

ASAD ASHRAF KAREL

All syntaxes for Statistics:

Bionomial, Poission, Normaml Distribution

```
In [ ]: 1 # Binomail Distribution
2 st.binom.cdf/sf/pmf(x,n,p)
3
4 # Poisson Distribution
5 st.poisson.cdf/sf/pmf(x,mean)
6
7 # Normal Distribution
8 st.norm.cdf/sf/pmf(x,mean,std)
9
10 st.norm.isf(0.7,mean,std) # Exmple: At least 70% customers have balance in t
11 st.norm.ppf(0.3,mean,std)
12
```

Calculation for ZScore value

```
In [ ]: 1 z1 = mean + std
2 z2 = mean - std
3
4 st.norm.cdf(z1,mean,std) - st.norm.cdf(z2,mean,std)
5 # 0.68
6
7 z1 = mean + 2*std
8 z2 = mean - 2*std
9
10 st.norm.cdf(z1,mean,std) - st.norm.cdf(z2,mean,std)
11 # 0.95
12
13 z1 = mean + 3*std
14 z2 = mean - 3*std
15
16 st.norm.cdf(z1,mean,std) - st.norm.cdf(z2,mean,std)
17 # 0.99
```

Sampling Error

Difference between the Pop Mean and the Sample Mean

```
In [ ]: 1 np.mean(population) - np.mean(sample)
```

Confidence Interval with z distribution

```
In [ ]: 1 # Std Error:
2 Std_Error = sigma/np.sqrt(n)
3
4 # OR:
5 Std_Error = st.sem(sample) # np.std(mins, ddof = 1)/np.sqrt(n)
6
```

Z_values Estimation:

```
In [ ]: 1 st.norm.isf(0.95) # for 95% area
2
3 # With upper Limit:
4 st.norm.isf(0.95) , st.norm.ppf(0.95) , st.norm.interval(0.95) # Here we get
5
6 # Upperlimit - Lower Limit: (CONFIDENCE INTERVAL)
7 values(UL,LL) = Mean - (z_alpha * std) , Mean + (z_alpha * std)
8
9 # Direct Estimation of Confidence Interval:
10 st.norm.interval(0.95, mean)
11 ### OR m - st.norm.isf(0.025) , m + st.norm.isf(0.025) Mathematically
```

Confidence interval with t distribution

```
In [ ]: 1 # t_critical_values:
2 st.t.isf(0.99, df = n-1), st.t.ppf(0.99, df = n-1) # For 99%.
3
4 # CI for t distribution:
5 t_stat = st.t.ppf(0.99, df = n-1)
6 SE = s/np.sqrt(n)
7
8 CI = xbar - (SE*t_stat) , xbar + (SE*t_stat)
9
10 # Direct Estimation of Confidence Interval:
11 st.t.interval(0.95, dof= n-1, xbar, SE)
```

Hypothesis Technique

Complete process:

If population std is given:

```
In [ ]: 1
2 # Given data:
3 Mu =
4 n =
5 sigma =
6 xbar =
7 alpha =
8
9 ### Critical Value:
10 zcrit = stats.norm.isf(0.975).round(3) # 0.025
11
12 # Z_stats:
13 std_error=sigma/np.sqrt(n)
14
15 z_stat=(xbar-Mu)/std_error
16
17 print(zcrit , z_stat)
18
19 ## If z_stat > zcrit, we reject our Null hypothesis.
20
21 # With P_value calculations:
22 p_value=st.norm.cdf(z_stat)*2
23 print(p_value)
24
25 ## If P_value > alpha, we accept our Null hypothesis.
```

If proportion is given:

```
In [ ]: 1 zcrit = st.norm.isf(0.05)
2 print(zcrit)
3
4 # Proportion:
5 p1 = x/n
6 q = 1 - p
7
8 zstat = (p1 - p) / np.sqrt(p*q/n)
9 print(zstat)
10
11 if zcrit > zstat :
12     print('\nWe are failed to reject null hypothesis.')
13     print('Student internship is across 45%.')
14 else:
15     print('\nWe reject our null hypothesis.')
16     print('Student internship is about our expectations.')
```

For T test:

```
In [ ]: 1 # Test Statistic
2 xbar =
3 popmean =
4 std_error = sigma/np.sqrt(n)
5 # OR
6 std_error = stats.sem(data)
7
8 teststats = (xbar-popmean)/ std_error
9 print(teststats)
10
11 # Remember that in Python, the Manual PValue Calc is One Sided. Thus, we need
12 1 - stats.t.cdf(teststats, n-1)
13
14 # Since P Value > 0.05, then we Fail to Reject the Claim.
```

OR

```

In [ ]: 1 # First we go for sample std.dev if std.dev is not given:
        2 sd = np.std(sample, ddof=1)
        3
        4 std_error = sd/np.sqrt(n)
        5 # OR
        6 std_error = stats.sem(data)
        7
        8 t_stat = (xbar-Mean)/std_error
        9
       10 t_crit = st.t.isf(alpha/2,n-1) # If two tail test.
       11
       12 # If t_stat < t_crit, hence data comes in accept region.
       13
       14 P_value = st.t.sf(t_stat,n-1)*2 # 2 tail test
       15 # OR
       16 1 - stats.t.cdf(teststats, n-1)
       17
       18 # If Pvalue > alpha, hence we accept null hypotheis.
       19
       20 # One line code, rather doing all mathematics:
       21 st.ttest_1samp(sample,Mean)

```

To check the normality of data, whether it is normal or not:

```

In [ ]: 1 st.shapiro(data)
        2
        3 # If Pvalue>Alpha: Which means, data follows normal distribution. Hence we a
        4 # Hence we approach non parametric test.

```

```

In [ ]: 1 # Pvalue>Alpha
        2 # Now test for variance equality:
        3
        4 st.levene(data1,data2)
        5
        6 # H0: Population1 variance = Population2 variance
        7 # H1: Population1 variance != Population2 variance
        8
        9 # pvalue>0.05: accept null value.
       10 # That means variance are same.

```

One Sample T Test - Sample and Pop Mean is Given

```

In [ ]: 1 xbar =
2 std = np.std(x,ddof=1)
3 se = stats.sem(mins)
4 mew =
5 freedom = len(mins) - 1 # Degrees of Freedom: n-1
6 alpha = 0.05
7
8 # Test Statistic and PValue
9 teststats, pvalue = st.ttest_1samp(sample, mean) # Sample Data, Pop Mean Par
10
11 # Conclusion
12
13 if(pvalue>alpha):
14     print("Fail to Reject the Ho")
15 else:
16     print("Reject the Ho meaning that the Mean Value is Different from the H

```

Independent T-Test - Two Samples

```

In [ ]: 1 # If both samples are independent
2 teststats, pvalue = stats.ttest_ind(Weight_Male, Weight_Female)
3
4 # If both samples are related (pairplot)
5 teststats, pval = stats.ttest_rel(Marks_before, Marks_after)
6
7

```

If the data is not given:

```

In [ ]: 1 n1 = ; xbar1 = ; sd1 =
2 n2 = ; xbar2 = ; sd2 =
3 df = n1+n2-2
4
5 t = (xbar1 - xbar2)/np.sqrt((sd1**2/n1)+(sd2**2/n2)) ; print(t)
6
7 # OR
8
9 st.ttest_ind_from_stats(n1,xbar1,sd1,n2,xbar2,sd2)

```

ANOVA

One way Anova

```
In [ ]: 1 # Ho: Tha the Income is Same Across ALL Gyms
        2 # H1: That the Income is not Same Across ALL Gyms
        3
        4 stats.f_oneway(Gym_1, Gym_2, Gym_3)
        5
        6 # Since P Value is > 0.05, We Fail to Reject the Ho.
```

Two way Anova

```
In [ ]: 1 import statsmodels.api as sm
        2 from statsmodels.formula.api import ols
        3
        4 model = ols('M_Income~Gym', data = income).fit() # Ordinary Least Square Met
```

And

```
In [ ]: 1 # First Step: Fit Anova using OLS Method
        2 # Step02: Generate the Anova Table: Sum of Sq, Mean Sq, DF, F Test Statistic
        3
        4 # y~x: y is dependent variable where as ind var is x (3 Levels: H, S, Suv)
        5
        6 model = ols('Mean_Pressure~Car_Type', data = pressure).fit() # Fiting the Mo
        7
        8 print(sm.stats.anova_lm(model))
```

In case of rejecting null hypothesis we recommend to see all the changes happened into the probabilities, we do following.

```
In [ ]: 1 from statsmodels.stats.multicomp import pairwise_tukeyhsd
        2 print(pairwise_tukeyhsd(pressure["Mean_Pressure"], pressure["Car_Type"])) #
```

And

```
In [ ]: 1 import statsmodels.api as sm
        2 from statsmodels.formula.api import ols
        3
        4 mobile_model = ols("Qtty~Discount+Loc+Discount:Loc", data = sales).fit()
        5 print(sm.stats.anova_lm(mobile_model))
```

Chi square:

```
In [ ]: 1 chiteststats = df["Obs-Exp^2"].sum()/Exp_value
        2
        3 crit = stats.chi2.isf(0.01, df) alpha=99, 1-0.99
        4
        5 # If chiteststats > crit, we reject the Ho meaning that the No of Transactio
```

OR

```
In [ ]: 1 # One line code:
        2 stats.chisquare(df["Users"], df["Exp_Values"])
        3
        4 # If p_value < alpha, we reject our null hypothesis.
```

OR

```
In [ ]: 1 data = pd.crosstab(tips.sex, tips.day)
```

```
In [ ]: 1 # Contingency Table to generate Statistic
        2 # Test Stat Value, P Value, DF, Exp Value
        3 teststats, pvalue, df, exp_freq = stats.chi2_contingency(data)
        4
        5 print(teststats)
        6 print(pvalue)
        7 print(df) # Degrees of Freedom = (r-1)*(c-1) = (2-1)*(4-1) = 3
        8 print(exp_freq)
```

If expected value is not given:

```
In [ ]: 1 st.chisquare(df.Column_name)
```

OR

```
In [ ]: 1 df5["Obs"] = np.sum(df5["Obs"])/shape.df5[0]
        2
        3 st.chisquare(df5["Obs"], df5["Exp"]) # df5['Exp'] = Expected column, was cre
```

Machine Learning:

Simple Linear Regreesion:


```
In [ ]: 1 st.linregress(df.a, df.b) # a & b are columns of the dataframe
```

```
In [ ]: 1 from sklearn.linear_model import LinearRegression
2
3 LR = LinearRegression()
4 LR.fit(X, y)
5
6 print(f'Coefficients: {lin_reg.coef_}')
7 print(f'Intercept: {lin_reg.intercept_}')
8 print(f'R^2 score: {lin_reg.score(X, y)}')
```

OR

```
In [ ]: 1 xbar = df.x.mean()
2 ybar = df.u.mean()
3
4 num = np.sum((df.x-xbar)*(df.u-ybar))
5 deno = np.sum((df.x-xbar)**2)
6 slope = num/deno
7
8 print("Slope: ", slope)
9
10 Intercept = ybar - slope * xbar
11
12 print("Intercept: ", Intercept)
```

```
In [ ]: 1 import statsmodels.formula.api as sm
2
3 model = sm.ols('Exp~Income', data=food).fit()
4 model.summary()
```

OR : A different way of making a model.

```
In [ ]: 1 import statsmodels.api as sm
2 from statsmodels.api import OLS
3 life1 = life[['Hepatitis ', 'Diphtheria ', 'Polio', 'Expected']]
4 life11 = life1.iloc[:, [0, 1, 2, 3]]
5 x = life11.drop('Expected', axis=1)
6 x = sm.add_constant(x)
7 y = life11.Expected
8 model = OLS(y, x).fit()
9 model.summary()
```

For R2 value:

```
In [ ]: 1 model.rsquared
```

Root mean square & Mean absolute error:

```
In [ ]: 1 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
2 RMSE = np.sqrt(mean_squared_error(food.Exp, food.Pred_value))
3
4 MAE = mean_absolute_error(food.Exp, food.Pred_value)
```

OR / Complete model.

```
In [ ]: 1 import statsmodels.api as sm
2 from statsmodels.api import OLS
3
4 x = lungs.drop(['FEV', 'Smoke', 'Gender'], axis=1)
5 x = sm.add_constant(x)
6 y = lungs.FEV
7 mod = OLS(y,x).fit()
8
9 from sklearn.metrics import mean_absolute_error, mean_squared_error
10 RMSE = np.sqrt(mean_squared_error(y, mod.predict(x)))
11 MAE = mean_absolute_error(y, mod.predict(x))
12
13 print(RMSE, MAE)
```

OR

```
In [ ]: 1 # To create constant (intercept into the data)
2 x = energy2.drop(['Appliances', 'Press_mm_hg'], axis=1)
3 y = energy2.Appliances
4
5 import statsmodels.api as sm
6 x = sm.add_constant(x)
7
8 from sklearn.model_selection import train_test_split
9 x_train, x_test, y_train, y_test = train_test_split(x,y, test_size=0.20, ran
10
11 model = sm.OLS(y_train,x_train).fit()
12 model.summary()
```

The mean absolute percentage error (MAPE) is commonly used to measure the predictive accuracy of models.

```
In [ ]: 1 def mape(actual, pred):
2         actual, pred = np.array(actual), np.array(pred)
3         return np.mean(np.abs((actual - pred) / actual)) * 100
4 mape(y, mod.predict(x))
```

One line code for Linear Regression:

```
In [ ]: 1 def linear_regression(df,x,y):
2         from sklearn.linear_model import LinearRegression
3         LR = LinearRegression()
4         LR.fit(x,y)
5         pred = LR.predict(x)
6
7         from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_
8         n = df.shape[0]
9         k = df.shape[1]-1
10        factor = (n-1)/(n-k-1)
11
12        RMSE= np.sqrt(mean_squared_error(y,pred))
13        MAE=mean_absolute_error(y,pred)
14        R2=r2_score(y,pred)
15        Ajd_R2 = 1 - ((1-R2) * factor)
16
17        df1 = pd.DataFrame({'Data':['RMSE', 'MAE', 'R2', 'Ajd_R2'], 'Values':[RMSE,
18        return(df1)
19
20 linear_regression(df=premium,x=premium.drop('Premium',axis=1),y=premium.Prem
```

OR

```

In [ ]: 1 def Linear_regression(df,x,y):
2         from sklearn.linear_model import LinearRegression
3         LR = LinearRegression()
4         pred = LR.fit(x,y).predict(x)
5
6         from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_
7         n = df.shape[0]
8         k = df.shape[1]-1
9         factor = (n-1)/(n-k-1)
10
11         RMSE= np.sqrt(mean_squared_error(y,pred)); MAE=mean_absolute_error(y,pre
12         Ajd_R2 = 1 - ((1-R2) * factor)
13         df1 = pd.DataFrame({'Data':['RMSE','MAE','R2','Ajd_R2'], 'Values':[RMSE,
14                                 index=['RMSE','MAE','R2','Ajd_R2']})
15
16         def analysis():
17             global R2
18             if R2>=0.5:
19                 return 'The model is quite good for prediction.'
20             else:
21                 return 'Model is not much appropriate.'
22
23         return(df1[['Values']], analysis())
24
25 Linear_regression(df=premium,x=premium.drop('Premium',axis=1),y=premium.Prem

```

To check the accuracy of the model:

```

In [ ]: 1 from sklearn.model_selection import train_test_split
2         xtrain, xtest, ytrain, ytest = train_test_split(x, y, test_size = 0.2, random_state=42)
3
4         lin_reg.score(xtrain, ytrain) , lin_reg.score(xtest, ytest)

```

Analysis with Random Forest Regression:

```

In [ ]: 1 from sklearn.ensemble import RandomForestRegressor
2         from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error
3         RF = RandomForestRegressor()
4
5         pred1 = RF.fit(x,y).predict(x)
6
7         print(r2_score(y,pred1), np.sqrt(mean_squared_error(y,pred1)), mean_absolute_error(y,pred1))

```

BaggingRegressor method of Model building:

```
In [ ]: 1 from sklearn.ensemble import BaggingRegressor
2 Bagg = BaggingRegressor()
3
4 pred = Bagg.fit(X,y).predict(test_dummies) # test_dummies, is the pd.get_dim
5
6 Prediction = solution.copy() # Solution: DataFrame of test_data, for predict
7
8 Prediction['Item_Outlet_Sales'] = pred
9
10 Prediction.head()
```

Calculations for Fstat and pvalue

```
In [ ]: 1 yp = -11.6905 * premium.Mileage + 327.0860 # Coefficient value and Intercept
2 SSE = np.sum((premium.Premium - yp)**2)
3 SSR = np.sum((yp-premium.Premium.mean())**2)
4 n = premium.shape[0]
5 MSR = SSR/1 ; MSE = SSE / (n-1-1)
6 Fst = MSR/MSE
7
8 1 - st.f.cdf(Fst, 1, n-1-1)
```

OR

```
In [ ]: 1 tstat = -11.6905/6.006 # Coefficient value and std_error
2 df = premium.shape[0]-2
3 pvalue = (1- st.t.cdf(abs(tstat), df))*2
```

Calculations for Confidence Interval

```
In [ ]: 1 df = premium.shape[0]-2
2 alp = 0.025
3
4 tcrit = abs(st.t.ppf(alp,df))
5 b1 = -11.6905 ; se = 6.006
6
7 b1 - tcrit*se , b1 + tcrit*se
```

Assumptions

Residual plot:

```
In [ ]: 1 plt.rcParams['figure.figsize'] = [15,15]
        2
        3 res = model.resid ; fit = model.fittedvalues
        4
        5 sns.residplot(fit, res)
        6 plt.show()
```

Linearity:

```
In [ ]: 1 import statsmodels.api as sm
        2 sm.stats.diagnostic.linear_rainbow(model)
        3 # pvalue > alpha, hence we failed to reject the null hypo. Data is linear.
```

Probabilty plot:

```
In [ ]: 1 plt.rcParams['figure.figsize'] = [15,15]
        2
        3 # Probplot
        4 from scipy.stats import probplot
        5
        6 probplot(model.resid, plot = plt)
        7 plt.show()
        8
        9 # Inferece: Since most of the points are closer to the red line hence, the d
       10 # If the points deviated from the red line, we can conclude that the data is
```

Normality:

```
In [ ]: 1 # Shapiro Wilk Test of Normality: TGT & Residual.
        2
        3 stats.shapiro(model.resid) # Residuals are Normally Distributed.
        4 #stats.shapiro(premium["Premium"]) # Tgt Variable is normally Distributed.
        5 #pvalue>5% then Data is normal.
```

OR

```
In [ ]: 1 tstat, pvalue = st.jarque_bera(model.resid)
        2 #pvalue>5% then Data is normal.
```

Multicollinearity:

```
In [ ]: 1 from statsmodels.stats.outliers_influence import variance_inflation_factor
2
3 x = premium.drop('Premium', axis=1)
4
5 li = []
6
7 for i in range(x.shape[1]):
8     li.append(variance_inflation_factor(x.values, i))
9
10 pd.DataFrame({'columns':x.columns, 'VIF':li}).sort_values(by='VIF')
11
12 # As VIF increases, probability of applying regression gets low.
```

Autocorrelation:

```
In [ ]: 1 from statsmodels.stats.stattools import durbin_watson
2
3 durbin_watson(model.resid)
4
5 # If errors are autocorrelated: Seems feature is significant but it actually
```

Thumb Rule - Test Stats is between 1.5 to 2.5 then we consider it to be normal. Values beyond this range is a sign of Autocorrelation

Inf: Since it is between 1.5 - 2.5, we fail to reject the Ho.

Parameters: DW Test Stat = 2, No Correlation, 0-2: +ve ACorrelation, >2 to 4: -ve ACorrel.

```
In [ ]: 1 import statsmodels.tsa.api as smt
2
3 acf = smt.graphics.plot_acf(LR.resid, lags=31, alpha=0.05)
4 acf.show()
```

Heteroskedasticity:

```
In [ ]: 1 import statsmodels.stats.api as ssa
2
3 predictors = premium.drop("Premium", axis = 1 )
4
5 # residuals, predictors:: endog, exhog
6 # Generate 04 statistic - Test Stat of BP, PValue(BP), Ftest & Pvalue(Ftest)
7 teststats, pvalue, ftest, f_pvalue = ssa.het_breuschpagan(model.resid, predi
8
9 # pvalue > 0.05, Data is not suffering from heteroskedasticity and there is
```

Interaction Effect:

```
In [ ]: 1 premium["Int_Eng_Age"] = premium["Engine"]*premium["Age"]
2
3 finalmodel = ols("Premium~Int_Eng_Age+Engine", data = premium).fit()
4 # OR
5 finalmodel = ols("Premium~Engine:Age+Engine", data = premium).fit()
6
7 finalmodel.summary()
```

Feature Selection:

```
In [ ]: 1 np.abs(finalmodel.params[1:]).sort_values().plot(kind = "barh")
2 # graph will show, which feature is contributing well.
```

OR


```

In [ ]: 1 from sklearn.ensemble import RandomForestRegressor
2 rf = RandomForestRegressor() #Created a Machine
3
4 # Input and Output Variable
5 X = premium.drop(["Mileage", "Premium"], axis = 1)
6 y = premium.Premium
7
8 # Fit and Predict
9
10 pred = rf.fit(X,y).predict(X)
11
12 # To Look the data:
13 pd.DataFrame({"Features":rf.feature_importances_},
14               index = X.columns)
15
16 # IN RF, We already have a parameter called feature importance
17 # Plot will show a good idea of best features.
18 pd.DataFrame({"Features":rf.feature_importances_},
19               index = X.columns).sort_values(by = "Features").plot(kind = "ba

```

```

In [ ]: 1 # Calc the R2 score and RMSE
2 from sklearn.metrics import r2_score, mean_squared_error
3
4 #print("LR R2Score: ", r2_score(y, finalmodel.predict(premium.drop("Mileage"
5 print("RF R2Score: ", r2_score(y, pred))
6
7 #print("LR rmse: ", np.sqrt(mean_squared_error(y, finalmodel.predict(X))))
8 print("RF rmse: ", np.sqrt(mean_squared_error(y, pred)))

```

Recursive Feature Elimination:

```

In [ ]: 1 from sklearn.ensemble import RandomForestRegressor
2 rf = RandomForestRegressor()
3
4 from sklearn.feature_selection import RFE
5 rfe = RFE(estimator=RandomForestRegressor(),
6           n_features_to_select=2) # estimator-Model, n_features_to_select
7
8 newinput = rfe.fit_transform(X,y) # return the dataset with the correct set
9
10 # Finding the Column Names
11 colnos = rfe.get_support(indices = True) # returns the indices or the column
12
13 X.iloc[:, colnos].columns # this is the correct approach

```

OR

```
In [ ]: 1 from sklearn.feature_selection import RFE
2 from sklearn.linear_model import LinearRegression
3 lr = LinearRegression()
4
5 pred_lr = lr.fit(X,y).predict(X)
6
7 rfe = RFE(estimator = LinearRegression(),
8           n_features_to_select=2)
9
10 newinput = rfe.fit(X,y)
11 colnos = rfe.get_support(indices = True)
12 print(X.iloc[:, colnos].columns)
```

OR

Sequential Feature Selection:

```
In [ ]: 1 #pip install mlxtend - Sequential Feature Selector
2 from mlxtend.feature_selection import SequentialFeatureSelector as sfs
3
4 from sklearn.linear_model import LinearRegression
5 lr = LinearRegression()
6
7 fs = sfs(estimator=LinearRegression(), k_features="best", forward = True,
8         verbose = 2, scoring = "r2", cv = 14)
9
10 # fit the model and generate the feature names
11 X = premium.drop(["Mileage", "Premium"], axis = 1)
12 y = premium.Premium
13
14 sfsmodel = fs.fit(X, y)
15
16 # Generate the feature names
17 print("Features: ", sfsmodel.k_feature_names_)
18 print("R2 Score: ", sfsmodel.k_score_)
```

To get more precise model we do boosting:

```
In [ ]: 1 from sklearn.ensemble import GradientBoostingRegressor
2 gbm = GradientBoostingRegressor()
3
4 from mlxtend.feature_selection import SequentialFeatureSelector as sfs
5
6 from sklearn.linear_model import LinearRegression
7 lr = LinearRegression()
8
9 fs = sfs(estimator=gbm, k_features="best", forward = True,
10         verbose = 2, scoring = "r2", cv = 14)
11
12 # fit the model and generate the feature names
13 X = premium.drop(["Mileage", "Premium"], axis = 1)
14 y = premium.Premium
15
16 sfsmodel = fs.fit(X, y)
17
18 # Generate the feature names
19 print("Features: ", sfsmodel.k_feature_names_)
20 print("R2 Score: ", sfsmodel.k_score_)
```

Feature selector Estimator

```
In [ ]: 1 SFS = sfs(estimator = lr, verbose=0, k_features = "best",
2             cv = 5, scoring = "r2")
3
4 sfsmodel = SFS.fit(x_train,y_train)
5 print("Features: ", sfsmodel.k_feature_names_)
6 print("R Squared: ", sfsmodel.k_score_)
```

Cross Validation:

K Fold:

```
In [ ]: 1 from sklearn.model_selection import KFold
2 n=10
3 KF = KFold(n_splits=n, shuffle=True)
4
5 x = energy2.drop(["Appliances", "lights"], axis = 1)
6 y = energy2.Appliances
7
8 for i in range(n):
9     result = next(KF.split(x))
10    x_train = x.iloc[result[0]]
11    x_test = x.iloc[result[1]]
12    y_train = y.iloc[result[0]]
13    y_test = y.iloc[result[1]]
14
15 x_train.shape, x_test.shape, y_train.shape, y_test.shape
```

One line code for cross validation score:

```
In [ ]: 1 from sklearn.model_selection import cross_val_score
2
3 from sklearn.model_selection import KFold
4 n =5
5
6 # shuffle the data after every draw
7 # n splits - how many splits
8
9 #kf_ = KFold(n_splits=n, shuffle = True)
10 cross_val_score(estimator=lrfit, X = x_train, y = y_train, scoring = "r2", c
```

OR

```
In [ ]: 1 from sklearn.model_selection import cross_val_score
2
3 from sklearn.model_selection import KFold
4 n =5
5
6 x = df.drop('quality', axis=1)
7 y = df.quality
8
9 lrfit = LinearRegression()
10
11 cross_val_score(estimator=lrfit, X=x, y= y, scoring = "r2", cv = 2).mean()
```

To extract numerical columns only

```
In [ ]: 1 df.select_dtypes(np.numeric)
```

Overfitting and Underfitting

```
In [ ]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.metrics import mean_squared_error
3
4 x = fish.drop(['Weight', 'Species'], axis = 1)
5 y = fish['Weight']
6 xtrain, xtest, ytrain, ytest = train_test_split(x, y, test_size = 0.2, random_state=1)
7
8 from sklearn.linear_model import LinearRegression
9 lrmodel = LinearRegression()
10
11 pred_train = lrmodel.fit(xtrain, ytrain).predict(xtrain)
12 pred_test = lrmodel.fit(xtrain, ytrain).predict(xtest)
13
14 print("RMSE for Train: ", np.sqrt(mean_squared_error(ytrain, pred_train)))
15 print("RMSE for Test: ", np.sqrt(mean_squared_error(ytest, pred_test)))
```

Elastic regression:

```
In [ ]: 1 from numpy import mean
2 from numpy import std
3 from numpy import absolute
4 from pandas import read_csv
5 from sklearn.model_selection import cross_val_score
6 from sklearn.model_selection import RepeatedKFold
7 from sklearn.linear_model import ElasticNet
8
9 model = ElasticNet(alpha=1.0, l1_ratio=0.001)
10 # define model evaluation method
11 cv = RepeatedKFold(n_splits=10, n_repeats=3, random_state=1)
12 # evaluate model
13 scores = cross_val_score(model, x, y, scoring='neg_mean_absolute_error', cv=cv)
14 # force scores to be positive
15 scores = absolute(scores)
16 print('Mean MAE: %.3f (%.3f)' % (mean(scores), std(scores)))
```

GridSearch

```
In [ ]: 1 from sklearn.model_selection import GridSearchCV
2 alphas = [{'alpha': [0.1, 0.01, 0.02, 0.0005, 0.09, 0.001, 0.5, 0.99, 1, 2, 5]}]
3 grid = GridSearchCV(estimator = Ridge(normalize=True), param_grid= alphas, cv=5)
4
5 grid.fit(xtrain, ytrain)
6
7 grid.best_params_
```

OR

```
In [ ]: 1 params = {'alpha' : [0.1,0.5,1,1.5,2,2.5,3,3.5,4] , 'l1_ratio' : [0.5,0.55,0.56,0.57,0.58,0.59,0.6,0.61,0.62,0.63,0.64,0.65,0.66,0.67,0.68,0.69,0.7,0.71,0.72,0.73,0.74,0.75,0.76,0.77,0.78,0.79,0.8,0.81,0.82,0.83,0.84,0.85,0.86,0.87,0.88,0.89,0.9,0.91,0.92,0.93,0.94,0.95,0.96,0.97,0.98,0.99,1] }
2
3 from sklearn.model_selection import GridSearchCV
4
5 grid = GridSearchCV(ene, param_grid=params,cv=10)
6 grid.fit(xtrain,ytrain)
7
8 grid.best_params_
```

Lasso & Regression

```
In [ ]: 1 from sklearn.model_selection import GridSearchCV
2 from sklearn.linear_model import LassoCV, RidgeCV, ElasticNet
```

BaggingRegressor

```
In [ ]: 1 from sklearn.ensemble import BaggingRegressor
2 Bagg = BaggingRegressor()
3
4 pred = Bagg.fit(X,y).predict(test_dummies)
5
6 Prediction = solution.copy()
7
8 Prediction['Item_Outlet_Sales'] = pred
9
10 Prediction.head()
```

```
In [ ]: 1 Prediction.to_csv('Predictions.csv', index=False)
```

Awesome...

END

