

# SVM Regression

Haniyyah Hamid, Jered Hightower

10/23/2022

Data set used: <https://www.kaggle.com/datasets/dhirajnrne/california-housing-data>

## Loading packages

```
library(e1071)
library(MASS)
```

## Importing data

```
df <- read.csv("housing.csv")
str(df)
```

```
## 'data.frame': 20640 obs. of 10 variables:
## $ longitude : num -122 -122 -122 -122 -122 ...
## $ latitude : num 37.9 37.9 37.9 37.9 37.9 ...
## $ housing_median_age: num 41 21 52 52 52 52 52 52 42 52 ...
## $ total_rooms : num 880 7099 1467 1274 1627 ...
## $ total_bedrooms : num 129 1106 190 235 280 ...
## $ population : num 322 2401 496 558 565 ...
## $ households : num 126 1138 177 219 259 ...
## $ median_income : num 8.33 8.3 7.26 5.64 3.85 ...
## $ median_house_value: num 452600 358500 352100 341300 342200 ...
## $ ocean_proximity : chr "NEAR BAY" "NEAR BAY" "NEAR BAY" "NEAR BAY" ...
```

## Clean up data

We will remove unnecessary columns. We will reduce the number of rows as well to be 10k to make tuning faster.

```
df <- df[,c(3, 7, 8, 9)]
df <- head(df, - 10000)
df$income <- factor(df$income)
str(df)
```

```
## 'data.frame': 10640 obs. of 4 variables:
## $ housing_median_age: num 41 21 52 52 52 52 52 52 42 52 ...
## $ households : num 126 1138 177 219 259 ...
## $ median_income : num 8.33 8.3 7.26 5.64 3.85 ...
## $ median_house_value: num 452600 358500 352100 341300 342200 ...
```

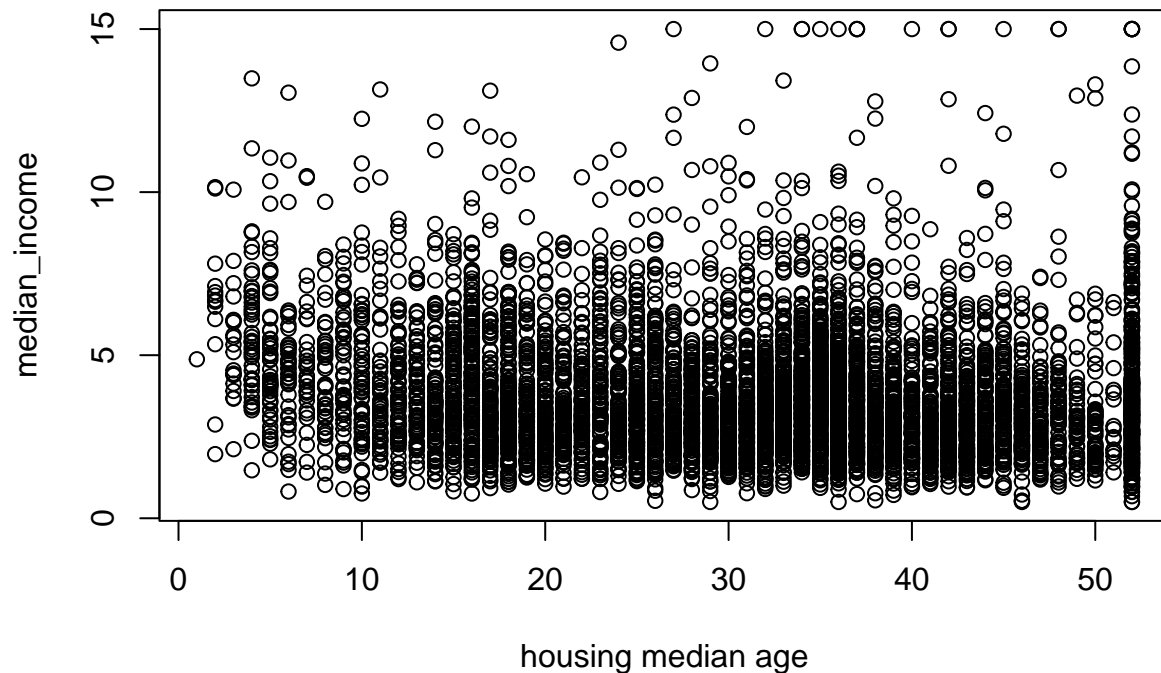
```
#head(df)
```

Divide into train, test, validate

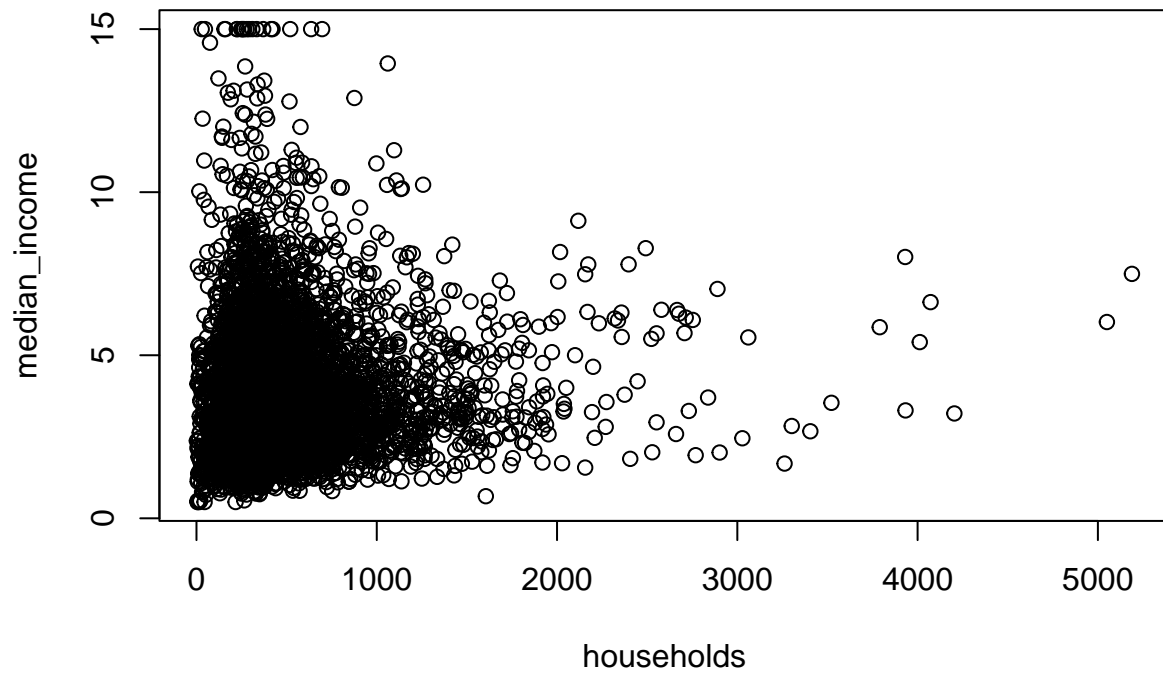
```
set.seed(1234)
spec <-c(train=.6, test=.2, validate=.2)
i <- sample(cut(1:nrow(df),nrow(df)*cumsum(c(0,spec))), labels=names(spec))
train <- df[i=="train",]
test <- df[i=="test",]
vald <- df[i=="validate",]
```

Plotting and statistically exploring the training data

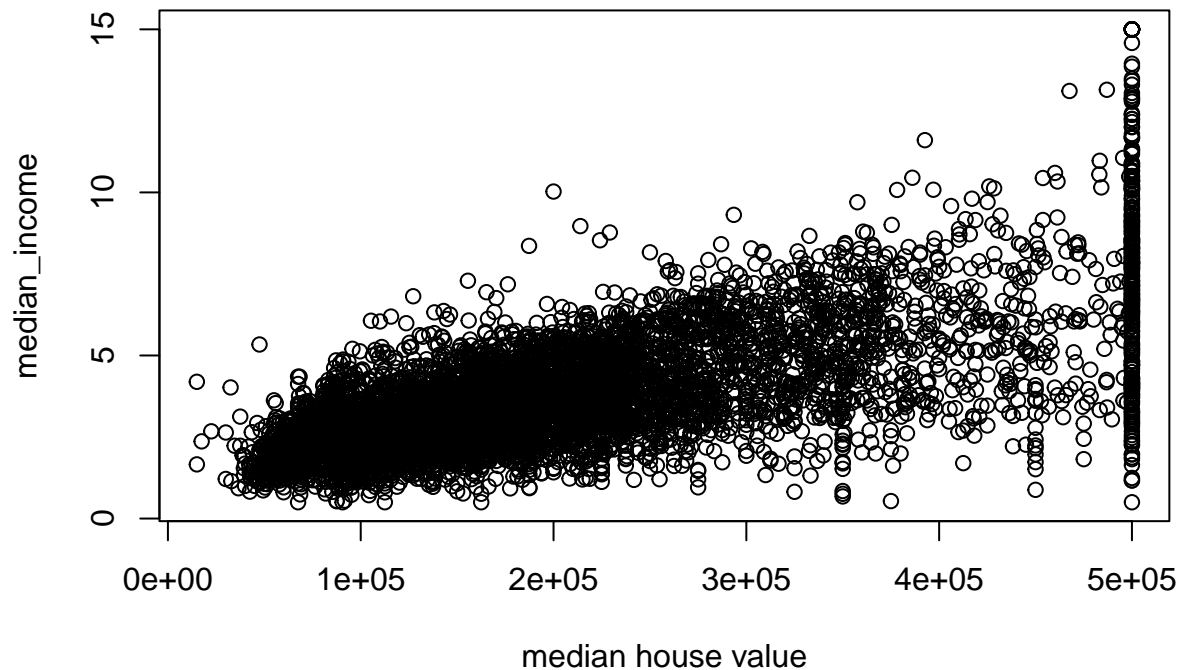
```
plot(train$median_income~train$housing_median_age, xlab="housing median age", ylab="median_income")
```



```
plot(train$median_income~train$households, xlab="households", ylab="median_income")
```



```
plot(train$median_income~train$median_house_value, xlab="median house value", ylab="median_income")
```



```
summary(df)
```

```
## housing_median_age  households  median_income  median_house_value
## Min.   : 1.00      Min.    : 2.0    Min.    : 0.4999  Min.    : 14999
## 1st Qu.:22.00     1st Qu.: 272.0  1st Qu.: 2.4666  1st Qu.:121600
## Median :33.00     Median : 390.0  Median : 3.4111  Median :181050
## Mean   :31.03     Mean     :482.8  Mean    : 3.8111  Mean    :207878
## 3rd Qu.:40.00     3rd Qu.: 574.0  3rd Qu.: 4.6736  3rd Qu.:264750
## Max.   :52.00     Max.     :6082.0  Max.    :15.0001  Max.    :500001
```

Of the three graphs, the relationship between median income and the median house value seemed to be the most linear. The graph of median income vs. households seems to all be clustered in the bottom left corner of the graph. While the graph of the median income vs. housing median age yields a graph that is distributed across the x axis, making it difficult to determine a particular pattern.

## Try linear regression

```
lm1 <- lm(median_income~., data=train)
pred <- predict(lm1, newdata=test)
cor_lm1 <- cor(pred, test$median_income)
mse_lm1 <- mean((pred-test$median_income)^2)
```

## Try a linear kernel

```
svm1 <- svm(median_income~., data=train, kernel="linear", cost=10, scale=TRUE)
summary(svm1)
```

```
##
## Call:
## svm(formula = median_income ~ ., data = train, kernel = "linear",
##      cost = 10, scale = TRUE)
##
##
## Parameters:
##   SVM-Type:  eps-regression
##   SVM-Kernel: linear
##      cost:   10
##    gamma: 0.3333333
##   epsilon: 0.1
##
##
## Number of Support Vectors: 5472
```

```
pred <- predict(svm1, newdata=test)
cor_svm1 <- cor(pred, test$median_income)
mse_svm1 <- mean((pred - test$median_income)^2)
```

## Tune

```
tune_svm1 <- tune(svm, median_income~., data=vald, kernel="linear", ranges=list(cost=c(0.001, 0.01, 0.1),
summary(tune_svm1)
```

```
##
## Parameter tuning of 'svm':
##
## - sampling method: 10-fold cross validation
##
## - best parameters:
##   cost
##   100
##
## - best performance: 1.934216
##
## - Detailed performance results:
##   cost      error dispersion
## 1 1e-03 2.293132 0.4154893
## 2 1e-02 1.940273 0.3596556
## 3 1e-01 1.934894 0.3506409
## 4 1e+00 1.934530 0.3489227
## 5 5e+00 1.934409 0.3488024
## 6 1e+01 1.934442 0.3487704
## 7 1e+02 1.934216 0.3487463
```

## Evaluate on best linear SVM

```
pred <- predict(tune_svm1$best.model, newdata=test)
cor_svm1_tune <- cor(pred, test$median_income)
mse_svm1_tune <- mean((pred - test$median_income)^2)
```

## Try a polynomial kernel

```
svm2 <- svm(median_income~., data=train, kernel="polynomial", cost=10, scale=TRUE)
summary(svm2)
```

```
##
## Call:
## svm(formula = median_income ~ ., data = train, kernel = "polynomial",
##      cost = 10, scale = TRUE)
##
##
## Parameters:
##   SVM-Type:  eps-regression
##   SVM-Kernel: polynomial
##      cost:   10
##   degree:    3
##   gamma:     0.3333333
##   coef.0:    0
##   epsilon:   0.1
##
##
## Number of Support Vectors:  5612
```

```
pred <- predict(svm2, newdata=test)
cor_svm2 <- cor(pred, test$median_income)
mse_svm2 <- mean((pred - test$median_income)^2)
```

## Try a radial kernel

```
svm3 <- svm(median_income~., data=train, kernel="radial", cost=10, gamma=1, scale=TRUE)
summary(svm3)
```

```
##
## Call:
## svm(formula = median_income ~ ., data = train, kernel = "radial",
##      cost = 10, gamma = 1, scale = TRUE)
##
##
## Parameters:
##   SVM-Type:  eps-regression
##   SVM-Kernel: radial
```

```
##      cost: 10
##      gamma: 1
##      epsilon: 0.1
##
##
## Number of Support Vectors: 5378
```

```
pred <- predict(svm3, newdata=test)
cor_svm3 <- cor(pred, test$median_income)
mse_svm3 <- mean((pred - test$median_income)^2)
```

## Tune hyperparameters

```
set.seed(1234)
tune.out <- tune(svm, median_income~., data=vald, kernel="radial", ranges=list(cost=c(0.1, 1, 10, 100, 1000), gamma=c(0.01, 0.1, 1, 10, 100, 1000)))
summary(tune.out)
```

```
##
## Parameter tuning of 'svm':
##
## - sampling method: 10-fold cross validation
##
## - best parameters:
##   cost gamma
##   1      0.5
##
## - best performance: 1.886606
##
## - Detailed performance results:
```

	cost	gamma	error	dispersion
## 1	1e-01	0.5	2.081685	0.5409546
## 2	1e+00	0.5	1.886606	0.4150534
## 3	1e+01	0.5	1.888209	0.3743815
## 4	1e+02	0.5	2.006507	0.3468630
## 5	1e+03	0.5	2.642792	0.9182706
## 6	1e-01	1.0	2.145599	0.5637560
## 7	1e+00	1.0	1.900676	0.4150506
## 8	1e+01	1.0	1.991707	0.3633135
## 9	1e+02	1.0	2.363728	0.5091026
## 10	1e+03	1.0	4.306932	1.4651069
## 11	1e-01	2.0	2.292982	0.6115231
## 12	1e+00	2.0	1.964140	0.4065697
## 13	1e+01	2.0	2.136951	0.3939967
## 14	1e+02	2.0	2.883890	0.6768614
## 15	1e+03	2.0	4.711954	0.8884599
## 16	1e-01	3.0	2.404436	0.6513568
## 17	1e+00	3.0	2.019773	0.4075842
## 18	1e+01	3.0	2.292831	0.4146292
## 19	1e+02	3.0	3.025913	0.5804187
## 20	1e+03	3.0	6.440873	2.9989875

```
## 21 1e-01    4.0 2.510755  0.6839634
## 22 1e+00    4.0 2.062424  0.3988603
## 23 1e+01    4.0 2.376882  0.4352604
## 24 1e+02    4.0 3.332062  0.5748096
## 25 1e+03    4.0 7.273798  2.2688460
```

```
svm4 <- svm(median_income~., data=train, kernel="radial", cost=1, gamma=0.5, scale=TRUE)
summary(svm4)
```

```
##
## Call:
## svm(formula = median_income ~ ., data = train, kernel = "radial",
##      cost = 1, gamma = 0.5, scale = TRUE)
##
##
## Parameters:
##   SVM-Type:  eps-regression
##   SVM-Kernel: radial
##      cost:   1
##      gamma:  0.5
##   epsilon:  0.1
##
##
## Number of Support Vectors:  5403
```

```
pred <- predict(svm4, newdata=test)
cor_svm4 <- cor(pred, test$median_income)
mse_svm4 <- mean((pred - test$median_income)^2)
```

## Comparing statistics of each of the SVM kernels

First, the correlations of each kernel.

```
print(paste("cor_lm1 = ", cor_lm1))
```

```
## [1] "cor_lm1 =  0.740081889222345"
```

```
print(paste("cor_svm1_tune = ", cor_svm1_tune))
```

```
## [1] "cor_svm1_tune =  0.740068040245718"
```

```
print(paste("cor_svm2 = ", cor_svm2))
```

```
## [1] "cor_svm2 =  0.663007515753886"
```

```
print(paste("cor_svm3 = ", cor_svm3))
```

```
## [1] "cor_svm3 =  0.737902959957496"
```



```
print(paste("cor_svm4 = ", cor_svm4))
```

```
## [1] "cor_svm4 = 0.747757223387671"
```

We see the greatest correlation was found with the radial SVM kernel. Meaning the radial decision boundary probably yielded the best correlation of the 3 kernels. The worst was the polynomial SVM kernel.

Now, the mean standard errors of each kernel.

```
print(paste("mse_lm1 = ", mse_lm1))
```

```
## [1] "mse_lm1 = 1.77490141707553"
```

```
print(paste("mse_svm1_tune = ", mse_svm1_tune))
```

```
## [1] "mse_svm1_tune = 1.77162998450893"
```

```
print(paste("mse_svm2 = ", mse_svm2))
```

```
## [1] "mse_svm2 = 2.23841248630074"
```

```
print(paste("mse_svm3 = ", mse_svm3))
```

```
## [1] "mse_svm3 = 1.78557131943635"
```

```
print(paste("mse_svm4 = ", mse_svm4))
```

```
## [1] "mse_svm4 = 1.72927586428391"
```

We see the lowest MSE was found when performing the radial SVM kernel. The highest MSE of the kernels was the polynomial SVM kernel.

Therefore, we can assume that because of all the given results, the radial SVM kernel performed the best of the 3 kernels.