots that show the data for each trim level. You can
ggplot(data = sclass350) +
  geom_point(mapping = aes(x = mileage, y = price), color='darkgrey')+
  ggtitle("Price Given Mileage with 350 Trim")
```



```
ggplot(data = sclass65AMG) +  
  geom_point(mapping = aes(x = mileage, y = price), color='darkgrey') +  
  ggtitle("Price Given Mileage with 65 Trim")
```



#According to these plots, there seems to be more of an even scatter with 65 AMG and more of a break be

#Train-test split for 350 AMG trim level: We are looking for a model that gives the best explanation wi

```
N = nrow(sclass350)
N_train = floor(0.8*N)
N_test = N - N_train

train_ind = sort(sample.int(N, N_train, replace=FALSE))
```

```
D_train = sclass350[train_ind,]
D_train = arrange(D_train, mileage)
D_test = sclass350[-train_ind,]
```

```
y_train = D_train$price
X_train = data.frame(mileage=jitter(D_train$mileage))
X_test = data.frame(mileage=jitter(D_test$mileage))
y_test = D_test$price
```

#Running KNN test for various values of K to determine the optimal predictive model. This will start at

#KNN = 2

#The first step is to make predictions to the data frame.

```
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn2 = knn.reg(train = X_train, test = X_test, y = y_train, k=2)
```

#rmse calculation

```
rmse = function(y, ypred) {  
  sqrt(mean(data.matrix((y-ypred)^2)))  
}
```

```
ypred_lm1 = predict(lm1, X_test)  
ypred_lm2 = predict(lm2, X_test)  
ypred_knn2 = knn2$pred
```

```
rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

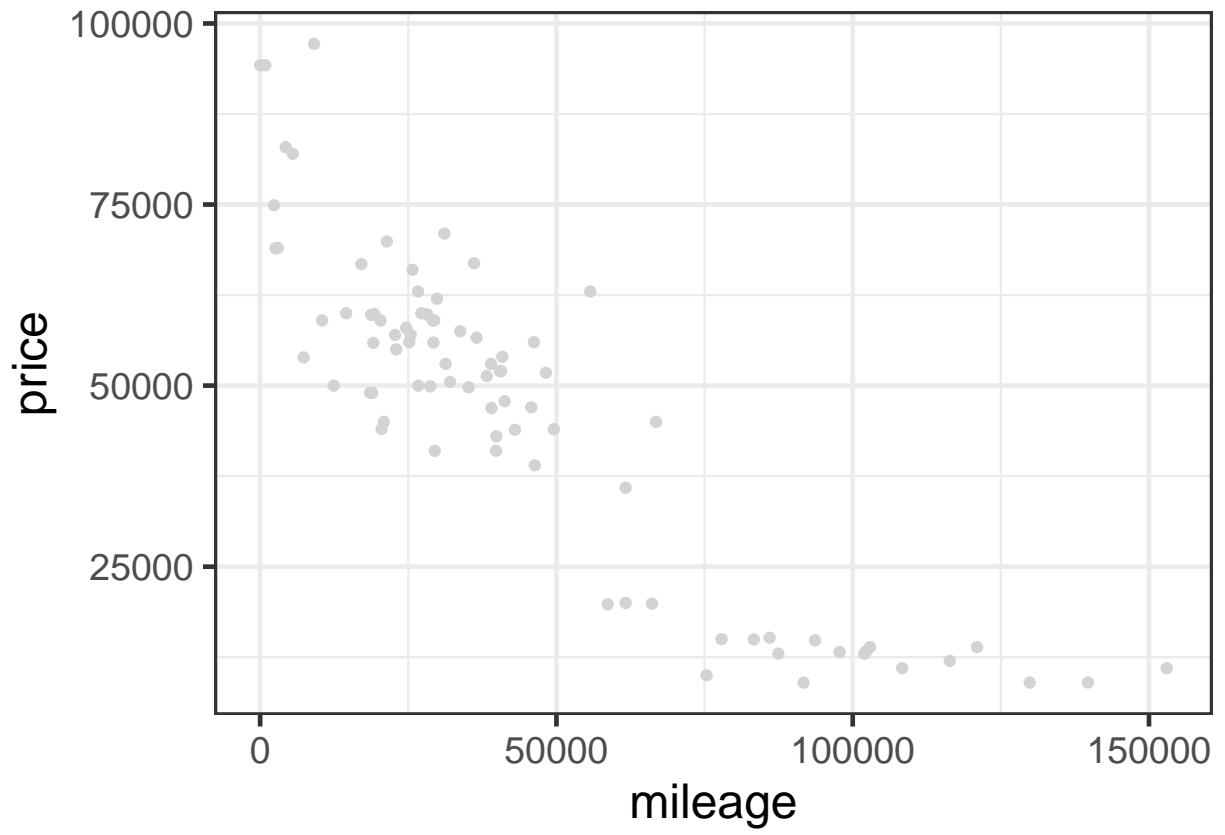
```
rmse(y_test, ypred_knn2)
```

```
## [1] 10582.86
```

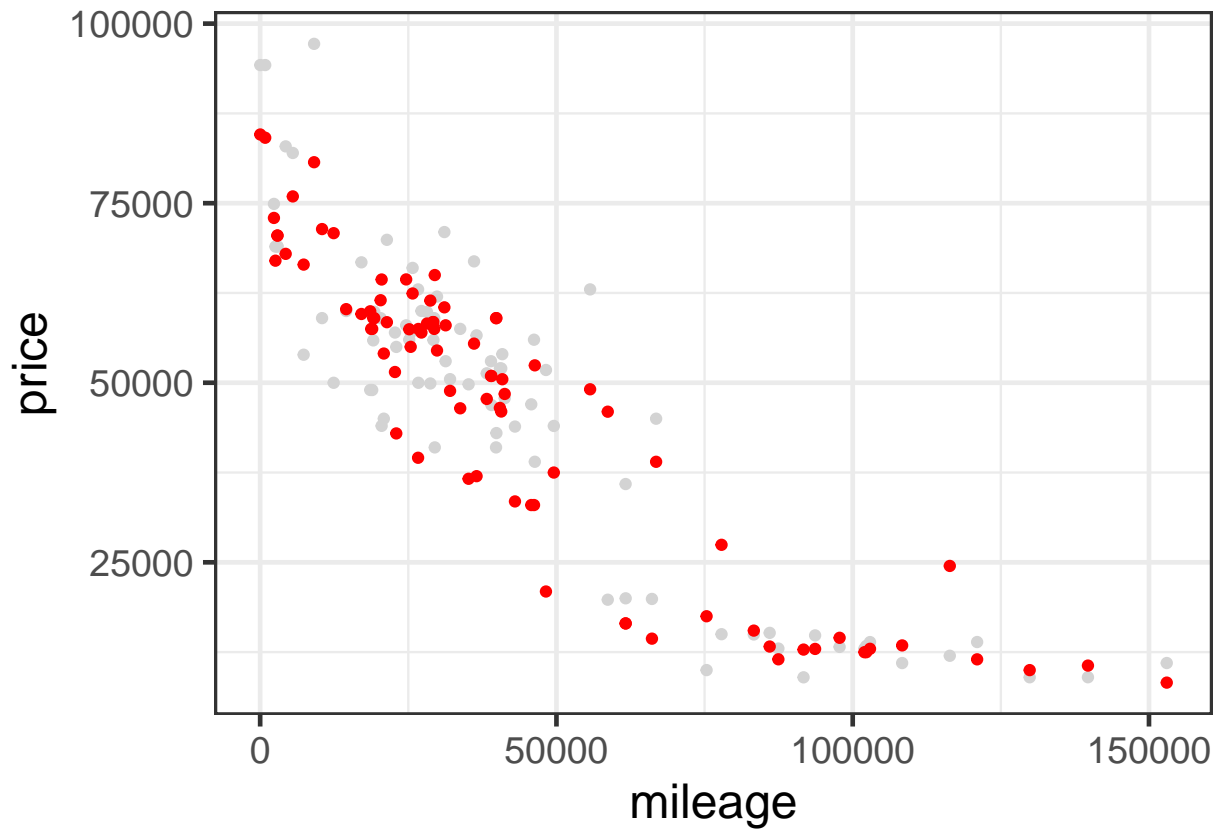
```
#attach predictions to data frame
```

```
D_test$ypred_lm2 = ypred_lm2  
D_test$ypred_knn2 = ypred_knn2
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

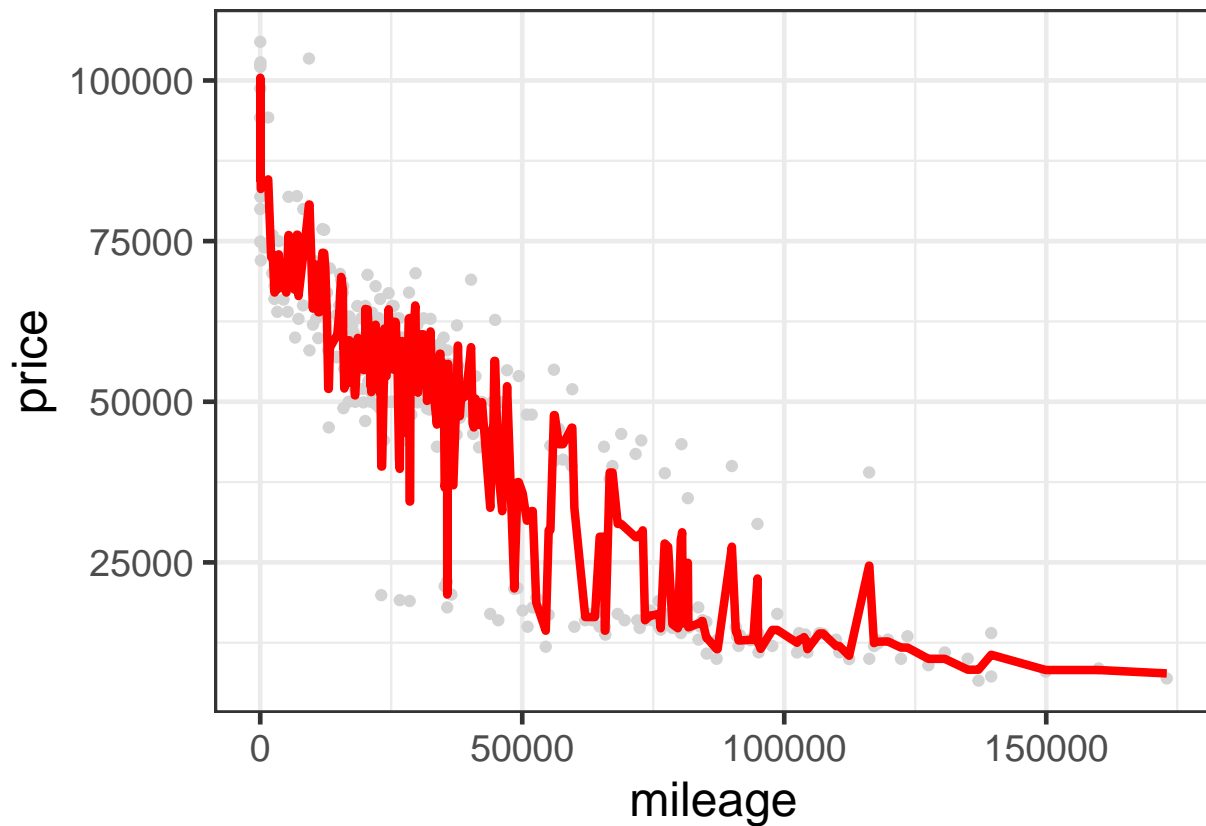



```
p_test + geom_point(aes(x = mileage, y = ypred_knn2), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 2)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 3
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn3 = knn.reg(train = X_train, test = X_test, y = y_train, k=3)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn3 = knn3$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn3)
```

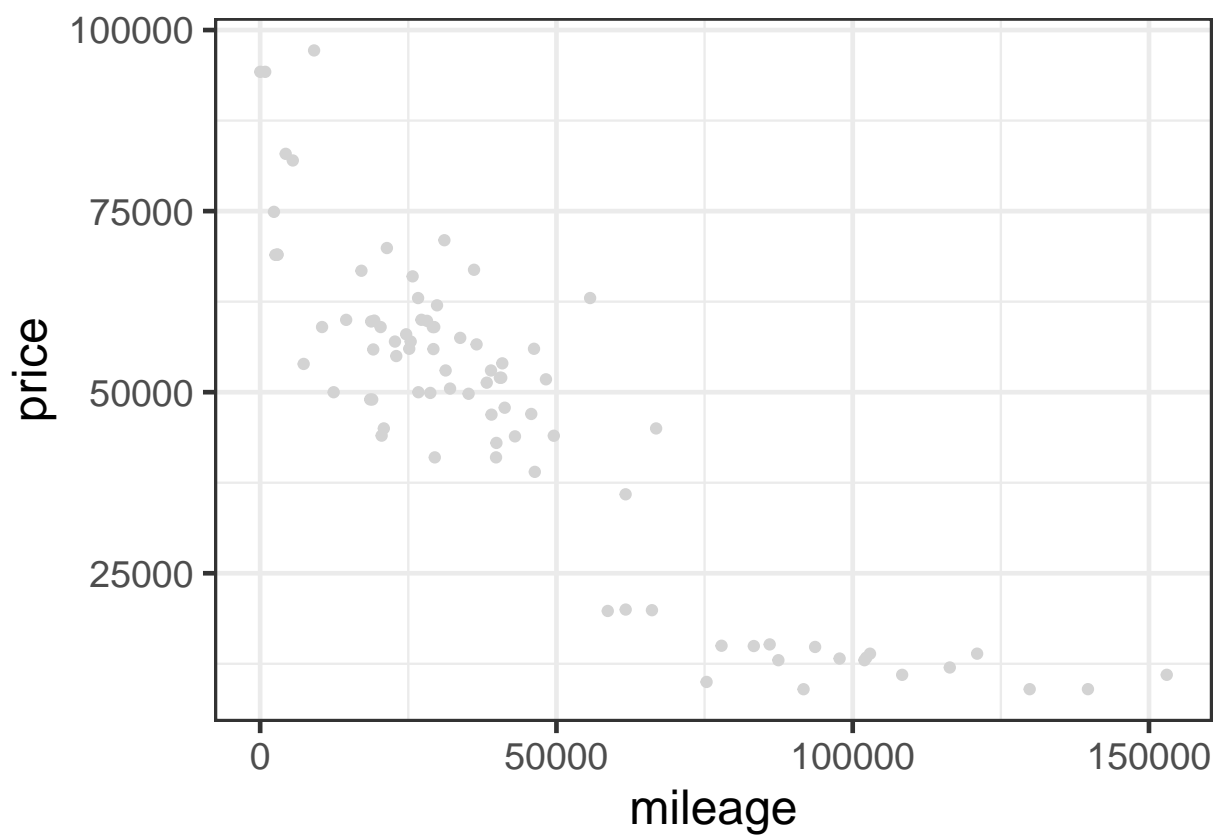
```
## [1] 9185.008
```

```
#attach predictions to data frame
```

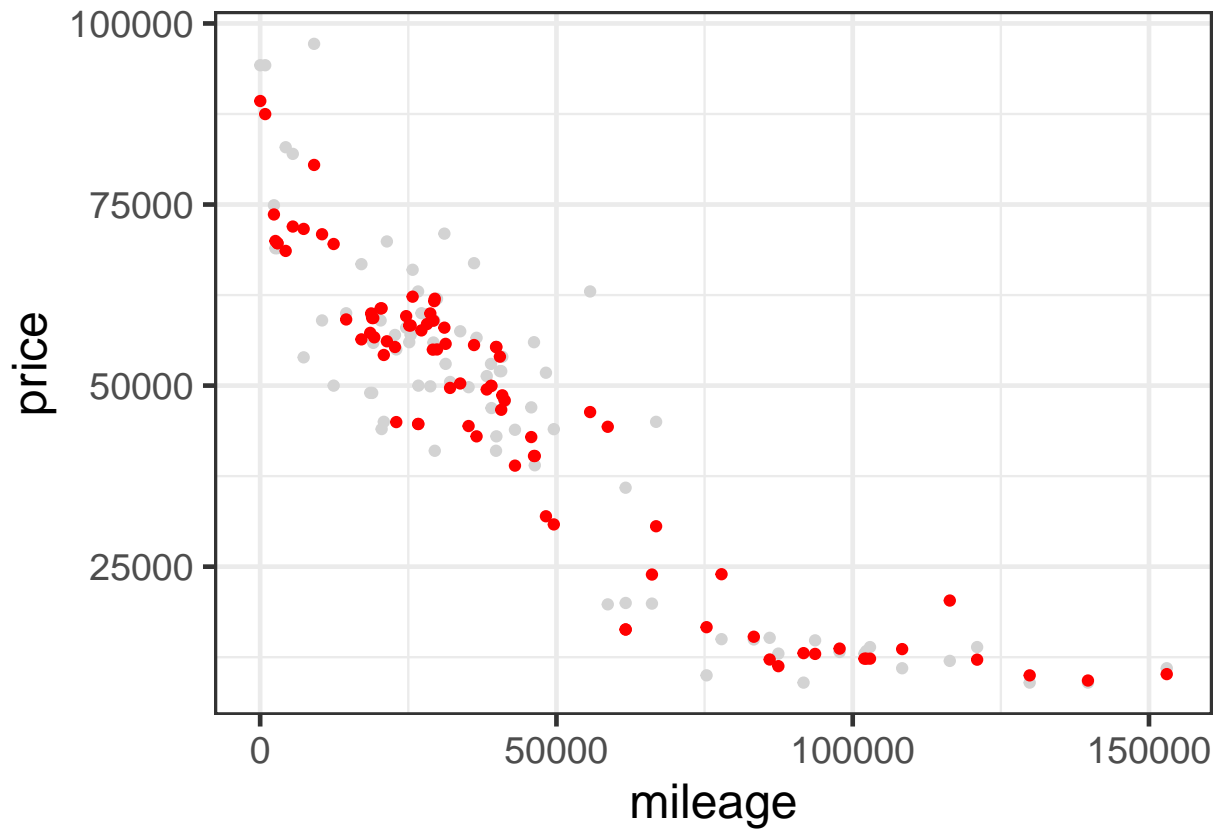
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn3 = ypred_knn3
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

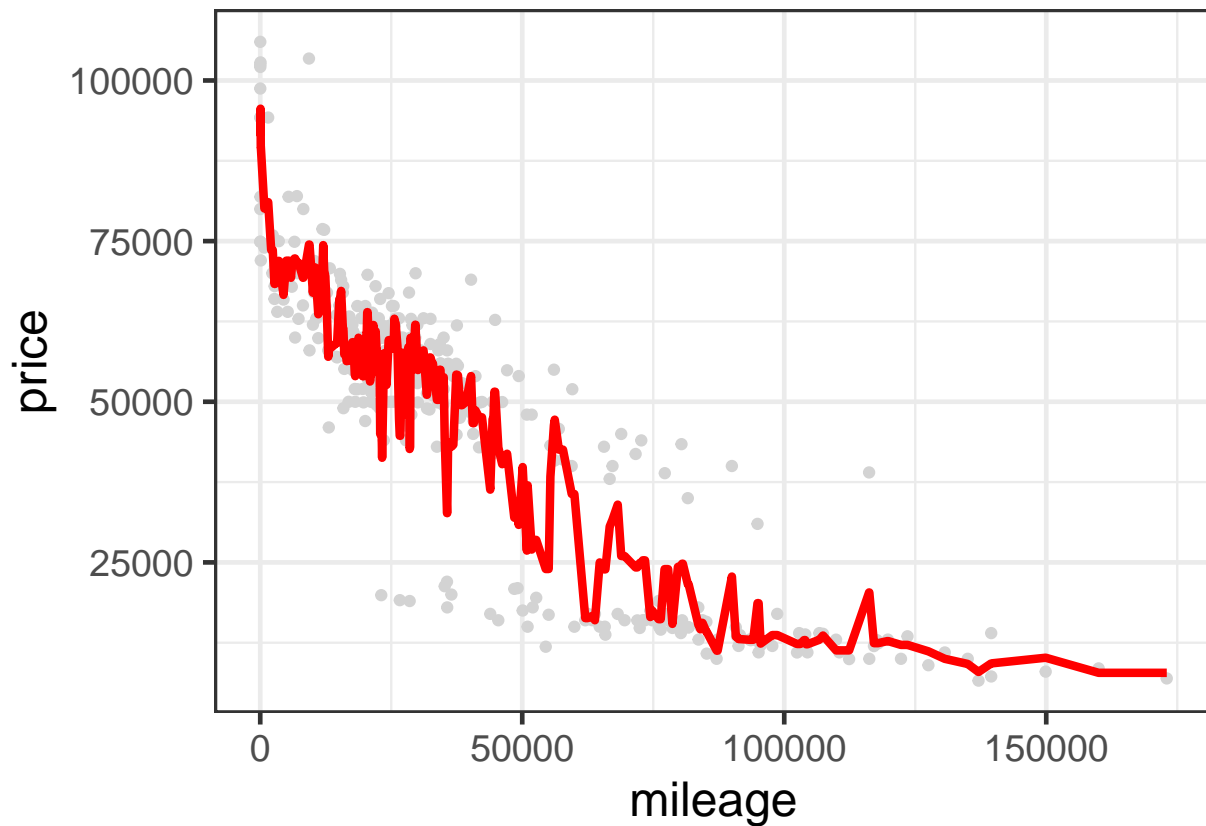


```
p_test + geom_point(aes(x = mileage, y = ypred_knn3), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 3)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 4
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn4 = knn.reg(train = X_train, test = X_test, y = y_train, k=4)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn4 = knn4$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn4)
```

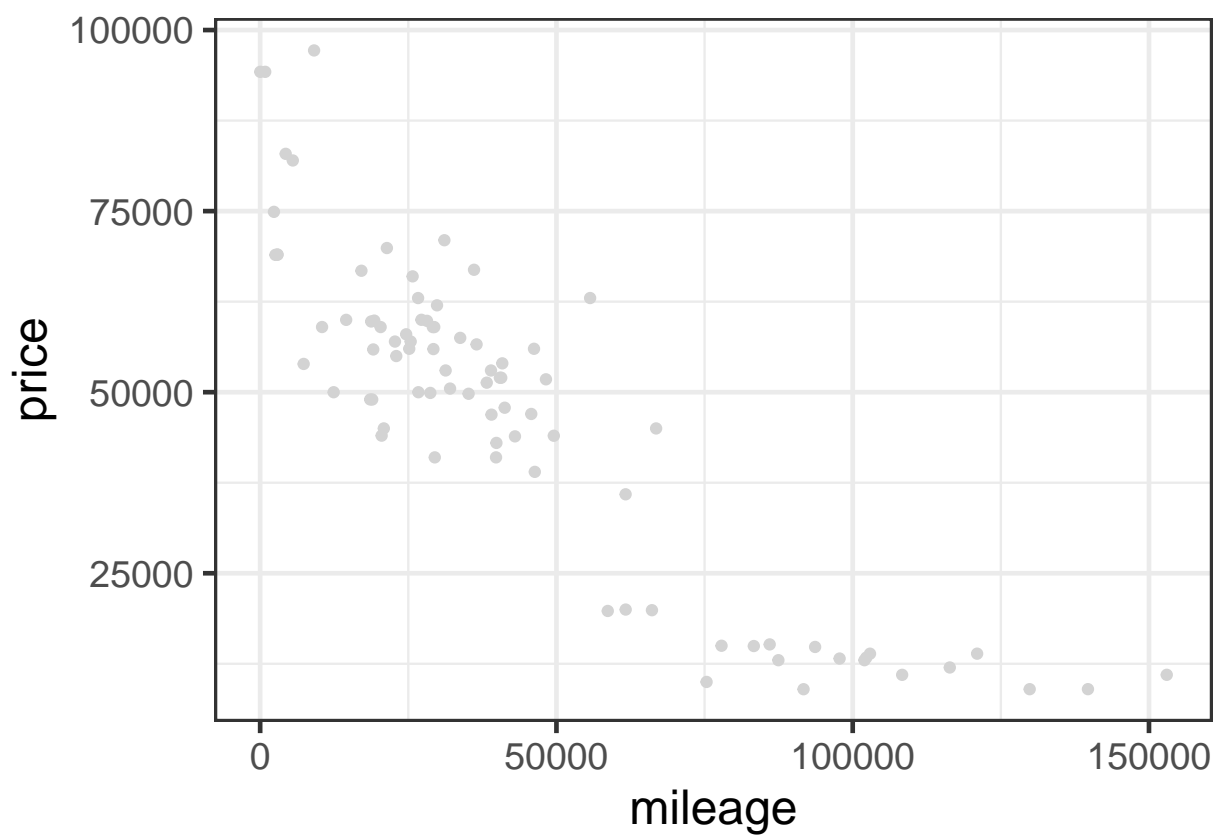
```
## [1] 9027.238
```

```
#attach predictions to data frame
```

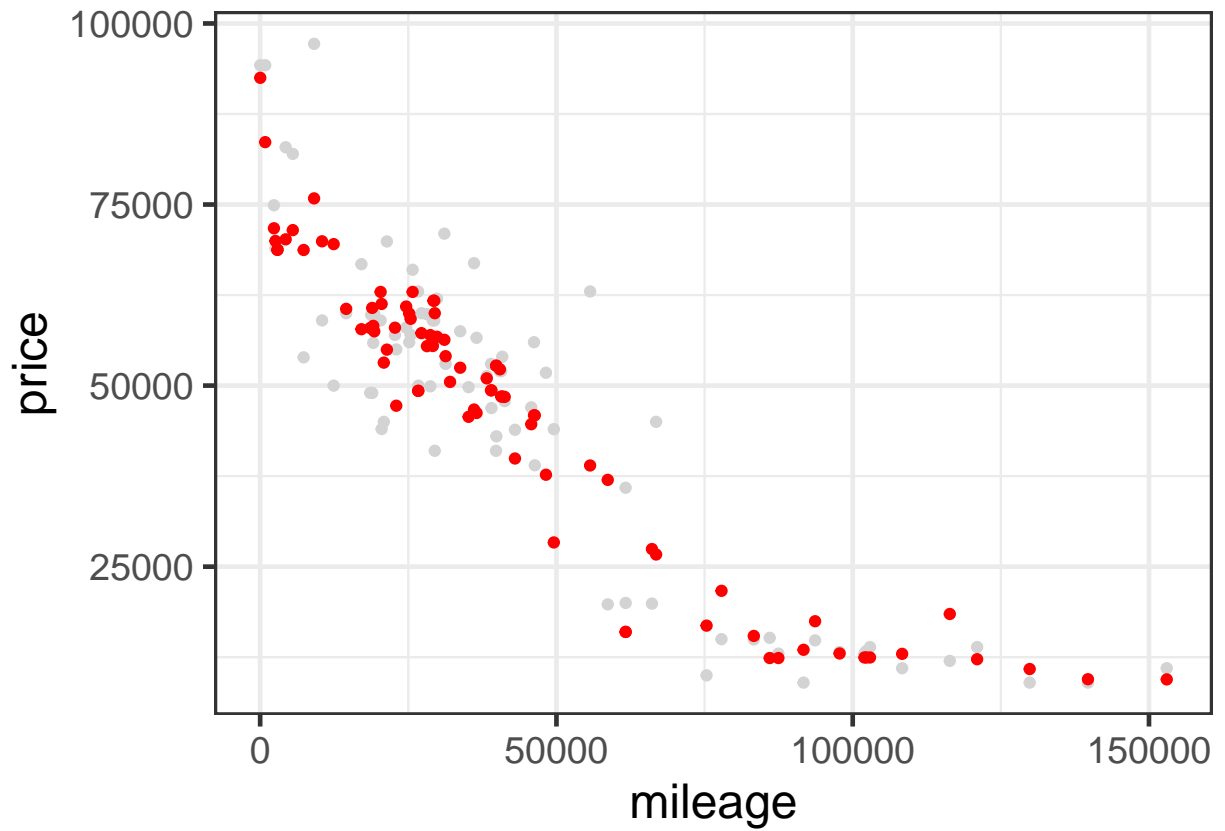
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn4 = ypred_knn4
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

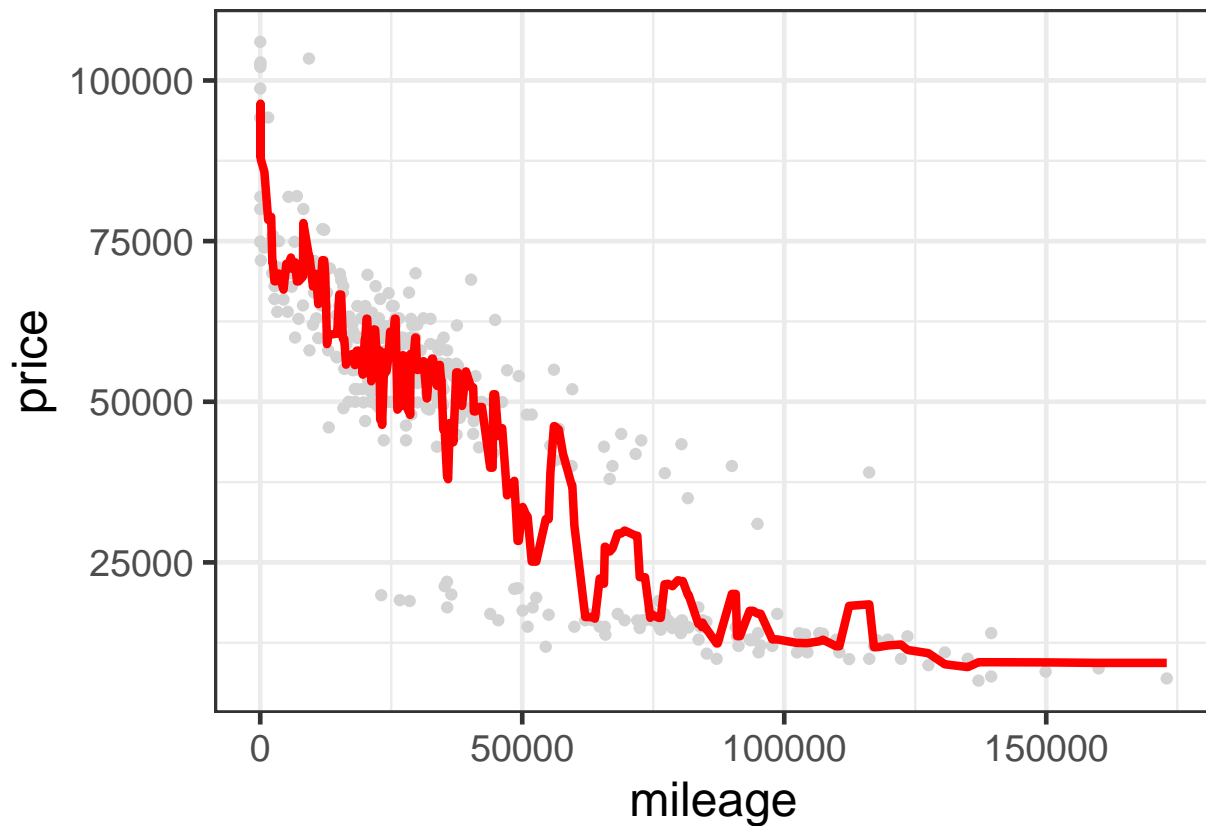


```
p_test + geom_point(aes(x = mileage, y = ypred_knn4), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 4)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```

```
#K = 5
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn5 = knn.reg(train = X_train, test = X_test, y = y_train, k=5)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn5 = knn5$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn5)
```

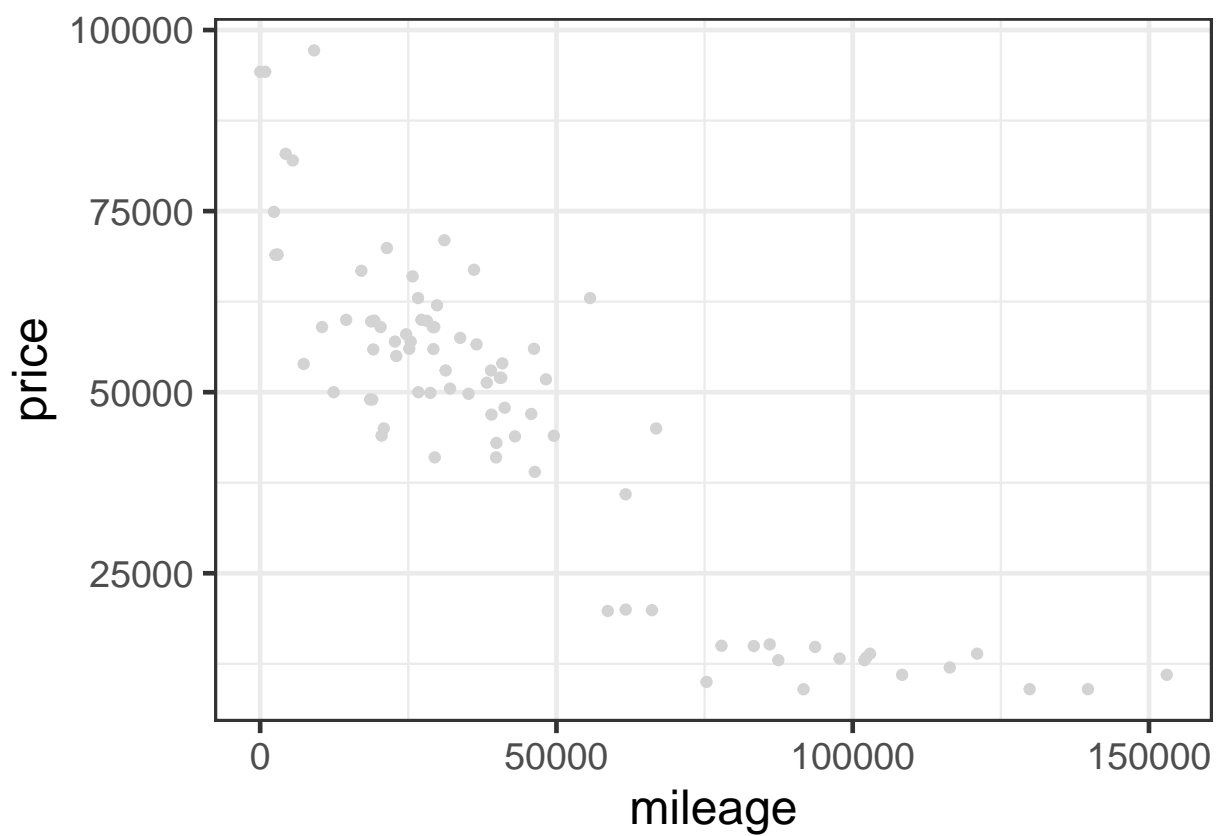
```
## [1] 9152.831
```

```
#attach predictions to data frame
```

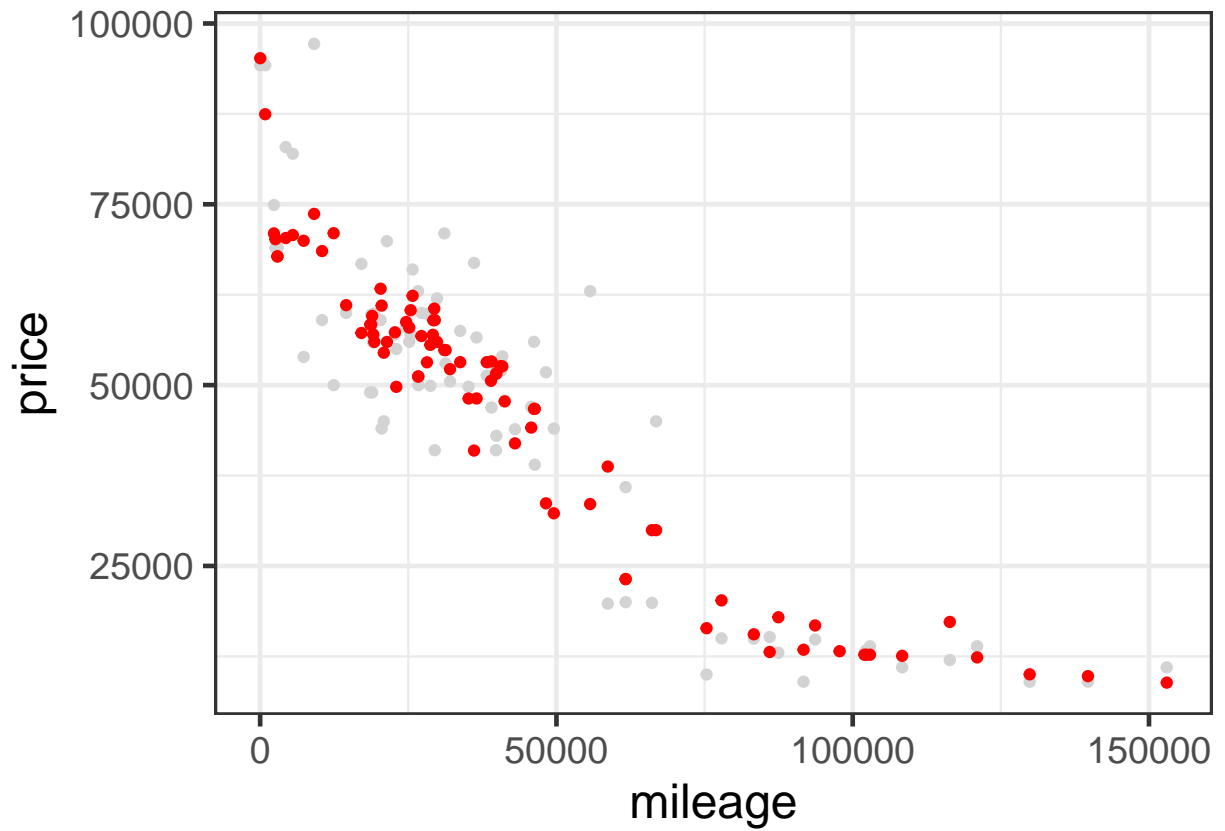
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn5 = ypred_knn5
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

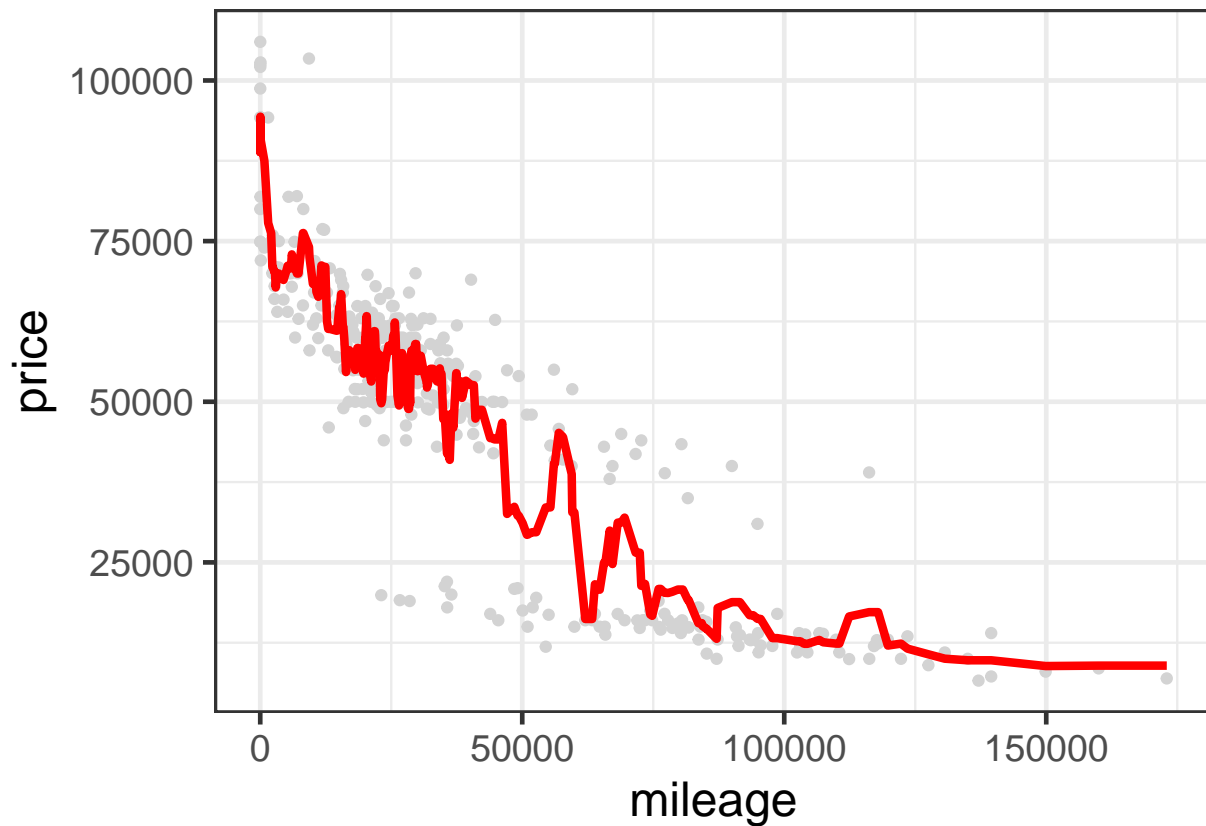


```
p_test + geom_point(aes(x = mileage, y = ypred_knn5), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 5)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 6
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn6 = knn.reg(train = X_train, test = X_test, y = y_train, k=6)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn6 = knn6$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn6)
```

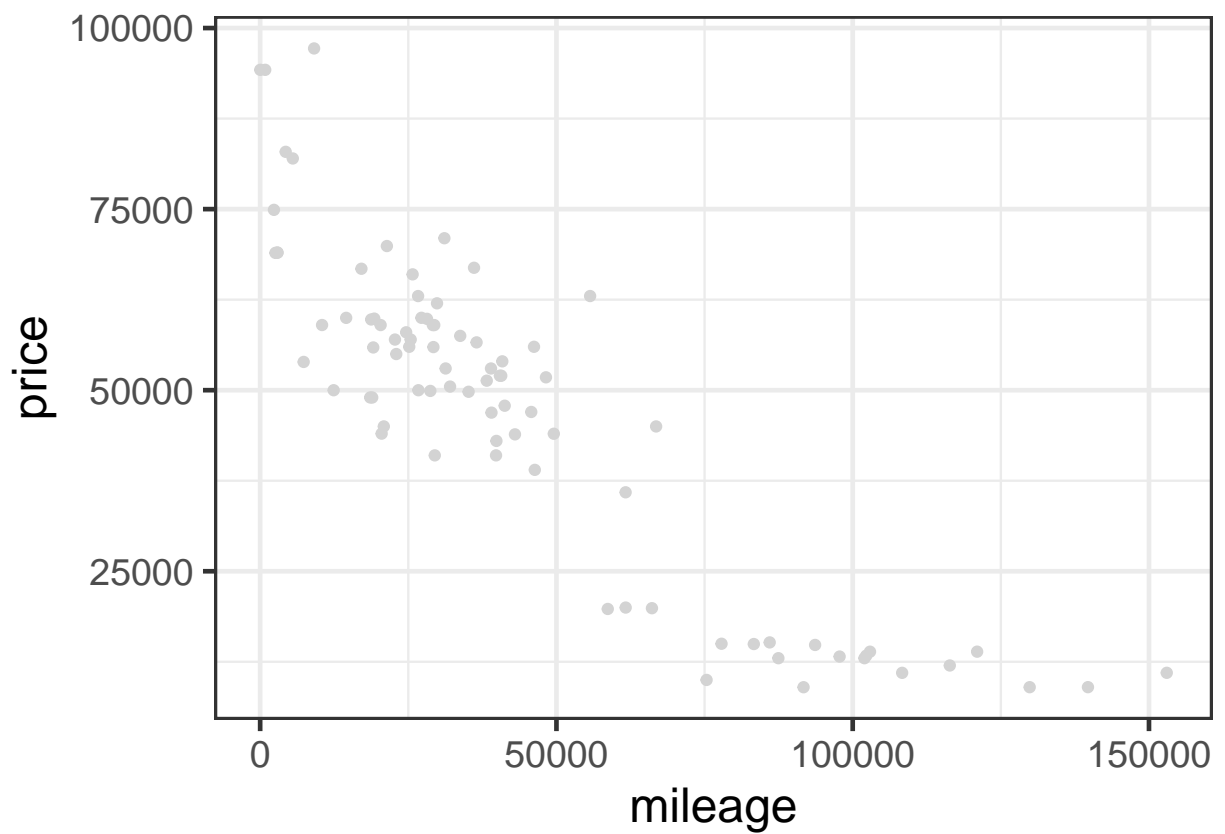
```
## [1] 9265.592
```

```
#attach predictions to data frame
```

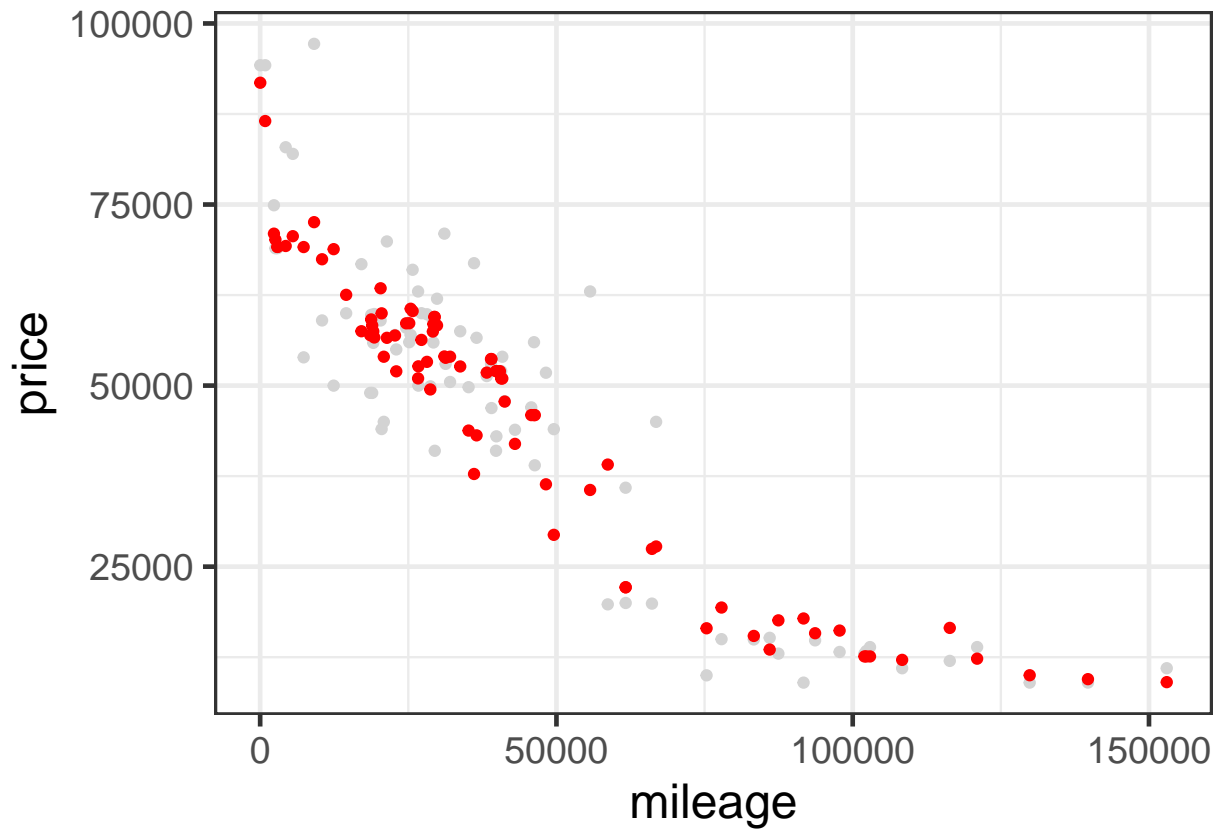
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn6 = ypred_knn6
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

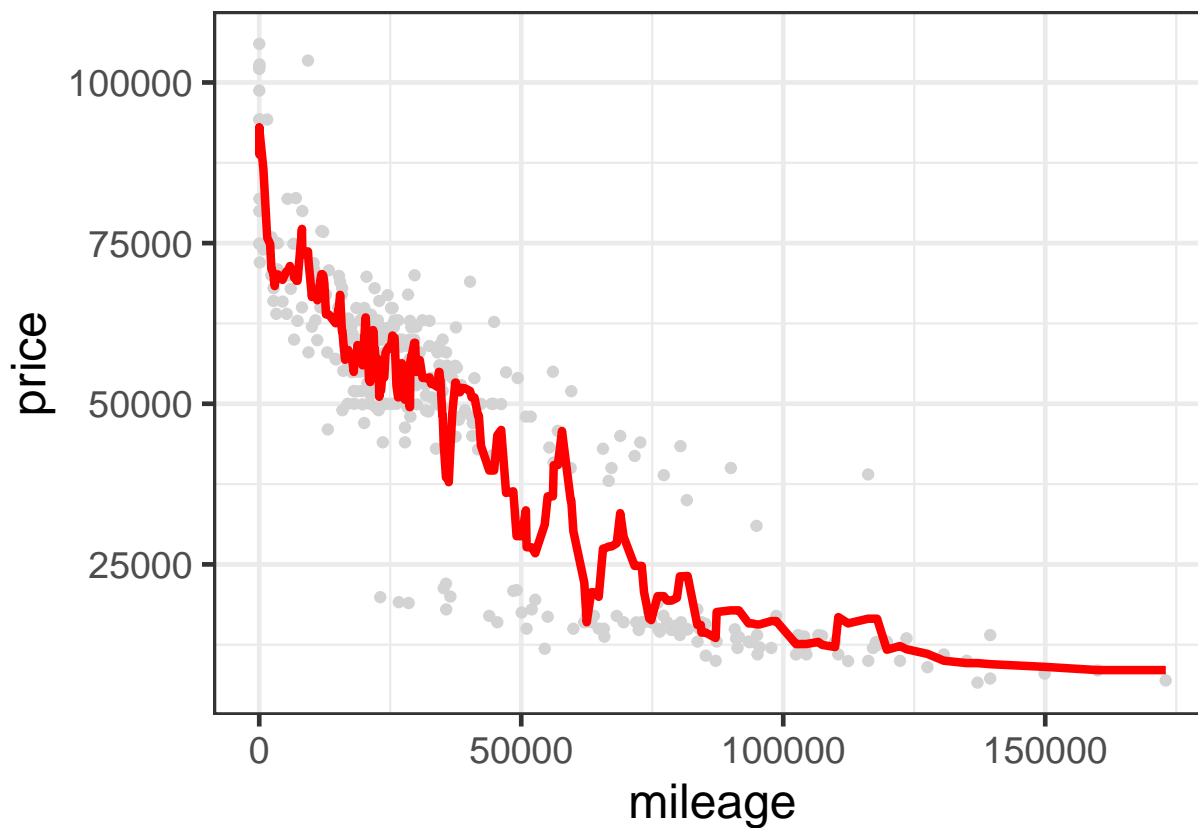


```
p_test + geom_point(aes(x = mileage, y = ypred_knn6), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 6)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 7
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn7 = knn.reg(train = X_train, test = X_test, y = y_train, k=7)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn7 = knn7$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn7)
```

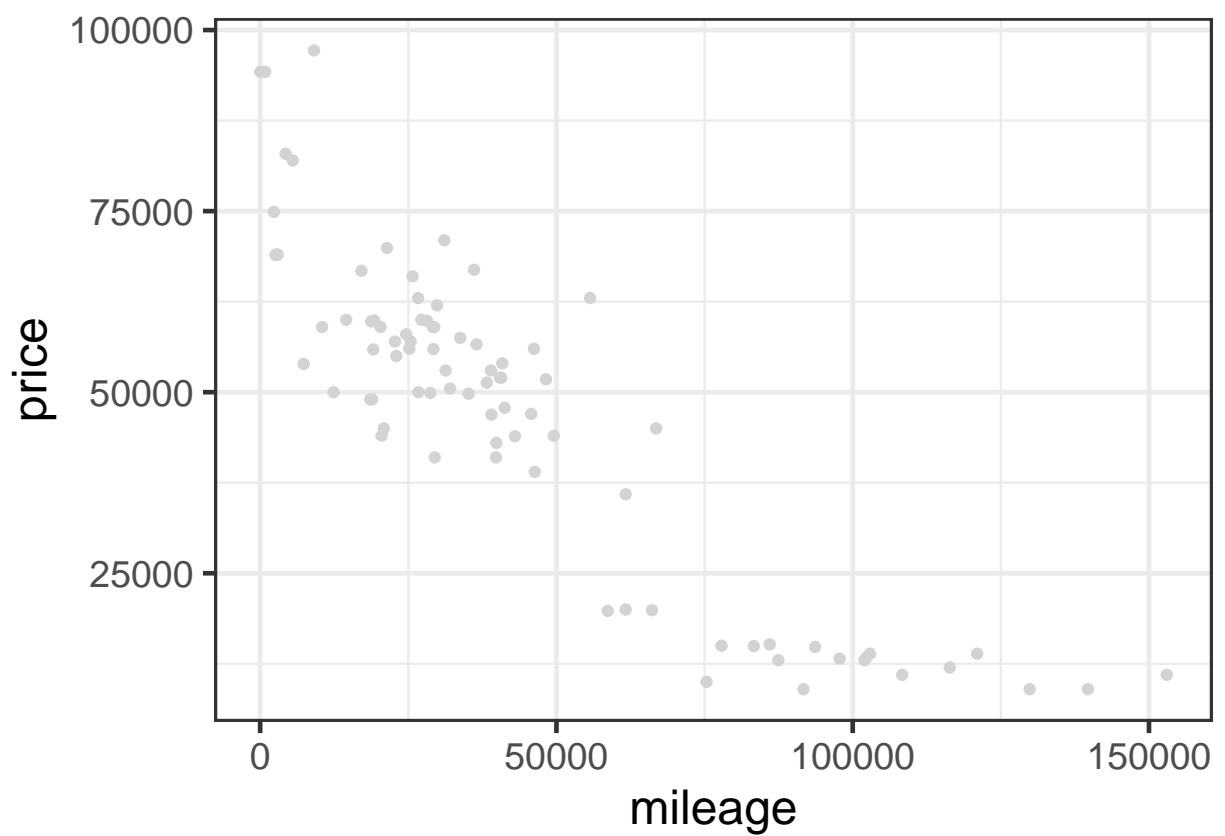
```
## [1] 9177.549
```

```
#attach predictions to data frame
```

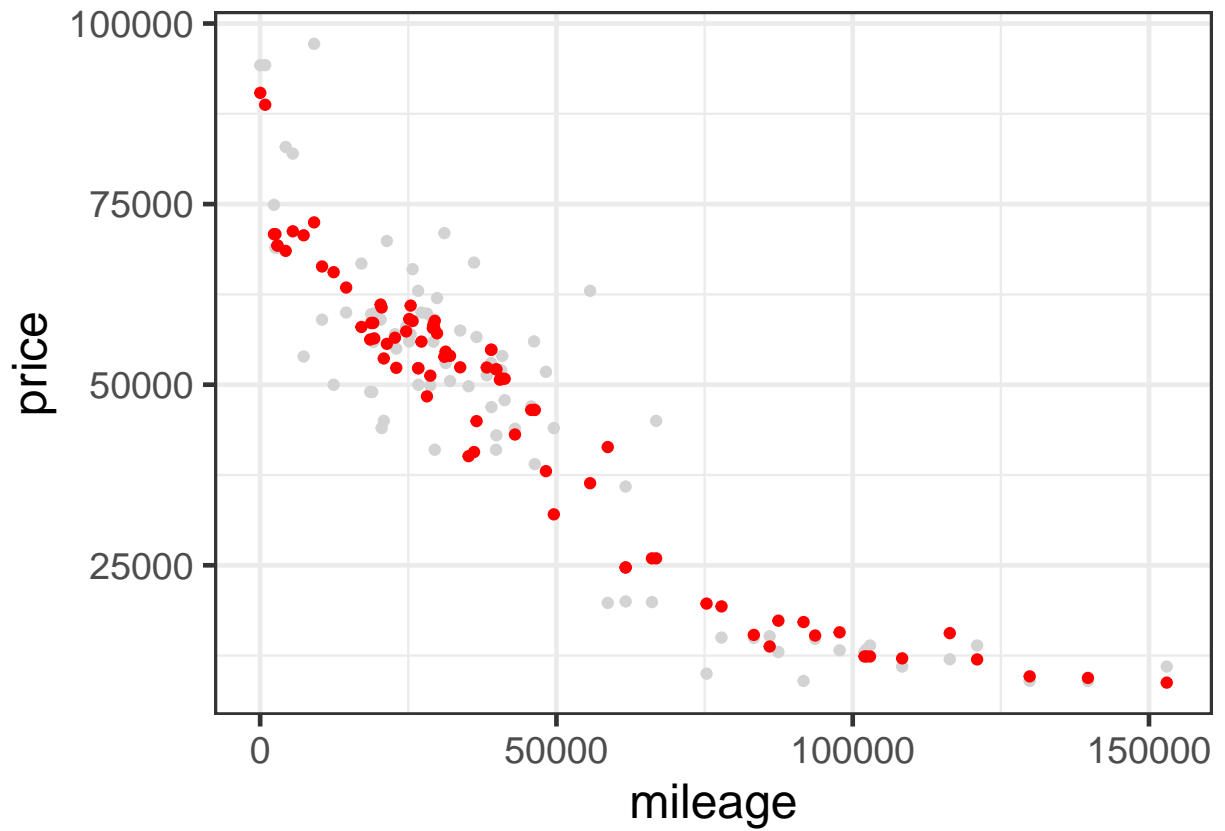
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn7 = ypred_knn7
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

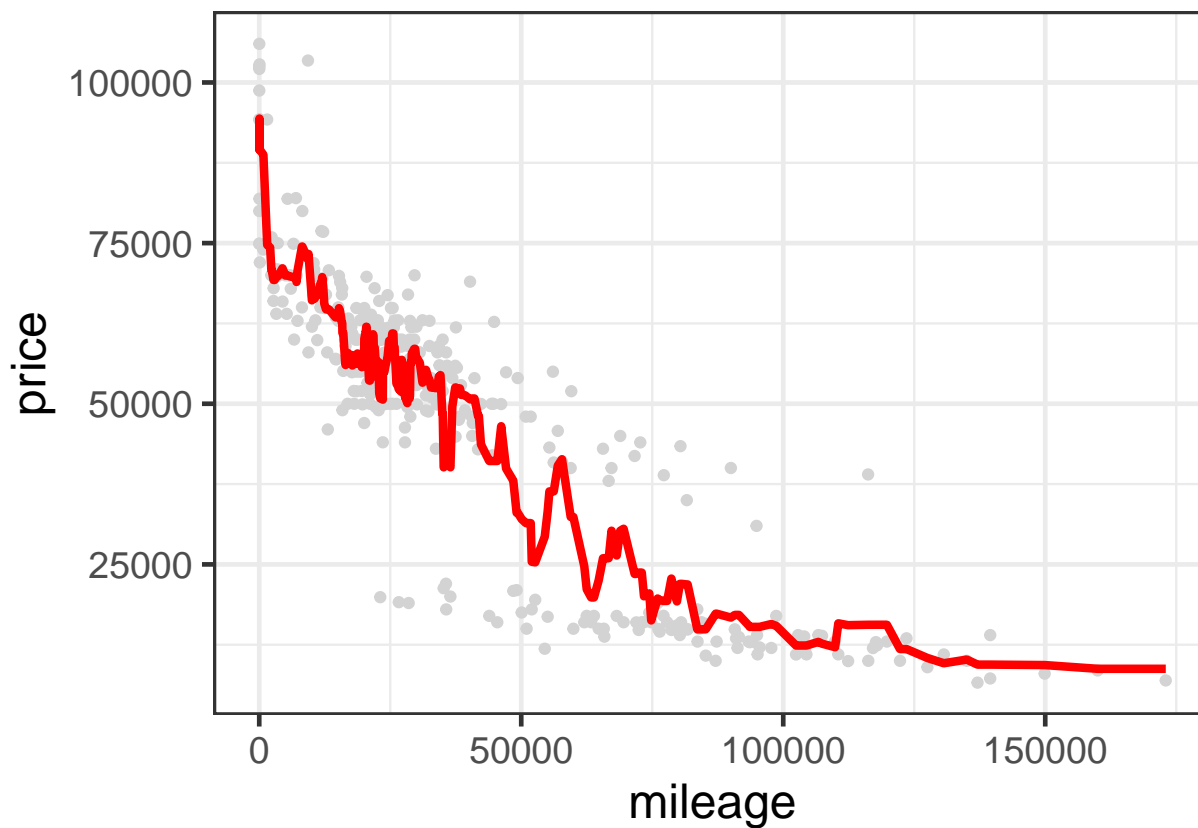


```
p_test + geom_point(aes(x = mileage, y = ypred_knn7), color='red')
```

```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 7)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 8
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn8 = knn.reg(train = X_train, test = X_test, y = y_train, k=8)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn8 = knn8$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn8)
```

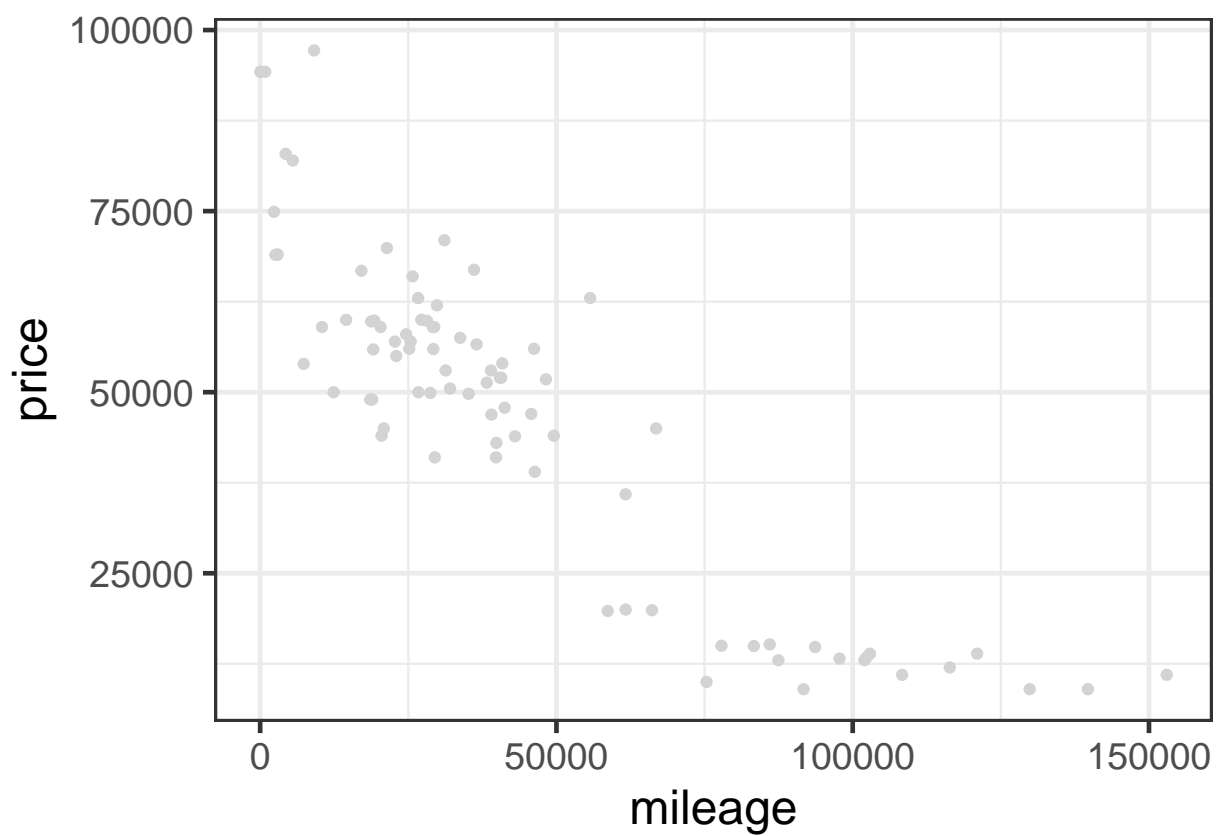
```
## [1] 9169.083
```

```
#attach predictions to data frame
```

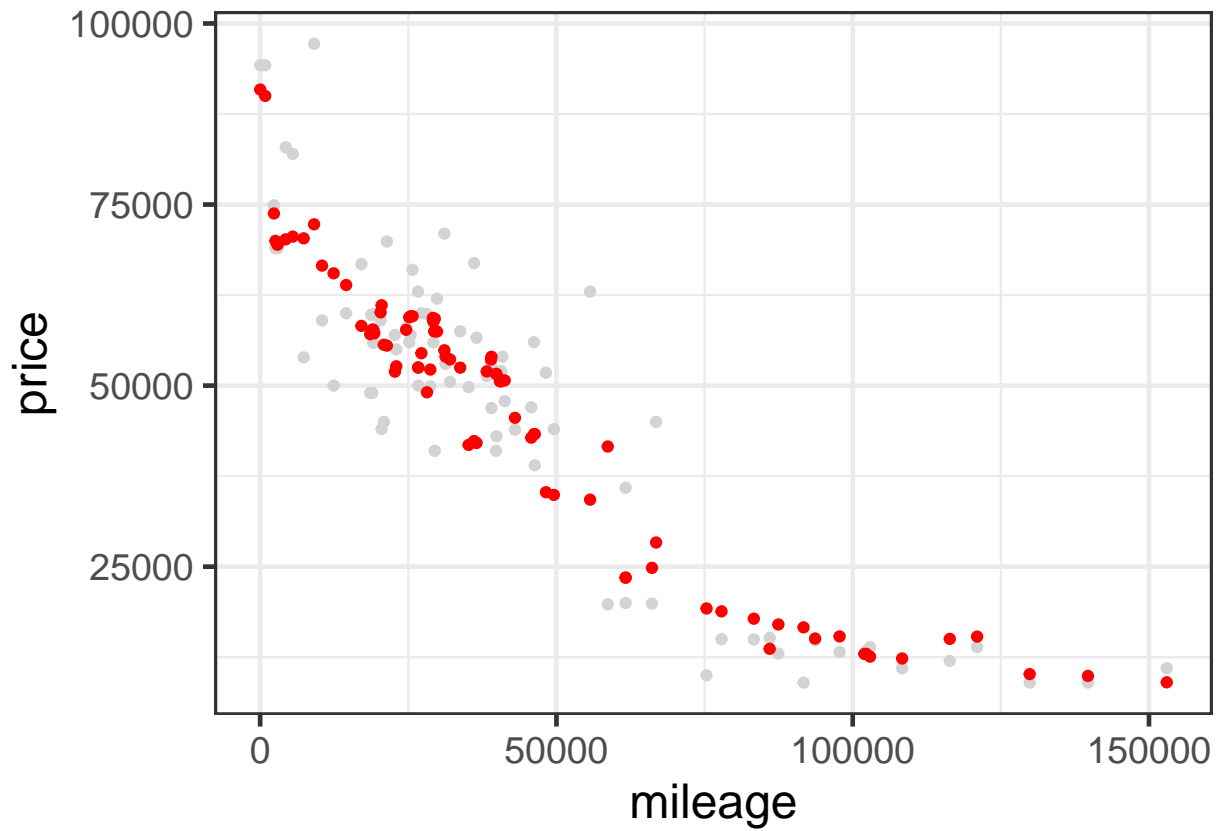
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn8 = ypred_knn8
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

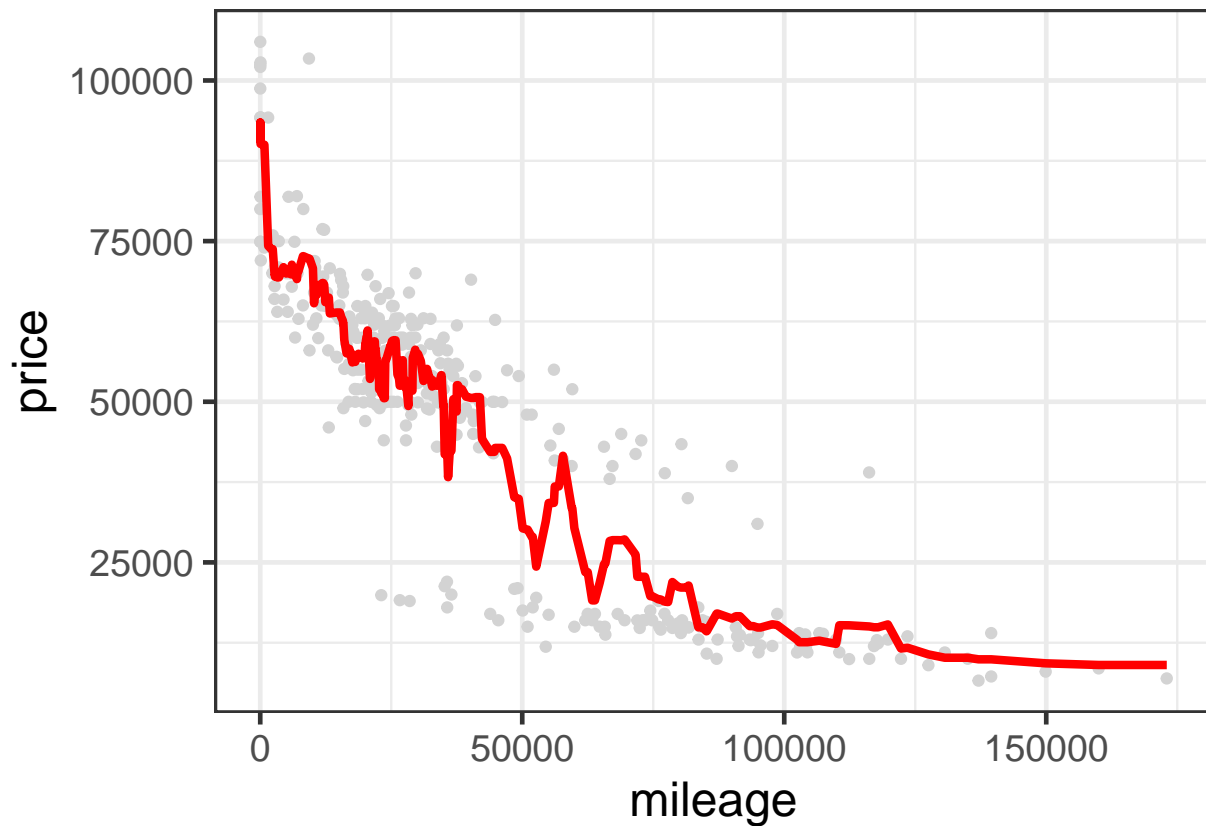


```
p_test + geom_point(aes(x = mileage, y = ypred_knn8), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 8)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 9
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn9 = knn.reg(train = X_train, test = X_test, y = y_train, k=9)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn9 = knn9$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn9)
```

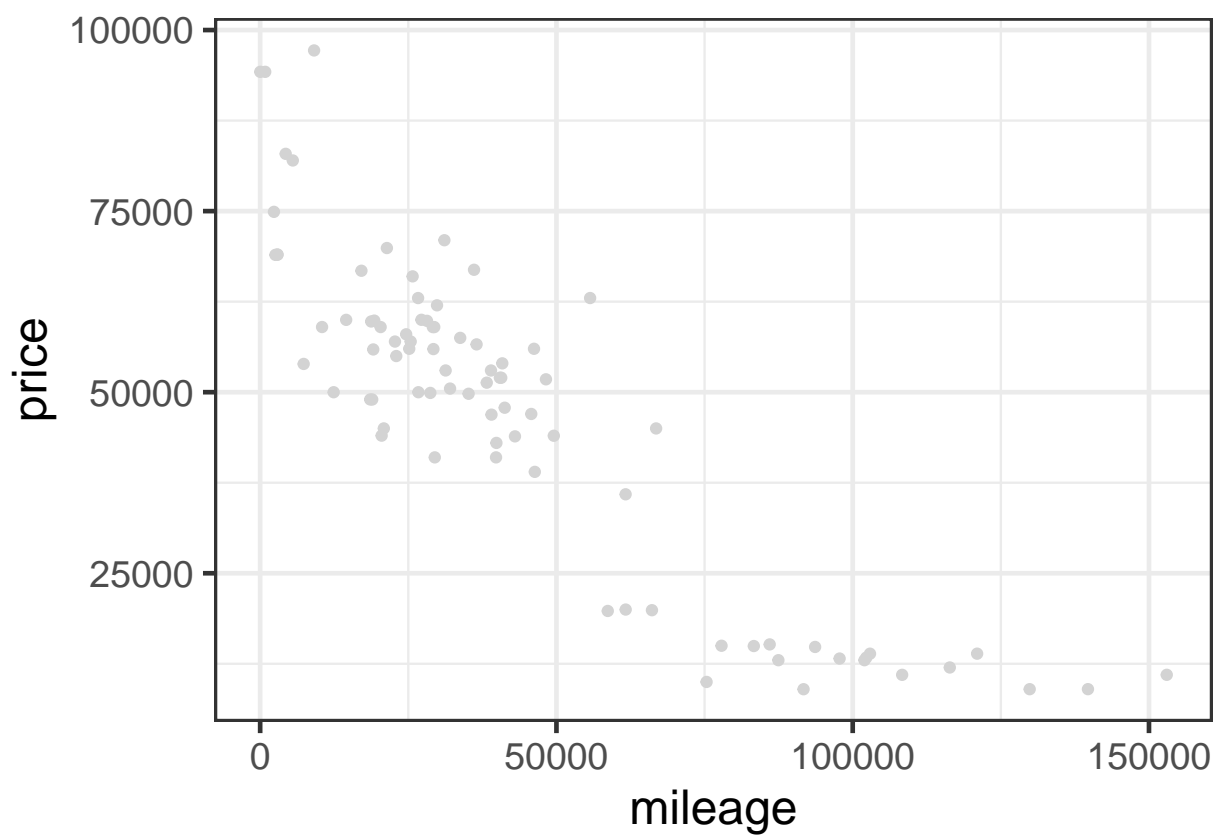
```
## [1] 9322.982
```

```
#attach predictions to data frame
```

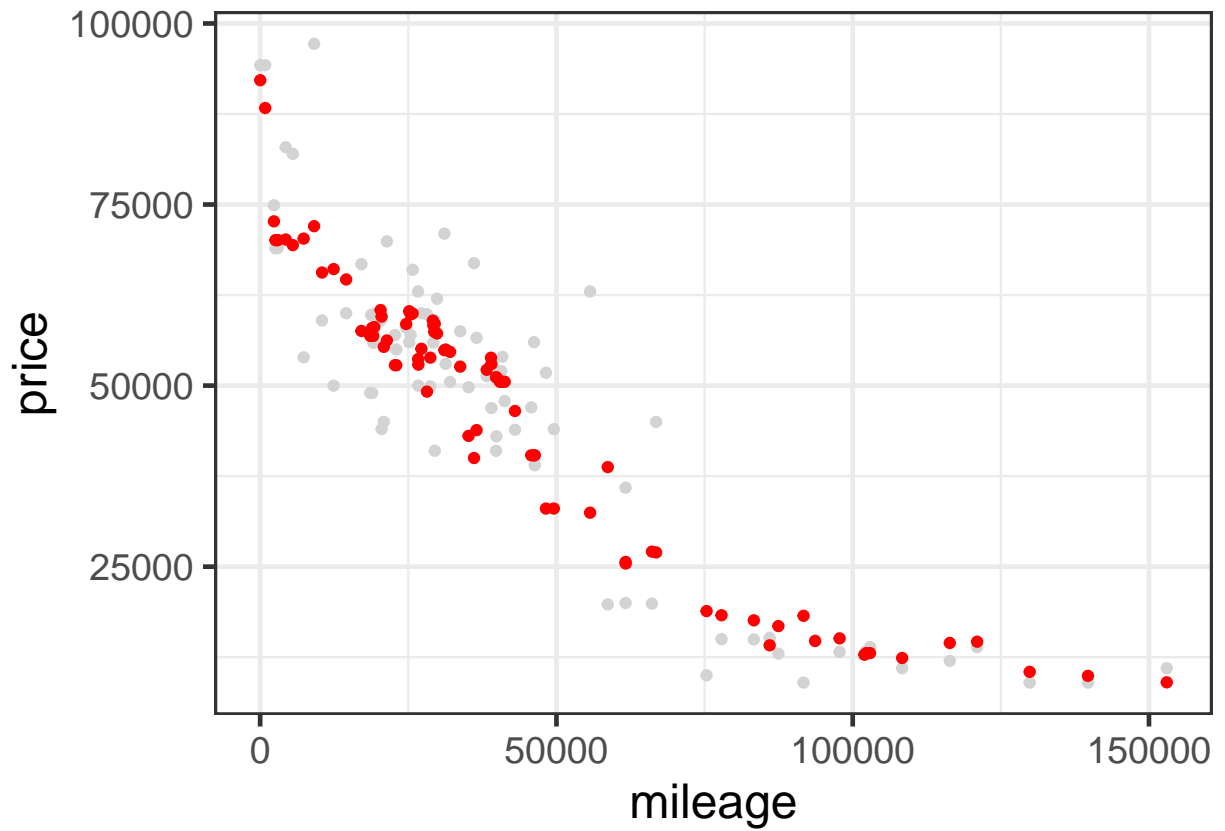
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn9 = ypred_knn9
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

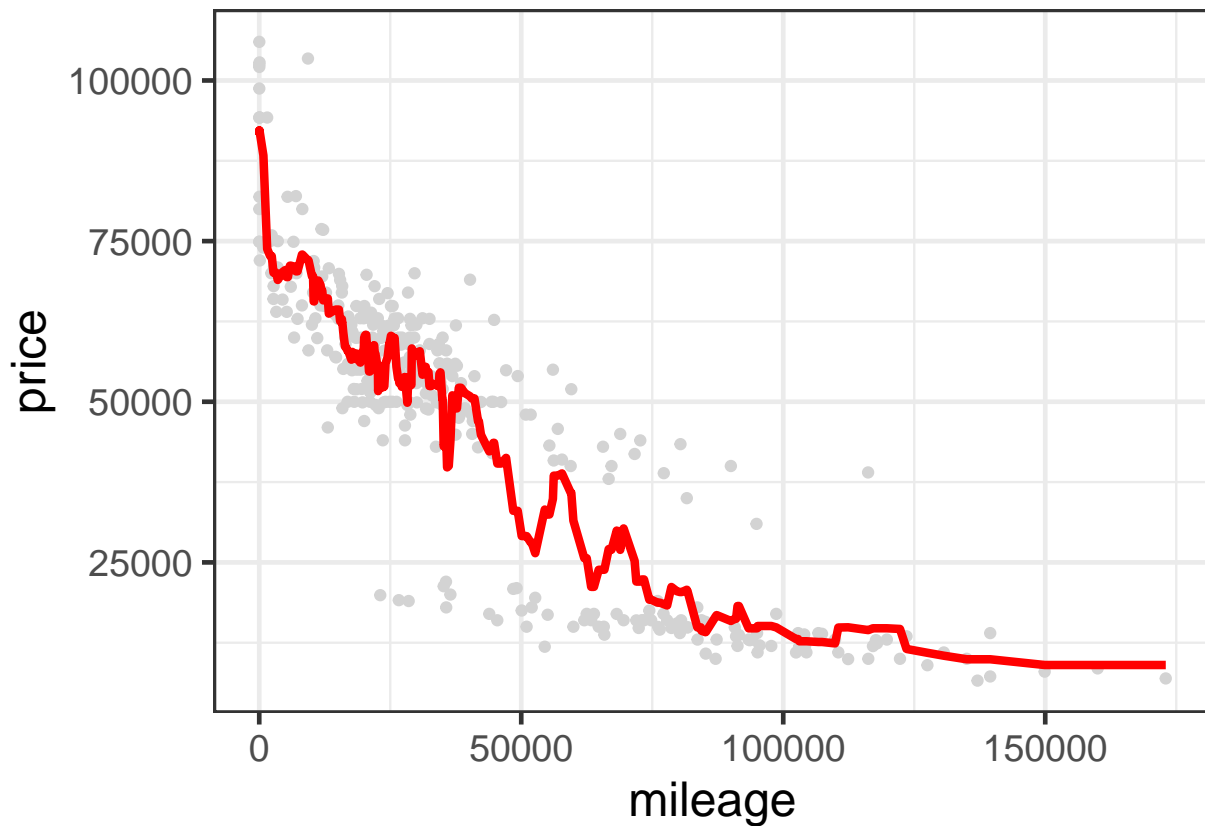


```
p_test + geom_point(aes(x = mileage, y = ypred_knn9), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 9)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 6
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn6 = knn.reg(train = X_train, test = X_test, y = y_train, k=6)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn6 = knn6$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```



```
rmse(y_test, ypred_knn6)
```

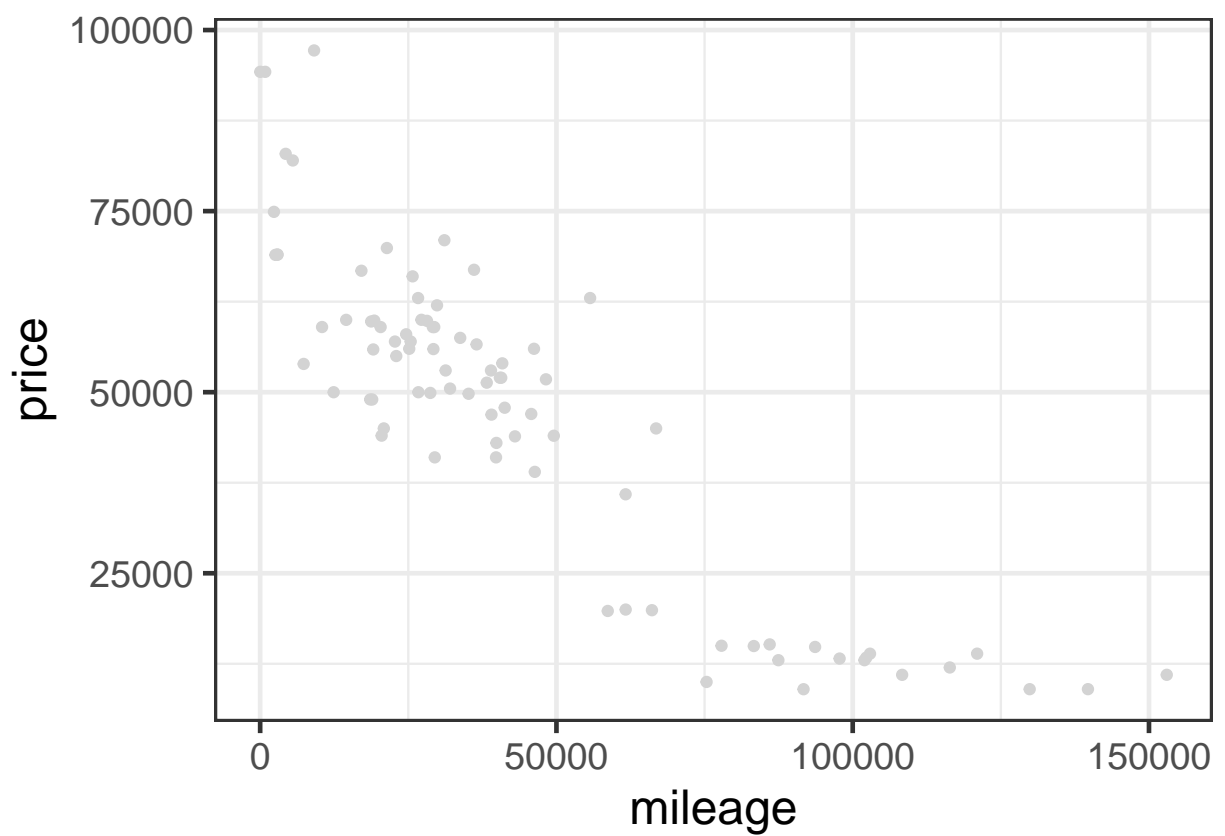
```
## [1] 9265.592
```

```
#attach predictions to data frame
```

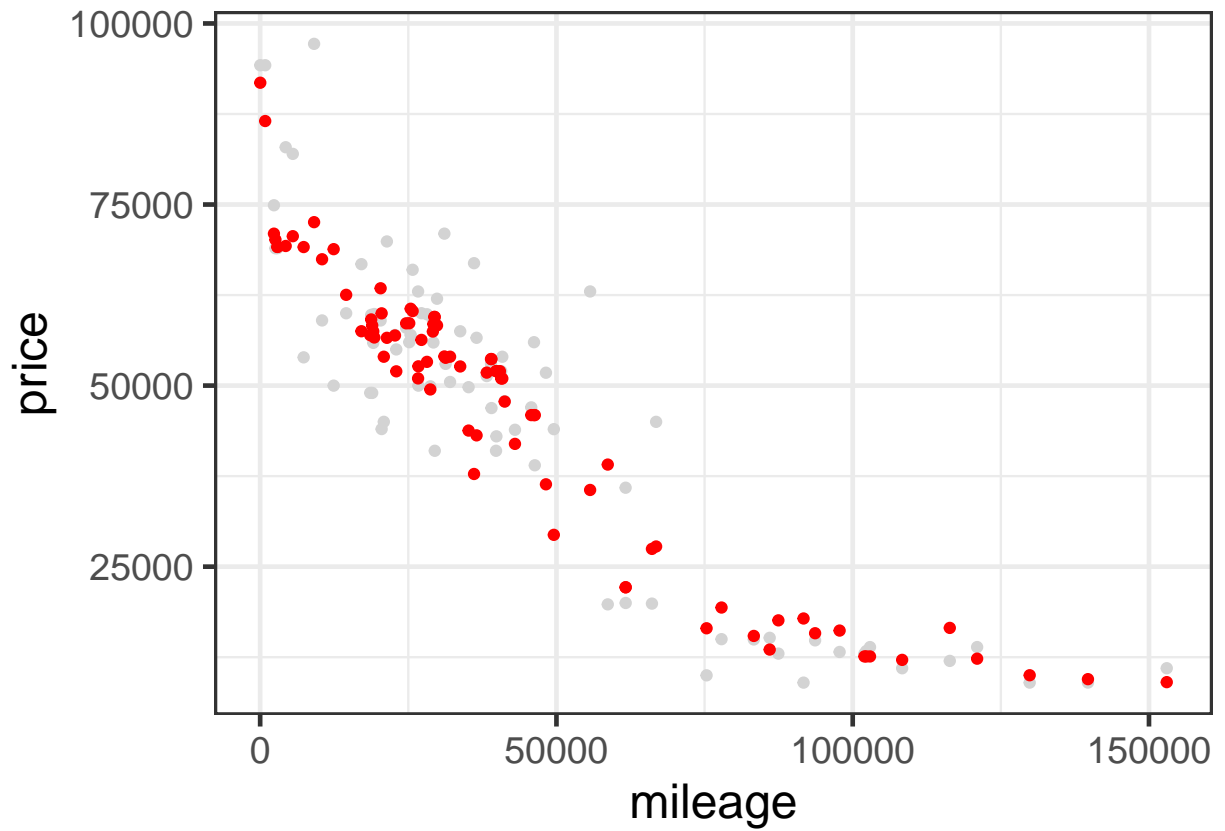
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn6 = ypred_knn6
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

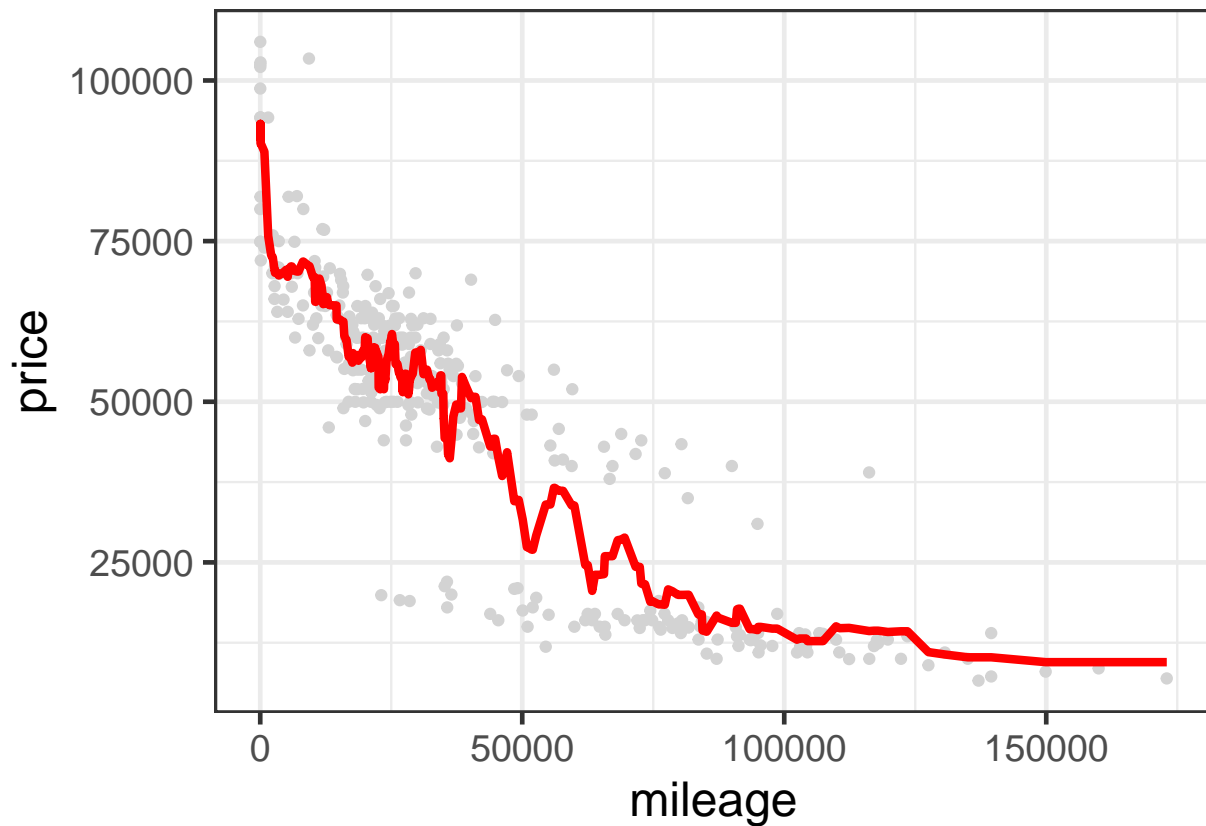


```
p_test + geom_point(aes(x = mileage, y = ypred_knn6), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 10)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 15
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn15 = knn.reg(train = X_train, test = X_test, y = y_train, k=15)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn15 = knn15$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn15)
```

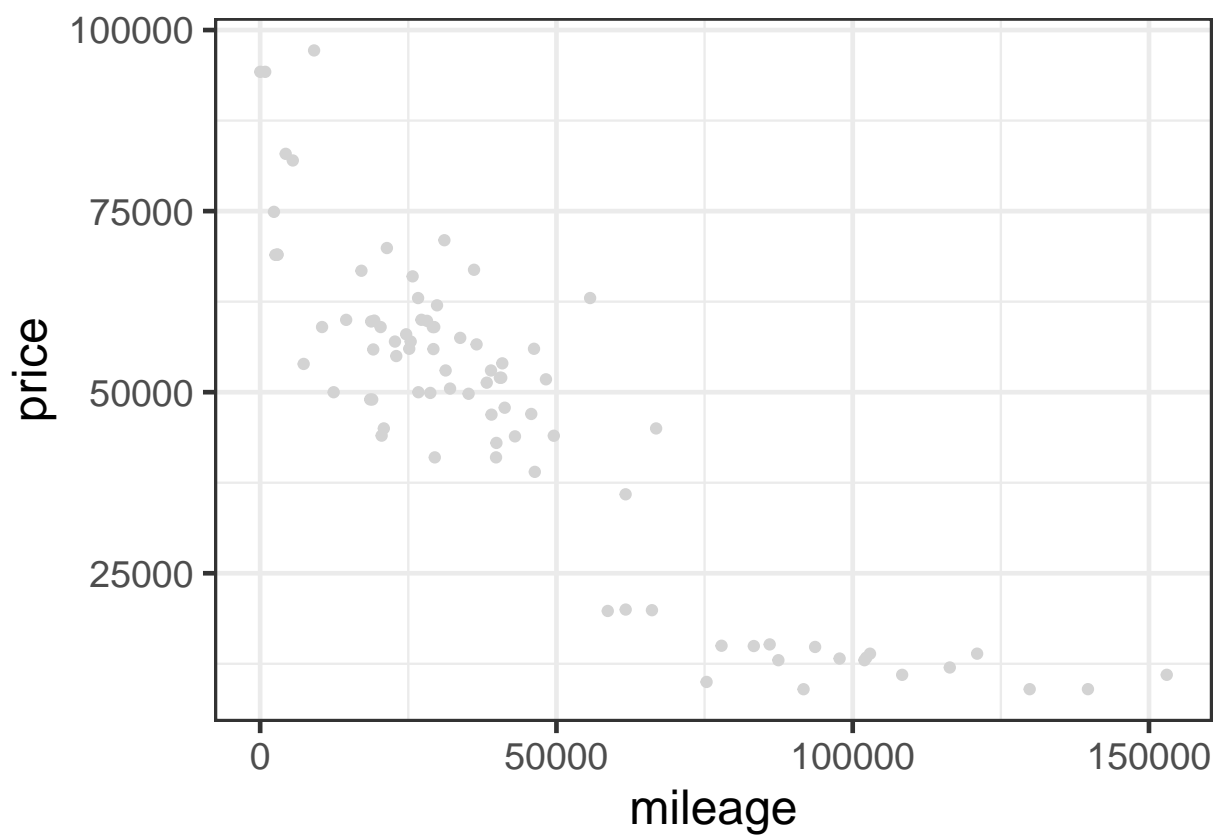
```
## [1] 8851.994
```

```
#attach predictions to data frame
```

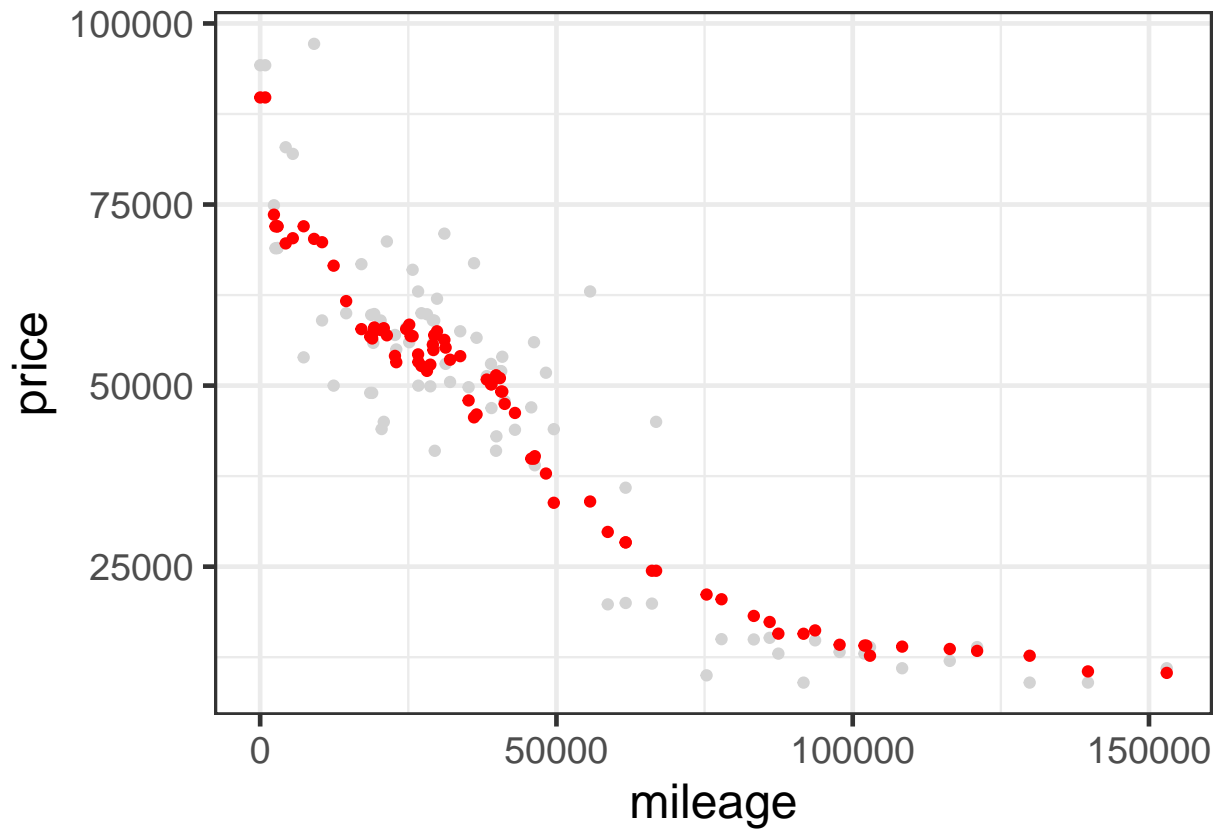
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn15 = ypred_knn15
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

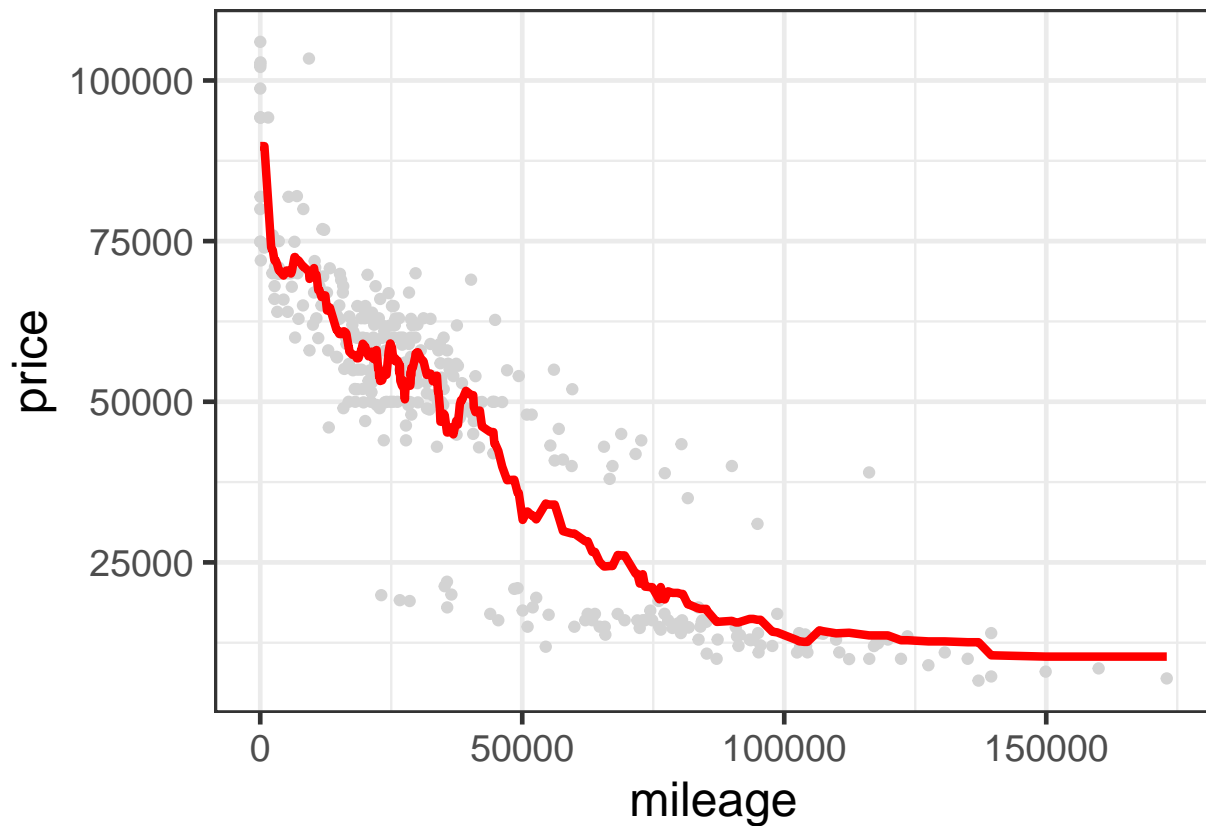


```
p_test + geom_point(aes(x = mileage, y = ypred_knn15), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 15)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 30
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn30 = knn.reg(train = X_train, test = X_test, y = y_train, k=30)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn30 = knn30$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn30)
```

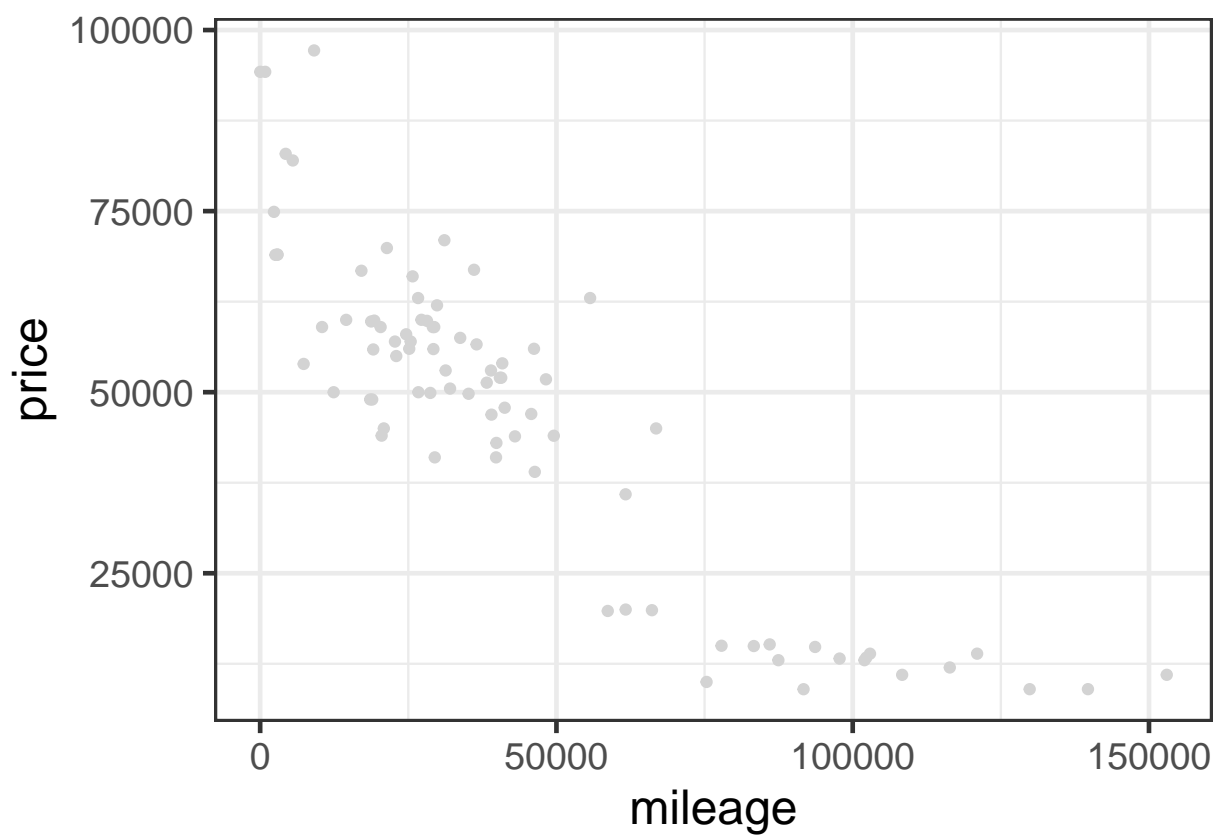
```
## [1] 9213.149
```

```
#attach predictions to data frame
```

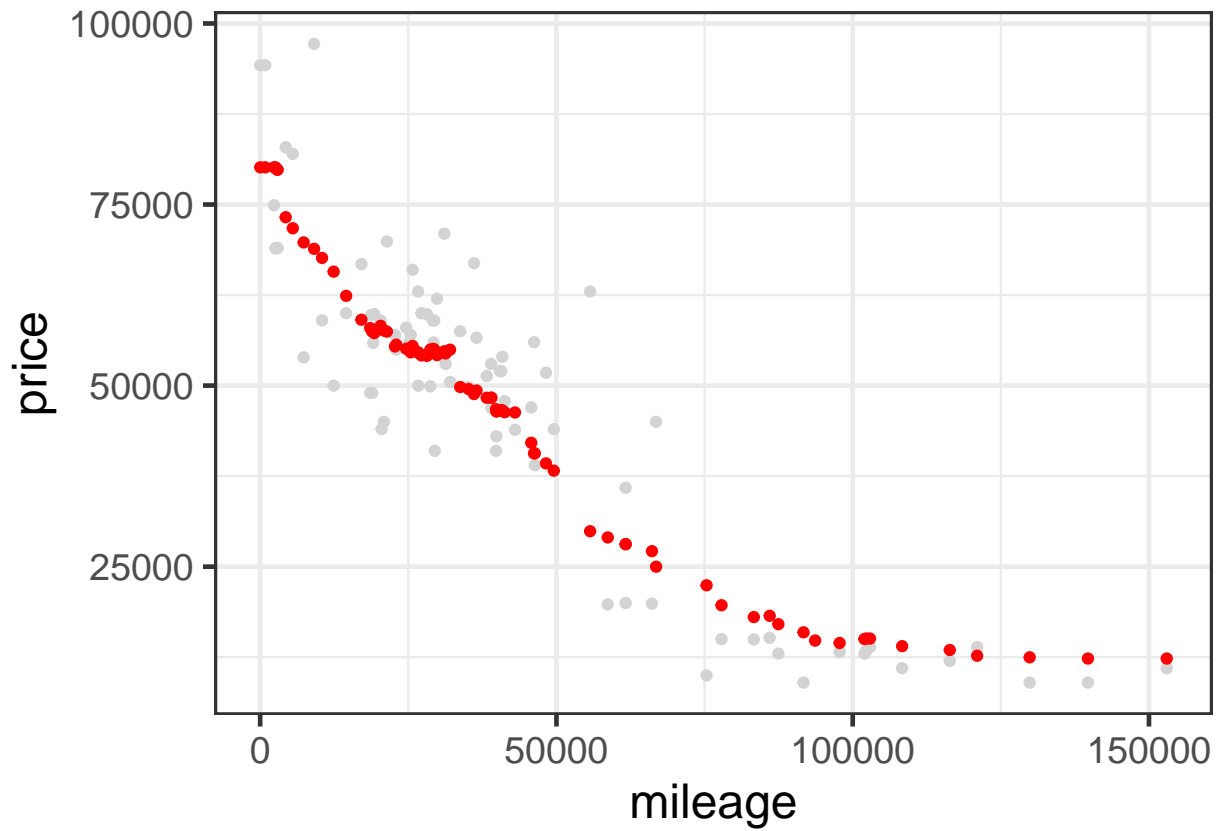
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn30 = ypred_knn30
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

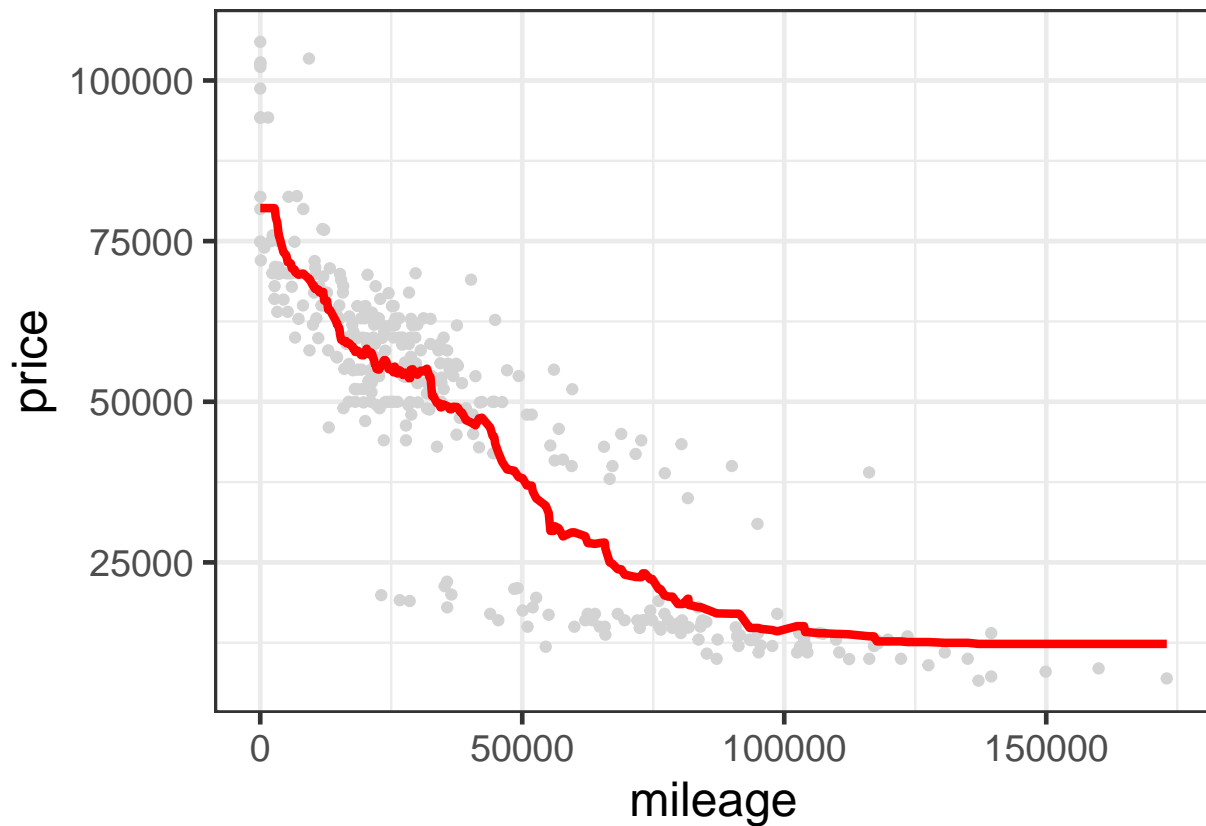


```
p_test + geom_point(aes(x = mileage, y = ypred_knn30), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 30)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```

```
#K = 50
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn50 = knn.reg(train = X_train, test = X_test, y = y_train, k=50)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn50 = knn50$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn50)
```

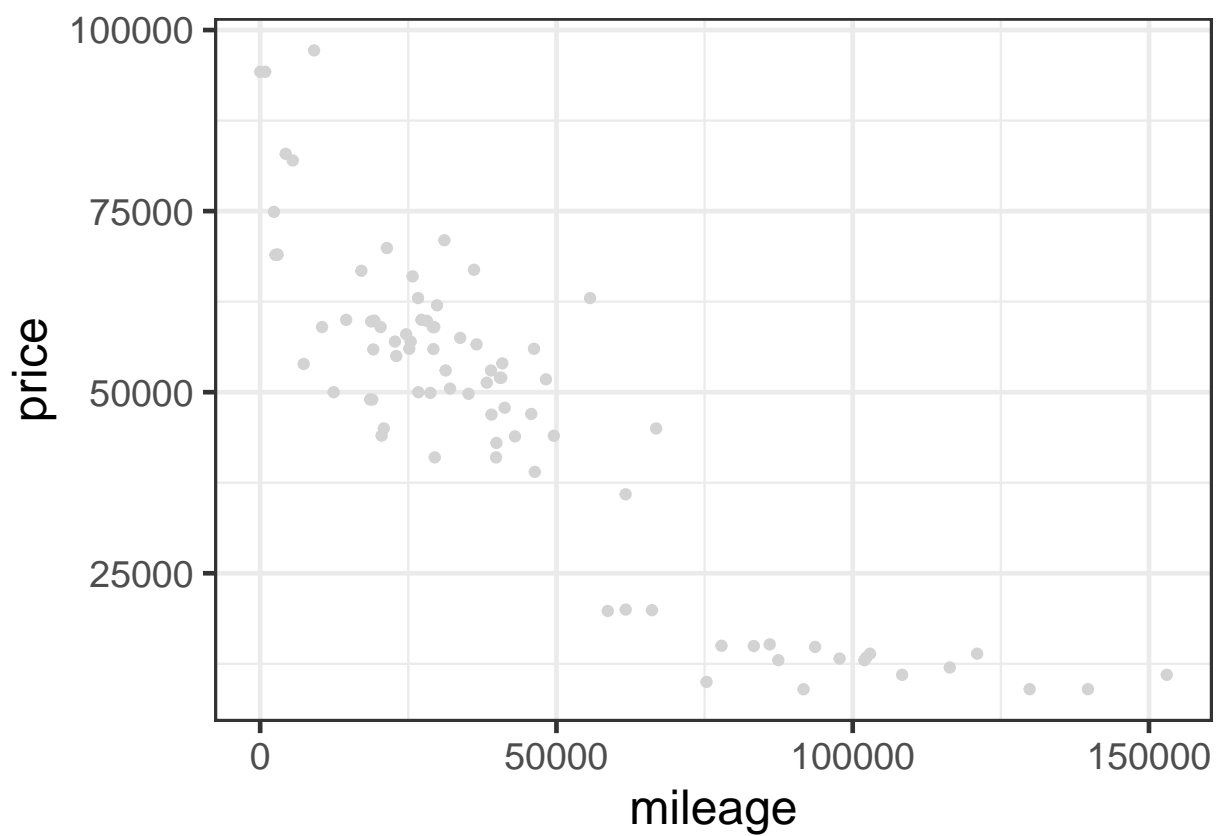
```
## [1] 8842.619
```

```
#attach predictions to data frame
```

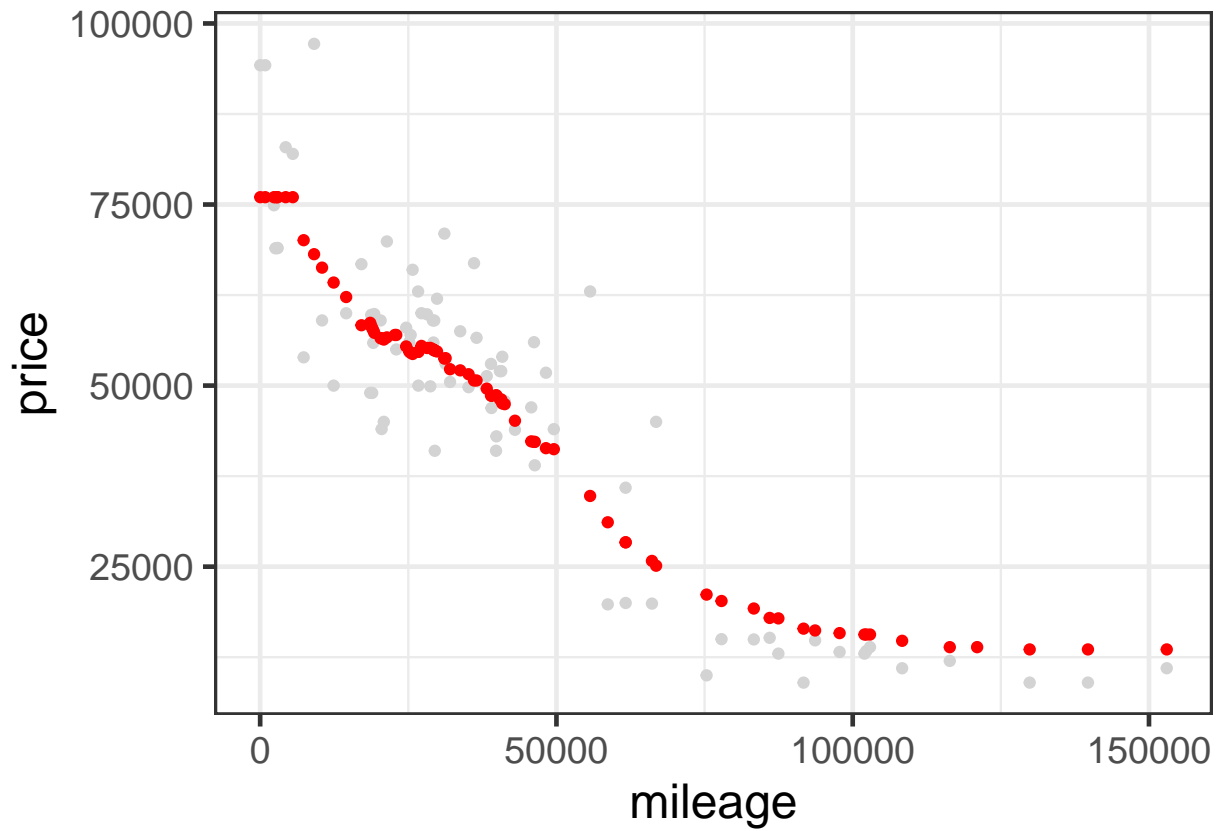
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn50 = ypred_knn50
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

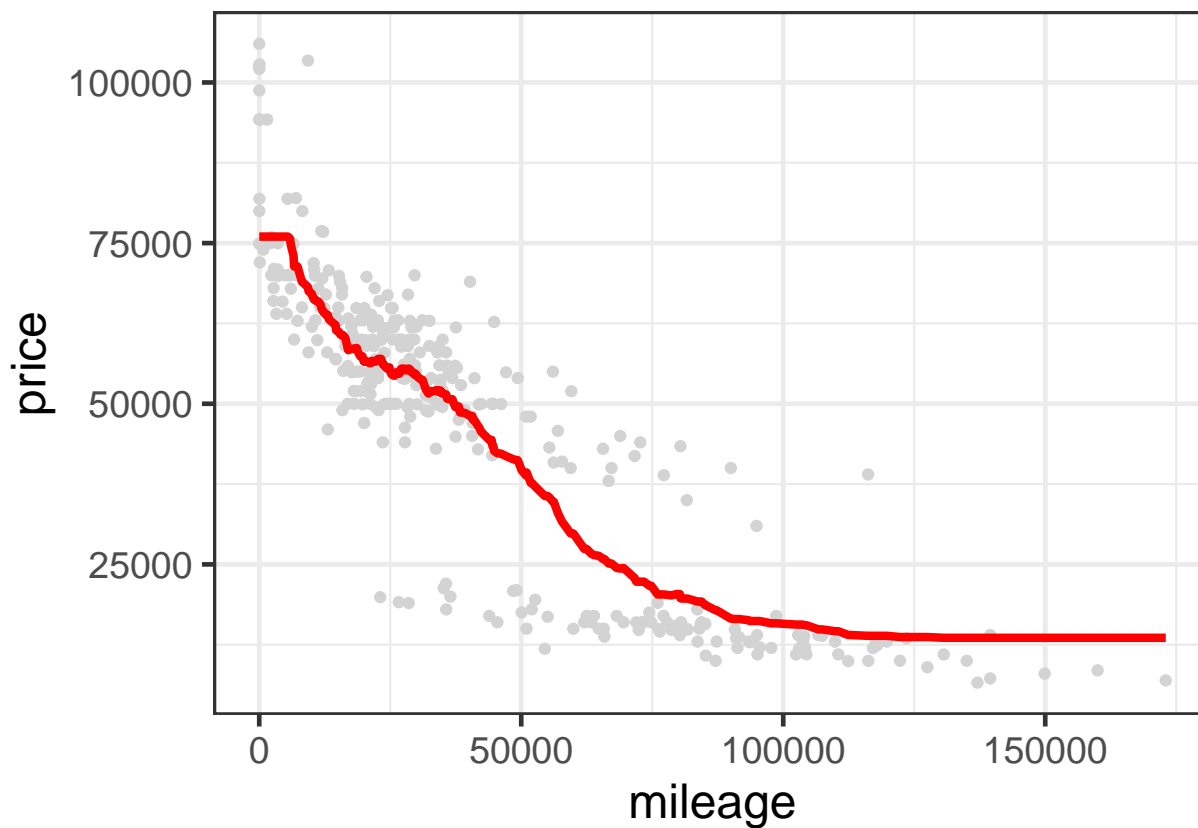


```
p_test + geom_point(aes(x = mileage, y = ypred_knn50), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 50)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 70
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn70 = knn.reg(train = X_train, test = X_test, y = y_train, k=70)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn70 = knn70$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn70)
```

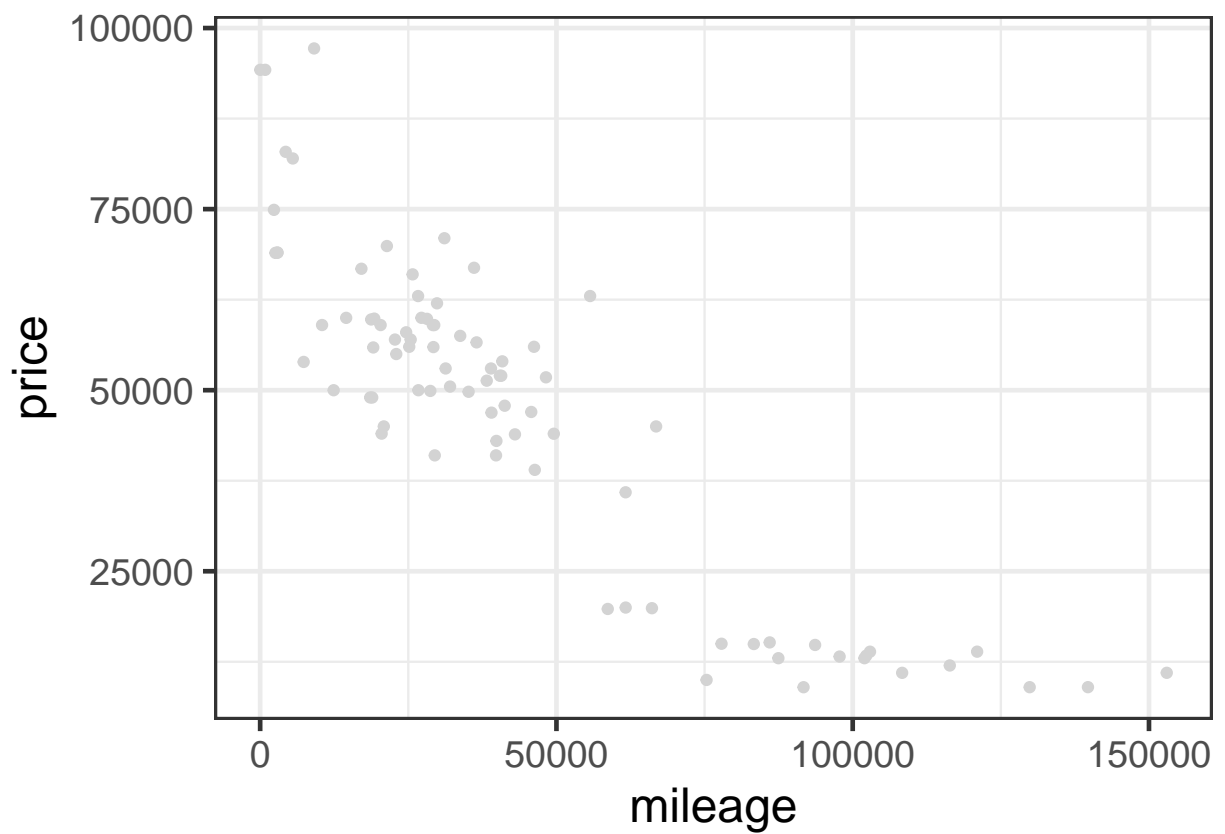
```
## [1] 9110.669
```

```
#attach predictions to data frame
```

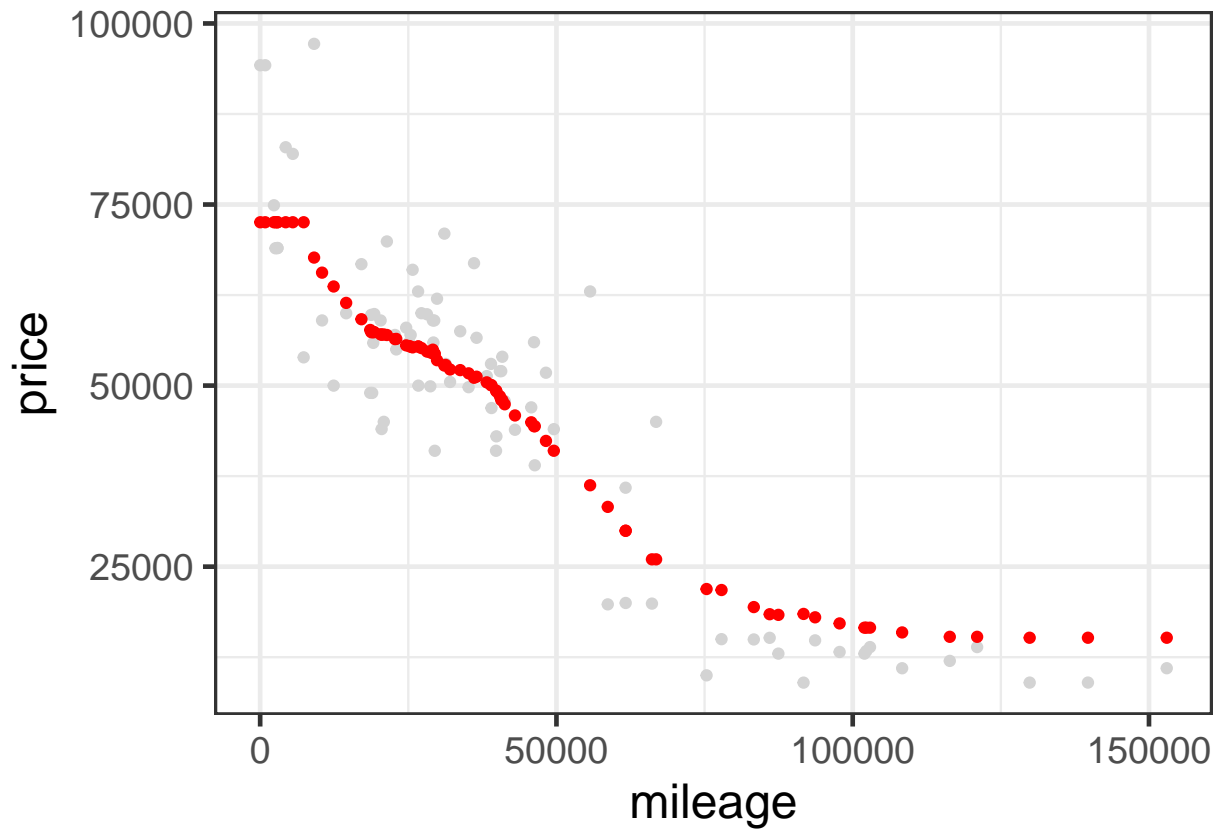
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn70 = ypred_knn70
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

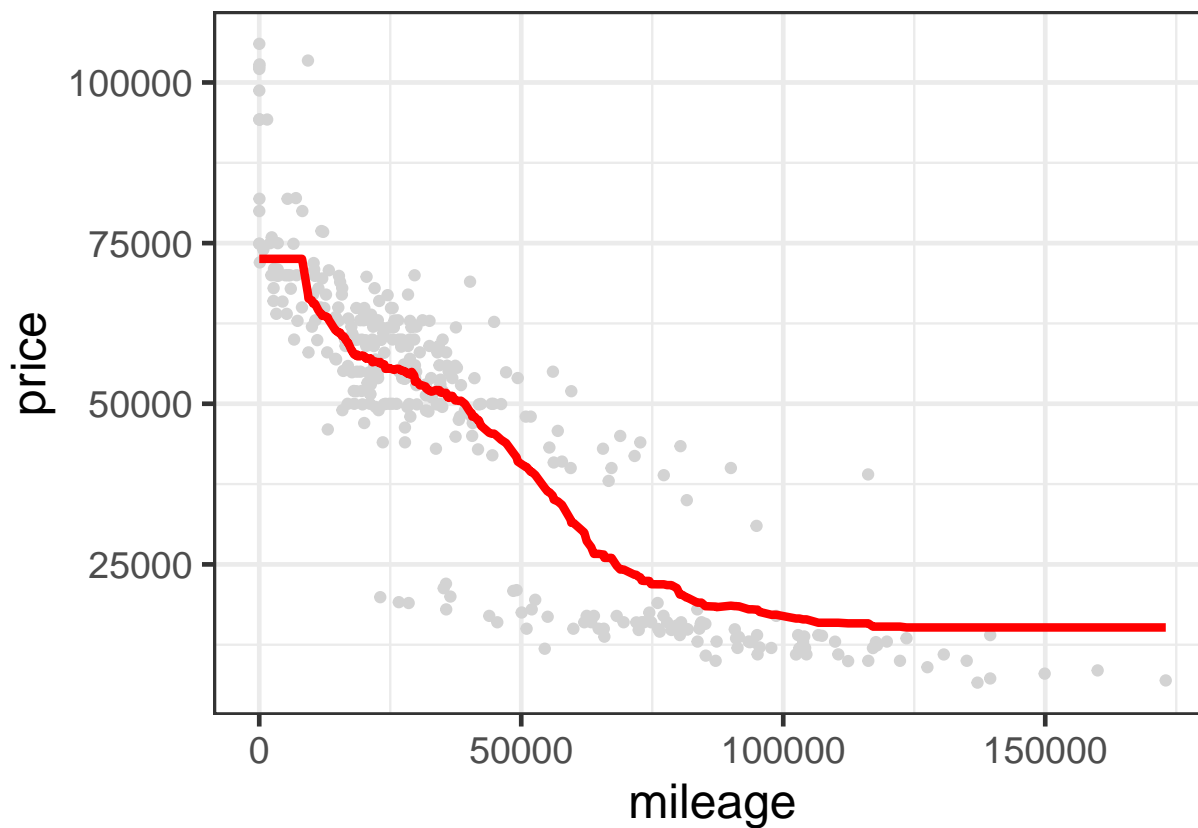


```
p_test + geom_point(aes(x = mileage, y = ypred_knn70), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 70)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 80
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn80 = knn.reg(train = X_train, test = X_test, y = y_train, k=80)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn80 = knn80$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn80)
```

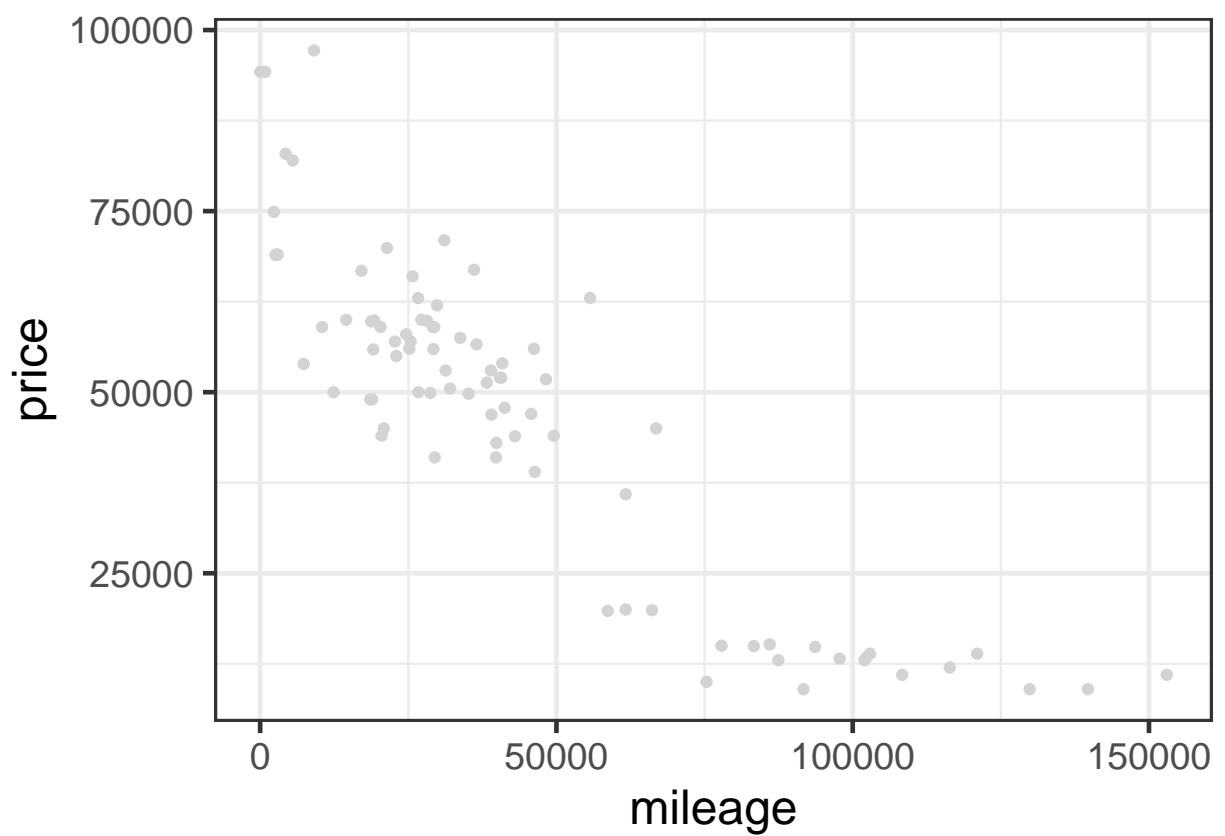
```
## [1] 9213.168
```

```
#attach predictions to data frame
```

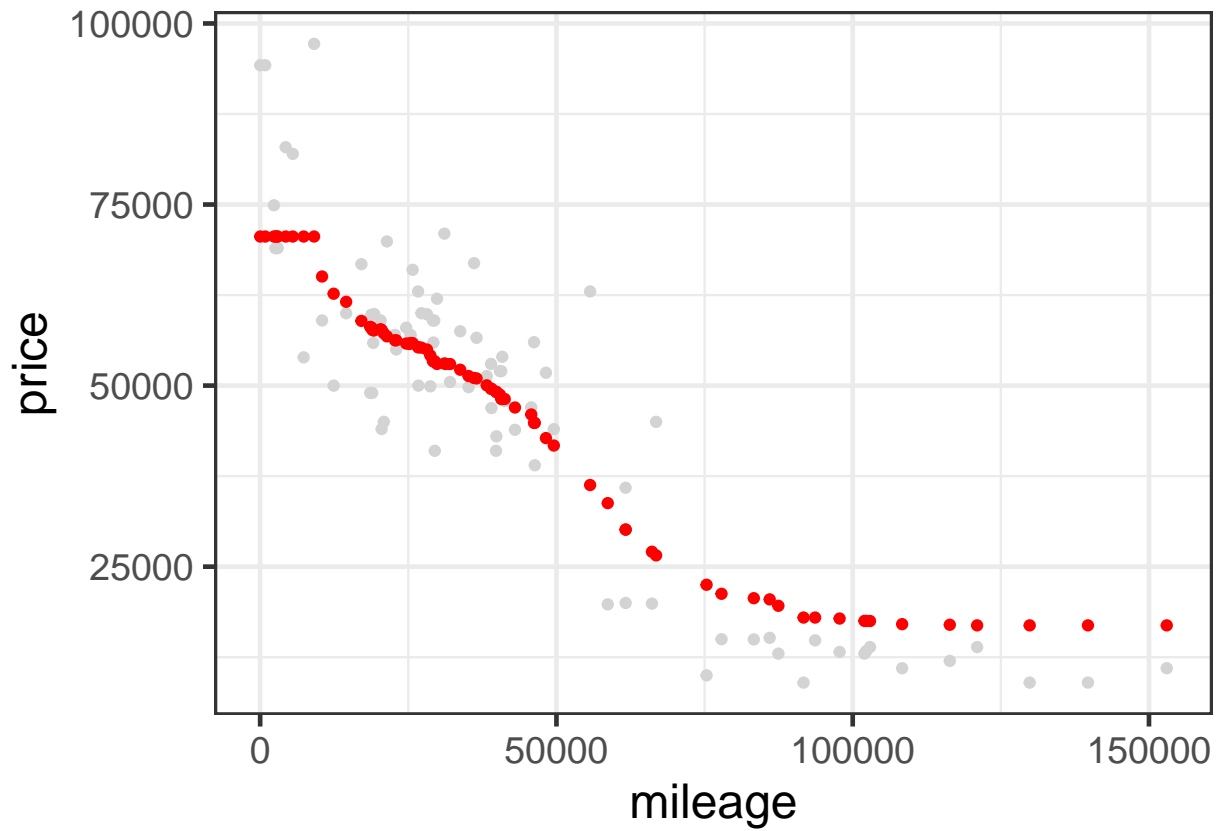
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn80 = ypred_knn80
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

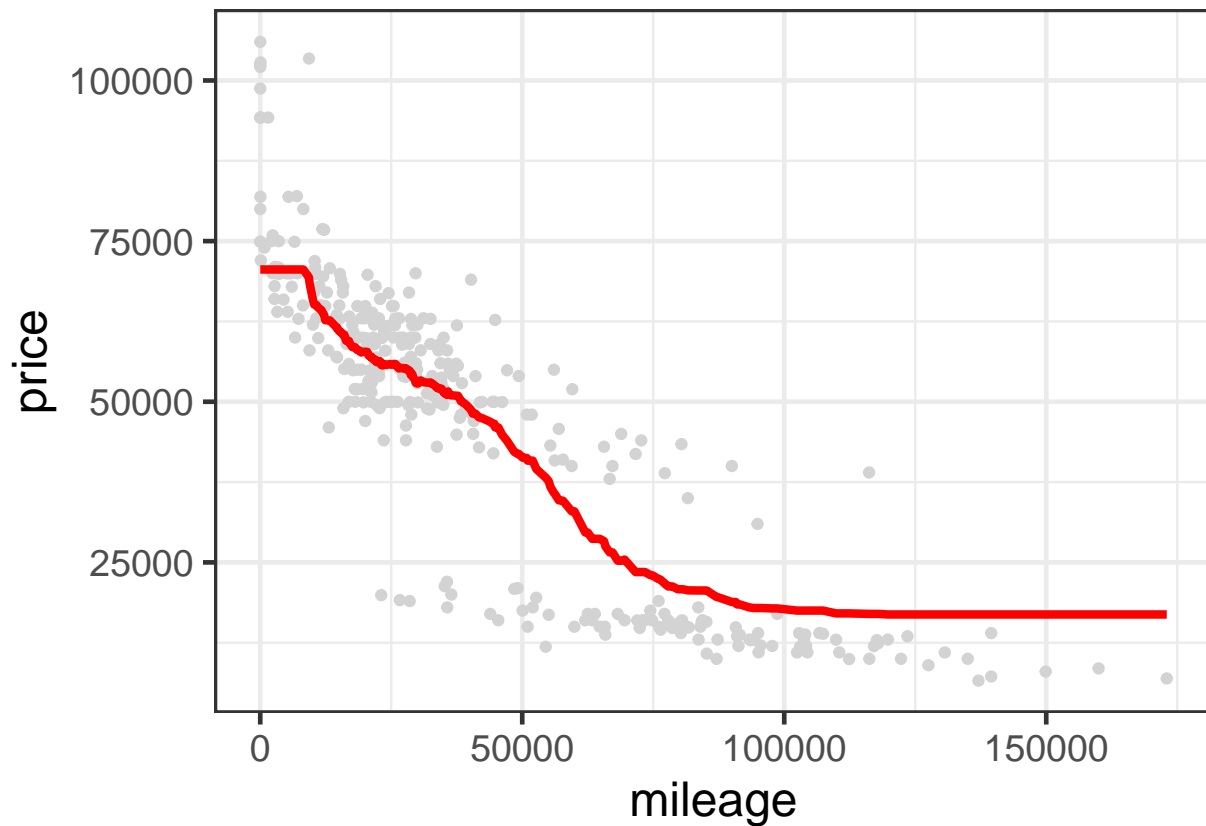


```
p_test + geom_point(aes(x = mileage, y = ypred_knn80), color='red')
```

```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 80)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 90
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn90 = knn.reg(train = X_train, test = X_test, y = y_train, k=90)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn90 = knn90$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn90)
```

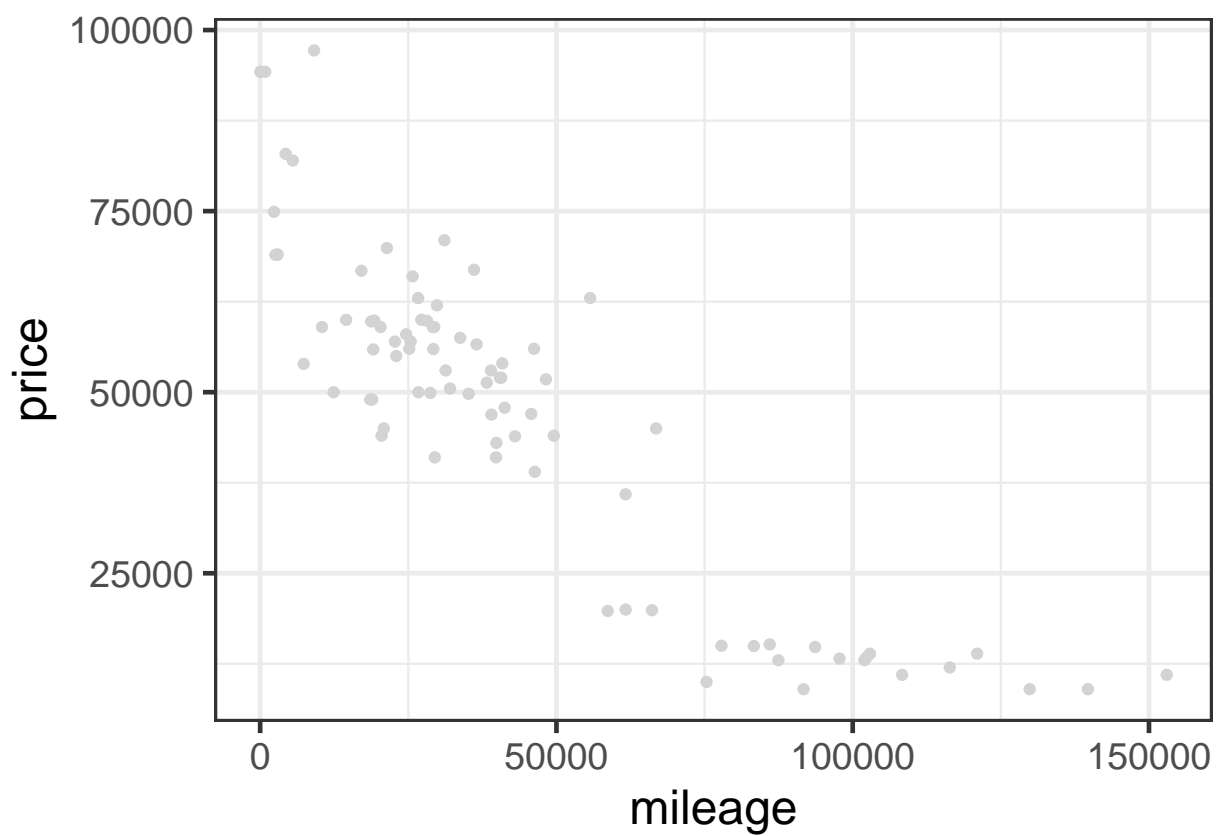
```
## [1] 9441.766
```

```
#attach predictions to data frame
```

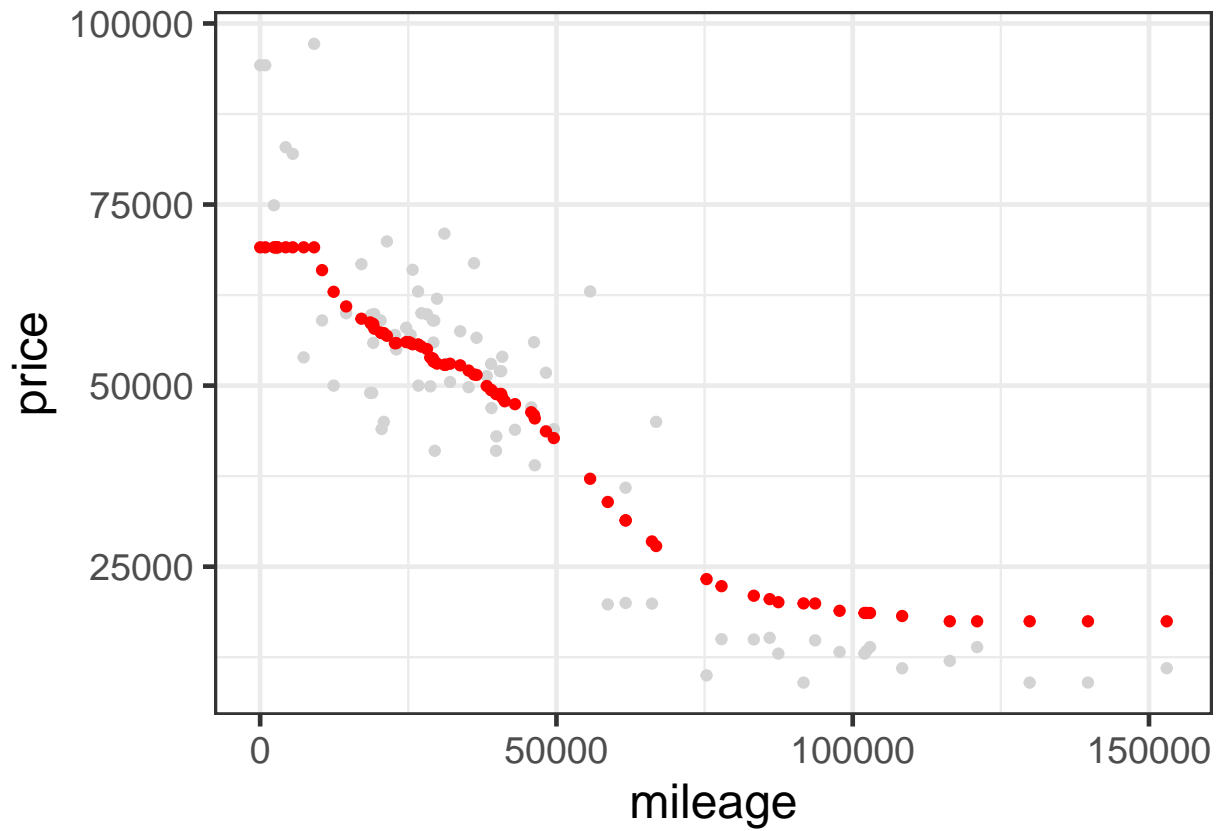
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn90 = ypred_knn90
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

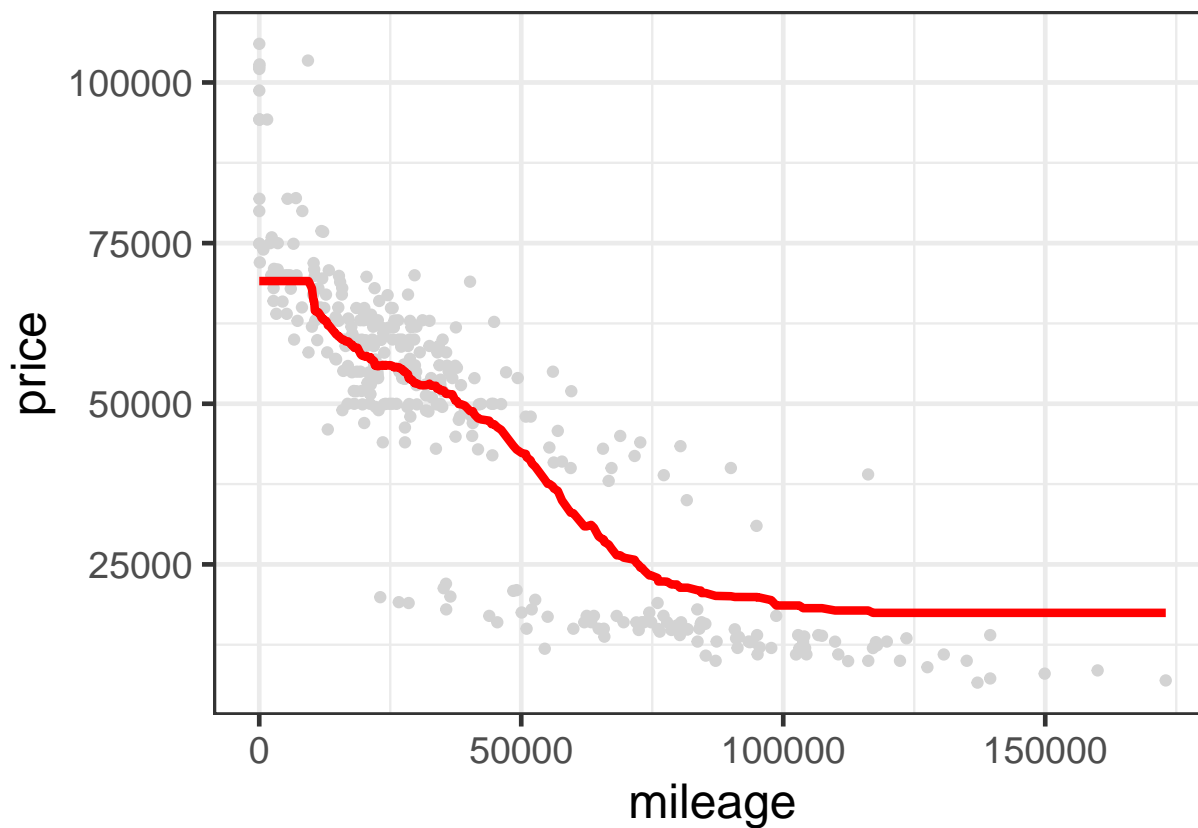


```
p_test + geom_point(aes(x = mileage, y = ypred_knn90), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 90)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K=100
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn100 = knn.reg(train = X_train, test = X_test, y = y_train, k=100)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn100 = knn100$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn100)
```

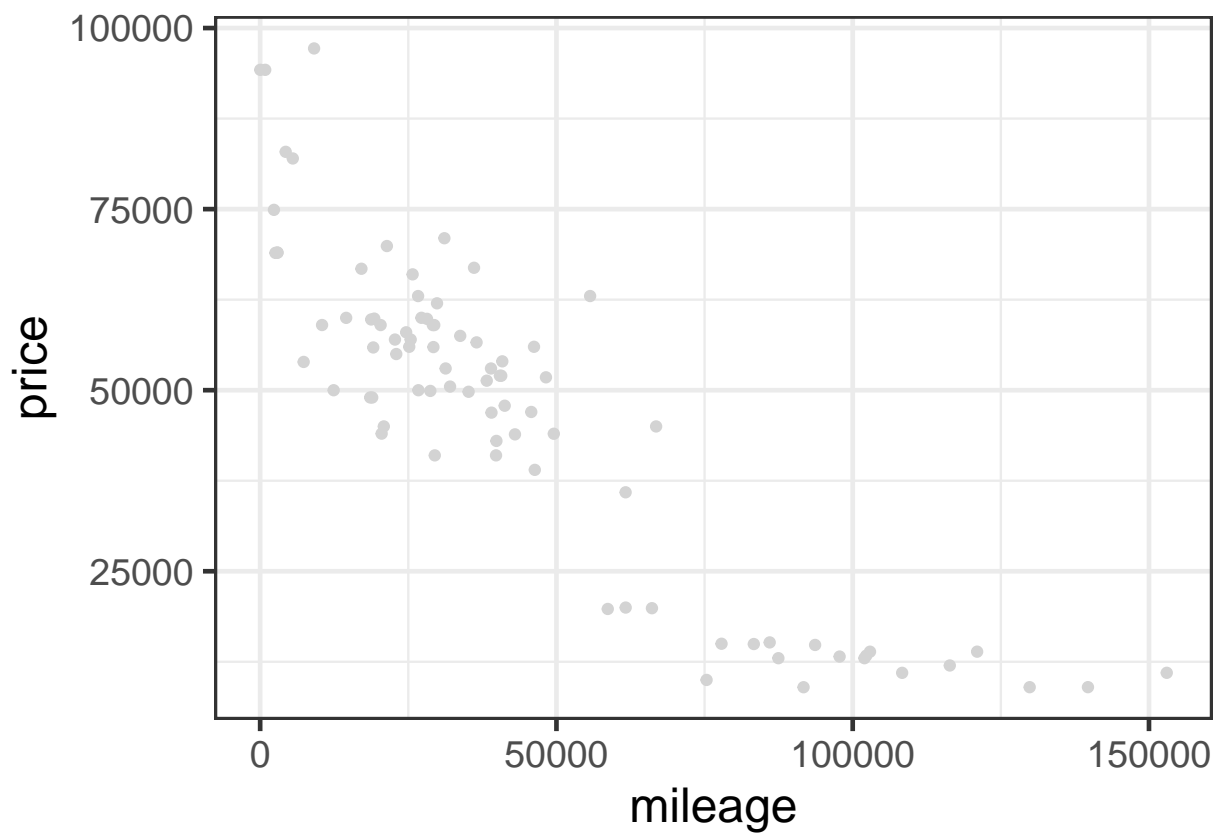
```
## [1] 9745.235
```

```
#attach predictions to data frame
```

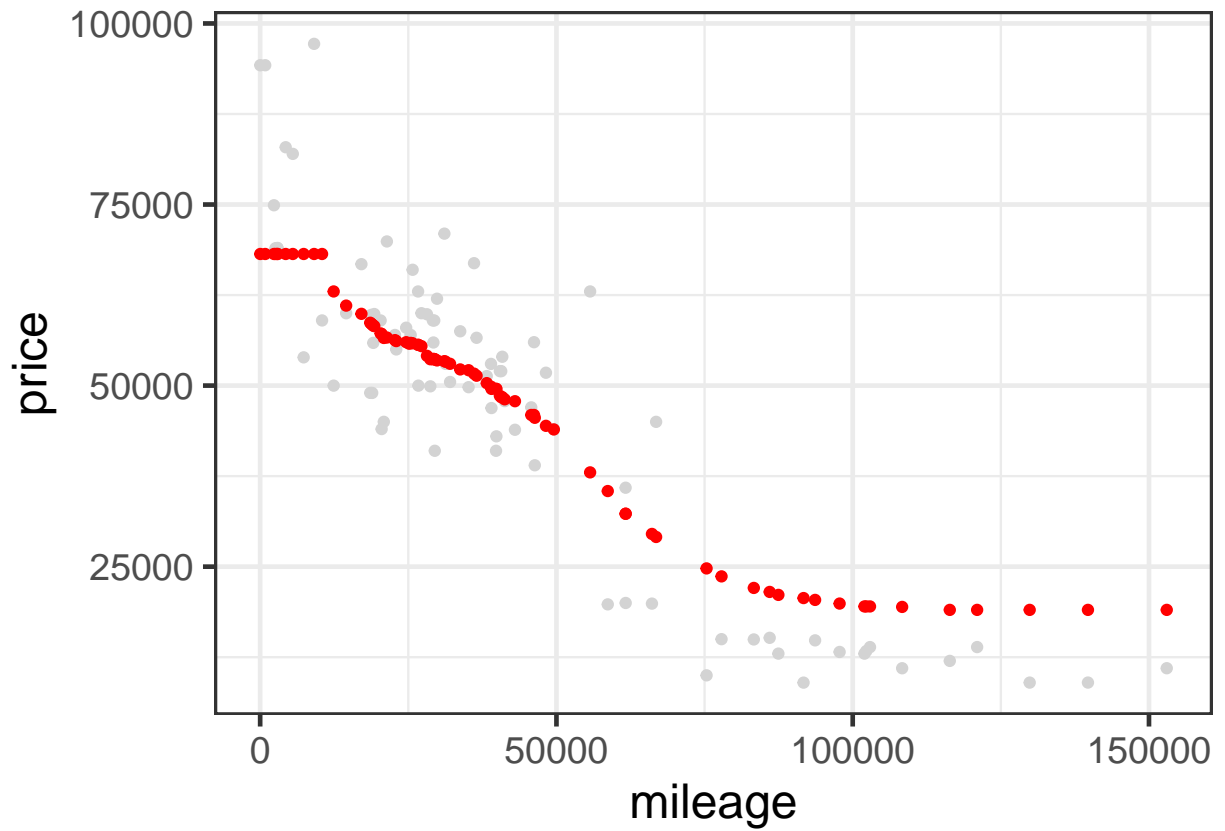
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn100 = ypred_knn100
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

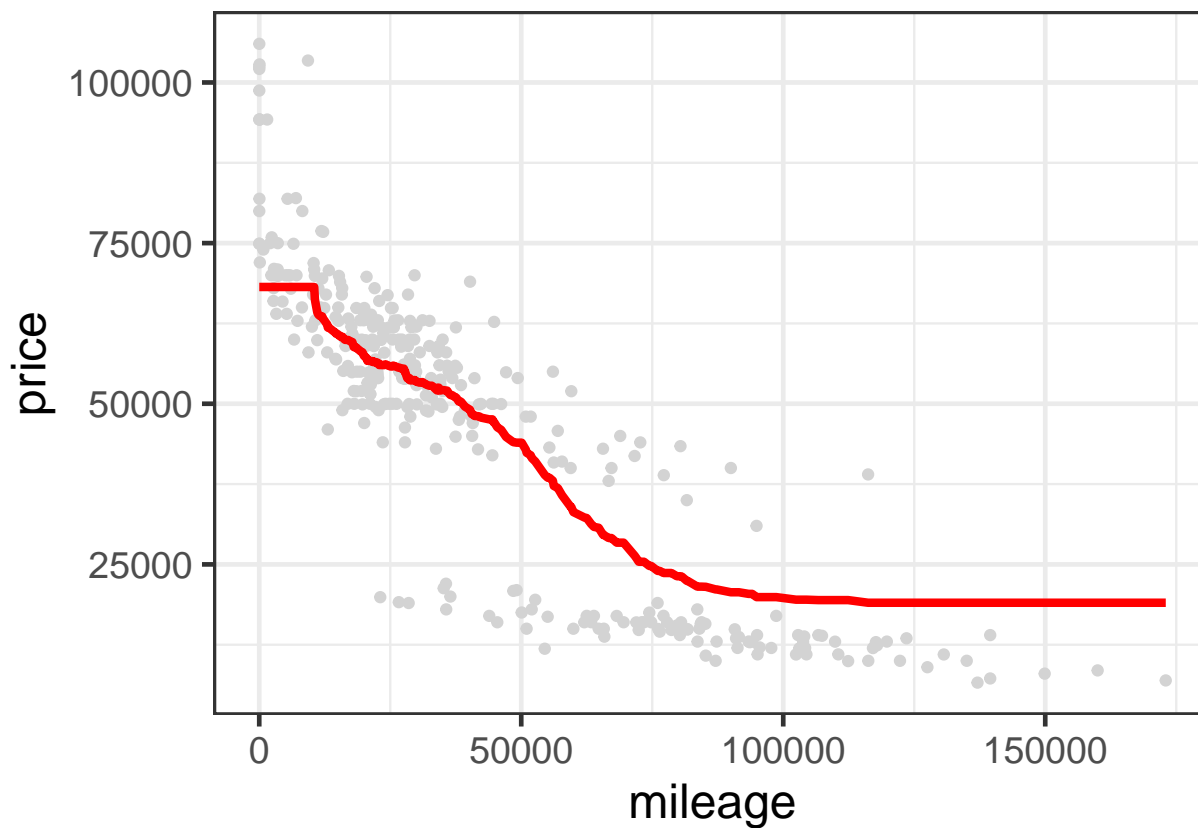


```
p_test + geom_point(aes(x = mileage, y = ypred_knn100), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 100)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 120
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn120 = knn.reg(train = X_train, test = X_test, y = y_train, k=120)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn120 = knn120$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```



```
rmse(y_test, ypred_knn120)
```

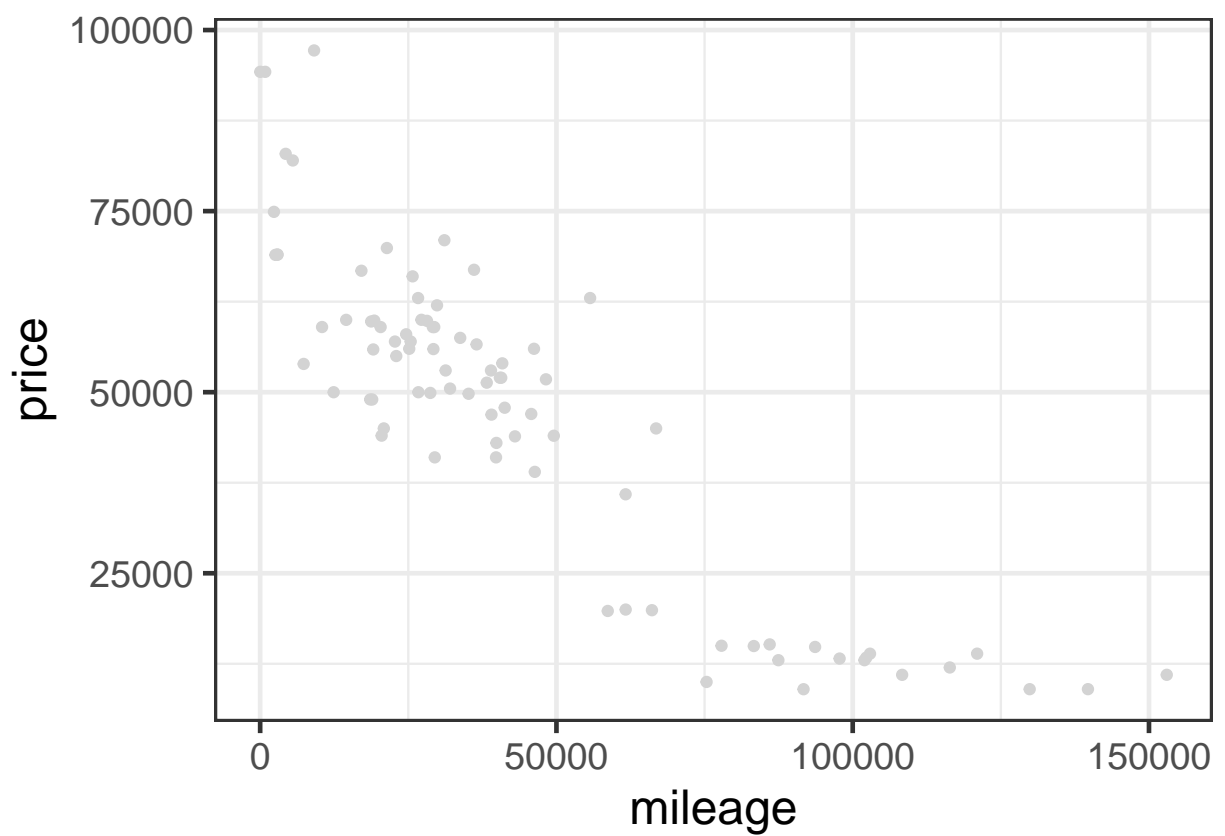
```
## [1] 10637.04
```

```
#attach predictions to data frame
```

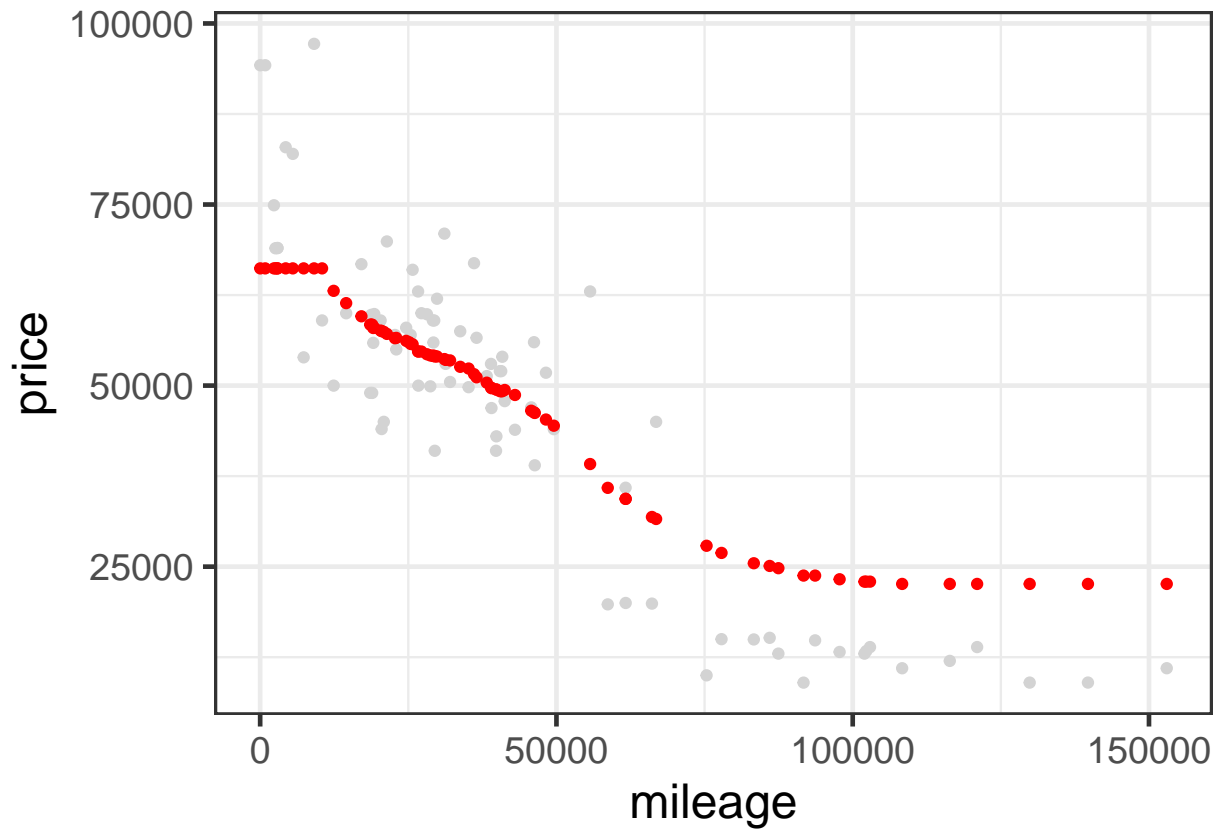
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn120 = ypred_knn120
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

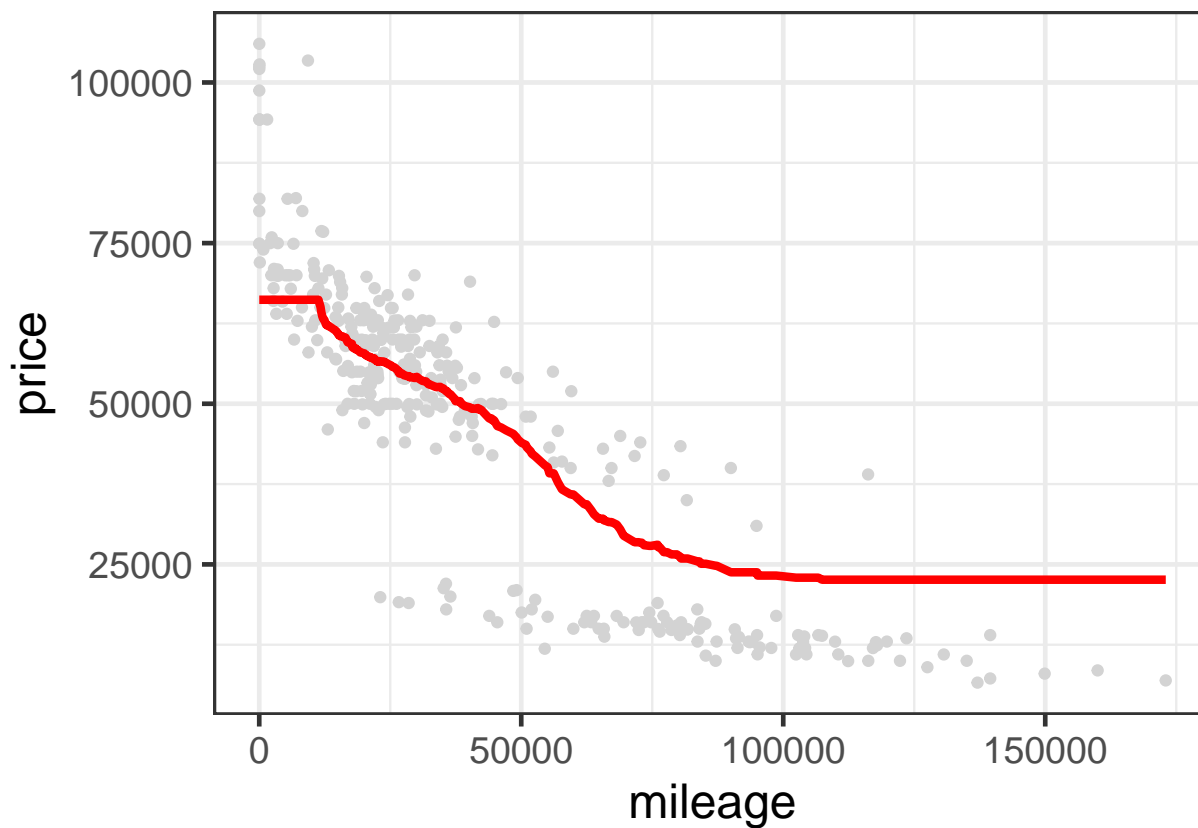


```
p_test + geom_point(aes(x = mileage, y = ypred_knn120), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 120)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn140 = knn.reg(train = X_train, test = X_test, y = y_train, k=140)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn140 = knn140$pred
```

```
rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

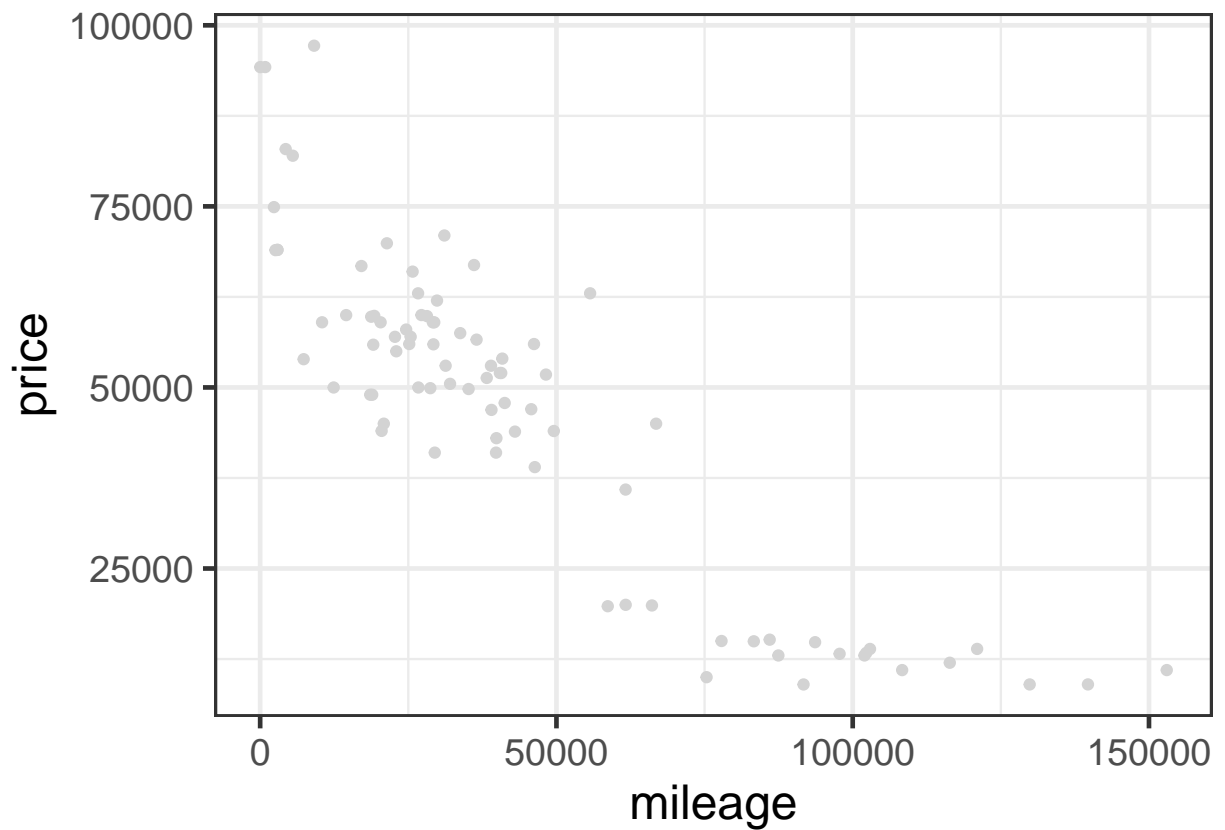
```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn140)
```

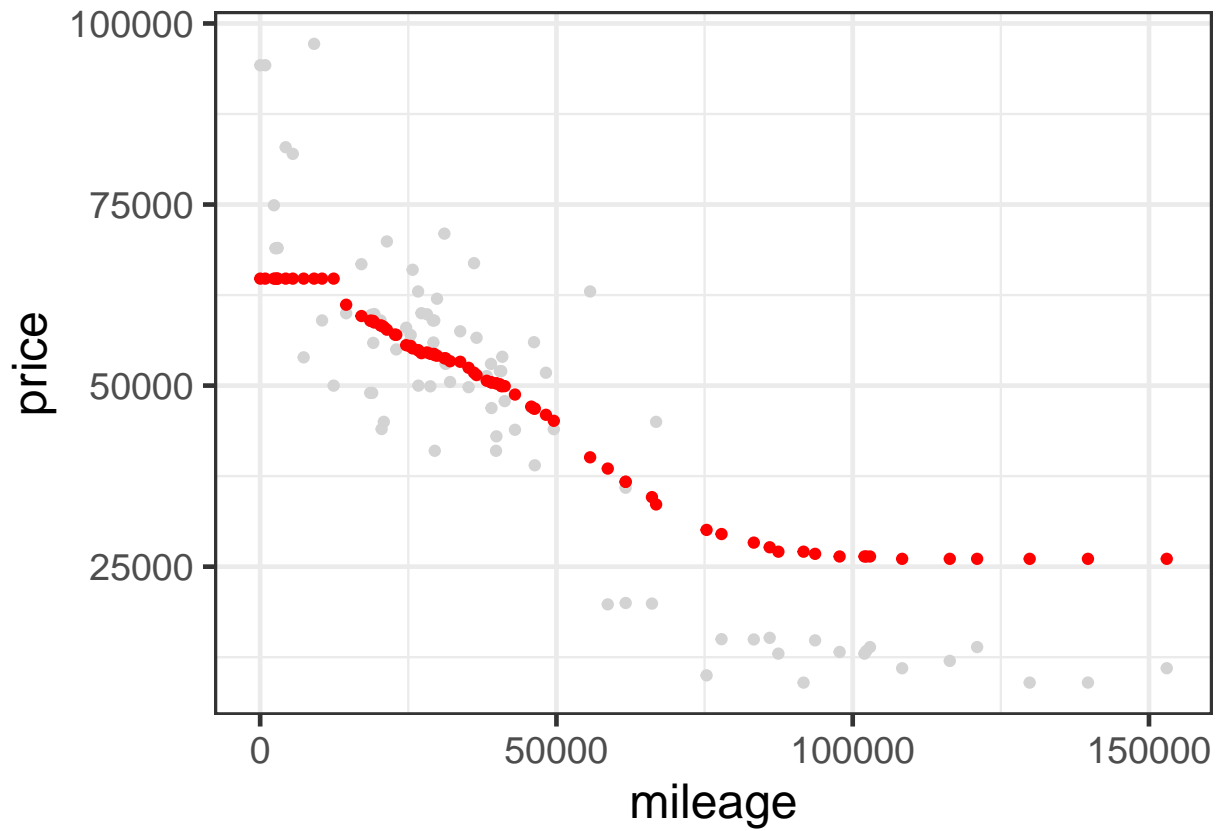
```
## [1] 11684.14
```

```
#attach predictions to data frame  
D_test$ypred_lm2 = ypred_lm2  
D_test$ypred_knn140 = ypred_knn140
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

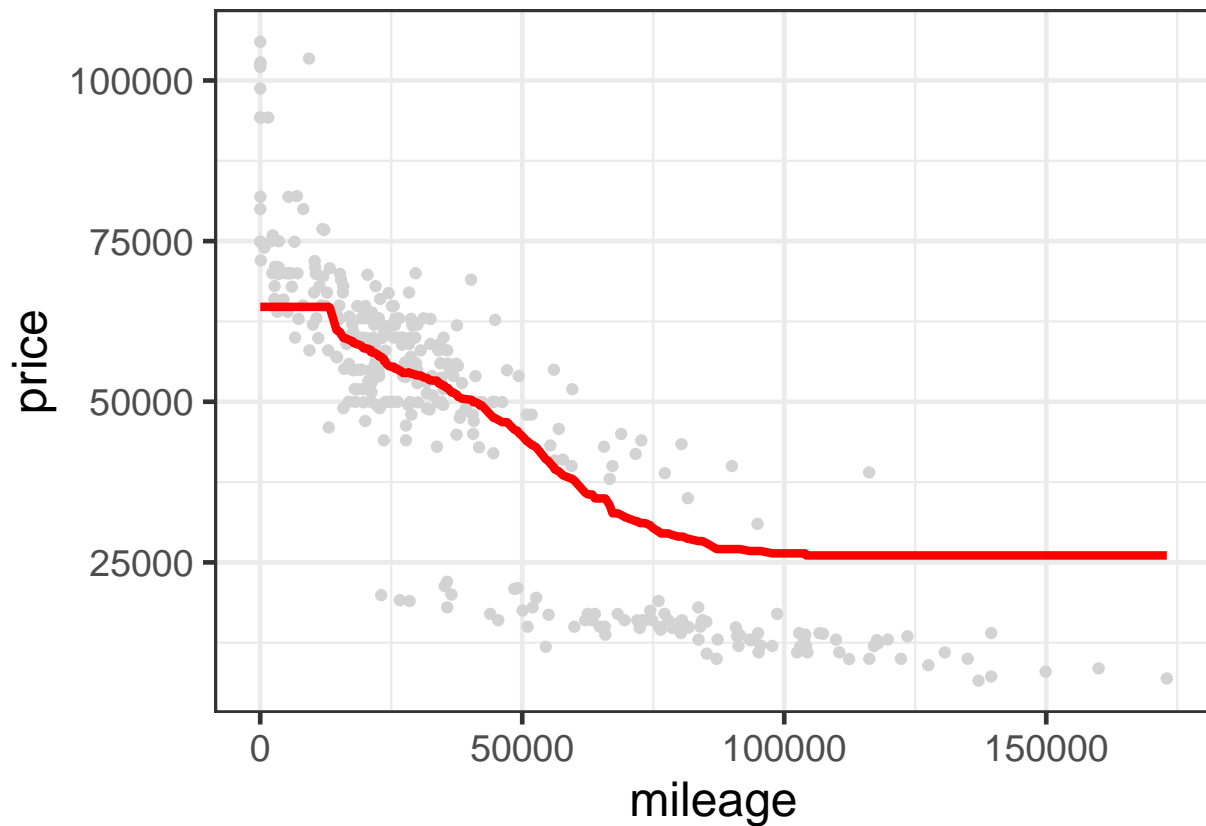


```
p_test + geom_point(aes(x = mileage, y = ypred_knn140), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 140)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn150 = knn.reg(train = X_train, test = X_test, y = y_train, k=150)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn150 = knn150$pred
```

```
rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

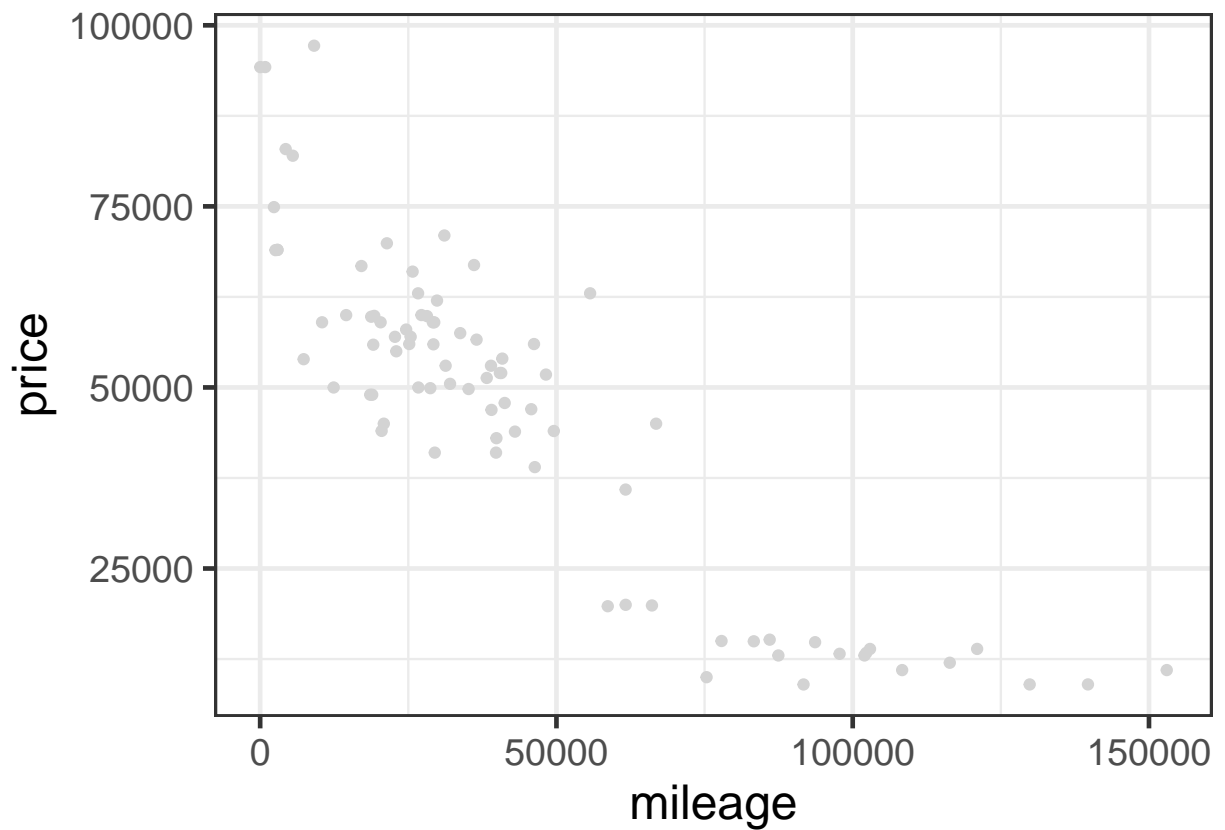
```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn150)
```

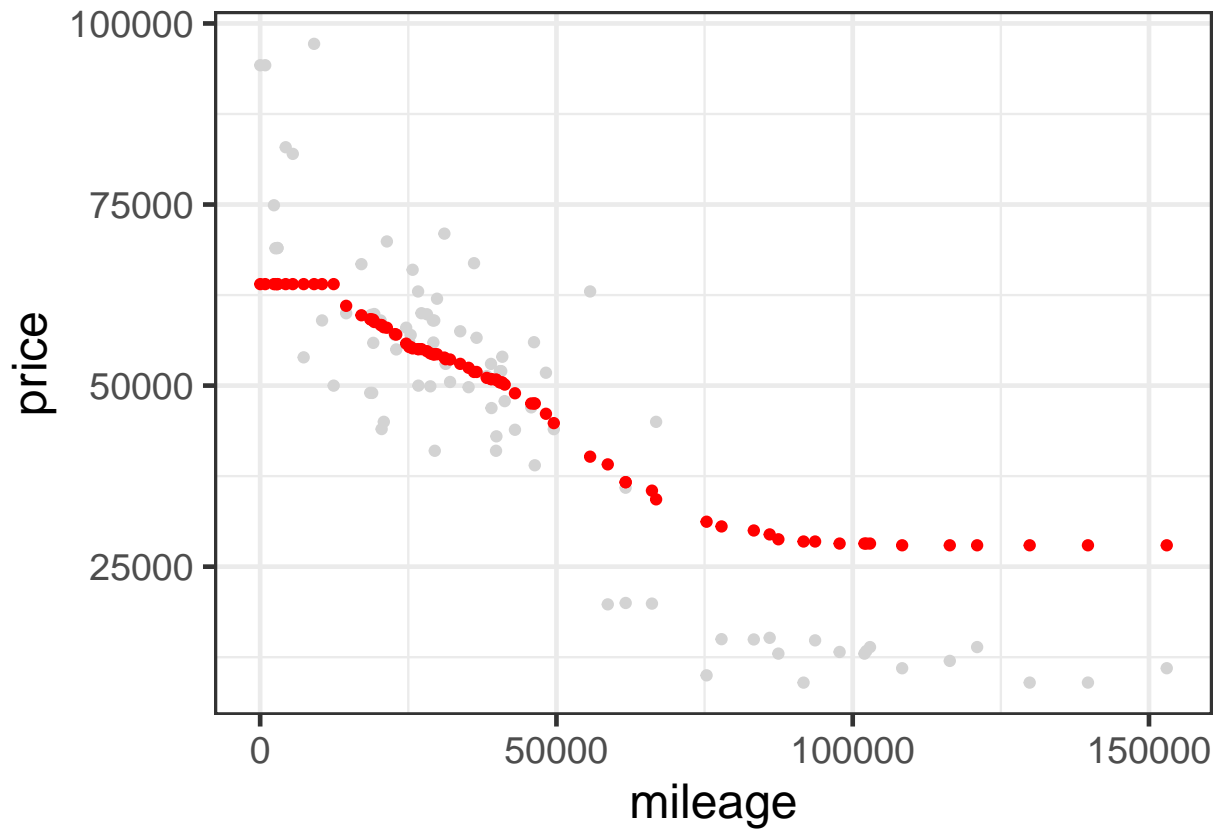
```
## [1] 12227.51
```

```
#attach predictions to data frame  
D_test$ypred_lm2 = ypred_lm2  
D_test$ypred_knn150 = ypred_knn150
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

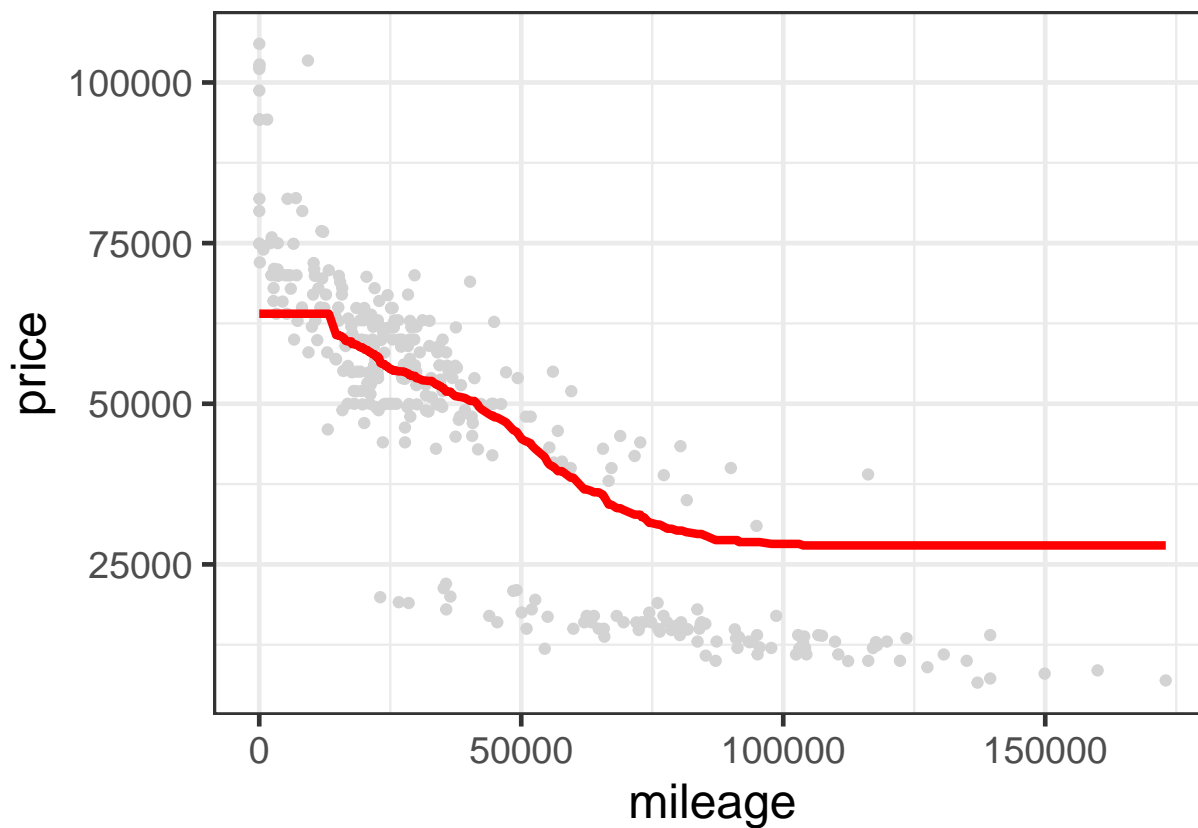


```
p_test + geom_point(aes(x = mileage, y = ypred_knn150), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 150)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```

```
#K=175
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn175 = knn.reg(train = X_train, test = X_test, y = y_train, k=175)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn175 = knn175$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn175)
```

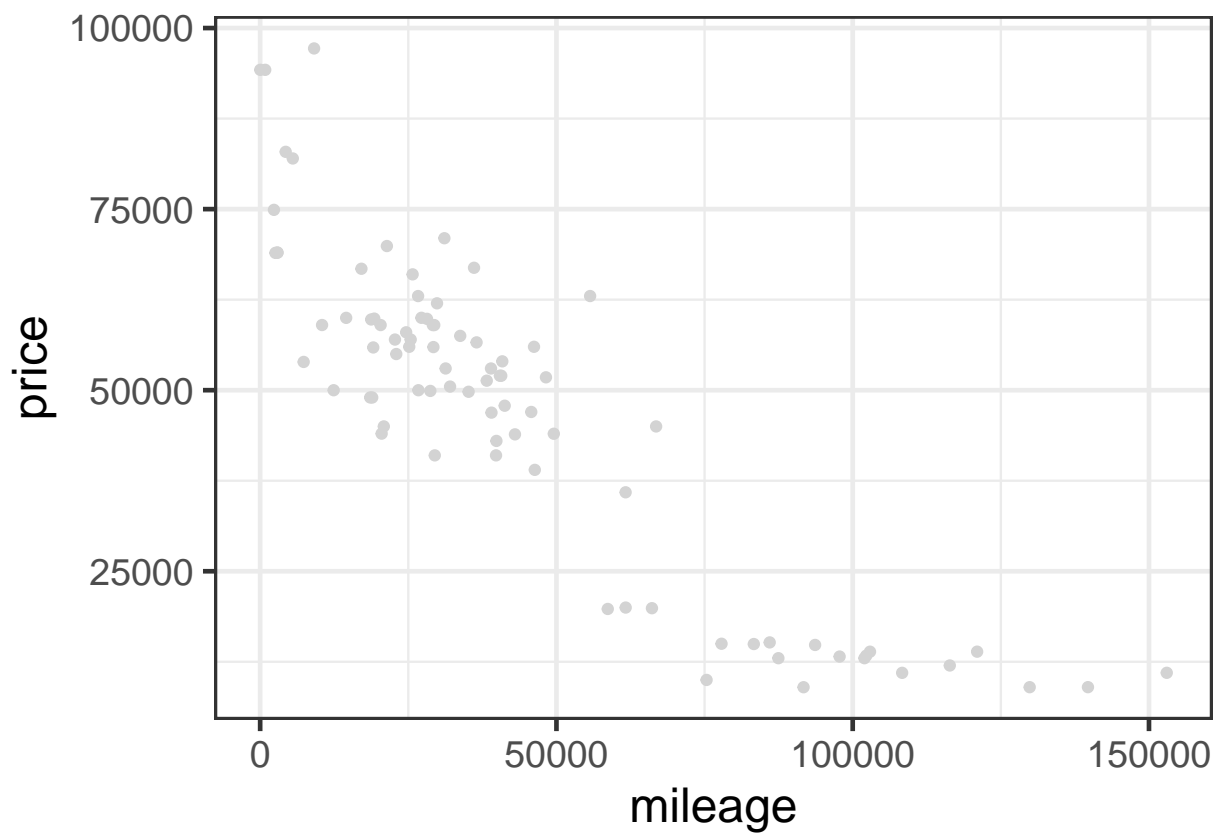
```
## [1] 13454.04
```

```
#attach predictions to data frame
```

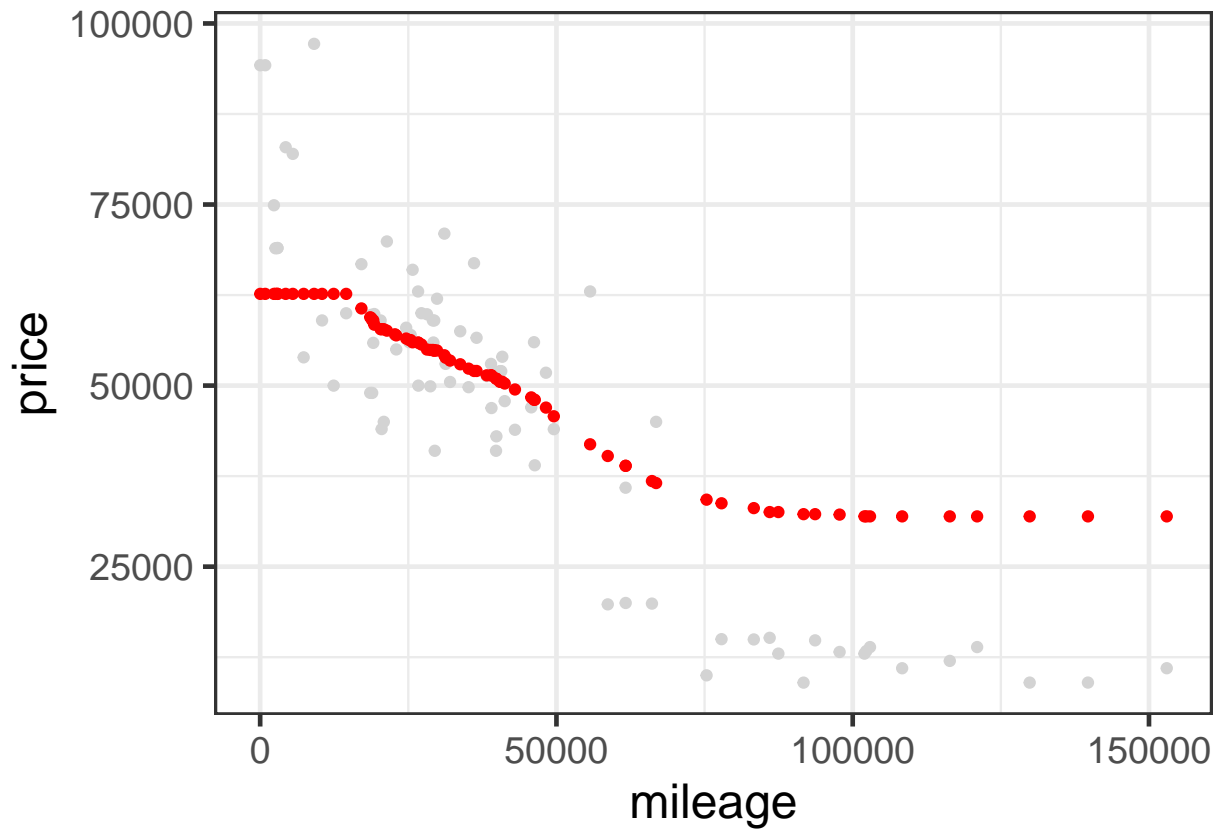
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn175 = ypred_knn175
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

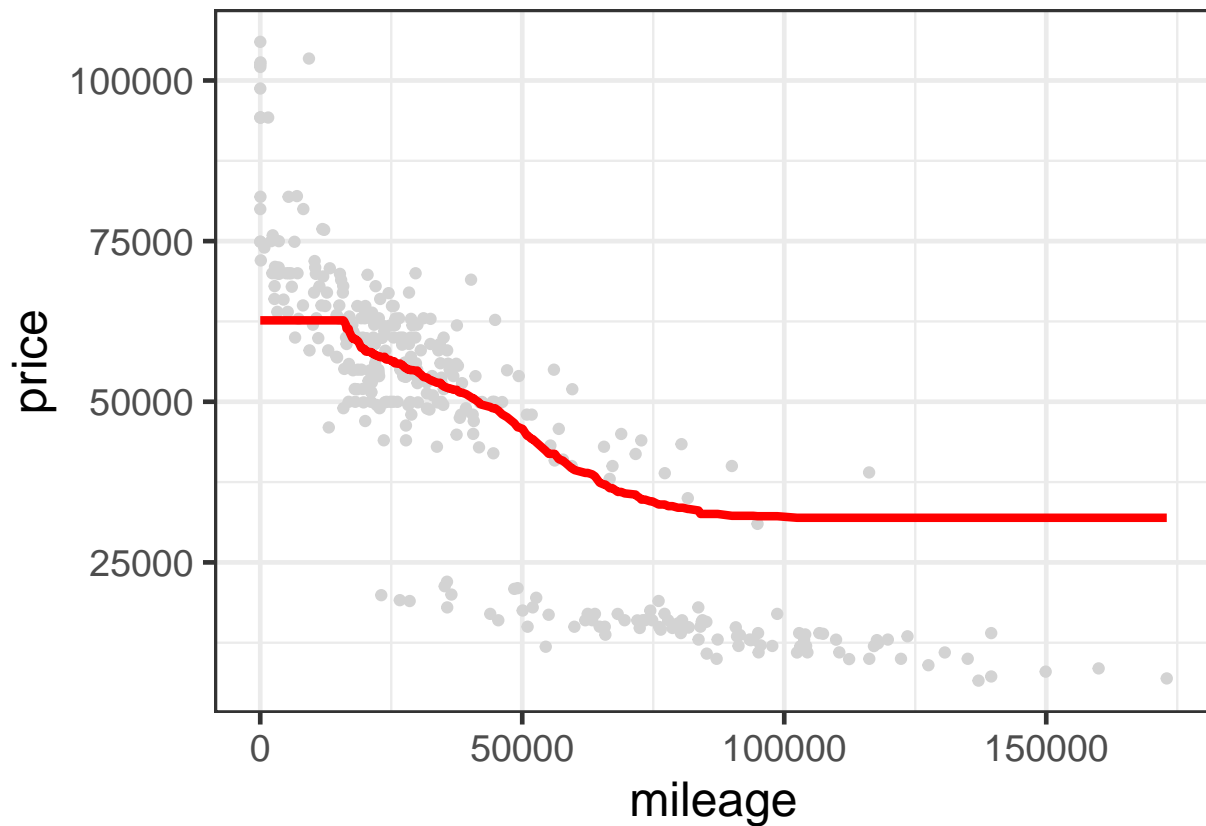


```
p_test + geom_point(aes(x = mileage, y = ypred_knn175), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 175)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K=190
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn190 = knn.reg(train = X_train, test = X_test, y = y_train, k=190)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn190 = knn190$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn190)
```

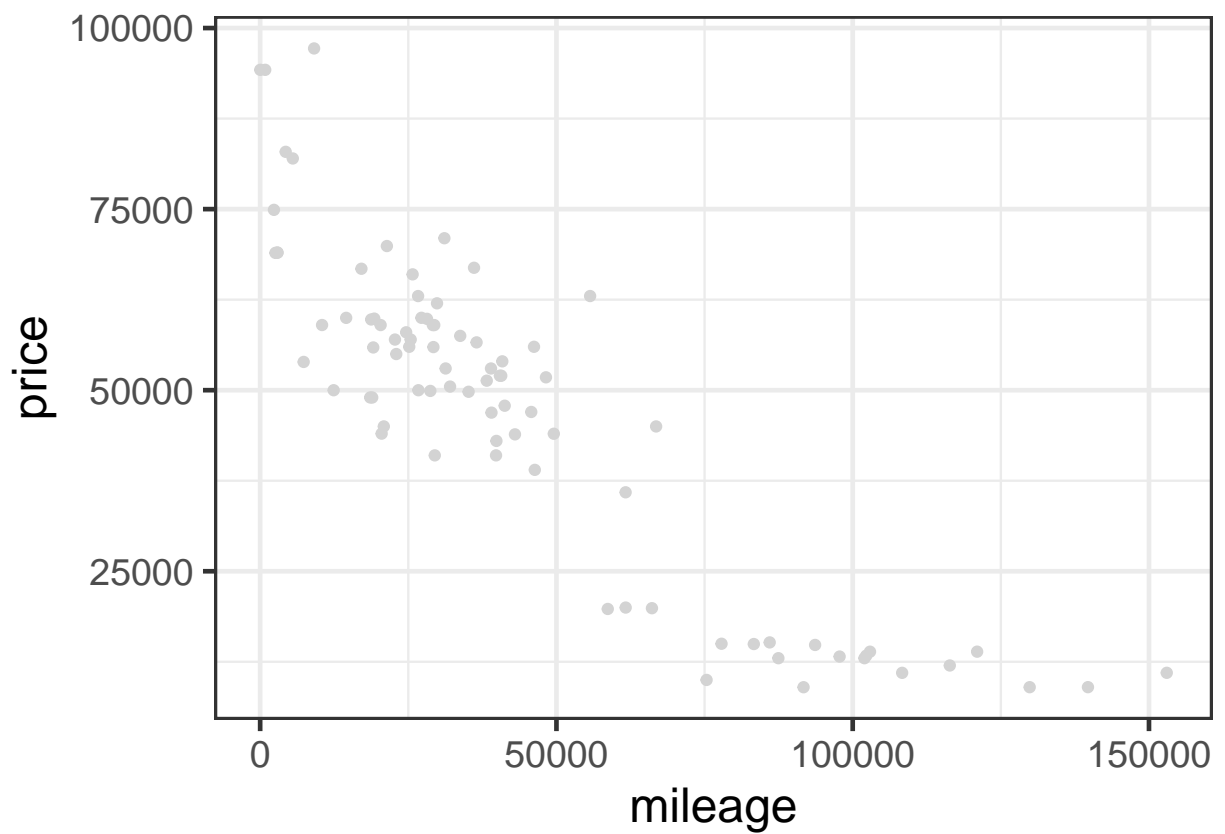
```
## [1] 14036.02
```

```
#attach predictions to data frame
```

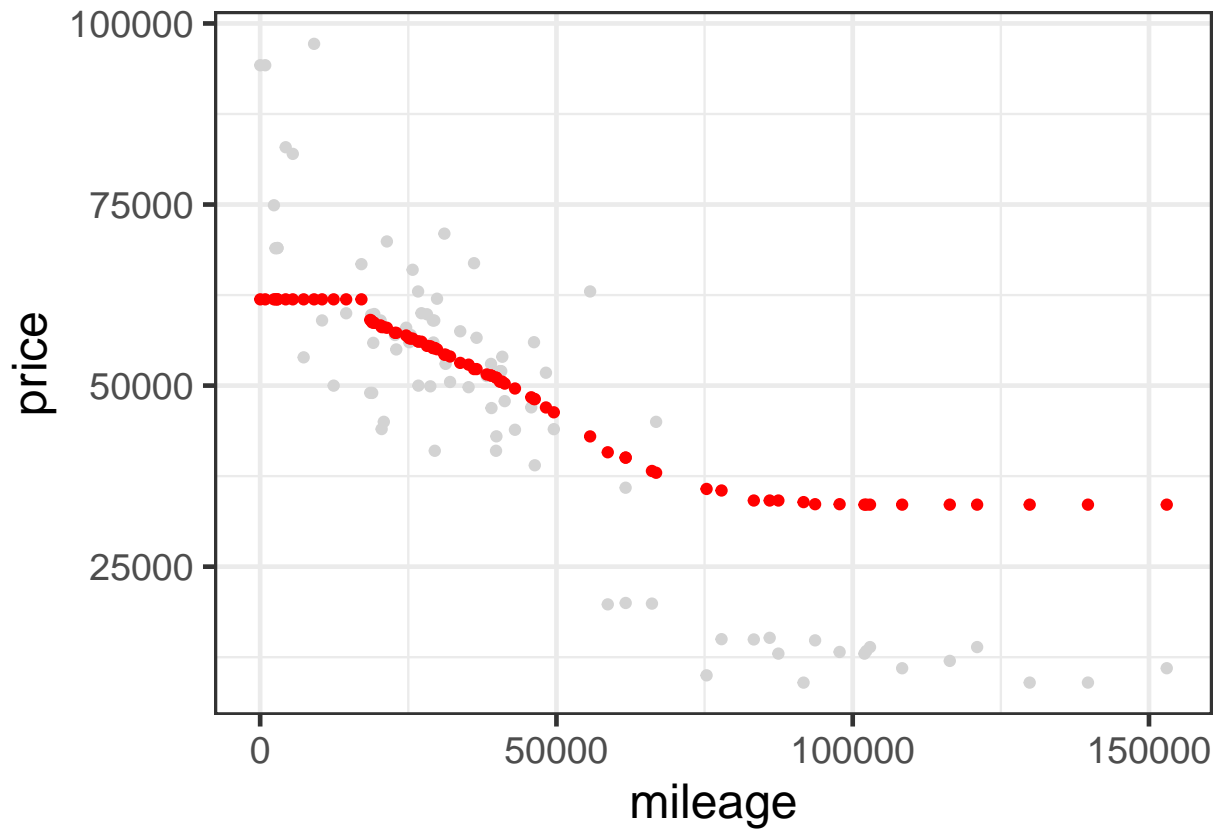
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn190 = ypred_knn190
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

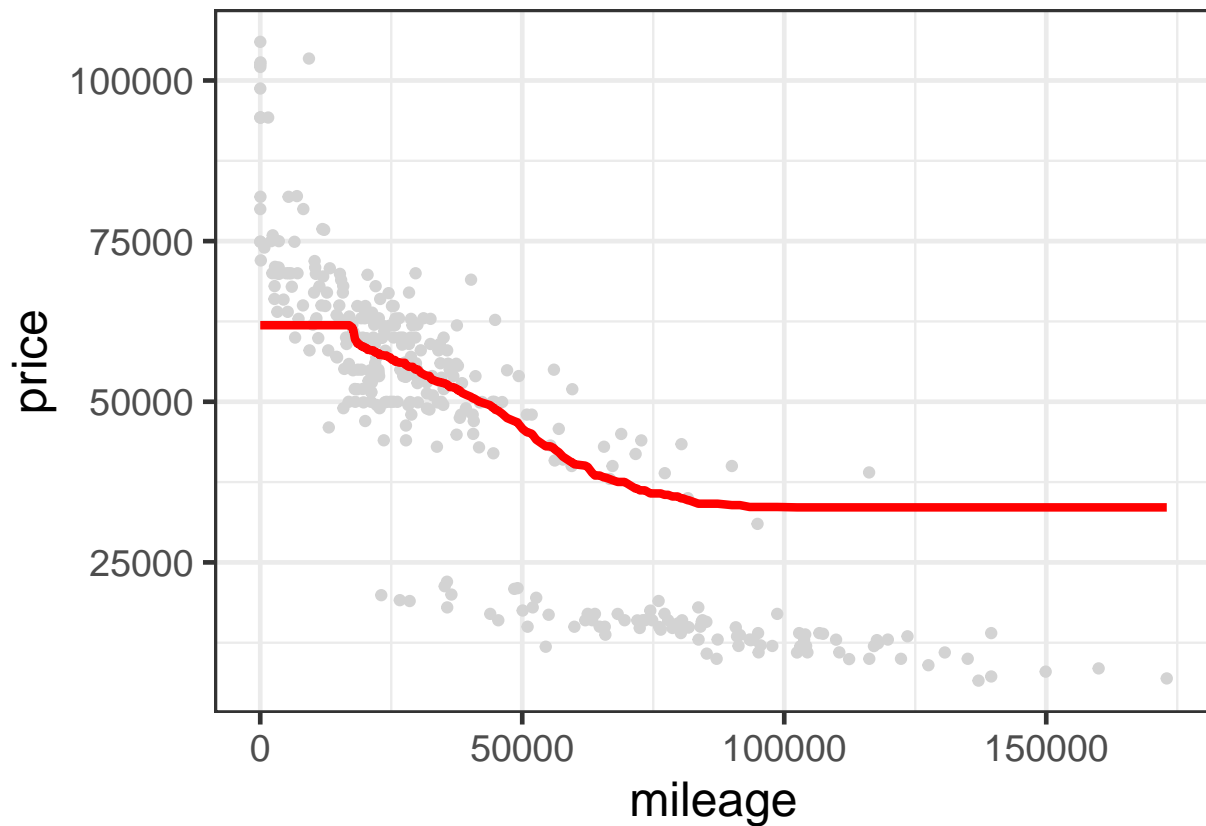


```
p_test + geom_point(aes(x = mileage, y = ypred_knn190), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 190)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K=210
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn210 = knn.reg(train = X_train, test = X_test, y = y_train, k=210)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn210 = knn210$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn210)
```

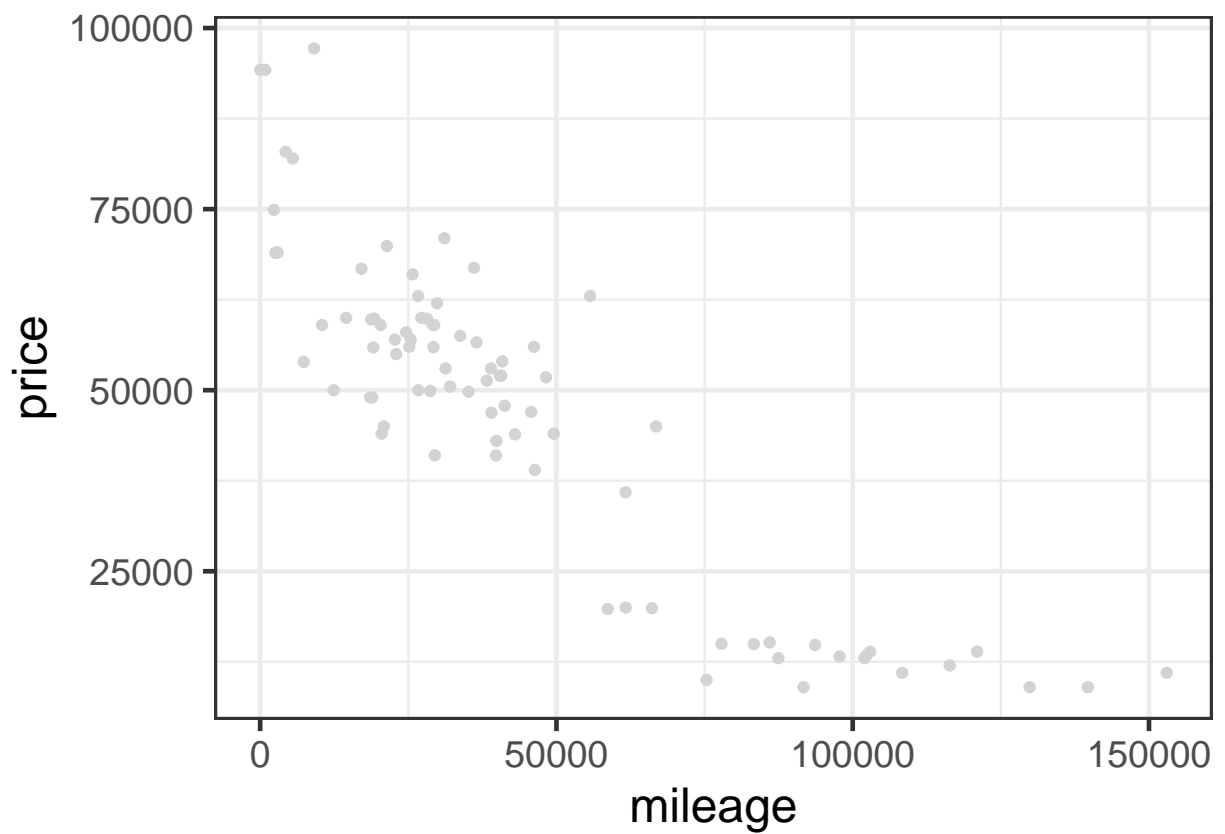
```
## [1] 14965.83
```

```
#attach predictions to data frame
```

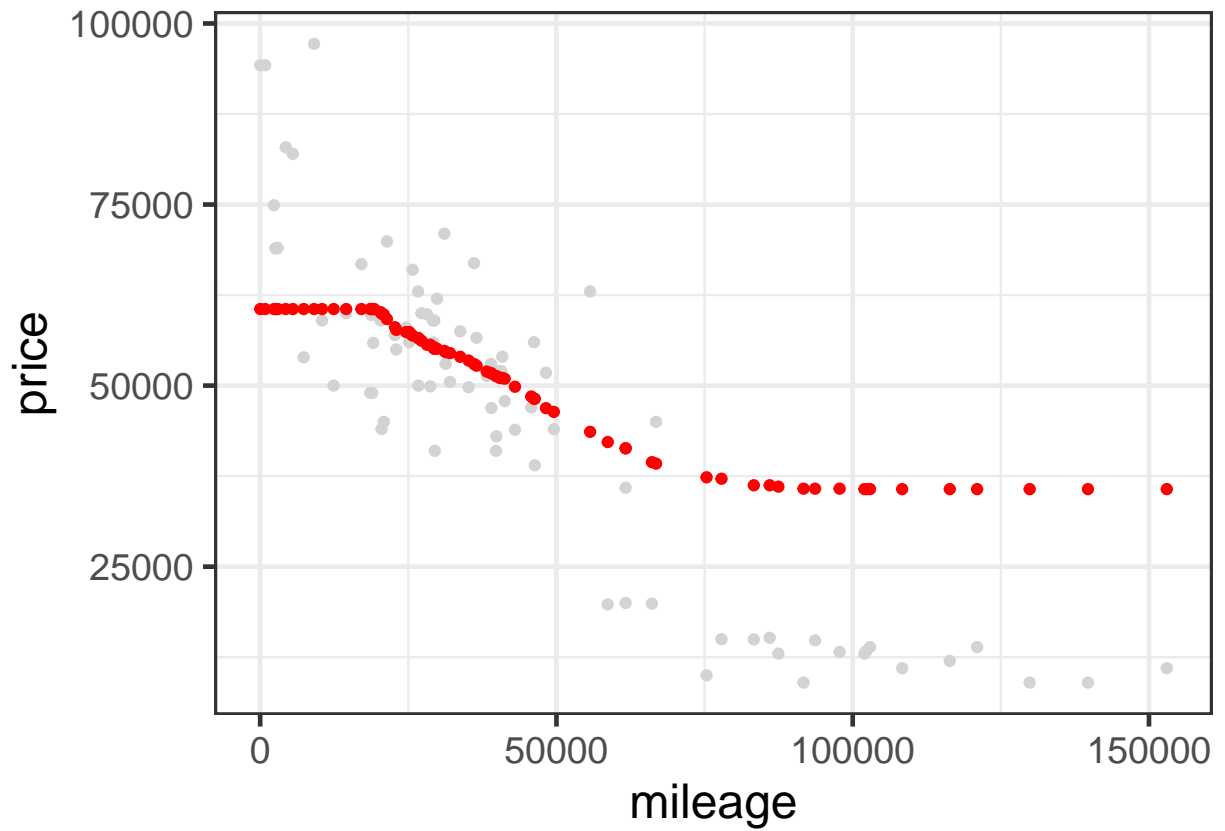
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn210 = ypred_knn210
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

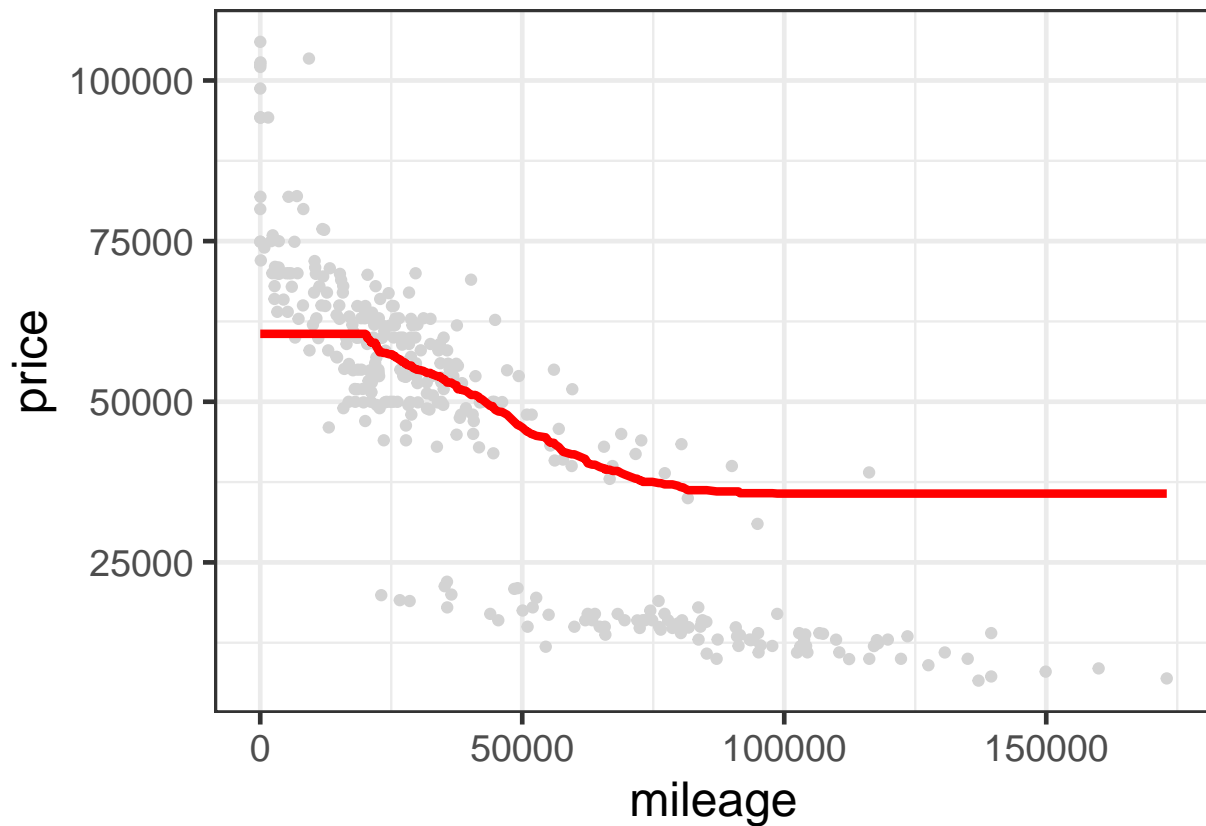


```
p_test + geom_point(aes(x = mileage, y = ypred_knn210), color='red')
```

```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 210)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K=230
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn230 = knn.reg(train = X_train, test = X_test, y = y_train, k=230)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn230 = knn230$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn230)
```

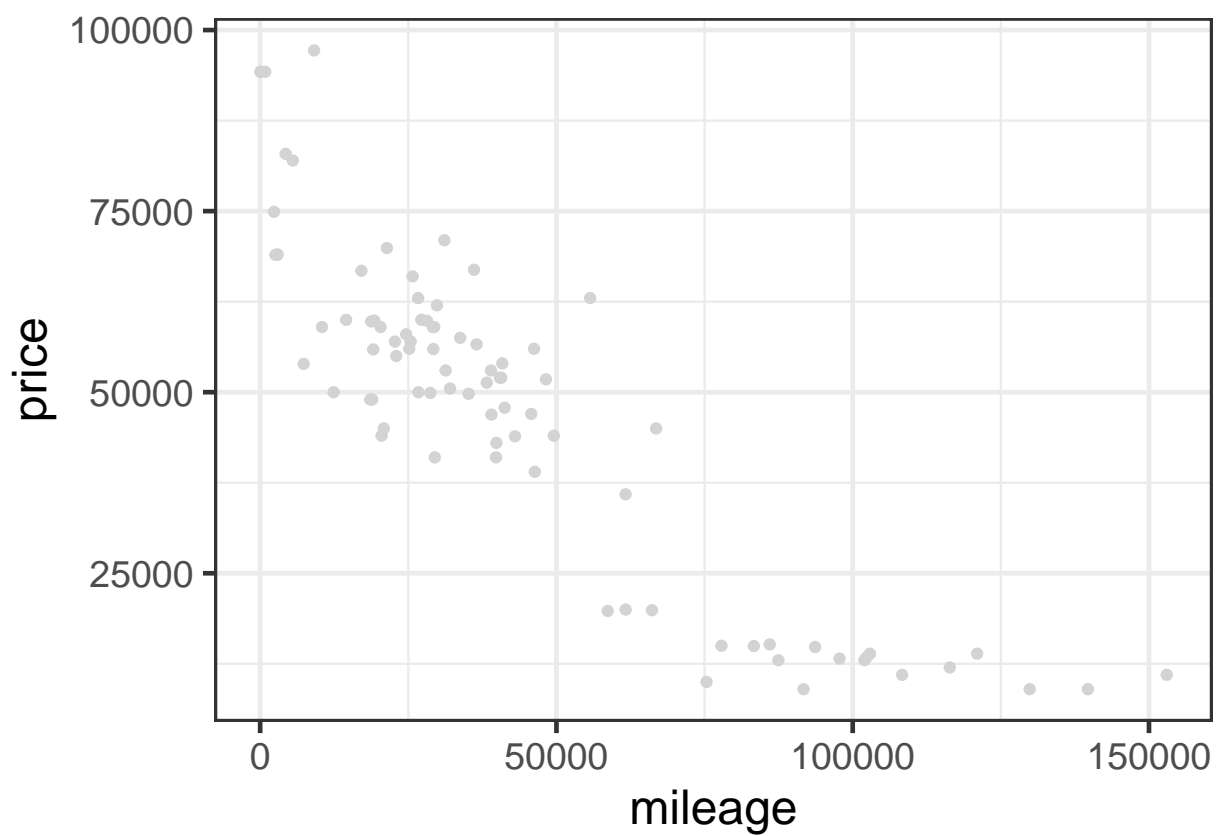
```
## [1] 15735.03
```

```
#attach predictions to data frame
```

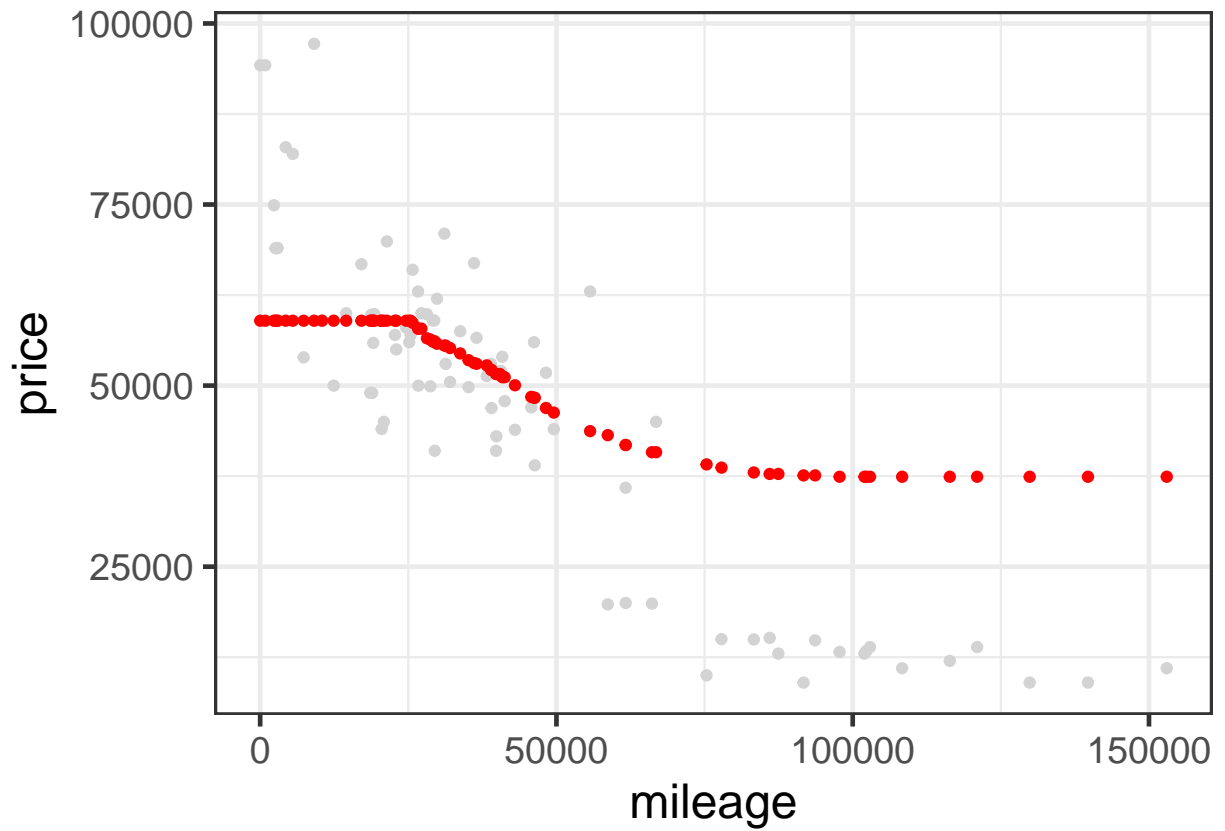
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn230 = ypred_knn230
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

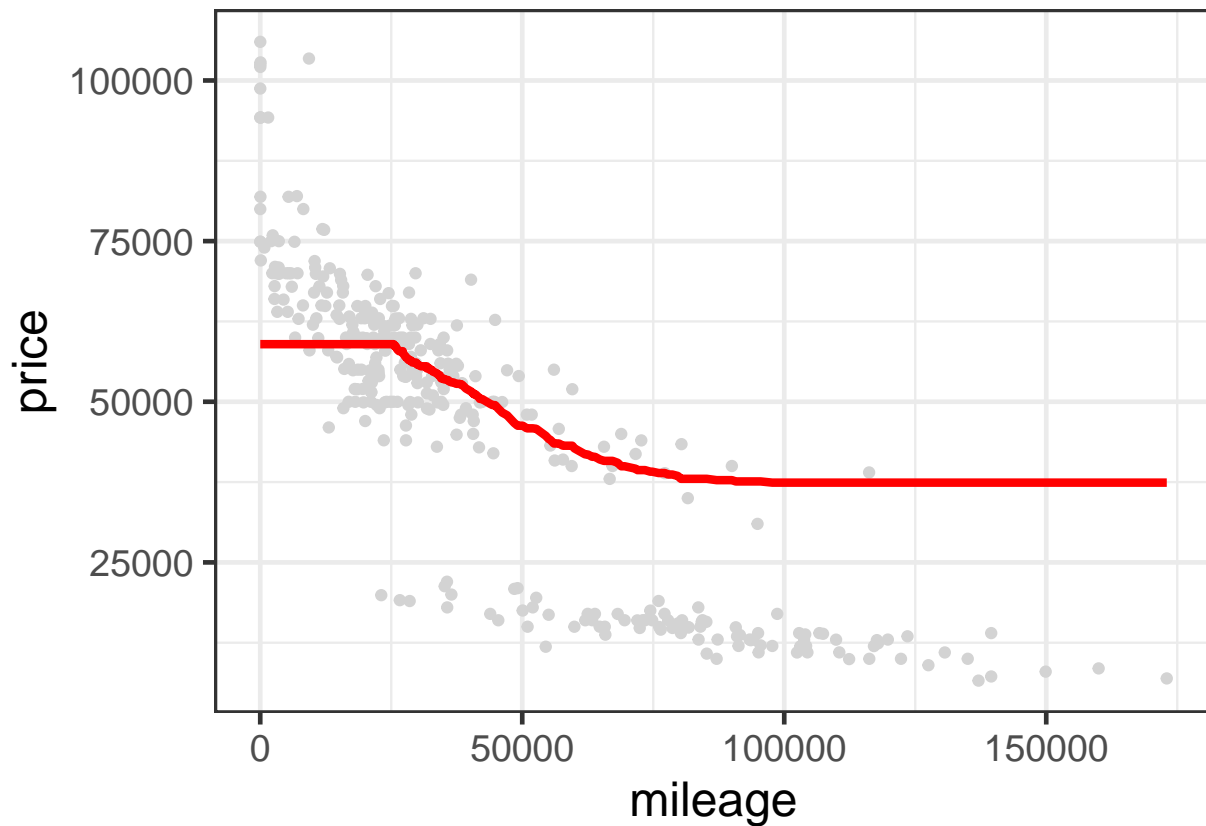


```
p_test + geom_point(aes(x = mileage, y = ypred_knn230), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 230)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K=250
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn250 = knn.reg(train = X_train, test = X_test, y = y_train, k=250)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn250 = knn250$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```

```
rmse(y_test, ypred_knn250)
```

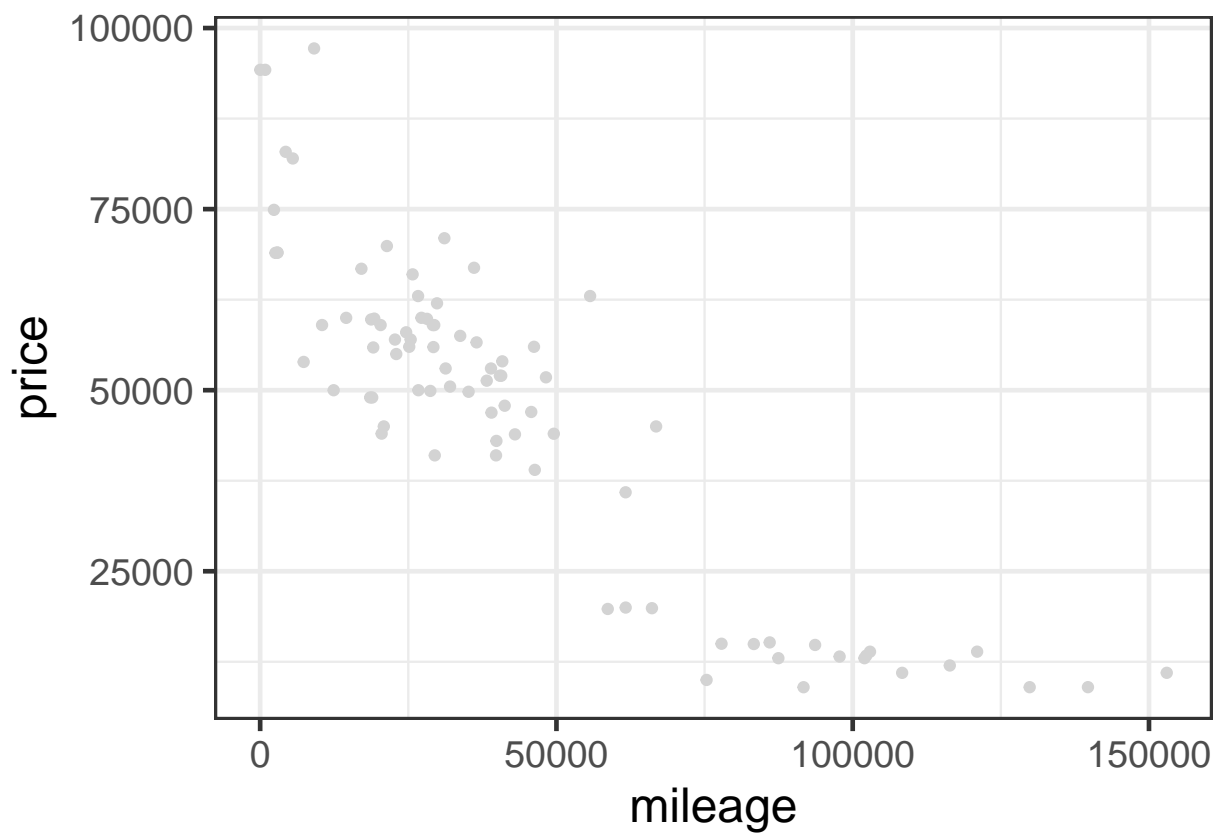
```
## [1] 16615.06
```

```
#attach predictions to data frame
```

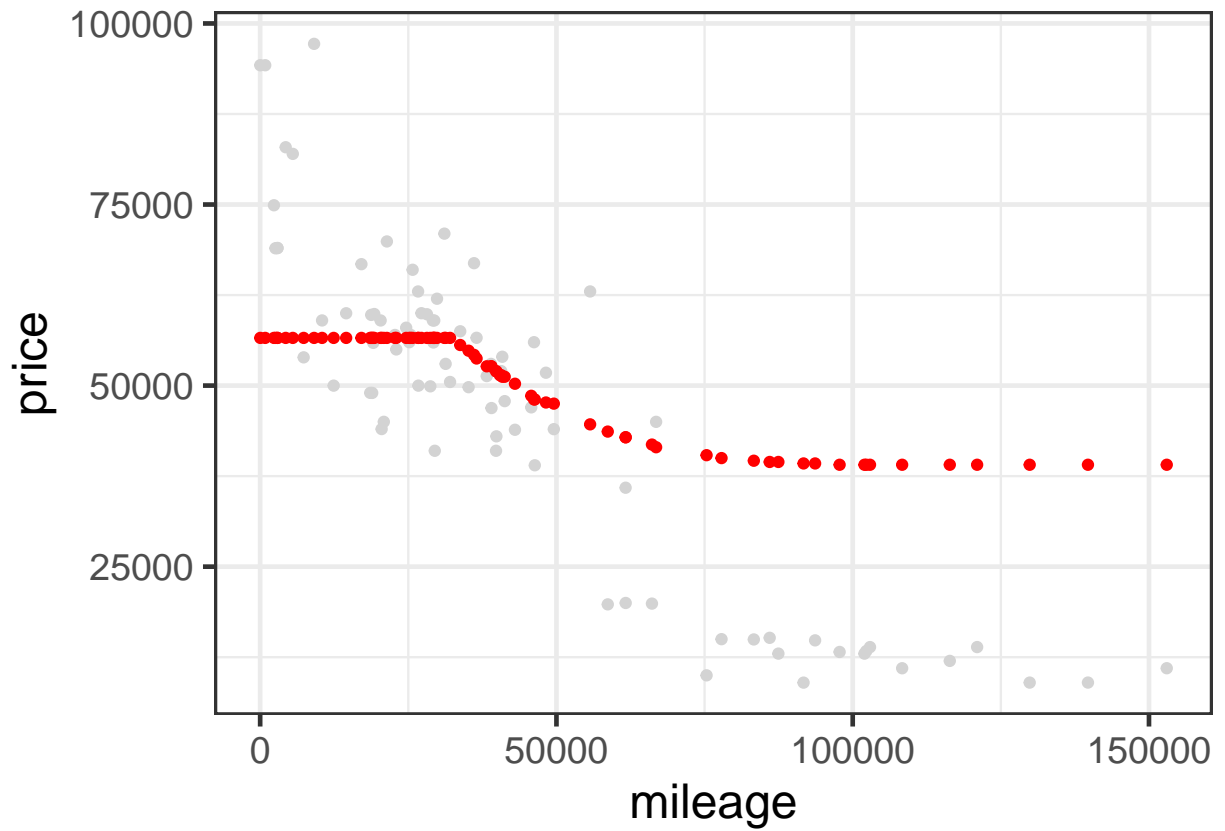
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn250 = ypred_knn250
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

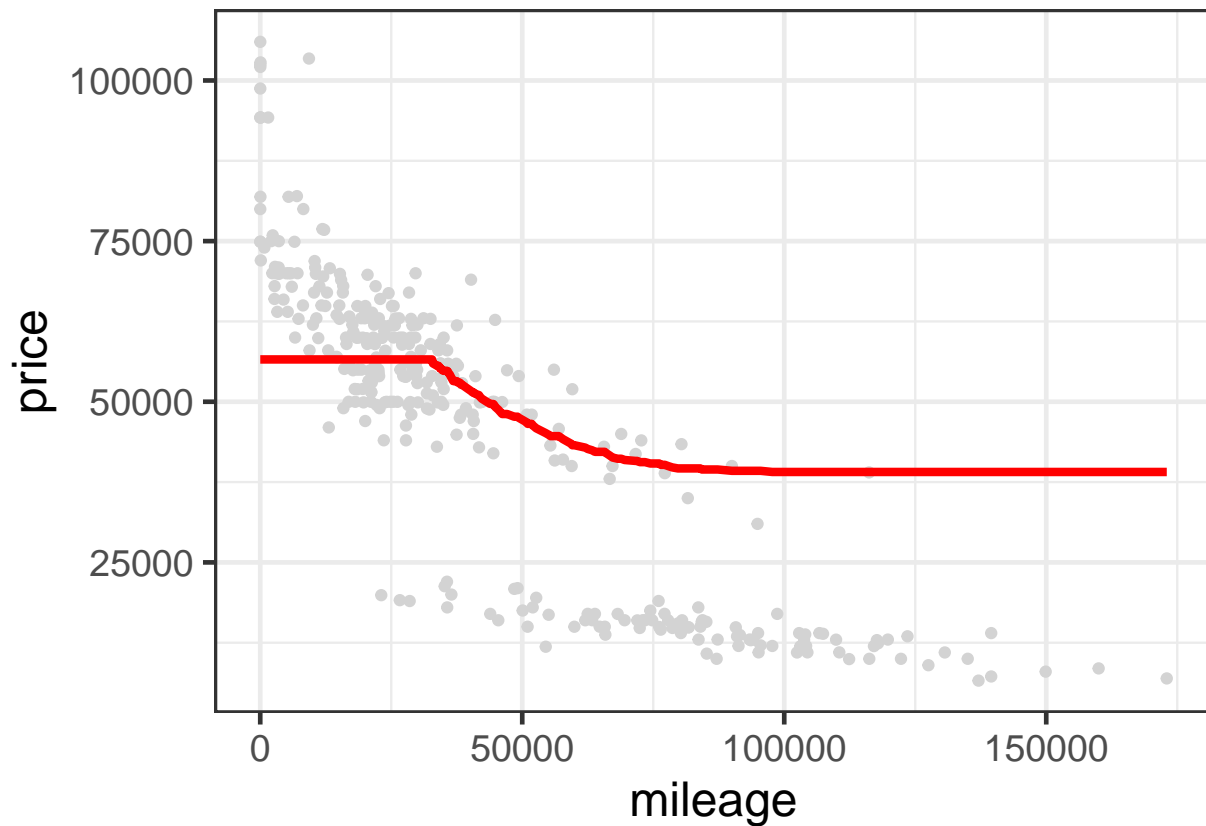


```
p_test + geom_point(aes(x = mileage, y = ypred_knn250), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 250)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K=300
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn300 = knn.reg(train = X_train, test = X_test, y = y_train, k=300)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn300 = knn300$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 10346.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9556.145
```



```
rmse(y_test, ypred_knn300)
```

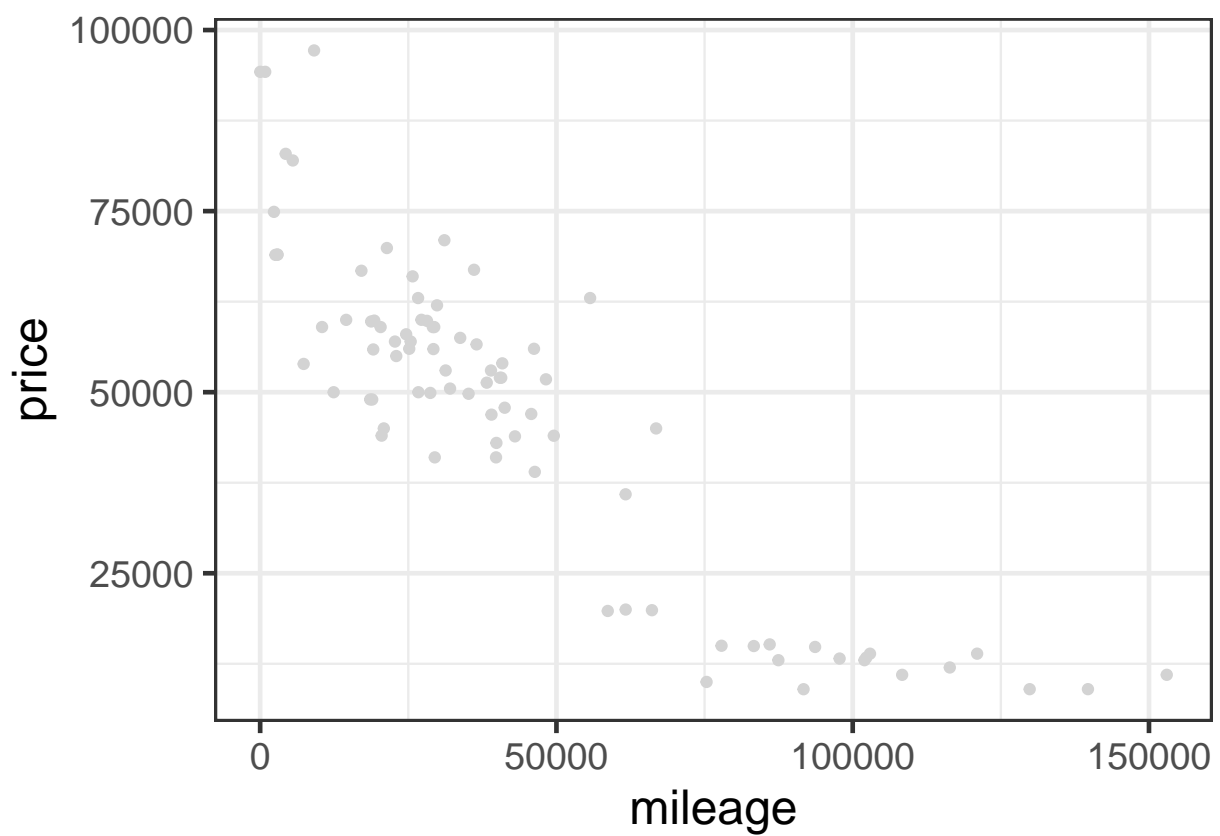
```
## [1] 19512.26
```

```
#attach predictions to data frame
```

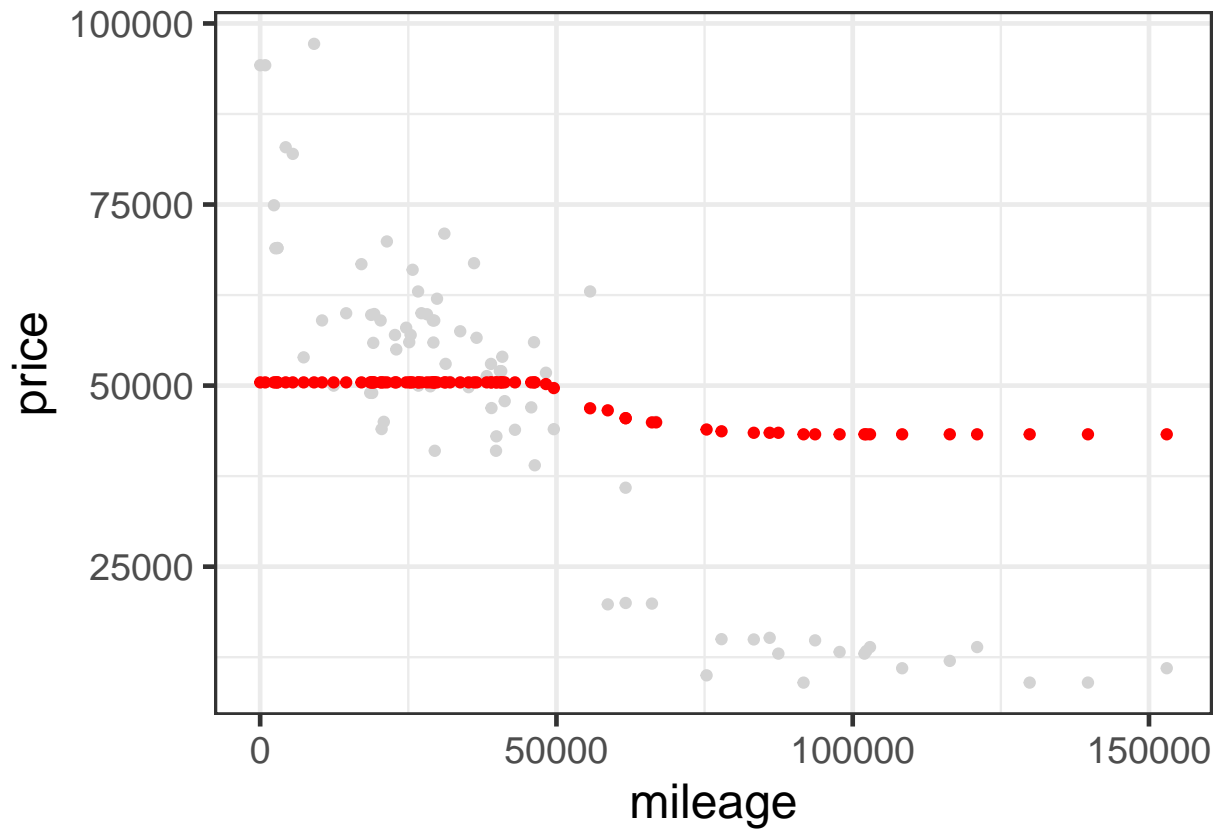
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn300 = ypred_knn300
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

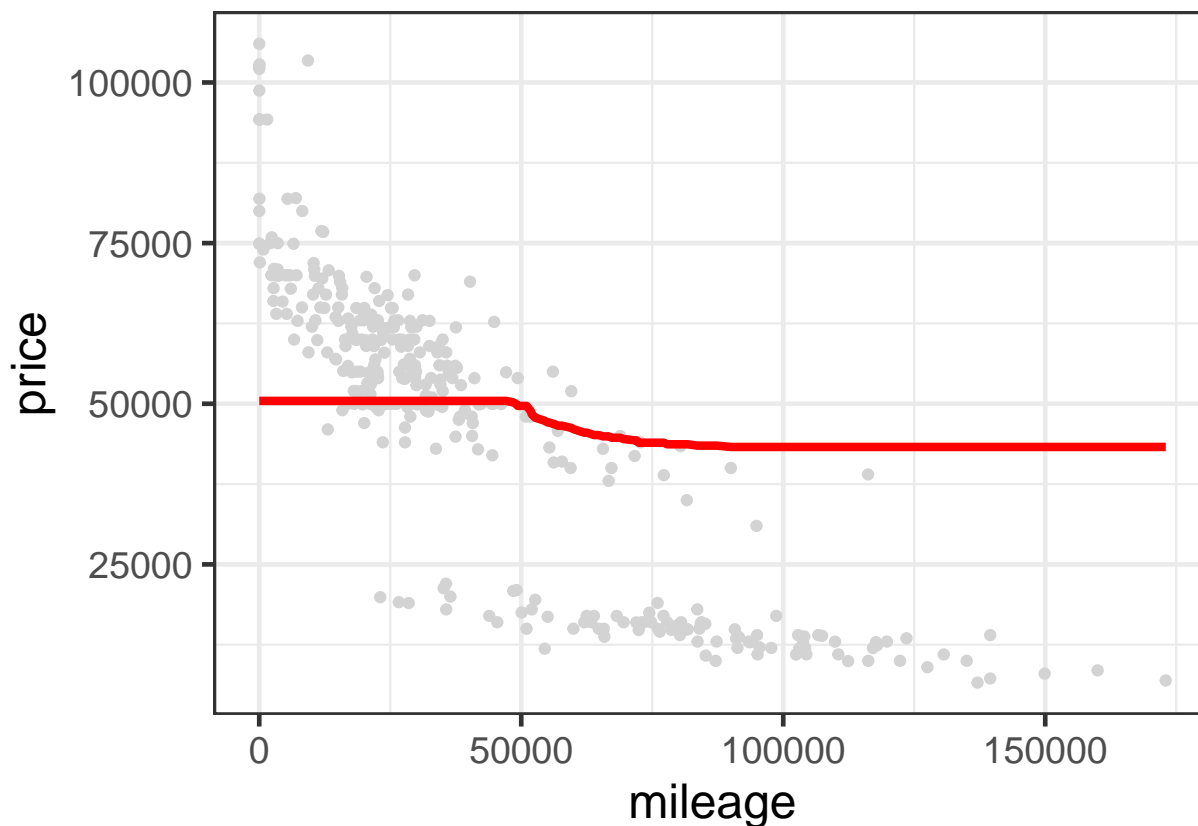


```
p_test + geom_point(aes(x = mileage, y = ypred_knn300), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 300)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



#RMSE plot vs K nearest neighbors 350 trim level.

```
N = nrow(sclass350)
N_train = floor(0.8*N)
N_test = N - N_train

train_ind = sort(sample.int(N, N_train, replace=FALSE))

D_train = sclass350[train_ind,]
D_train = arrange(D_train, mileage)
D_test = sclass350[-train_ind,]

y_train = D_train$price
X_train = data.frame(mileage=jitter(D_train$mileage))
X_test = data.frame(mileage=jitter(D_test$mileage))
y_test = D_test$price
library(foreach)
```

```
## Warning: package 'foreach' was built under R version 3.6.3
```

```
##
```

```
## Attaching package: 'foreach'
```

```
## The following objects are masked from 'package:purrr':
```

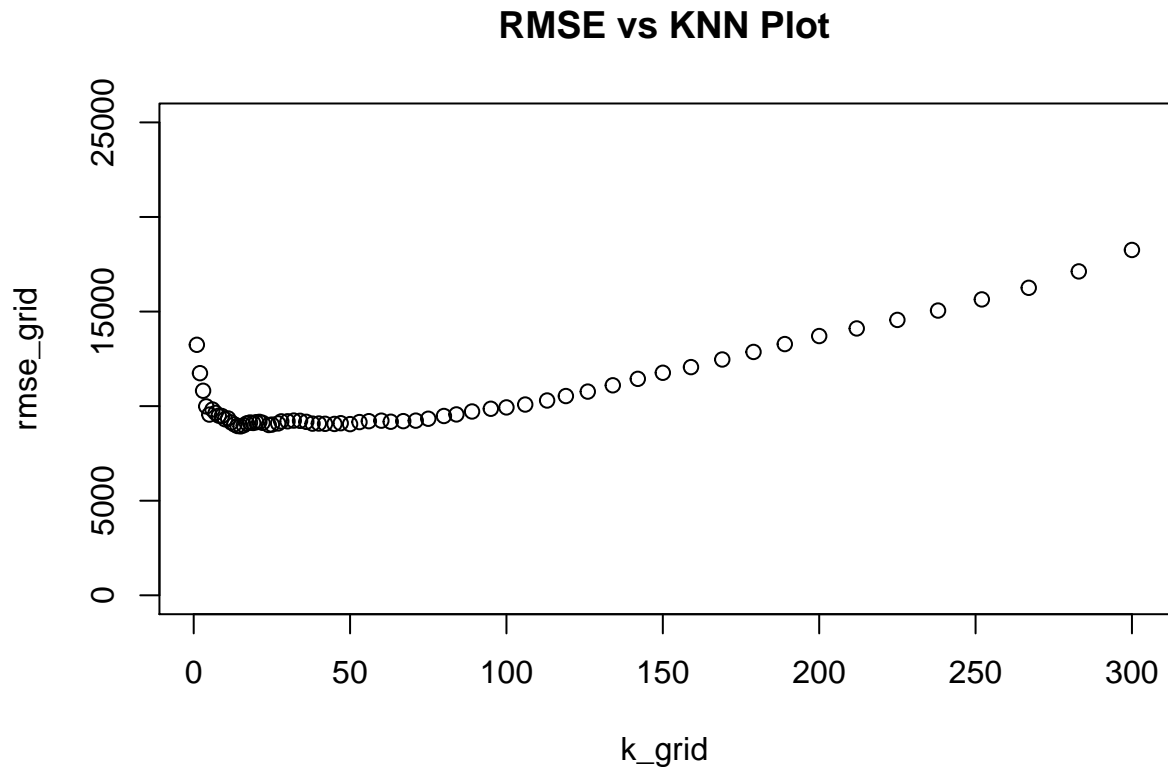
```
##
```

```
## accumulate, when
```

```

k_grid = exp(seq(log(1), log(300), length=100)) %>% round %>% unique
rmse_grid = foreach(K = k_grid, .combine='c') %do% {
  knn_model = knn.reg(train = X_train, test = X_test, y = y_train, k=K)
  rmse(y_test, knn_model$pred)
}
rmse_plot = plot(k_grid, rmse_grid, ylim=c(0,25000), main="RMSE vs KNN Plot")
abline(rmse_plot)

```



```

#optimal K value for 350 AMG trim level

#K = 30
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn30 = knn.reg(train = X_train, test = X_test, y = y_train, k=30)

#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}

ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn30 = knn30$pred

```

```
rmse(y_test, ypred_lm1)
```

```
## [1] 10425.46
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 9695.712
```

```
rmse(y_test, ypred_knn30)
```

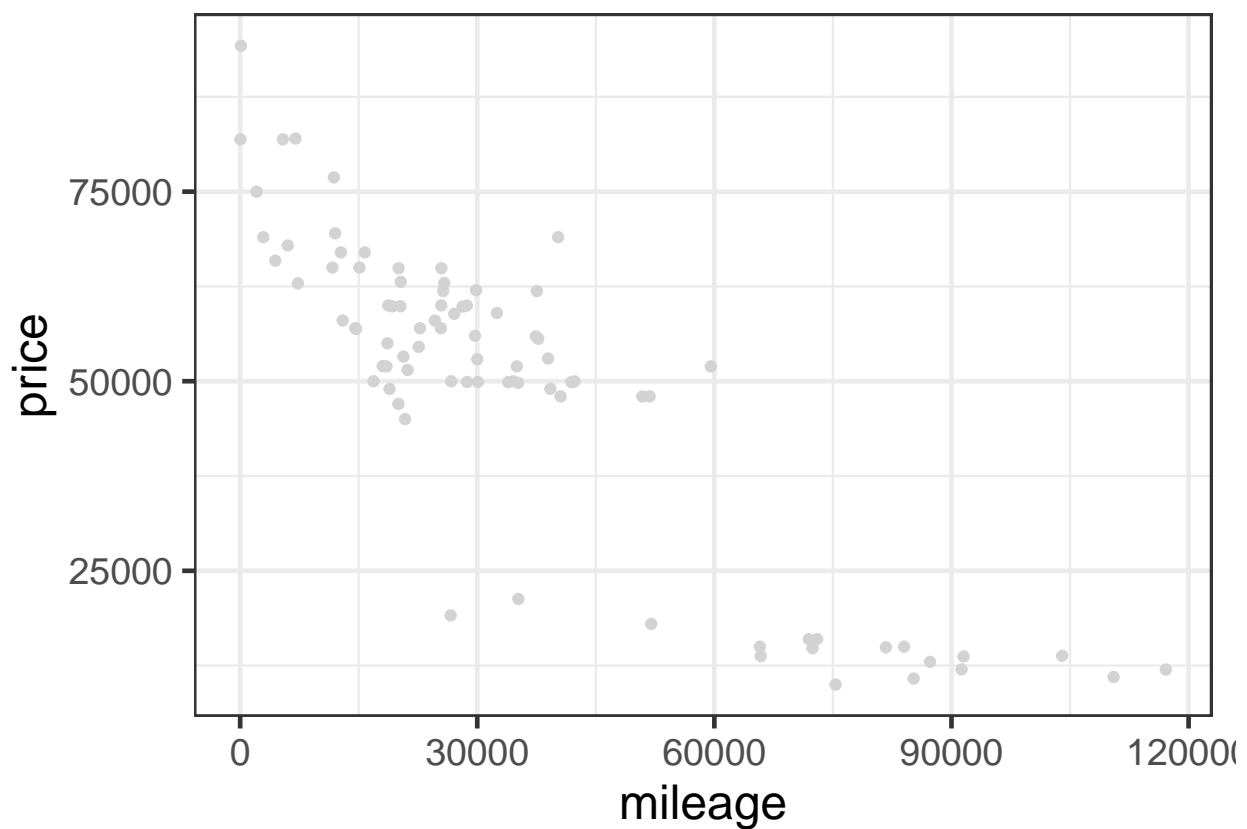
```
## [1] 9198.996
```

```
#attach predictions to data frame
```

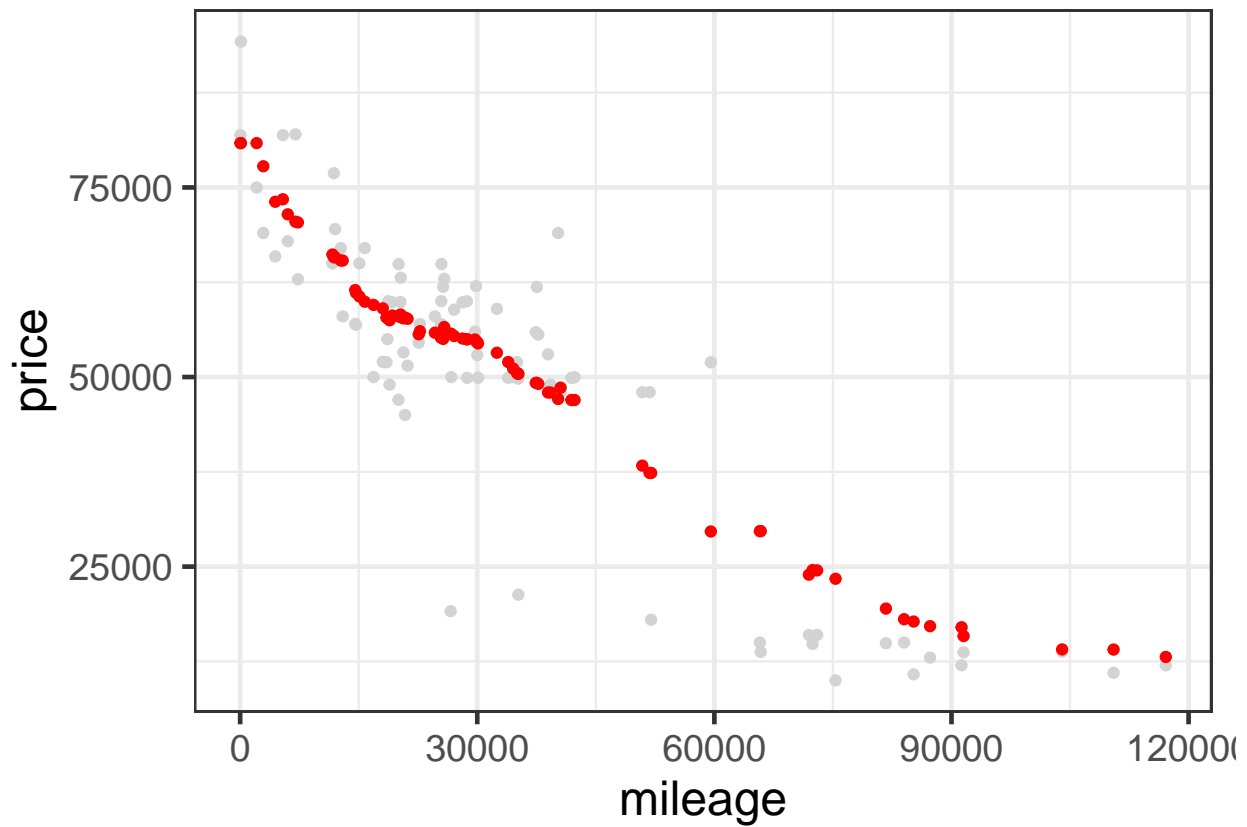
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn30 = ypred_knn30
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

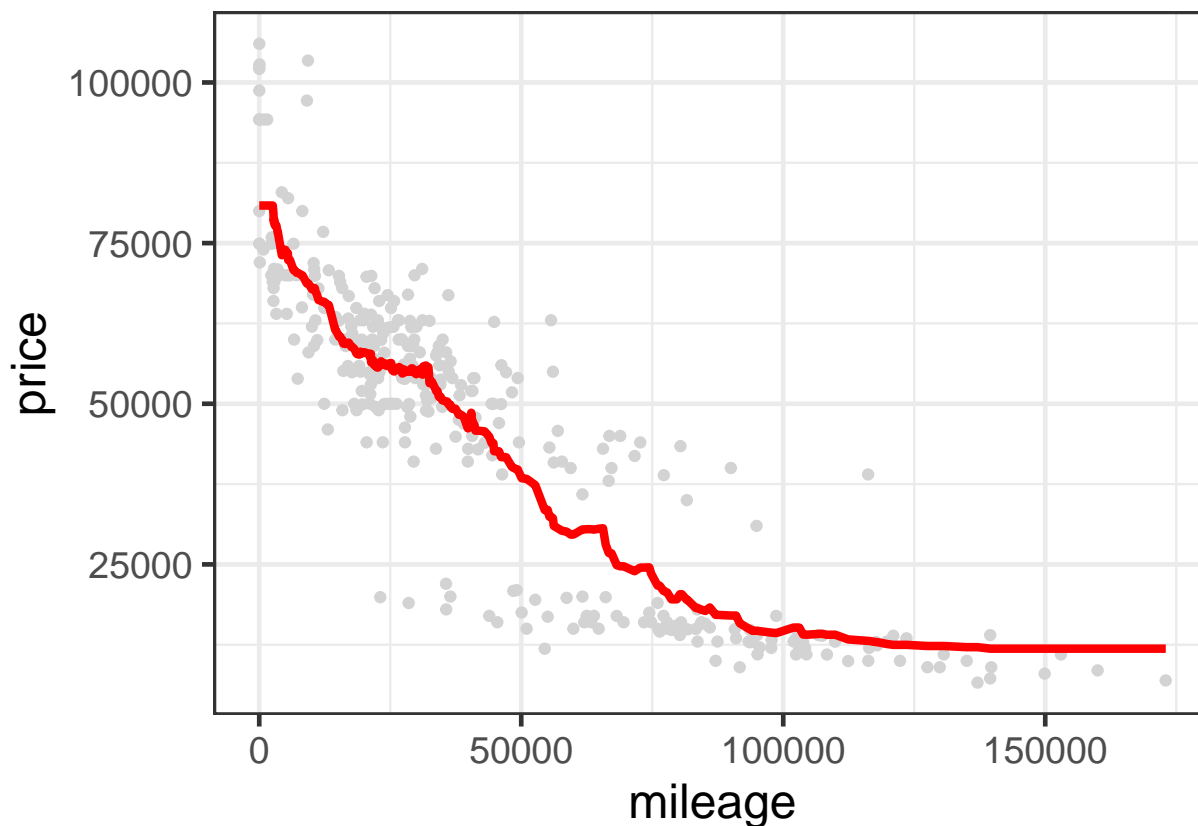


```
p_test + geom_point(aes(x = mileage, y = ypred_knn30), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 30)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



#The value of K=30 is optimal for the 350 AMG trim level with an out of sample RMSE of 25924.26. This v

#Train-test split for 65 AMG trim level: We are looking for a model that gives the best explanation wit

```
N = nrow(sclass65AMG)
N_train = floor(0.8*N)
N_test = N - N_train

train_ind = sort(sample.int(N, N_train, replace=FALSE))
```

```
D_train = sclass65AMG[train_ind,]
D_train = arrange(D_train, mileage)
D_test = sclass65AMG[-train_ind,]
```

```
y_train = D_train$price
X_train = data.frame(mileage=jitter(D_train$mileage))
X_test = data.frame(mileage=jitter(D_test$mileage))
y_test = D_test$price
```

#Running KNN test for various values of K to determine the optimal predictive model.This will start at

```
#KNN = 2
```

#The first step is to make predictions to the data frame.

```
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn2 = knn.reg(train = X_train, test = X_test, y = y_train, k=2)
```

```
#rmse calculation
```

```
rmse = function(y, ypred) {  
  sqrt(mean(data.matrix((y-ypred)^2)))  
}
```

```
ypred_lm1 = predict(lm1, X_test)  
ypred_lm2 = predict(lm2, X_test)  
ypred_knn2 = knn2$pred
```

```
rmse(y_test, ypred_lm1)
```

```
## [1] 46463.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 34814.4
```

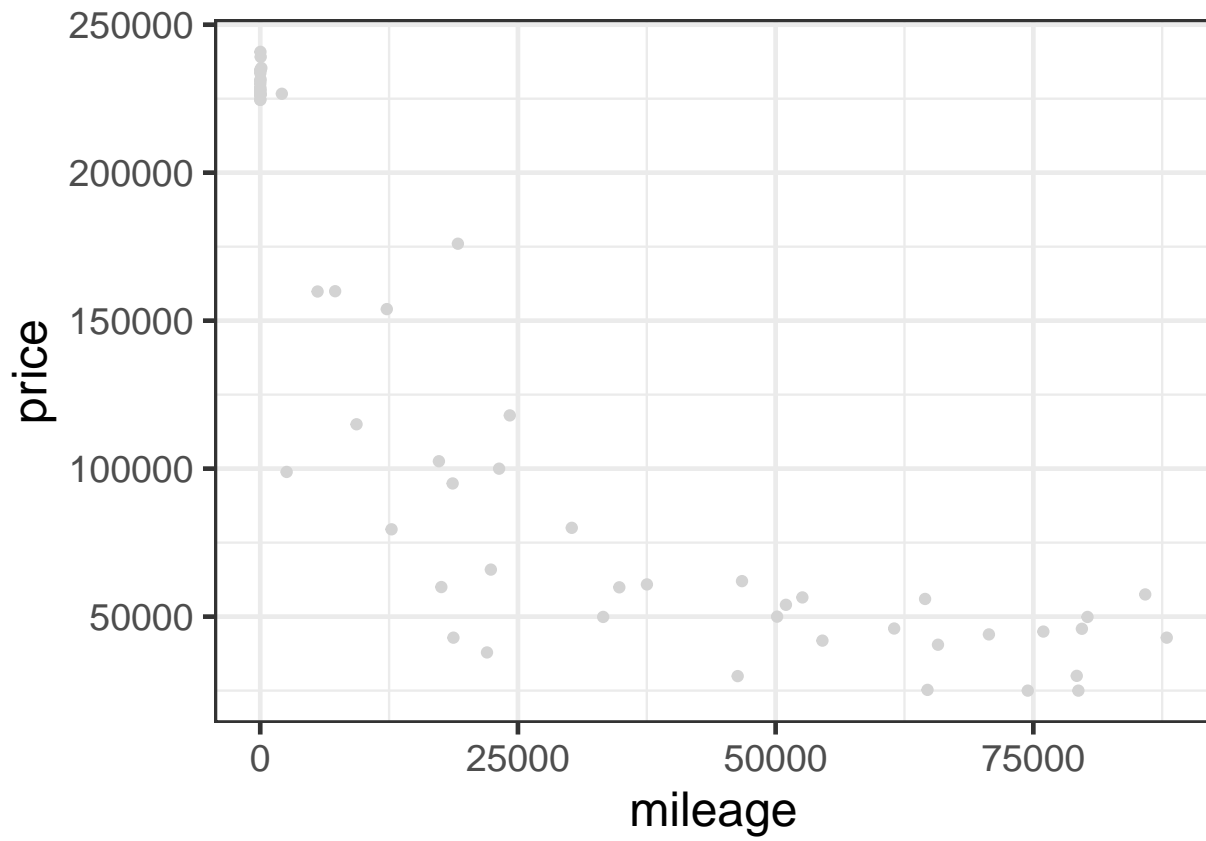
```
rmse(y_test, ypred_knn2)
```

```
## [1] 26213.38
```

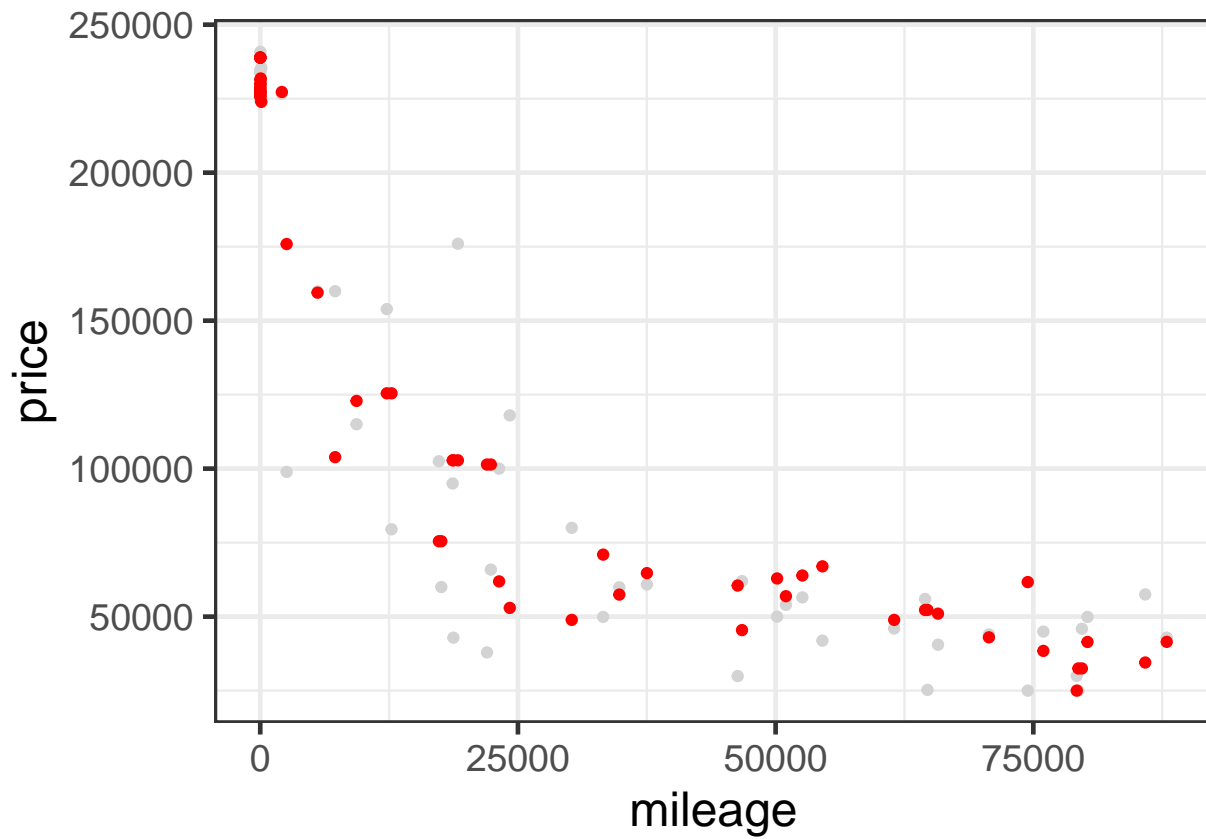
```
#attach predictions to data frame
```

```
D_test$ypred_lm2 = ypred_lm2  
D_test$ypred_knn2 = ypred_knn2
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

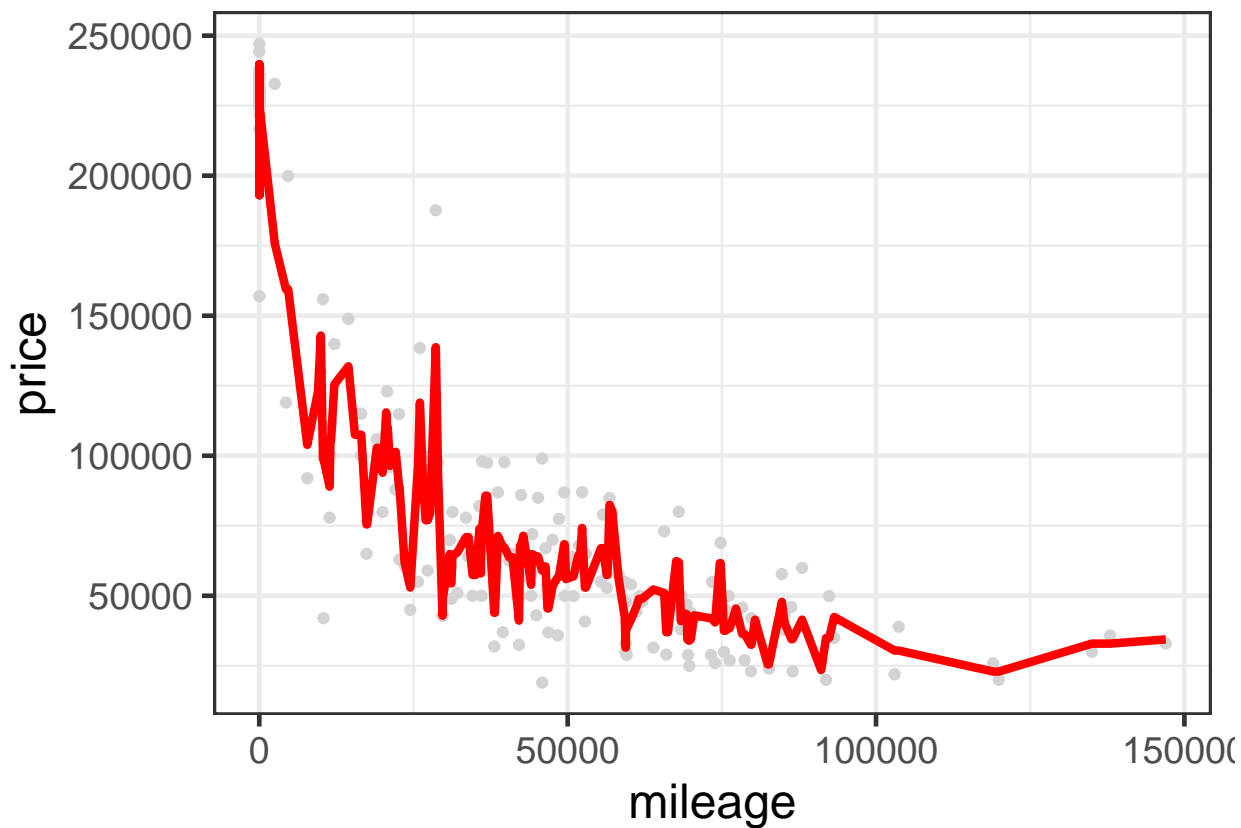



```
p_test + geom_point(aes(x = mileage, y = ypred_knn2), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 2)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 3
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn3 = knn.reg(train = X_train, test = X_test, y = y_train, k=3)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn3 = knn3$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 46463.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 34814.4
```

```
rmse(y_test, ypred_knn3)
```

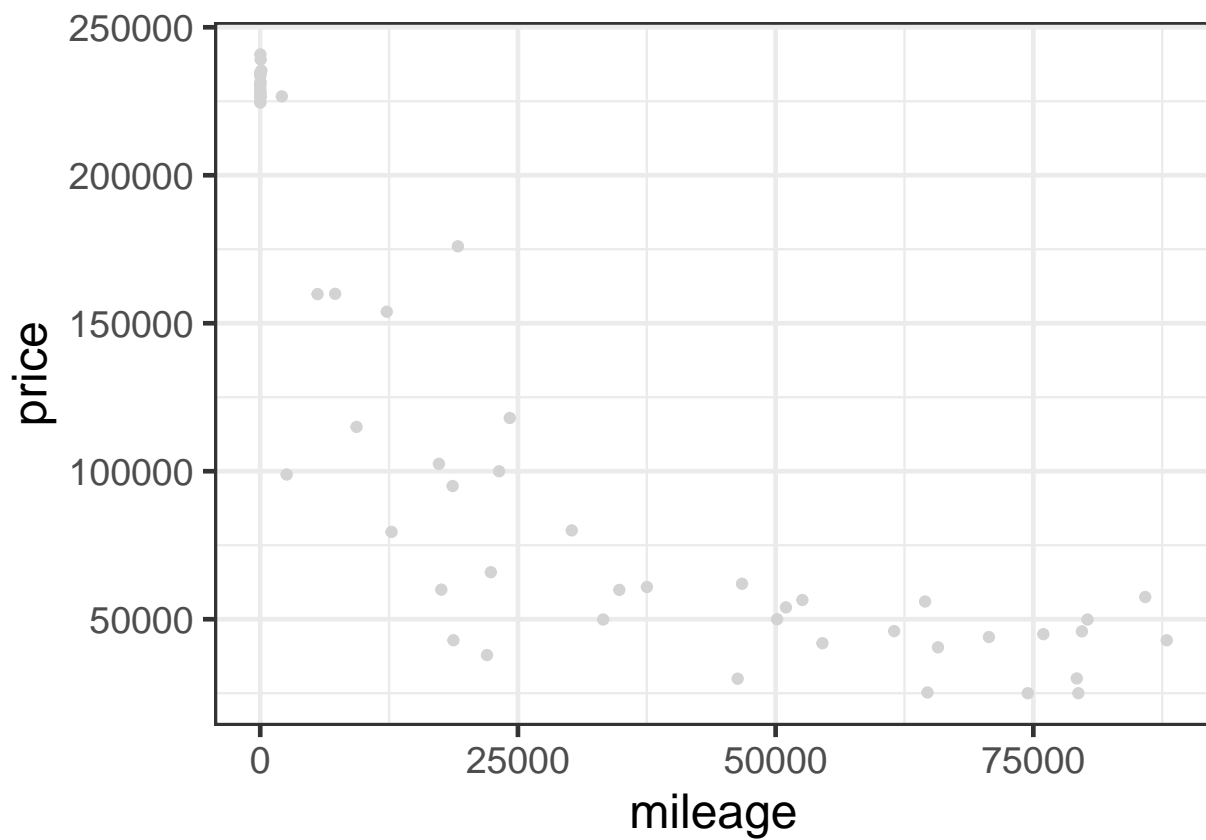
```
## [1] 24130.88
```

```
#attach predictions to data frame
```

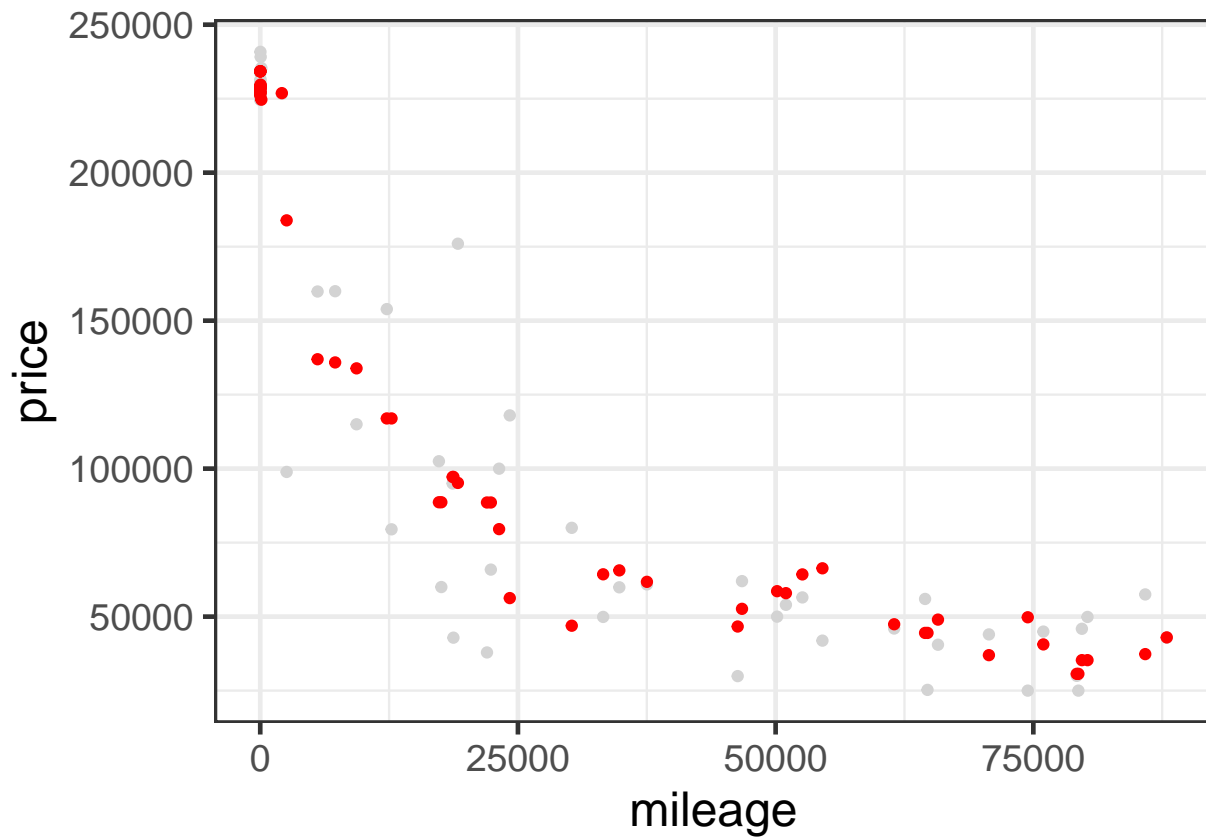
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn3 = ypred_knn3
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

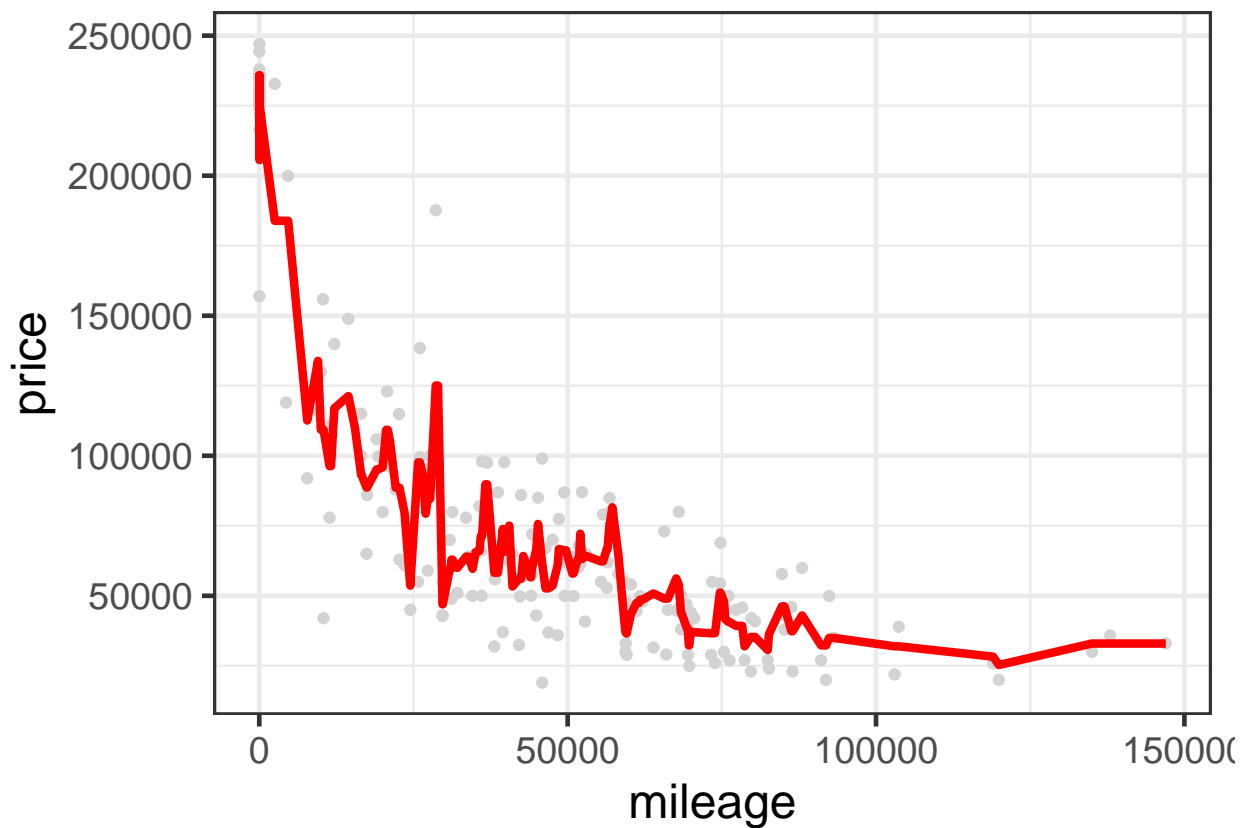


```
p_test + geom_point(aes(x = mileage, y = ypred_knn3), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 3)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 4
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn4 = knn.reg(train = X_train, test = X_test, y = y_train, k=4)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn4 = knn4$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 46463.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 34814.4
```

```
rmse(y_test, ypred_knn4)
```

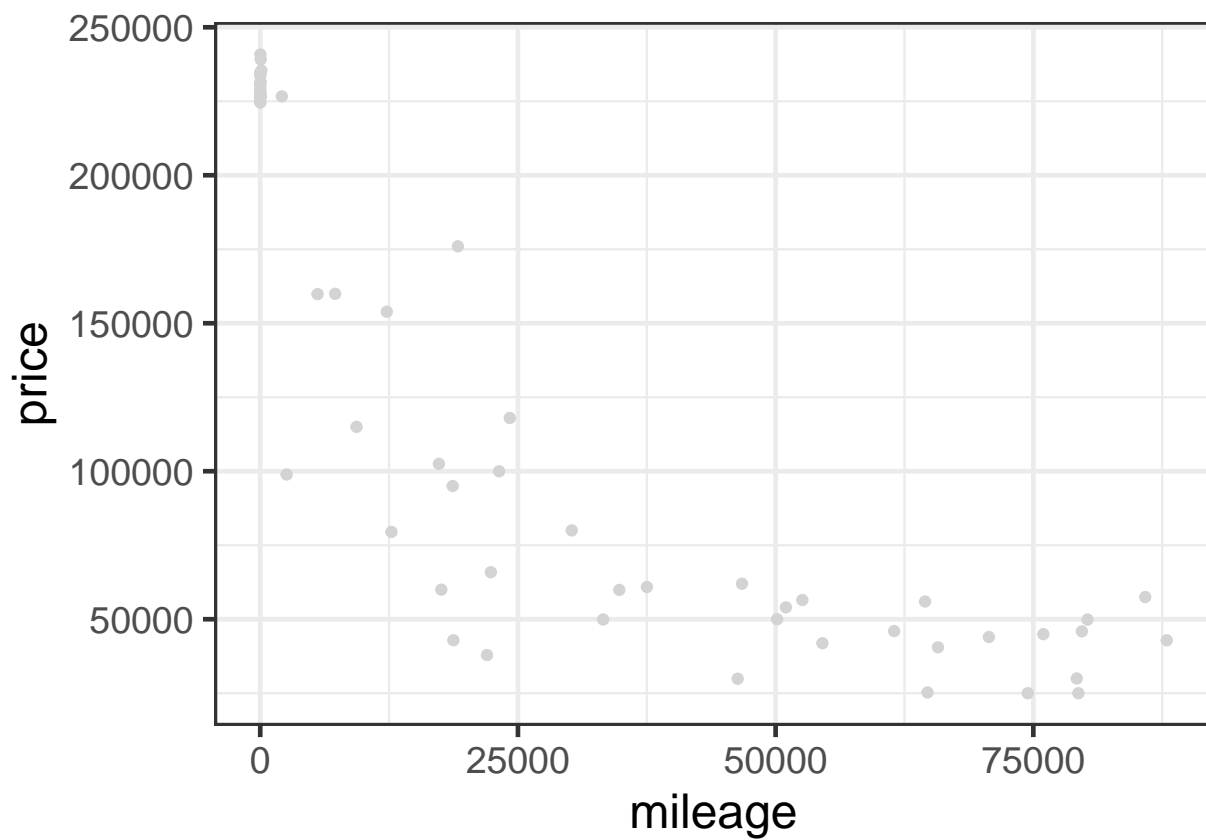
```
## [1] 24286.23
```

```
#attach predictions to data frame
```

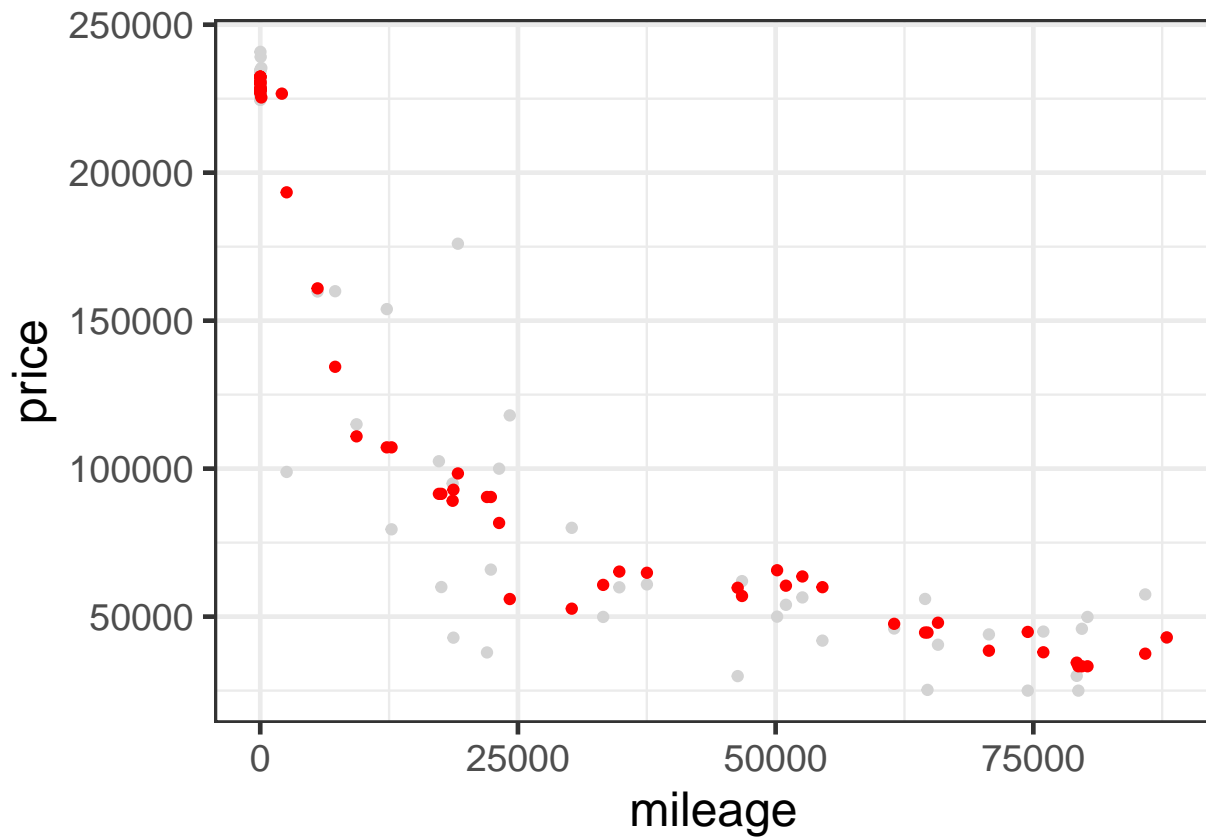
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn4 = ypred_knn4
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

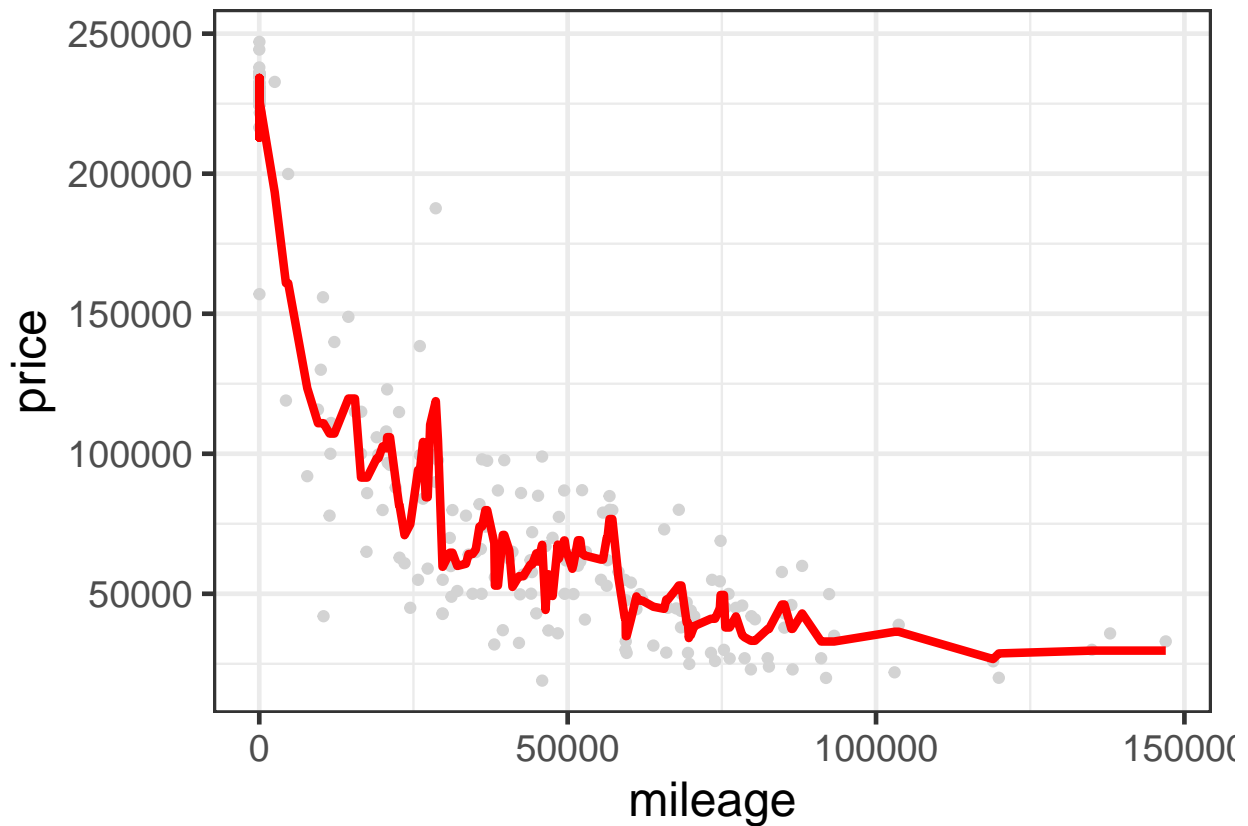


```
p_test + geom_point(aes(x = mileage, y = ypred_knn4), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 4)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```

```
#K = 5
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn5 = knn.reg(train = X_train, test = X_test, y = y_train, k=5)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn5 = knn5$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 46463.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 34814.4
```

```
rmse(y_test, ypred_knn5)
```

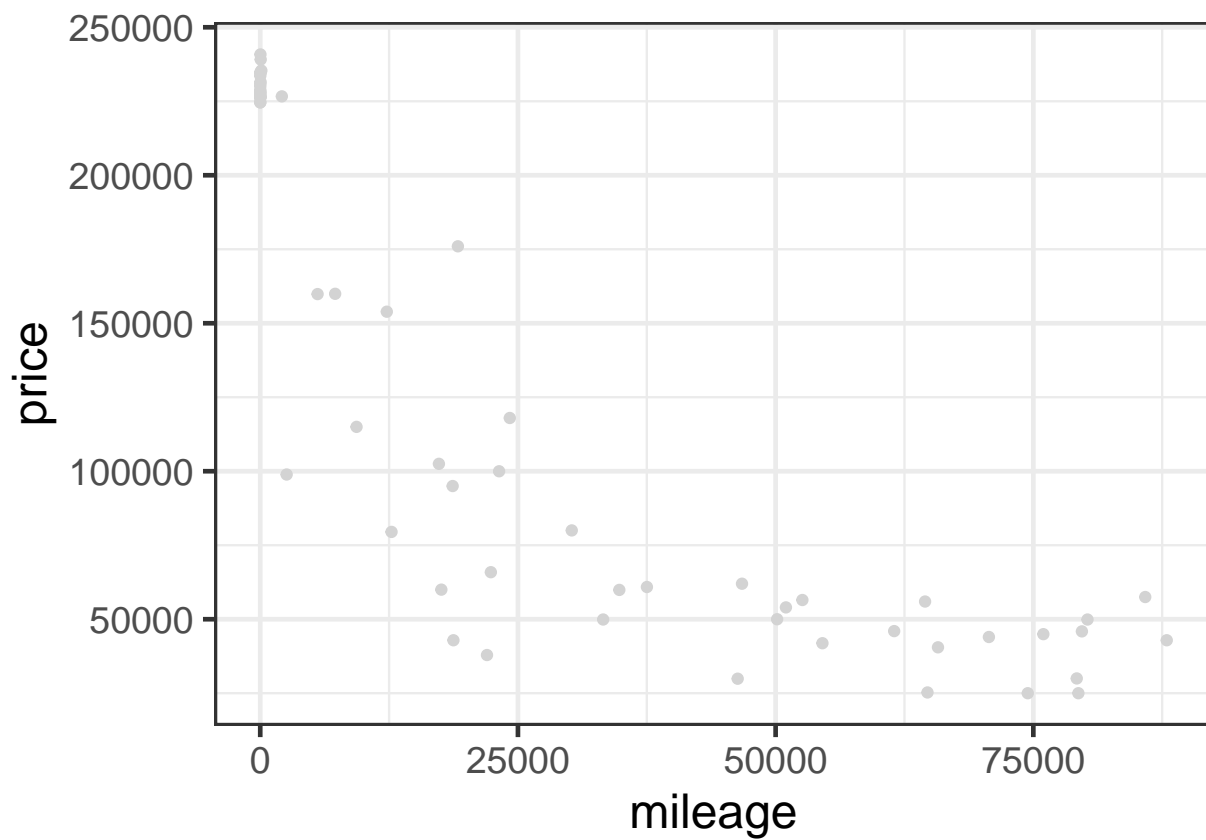
```
## [1] 24550.1
```

```
#attach predictions to data frame
```

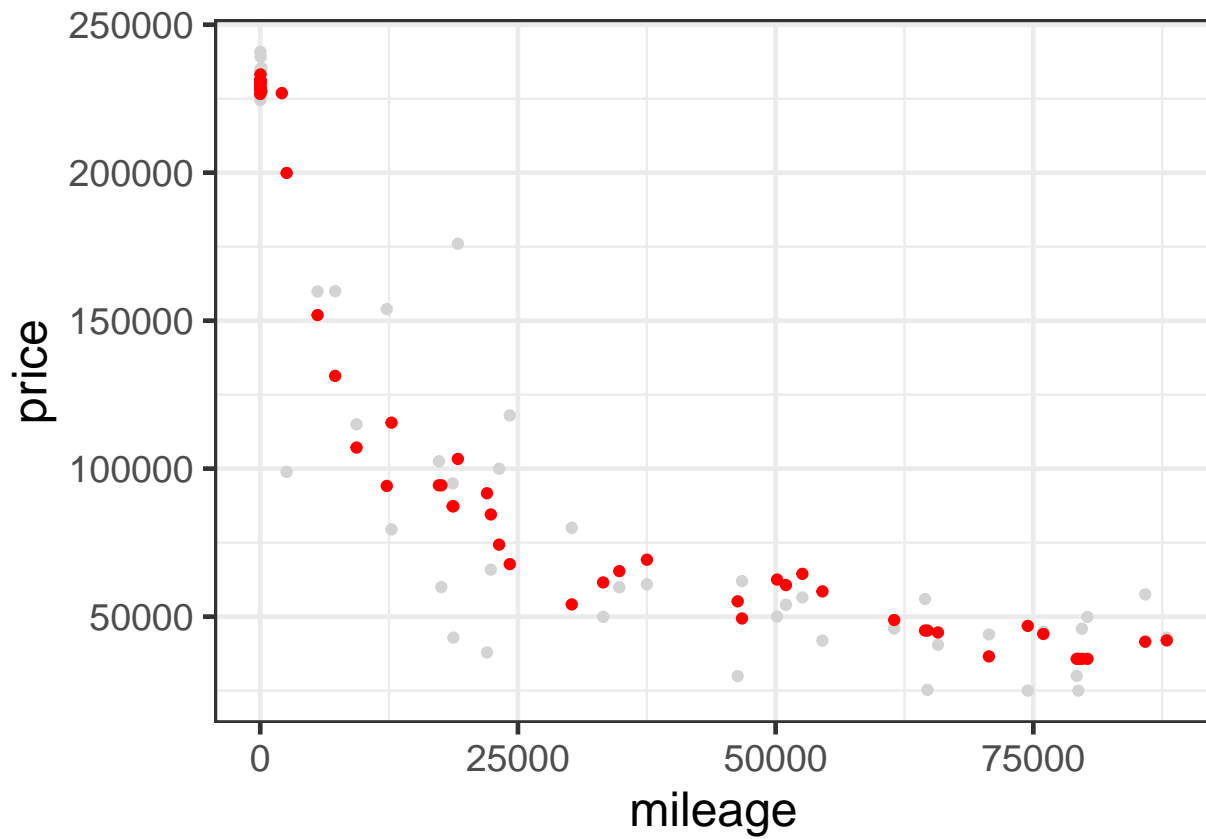
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn5 = ypred_knn5
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

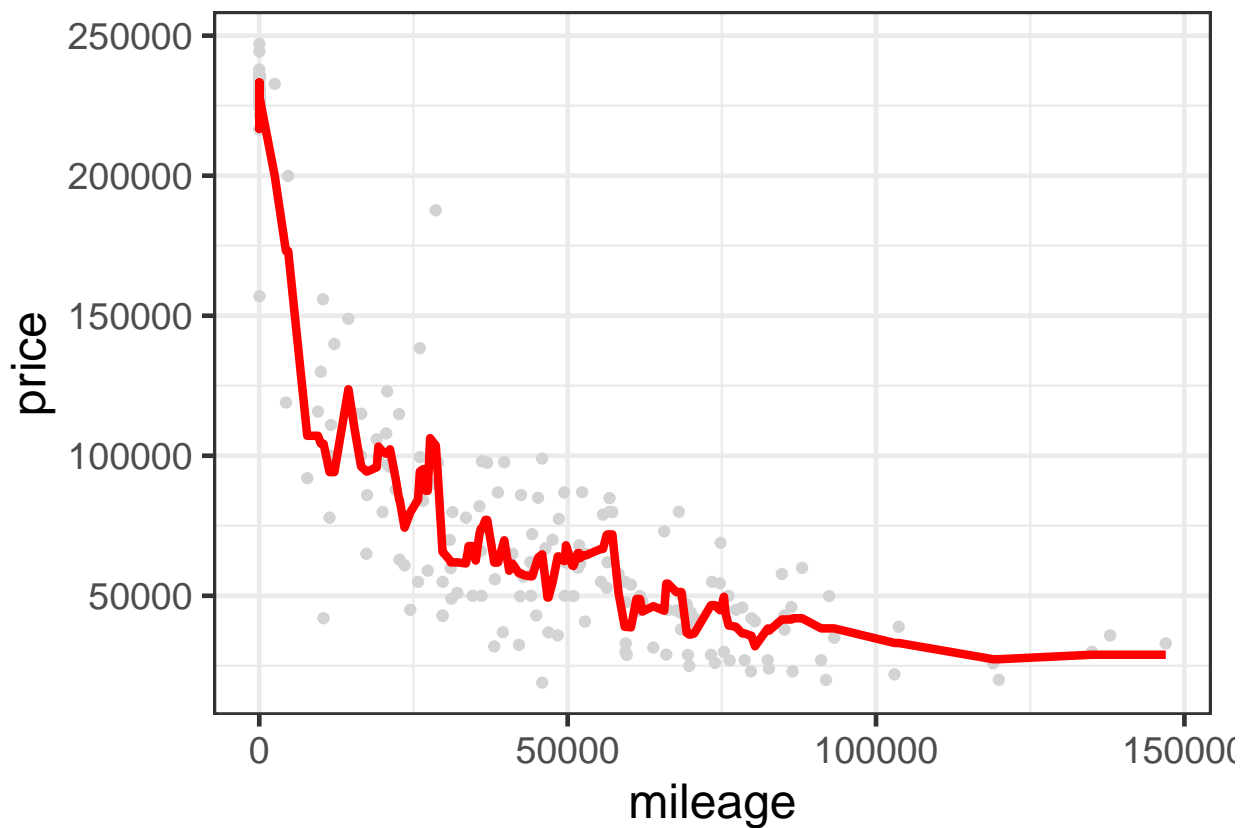


```
p_test + geom_point(aes(x = mileage, y = ypred_knn5), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 5)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 6
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn6 = knn.reg(train = X_train, test = X_test, y = y_train, k=6)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn6 = knn6$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 46463.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 34814.4
```

```
rmse(y_test, ypred_knn6)
```

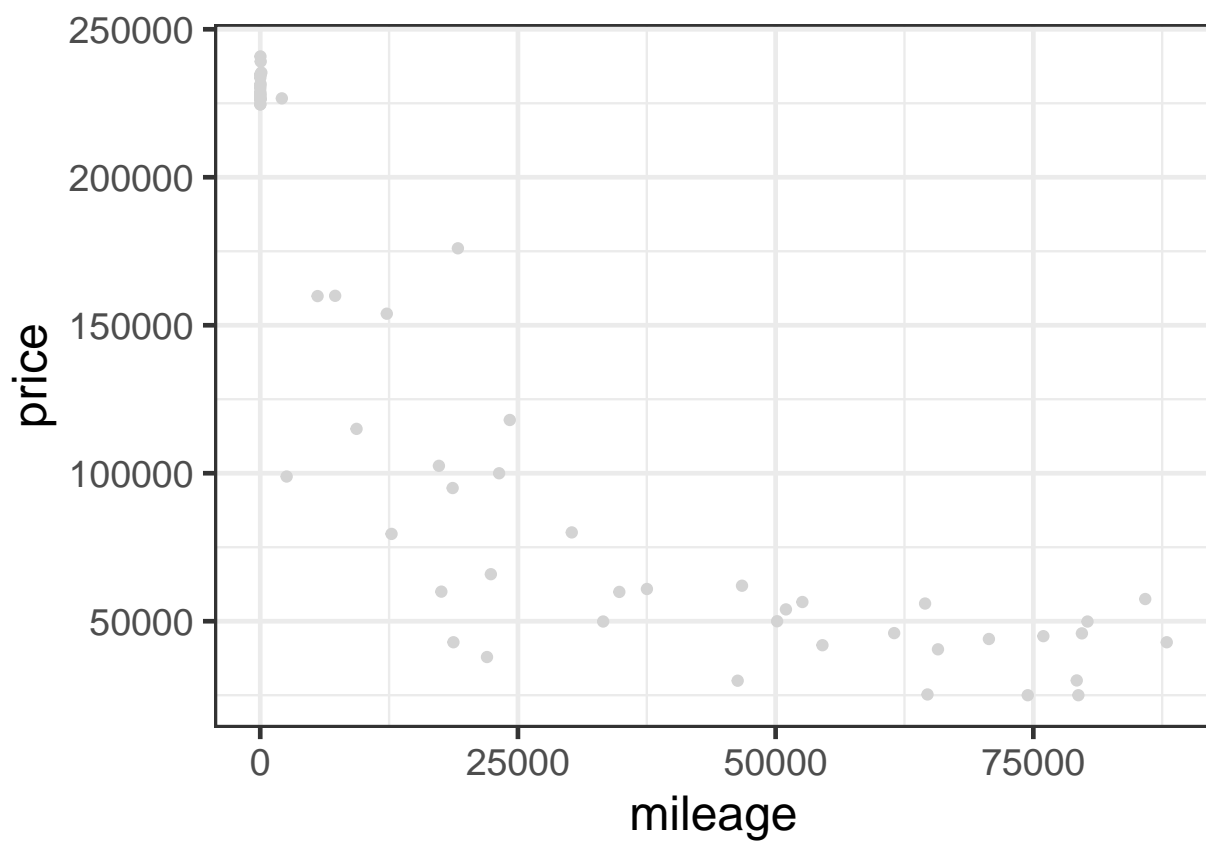
```
## [1] 24405.25
```

```
#attach predictions to data frame
```

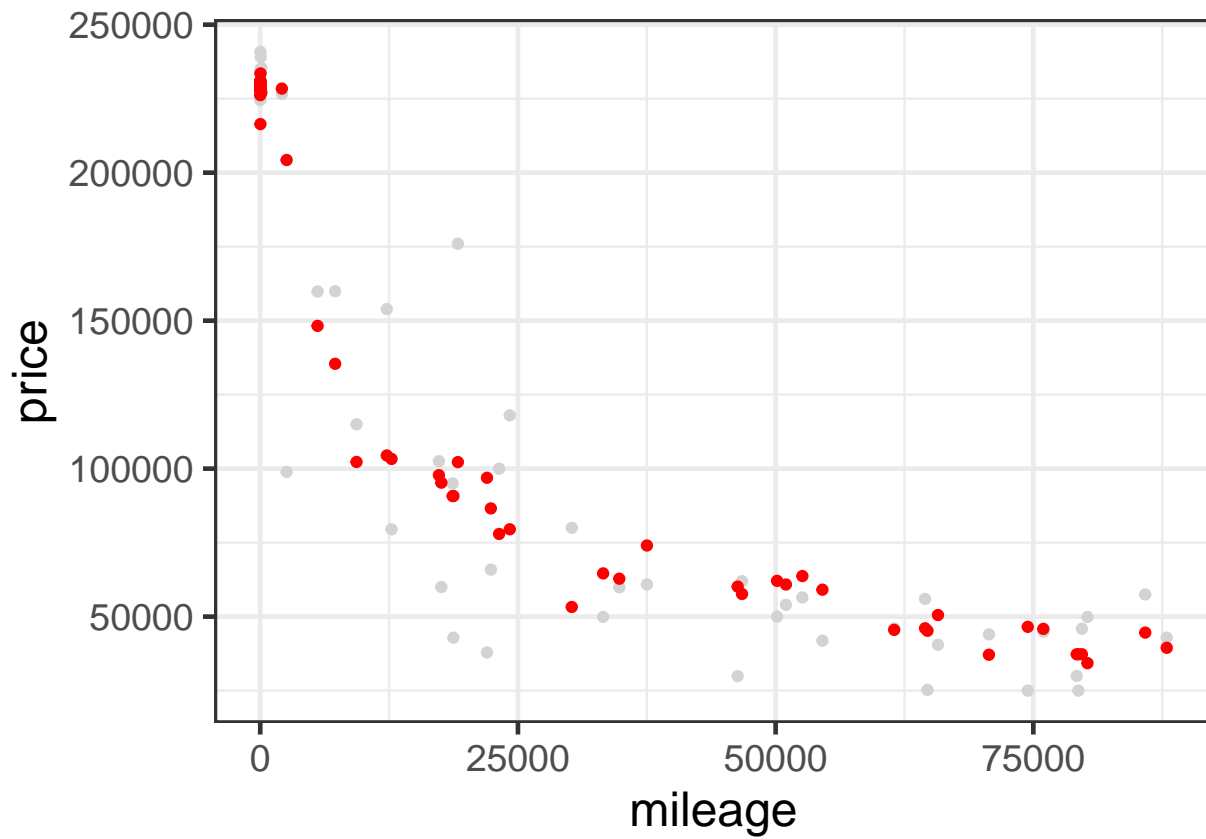
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn6 = ypred_knn6
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

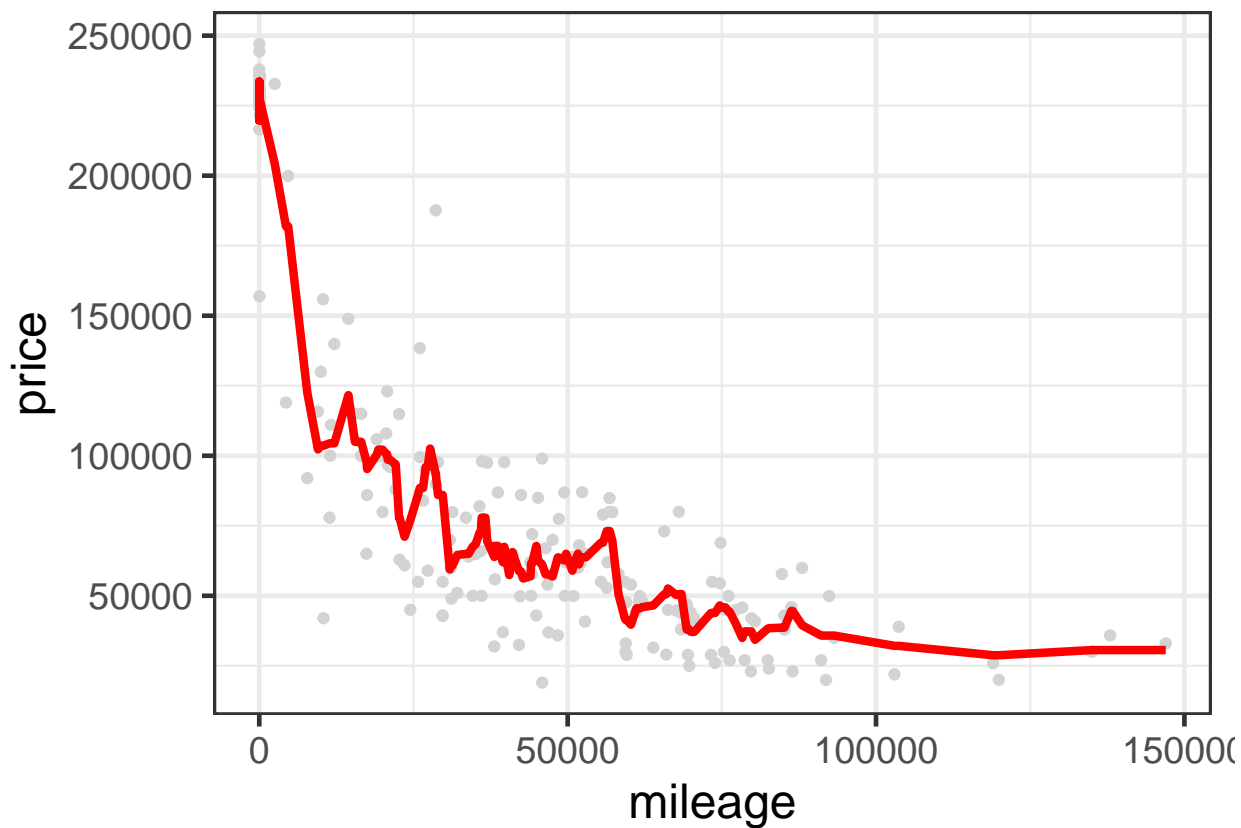


```
p_test + geom_point(aes(x = mileage, y = ypred_knn6), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 6)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 7
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn7 = knn.reg(train = X_train, test = X_test, y = y_train, k=7)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn7 = knn7$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 46463.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 34814.4
```

```
rmse(y_test, ypred_knn7)
```

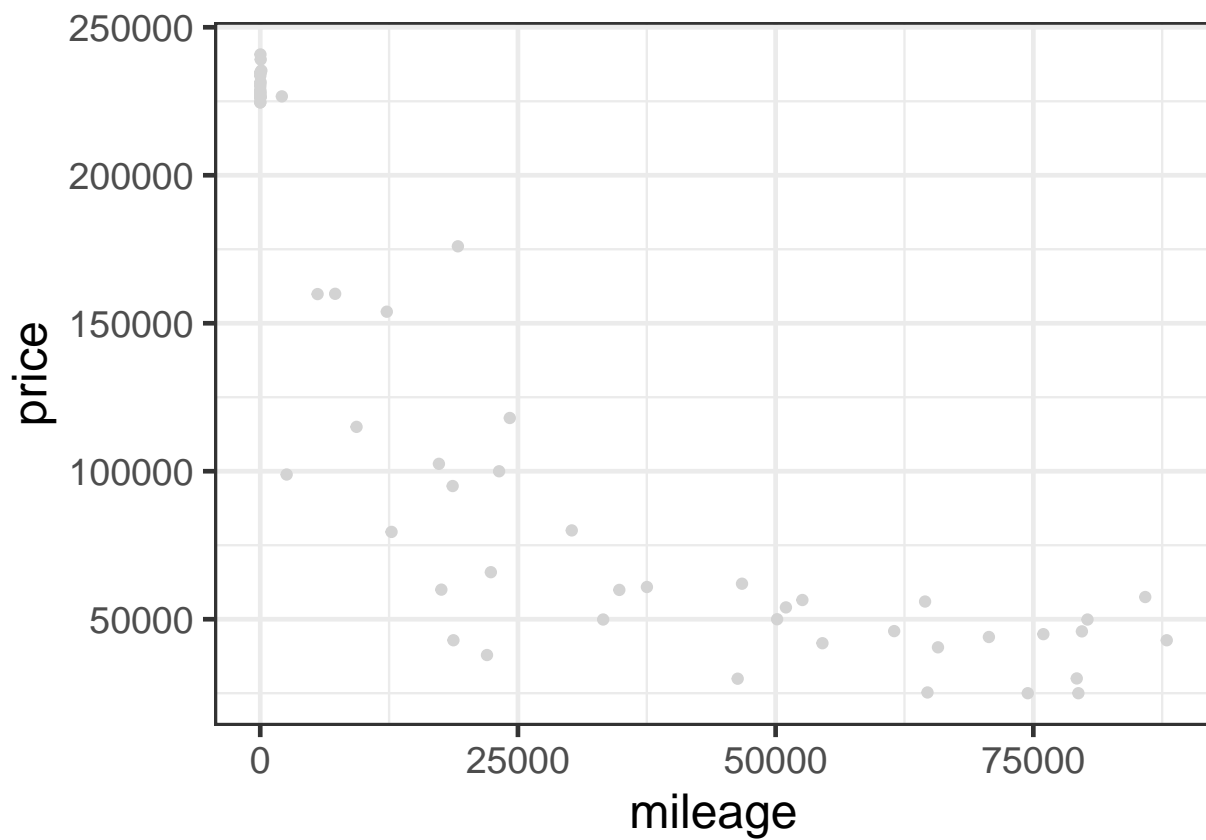
```
## [1] 25222.36
```

```
#attach predictions to data frame
```

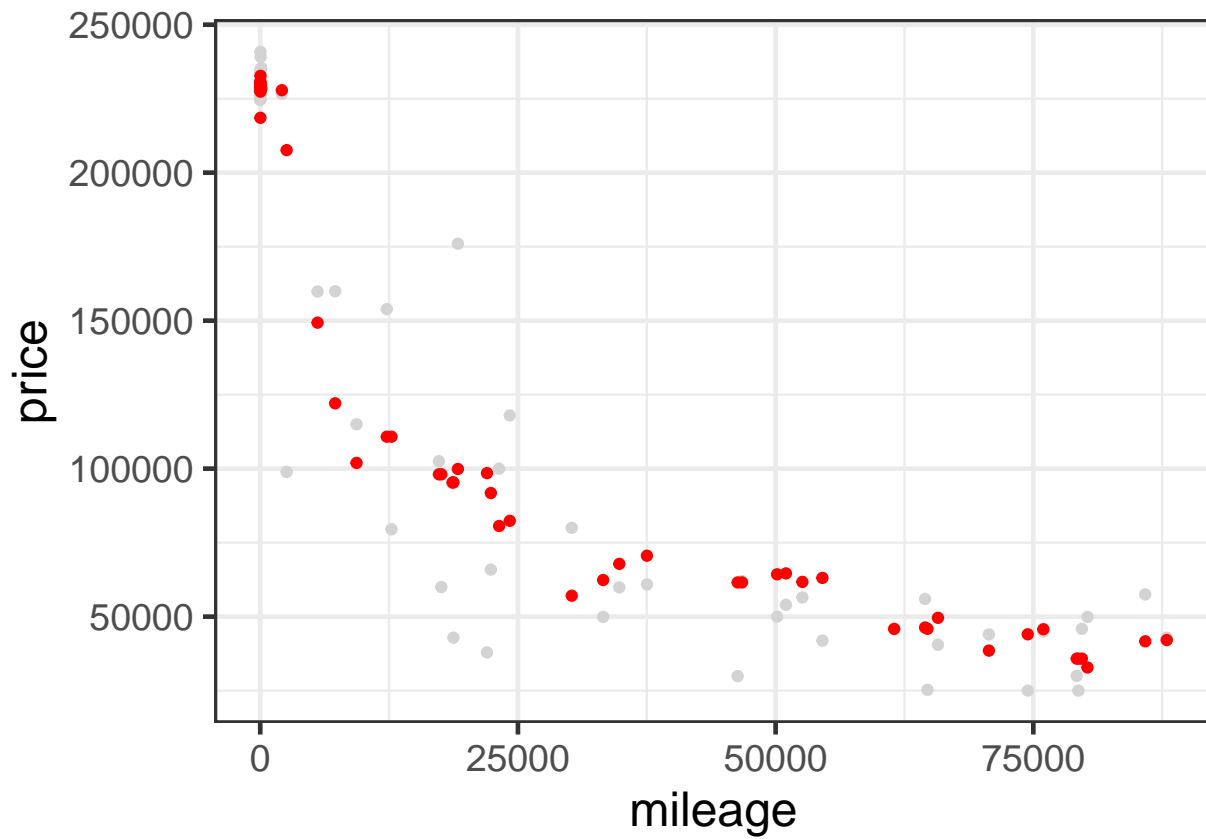
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn7 = ypred_knn7
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

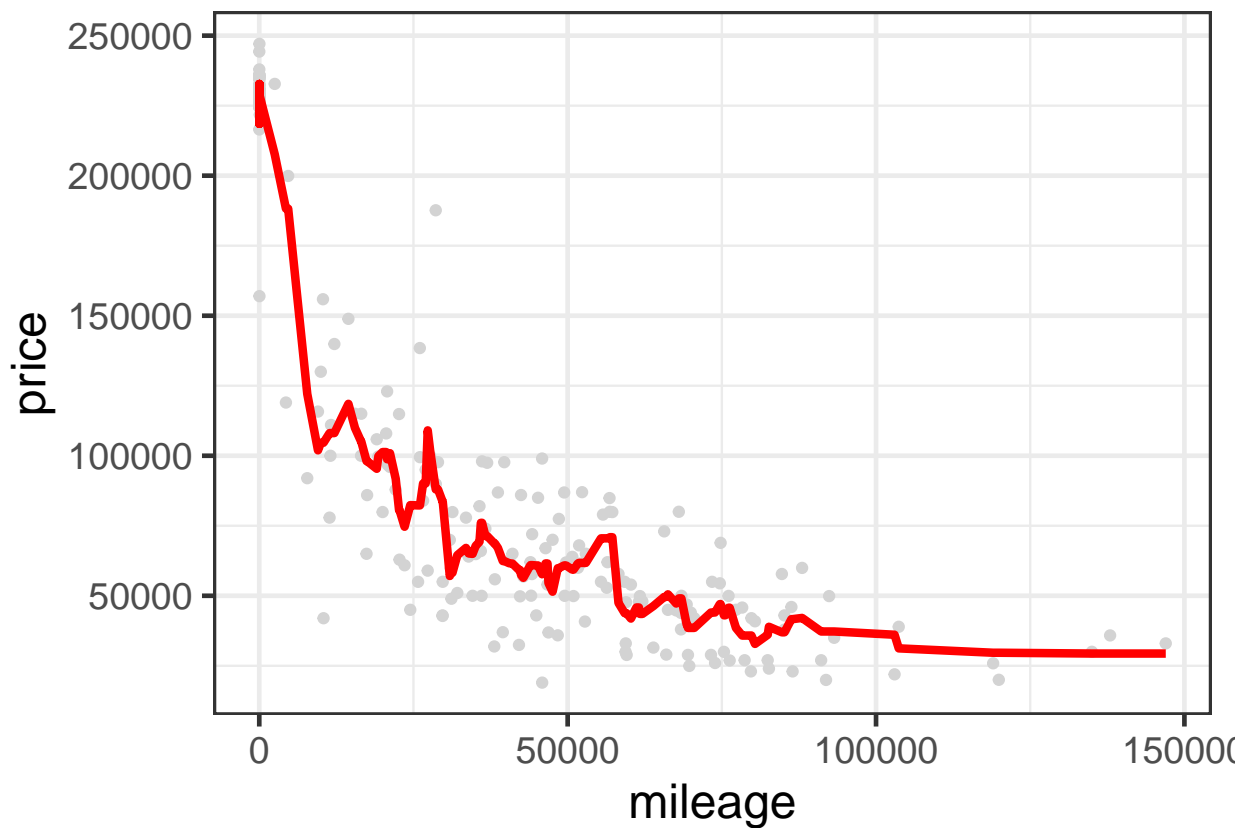


```
p_test + geom_point(aes(x = mileage, y = ypred_knn7), color='red')
```

```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 7)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 15
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn15 = knn.reg(train = X_train, test = X_test, y = y_train, k=15)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn15 = knn15$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 46463.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 34814.4
```

```
rmse(y_test, ypred_knn15)
```

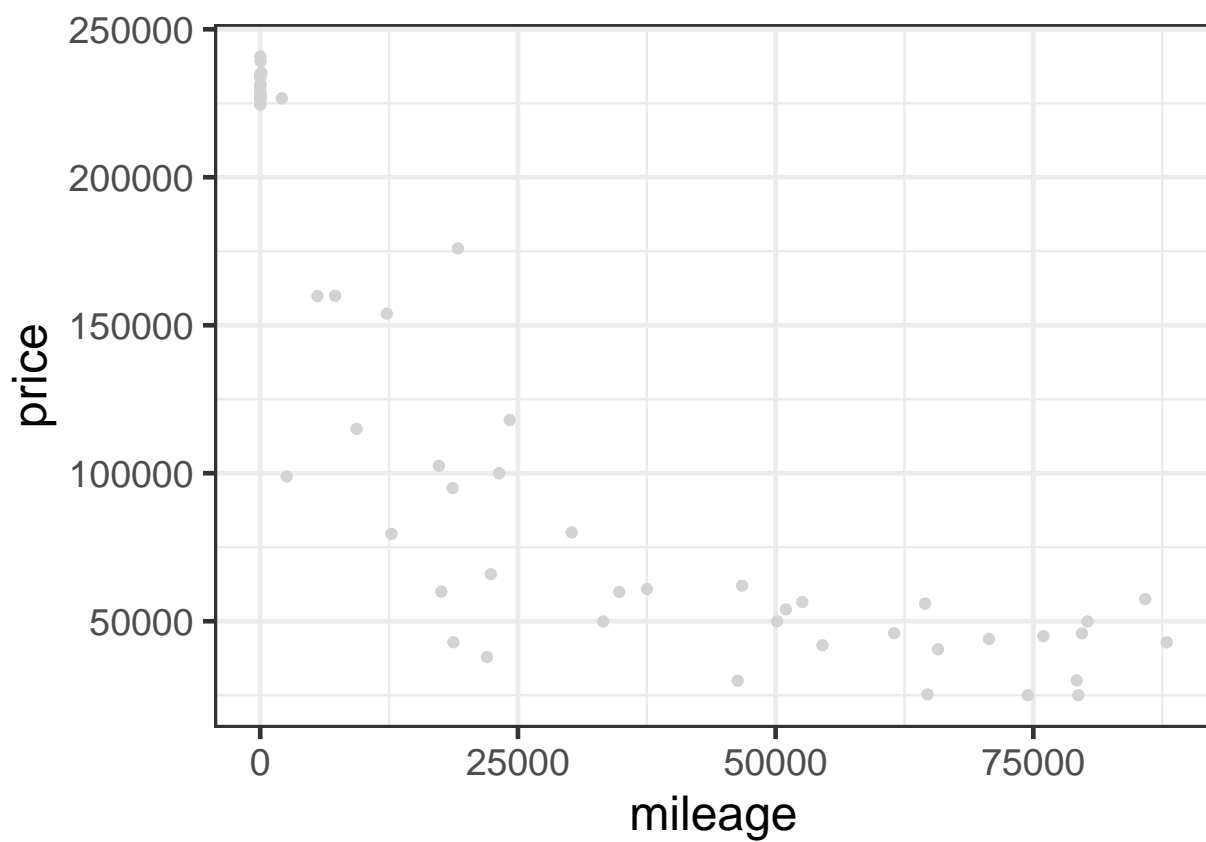
```
## [1] 25562.88
```

```
#attach predictions to data frame
```

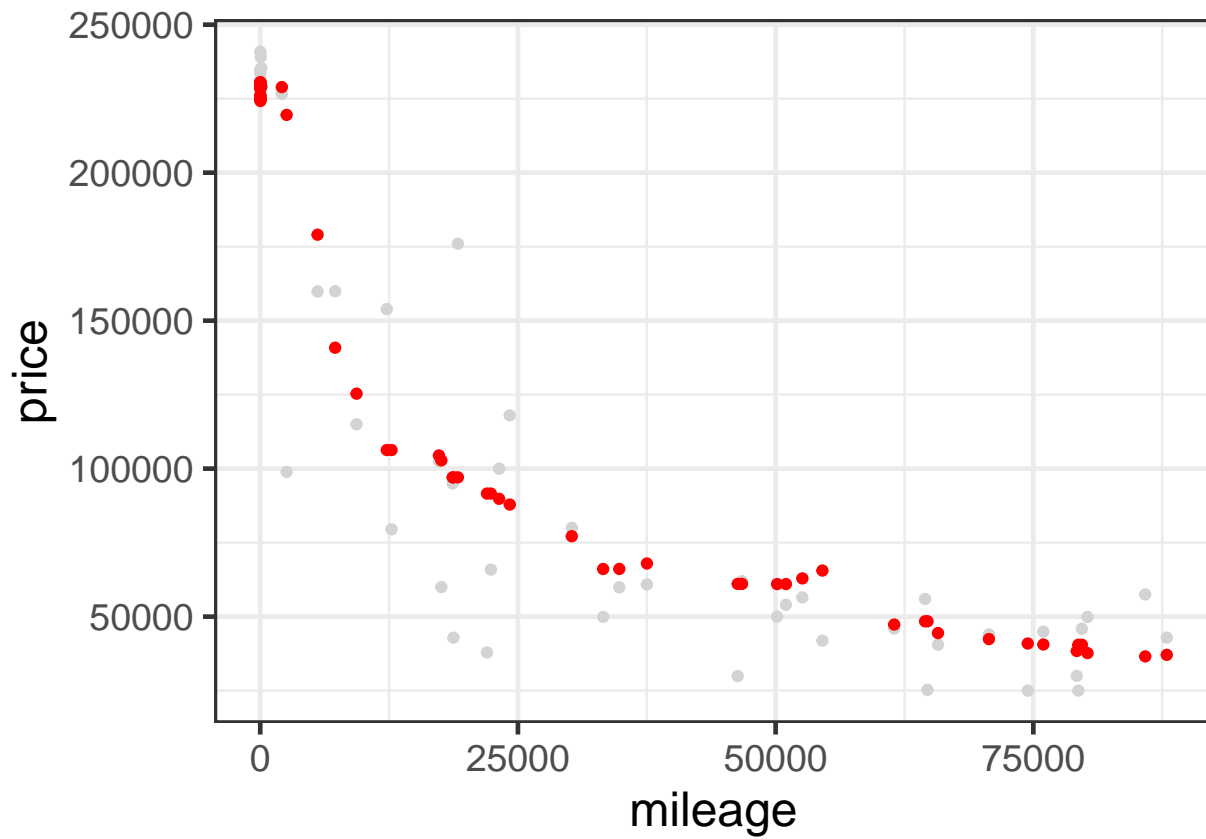
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn15 = ypred_knn15
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

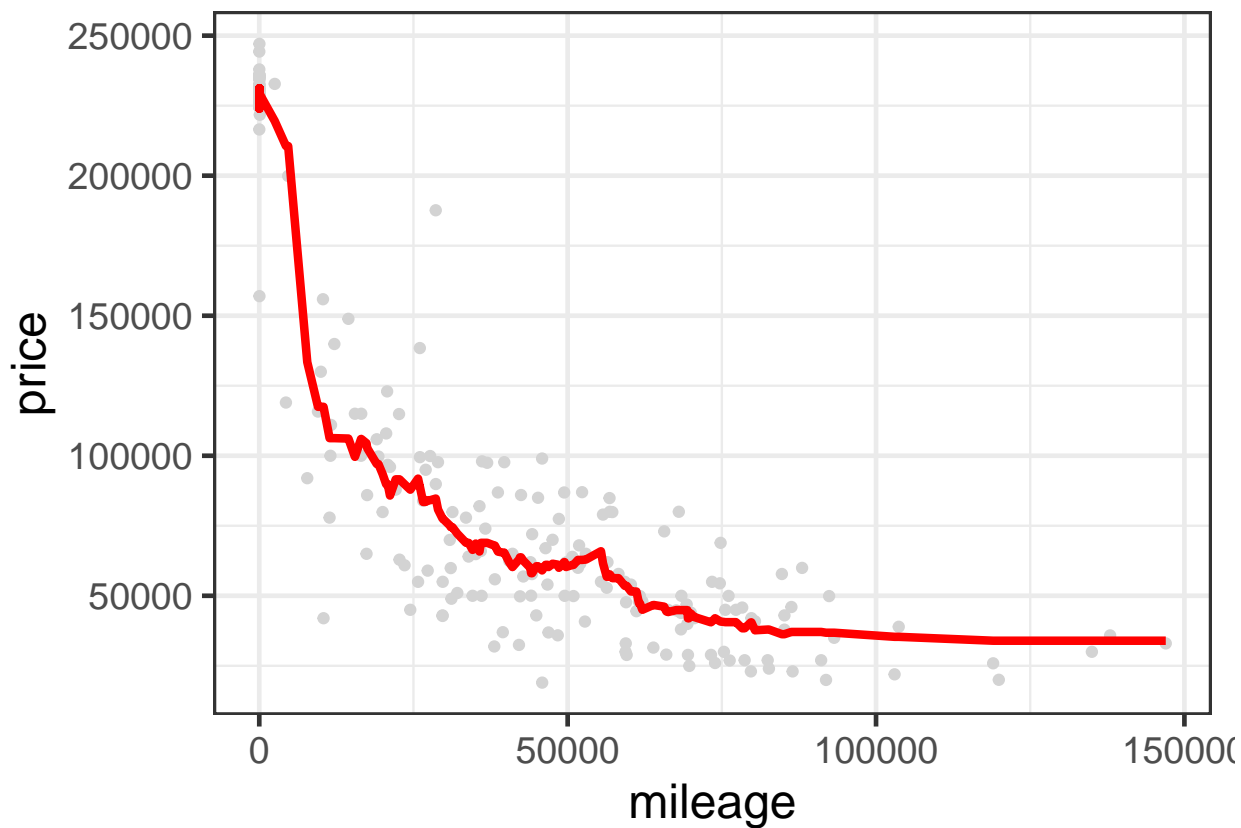


```
p_test + geom_point(aes(x = mileage, y = ypred_knn15), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 15)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 30
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn30 = knn.reg(train = X_train, test = X_test, y = y_train, k=30)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn30 = knn30$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 46463.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 34814.4
```

```
rmse(y_test, ypred_knn30)
```

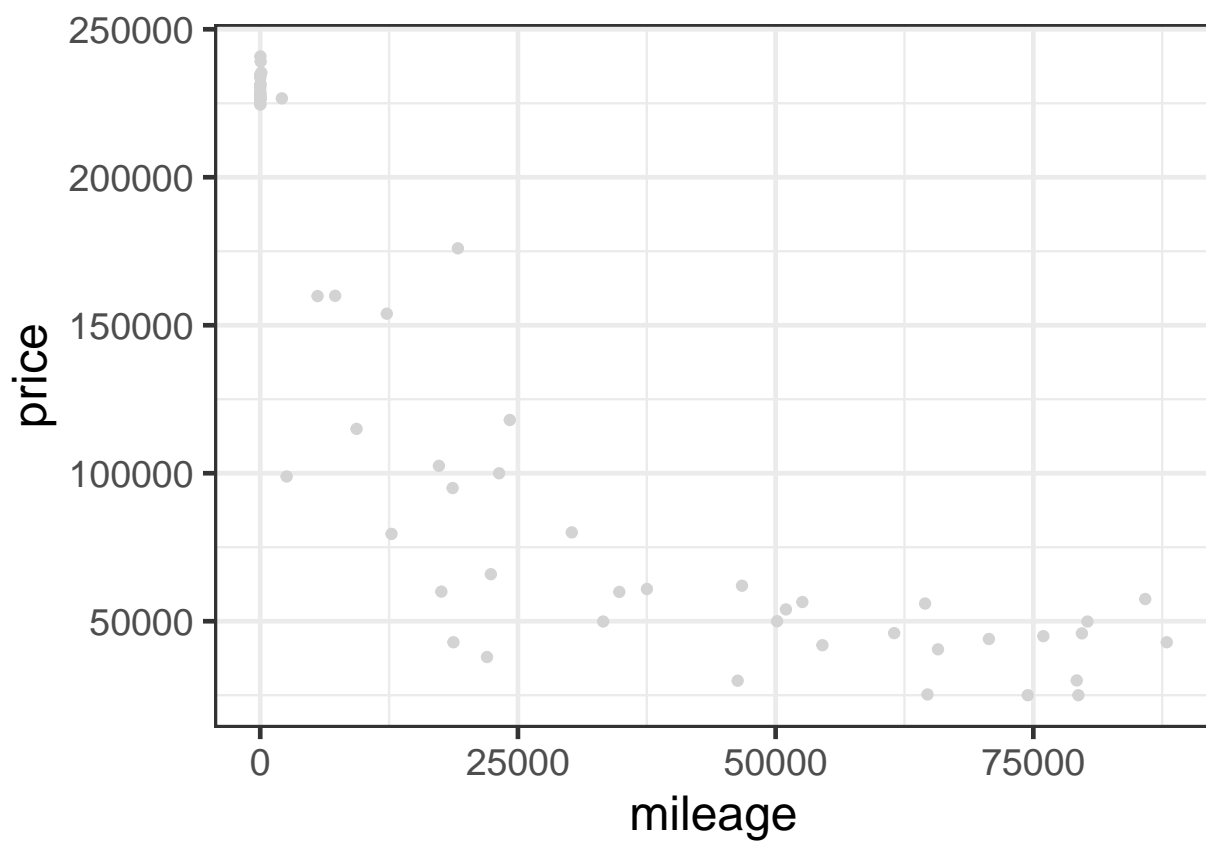
```
## [1] 27194.3
```

```
#attach predictions to data frame
```

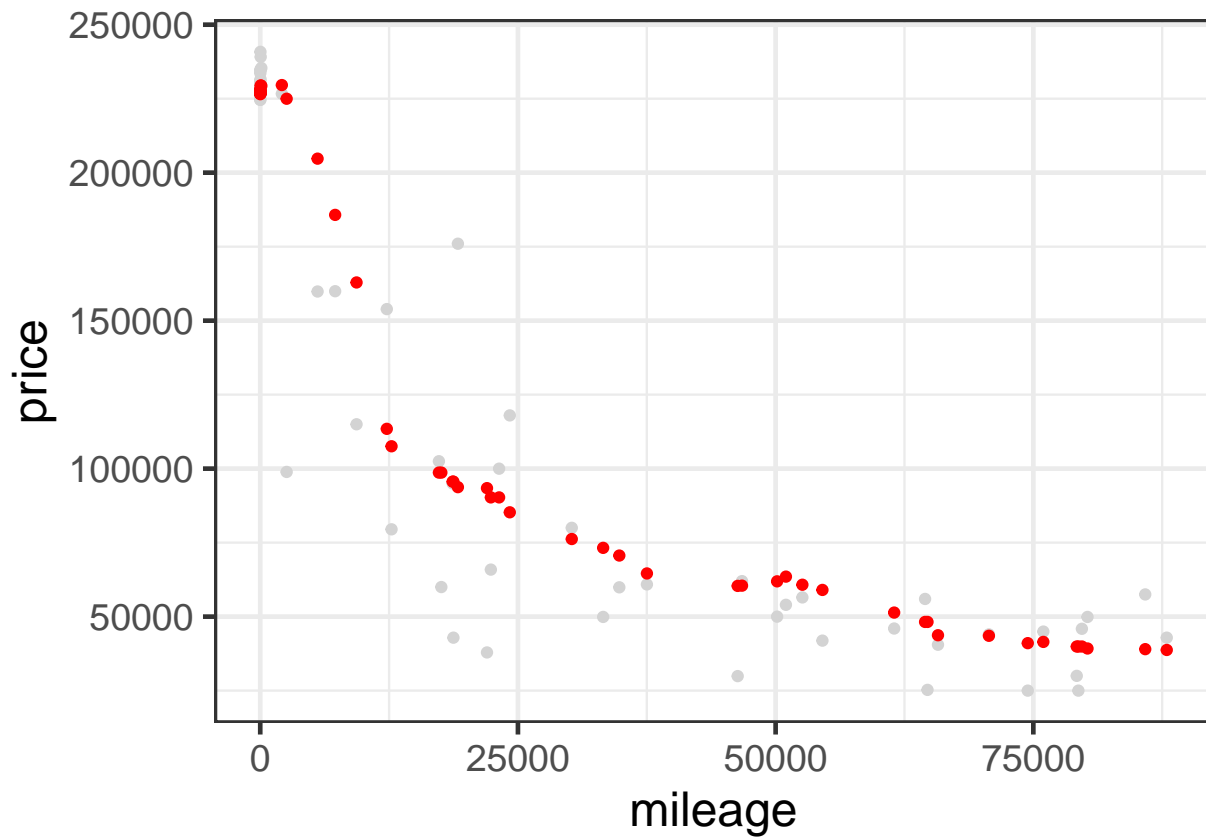
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn30 = ypred_knn30
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

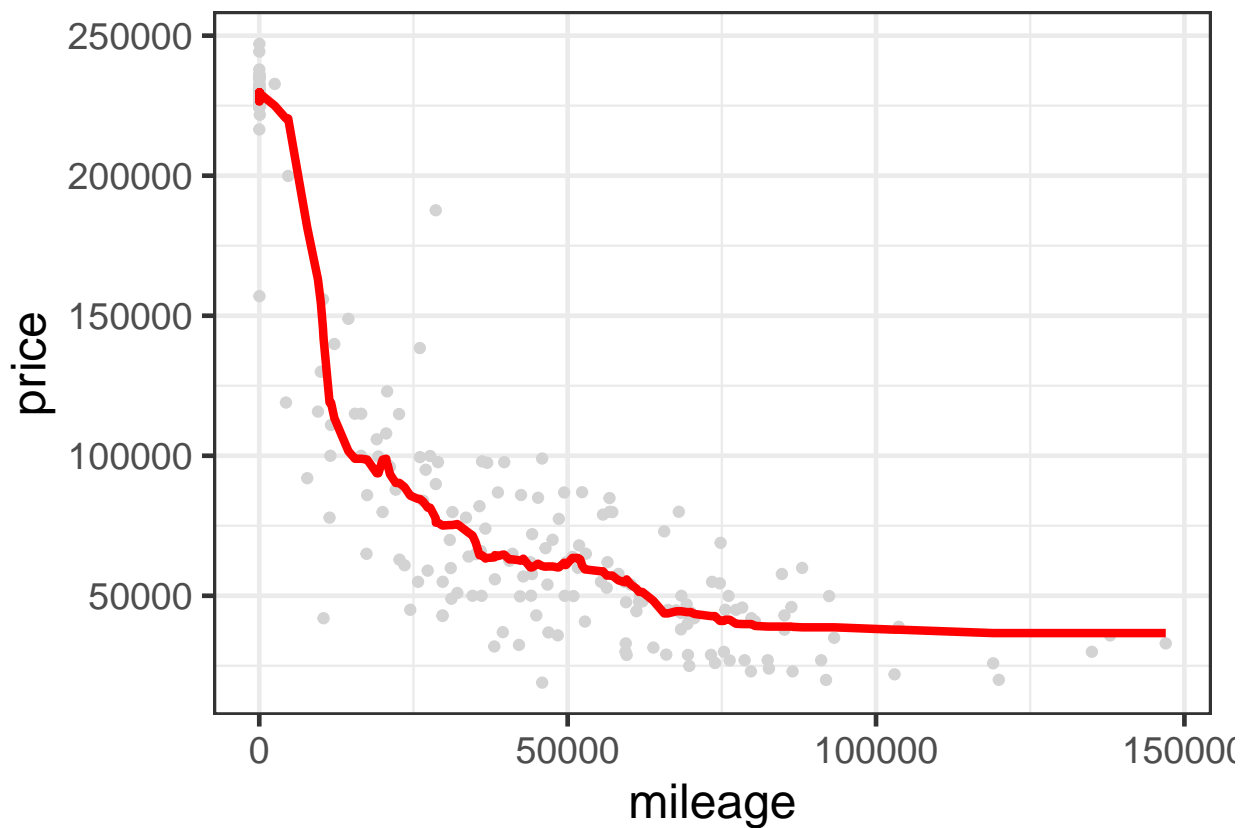


```
p_test + geom_point(aes(x = mileage, y = ypred_knn30), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 30)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 50
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn50 = knn.reg(train = X_train, test = X_test, y = y_train, k=50)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn50 = knn50$pred
```

```
rmse(y_test, ypred_lm1)
```

```
## [1] 46463.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 34814.4
```



```
rmse(y_test, ypred_knn50)
```

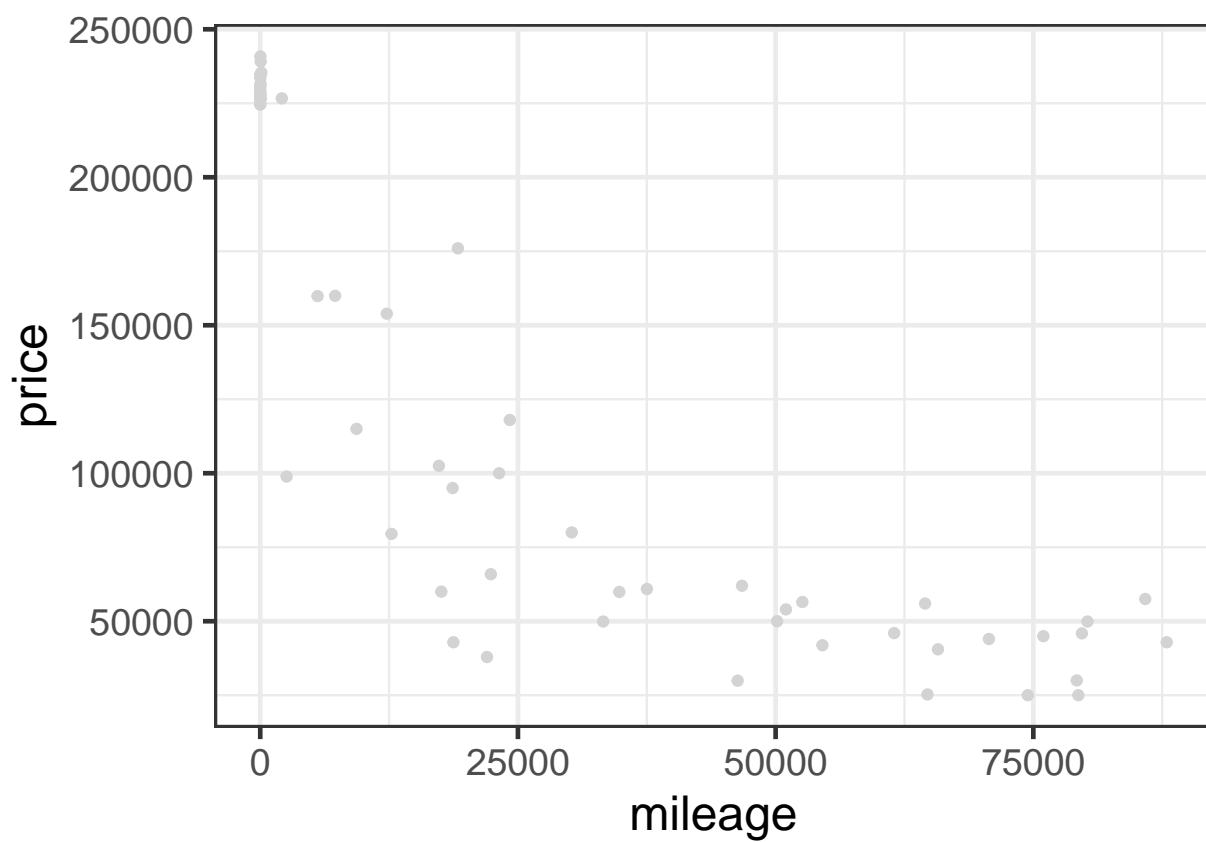
```
## [1] 29611.44
```

```
#attach predictions to data frame
```

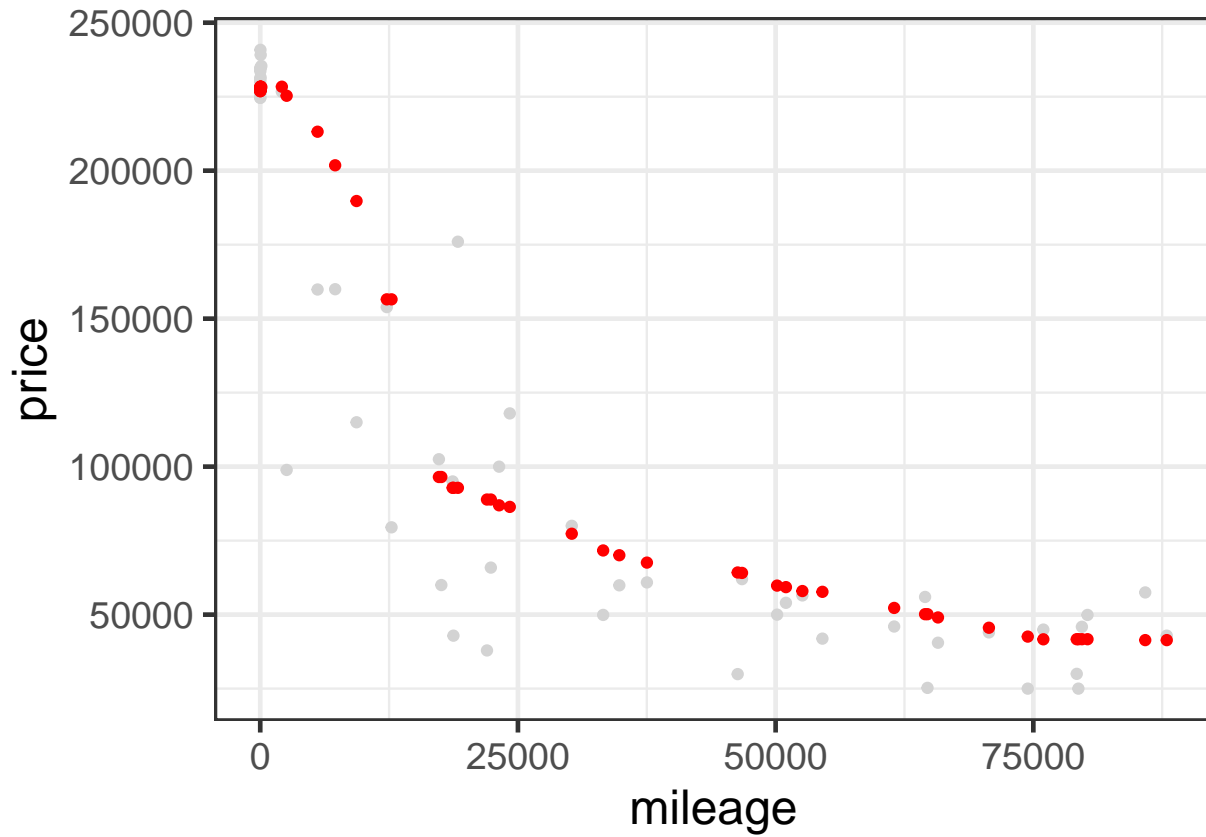
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn50 = ypred_knn50
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

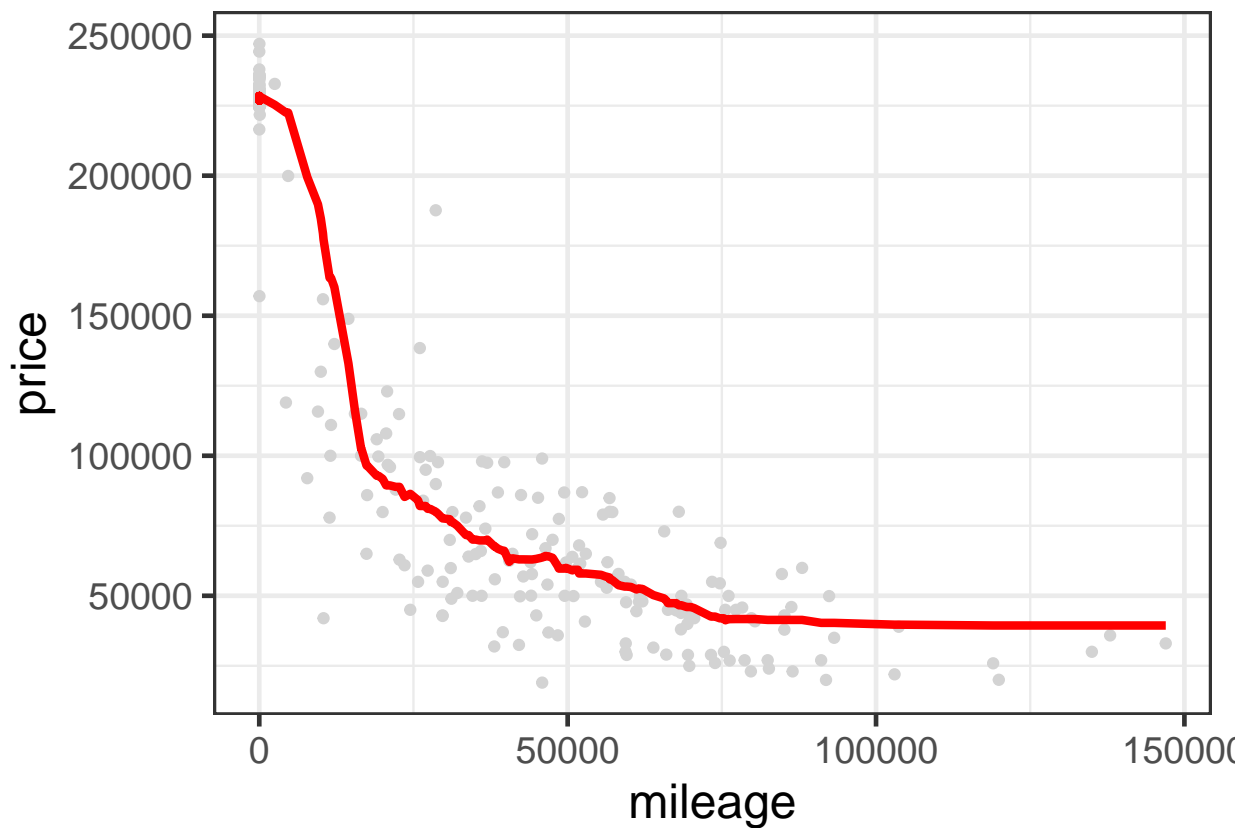


```
p_test + geom_point(aes(x = mileage, y = ypred_knn50), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 50)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 60
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn60 = knn.reg(train = X_train, test = X_test, y = y_train, k=60)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn60 = knn60$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 46463.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 34814.4
```

```
rmse(y_test, ypred_knn60)
```

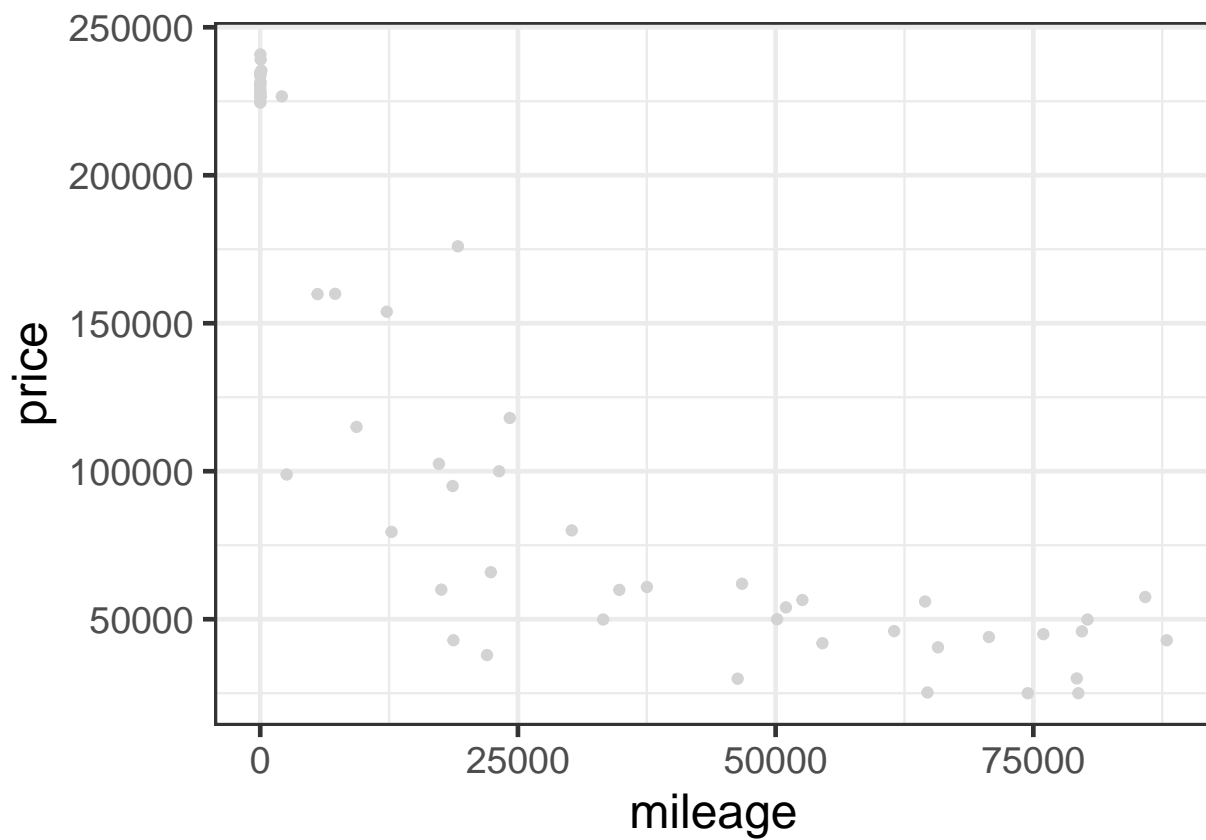
```
## [1] 31321.59
```

```
#attach predictions to data frame
```

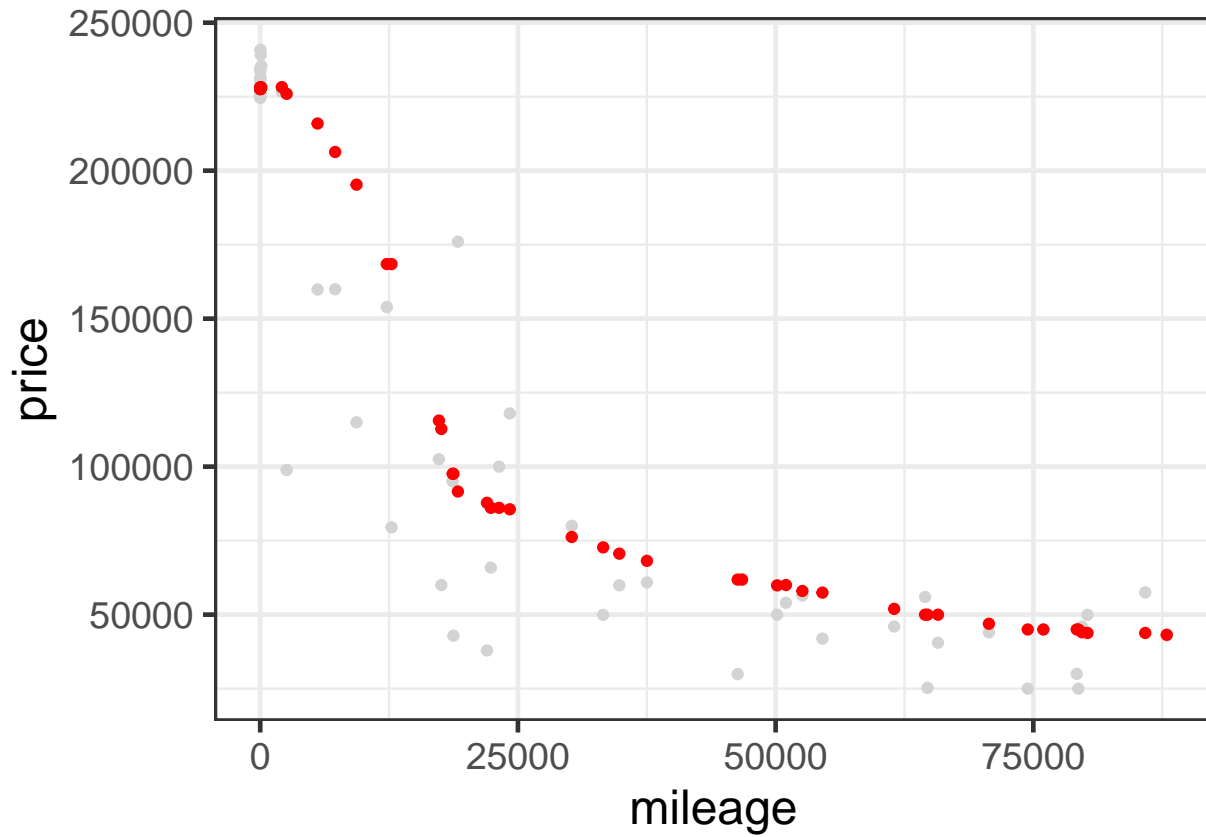
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn60 = ypred_knn60
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

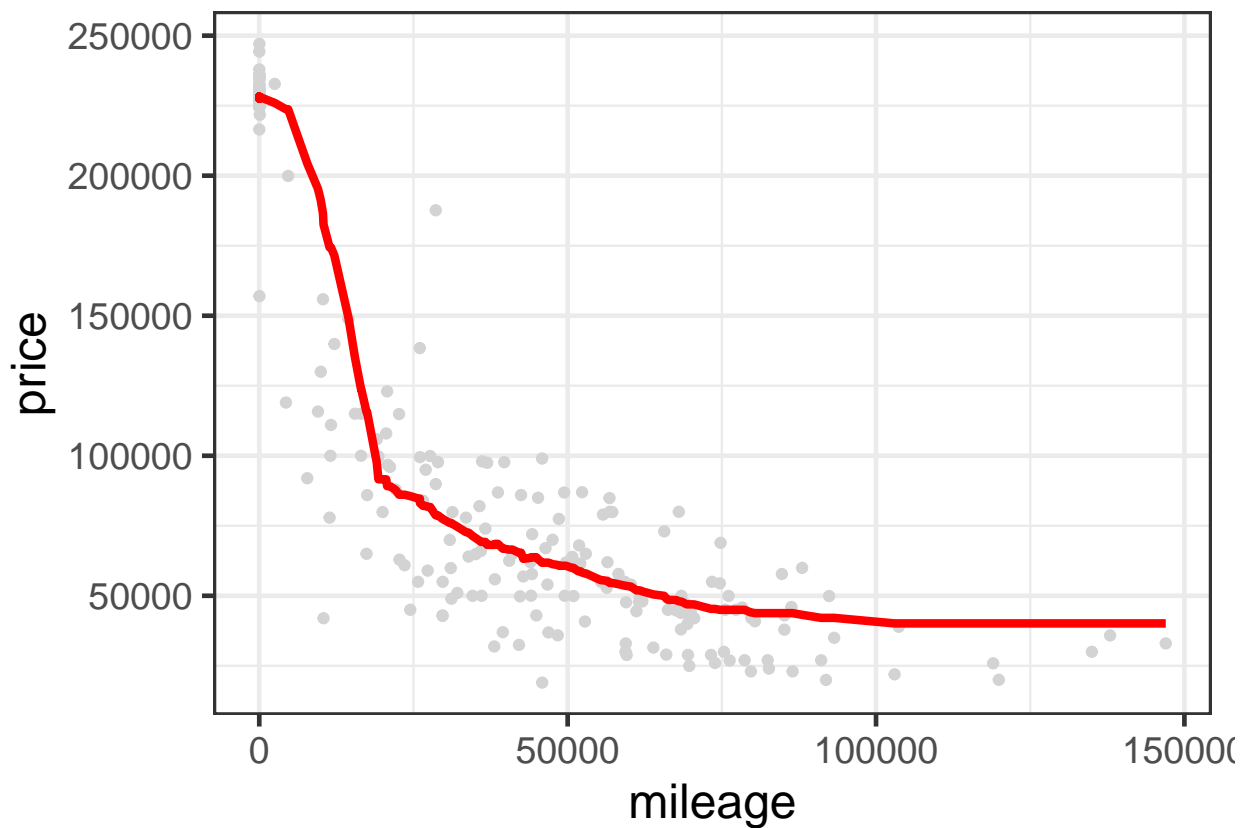


```
p_test + geom_point(aes(x = mileage, y = ypred_knn60), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 60)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



```
#K = 80
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn80 = knn.reg(train = X_train, test = X_test, y = y_train, k=80)
```

```
#rmse
rmse = function(y, ypred) {
  sqrt(mean(data.matrix((y-ypred)^2)))
}
```

```
ypred_lm1 = predict(lm1, X_test)
ypred_lm2 = predict(lm2, X_test)
ypred_knn80 = knn80$pred

rmse(y_test, ypred_lm1)
```

```
## [1] 46463.64
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 34814.4
```

```
rmse(y_test, ypred_knn80)
```

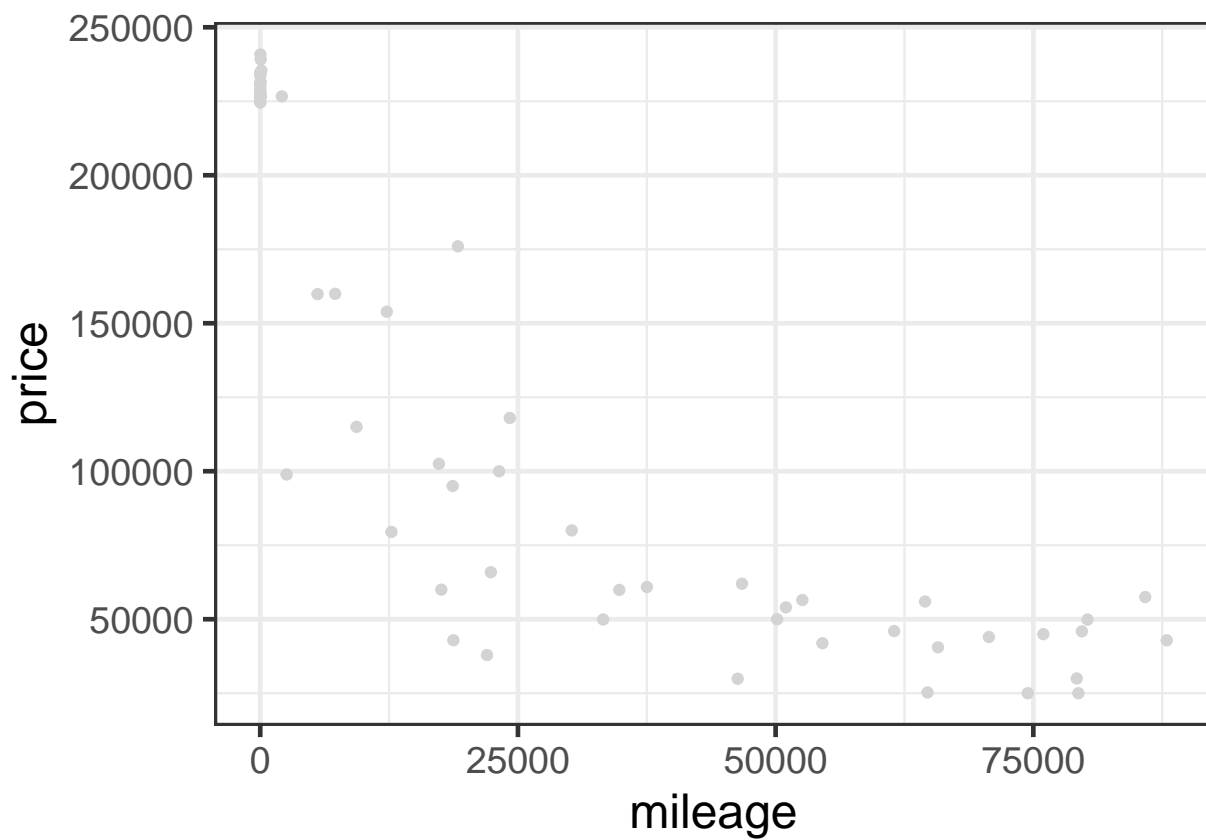
```
## [1] 35630.38
```

```
#attach predictions to data frame
```

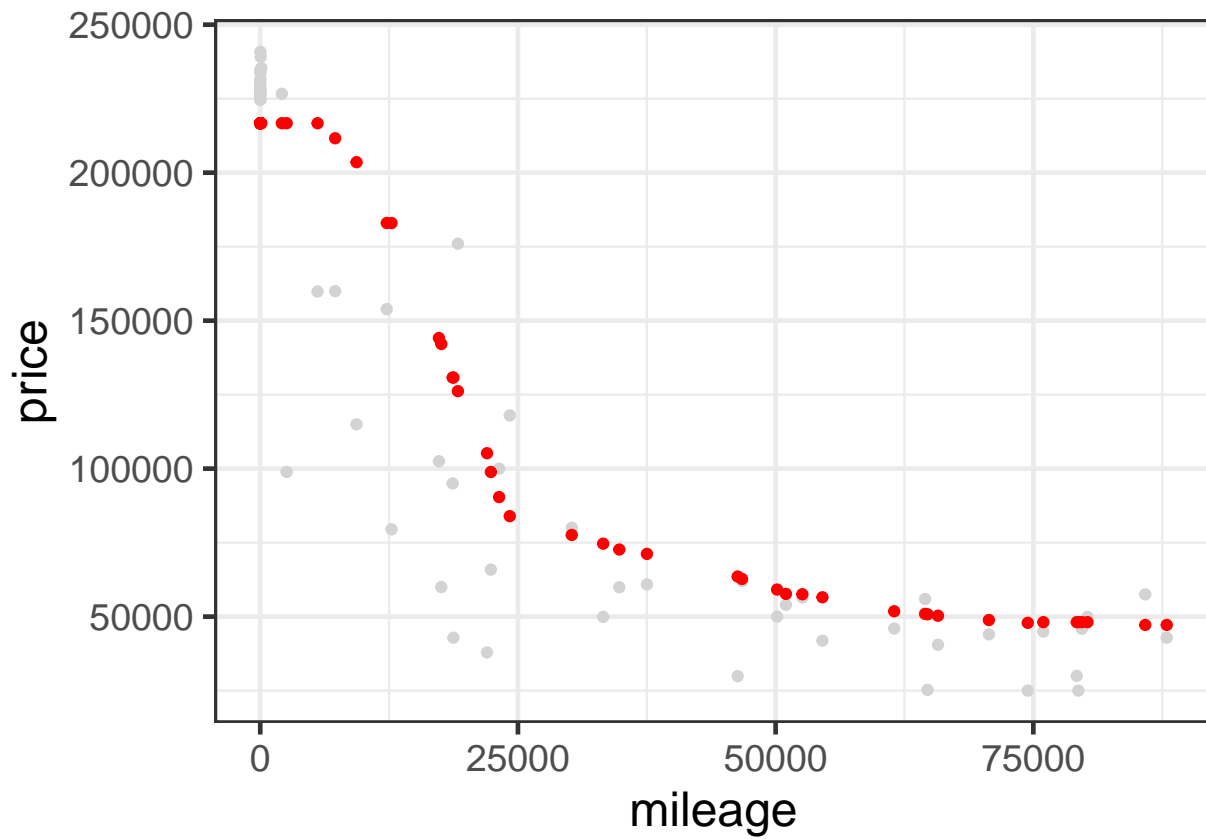
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn80 = ypred_knn80
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

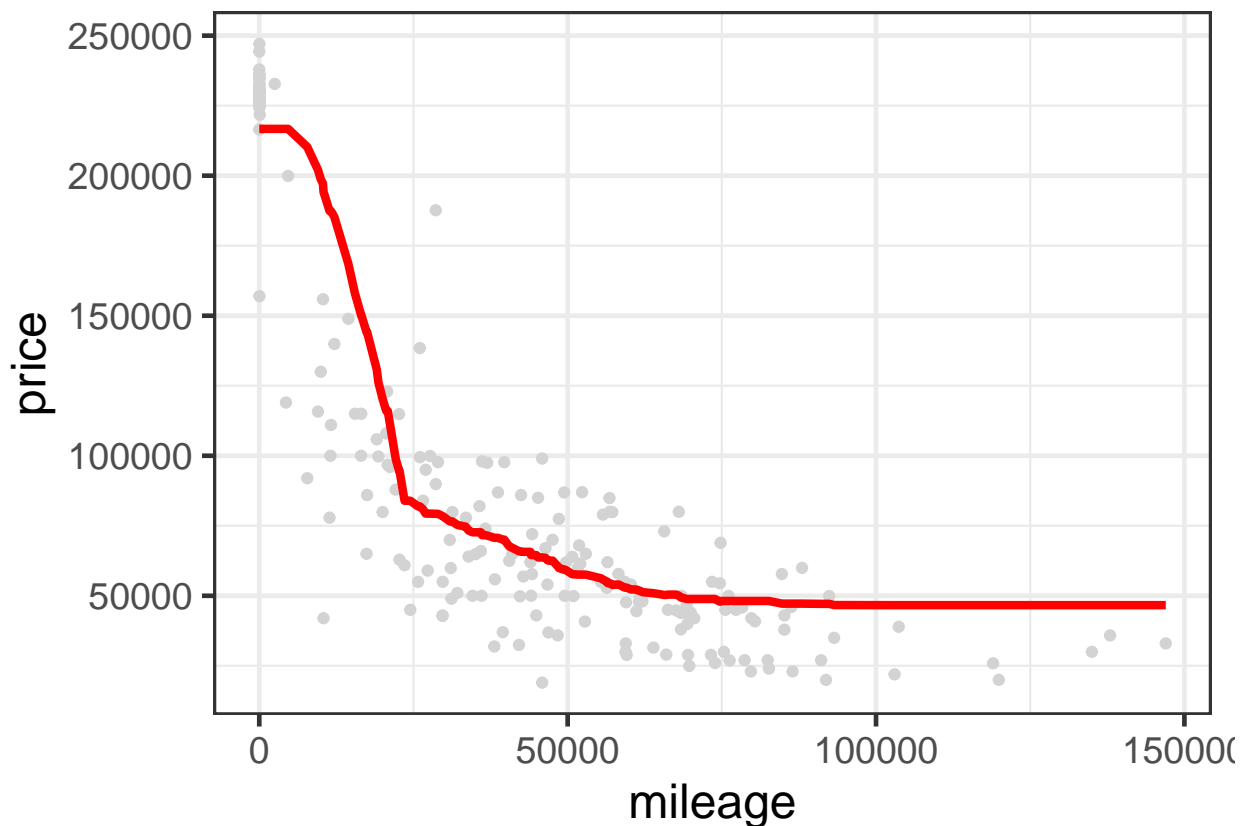


```
p_test + geom_point(aes(x = mileage, y = ypred_knn80), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 80)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```

```
#RMSE plot vs K nearest neighbors 65 AMG trim level.
sclass65AMG = subset(sclass, trim == '65 AMG')
N = nrow(sclass65AMG)
N_train = floor(0.8*N)
N_test = N - N_train

train_ind = sort(sample.int(N, N_train, replace=FALSE))

D_train = sclass65AMG[train_ind,]
D_train = arrange(D_train, mileage)
D_test = sclass65AMG[-train_ind,]

y_train = D_train$price
X_train = data.frame(mileage=jitter(D_train$mileage))
X_test = data.frame(mileage=jitter(D_test$mileage))
y_test = D_test$price
library(foreach)

#Optimal value for K for 65 AMG trim level is
#K = 80
lm1 = lm(price ~ mileage, data=D_train)
lm2 = lm(price ~ poly(mileage, 2), data=D_train)
knn80 = knn.reg(train = X_train, test = X_test, y = y_train, k=80)

#rmse
```

```
rmse = function(y, ypred) {  
  sqrt(mean(data.matrix((y-ypred)^2)))  
}
```

```
ypred_lm1 = predict(lm1, X_test)  
ypred_lm2 = predict(lm2, X_test)  
ypred_knn80 = knn80$pred
```

```
rmse(y_test, ypred_lm1)
```

```
## [1] 42407.7
```

```
rmse(y_test, ypred_lm2)
```

```
## [1] 30343.78
```

```
rmse(y_test, ypred_knn80)
```

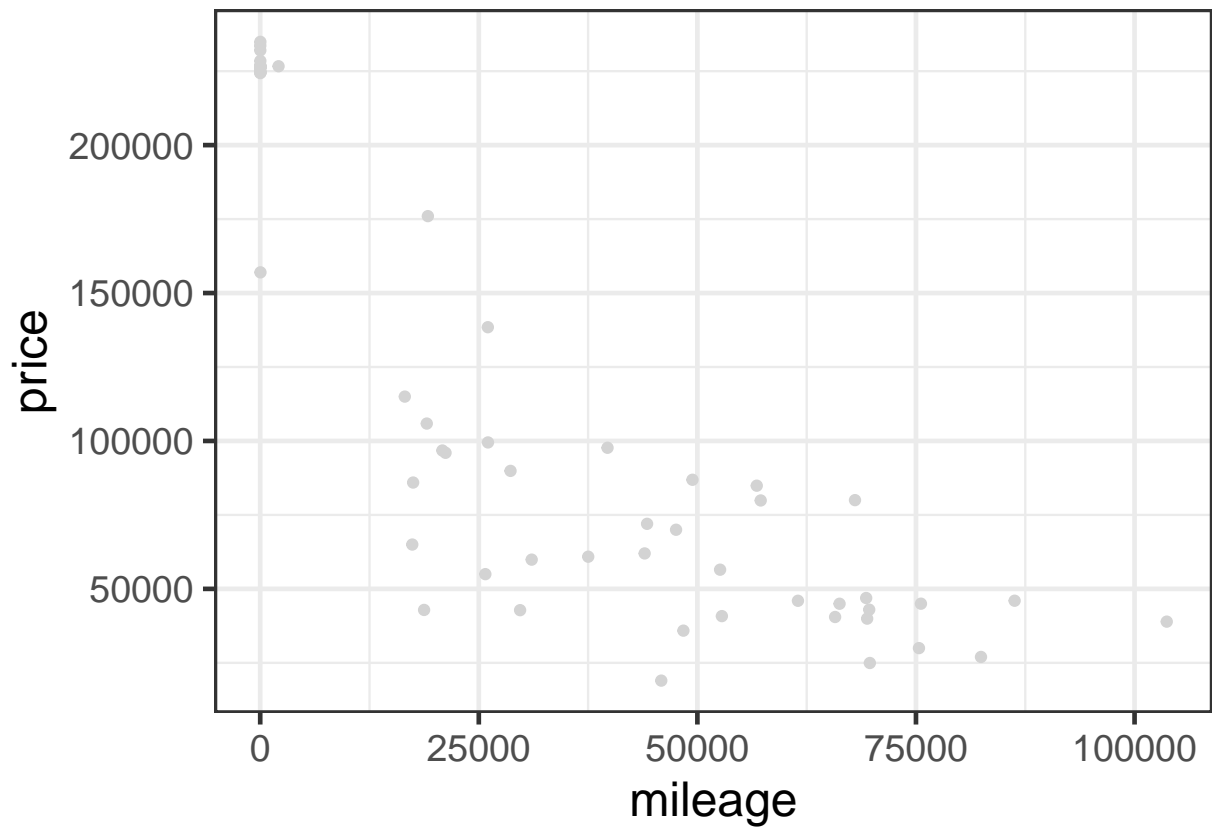
```
## [1] 26022.42
```

```
#attach predictions to data frame
```

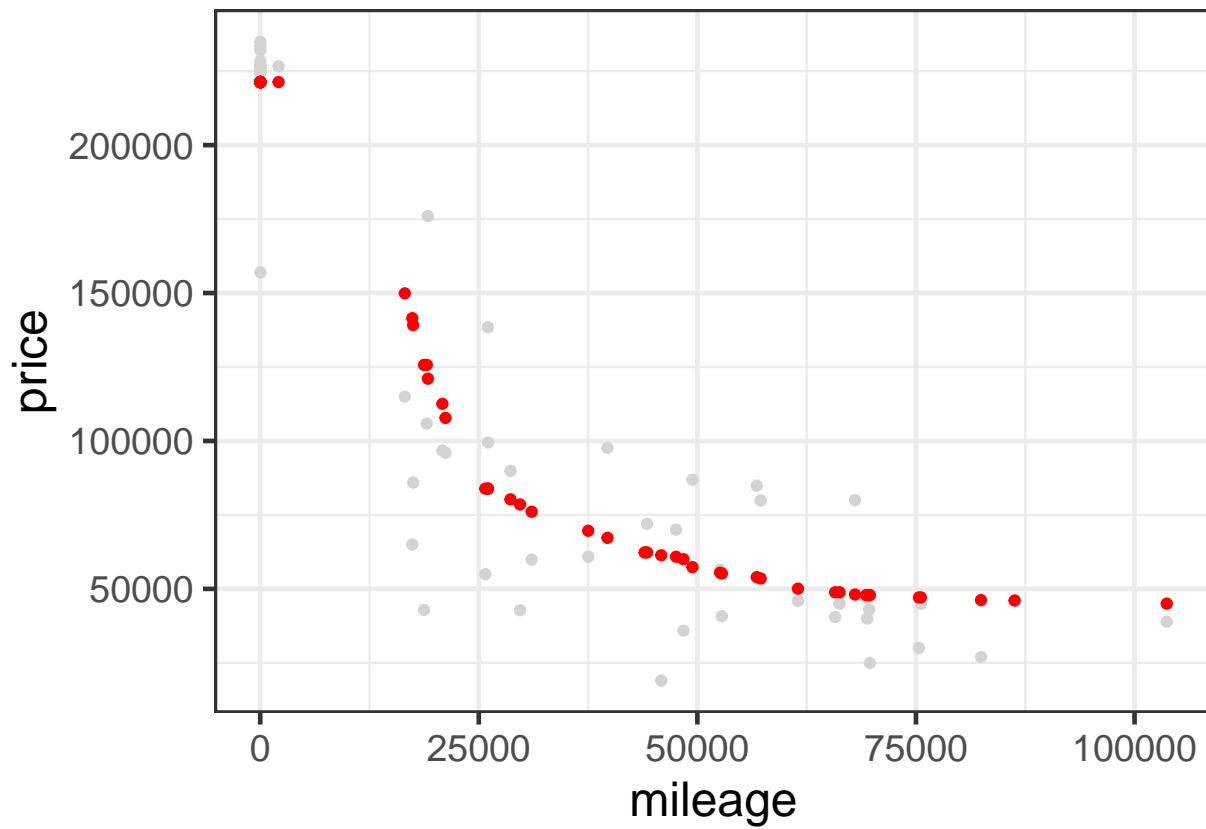
```
D_test$ypred_lm2 = ypred_lm2
```

```
D_test$ypred_knn80 = ypred_knn80
```

```
p_test = ggplot(data = D_test) +  
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +  
  theme_bw(base_size=18)  
p_test
```

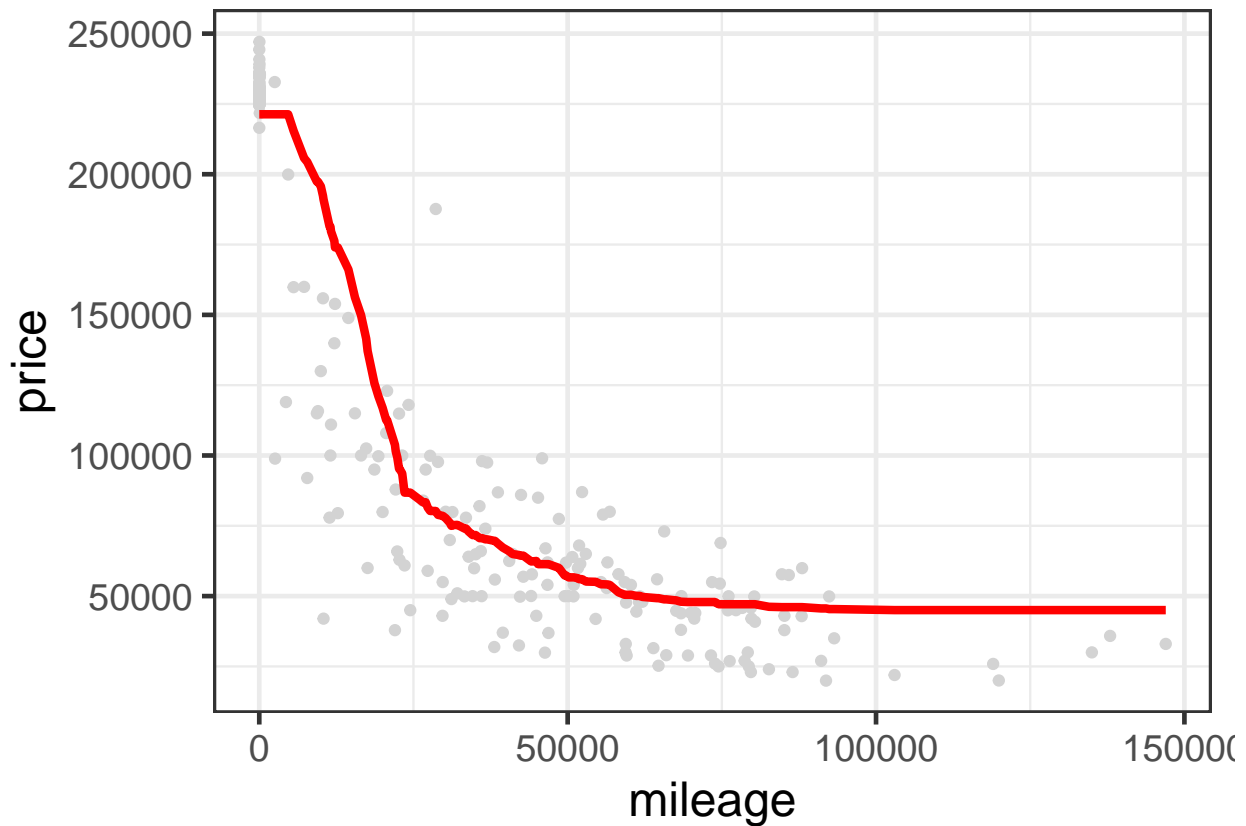


```
p_test + geom_point(aes(x = mileage, y = ypred_knn80), color='red')
```



```
#KNN variances
knn_model = knn.reg(X_train, X_train, y_train, k = 80)

D_train$ypred = knn_model$pred
p_train = ggplot(data = D_train) +
  geom_point(mapping = aes(x = mileage, y = price), color='lightgrey') +
  theme_bw(base_size=18)
p_train + geom_path(mapping = aes(x=mileage, y=ypred), color='red', size=1.5)
```



#The optimal value of K for the 65 AMG trim level is K=80. With this value of K, the rmse is 32901.62, and the rmse for the 350 trim level is 32901.62.

#Conclusion: The 65 AMG level has an optimal value of K that is higher than the 350 trim level optimal value.

Problem 2 - Saratoga Houses

Question

This analysis tries to examine and identify key indicators in home prices in Saratoga, New York. Homes have a nearly infinite number of attributes that may contribute to the value of the home, and this report seeks to find the most relevant for predicting what the value of a home is.

Data

```
library(tidyverse)
library(mosaic)
data(SaratogaHouses)

summary(SaratogaHouses)
```

```
##      price      lotSize      age      landValue
## Min.   : 5000   Min.   : 0.0000   Min.   : 0.00   Min.   : 200
```

```
## 1st Qu.:145000 1st Qu.: 0.1700 1st Qu.: 13.00 1st Qu.: 15100
## Median :189900 Median : 0.3700 Median : 19.00 Median : 25000
## Mean :211967 Mean : 0.5002 Mean : 27.92 Mean : 34557
## 3rd Qu.:259000 3rd Qu.: 0.5400 3rd Qu.: 34.00 3rd Qu.: 40200
## Max. :775000 Max. :12.2000 Max. :225.00 Max. :412600
## livingArea pctCollege bedrooms fireplaces bathrooms
## Min. : 616 Min. :20.00 Min. :1.000 Min. :0.0000 Min. :0.0
## 1st Qu.:1300 1st Qu.:52.00 1st Qu.:3.000 1st Qu.:0.0000 1st Qu.:1.5
## Median :1634 Median :57.00 Median :3.000 Median :1.0000 Median :2.0
## Mean :1755 Mean :55.57 Mean :3.155 Mean :0.6019 Mean :1.9
## 3rd Qu.:2138 3rd Qu.:64.00 3rd Qu.:4.000 3rd Qu.:1.0000 3rd Qu.:2.5
## Max. :5228 Max. :82.00 Max. :7.000 Max. :4.0000 Max. :4.5
## rooms heating fuel
## Min. : 2.000 hot air :1121 gas :1197
## 1st Qu.: 5.000 hot water/steam: 302 electric: 315
## Median : 7.000 electric : 305 oil : 216
## Mean : 7.042
## 3rd Qu.: 8.250
## Max. :12.000
## sewer waterfront newConstruction centralAir
## septic : 503 Yes: 15 Yes: 81 Yes: 635
## public/commercial:1213 No :1713 No :1647 No :1093
## none : 12
##
##
##
```

There are 16 variables in a dataset of 1,728 observations. The objective is to predict the variable “Price” using some combination of the other 15 variables.

Method

```
library(mosaic)

n = nrow(SaratogaHouses)
n_train = round(0.8*n)
n_test = n - n_train

rmse = function(y, yhat) {
  sqrt( mean( (y - yhat)^2 ) )
}

rmse_vals = do(100)*[

  train_cases = sample.int(n, n_train, replace=FALSE)
  test_cases = setdiff(1:n, train_cases)
  saratoga_train = SaratogaHouses[train_cases,]
  saratoga_test = SaratogaHouses[test_cases,]

  lm1 = lm(price ~ lotSize + bedrooms + bathrooms, data=saratoga_train)
  lm2 = lm(price ~ . - sewer - waterfront - landValue - newConstruction, data=saratoga_train)
  lm3 = lm(price ~ (. - sewer - waterfront - landValue - newConstruction)^2, data=saratoga_train)
```

```

lm4 = lm(price ~ . - sewer - landValue - newConstruction, data=saratoga_train)
lm5 = lm(price ~ . - sewer - waterfront - landValue, data=saratoga_train)
lm6 = lm(price ~ . - sewer - landValue, data=saratoga_train)
lm7 = lm(price ~ lotSize + age + livingArea + pctCollege +
          bedrooms + fireplaces + bathrooms + rooms + heating + fuel +
          centralAir + lotSize:heating + livingArea:rooms + newConstruction + livingArea:newConstruction, data=saratoga_train)
lm8 = lm(price ~ (lotSize + age + livingArea + pctCollege +
          bedrooms + fireplaces + bathrooms + rooms + heating + fuel +
          centralAir + lotSize:heating + livingArea:rooms + newConstruction + livingArea:newConstruction)^2, data=saratoga_train)
lm9 = lm(price ~ . - sewer - landValue - rooms, data=saratoga_train)
lm10 = lm(price ~ . - sewer - landValue - bathrooms - bedrooms, data=saratoga_train)
lm11 = lm(price ~ (. - sewer - landValue)^2, data=saratoga_train)

yhat_test1 = predict(lm1, saratoga_test)
yhat_test2 = predict(lm2, saratoga_test)
yhat_test3 = predict(lm3, saratoga_test)
yhat_test4 = predict(lm4, saratoga_test)
yhat_test5 = predict(lm5, saratoga_test)
yhat_test6 = predict(lm6, saratoga_test)
yhat_test7 = predict(lm7, saratoga_test)
yhat_test8 = predict(lm8, saratoga_test)
yhat_test9 = predict(lm9, saratoga_test)
yhat_test10 = predict(lm10, saratoga_test)
yhat_test11 = predict(lm11, saratoga_test)

c(rmse(saratoga_test$price, yhat_test1),
  rmse(saratoga_test$price, yhat_test2),
  rmse(saratoga_test$price, yhat_test3),
  rmse(saratoga_test$price, yhat_test4),
  rmse(saratoga_test$price, yhat_test5),
  rmse(saratoga_test$price, yhat_test6),
  rmse(saratoga_test$price, yhat_test7),
  rmse(saratoga_test$price, yhat_test8),
  rmse(saratoga_test$price, yhat_test9),
  rmse(saratoga_test$price, yhat_test10),
  rmse(saratoga_test$price, yhat_test11))
}

```

To find the best model for predicting home price, we create a train-test split, and use the training data to create models and test their accuracy by calculating the Root Mean Square Error (RMSE) of each model.

Results

```
rmse_vals
```

##	V1	V2	V3	V4	V5	V6	V7	V8
## 1	78394.63	67536.23	69939.49	65246.67	67460.27	65160.95	67587.02	70561.00
## 2	74855.03	61709.06	63743.08	60870.53	61495.59	60665.60	60838.45	98375.30
## 3	79003.51	67938.71	68380.35	67587.67	67617.61	67327.69	67592.28	66853.64
## 4	83257.85	71764.64	76215.13	67467.14	71665.04	67425.29	71530.80	84258.55
## 5	87205.86	72407.08	82180.75	67912.58	72358.55	67921.34	73021.63	115431.85

## 6	75602.58	67467.93	66437.47	66379.06	67390.44	66321.64	66929.12	64389.28
## 7	81742.16	68648.51	72987.88	69026.13	68304.43	68756.19	68537.79	73244.74
## 8	74332.27	61398.06	66405.32	60660.79	61563.11	60782.07	60630.61	71088.41
## 9	70535.02	62820.18	67743.71	60632.36	63037.30	60934.47	62454.54	78018.21
## 10	81092.11	68662.29	74235.22	65568.03	68322.04	65295.84	67778.32	73816.76
## 11	76461.01	63743.64	118631.18	63824.44	63712.56	63793.13	63177.55	112698.22
## 12	86106.10	74821.02	74817.69	71074.01	74707.36	71005.40	74193.87	75351.64
## 13	82731.24	66505.70	68428.86	66010.71	66482.86	66029.48	66003.96	78105.40
## 14	72273.43	65189.86	64860.65	61425.79	65171.32	61452.06	65319.61	68659.04
## 15	79865.81	69810.29	67084.38	70818.99	69734.66	70709.27	69157.48	85241.26
## 16	73716.57	62699.53	64321.26	62193.79	62632.22	62095.11	64025.42	68499.87
## 17	86010.91	79563.36	79265.14	77439.42	79368.35	77274.06	78949.39	87457.59
## 18	71489.57	62952.21	61250.99	63501.25	62897.05	63475.45	63514.47	65363.35
## 19	66267.91	57867.18	56912.73	56984.53	57595.30	56754.22	57786.63	66326.04
## 20	84313.74	73546.90	73524.07	71594.00	73378.00	71460.79	73342.96	76835.20
## 21	85071.85	70403.50	68426.64	65363.17	70152.24	65164.01	70099.42	71479.13
## 22	79264.15	70826.35	71769.60	68943.26	70555.06	68712.91	70255.52	71207.43
## 23	75222.71	70270.52	68593.61	67238.32	69959.81	66979.73	69394.88	99056.25
## 24	73515.56	62864.24	63790.67	63483.99	62808.95	63466.45	64057.05	67401.36
## 25	81001.13	68819.45	69062.88	67956.45	68727.89	67885.51	68749.99	88619.95
## 26	77892.26	62669.66	61377.28	62157.48	62310.92	61865.91	61871.78	75084.77
## 27	77190.60	66171.25	75414.78	65432.17	65977.35	65252.44	65271.18	80158.73
## 28	79945.74	72638.72	73956.68	71619.31	72432.63	71429.44	71919.36	84417.64
## 29	78969.66	64622.61	65940.41	64863.96	64506.56	64754.67	63699.25	76647.29
## 30	77979.31	70235.22	71642.73	67060.33	69969.43	66832.45	70323.69	74388.64
## 31	78787.82	68734.13	70442.48	65368.21	68513.29	65187.60	67882.22	73666.18
## 32	73897.35	62501.05	60363.24	59399.82	62248.74	59193.50	61908.77	65229.87
## 33	82318.40	68411.88	68099.76	67865.52	68173.63	67667.52	68350.59	85207.72
## 34	72733.73	64282.36	64350.92	63449.59	63965.45	63182.04	63426.15	65270.37
## 35	72433.86	65123.70	68404.49	60719.53	64916.82	60562.86	65987.59	86910.34
## 36	81081.80	70909.40	68598.44	70121.45	70780.10	70026.36	70899.52	78139.19
## 37	77128.53	67747.53	70152.88	66172.14	67665.40	66113.56	67324.60	74699.51
## 38	74904.48	61490.75	68423.54	60546.39	61515.42	60576.74	61120.53	80322.26
## 39	76164.18	63299.07	61900.97	61186.33	62978.79	60912.58	62119.89	64548.04
## 40	80022.09	71519.00	75992.89	66321.98	71353.17	66204.35	72334.16	77925.39
## 41	79788.41	70169.69	67849.88	68612.61	69867.32	68363.31	69060.19	85525.13
## 42	71869.36	64250.02	68318.01	61795.43	64089.58	61684.32	63908.34	75640.46
## 43	75310.40	62566.37	67571.38	62620.50	62495.46	62552.77	62166.92	64898.56
## 44	81686.50	70046.35	70790.57	69744.43	70117.10	69814.47	69844.23	71363.37
## 45	72334.15	59298.21	65344.76	59113.93	59237.68	59043.93	59690.46	67961.77
## 46	80326.09	69253.32	72239.91	68246.24	69011.58	68067.08	68246.73	388943.06
## 47	80869.38	66991.04	63955.82	64873.21	66853.86	64762.70	66536.65	64300.20
## 48	76438.56	67922.51	70895.34	67382.87	67987.03	67457.76	68004.15	79617.46
## 49	77738.96	69272.29	72838.43	66087.17	69122.72	65943.69	68867.61	77835.01
## 50	76948.44	64475.84	66359.30	64061.28	64453.52	64009.96	64880.20	64213.84
## 51	67739.32	54268.06	67826.81	53713.81	54221.63	53674.77	55426.28	65641.93
## 52	77600.79	65041.58	64048.99	62166.69	64928.20	62081.87	64808.25	65650.31
## 53	71929.34	64222.66	93592.71	63026.34	64269.59	63103.48	63982.79	75621.87
## 54	76133.44	65688.22	63586.93	63812.69	65456.50	63632.16	65423.04	87361.39
## 55	68690.90	59347.48	77930.10	57982.21	59382.52	58041.14	60986.76	70718.38
## 56	73745.46	62089.04	63030.79	60050.03	61938.84	59928.23	61223.50	66775.25
## 57	68100.01	60221.60	65667.67	59498.77	60025.66	59330.01	59935.17	75687.27
## 58	75490.93	62583.37	61326.61	60791.36	62375.70	60600.26	62390.26	67768.96
## 59	78614.12	68762.92	96535.60	67771.06	68684.96	67737.04	68665.93	73188.77

## 60	90291.96	78524.77	78789.38	74634.47	78262.89	74427.98	77989.11	77463.42
## 61	81063.34	71325.63	82860.33	69430.17	71031.51	69189.12	71014.59	87128.89
## 62	77993.48	63672.68	69305.80	62417.42	63445.38	62174.15	63425.19	77648.78
## 63	85037.15	72535.38	73401.45	72753.41	72347.21	72581.30	72495.50	71847.44
## 64	78717.35	69765.74	74173.77	67226.88	69506.91	67005.29	68918.49	76325.67
## 65	71576.62	61877.52	60951.40	61929.93	61788.91	61827.89	61069.34	78869.17
## 66	77333.95	63829.17	77475.39	64559.13	63615.07	64346.25	63227.77	80873.90
## 67	71450.76	63804.26	65115.63	62771.26	63517.92	62508.31	62972.17	66825.64
## 68	76993.54	64674.95	68536.92	64780.30	64754.00	64829.42	64694.72	81162.79
## 69	84654.72	75335.45	77153.96	73047.07	75041.24	72808.98	76473.76	84650.35
## 70	86500.11	74881.24	75116.66	74499.29	74573.90	74245.61	74001.44	79057.38
## 71	81083.69	72793.68	265343.46	72383.50	72557.64	72175.89	72014.83	218817.61
## 72	84365.56	73612.40	73334.51	73334.45	73639.12	73368.41	73478.33	82683.42
## 73	72872.81	63031.00	64430.48	61427.42	63022.82	61406.04	63115.59	79114.09
## 74	81763.92	72080.25	74417.62	67826.31	71800.95	67601.31	71954.68	74921.43
## 75	78199.23	71137.95	70128.29	71031.06	71592.88	71397.65	71816.76	72039.20
## 76	78917.08	67544.78	76210.75	67021.52	67208.38	66745.11	67161.69	91941.03
## 77	79192.34	64255.50	64308.50	64114.51	64064.71	63937.03	62975.31	65758.84
## 78	82355.20	71551.39	69075.47	68804.37	71448.94	68756.08	70857.05	89981.62
## 79	78647.80	69916.39	68374.27	67442.53	69654.32	67252.53	68989.15	68477.84
## 80	78541.21	65611.09	66028.70	64675.75	65519.80	64596.57	64978.06	75473.91
## 81	78511.33	64720.87	65100.73	62838.73	64619.09	62763.51	64500.82	71783.62
## 82	81158.30	68022.13	72456.61	63988.23	67799.73	63831.24	67597.48	75697.79
## 83	82701.67	71133.23	77453.50	67358.68	71208.53	67464.08	71817.05	137286.38
## 84	79261.50	67073.02	65098.41	65790.50	66874.84	65650.08	66340.83	64285.85
## 85	73219.74	60933.43	62851.53	60113.82	60726.56	59962.74	60489.86	61506.05
## 86	70148.46	57440.97	61474.21	56592.61	57371.74	56497.17	57126.74	71842.42
## 87	82243.87	72681.13	76269.27	72483.45	72536.00	72322.14	72722.15	89867.90
## 88	80037.98	70284.68	143923.33	69274.65	70243.82	69201.96	70210.24	89877.90
## 89	79459.27	66135.14	65222.87	65521.97	66055.97	65438.30	65625.67	76844.05
## 90	78148.33	60847.32	89049.76	60491.28	60602.14	60269.86	62521.08	95377.14
## 91	76230.14	70498.04	69988.75	70287.90	70487.23	70275.12	69641.94	147540.45
## 92	89084.73	77123.97	73422.60	72617.73	76985.24	72536.39	76289.29	78565.23
## 93	77095.59	67062.37	67530.36	62475.32	67121.45	62603.41	66432.61	73365.00
## 94	74309.54	61972.37	74665.94	59741.40	61699.09	59553.82	61201.68	78796.50
## 95	75540.06	67514.49	67027.10	68216.20	67638.24	68427.95	67501.76	78955.98
## 96	67105.45	58830.46	70377.99	58847.55	58549.06	58615.76	58036.56	65080.20
## 97	71261.04	59621.76	63308.76	57198.73	59840.21	57443.67	59460.15	62172.32
## 98	85307.77	76751.66	96012.96	73288.73	76623.46	73202.70	77114.24	78644.42
## 99	84246.32	72720.25	70459.02	69064.54	72473.91	68876.69	71872.83	74658.96
## 100	75252.77	62951.34	60078.14	61284.90	62955.71	61286.11	62622.69	65305.88
##	V9	V10	V11					
## 1	65716.49	67069.59	315022.04					
## 2	61026.54	60984.55	84323.33					
## 3	67218.43	68741.90	3554456.98					
## 4	67949.36	68619.01	224448.89					
## 5	68076.52	68135.42	235200.09					
## 6	66706.78	66230.74	944316.50					
## 7	69275.24	70216.31	80547.96					
## 8	60862.48	61663.36	66627.45					
## 9	61096.42	61957.07	283937.50					
## 10	65718.79	66647.86	86061.71					
## 11	63791.39	64777.90	119777.55					
## 12	71147.97	71051.50	192711.99					

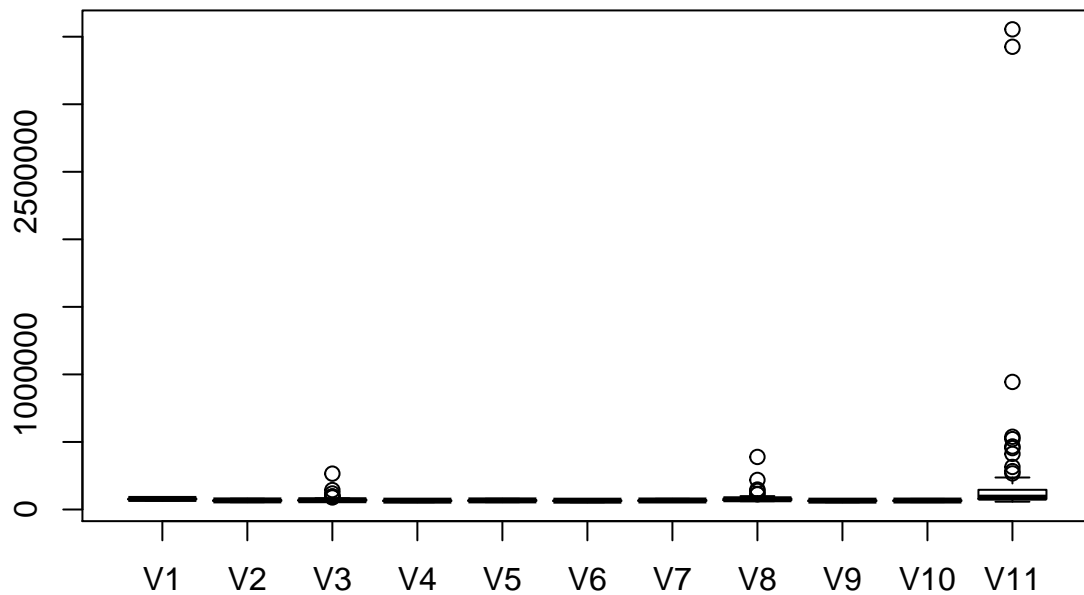
## 13	66327.54	67469.83	83813.85
## 14	61563.51	62122.65	125204.09
## 15	70806.50	72159.08	71559.96
## 16	62150.13	62490.31	90830.21
## 17	77485.35	78242.05	162340.15
## 18	63379.34	62444.59	69977.90
## 19	56889.49	57951.99	71479.66
## 20	71688.25	72950.96	105570.12
## 21	65111.15	65415.83	178545.19
## 22	69073.46	70181.45	80594.79
## 23	67242.21	68167.85	87721.05
## 24	63581.51	64777.37	70711.33
## 25	68088.24	69301.87	80121.49
## 26	61628.33	63349.73	83456.07
## 27	65205.89	67158.77	118839.34
## 28	71695.43	72830.74	75241.71
## 29	65120.08	65056.92	79330.74
## 30	67307.56	67590.27	87915.25
## 31	65234.71	65821.19	68591.41
## 32	59266.43	60046.51	81535.44
## 33	67957.00	68669.01	85779.51
## 34	63051.58	64071.46	107646.79
## 35	60920.83	61529.70	87416.23
## 36	70460.83	71686.73	81370.64
## 37	66432.07	66209.60	101462.46
## 38	60535.33	61269.37	85995.44
## 39	61038.11	62414.73	114041.38
## 40	65837.56	67729.73	153080.38
## 41	68683.62	68825.29	69873.16
## 42	61847.08	61795.60	168512.91
## 43	62991.04	64091.98	72762.87
## 44	70233.36	71811.89	75559.57
## 45	59219.19	60705.42	67381.69
## 46	67931.05	68262.67	87162.64
## 47	64869.31	65720.26	65765.02
## 48	67410.95	68268.82	70499.43
## 49	66374.85	66195.30	115517.87
## 50	64241.16	66182.76	79755.47
## 51	53673.35	54603.18	67796.03
## 52	62442.49	63400.41	179949.66
## 53	63565.01	64342.38	119523.82
## 54	63576.34	64549.34	74883.31
## 55	58335.68	57461.31	81572.33
## 56	60323.26	60445.60	80865.87
## 57	59549.56	61000.13	82774.00
## 58	60946.97	61892.40	125741.43
## 59	68089.33	68696.49	99085.59
## 60	74662.14	75309.49	78222.68
## 61	69447.98	70484.70	85633.65
## 62	62433.30	63517.67	267223.83
## 63	72754.70	74023.01	110870.67
## 64	66946.77	66611.43	83292.71
## 65	61941.32	63076.40	66596.22
## 66	64269.82	65004.53	81254.90

```
## 67 62115.86 63156.85 72532.77
## 68 65011.90 66721.56 84128.65
## 69 73185.47 74466.81 120617.47
## 70 74642.92 75278.53 99745.62
## 71 72491.45 73655.49 236733.53
## 72 73464.26 73909.43 140018.04
## 73 61283.29 61647.33 76427.36
## 74 67849.89 69111.88 183001.79
## 75 71597.68 72454.35 77204.94
## 76 66583.88 68161.90 81823.07
## 77 63754.24 63970.47 63251.51
## 78 68457.12 69690.46 454382.08
## 79 67631.76 69030.80 466596.74
## 80 65073.18 66075.47 204485.82
## 81 62771.61 63690.30 92006.53
## 82 64130.26 65035.40 162723.17
## 83 67629.69 68400.35 159687.08
## 84 65768.59 66108.49 105485.27
## 85 60298.14 61149.78 107955.84
## 86 56613.13 57328.11 66184.48
## 87 72702.64 73630.40 520908.44
## 88 69677.50 70032.32 143600.28
## 89 65873.06 65715.63 75220.83
## 90 60711.54 61340.74 148152.73
## 91 70110.78 70812.35 69064.92
## 92 72749.69 73620.94 411948.07
## 93 62505.14 62782.87 87332.27
## 94 59724.54 59320.17 143089.69
## 95 68333.25 68146.16 538394.63
## 96 58825.49 59169.80 69733.97
## 97 57367.42 57240.29 92606.15
## 98 73492.36 74351.26 105779.05
## 99 69376.94 69946.41 3425624.41
## 100 61350.59 61566.06 57676.83
```

```
colMeans(rmse_vals)
```

```
##      V1      V2      V3      V4      V5      V6      V7      V8
## 77900.13 67031.43 73300.84 65464.26 66897.32 65358.43 66761.65 82281.85
##      V9      V10     V11
## 65535.80 66282.03 203798.06
```

```
boxplot(rmse_vals)
```



```
Xtrain = model.matrix(~ . - (price + sewer + landValue) - 1, data=saratoga_train)
Xtest = model.matrix(~ . - (price + sewer + landValue) - 1, data=saratoga_test)
ytrain = saratoga_train$price
ytest = saratoga_test$price
```

```
scale_train = apply(Xtrain, 2, sd)
Xtilde_train = scale(Xtrain, scale = scale_train)
Xtilde_test = scale(Xtest, scale = scale_train)
```

```
head(Xtrain, 2)
```

```
##      lotSize age livingArea pctCollege bedrooms fireplaces bathrooms rooms
## 80      2.12 32      1440         52         3           1          1.0      8
## 1086    0.19 16      1480         54         2           1          1.5      6
##      heatinghot air heatinghot water/steam heatingelectric fuelelectric fueloil
## 80              1              0              0              0              1
## 1086            1              0              0              0              0
##      waterfrontNo newConstructionNo centralAirNo
## 80              1              1              1
## 1086            1              1              1
```

```
library(FNN)
K = 10
```

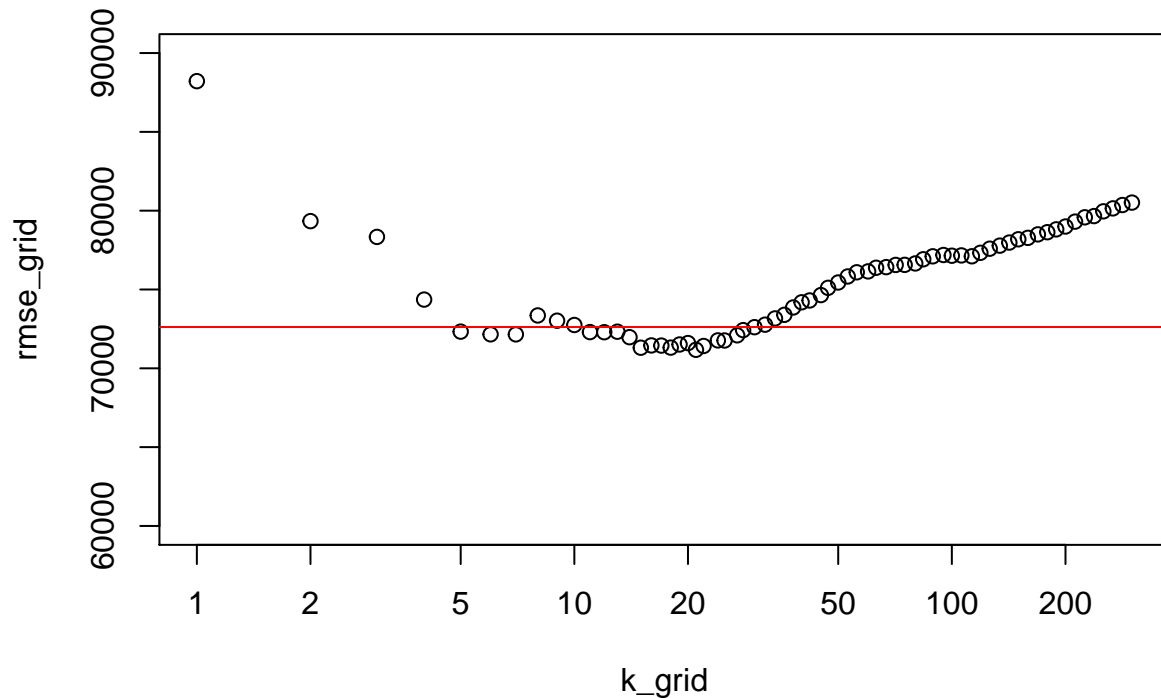
```
knn_model = knn.reg(Xtilde_train, Xtilde_test, ytrain, k=K)
```

```
rmse(ytest, knn_model$pred)
```

```
## [1] 72749.14
```

```
library(foreach)
k_grid = exp(seq(log(1), log(300), length=100)) %>% round %>% unique
rmse_grid = foreach(K = k_grid, .combine='c') %do% {
  knn_model = knn.reg(Xtilde_train, Xtilde_test, ytrain, k=K)
  rmse(ytest, knn_model$pred)
}

plot(k_grid, rmse_grid, log='x', ylim=c(60000, 90000))
abline(h=rmse(ytest, yhat_test6), col="red")
```



After running 100 train-tests splits, we find the average RMSE for each model. We find that Linear Model #6 is the best model because it has the lowest average RMSE. This model is a modification of the original, medium sized model which factors in Lot Size, age, living area, percent of college students in the area, number of bedrooms, number of fireplaces, number of bathrooms, number of rooms, heating type, fuel, and central air system. The modification adds the variables Waterfront and New Construction. Meanwhile, we exclude the land value and the sewer system from the model. We assume that some of the price already reflects the land value, so it doesn't explain any additional difference in home prices. We exclude the sewer system variable, because most home-buyers are unlikely to closely examine the type of sewer system in a potential home.

We observe that including interactions between variables makes the model more inaccurate (LM-11). We compare the model to a K-nearest neighbors model. We find that the best of the linear models (LM-6) is still better than all of the K-nearest neighbors models generated.

Conclusion

An improved model for predicting home prices is a linear model that includes Lot Size, age, living area, percent of college students in the area, number of bedrooms, number of fireplaces, number of bathrooms, number of rooms, heating type, fuel, central air system, Waterfront, and New Construction. This linear model excludes interactions and is slightly more accurate than the most accurate k-nearest neighbors models.

Problem 3 - Viral News

Regression

Testing different models out of a set of potentially meaningful predictors for number of shares (virality).

```
require(MASS)
```

```
## Loading required package: MASS
```

```
##
```

```
## Attaching package: 'MASS'
```

```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
##      select
```

```
require(leaps)
```

```
## Loading required package: leaps
```

```
## Warning: package 'leaps' was built under R version 3.6.3
```

```
#create set of linear models
```

```
online_news <- read_csv("SDS323/data/online_news.csv")
```

```
## Parsed with column specification:
```

```
## cols(
```

```
##   .default = col_double(),
```

```
##   url = col_character()
```

```
## )
```

```
## See spec(...) for full column specifications.
```

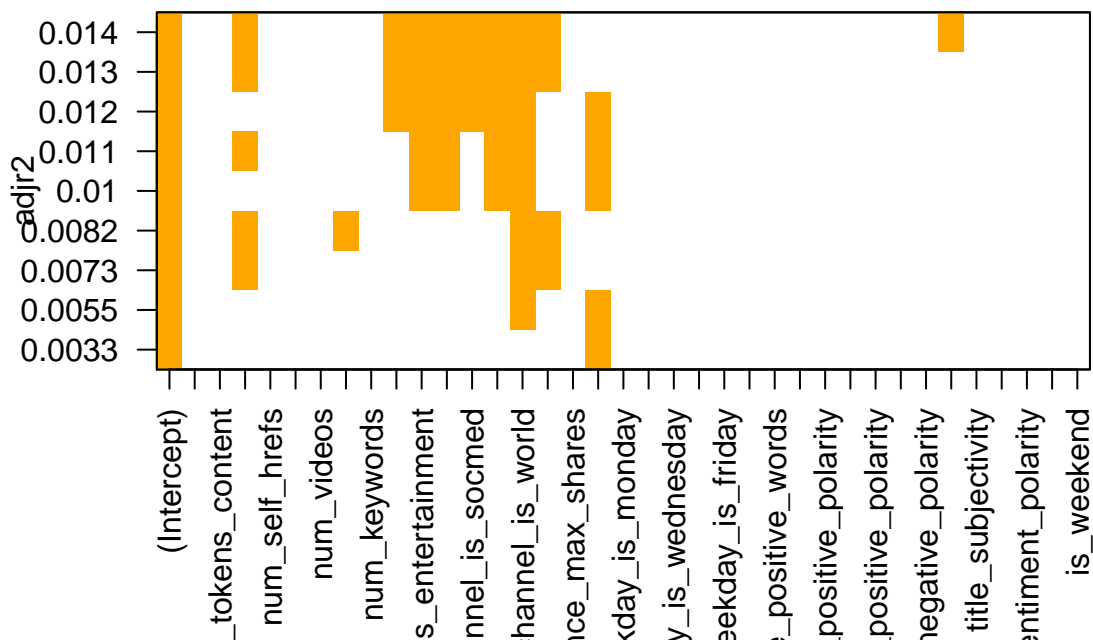
```
mysubsets <- regsubsets(shares ~ . - url, data=online_news)
```

```
## Warning in leaps.setup(x, y, wt = wt, nbest = nbest, nvmax = nvmax, force.in =
```

```
## force.in, : 2 linear dependencies found
```

```
## Reordering variables and trying again:
```

```
#rank by adjusted R-squared
plot(mysubsets, scale="adjr2", col="orange")
```



Resulting Model

```
model1 <- lm(shares ~ n_tokens_content + num_hrefs + num_self_hrefs + num_imgs + num_videos + average_token_length + self_reference_avg_share + avg_negative_polarity, data = online_news)
summary(model1)
```

```
##
## Call:
## lm(formula = shares ~ n_tokens_content + num_hrefs + num_self_hrefs +
##      num_imgs + num_videos + average_token_length + self_reference_avg_share +
##      avg_negative_polarity, data = online_news)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -25236  -2304  -1587   -398  838592
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    4.685e+03  3.217e+02  14.563  < 2e-16 ***
## n_tokens_content -5.724e-01  1.428e-01  -4.009  6.11e-05 ***
## num_hrefs       5.494e+01  6.194e+00   8.869  < 2e-16 ***
## num_self_hrefs  -7.070e+01  1.674e+01  -4.223  2.42e-05 ***
## num_imgs       4.700e+01  7.786e+00   6.036  1.59e-09 ***
```

```
## num_videos          5.179e+01  1.453e+01   3.564 0.000365 ***
## average_token_length -5.710e+02  7.473e+01  -7.641 2.21e-14 ***
## self_reference_avg_sharess 2.648e-02  2.409e-03  10.992 < 2e-16 ***
## avg_negative_polarity   -3.112e+03  4.877e+02  -6.382 1.77e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11570 on 39635 degrees of freedom
## Multiple R-squared:  0.009373,    Adjusted R-squared:  0.009173
## F-statistic: 46.88 on 8 and 39635 DF,  p-value: < 2.2e-16
```

```
adjr2_model <- lm(shares ~ num_hrefs + data_channel_is_lifestyle + data_channel_is_entertainment + data_channel_is_socmed + data_channel_is_tech + data_channel_is_world + self_reference_min_shares + max_negative_polarity, data = online_news)
summary(adjr2_model)
```

```
##
## Call:
## lm(formula = shares ~ num_hrefs + data_channel_is_lifestyle +
##     data_channel_is_entertainment + data_channel_is_bus + data_channel_is_socmed +
##     data_channel_is_tech + data_channel_is_world + self_reference_min_shares +
##     max_negative_polarity, data = online_news)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -30123  -2102  -1454   -397  837044
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    5.098e+03  1.836e+02  27.760 < 2e-16 ***
## num_hrefs       3.694e+01  5.184e+00   7.127 1.05e-12 ***
## data_channel_is_lifestyle -2.174e+03  2.924e+02  -7.434 1.08e-13 ***
## data_channel_is_entertainment -2.738e+03  2.027e+02 -13.507 < 2e-16 ***
## data_channel_is_bus      -2.621e+03  2.092e+02 -12.529 < 2e-16 ***
## data_channel_is_socmed    -2.250e+03  2.815e+02  -7.992 1.37e-15 ***
## data_channel_is_tech      -2.631e+03  2.016e+02 -13.047 < 2e-16 ***
## data_channel_is_world     -3.385e+03  1.960e+02 -17.268 < 2e-16 ***
## self_reference_min_shares   2.965e-02  2.947e-03  10.060 < 2e-16 ***
## max_negative_polarity     -1.229e+03  6.135e+02  -2.003  0.0451 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11550 on 39634 degrees of freedom
## Multiple R-squared:  0.01378,    Adjusted R-squared:  0.01356
## F-statistic: 61.55 on 9 and 39634 DF,  p-value: < 2.2e-16
```

Creating the testing and training sets

```
#random sample from data for train set
random_sample<-sample(seq_len(nrow(online_news)), size = 35000)
head(random_sample)
```

```
## [1] 561 9271 34285 14343 19517 6176
```



```

train<-online_news[random_sample,]
#use the rest for test
test<-online_news[-random_sample,]
#training model
training_adj2_model <- lm(shares ~ n_tokens_content + num_hrefs + num_self_hrefs + num_imgs + num_videos +
summary(training_adj2_model)

```

```

##
## Call:
## lm(formula = shares ~ n_tokens_content + num_hrefs + num_self_hrefs +
##      num_imgs + num_videos + average_token_length + self_reference_avg_sharess +
##      avg_negative_polarity, data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -25220  -2329  -1591   -400  838586
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    4.791e+03  3.446e+02  13.905 < 2e-16 ***
## n_tokens_content -7.736e-01  1.510e-01  -5.122 3.03e-07 ***
## num_hrefs       5.309e+01  6.526e+00   8.136 4.22e-16 ***
## num_self_hrefs  -6.964e+01  1.784e+01  -3.904 9.47e-05 ***
## num_imgs       5.286e+01  8.203e+00   6.444 1.18e-10 ***
## num_videos     5.702e+01  1.535e+01   3.714 0.000204 ***
## average_token_length -5.793e+02  7.991e+01  -7.250 4.26e-13 ***
## self_reference_avg_sharess 2.645e-02  2.495e-03  10.602 < 2e-16 ***
## avg_negative_polarity -3.269e+03  5.171e+02  -6.323 2.60e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 11560 on 34991 degrees of freedom
## Multiple R-squared:  0.01004,    Adjusted R-squared:  0.009815
## F-statistic: 44.36 on 8 and 34991 DF,  p-value: < 2.2e-16

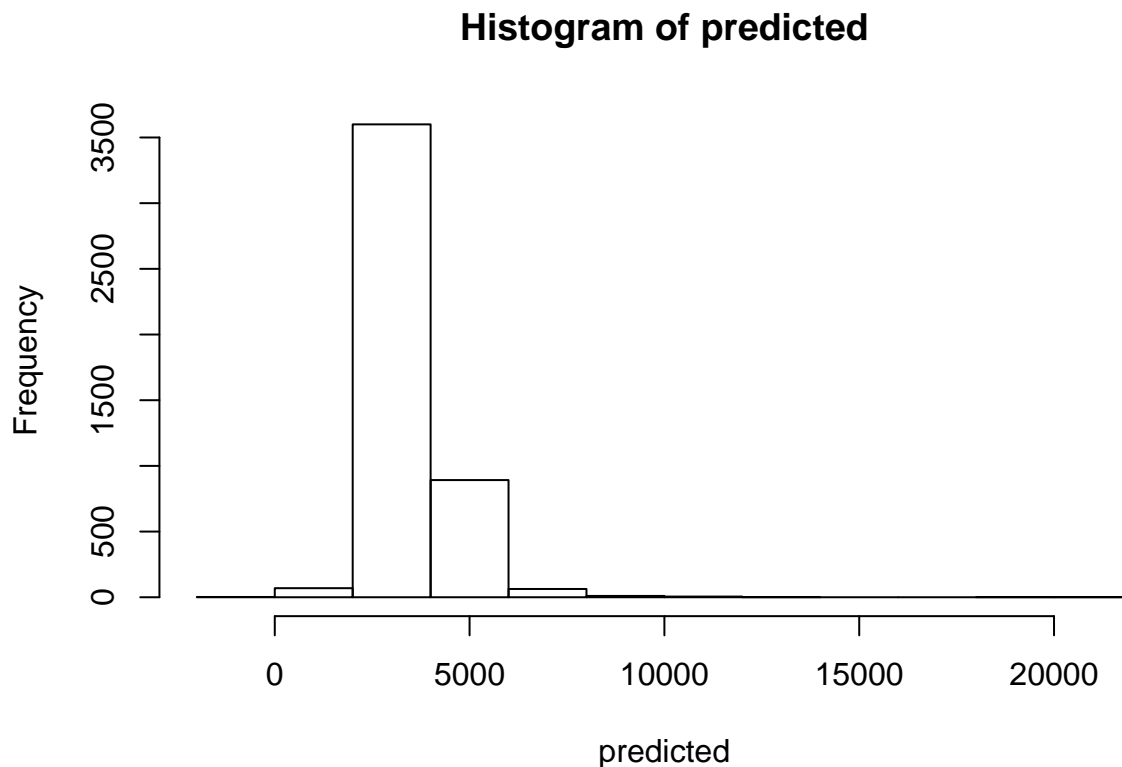
```

Creating set of model predictions on the testing set

```

#generate results
predicted<-predict.lm(training_adj2_model,newdata=test)
hist(predicted)

```



Covertng to binary (viral or not) form:

```
bin_predicted = ifelse(predicted >= 1400,1,0)
bin_actual = ifelse(test$shares >= 1400,1,0)
#creating confusion matrix
conf_viral = table(bin_actual, bin_predicted)
conf_viral
```

```
##          bin_predicted
## bin_actual    0     1
##          0     4 2175
##          1     2 2463
```

Finding average values:

```
correct_neg <- mean(c(2,1,2,2,2))
false_neg <- mean(c(3,1,0,1,0))
correct_pos <- mean(c(2478,2473,2450,2453,2482))
false_pos <- mean(c(2161,2169,2192,2188,2160))
```

Converting average values into confusion matrix:

```
library(knitr)
library(kableExtra)
```

Table 1: Confusion Matrix

Actually_Viral	Predicted_Not	Predicted_Viral
No	1.8	2174.0
Yes	1.0	2467.2

```
## Warning: package 'kableExtra' was built under R version 3.6.3
```

```
##
```

```
## Attaching package: 'kableExtra'
```

```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
##      group_rows
```

```
avg_conf_viral <- data.frame(
  Actually_Viral = c("No", "Yes"),
  Predicted_Not = c(correct_neg,false_neg),
  Predicted_Viral = c(false_pos,correct_pos)
)
kable(avg_conf_viral, caption = "Confusion Matrix") %>% row_spec(0,bold=FALSE) %>% kable_styling()
```

```
#total error rate
```

```
(false_neg+false_pos)/nrow(test) #46.83%
```

```
## [1] 0.4683463
```

```
#false positive rate
```

```
false_pos/(false_pos+correct_neg) #99.91%
```

```
## [1] 0.9991727
```

```
#false negative rate
```

```
false_neg/(false_neg+correct_pos) #0.04%
```

```
## [1] 0.0004051536
```

This linear model just barely outperforms the null model if the null model were to predict all articles to go viral. This means that the linear probability model we derived is not very useful, but that result is what we expected with an adjusted R-squared of only 0.0106. The model's total error rate is 46.51%, false positive rate is 99.91%, and false negative rate is 0.04%. This linear model predicts all but a few of the articles to viral. The reason that our error rate is below 50% (the expected error rate for random guesses) is there are simply more viral articles than not in the population.

Classification

Create a new binary variable (viral or not):

```
online_news$viral = ifelse(online_news$shares > 1400, 1, 0)
```

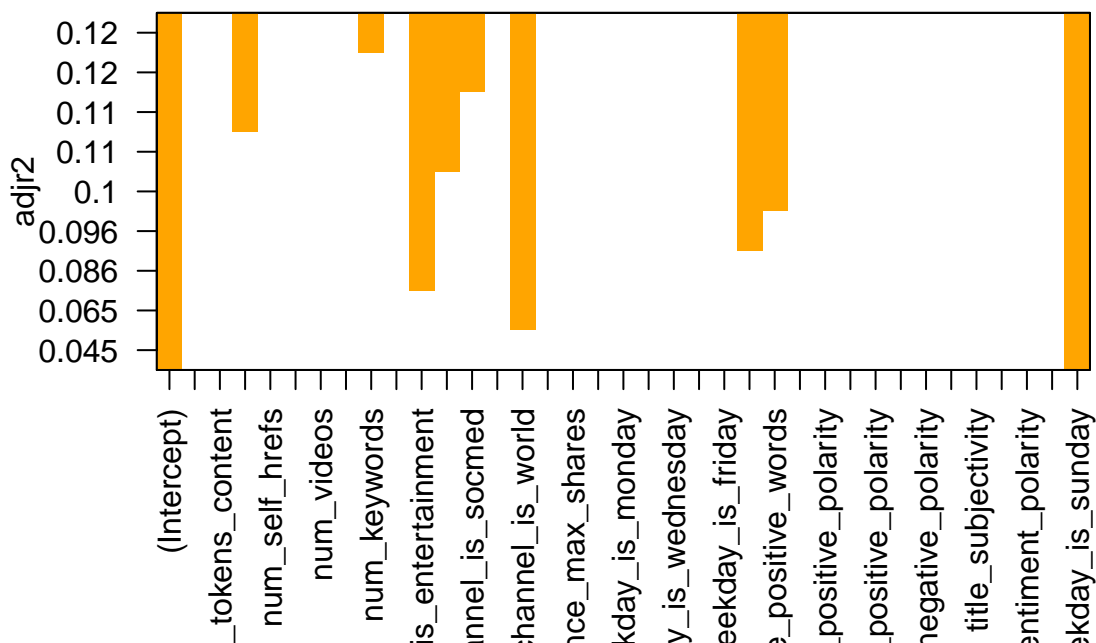
Test possible predictor combinations for predicting this new variable:

```
require(MASS)
require(leaps)
#create set of linear models - remove weekend b/c colinear
mybinsubsets <- regsubsets(viral ~ . - url - is_weekend, data=online_news)
```

```
## Warning in leaps.setup(x, y, wt = wt, nbest = nbest, nvmax = nvmax, force.in =
## force.in, : 1 linear dependencies found
```

```
## Reordering variables and trying again:
```

```
#rank by adjusted R-squared
plot(mybinsubsets, scale="adjr2", col="orange")
```



Resulting Model

```
binary_model <- lm(viral ~ num_hrefs + num_keywords + data_channel_is_entertainment + data_channel_is_b
```

```
## Warning: In lm.fit(x, y, offset = offset, singular.ok = singular.ok, ...) :
## extra argument 'family' will be disregarded
```

```
summary(binary_model)
```

```
##
## Call:
## lm(formula = viral ~ num_hrefs + num_keywords + data_channel_is_entertainment +
##     data_channel_is_bus + data_channel_is_socmed + data_channel_is_world +
##     weekday_is_saturday + weekday_is_sunday + global_rate_positive_words,
##     data = online_news, family = "binomial")
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0693 -0.4496 -0.2795  0.4535  0.7512
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    0.4528646   0.0120051  37.723 < 2e-16 ***
## num_hrefs       0.0026449   0.0002156  12.270 < 2e-16 ***
## num_keywords    0.0104067   0.0013172   7.901 2.84e-15 ***
## data_channel_is_entertainment -0.2074708   0.0069628 -29.797 < 2e-16 ***
## data_channel_is_bus    -0.0837258   0.0073790 -11.346 < 2e-16 ***
## data_channel_is_socmed   0.1264028   0.0108247  11.677 < 2e-16 ***
## data_channel_is_world   -0.2352816   0.0066583 -35.336 < 2e-16 ***
## weekday_is_saturday     0.2105378   0.0100605  20.927 < 2e-16 ***
## weekday_is_sunday       0.1658719   0.0095576  17.355 < 2e-16 ***
## global_rate_positive_words  0.1234263   0.1437590   0.859  0.391
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4797 on 39634 degrees of freedom
## Multiple R-squared:  0.07963, Adjusted R-squared:  0.07942
## F-statistic: 381 on 9 and 39634 DF, p-value: < 2.2e-16
```

Creating the testing and training sets

```
#random sample from data for train set
random_bin_sample<-sample(seq_len(nrow(online_news)), size = 35000)
head(random_bin_sample)
```

```
## [1] 18554 37294 1669 6594 2725 28383
```

```
bin_train<-online_news[random_bin_sample,]
#use the rest for test
bin_test<-online_news[-random_bin_sample,]
#training model
training_binary_model <- lm(viral ~ num_hrefs + num_keywords + data_channel_is_entertainment + data_chan
```

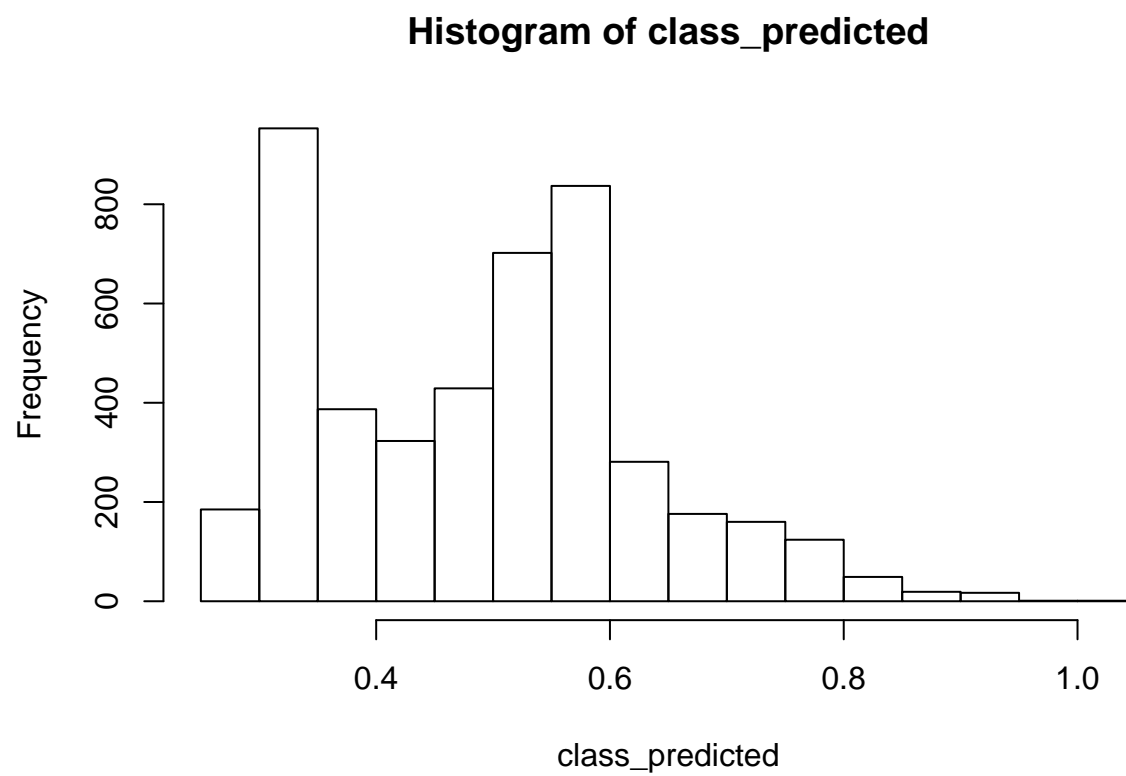
```
## Warning: In lm.fit(x, y, offset = offset, singular.ok = singular.ok, ...) :
## extra argument 'family' will be disregarded
```

```
summary(training_binary_model)
```

```
##
## Call:
## lm(formula = viral ~ num_hrefs + num_keywords + data_channel_is_entertainment +
##     data_channel_is_bus + data_channel_is_socmed + data_channel_is_world +
##     weekday_is_saturday + weekday_is_sunday + global_rate_positive_words,
##     data = bin_train, family = "binomial")
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0572 -0.4479 -0.2816  0.4542  0.7417
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    0.4572753   0.0127921   35.747 < 2e-16 ***
## num_hrefs       0.0026699   0.0002289   11.665 < 2e-16 ***
## num_keywords    0.0101978   0.0014043    7.262 3.91e-13 ***
## data_channel_is_entertainment -0.2046706   0.0074274 -27.556 < 2e-16 ***
## data_channel_is_bus    -0.0858602   0.0078699 -10.910 < 2e-16 ***
## data_channel_is_socmed   0.1221070   0.0115141  10.605 < 2e-16 ***
## data_channel_is_world   -0.2353292   0.0070884 -33.199 < 2e-16 ***
## weekday_is_saturday    0.2051974   0.0107480  19.092 < 2e-16 ***
## weekday_is_sunday     0.1661671   0.0101645  16.348 < 2e-16 ***
## global_rate_positive_words  0.0336517   0.1532106   0.220  0.826
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4801 on 34990 degrees of freedom
## Multiple R-squared:  0.07803,    Adjusted R-squared:  0.07779
## F-statistic: 329 on 9 and 34990 DF,  p-value: < 2.2e-16
```

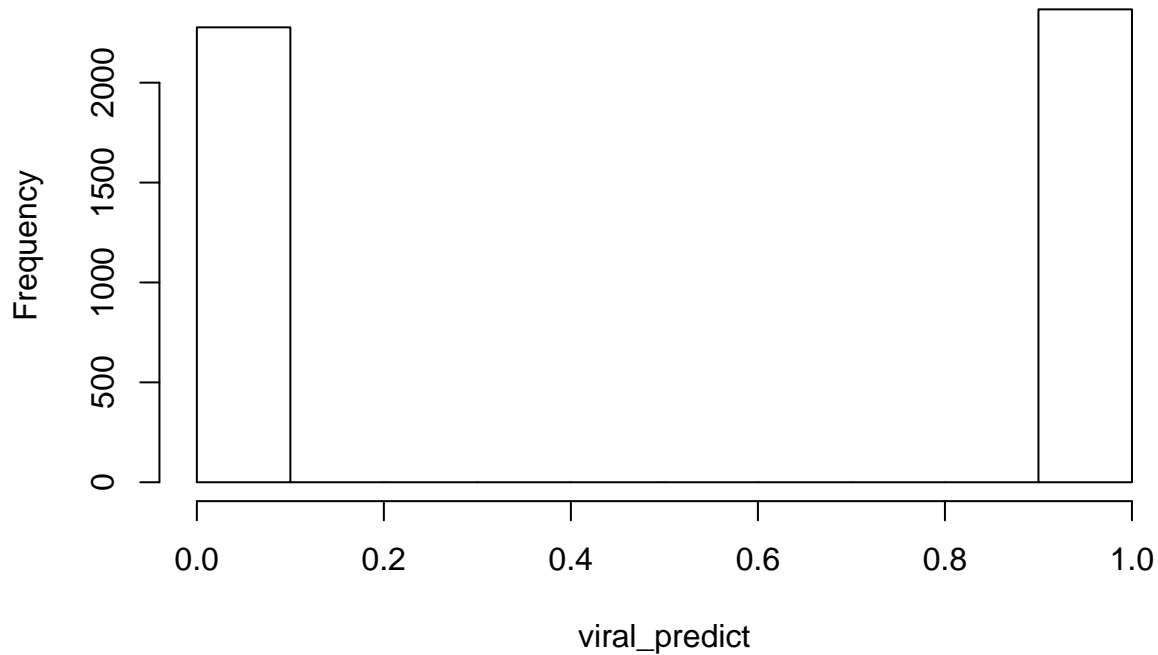
Creating set of probability predictions on the testing set

```
#generate results
class_predicted<-predict.glm(training_binary_model,newdata=bin_test,type = "response")
hist(class_predicted)
```



```
viral_predict = ifelse(class_predicted >= .5,1,0)  
hist(viral_predict)
```

Histogram of viral_predict



Comparing predictions to actual virality:

```
viral_actual = ifelse(bin_test$viral==1,1,0)
#creating confusion matrix
class_conf_viral = table(viral_actual, viral_predict)
class_conf_viral
```

```
##          viral_predict
## viral_actual    0    1
##           0 1448  885
##           1  829 1482
```

Finding average values:

```
class_correct_neg <- mean(c(1389,1365,1348,1419,1450))
class_false_neg <- mean(c(818,834,771,821,784))
class_correct_pos <- mean(c(1487,1486,1550,1479,1439))
class_false_pos <- mean(c(950,959,975,925,971))
```

Converting average values into confusion matrix:

```
library(knitr)
library(kableExtra)
bin_avg_conf_viral <- data.frame(
  Actually_Viral = c("No", "Yes"),
```


Table 2: Confusion Matrix

Actually_Viral	Predicted_Not	Predicted_Viral
No	1394.2	956.0
Yes	805.6	1488.2

```

Predicted_Not = c(class_correct_neg,class_false_neg),
Predicted_Viral = c(class_false_pos,class_correct_pos)
)
kable(bin_avg_conf_viral, caption = "Confusion Matrix") %>% row_spec(0,bold=FALSE) %>% kable_styling()

```

```

#total error rate
(class_false_neg+class_false_pos)/nrow(bin_test) #37.93%

```

```
## [1] 0.3793282
```

```

#false positive rate
class_false_pos/(class_false_pos+class_correct_neg) #40.68%

```

```
## [1] 0.4067739
```

```

#false negative rate
class_false_neg/(class_false_neg+class_correct_pos) #35.12%

```

```
## [1] 0.3512076
```

The classification model is more successful than the linear model and far outperformed the null model as both the false negative and false positive error rates are below 50%. This model may be useful, but it is not very accurate or predictive with an adjusted R-squared of 0.0810. The model's total error rate is 37.93%, false positive rate is 40.68%, and false negative rate is 35.12%. The classification model is more evenly split in its predictions that the regression model, predicting more articles to go viral than not which mirrors the data set.

#Summary ### The Classification approach (threshold first and regress second) performs better than the Regression approach (regress first and threshold second). This approach fits the problem better because Mashable is look for a yes or no response on if each article they publish goes viral, which ultimately comes down to a set of probabilities. By framing this question as a logistic regression, each predictor is given the chance to affect the probability that an article with certain features will reach 1400 shares (going viral) rather than just attempting to predict how many shares an article will have. The classification approach proved to deliver a more useful and accurate model, with a total error rate of 37.93% compared to 46.51% and adjusted R-squared of 0.0810 compared to 0.0106. The Regression model predicted that nearly all articles would go viral, making it nearly useless.