

Linear Regression

Objectives

After completing this module, you should be able to:

Work with Linear Regression

Build Univariate Linear Regression Model

Understand Gradient Descent Algorithm

Implement Linear Regression with sklearn

Multivariate Linear Regression

Dummy Variables

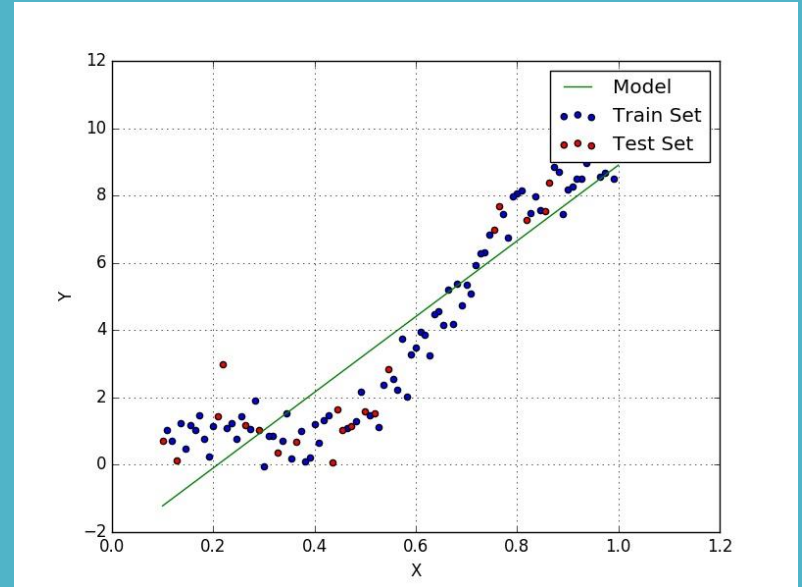
Feature Scaling

Boston Housing Prizes Prediction

Cross Validation

Regularization: Lasso & Ridge

WHAT IS LINEAR REGRESSION?



- A Supervised Learning Algorithm that learns from a set of training samples
- It estimates relationship between a dependent variable (target/label) and one or more independent variable (predictors).

Linear Regression

Univariate
Linear
Regression

$$y = m_1x_1 + c$$

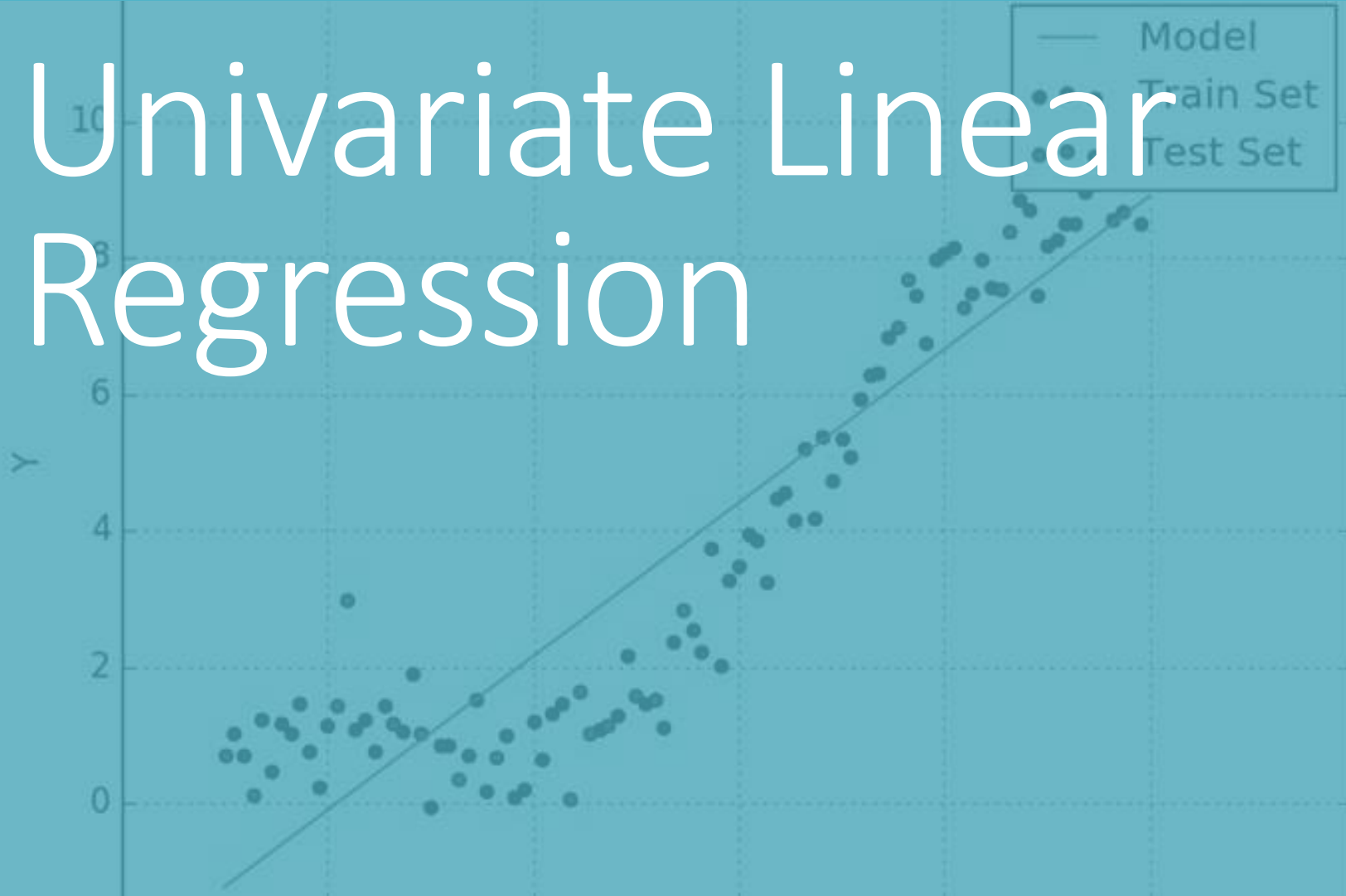
Multivariate
Linear
Regression

$$y = m_1x_1 + m_2x_2 + m_3x_3 + \dots + m_nx_n + c$$

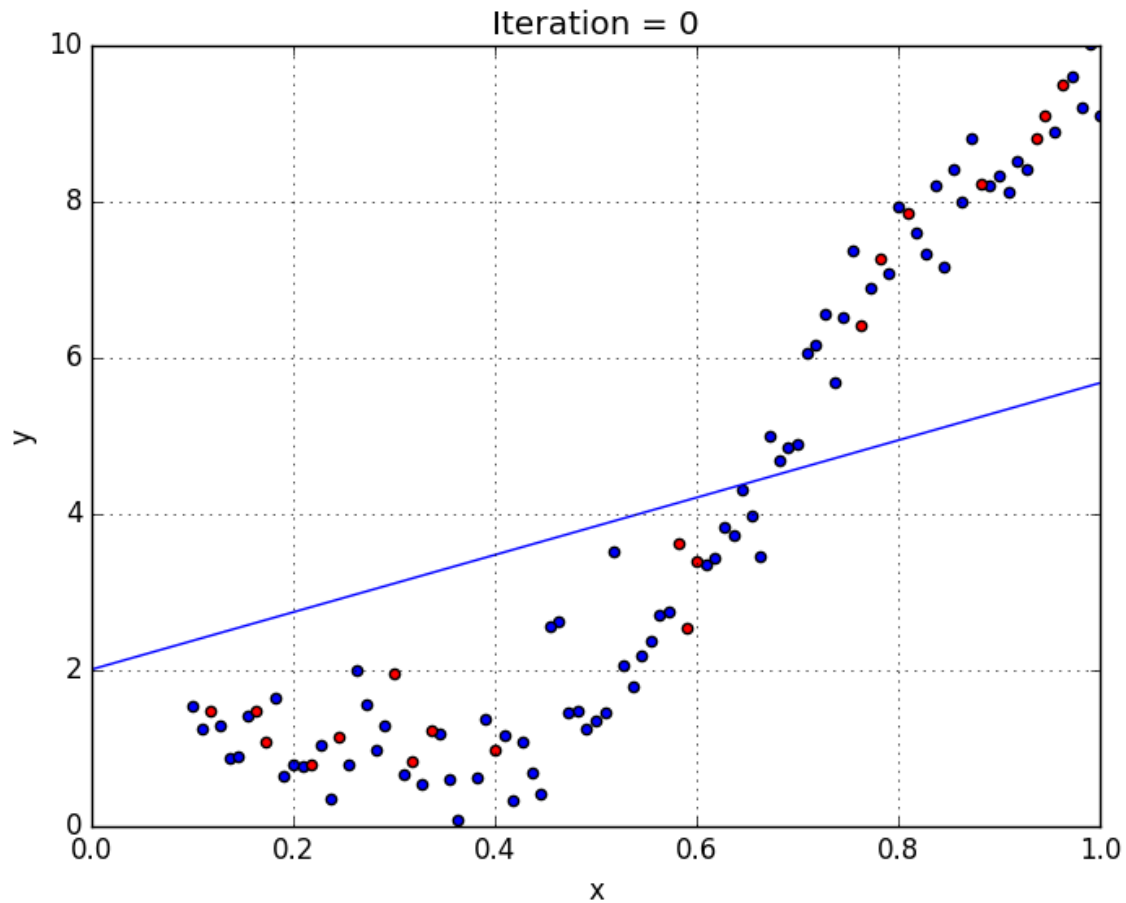
Polynomial
Linear
Regression

$$y = m_1x_1 + m_2x_1^2 + m_3x_1^3 + \dots + m_nx_1^n + c$$

Univariate Linear Regression

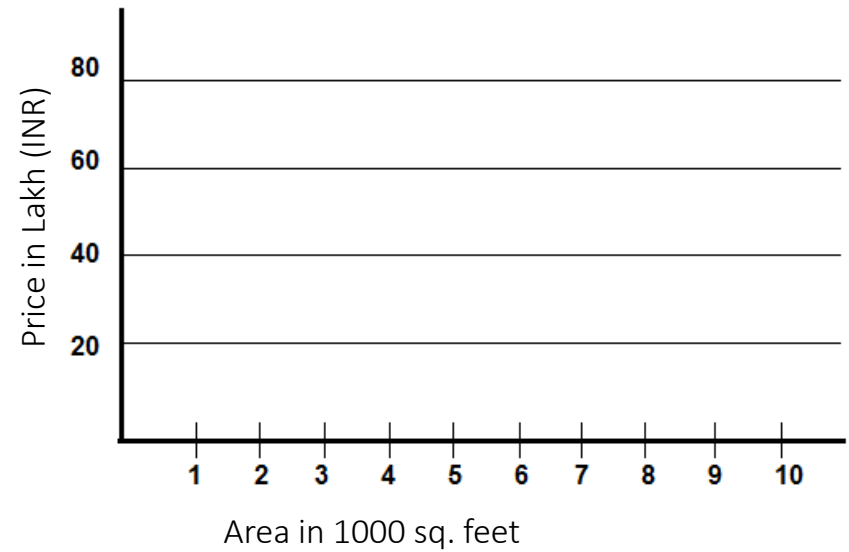


During the training period the regression line is getting more fit.



Housing Prices Prediction

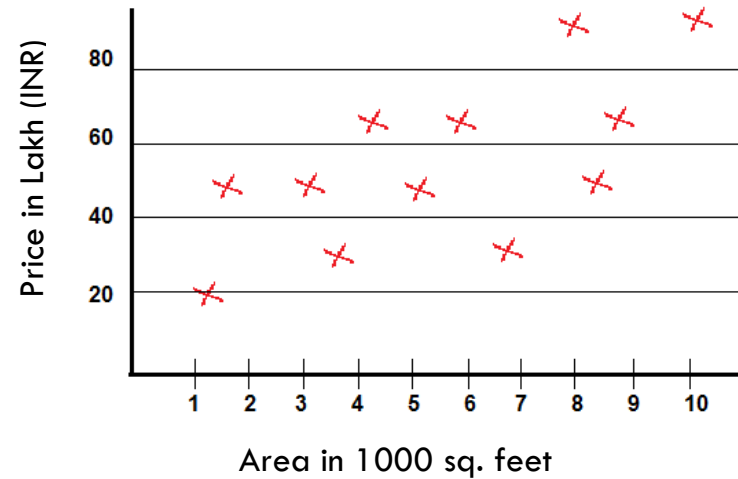
Area (sq ft)	Price In INR
1200	20,00,000
1800	42,00,000
3200	44,00,000
3800	25,00,000
4200	62,00,000



Housing Prices Prediction

y: Dependent Variable, criterion variable, or regressand.
x: Independent variable, predictor variables or regressors.

Area (sq ft) (x)	Price In INR (y)
1200	20,00,000
1800	42,00,000
3200	44,00,000
3800	25,00,000
4200	62,00,000



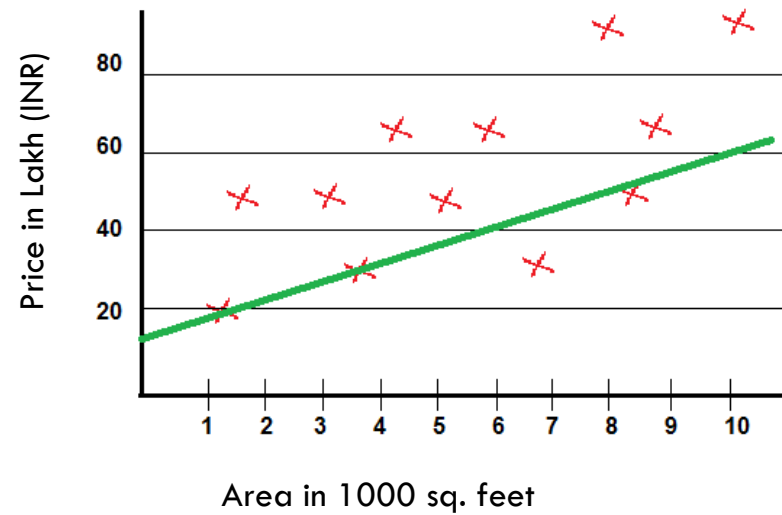
Housing Prices Prediction

$$\hat{y} = mx + c$$

\hat{y} = Value predicted by current Algorithm

Linear Regression in one Variable

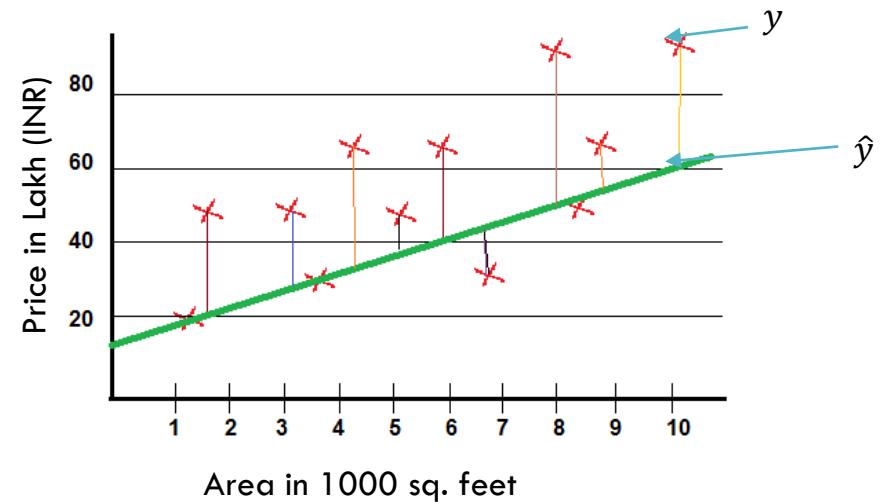
Area (sq ft)	Price In INR
1200	20,00,000
1800	42,00,000
3200	44,00,000
3800	25,00,000
4200	62,00,000



Housing Prices Prediction

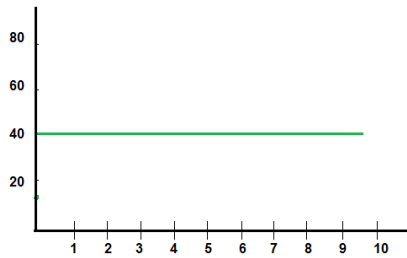
minimize
 $(y - \hat{y})$

Predictor
 $\hat{y} = mx + c$

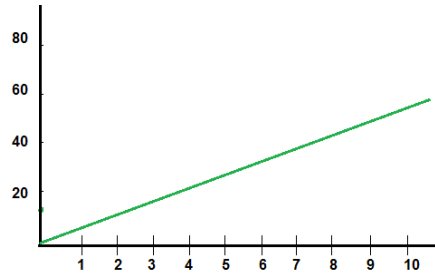


Variables affecting Regression Equation

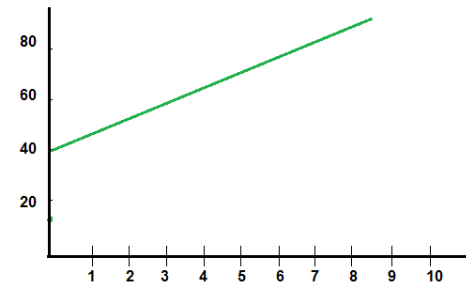
$$m = 0$$
$$c = 40$$



$$m = 0.8$$
$$c = 0$$



$$m = 0.8$$
$$c = 40$$



$$\hat{y} = mx + c$$

Housing Prices Prediction

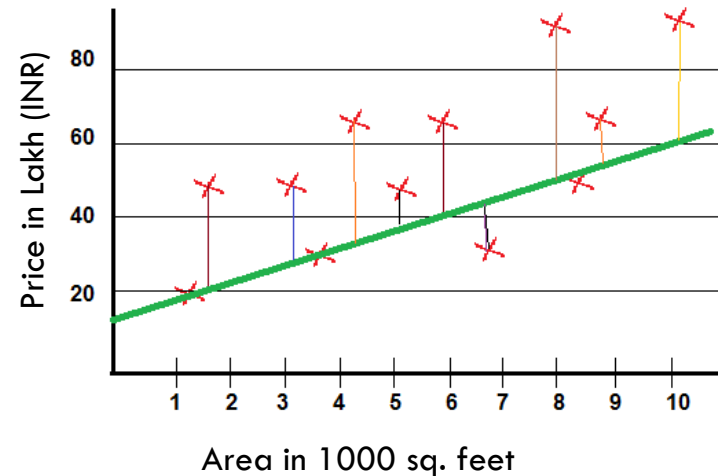
Predictor

$$\hat{y} = mx + c$$

Cost Function

$$J = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$j(m_i, c) = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



Regression Equation:

$$\hat{y} = mx + c$$

Parameters

$$m_i, c$$

Cost Function:

$$j(m_i, c) = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Goal

$$\underset{m_i, c}{\text{minimize}} J(m_i, c)$$

Gradient Descent Algorithm

Repeat Until converge

$$w_j := w_j - lr \frac{\partial}{\partial w} J(w_j)$$

simultaneously update, $j = 0, j = 1$
where, w =parameter (coefficient & constant)

Learning Rate lr

Learning Rate lr controls how big step we take while updating our parameter w .

If lr is too small, gradient descent can be slow.

If lr is too big, gradient descent can overshoot the minimum, it may fail to converge

Gradient Descent Algorithm

Repeat Until converge

$$w_j := w_j - lr \frac{\partial}{\partial w} J(w_j)$$

simultaneously update, $j = 0, j = 1$
where, w =parameter (coefficient &
constant)

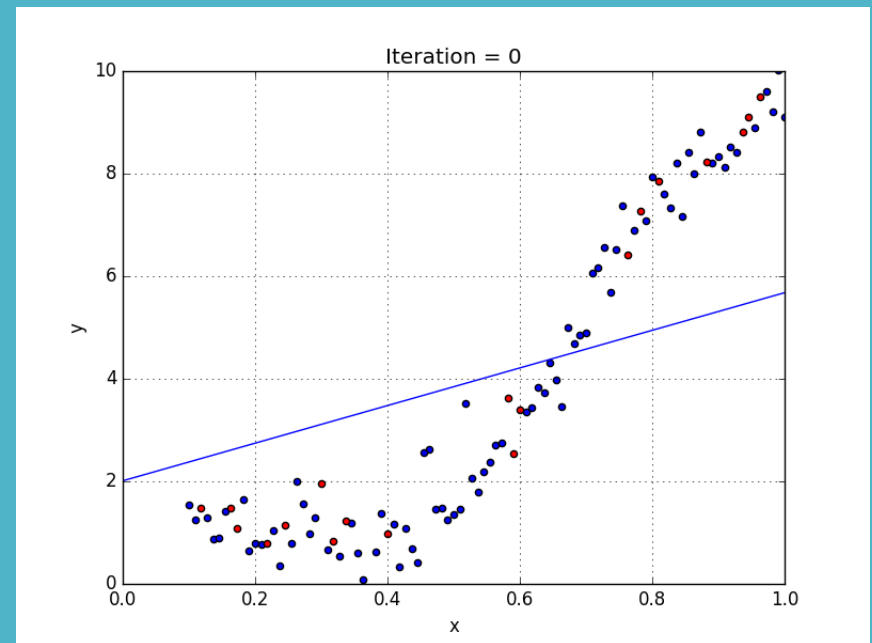
Linear Regression Model

$$\hat{y} = mx + c$$

$$j(m, c) = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Univariate Linear Regression

Linear Regression Process Visualization



Objective of Linear Regression

- Establish If there is a relationship between two variables.
Examples – relationship between housing process and area of house, no of hours of study and the marks obtained, income and spending etc.
- Prediction of new possible values
Based on the area of house predicting the house prices in a particular month; based on number of hour studied predicting the possible marks. Sales in next 3months etc.

LINEAR REGRESSION USE CASES

Real Estate

- To model residential home prices as a function of the home's living area, bathrooms, number of bedrooms, lot size.

Medicine

- To analyze the effect of a proposed radiation treatment on reducing tumor sizes based on patient attributes such as age or weight.

Demand Forecasting

- To predict demand for goods and services. For example, restaurant chains can predict the quantity of food depending on weather.

Marketing

- To predict company's sales based on previous month's sales and stock prices of a company.

Multiple Linear Regression

One Hot Encoding

When some inputs are categories (e.g. gender) rather than numbers (e.g. age) we need to represent the category values as numbers so they can be used in our linear regression equations.

Dummy Variables

Dummy Variables



Salary	Credit Score	Age	State	New York	California
192,451	485	42	New York	1	0
118,450	754	35	California	0	1
258,254	658	28	California	0	1
200,123	755	48	New York	1	0
152,485	654	52	California	0	1

Encoding Categorical Data

```
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import OneHotEncoder
labelencoder = LabelEncoder()
#considering X is dataset from above slide
# 3 is the index number of state
X[:, 3] = labelencoder.fit_transform(X[:, 3])
onehotencoder = OneHotEncoder(categorical_features = [3])
X = onehotencoder.fit_transform(X).toarray()
```

Avoiding the Dummy variable trap

```
X=X[:,1:]
```

NOTE : if you have n dummy variables remove one dummy variable to avoid the dummy variable trap. However the linear regression model that is built in R and Python takes care of this. But there is no harm in removing it by ourselves

Feature Scaling

(Mean Normalization) Standardization	(Min-Max Normalization) Normalization
$X_{stand} = \frac{x - mean(x)}{standard_deviation(x)}$	$X_{norm} = \frac{x - min(x)}{max(x) - min(x)}$

Standard Scale using sklearn

```
from sklearn.preprocessing import StandardScaler  
sc_x = StandardScaler()  
sc_y = StandardScaler()  
X_std = sc_x.fit_transform(X)  
y_std = sc_y.fit_transform(y)
```

Cross Validation

Split 1	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 1
Split 2	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 2
Split 3	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 3
Split 4	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 4
Split 5	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Metric 5



Cross-validation and model performance

- 5 folds = 5-fold CV
- 10 folds = 10-fold CV
- k folds = k-fold CV
- More folds = More computationally expensive

Cross-validation in scikit-learn

```
In [1]: from sklearn.model_selection import cross_val_score
In [2]: alg = linear_model.LinearRegression()
In [3]: cv_results = cross_val_score(alg, X, y, cv=5)
In [4]: print(cv_results)
[ 0.63919994  0.71386698  0.58702344  0.07923081 -0.25294154]
In [5]: numpy.mean(cv_results)
Out[5]: 0.35327592439587058
```

Overfitting & Generalisation

As we train our model with more and more data the it may start to fit the training data more and more accurately, but become worse at handling test data that we feed to it later.

This is known as “over-fitting” and results in an increased generalization error.

Large coefficients lead to overfitting

Penalizing large coefficients: Regularization

How to minimize?

- To minimize the generalization error we should
- Collect as much sample data as possible.
- Use a random subset of our sample data for training.
- Use the remaining sample data to test how well our model copes with data it was not trained with.

L1 Regularisation (Lasso)

(Least Absolute Shrinkage and Selection Operator)

- Having a large number of samples (n) with respect to the number of dimensionality (d) increases the quality of our model.
- One way to reduce the effective number of dimensions is to use those that most contribute to the signal and ignore those that mostly act as noise.
- L1 regularization achieves this by adding a penalty that results in the weight for the dimensions that act as noise becoming 0.
- L1 regularization encourages a sparse vector of weights in which few are non-zero and many are zero.

L1 Regularisation (Lasso)

Depending on the regularization strength, certain weights can become zero, which makes the LASSO also useful as a supervised feature selection technique:

$$j(w_i) = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|w_i\|$$

A limitation of the LASSO is that it selects at most n variables if $m > n$.

Lasso regression in scikit-learn

```
In [1]: from sklearn.linear_model import Lasso
In [2]: X_train, X_test, y_train, y_test = train_test_split(X,
y,test_size = 0.3, random_state=42)
In [3]: lasso = Lasso(alpha=0.1, normalize=True)
In [4]: lasso.fit(X_train, y_train)
In [5]: lasso_pred = lasso.predict(X_test)
In [6]: lasso.score(X_test, y_test)
Out[6]: 0.59502295353285506
```

L2 Regularisation (Ridge)

- Another way to reduce the complexity of our model and prevent overfitting to outliers is L2 regression, which is also known as ridge regression.
- In L2 Regularization we introduce an additional term to the cost function that has the effect of penalizing large weights and thereby minimizing this skew.

L2 Regularisation (Ridge)

Ridge regression is an L2 penalized model where we simply add the squared sum of the weights to our least-squares cost function:

$$j(\mathbf{w}_i) = \frac{1}{2n} \sum_{i=1}^n (\mathbf{y}_i - \hat{\mathbf{y}}_i)^2 + \lambda \|\mathbf{w}_i\|^2$$

By increasing the value of the hyperparameter λ , we increase the regularization strength and shrink the weights of our model.

Ridge regression in scikit-learn

```
In [1]: from sklearn.linear_model import Ridge
In [2]: X_train, X_test, y_train, y_test = train_test_split(X,
y,test_size = 0.3, random_state=42)
In [3]: ridge = Ridge(alpha=0.1, normalize=True)
In [4]: ridge.fit(X_train, y_train)
In [5]: ridge_pred = ridge.predict(X_test)
In [6]: ridge.score(X_test, y_test)
Out[6]: 0.69969382751273179
```

L1 & L2 Regularisation (Elastic Net)

L1 Regularisation minimises the impact of dimensions that have low weights and are thus largely “noise”.

L2 Regularisation minimise the impacts of outliers in our training data.

L1 & L2 Regularisation can be used together and the combination is referred to as Elastic Net regularisation.

Because the differential of the error function contains the sigmoid which has no inverse, we cannot solve for w and must use gradient descent.

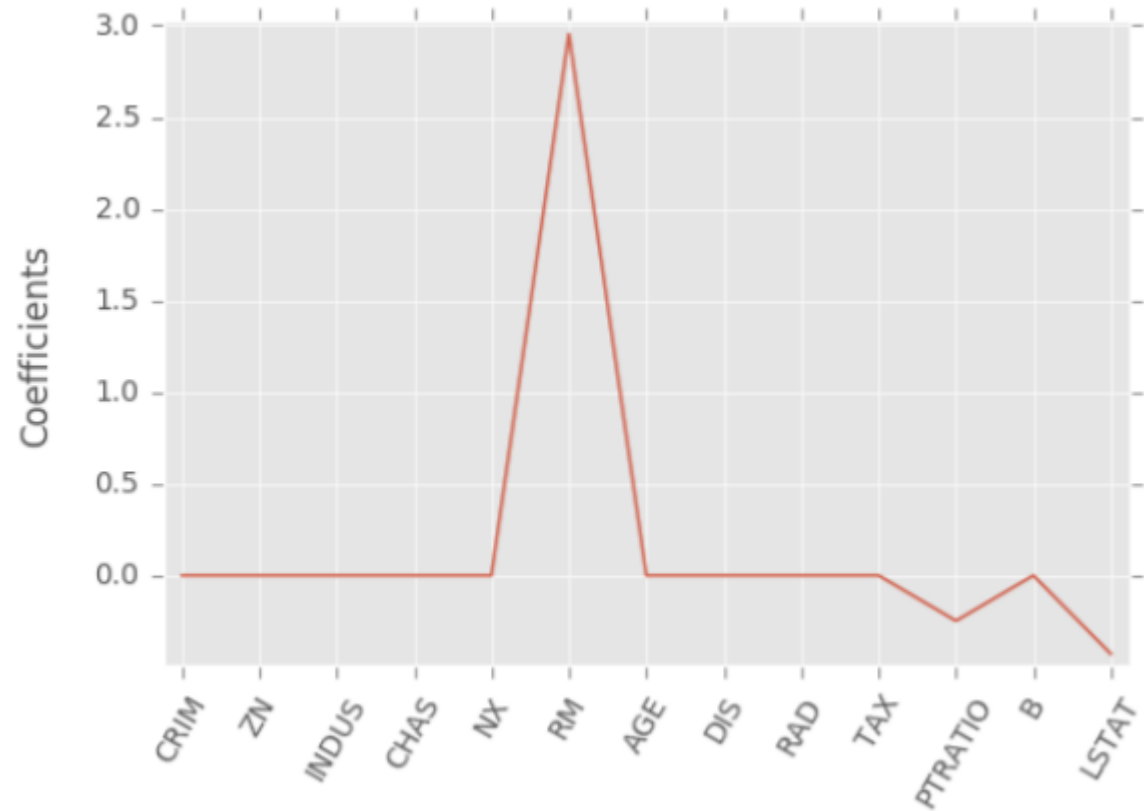
Lasso regression for feature selection

- Can be used to select important features of a dataset
- Shrinks the coefficients of less important features to exactly 0.

Lasso regression for feature selection

```
In [1]: from sklearn.linear_model import Lasso
In [2]: names = boston.drop('MEDV', axis=1).columns
In [3]: lasso = Lasso(alpha=0.1)
In [4]: lasso_coef = lasso.fit(X, y).coef_
In [5]: plt.plot(range(len(names)), lasso_coef)
In [6]: plt.xticks(range(len(names)), names, rotation=60)
In [7]: plt.ylabel('Coefficients')
In [8]: plt.show()
```


Lasso regression for feature selection



Summary

This module covered the following topics:

Linear Regression

Gradient Descent Algorithm

Multiple Linear Regression

Data Manipulation

Regularization

Thank you