

## Regression

- Regression is the process of predicting a continuous value.
- $x$  : Independent variable
- $y$  : Dependent variable  $\rightarrow$  should be continuous
- Types of regression model:
  1. Simple ~~Linear~~ Regression
    - $\rightarrow$  Simple Linear Regression
    - $\rightarrow$  Simple Non-linear Regression.
  2. Multiple Regression
    - $\rightarrow$  Multiple Linear Regression.
    - $\rightarrow$  Multiple Non-linear Regression.
- Applications of Regression
  1. Sales forecasting
  2. Satisfaction analysis
  3. Price estimation
  4. Employment income

## • Regression Algorithms

1. Ordinal regression
2. Poisson regression
3. Fast forest quantile regression
4. Linear, Polynomial, Lasso, Stepwise, Ridge regression.
5. Bayesian linear regression.
6. Neural Network Regression.
7. Decision forest Regression.
8. KNN (K-nearest neighbors).

## • Simple Linear Regression

- we can plot scatter plot b/w Dependent & independent variable to find the correlation b/w them.

### • Fitting line

$$\hat{y} = \theta_0 + \theta_1 x_1$$

response or target variable →  $\hat{y}$  = dependent variable  
or  
a single predictor

$\theta_0$  &  $\theta_1$  → Coefficient of equations.

- The objective of linear regression is to minimize the MSE (Mean Squared Error).
- How to calculate  $\theta_0$  &  $\theta_1$  mathematically?

$$\hat{y} = \theta_0 + \theta_1 x_i$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\theta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

### Pros of Linear Regression

- very fast
- No Parameter Tuning
- Easy to understand, and highly interpretable.

### Model Evaluation Approaches

- Train and Test on the same dataset
- Train / Test split.

- what is an error of one model?

Error: measure of how far the data is from the fitted regression line.

- ↳ Mean Absolute Error

$$MAE = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

- ↳ mean squared Error

$$MSE = \frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2$$

- ↳ Root Mean Squared Error

$$RMSE = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

- ↳ Relative Absolute Error

$$RAE = \frac{\sum_{j=1}^n |y_j - \hat{y}_j|}{\sum_{j=1}^n |y_j - \bar{y}|}$$

- ↳ Relative Square Error

$$RSE = \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{\sum_{j=1}^n (y_j - \bar{y})^2}$$

$$R^2 = 1 - RSE$$

## • Fuel consumption

### • Reading the data

```
df = pd.read_csv("FuelConsumption.csv")
```

### • Data Exploration

```
# summarize the data
```

```
df.describe()
```

```
# let's select some features to explore more
```

```
cdf = df[['ENGINESIZE', 'CYLINDERS', 'FUELCONSUMPTION-COMB', 'CO2EMISSIONS']]
```

```
# we can plot each of these features:
```

```
viz = cdf[[ ]]
```

```
viz.hist()
```

```
plt.show()
```

```
# Now, let's plot each of these features vs  
the EMISSION, to see how linear is their  
relation:
```

```
plt.scatter(cdf.FUELCONSUMPTION-COMB,  
            cdf.CO2EMISSIONS, color='blue')
```

```
plt.xlabel("FUELCONSUMPTION-COMB")
```

```
plt.ylabel("Emission")
```

```
plt.show()
```

# Let's split our dataset into train and test sets, 80% of the entire data for training, and the 20% for testing. We can create a mask to select random rows using np.random.rand() function:

```
msk = np.random.rand(len(df)) < 0.8  
train = df[msk]  
test = df[~msk]
```

# modeling

Using sklearn package to model data.

```
from sklearn import linear_model  
regr = linear_model.LinearRegression()  
train_x = np.asarray(train[['ENGINESIZE']])  
train_y = np.asarray(train[['CO2EMISSIONS']])  
regr.fit(train_x, train_y)
```

# The coefficients

```
print("Coefficients: ", regr.coef_ )  
print("Intercept: ", regr.intercept_ )
```

## # PLOT outputs

We can plot the fit line over the data:

```
plt.scatter(train.ENGINESIZE, train.CO2EMISSIONS,  
           color='blue')  
plt.plot(train_x, regr.coef_[0][0]*train_x  
         + regr.intercept_[0], '-r')  
plt.xlabel("Engine Size")  
plt.ylabel("Emission")
```

## • Evaluation

```
from sklearn.metrics import r2_score  
test_x = np.asarray(test[['ENGINESIZE']])  
test_y = np.asarray(test[['CO2EMISSIONS']])  
test_y_hat = regr.predict(test_x)  
  
print("Mean absolute error: %.2f" %  
      np.mean(np.absolute(test_y_hat - test_y)))  
print("Residual sum of squares (MSE): %.2f" %  
      np.mean((test_y_hat - test_y) ** 2))  
print("R2-score: %.2f" % r2_score(test_y_hat,  
                                 test_y))
```

05/06/2020

## Introduction to Regression

[DATACAMP]

### # Simple linear regression

```
# Import numpy and pandas
```

```
import numpy as np
```

```
import pandas as pd
```

```
# Read the csv file into a dataframe: df
```

```
df = pd.read_csv("gapminder.csv")
```

```
# Create arrays for features and target variable
```

```
y = df['life'].values
```

```
x = df['fertility'].values
```

```
# Print the dimension of X & y before reshaping.
```

```
print("Dimension of y before reshaping: {}".format(y.shape))
```

```
print("Dimension of x before reshaping: {}".format(x.shape))
```

```
# Reshape x & y
```

```
y = y.reshape(-1, 1)
```

```
x = x.reshape(-1, 1)
```

```
# Print the dimension of x & y after reshaping
```

## # Exploring Gapminder data

→ df.info()

→ df.describe()

→ df.head()

## # Seaborn heatmap function

sns. heatmap(df.corr(), square=True, cmap='RdYlGn')

'RdYlGn')

06/06/2020

[DATACAMP]

## Regression Mechanics

- $y = ax + b$ 
  - $y$  = target
  - $x$  = single feature
  - $a, b$  = parameters of model
- How do we choose  $a$  &  $b$ ?
- Define an error function for any given line
  - choose the line that minimizes the error function.

### The loss function

- ordinary least squares (OLS):
  - Minimize sum of squares of residuals.

### Linear Regression in higher dimension

$$y = a_1x_1 + a_2x_2 + b$$

- To fit a linear regression model we:
  - Need to specify 3 variables  
 $a_1, a_2$  &  $b$

- In higher dimensions:
  - Must specify coefficient for each feature and one variable  $b$ .
- Scikit-learn API works exactly the same way:
  - Pass 2 arrays: features and target
- Linear regression on all features

```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
```

$x_{\text{train}}, x_{\text{test}}, y_{\text{train}}, y_{\text{test}} = \text{train\_test\_split}(x, y,$   
 $\text{test\_size} = 0.3, \text{random\_state} = 42)$

```
reg_all = LinearRegression()
```

```
reg_all.fit(x_train, y_train)
```

```
y_pred = reg_all.predict(x_test)
```

```
reg_all.score(x_test, y_test)
```

## Fit & Predict for Regression [working on gapminder dataset]

```
# Import linear regression  
from sklearn.linear_model import LinearRegression  
  
# Create the regressor: reg  
reg = LinearRegression()  
  
# Create the prediction space  
prediction_space = np.linspace(min(X_fertility),  
                               max(X_fertility)).reshape(-1,1)  
  
# Fit the model to the data  
reg.fit(X_fertility, y)  
  
# Compute predictions over the prediction space:  
y_pred = reg.predict(prediction_space)  
  
# Print R2  
print('R^2 score:', reg.score(X_fertility, y))  
  
# Plot regression line  
plt.plot(prediction_space, y_pred,  
          color='black', linewidth=3)  
plt.show()
```

## train/test-split for Regression

```
# Import necessary modules
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import train_test_split
```

```
# Create training & test sets.
```

```
X_train, X_test, y_train, y_test =
train_test_split(X, y, test_size=0.3,
random_state=42)
```

```
# Create the regressor: reg_all
```

```
reg_all = LinearRegression()
```

```
# Fit the regressor to the training data
reg_all.fit(X_train, y_train)
```

```
# Predict on the test data: y_pred
```

```
y_pred = reg_all.predict(X_test)
```

```
# Compute & print R2 & RMSE
```

```
print("R2: {}".format(reg_all.score(X_test, y_test)))
```

```
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
```

```
print("Root Mean Squared Error: {}".format(rmse))
```

## ⇒ Cross-validation motivation

- Model performance is dependent on way the data is split.
- Not representative of the model's ability to generalize.
- Solution: Cross-validation!

## Cross-validation basics

split 1	fold 1	fold 2	fold 3	fold 4	fold 5	metric 1
split 2	1	2	3	4	5	metric 2
split 3	1	2	fold 3	4	5	3
split 4	1	2	3	fold 4	5	4
split 5	1	2	3	4	fold 5	5

Training data      Test data

- 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

```
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LinearRegression
reg = LinearRegression()
cv_results = cross_val_score(reg, X, y, cv=5)
print(cv_results)
```

### • 5-fold cross validation

# Import the necessary modules

```
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LinearRegression
```

# Create a LinearRegression object : reg

```
reg = LinearRegression()
```

# Compute 5-fold cross-validation scores:

cv\_scores

```
cv_scores = cross_val_score(reg, X, y, cv=5)
```

# Print the 5-fold cross-validation scores

```
print(cv_scores)
```

```
print("Average 5-fold CV scores: {}".format(np.mean(cv_scores)))
```

time

# To see how much time a single cross-validation score takes or  $K$ -CV takes.

% timeit cross\_val\_score (reg, X, y, cv=

$K$ )

Here  $K \rightarrow$  can be

1, 2, 3, ...,  $K$ .

import timeit

timeit

and you get number of times each fold

is used. For example, if we have 5 folds

Quantity  $\rightarrow$  Time  $\rightarrow$  5

for example, if we have 5 folds, quantity will be 5

Time  $\rightarrow$  5 \* Time of each fold

Example: If I want to calculate time for

cross\_val\_score (reg, X, y, cv=5)

Time  $\rightarrow$  5 \* Time of each fold

Time of each fold = Time of one fold

Time of one fold = Time of one iteration

Iteration  $\rightarrow$  Number of iterations

Number of iterations = Number of folds

Number of folds = 5

Time of one iteration = Time of one fold

Time of one fold = Time of one fold

## Multiple Linear Regression [L18M]

### 1. Examples of Multiple Linear Regression

- Independent variable effectiveness on prediction

- Does revision time, test anxiety, lecture attendance and gender have any effect on the exam performance of students?

### 2. Predicting Impact of changes

- How much does blood pressure go up & down for every unit increase or decrease in the BMI of a patient?

### 3. Estimating multiple linear Regression Parameters

- How to estimate  $\theta$ ?

- Ordinary Least Squares

- Linear Algebra Operations

- Takes a long time for large datasets ( $10k + \text{rows}$ )

- An optimization algorithm

- Gradient Descent

- Proper approach if you have a very large dataset.

## Q. What is Polynomial Regression?

- Some curvy data can be modelled by a polynomial regression.
- for e.g.  $\hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$

08/06/2020

## 2. Regularized Regression

[DATACAMP]

### 2. why Regularize:

- Recall: Linear Regression minimizes a loss function.
- It chooses a coefficient for each feature variable.
- Large coefficients can lead to overfitting.
- Penalize large coefficients: regularization.

#### [1] Ridge Regression

- Loss function = OLS loss function +  $\alpha \cdot \sum_{i=1}^n a_i^2$
- Alpha: Parameter we need to choose
- Picking alpha here is similar to picking  $k$  in KNN.
- Alpha controls model complexity
  - Alpha = 0: We get back OLS (can lead to overfitting).
  - very high alpha: Can lead to underfitting

## Ridge Regression in scikit-learn

```
from sklearn.linear_model import Ridge  
x_train, x_test, y_train, y_test = train_test_split  
(x, y, test_size=0.3, random_state=42)  
ridge = Ridge(alpha=0.1, normalize=True)  
ridge.fit(x_train, y_train)  
ridge_pred = ridge.predict(x_test)  
ridge.score(x_test, y_test)  
↳ output = 0.69969.
```

## [2] Lasso Regression

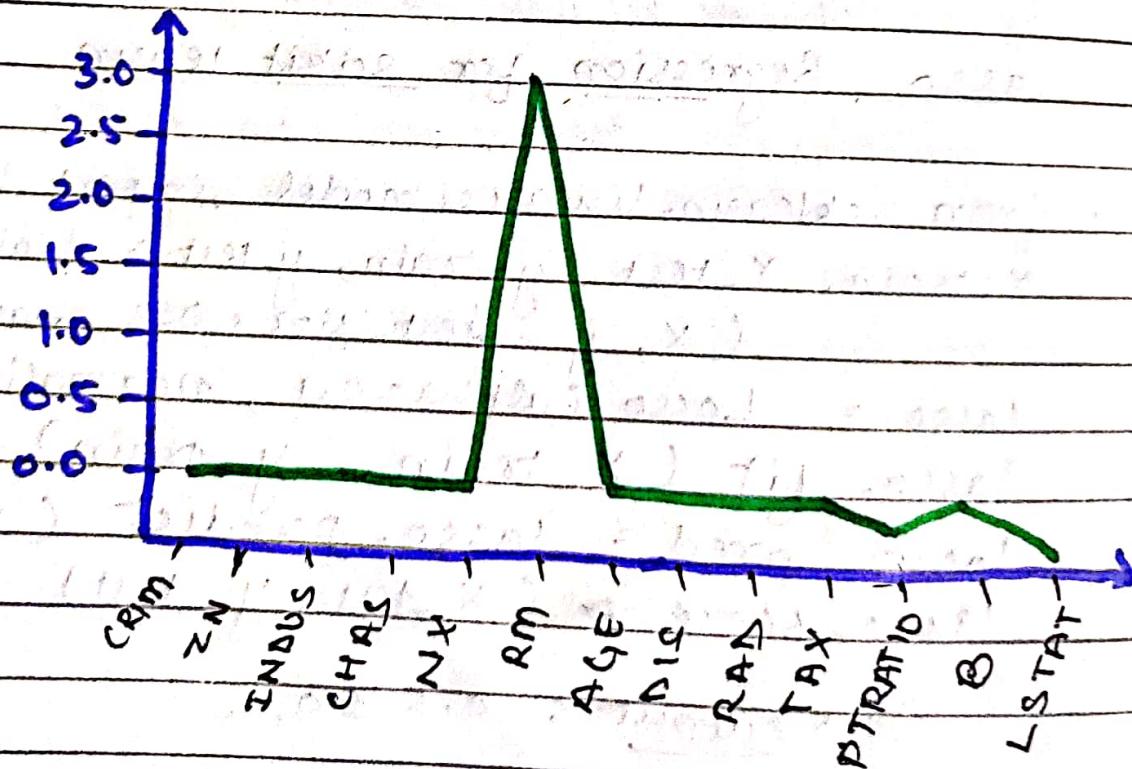
- loss function = OLS loss function +  
$$\lambda \sum_{i=1}^n |a_i|$$

## Lasso Regression for scikit-learn

- from sklearn.linear\_model import Lasso  
x\_train, x\_test, y\_train, y\_test = train\_test\_split  
(x, y, test\_size=0.3, random\_state=42)  
lasso = Lasso(alpha=0.1, normalize=True)  
lasso.fit(x\_train, y\_train)  
lasso\_pred = lasso.predict(x\_test)  
lasso.score(x\_test, y\_test)  
↳ output = 0.5950.

## 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.
- from sklearn.linear\_model import Lasso  
names = boston.drop('MEDV', axis=1).columns  
lasso = Lasso(alpha=0.1)  
lasso.coef\_ = lasso.fit(X, y).coef\_  
- = plt.plot(range(len(names)), lasso.coef\_)  
- = plt.xticks(range(len(names)), names,  
rotation=60)  
- = plt.ylabel('coefficients')  
plt.show()



### FINE-TUNING YOUR MODEL

HOW GOOD IS YOUR MODEL?

#### + Classification metrics

- Measuring model performance "with accuracy":
  - fraction of correctly classified samples.
  - Not always a useful metric.

check imbalanced examples: emails

- Spam classification
  - 99% of emails are real, 1% of emails are spam.
- Could build a classifier that predicts all emails as real
  - 99% accurate!
  - But horrible at actually classifying spam.
  - Fails at its original purpose.
- Need more nuanced metrics.

## Diagnosing classification predictions

- confusion matrix

	Predicted: spam Email	Predicted: Real Email
Actual: spam Email	TRUE POSITIVE	FALSE NEGATIVE
Actual: Real Email	FALSE POSITIVE	TRUE NEGATIVE

- Accuracy:

$$\frac{TP + TN}{TP + TN + FP + FN}$$

Metrics from the confusion matrix

- Precision:

$$\frac{TP}{TP + FP}$$

- Recall:

$$\frac{TP}{TP + FN}$$

- F1 score:  $\frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$

- High Precision: Not many real emails predicted as spam.
- High Recall: Predicted most spam emails correctly.

### confusion matrix in scikit learn

```
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
```

```
knn = KNeighborsClassifier(n_neighbors = 8)
```

```
X_train, X_test, y_train, y_test = train_test_split
(X, y, test_size = 0.4, random_state=42)
```

```
knn.fit(X_train, y_train)
```

```
y_pred = knn.predict(X_test)
```

```
print(confusion_matrix(y_test, y_pred))
```

↳ output

```
[ [52  7]
  [3 112]]
```

```
print(classification_report(y_test, y_pred))
```

↳ output

	Precision	Recall	f1-score	Support
0	0.95	0.88	0.91	59
1	0.94	0.97	0.96	115
avg / total	0.94	0.94	0.94	174