

# Assignment

- The training job you used, we started based on a pre-trained model with Abalone data
- Did that really make a difference? What if you use the data I gave you and start a new training from scratch and use LightGBM algorithm to train a new model?  
(<https://docs.aws.amazon.com/sagemaker/latest/dg/lightgbm.html>)
- You need to write code (not console) in the notebook to do this assignment
- Use the same hyperparameters you used during the class in JS to train from scratch and compare at least the following items:
  - Training time
  - Quality of model against validation data
  - Residual values
- Write a report to explain the steps you took and what you learned out of this activity
- Submit the completed notebook and report to blackboard. The notebook should have all the cells run already and I want to see the result of each cell
- You just open the file in the presentation day and explain what happens in each cell. You do not run those cells in the presentation day. We just talk about results and reports.

## Comparison

After using the same hyper parameters and re-running a training from scratch by using LightGBM algorithm, it was seen that the both pretrained model and the model trained on notebook takes approximately same time. The pretrained model takes equivalent to 2 minutes as shown below –

The screenshot displays the Amazon SageMaker Studio interface. The left sidebar contains navigation options: Home, Data, Data Wrangler, Feature Store, Clusters, AutoML, Experiments, Notebook jobs, Pipelines, Models, Deployments, and SageMaker JumpStart. The main panel shows the 'TRAINING JOB' for 'smjs-c-lgb-regression-model-20230403-210228'. The training status is 'Complete', with a message stating: 'The training job that fitted the model to your data is complete. It created a new model and uploaded it to Amazon S3. From here you can see information about the model and deploy the model to an endpoint. You can also access this model from the AWS SageMaker console.' The training time is listed as '~2 minutes'. The output path is 's3://sagemaker-ca-central-1-224670572127/smjs-c-lgb-regression-model-20230403-210228/output/model.tar.gz'. The instance settings show 'Instance type: m5.xlarge' and 'Number of instances: 1'. The bottom status bar indicates 'Simple' mode, '0' errors, '6' warnings, and the time '6:54 PM 2023-04-03'.

The screenshot shows the Amazon SageMaker Studio interface. The left sidebar contains navigation options: Home, Data, Data Wrangler, Feature Store, Clusters, AutoML, Experiments, Notebook jobs, Pipelines, Models, Deployments, and SageMaker JumpStart. The main panel displays the 'TRAINING JOB' details for 'smjs-c-lgb-regression-model-20230403-210228'. The 'Instance Settings' tab is active, showing a list of hyperparameters:

| Hyper Parameter         | Value |
|-------------------------|-------|
| Bagging_fraction        | 0.53  |
| Bagging_freq            | 5     |
| Boosting                | gbdt  |
| Early_stopping_rounds   | 30    |
| Feature_fraction        | 0.74  |
| Feature_fraction_bynode | 1     |
| Is_unbalance            | False |
| Lambda_11               | 0     |
| Lambda_12               | 0     |
| Learning_rate           | 0.009 |
| Max_bin                 | 255   |
| Max_delta_step          | 0     |
| Max_depth               | 11    |
| Evaluation metric       | auto  |
| Min_data_in_leaf        | 26    |
| Min_gain_to_split       | 0     |

The bottom status bar shows 'Simple' mode, 0/6/6 metrics, and a temperature of 12°C Cloudy. The system clock indicates 6:54 PM on 2023-04-03.

The screenshot shows the same Amazon SageMaker Studio interface, but with the 'Security Settings' tab selected. It displays the following settings:

| Security Setting       | Value  |
|------------------------|--------|
| Learning_rate          | 0.009  |
| Max_bin                | 255    |
| Max_delta_step         | 0      |
| Max_depth              | 11     |
| Evaluation metric      | auto   |
| Min_data_in_leaf       | 26     |
| Min_gain_to_split      | 0      |
| Num_boost_round        | 5000   |
| Num_leaves             | 67     |
| Num_threads            | 0      |
| Tree_learner           | serial |
| Tweedie_variance_power | 1.5    |
| Use_dask               | False  |
| Verbosity              | 1      |

Below the hyperparameters, the 'Security Settings' section is expanded, showing:

- Network Isolation: enabled
- IAM Role: Default SageMaker Role (arn:aws:iam::224670572127:role/fast-ai-academic-12-Student-Azure)
- Volume encryption: No custom KMS key applied.
- Output encryption: No custom KMS key applied.
- VPC Settings: No custom VPC settings applied.

The 'Deploy Model' button is visible at the bottom of the main panel. The bottom status bar and system clock are identical to the first screenshot.

Whereas the model trained on notebook also takes nearly 101 training seconds which is equivalent to 2 minutes.

For algorithm specific hyper-parameters, we start by fetching python dictionary of the training hyper-parameters that the algorithm accepts with their default values.

```
In [8]: from sagemaker import hyperparameters

# Retrieve the default hyper-parameters for fine-tuning the model
hyperparameters = hyperparameters.retrieve_default(
    model_id=train_model_id, model_version=train_model_version
)
print(hyperparameters)

{'num_boost_round': '5000', 'early_stopping_rounds': '30', 'metric': 'auto', 'learning_rate': '0.009', 'num_leaves': '67', 'feature_fraction': '0.74', 'bagging_fraction': '0.53', 'bagging_freq': '5', 'max_depth': '11', 'min_data_in_leaf': '26', 'max_delta_step': '0.0', 'lambda_l1': '0.0', 'lambda_l2': '0.0', 'boosting': 'gbdt', 'min_gain_to_split': '0.0', 'tree_learner': 'serial', 'feature_fraction_bynode': '1.0', 'is_unbalance': 'False', 'max_bin': '255', 'tweedie_variance_power': '1.5', 'num_threads': '0', 'verbosity': '1', 'use_dask': 'False'}
```

```
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[473]#011train's rmse: 0.0526646#011val's rmse: 0.0525385
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[474]#011train's rmse: 0.0526473#011val's rmse: 0.0525326
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[475]#011train's rmse: 0.052632#011val's rmse: 0.0525453
Early stopping, best iteration is:
[445]#011train's rmse: 0.0530713#011val's rmse: 0.0524256
INFO:root:Saving model...
INFO:root:Info file not found at '_input_model_extracted/__models_info__.json'.
2023-04-03 22:27:53,836 sagemaker-training-toolkit INFO      Reporting training SUCCESS

2023-04-03 22:28:38 Uploading - Uploading generated training model
2023-04-03 22:28:38 Completed - Training job completed
Training seconds: 101
Billable seconds: 101
```

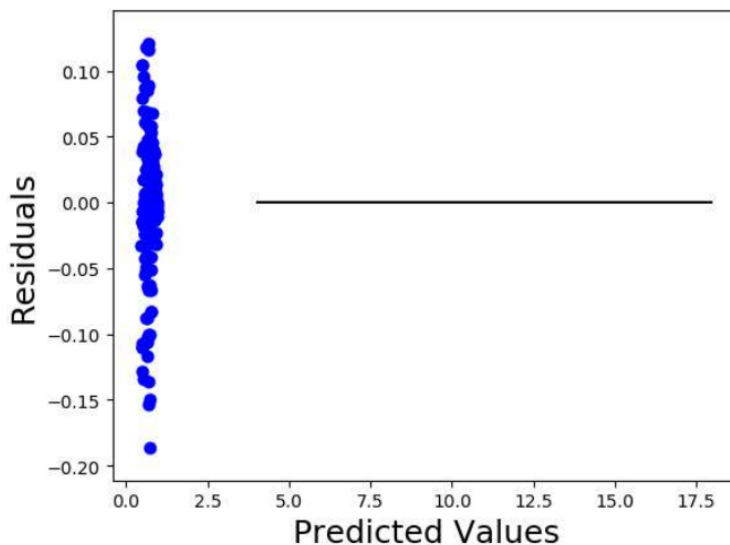
When we compare with the quality of model with respect to the validation data as well as with the residual values, it is noticed that with pretrained model, the best training rmse value is 0.0530 and best validation rmse is 0.0524. The pretrained model's R2 Score is 0.83478, mean squared error is 0.00327 and mean absolute error is 0.04144, which is as shown below.



```
In [6]: # Visualization: a residual plot to compare the model predictions and ground truth targets. For each example, the residual value
# is the subtraction between the prediction and ground truth target.
# We can see that the points in the residual plot are randomly dispersed around the horizontal axis  $y = 0$ ,
# which indicates the fitted regression model is appropriate for the ABALONE data

residuals = ground_truth_label.values[:, 0] - model_predictions
plt.scatter(model_predictions, residuals, color="blue", s=40)
plt.hlines(y=0, xmin=4, xmax=18)
```

```
plt.xlabel('Predicted Values', fontsize=18)
plt.ylabel('Residuals', fontsize=18)
plt.show()
```



```
In [7]: # Evaluate the model predictions quantitatively.
eval_r2_score = r2_score(ground_truth_label.values, model_predictions)
eval_mse_score = mean_squared_error(ground_truth_label.values, model_predictions)
eval_mae_score = mean_absolute_error(ground_truth_label.values, model_predictions)
print (
    f"{bold}Evaluation result on test data{unbold}:{newline}"
    f"{bold}{r2_score.__name__}{unbold}: {eval_r2_score}{newline}"
    f"{bold}{mean_squared_error.__name__}{unbold}: {eval_mse_score}{newline}"
    f"{bold}{mean_absolute_error.__name__}{unbold}: {eval_mae_score}{newline}"
)

Evaluation result on test data:
r2_score: 0.8347803463244985
mean_squared_error: 0.0032756969897500972
mean_absolute_error: 0.041440139984582014
```

When we compare with the quality of model with respect to the validation data as well as with the residual values, it is noticed that with model trained on notebook gives identical values with respect to pretrained model. The best training rmse value is 0.05307 and best validation rmse is 0.05242. The pretrained model's R2 Score is 0.83478, mean squared error is 0.00327 and mean absolute error is 0.04144, which is as shown below.



Early stopping, best iteration is:  
[445]#011train's rmse: 0.0530713#011val's rmse: 0.0524256  
INFO:root:Saving model...  
INFO:root:Info file not found at '\_input\_model\_extracted/\_models\_info\_.json'.  
2023-04-03 22:27:53,836 sagemaker-training-toolkit INFO Reporting training SUCCESS

The screenshot shows the AWS CloudWatch console interface. The left sidebar contains navigation options: Dashboards, Alarms, Logs, Metrics, X-Ray traces, Events, Application monitoring, and Insights. The 'Logs' section is selected, and the 'Log groups' tab is active. The main panel displays a list of log events for the log group `/aws/sagemaker/TrainingJobs`. The events are filtered by the search term `LightGBM`. The log messages show the training progress of a LightGBM model, including RMSE values for training and validation sets. The final log message indicates that the training was successful and that the best iteration was found.

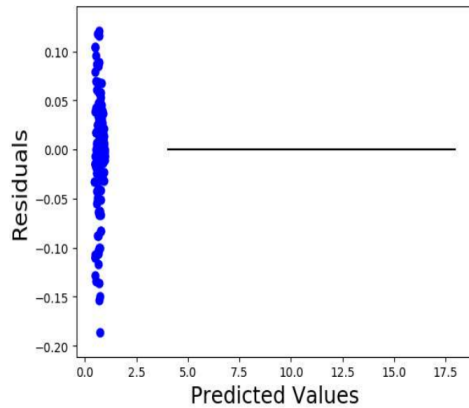
| Timestamp                     | Message  |
|-------------------------------|--|
| 2023-04-03T12:27:53.720-04:00 | [373]#011train's rmse: 0.054330#011val's rmse: 0.052899                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [379]#011train's rmse: 0.054309#011val's rmse: 0.052842                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [380]#011train's rmse: 0.054076#011val's rmse: 0.052843                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [381]#011train's rmse: 0.054025#011val's rmse: 0.052825                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [382]#011train's rmse: 0.054049#011val's rmse: 0.052804                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [383]#011train's rmse: 0.054036#011val's rmse: 0.052876                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [384]#011train's rmse: 0.053996#011val's rmse: 0.052844                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [385]#011train's rmse: 0.053975#011val's rmse: 0.052836                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [386]#011train's rmse: 0.053957#011val's rmse: 0.052833                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [387]#011train's rmse: 0.053937#011val's rmse: 0.052837                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [388]#011train's rmse: 0.053919#011val's rmse: 0.052807                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [389]#011train's rmse: 0.053902#011val's rmse: 0.052799                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |
| 2023-04-03T12:27:53.720-04:00 | [390]#011train's rmse: 0.053879#011val's rmse: 0.052793                    |
| 2023-04-03T12:27:53.720-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf |

This screenshot shows the same AWS CloudWatch console interface as the previous one, but with more log events displayed. The final log message indicates that the training was successful and that the best iteration was found. The log messages show the training progress of a LightGBM model, including RMSE values for training and validation sets. The final log message indicates that the training was successful and that the best iteration was found.

| Timestamp                     | Message  |
|-------------------------------|--|
| 2023-04-03T12:27:53.723-04:00 | [440]#011train's rmse: 0.052259#011val's rmse: 0.052498                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [441]#011train's rmse: 0.052239#011val's rmse: 0.052491                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [442]#011train's rmse: 0.052207#011val's rmse: 0.052465                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [443]#011train's rmse: 0.052190#011val's rmse: 0.052473                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [444]#011train's rmse: 0.052179#011val's rmse: 0.052463                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [445]#011train's rmse: 0.052155#011val's rmse: 0.052464                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [446]#011train's rmse: 0.052137#011val's rmse: 0.052494                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [447]#011train's rmse: 0.052129#011val's rmse: 0.052495                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [448]#011train's rmse: 0.052126#011val's rmse: 0.052501                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [449]#011train's rmse: 0.052124#011val's rmse: 0.052493                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [470]#011train's rmse: 0.052144#011val's rmse: 0.052828                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [471]#011train's rmse: 0.052075#011val's rmse: 0.052515                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [472]#011train's rmse: 0.052049#011val's rmse: 0.052574                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [473]#011train's rmse: 0.052046#011val's rmse: 0.052595                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [474]#011train's rmse: 0.052047#011val's rmse: 0.052526                            |
| 2023-04-03T12:27:53.723-04:00 | [LightGBM] [warning] No further splits with positive gain, best gain: -inf         |
| 2023-04-03T12:27:53.723-04:00 | [475]#011train's rmse: 0.052030#011val's rmse: 0.052543                            |
| 2023-04-03T12:27:53.724-04:00 | Early stopping, best iteration is:   |
| 2023-04-03T12:27:53.724-04:00 | [445]#011train's rmse: 0.0530713#011val's rmse: 0.0524256                          |
| 2023-04-03T12:27:53.724-04:00 | INFO:root:Saving model...  |
| 2023-04-03T12:27:53.724-04:00 | INFO:root:Info file not found at '_input_model_extracted/_models_info_.json'.      |
| 2023-04-03T12:27:53.836-04:00 | 2023-04-03 22:27:53,836 sagemaker-training-toolkit INFO Reporting training success |

```
In [22]: # Visualization: a residual plot to compare the model predictions and ground truth targets. For each example, the residual value
# is the subtraction between the prediction and ground truth target.
# We can see that the points in the residual plot are randomly dispersed around the horizontal axis y = 0,
# which indicates the fitted regression model is appropriate for the ABALONE data

residuals = ground_truth_label.values[:, 0] - model_predictions
plt.scatter(model_predictions, residuals, color="blue", s=40)
plt.hlines(y=0, xmin=4, xmax=18)
plt.xlabel("Predicted Values", fontsize=18)
plt.ylabel("Residuals", fontsize=18)
plt.show()
```



```
In [23]: # Evaluate the model predictions quantitatively.
eval_r2_score = r2_score(ground_truth_label.values, model_predictions)
eval_mse_score = mean_squared_error(ground_truth_label.values, model_predictions)
eval_mae_score = mean_absolute_error(ground_truth_label.values, model_predictions)
print(
    f"{bold}Evaluation result on test data{unbold}:{newline}"
    f"{bold}{r2_score.__name__}{unbold}: {eval_r2_score}{newline}"
    f"{bold}{mean_squared_error.__name__}{unbold}: {eval_mse_score}{newline}"
    f"{bold}{mean_absolute_error.__name__}{unbold}: {eval_mae_score}{newline}"
)
```

```
Evaluation result on test data:
r2_score: 0.8347803463244985
mean_squared_error: 0.0032756969897500972
mean_absolute_error: 0.041440139984582014
```