

Regularization

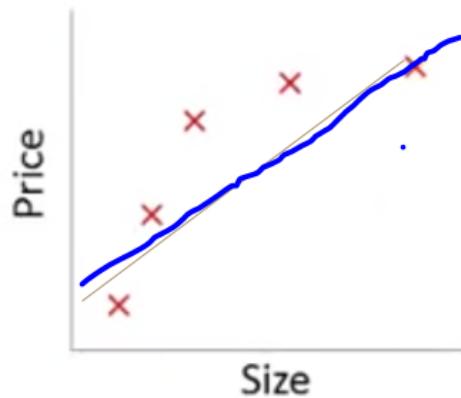
~Abhishek Kumar

Scope

- Problem of Overfitting
- What is regularization?
- Why do we need regularization?
- Types of regularization
- When to use them?
- Model Evaluation

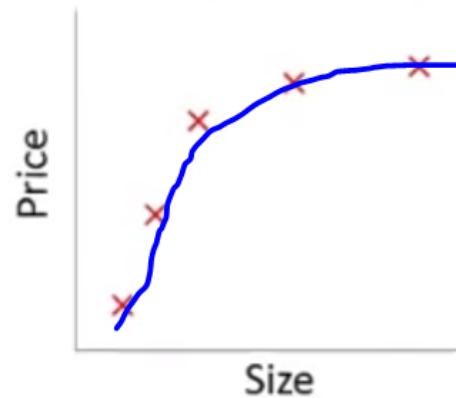
Overfitting

Example: Linear regression (housing prices)



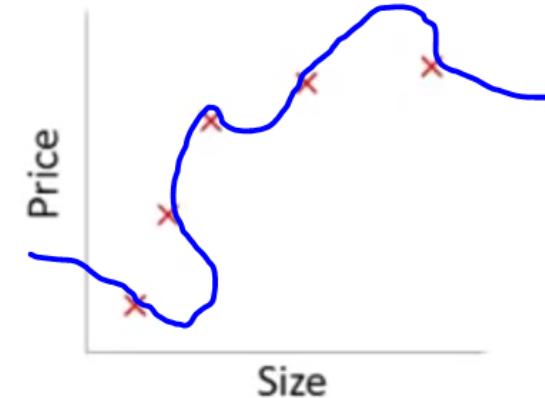
$$mx + c$$

“Underfitting”
High Bias



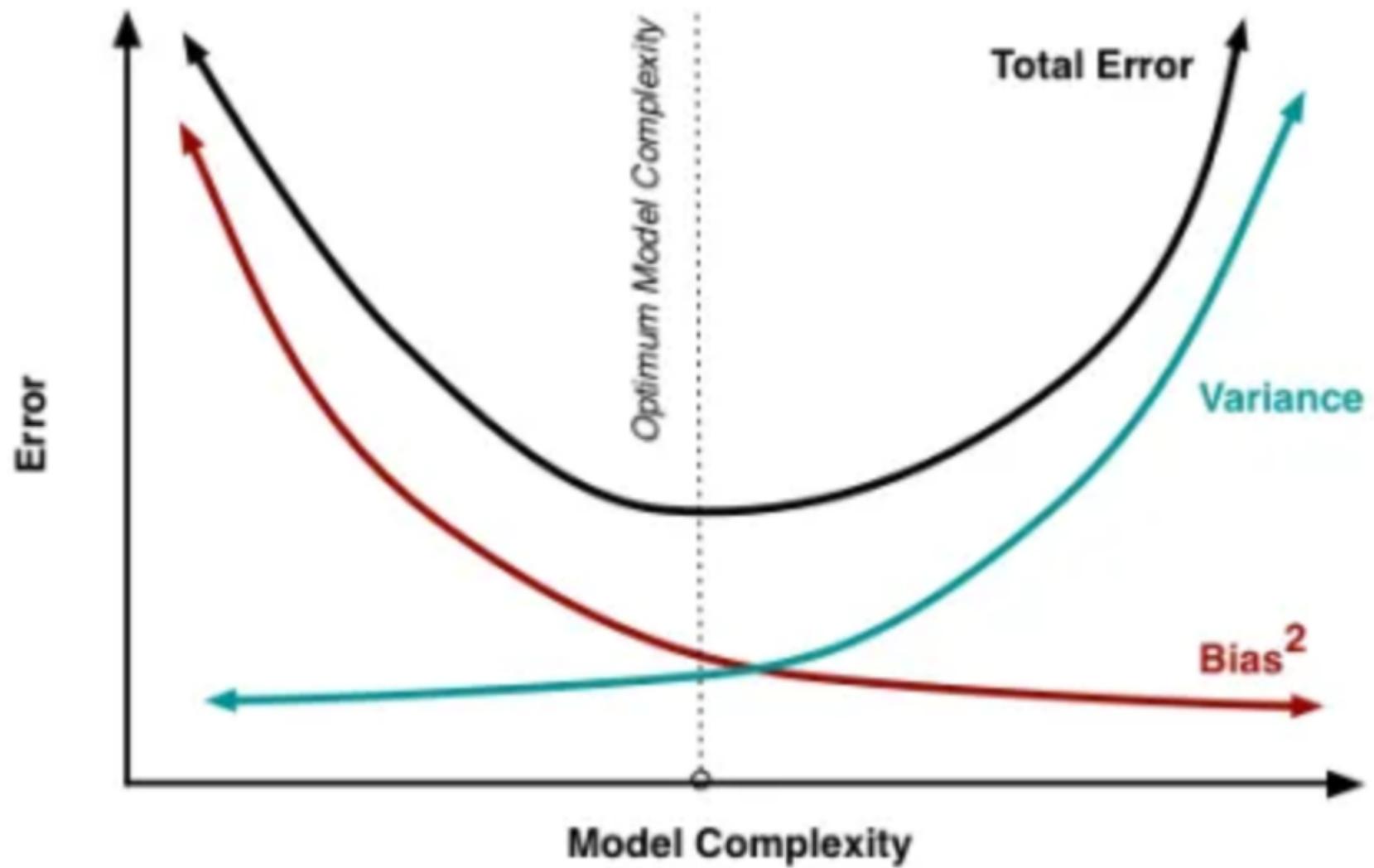
$$m_1x + m_2x^2 + c$$

“Just right”



$$m_1x + m_2x^2 + m_3x^3 + m_4x^4 + c$$

“Overfitting”
High Variance



Bias & Variance

- High Bias: It means the algorithm has a strong preconception / bias that housing prices are going to vary linearly with their size despite the evidence based on previous data.
- High Variance: If we fit a high order polynomial then hypothesis can fit to any function. So in this case our possible hypothesis is too large or too variable.

Regularization

- If we have too many features, the learned hypothesis may fit the training set very well(cost function almost zero) , but fail to generalize new examples.
- When we have a high dimensional data set, it would be highly inefficient to use all the variables since some of them might be imparting redundant information.

Addressing overfitting:

x_1 = size of house

x_2 = no. of bedrooms

x_3 = no. of floors

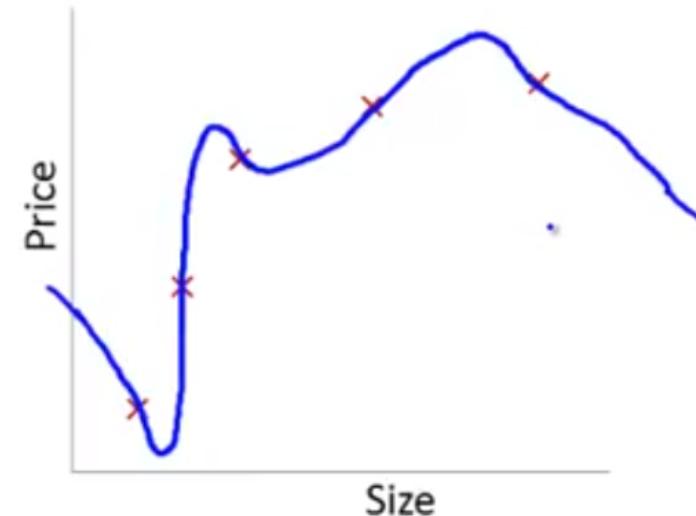
x_4 = age of house

x_5 = average income in neighborhood

x_6 = kitchen size

:

x_{100}



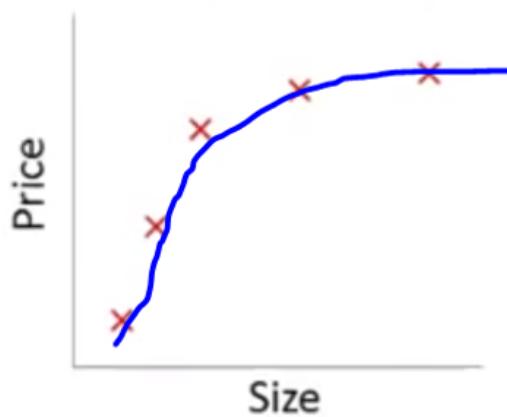
Addressing Overfitting

- Reduce number of features
 - Manually select which features to keep
 1. Business Understanding & Domain knowledge
 2. Forward Selection
 3. Backward Elimination
 - Regularization

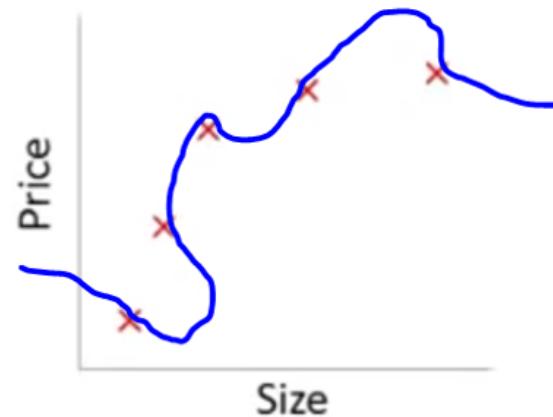
Regularization

- Keep all /some of the features, but reduce magnitude/value of parameter m_j .
- Works well when we have lot of features, each of which contributes a bit to predicting y .

Cost Function



$$m_1x + m_2x^2 + c$$



$$m_1x + m_2x^2 + m_3x^3 + m_4x^4 + c$$

We penalize and make m_3 and m_4 really small

“Everything should be made simple as possible, but not simpler – Albert Einstein”

- Small values of parameters m_1, m_2, \dots, m_k
- “Simple” hypothesis
- Less prone to overfitting

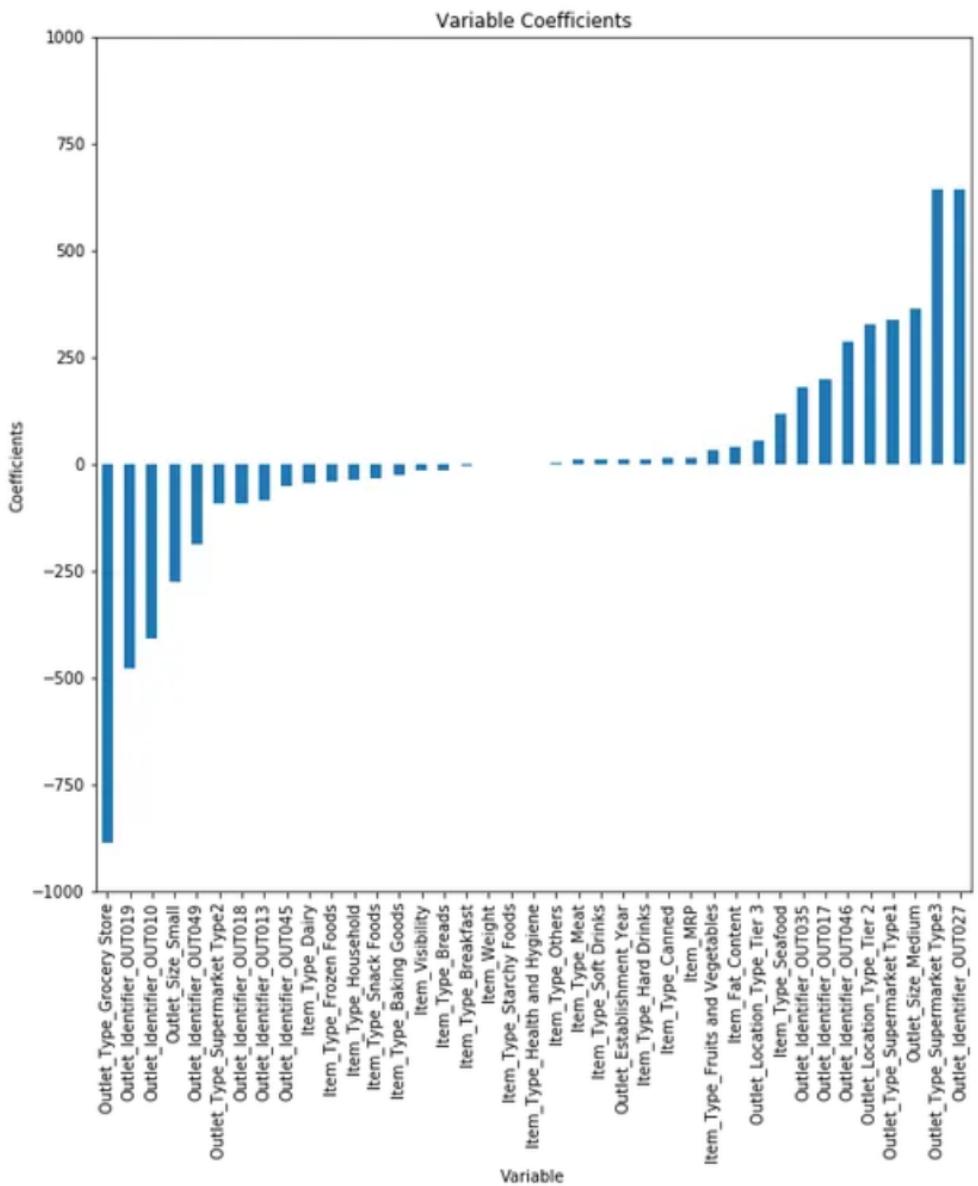
$$J(m) = \frac{1}{2n} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^k m_j^p$$

Regularization parameter

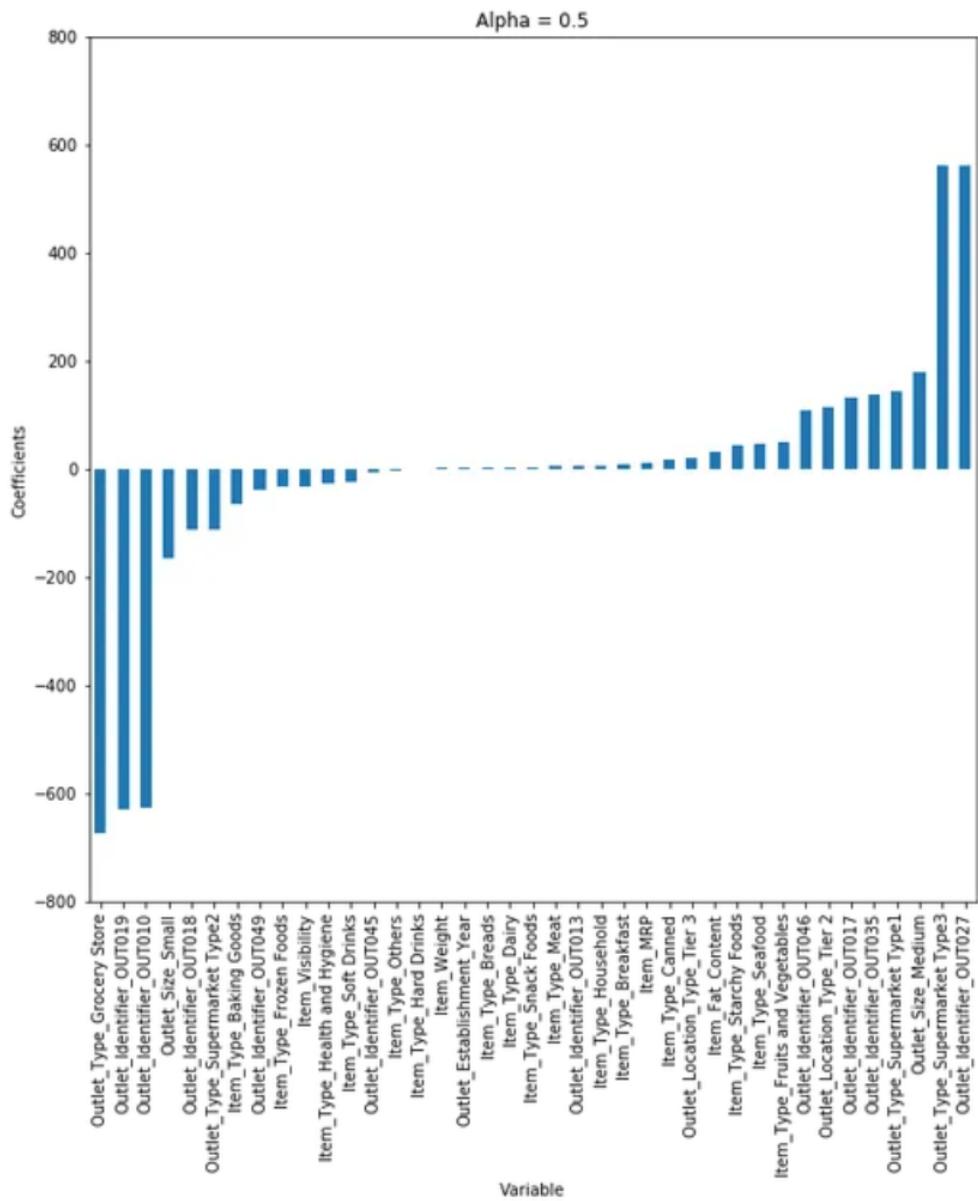
- $\min J(m)$ Regularization term
- Lambda controls the trade off between fitting the training set well and setting the value of parameters m_j small and therefore keeping the hypothesis simple to avoid overfitting.

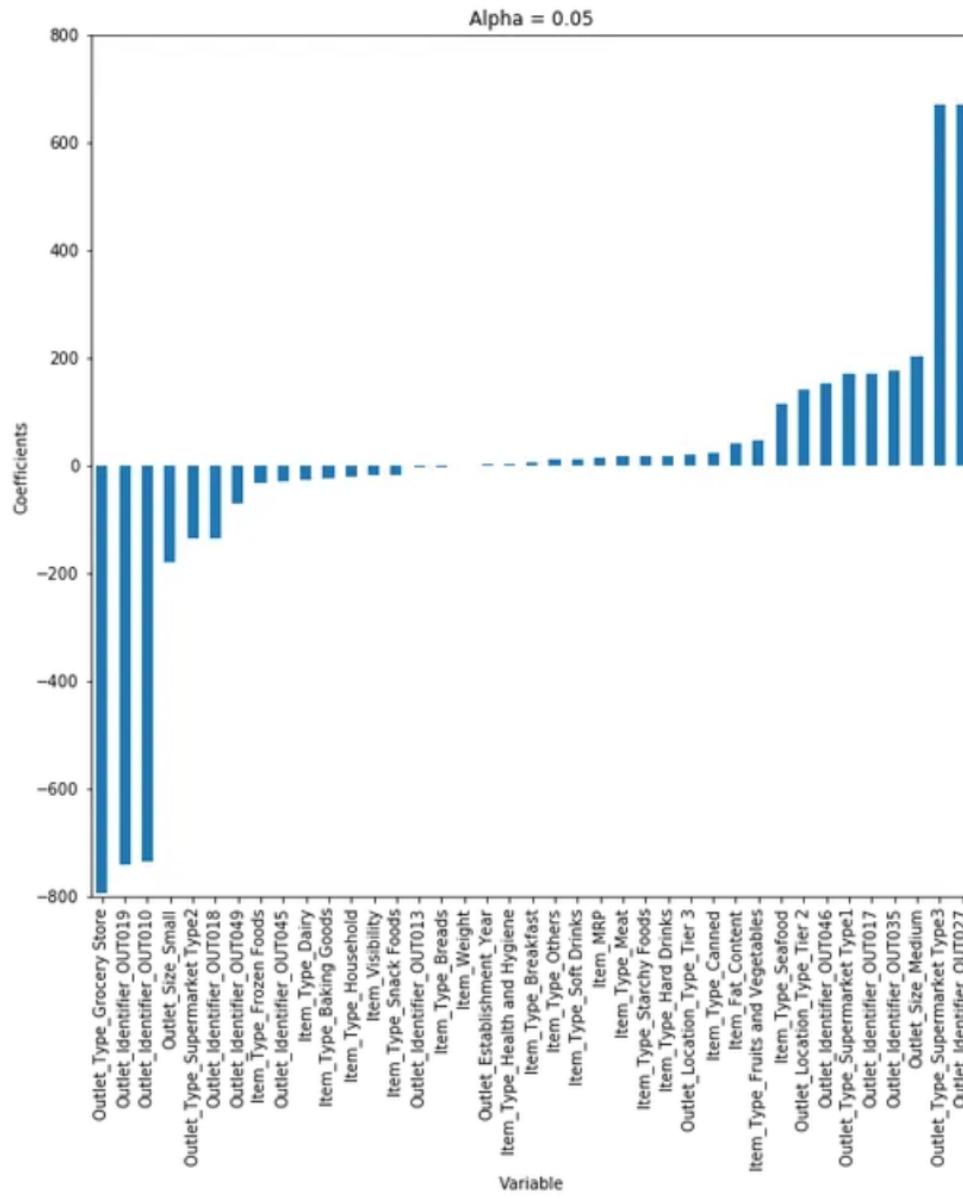
What happens if we set lamda very large?

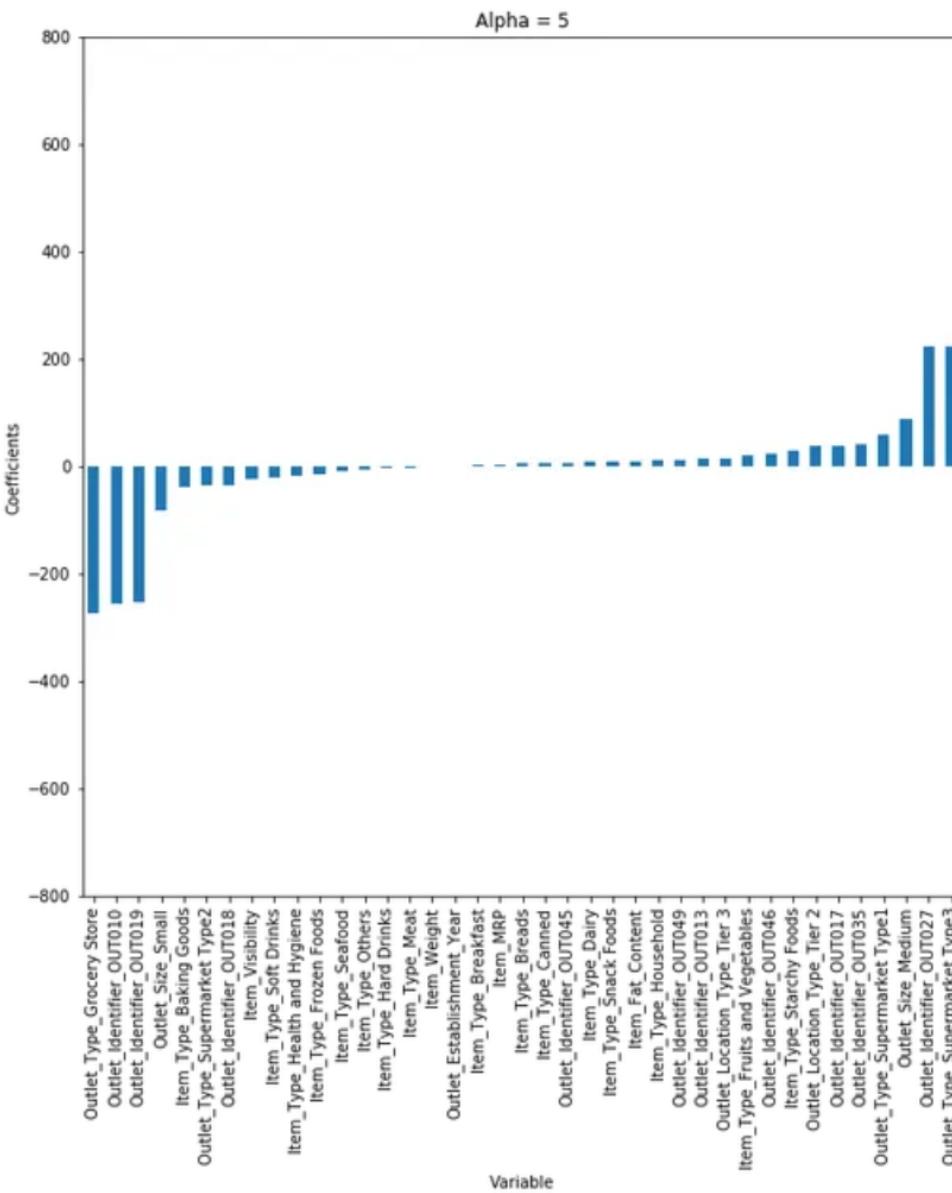
- It means we end up penalizing parameters m_j very heavily.
- Horizontal flat line a result, “Underfitting”
- So lamda should be choose very carefully.

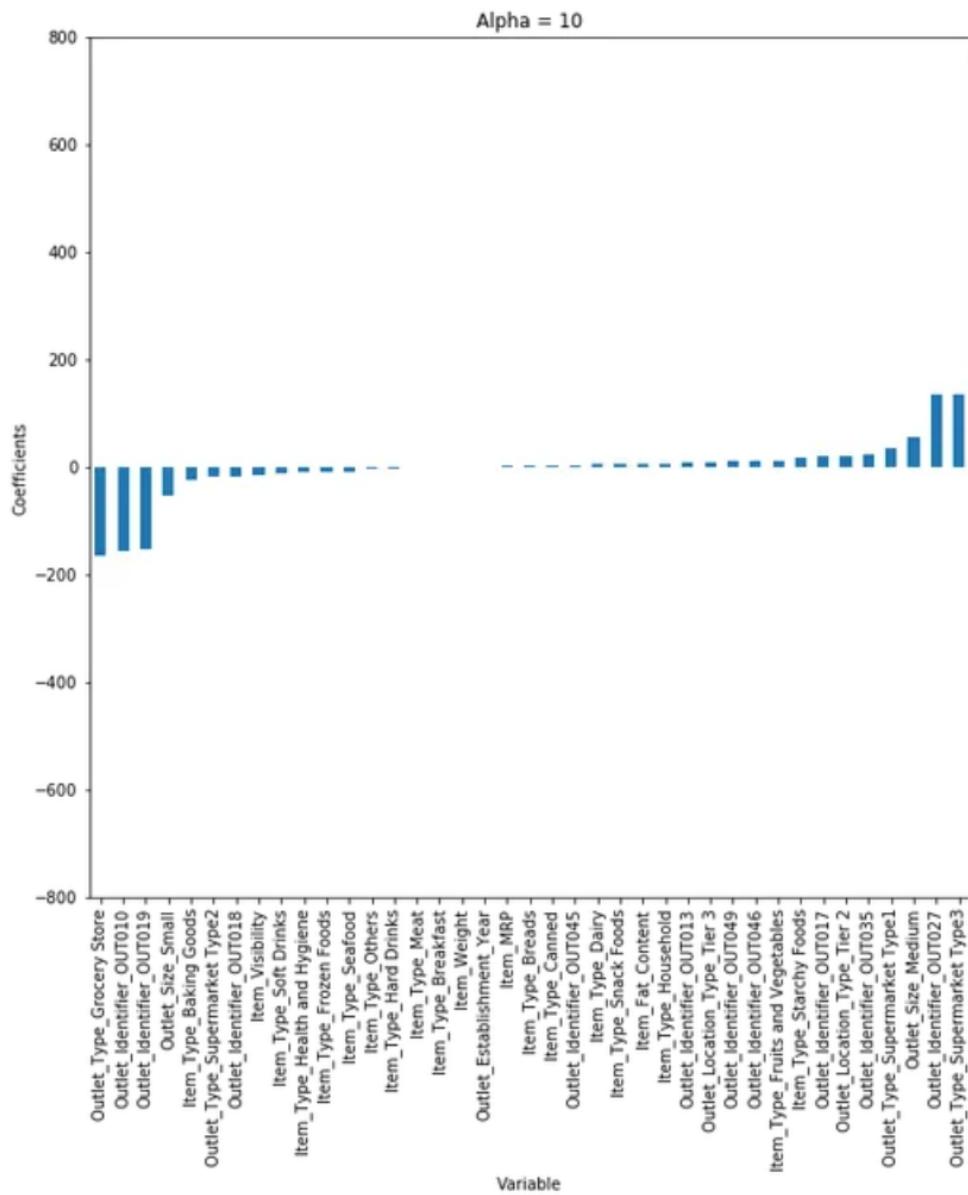


Ridge Regression









Ridge Regression

- You can see that, as we increase the value of alpha, the magnitude of the coefficients decreases, where the values reaches to zero but not absolute zero.
- But if you calculate R-square for each alpha, we will see that the value of R-square will be maximum at alpha=0.05.
- So we have to choose it wisely by iterating it through a range of values and using the one which gives us lowest error.

Ridge Regression

```
from sklearn.linear_model import  
Ridge  
## training the model  
ridgeReg = Ridge(alpha=0.05,  
normalize=True)  
ridgeReg.fit(x_train,y_train)  
pred = ridgeReg.predict(x_cv)
```

- Higher the values of alpha, bigger is the penalty and therefore the magnitude of coefficients are reduced.

Ridge Regression

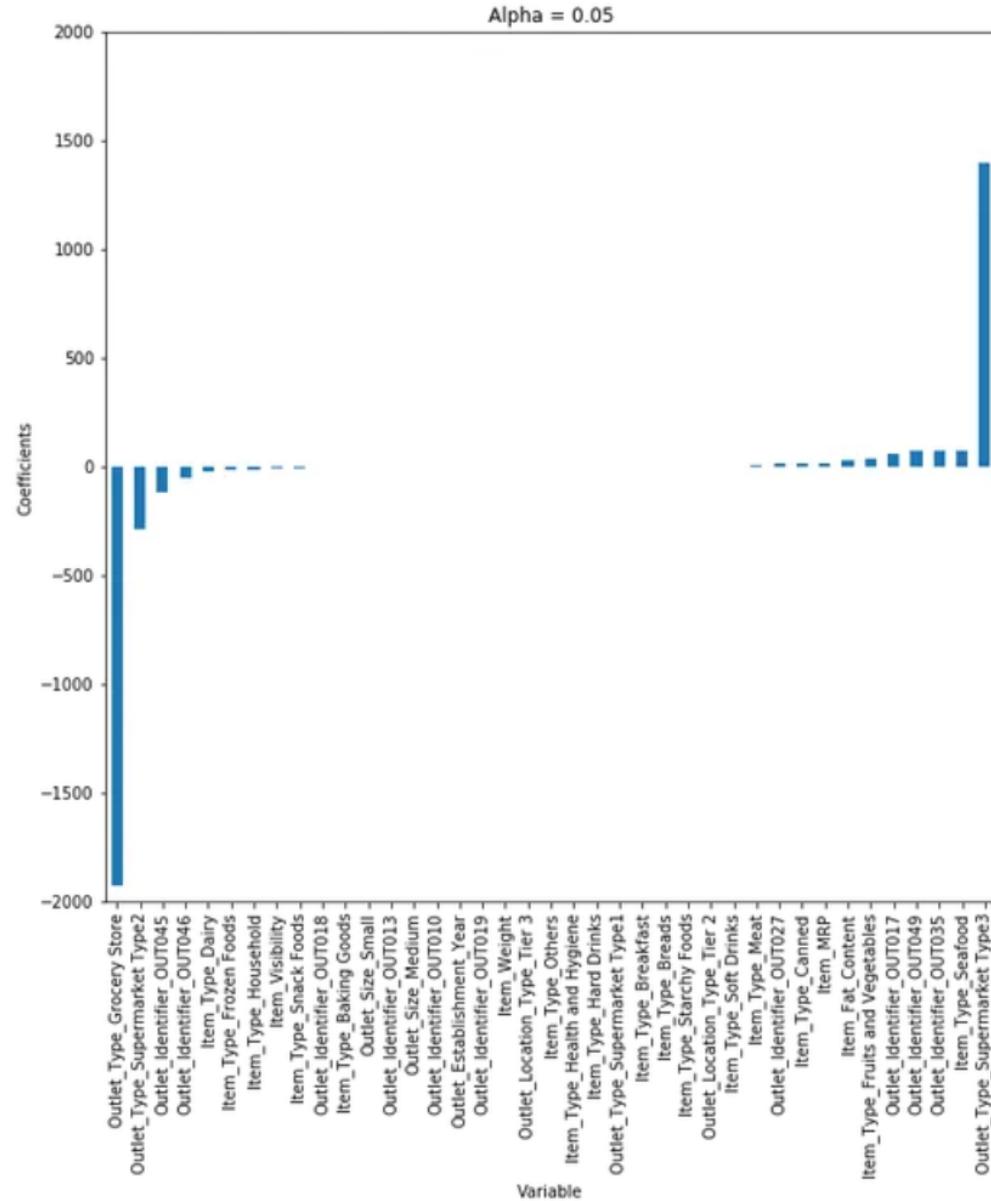
$$J(m) = \frac{1}{2n} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^k m_j^p$$

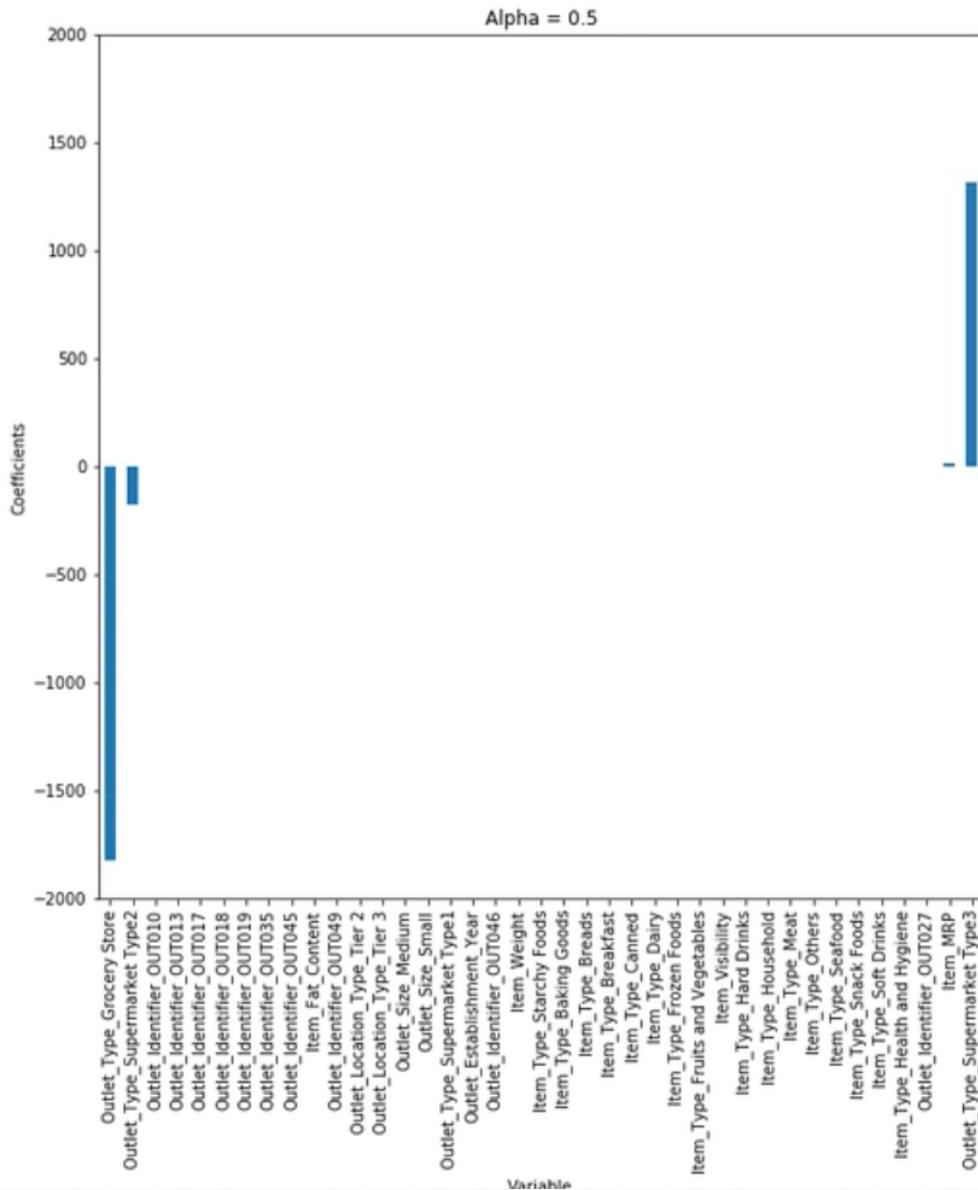
$$p = 2$$

It uses **L-2 regularization** technique.

Lasso regression

- LASSO (Least Absolute Shrinkage Selector Operator)





Any differences between ridge and lasso?

Lasso regression

- We can see that as we increased the value of alpha, coefficients were approaching towards zero.
- But if you see in case of lasso, even at smaller alpha's, our coefficients are reducing to absolute zeroes.
- Therefore, lasso selects the only some feature while reduces the coefficients of others to zero.
- This property is known as **feature selection** and which is absent in case of ridge.

Lasso regression

- It uses L1 regularization technique $p = 1$

$$J(m) = \frac{1}{2n} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^k m_j^p$$

- It is generally used when we have more number of features, because it automatically does feature selection.

Code

```
from sklearn.linear_model import Lasso
lassoReg = Lasso(alpha=0.3, normalize=True)
lassoReg.fit(x_train,y_train)
pred = lassoReg.predict(x_cv)
```

Which one to use?

- An example where we have a large dataset, lets say it has 10,000 features.
 - Ridge : model complex
 - Lasso: loose some information

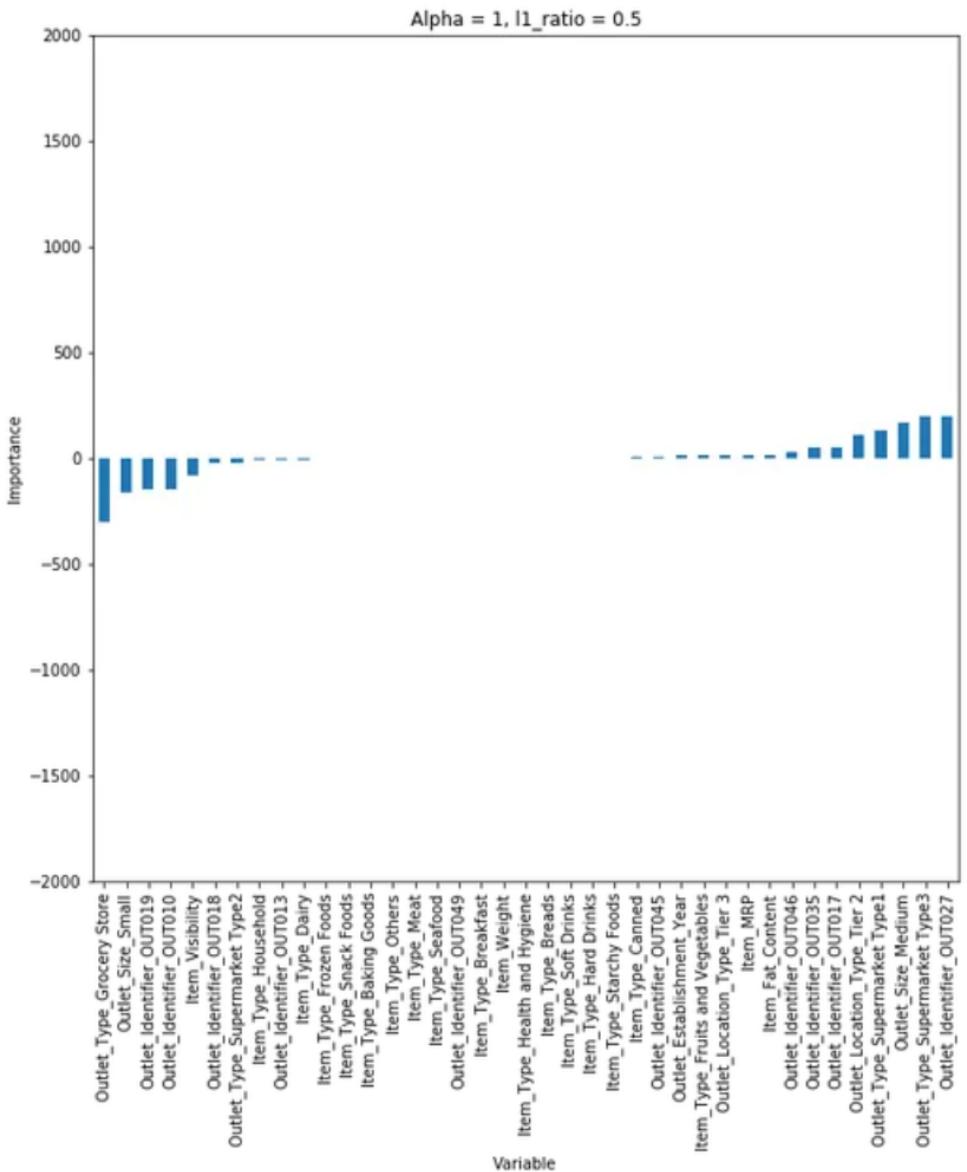
Elastic Net Regression

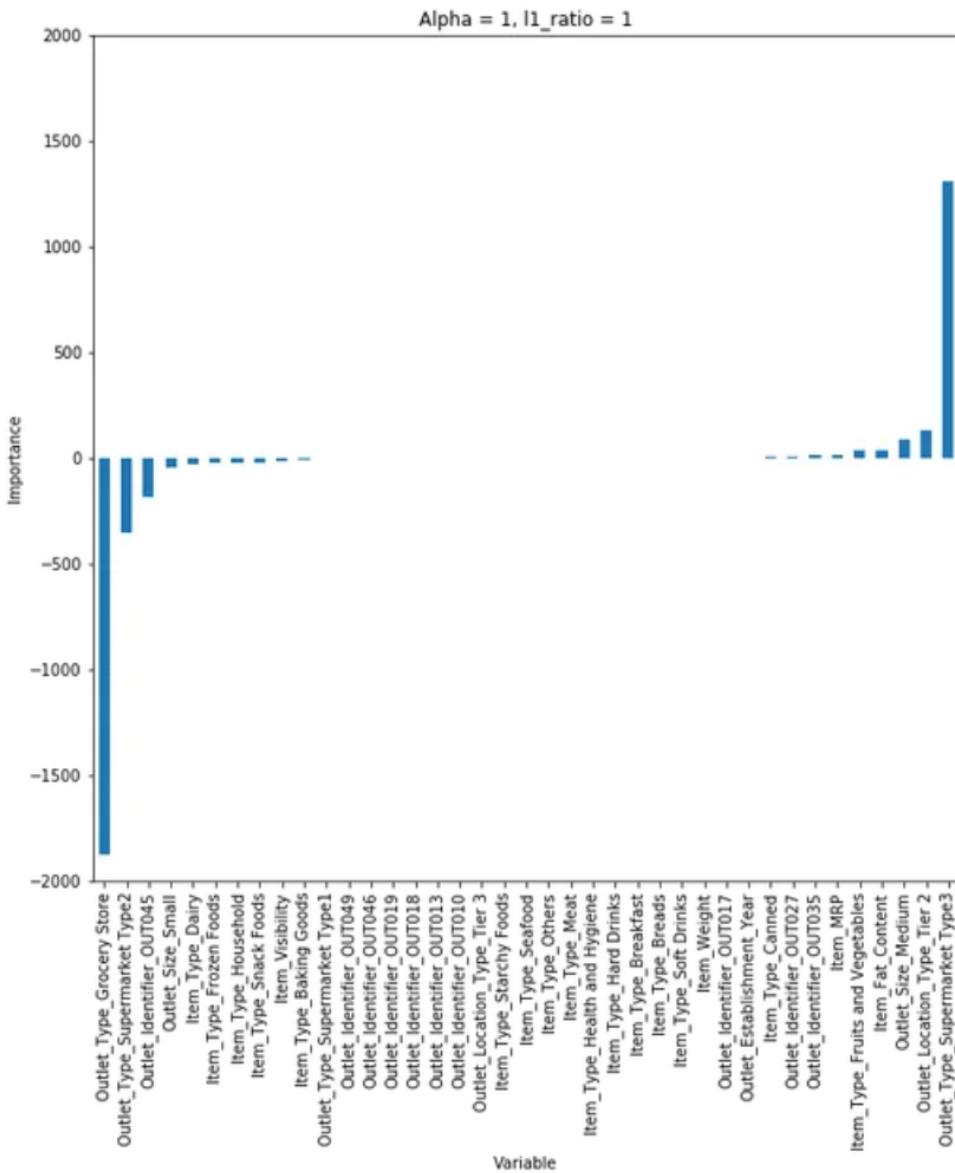
- ```
from sklearn.linear_model import
ElasticNet
ENreg = ElasticNet(alpha=1,
l1_ratio=0.5, normalize=False)
ENreg.fit(x_train,y_train)
pred_cv = ENreg.predict(x_cv)
```
- L1 and L2 regularization (you can implement both Ridge and Lasso by tuning the parameters)
  - We need to define alpha and l1\_ratio

# Elastic Net Regression

---

- $\text{Alpha} = a + b$       and       $\text{l1\_ratio} = a / (a+b)$
- Let  $\text{alpha}$  (or  $a+b$ ) = 1, and now consider the following cases:
  1. If  $\text{l1\_ratio} = 1$ , therefore if we look at the formula of  $\text{l1\_ratio}$ , we can see that  $\text{l1\_ratio}$  can only be equal to 1 if  $a=1$ , which implies  $b=0$ . Therefore, it will be a lasso penalty.
  2. Similarly if  $\text{l1\_ratio} = 0$ , implies  $a=0$ . Then the penalty will be a ridge penalty.
  3. For  $\text{l1\_ratio}$  between 0 and 1, the penalty is the combination of ridge and lasso.





# Evaluate the model

---

- How accurate do you think the model is?
- R square and adjusted R- square

$$R - Square = 1 - \frac{\sum(Y_{actual} - Y_{predicted})^2}{\sum(Y_{actual} - Y_{mean})^2}$$

# R square

---

- It determines how much of the total variation in Y (dependent variable) is explained by the variation in X (independent variable).
- The value of R-square is always between 0 and 1,
- where 0 means that the model does not model explain any variability in the target variable (Y)
- and 1 meaning it explains full variability in the target variable.

# Adjusted R-square

---

- The only drawback of  $R^2$  is that if new predictors ( $X$ ) are added to our model,  $R^2$  only increases or remains constant but it never decreases.
- We can not judge that by increasing complexity of our model, are we making it more accurate?
- The adjusted R-Square only increases if the new term improves the model accuracy.

---

*“Knowledge is the treasure and practice is the key to it”*

# Discussion

---



# Thank you!!!!!!

---

