

Practical#2

Name: Saurin Anilkumar Prajapati

Roll No.: 19BCE239

Course Code and Name: 2CS702 Big Data Analytics

Batch: D1

AIM:

Learning limitation of data analytics by applying Machine Learning Techniques on large amount of data. Write a program to read data set from any online website, excel file and CSV file and to perform

- a) Linear regression and logistic regression on iris dataset.
- b) K-means clustering.

□ Students will learn the limitation of platform and algorithm

Data:

Tesla stock pricing dataset on Kaggle

(for small dataset → in KBs)

(for mid sized dataset → in MBs)

(for large dataset → 1.7 GB)

Algorithms Used:

Lazy regressor and XGBoost

```
#Ready for model selection!

if REGRESSOR == True:
    try:
        reg = LazyRegressor(verbose=2, ignore_warnings=False, custom_metric=None)
        models, predictions = reg.fit(X_train, X_test, y_train, y_test)
        print(f'Project: {PROJECT_NAME}')
        print(PROJECT_NAME)
        print(f'Target: {target}')
```

```

        print(target)
        target_std = y_train.std()
        print(f'Target Standard Deviation: {target_std}')
        print(models)
        models['project'] = PROJECT_NAME
        models['target'] = target
        models['target_std'] = target_std
        #rename index of
        models.to_csv(f'{PARAM_DIR}/regression_results_{target_str}.csv', mode='a', header=True, index=True)
    except:
        print('Issue during lazypredict analysis')
    else:
        #TODO: remove this
        try:
            clf = LazyClassifier(verbose=2, ignore_warnings=False, custom_metric=None)
            models, predictions = clf.fit(X_train, X_test, y_train, y_test)
            print(f'Project: {PROJECT_NAME}')
            print(PROJECT_NAME)
            print(f'Target: {target}')
            print(target)
            print(f'Target Standard Deviation: {y_train.std()}')
            print(models)
            models.to_csv(f'{PARAM_DIR}/classification_results.csv', mode='a', header=False)
        except:
            print('Issue during lazypredict analysis')

model_name = 'tabnet'

# FastAI + pre-trained TabNet
#=====
learn = None
i = 0
while True:
    try:
        del learn
    except:
        pass
    try:
        learn = 0
        model = TabNetModel(get_emb_sz(to), len(to.cont_names), dls.c, n_d=64, n_a=64,
n_steps=5, virtual_batch_size=256)
        # save the best model so far, determined by early stopping
        cbs = [SaveModelCallback(monitor='_rmse', comp=np.less, fname=f'{model_name}_{PROJECT_NAME}_{target_str}_best'), EarlyStoppingCallback()]
        learn = Learner(dls, model, loss_func=MSELossFlat(), metrics=rmse, cbs=cbs)
        #learn = get_learner(to)
        if(learn != 0):
            break
        if i > 50:
            break
    except:
        i += 1

```

```

        print('Error in FastAI TabNet')
        traceback.print_exc()
        continue
    try:
        #display learning rate finder results
        x = learn.lr_find()
    except:
        pass
    if AUTO_ADJUST_LEARNING_RATE == True:
        FASTAI_LEARNING_RATE = x.valley
    print(f'LEARNING RATE: {FASTAI_LEARNING_RATE}')
    try:
        if i < 50:
            learn.fit_one_cycle(20, FASTAI_LEARNING_RATE)
            plt.figure(figsize=(10, 10))
            try:
                ax = learn.show_results()
                plt.show(block=True)
            except:
                print('Could not show results')
                pass
    except:
        print('Could not fit model')
        traceback.print_exc()
        pass

=====

#fit an xgboost model
=====
if REGRESSOR == True:
    xgb = XGBRegressor()
else:
    xgb = XGBClassifier()
try:
    xgb = XGBRegressor()
    xgb.fit(X_train, y_train)
    y_pred = xgb.predict(X_test)
    print('XGBoost Predictions vs Actual=====')
    print(pd.DataFrame({'actual': y_test, 'predicted': y_pred}).head())
    print('XGBoost RMSE: ', np.sqrt(mean_squared_error(y_test, y_pred)))
    #save feature importance plot to file
    plot_importance(xgb)
    plt.title(f'XGBoost Feature Importance for {PROJECT_NAME} | Target : {target}', wr
ap=True)
    plt.tight_layout()
    plt.show()
    plt.savefig(f'{PARAM_DIR}/xgb_feature_importance_{target_str}.png')
    fi_df = pd.DataFrame([xgb.get_booster().get_score()]).T
    fi_df.columns = ['importance']
    #create a column based off the index called feature
    fi_df['feature'] = fi_df.index
    #create a dataframe of feature importance
    fi_df = fi_df[['feature', 'importance']]

```

```

fi_df.to_csv(f'{PARAM_DIR}/xgb_feature_importance_{target_str}.csv', index=False)
#xgb_fi = pd.DataFrame(xgb.feature_importances_, index=X_train.columns, columns=
['importance'])
#xgb_fi.to_csv(f'{PARAM_DIR}/xgb_feature_importance_{target_str}.csv')
#print('XGBoost AUC: ', roc_auc_score(y_test, y_pred))
except:
    traceback.print_exc()
    print('XGBoost failed')

```

On Kaggle:

Small data	367ms
Mid data	34283ms
Large data	700000+ms(encountered error - session ended)

Thank You.