# Part 2 Math grades for both schools

```
In [2]:  import pandas as pd                    #Importing the required package
         import numpy as np
         import matplotlib.pyplot as plt
         import seaborn as sns
         from sklearn.preprocessing import MinMaxScaler
         from sklearn.model_selection import KFold, cross_val_predict
         from sklearn.pipeline import Pipeline
         from sklearn.metrics import mean_squared_error
         from sklearn.metrics import r2_score
         from sklearn.model_selection import GridSearchCV
         from timeit import default_timer as timer
         from sklearn.linear_model import LinearRegression
         import warnings
         warnings.simplefilter(action='ignore', category=FutureWarning)
         warnings.simplefilter(action='ignore', category=DeprecationWarning)
         warnings.simplefilter(action='ignore', category=UserWarning)
```

```
In [ ]:  data = pd.read_csv('student-mat.csv',sep=';')
         data.head()
```

```
In [ ]:  data = data.drop(['G1','G2'],axis=1) #dropping columns G1 and G2 in order to get a m
         data.head()
```

```
In [ ]:  print("Number of rows:", len(data)) #Looking for number of rows in total
```

```
In [ ]:  data.isna().sum() #Checking for Null Values
```

```
In [ ]:  data.info()               #checking for Data types of each variables
```

```
In [ ]:  for var in ['traveltime','studytime','Medu','Fedu','famrel','freetime','goout','Dalc
             data[var] = data[var].astype('category')
         data.info()
```

```
In [ ]:  num_col = data._get_numeric_data().columns          #Getting the numerical dataset fr
         cat_col = list(set(data.columns)-set(num_col))       #By removing the numerical datase
         num_col = list(num_col)[0:3]
```

```
In [ ]:  nominal_col = ['higher','Fjob','address','guardian','school','paid','sex','Mjob','ac
         ordinal_col = np.setdiff1d(cat_col,nominal_col)
```

```
In [ ]:  for col in ordinal_col:
             data[col] = data[col].astype(int)               # Changing Ordinal data to int
```

```
In [ ]:  for col in data[cat_col]:
             print(col,data[col].unique())                   # Checking for erroneous data in
```

In [ ]:
```python
data[num_col].describe()
```

In [ ]:
```python
for i in data[num_col]:
    print(i,data[i].unique())
```

In [ ]:
```python
data['G3'].describe()
```

In [ ]:
```python
data['G3'].unique()
```

# Exploratory Data Analytics

In [ ]:
```python
sns.pairplot(data, hue = 'school')
```

In [ ]:
```python
corrs = data.corr()
fig, ax = plt.subplots(figsize=(15,10))
sns.heatmap(corrs, annot=True, fmt='.2f', ax=ax)
plt.show()
```

**Based on observation, it is shown the following:**

Fedu and Medu has moderate correlation

Walc and Dalc has moderate correlation

Walc and goout has weak correlation

In [ ]:
```python
sns.set(rc={'figure.figsize':(15,5)})
sns.swarmplot(data=data, x='school', y='age')
```

**Based on observation, it is shown that:**

Students in GP are relatively younger than MS

In [ ]:
```python
ax = sns.boxplot(data=data, x='school', y='absences')
ax.set_ylim([0, 35])
```

**Based on observation, it is shown that:**

GP has more absentees comapred to MS

In [ ]:
```python
data.groupby(by='school').mean()
```

In [ ]:
```python
sns.boxplot(data=data, x="school", y="G3",
            palette=["m", "g"]
            )
sns.set(rc={'figure.figsize':(8,5)})
```

**Based on Observation on the boxplot, it is shown that:**

In terms of median, 75th and 25th percentile, students in GP are relatively doing better than students in MS

# Reducing Skewness in numerical data

In [ ]:
```python
skew_limit = 0.75

skew_cols = (data[num_col].skew()
            .sort_values(ascending=False)
            .to_frame()
            .rename(columns={0:'Skew'})
            .query('abs(Skew) > {}'.format(skew_limit)))

skew_cols
```

In [ ]:
```python
field = "absences"
fig, (ax_original, ax_log1p) = plt.subplots(1, 2, figsize=(15, 5))

data[field].hist(ax=ax_original)

# Apply a log transformation (numpy syntax) to this column
data[field].apply(np.log1p).hist(ax=ax_log1p)

# Formatting of titles etc. for each subplot
ax_original.set(title='before np.log1p', ylabel='frequency', xlabel='value')
ax_log1p.set(title='after np.log1p', ylabel='frequency', xlabel='value')
fig.suptitle('Field "{}"'.format(field));
print('pop_orignial skewness: ',data[field].skew())
print('pop_log1p skewness: ',data[field].apply(np.log1p).skew())
# fall in skewness after log1p
```

In [ ]:
```python
# apply log1p across all numerical columns

for col in data[skew_cols.index]:
# drop famrel as skewness increase after log1p
    data[col] = data[col].apply(np.log1p)
new_skew_cols = (data[skew_cols.index].skew()
            .sort_values(ascending=False)
            .to_frame()
            .rename(columns={0:'Skew'}))

skew_cols['New_skew'] = new_skew_cols
skew_cols['Difference'] = abs(skew_cols['New_skew']) - abs(skew_cols['Skew'])
skew_cols
```

# MiniMax scale numerical data

In [ ]:
```python
MinM = MinMaxScaler()
for col in num_col:
    data[col] = MinM.fit_transform(data[[col]])
```

# One hot encoding of categorical Data

In [ ]:
```python
data = pd.get_dummies(data, columns=nominal_col,drop_first=True)
data.head()
# df is the transformed dataest for ML, df is the orignal dataset
```

# Polynomial Features

```python
In [ ]:    from sklearn.preprocessing import PolynomialFeatures
```

```python
In [ ]:    feature_cols = list(filter(lambda x: x!= 'G3', data.columns))

           polyf = PolynomialFeatures(degree=2, include_bias=False,)
           X_pf = polyf.fit_transform(data[feature_cols])
```

# Train Test Split

```python
In [ ]:    from sklearn.model_selection import train_test_split

           train,test = train_test_split(data,
                                         train_size = 0.7,
                                         test_size = 0.3,
                                         random_state = 42)

           feature_cols = list(filter(lambda x: x!= 'G3', data.columns))

           X_train = train[feature_cols]
           y_train = train['G3']

           X_test = test[feature_cols]
           y_test = test['G3']
```

```python
In [ ]:    # splitting of Polynomial feautures
           X_pf_train = X_pf[X_train.index]
           y_pf_train = data['G3'][X_train.index]

           X_pf_test = X_pf[X_test.index]
           y_pf_test = data['G3'][X_test.index]
```

# Model Kfold + Evaluation

```python
In [ ]:    # define root mean sq error func
           def rmse(ytrue, ypredicted):
               return np.sqrt(mean_squared_error(ytrue, ypredicted))

           folds = KFold(n_splits = 5, shuffle = True, random_state = 100)
```

```python
In [ ]:    import statsmodels.api as sm
           def diagnostic_plot(y_test, prediction):
               residual = y_test - prediction

               fig, axs = plt.subplots(2, 2,figsize=(15,10))
               fig.suptitle('Diagnostic plots')
               axs[0,0].plot([0, 15], [0, 20], ls="--", c="red", alpha=0.5)
               axs[0,0].scatter(y_test, prediction)
               axs[0,0].set_title('Truth-Prediction plot')
               axs[0,0].set_xlabel("Truth")
               axs[0,0].set_ylabel("Predictions")
```

```python
        sm.qqplot(residual,fit=True,line= 's', ax=axs[0,1])
        axs[0,1].set_title('QQ plot')
        axs[0,1].set_xlabel("Theoretical Quantiles")
        axs[0,1].set_ylabel("Sample Quantiles")

        sns.residplot(prediction,y_test,
                              lowess=True,
                              scatter_kws={'alpha': 0.5},
                              line_kws={'color': 'red', 'lw': 1, 'alpha': 0.8},
                              ax=axs[1,0])
        axs[1,0].set_title('Residual plot')
        axs[1,0].set_xlabel("Predicted")
        axs[1,0].set_ylabel("Residual")

        sqrt_standardized_residual=np.sqrt(np.abs(residual))
        sns.regplot(prediction, sqrt_standardized_residual,
                        scatter=True,
                        lowess=True,
                        line_kws={'color': 'red', 'lw': 1, 'alpha': 0.5},
                        ax=axs[1,1])
        axs[1,1].set_title('Scale-Location plot')
        axs[1,1].set_xlabel("Predicted")
        axs[1,1].set_ylabel("Sqrt Standarized residuals")
```

# Dummy Model

In [ ]:
```python
average_G3 = np.mean(y_test)
average_G3
```

In [ ]:
```python
start = timer()

dummy_prediction =[]
for row in range(len(y_test)):
    dummy_prediction.append(average_G3)

print("R2 score: ",r2_score(y_test, dummy_prediction))
print('RMSE: ', rmse(y_test,dummy_prediction))

end = timer()
dummy_time = end - start
print('Time Elapsed (sec):', dummy_time)
```

# Base Linear Regression

In [ ]:
```python
linear = LinearRegression()
linear_cv = GridSearchCV(estimator = linear,
                         param_grid = {},
                         cv = folds)
start = timer()
linear_cv.fit(X_train, y_train)

end = timer()
linear_time = end - start
print('Time Elapsed (sec):', linear_time)
```

```python
linear_prediction = linear_cv.predict(X_test)

print("R2 score: ",r2_score(y_test, linear_prediction))
print('RMSE: ', rmse(y_test,linear_prediction))
```

```python
diagnostic_plot(y_test, linear_prediction)
```

```python
print("\u0332".join('Feature_coefficient (Top)'),'\n')
features = []
coeffs = []
coeffs_abs = []
for feature, coeff in zip(X_test.columns, linear_cv.best_estimator_.coef_):
    features.append(feature)
    coeffs.append(coeff)
    coeffs_abs.append(abs(coeff))

pd.DataFrame({'Features': features, 'Coefficients': coeffs ,'abs_coeffs':coeffs_abs}
```

```python
print("\u0332".join('Feature_coefficient (Bottom)'),'\n')
pd.DataFrame({'Features': features, 'Coefficients': coeffs ,'abs_coeffs':coeffs_abs}
```

```python
lm = LinearRegression()
lm_pf_cv = GridSearchCV(estimator = lm,
                        param_grid = {},
                        cv = folds)
start = timer()
lm_pf_cv.fit(X_pf_train, y_pf_train)

end = timer()
lm_pf_time = end - start
print('Time Elapsed (sec):', lm_pf_time)
```

```python
lm_pf_prediction = lm_pf_cv.predict(X_pf_test)

print("R2 score: ",r2_score(y_pf_test, lm_pf_prediction))
print('RMSE: ', rmse(y_pf_test,lm_pf_prediction))
```

```python
diagnostic_plot(y_test, lm_pf_prediction)
```

## Lasso Regression

```python
alphas = np.geomspace(0.001, 10, 30)
```

```python
from sklearn.linear_model import LassoCV

start = timer()
lasso_estimator = LassoCV(alphas=alphas,
                          cv=5).fit(X_train,y_train)
lasso_prediction = lasso_estimator.predict(X_test)

print('Alpha param: ', lasso_estimator.alpha_)
print("R2 score: ",r2_score(y_test, lasso_prediction))
```

```python
print('RMSE: ', rmse(y_test,lasso_prediction))

end = timer()
lasso_time = end - start
print('Time Elapsed (sec):', lasso_time)
```

```python
print("Number of non-zero coeff: ", sum(lasso_estimator.coef_ != 0))
print("Mean coeff: ",sum(abs(lasso_estimator.coef_))/sum(lasso_estimator.coef_ != 0)
```

```python
print("\u0332".join('Feature_coefficient (Top)'),'\n')
features = []
coeffs = []
coeffs_abs = []
for feature, coeff in zip(X_test.columns, lasso_estimator.coef_):
    features.append(feature)
    coeffs.append(coeff)
    coeffs_abs.append(abs(coeff))

pd.DataFrame({'Features': features, 'Coefficients': coeffs ,'abs_coeffs':coeffs_abs}
```

```python
diagnostic_plot(y_test, lasso_prediction)
```

## Ridge Regression

```python
alphas = np.geomspace(1, 100, 30)
```

```python
from sklearn.linear_model import RidgeCV

start = timer()
ridge_estimator = RidgeCV(alphas=alphas,
                          cv=5).fit(X_train,y_train)
ridge_prediction = ridge_estimator.predict(X_test)

print('Alpha param: ', ridge_estimator.alpha_)
print("R2 score: ",r2_score(y_test, ridge_prediction))
print('RMSE: ', rmse(y_test,ridge_prediction))

end = timer()
ridge_time = end - start
print('Time Elapsed (sec):', ridge_time)
```

```python
print("Number of non-zero coeff: ", sum(ridge_estimator.coef_ != 0))
print("Mean coeff: ",sum(abs(ridge_estimator.coef_))/sum(ridge_estimator.coef_ != 0)
```

```python
print("\u0332".join('Feature_coefficient (Top)'),'\n')
features = []
coeffs = []
coeffs_abs = []
for feature, coeff in zip(X_test.columns, ridge_estimator.coef_):
    features.append(feature)
    coeffs.append(coeff)
    coeffs_abs.append(abs(coeff))

pd.DataFrame({'Features': features, 'Coefficients': coeffs ,'abs_coeffs':coeffs_abs}
```

In [ ]:
```python
diagnostic_plot(y_test, ridge_prediction)
```

# Elasticnet Regression

In [ ]:
```python
from sklearn.linear_model import ElasticNetCV
alphas = np.geomspace(0.01, 1, 30)
l1_ratios = np.linspace(0.1, 0.9, 9)

start = timer()
elasticNetCV = ElasticNetCV(alphas=alphas,
                            l1_ratio=l1_ratios,
                            cv=5).fit(X_train,y_train)
elasticNet_prediction = elasticNetCV.predict(X_test)

print('Alpha param: ', elasticNetCV.alpha_)
print('l1 ratio param: ', elasticNetCV.l1_ratio_)
print("R2 score: ",r2_score(y_test, elasticNet_prediction))
print('RMSE: ', rmse(y_test,elasticNet_prediction))

end = timer()
EN_time = end - start
print('Time Elapsed (sec):', EN_time)
```

In [ ]:
```python
print("Number of non-zero coeff: ", sum(elasticNetCV.coef_ != 0))
print("Mean coeff: ",sum(abs(elasticNetCV.coef_))/sum(elasticNetCV.coef_ != 0))
```

In [ ]:
```python
print("\u0332".join('Feature_coefficient (Top)'),'\n')
features = []
coeffs = []
coeffs_abs = []
for feature, coeff in zip(X_test.columns, elasticNetCV.coef_):
    features.append(feature)
    coeffs.append(coeff)
    coeffs_abs.append(abs(coeff))

pd.DataFrame({'Features': features, 'Coefficients': coeffs ,'abs_coeffs':coeffs_abs}
```

In [ ]:
```python
diagnostic_plot(y_test, elasticNet_prediction)
```

# Stepwise regression - Forward

In [ ]:
```python
# Starting with no features, features are added in iteratively based on the p-value
def forward_regression(X, y,
                       threshold_in,
                       verbose=False):
    initial_list = []
    included = list(initial_list)
    while True:
        changed=False
        exc = list(set(X.columns)-set(included))
        pval = pd.Series(index=exc)
        for column in exc:
            model = sm.OLS(y, sm.add_constant(pd.DataFrame(X[included+[column]]))).f
```

```python
                pval[column] = model.pvalues[column]
            best = pval.min()
            if best < threshold_in:
                best_feature = pval.idxmin()
                included.append(best_feature)
                changed=True
                if verbose:
                    print('Add  {:17} with p-value {:.6}'.format(best_feature, best))

            if not changed:
                break

    return included


# Starting with all features, features are removed iteratively based on the p-value
def backward_regression(X, y,
                            threshold_out,
                            verbose=False):
    included=list(X.columns)
    while True:
        changed=False
        model = sm.OLS(y, sm.add_constant(pd.DataFrame(X[included]))).fit()
        # use all coefs except intercept
        pvalues = model.pvalues.iloc[1:]
        worst_pval = pvalues.max() # null if pvalues is empty
        if worst_pval > threshold_out:
            changed=True
            worst_feature = pvalues.idxmax()
            included.remove(worst_feature)
            if verbose:
                print('Drop {:17} with p-value {:.6}'.format(worst_feature, worst_pv
        if not changed:
            break
    return included
```

In [ ]:
```python
start = timer()
forward_regression_features = forward_regression(X_train,y_train, threshold_in=0.05,
forward_regression_features
```

In [ ]:
```python
lm = LinearRegression()
forward = GridSearchCV(estimator = lm,
                        param_grid = {},
                        cv = folds)

# fit the model
forward.fit(X_train[forward_regression_features], y_train)

end = timer()
forward_time = end - start
print('Time Elapsed (sec):', forward_time)
```

In [ ]:
```python
forward_prediction = forward.predict(X_test[forward_regression_features])

print("R2 score: ",r2_score(y_test,forward_prediction))
print('RMSE: ', rmse(y_test,forward_prediction))
```

In [ ]:
```python
diagnostic_plot(y_test, forward_prediction)
```

# Stepwise regression - Backward

In [ ]:
```python
start = timer()
backward_regression_features = backward_regression(X_train,y_train, threshold_out=0.
```

In [ ]:
```python
lm = LinearRegression()
backward = GridSearchCV(estimator = lm,
                        param_grid = {},
                        cv = folds)

# fit the model
backward.fit(X_train[backward_regression_features], y_train)

end = timer()
backward_time = end - start
print('Time Elapsed (sec):', backward_time)
```

In [ ]:
```python
backwards_prediction = backward.predict(X_test[backward_regression_features])

print("R2 score: ",r2_score(y_test, backwards_prediction))
print('RMSE: ', rmse(y_test,backwards_prediction))
```

In [ ]:
```python
diagnostic_plot(y_test, backwards_prediction)
```

# Random Forest Regressor

In [ ]:
```python
from sklearn.ensemble import RandomForestRegressor
random_tree = RandomForestRegressor(criterion = 'mse')
```

In [ ]:
```python
n_features = len(data.columns)
max_features = ['auto', 'sqrt', 'log2']
max_depth = [3, 4]
min_samples_split = [2, 5, 10]
min_samples_leaf = [1, 2, 4]
param_test = {'max_features': max_features,
 'max_depth': max_depth,
 'min_samples_split': min_samples_split,
 'min_samples_leaf': min_samples_leaf}
print(param_test)
```

In [ ]:
```python
random_tree = GridSearchCV(estimator = random_tree,
                           param_grid = param_test,
                           cv = folds)
start = timer()
random_tree.fit(X_train, y_train)
end = timer()
rf_time = end - start
print('Time Elapsed (sec):', rf_time)
```

In [ ]:
```python
random_tree_prediction = random_tree.predict(X_test)
print('Best param: ', random_tree.best_params_)
```

```python
print("R2 score: ",r2_score(y_test, random_tree_prediction))
print('RMSE: ', rmse(y_test,random_tree_prediction))
```

In [ ]:
```python
diagnostic_plot(y_test, random_tree_prediction)
```

In [ ]:
```python
from sklearn.tree import export_graphviz
tree = random_tree.best_estimator_.estimators_[5]
export_graphviz(tree, out_file='small.dot',
                feature_names = X_train.columns)
```

In [ ]:
```python
from IPython.display import Image
Image(filename = 'small.png')
```

In [ ]:
```python
feat_df = pd.DataFrame({'Feature':X_train.columns, 'Importance':random_tree.best_est
feat_df = feat_df.sort_values(by = 'Importance', ascending=False)
feat_df.head()
```

In [ ]:
```python
g = sns.barplot(data=feat_df, x='Feature', y= 'Importance')
for item in g.get_xticklabels():
    item.set_rotation(90)
```

In [ ]:
```python
models_pred = [dummy_prediction, linear_prediction, lm_pf_prediction, lasso_predicti
time_elapsed = [dummy_time, linear_time, lm_pf_time, lasso_time, ridge_time, EN_time
rmse_vals = []
for pred in models_pred:
    rmse_vals.append(rmse(y_test, pred))

R2_score = []
for pred in models_pred:
    R2_score.append(r2_score(y_test, pred))

labels = ['Dummy', 'Linear', 'Linear + PF', 'Ridge', 'Lasso', 'ElasticNet','Stepwise

eval_df = pd.DataFrame({'RMSE':rmse_vals, 'R2 Score':R2_score, 'Time Elapsed':time_e
eval_df
```

In [ ]: