


## Loading all necessary Packages/Libraries

In [1]:

```
1 # Loading the iconic trio 
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 # Importing model_selection to get access to some dope functions like GridSearchCV()
7 from sklearn import model_selection
8
9 # from sklearn.externals import joblib
10
11 # Loading models
12 from sklearn import linear_model
13 from sklearn import svm
14 from sklearn import tree
15 from sklearn import ensemble
16
17 # custom
18 import helper
19
20 # Loading black for formatting codes
21 %load_ext blackcellmagic
```

## Styling Tables

In [2]:

```
1 %%HTML
2 <style type='text/css'>
3 table.dataframe th, table.dataframe td{
4     border: 3px solid purple !important;
5     color: solid black !important;
6 }
7 </style>
```

## Loading the Dataset

In [3]:

```
1 # Loading dataset
2 filename = "Clean_Akosombo_data.csv"
3 akosombo = helper.load_csv_data(filename)
```

Successfully loaded!

## Splitting the Dataset

In [4]:

```
1 # Splitting dataset
2 target_variable = "generation"
3 X, y, X_train, X_test, y_train, y_test = helper.split_data(akosombo, target_variable)
```

Data is splitted into X, y, X\_train, X\_test, y\_train, y\_test.

Shape Info of Features Training Set:

Number of datapoints (rows): 10001

Number of features (columns): 2

Shape Info of Features Test Set:

Number of datapoints (rows): 2501

Number of features (columns): 2

## Scaling the Dataset

In [5]:

```
1 # Data Scaling
2 X_train, X_test = helper.scale(X_train, X_test)
```

## Chosing Baseline Models and Training Models

In [6]:

```
1  # Instantiating baseline models
2  models = [
3      ("Linear Regression", linear_model.LinearRegression()),
4      ("Lasso", linear_model.Lasso()),
5      ("Ridge", linear_model.Ridge()),
6      ("SVR", svm.LinearSVR()),
7      ("Decision Tree", tree.DecisionTreeRegressor()),
8      ("Random Forest", ensemble.RandomForestRegressor()),
9  ]
10
11 model_names = []
12 accuracies = []
13
14 # Fitting models to Training Dataset and Scoring them on Test set
15 for dataset_name, dataset in [("Akosomba_Data", akosombo)]:
16     for model_name, model in models:
17         regressor_model = model
18         regressor_model.fit(X_train, y_train)
19
20         accuracy = regressor_model.score(X_test, y_test)
21         print(dataset_name, model_name, accuracy)
22
23         model_names.append(model_name)
24         accuracies.append(accuracy)
```

Akosomba\_Data Linear Regression 0.9697494269312791

Akosomba\_Data Lasso -0.00018218953746829136

Akosomba\_Data Ridge 0.9697837468881153

Akosomba\_Data SVR 0.9696773290140956

Akosomba\_Data Decision Tree 0.962324877287925

Akosomba\_Data Random Forest 0.9764811116031492

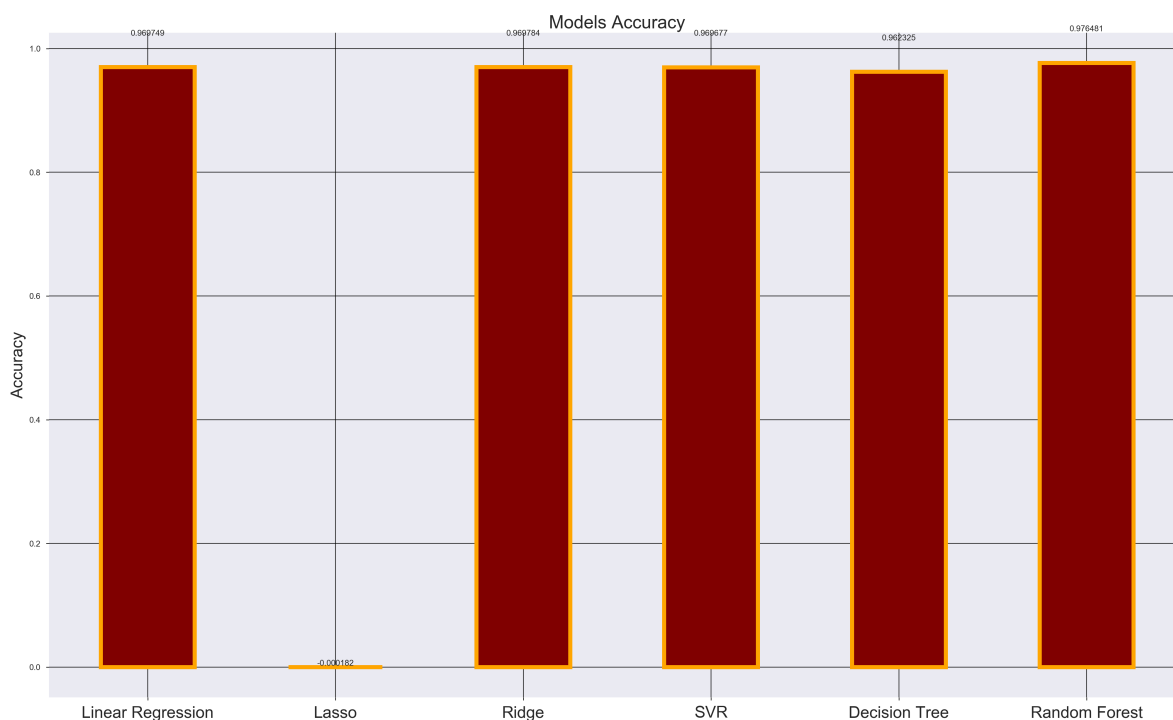
## Visualizing Models' Accuracy with Bar Charts

In [7]:

```

1  # Size in inches (width, height) & resolution(DPI)
2  plt.figure(figsize=(25, 15), dpi=200)
3
4  x_loc = np.arange(len(models)) # the x locations for the groups
5  width = 0.5 # bar width
6
7  # plotting the graphs with bar chart
8  models_graph = plt.bar(
9      x_loc, accuracies, width, color="maroon", edgecolor="orange", linewidth=5,
10 )
11
12 plt.title("Models Accuracy", fontsize=22)
13 plt.xticks(x_loc, model_names, fontsize=20)
14 plt.ylabel("Accuracy", fontsize=20)
15 plt.grid(b=True, which="both", axis="both", color="black", linewidth=0.8)
16
17 # adding model accuracy on top of every bar
18 def addLabel(models):
19     for model in models:
20         height = model.get_height()
21         plt.text(
22             model.get_x() + model.get_width() / 2.0,
23             1.05 * height,
24             "%f" % height,
25             ha="center",
26             va="bottom",
27         )
28
29
30 addLabel(models_graph)
31
32 plt.savefig('Bar_Charts_of_Models_and_their_Accuracy.png', dpi=300, transparent=True)
33
34 plt.show()

```



## Evaluating Models



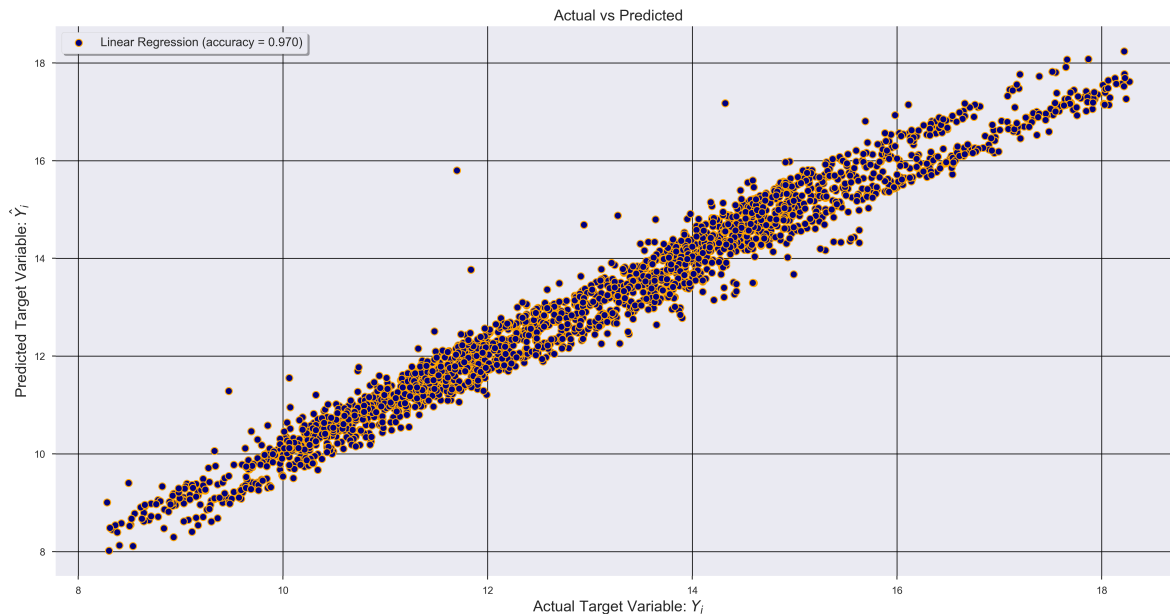
In [8]:

```

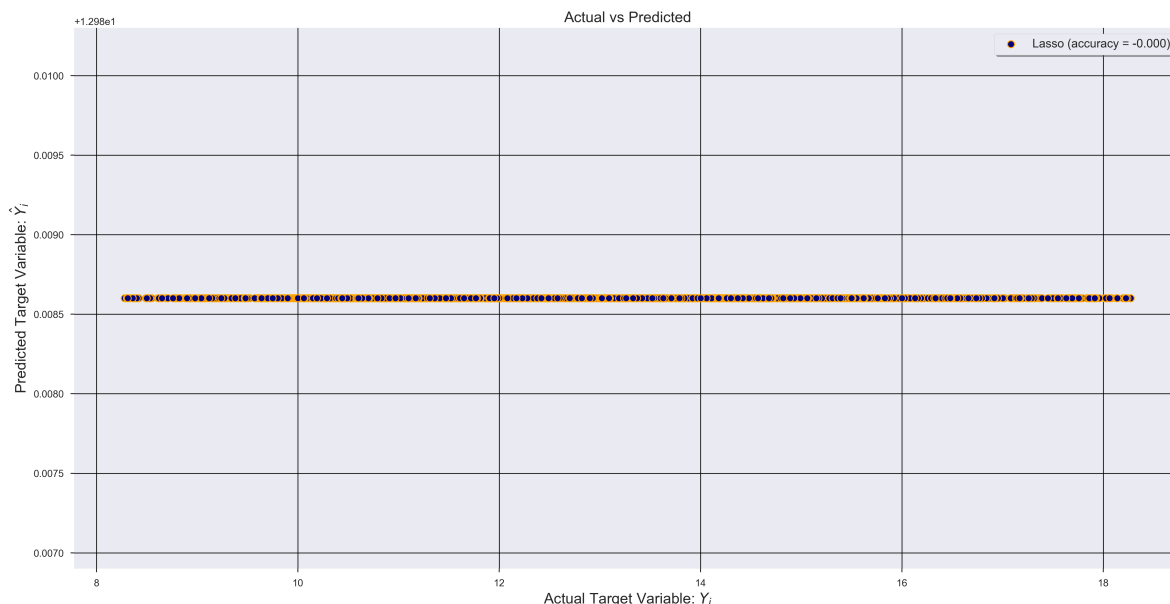
1 # Model Evaluation
2 for model_name, model in models:
3     helper.evaluate(X_test, y_test, model_name, model)

```

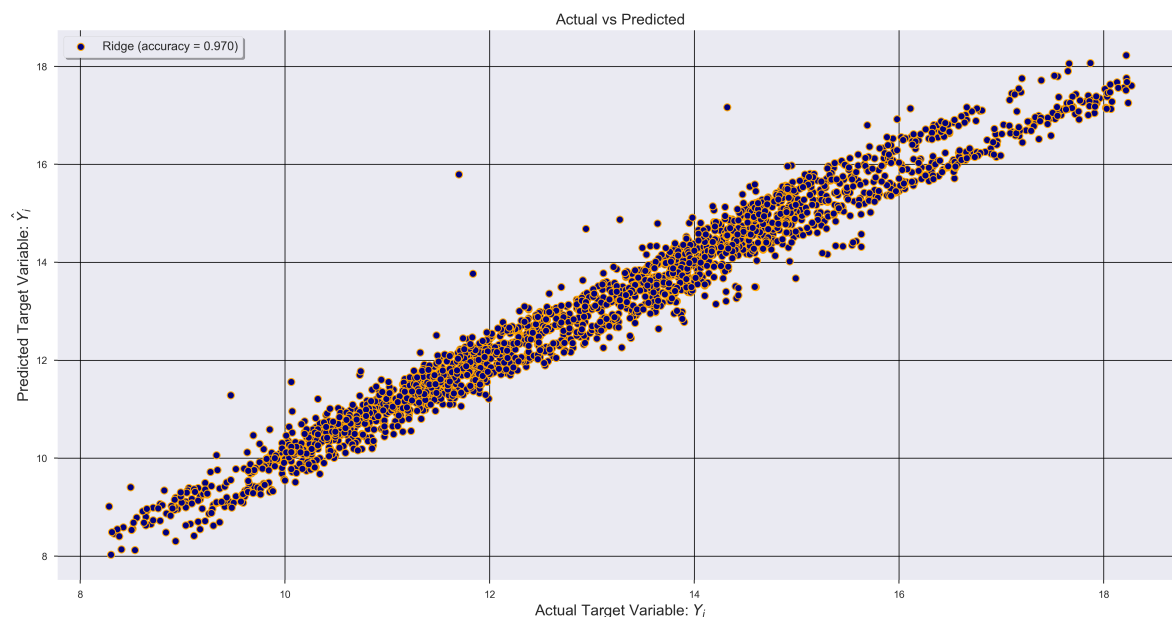
Linear Regression Mean Squared Error: 0.14625233556901102  
 Linear Regression Root Mean Squared Error: 0.38242951712571954  
 Linear Regression R2 Score: 0.969749426931279  
 Linear Regression Explained Variance Score: 0.969772234576705  
 Linear Regression Mean Absolute Error: 0.29263695465956613  
 Linear Regression Meadian Abosolute Error: 0.2677789934798991  
 Linear Regression Mean Squared Log Error: 0.0007141004739016008



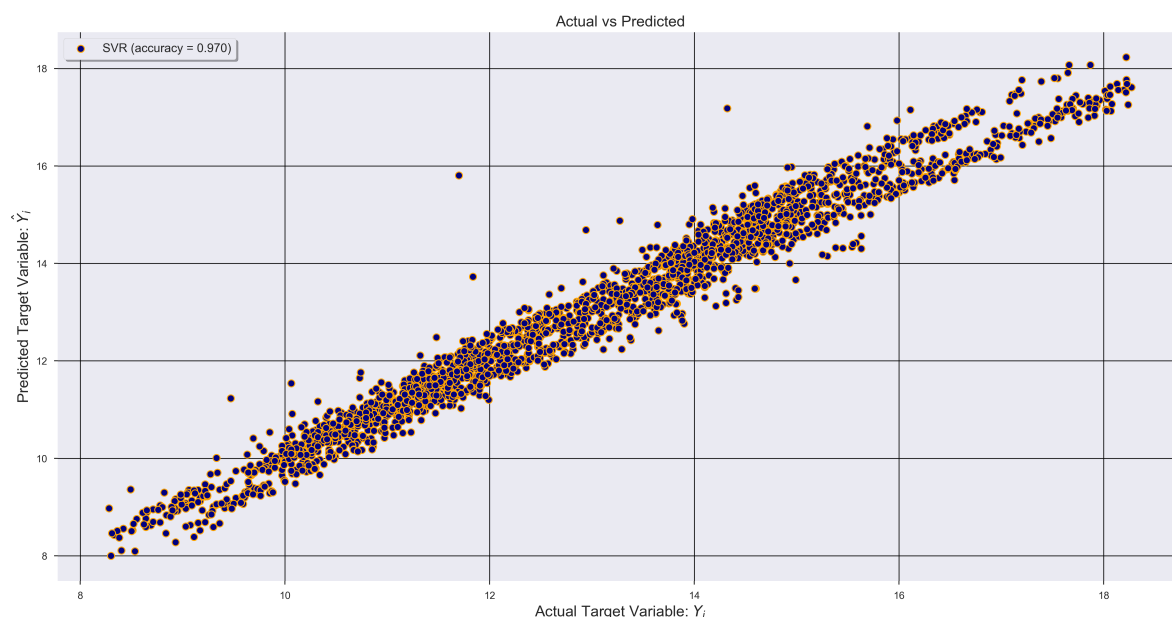
Lasso Mean Squared Error: 4.83557719326759  
 Lasso Root Mean Squared Error: 2.198994586911844  
 Lasso R2 Score: -0.00018218953746829136  
 Lasso Explained Variance Score: 0.0  
 Lasso Mean Absolute Error: 1.8671255680309624  
 Lasso Meadian Abosolute Error: 1.7113991700829914  
 Lasso Mean Squared Log Error: 0.0253423068715663



Ridge Mean Squared Error: 0.14608640899854475  
Ridge Root Mean Squared Error: 0.3822125181081132  
Ridge R2 Score: 0.9697837468881153  
Ridge Explained Variance Score: 0.9698066436041364  
Ridge Mean Absolute Error: 0.2926293335519423  
Ridge Meadian Abosolute Error: 0.26566972718727655  
Ridge Mean Squared Log Error: 0.0007123218460580986

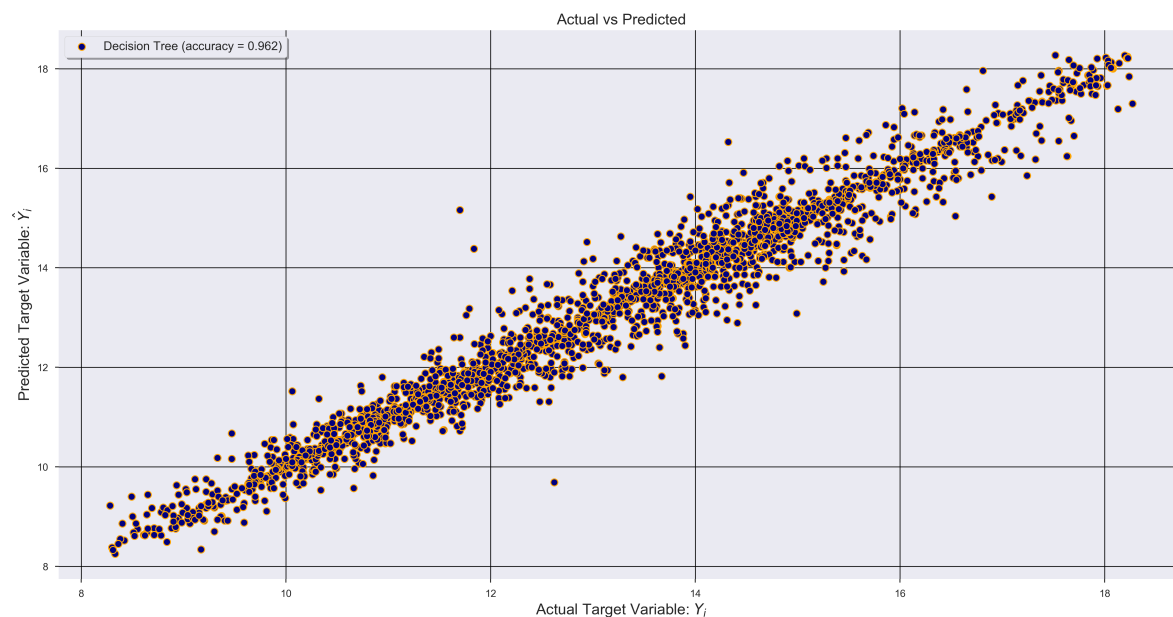


SVR Mean Squared Error: 0.14660090710693777  
SVR Root Mean Squared Error: 0.3828849789518228  
SVR R2 Score: 0.9696773290140956  
SVR Explained Variance Score: 0.9696893927526878  
SVR Mean Absolute Error: 0.2910085934765154  
SVR Meadian Abosolute Error: 0.26332581829105983  
SVR Mean Squared Log Error: 0.0007128432992722564

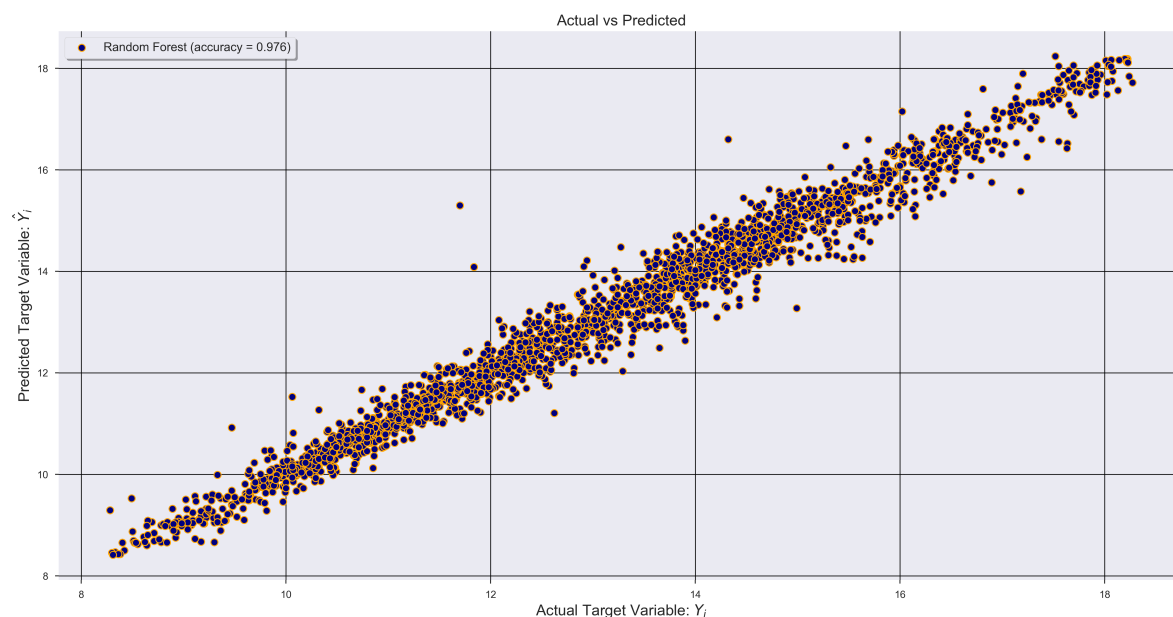


Decision Tree Mean Squared Error: 0.18214777872051172  
Decision Tree Root Mean Squared Error: 0.4267877443419758  
Decision Tree R2 Score: 0.962324877287925  
Decision Tree Explained Variance Score: 0.9623848910919239  
Decision Tree Mean Absolute Error: 0.2797403438624543

Decision Tree Meadian Abosolute Error: 0.155000000000000114  
Decision Tree Mean Squared Log Error: 0.0009010063568695894



Random Forest Mean Squared Error: 0.11370668417462113  
Random Forest Root Mean Squared Error: 0.337204217314406  
Random Forest R2 Score: 0.9764811116031491  
Random Forest Explained Variance Score: 0.9765071315241614  
Random Forest Mean Absolute Error: 0.23050171260067406  
Random Forest Meadian Abosolute Error: 0.151079999999999855  
Random Forest Mean Squared Log Error: 0.0005612420264805661



## Cross Validating Models

### Cross Validating with a single metric



In [9]:

```

1  # Splitting data into 10 folds
2  cv_kfold = model_selection.KFold(n_splits=10, shuffle=True, random_state=23)
3  scorer = "r2"
4
5  model_names = []
6  cv_mean_scores = []
7  cv_std_scores = []
8
9  for model_name, model in models:
10     regressor_model = model
11     model_scores = model_selection.cross_val_score(
12         regressor_model, X, y, cv=cv_kfold, scoring=scorer, n_jobs=-1, verbose=1,
13     )
14
15     print(
16         f"{model_name} Accuracy: %0.2f (+/- %0.2f)"
17         % (model_scores.mean(), model_scores.std() * 2)
18     )
19
20     model_names.append(model_name)
21     cv_mean_scores.append(model_scores.mean())
22     cv_std_scores.append(model_scores.std())

```

[Parallel(n\_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.  
 [Parallel(n\_jobs=-1)]: Done 6 out of 10 | elapsed: 1.1s remaining: 0.7s

[Parallel(n\_jobs=-1)]: Done 10 out of 10 | elapsed: 1.1s finished  
 [Parallel(n\_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.  
 [Parallel(n\_jobs=-1)]: Done 6 out of 10 | elapsed: 0.0s remaining: 0.0s

[Parallel(n\_jobs=-1)]: Done 10 out of 10 | elapsed: 0.0s finished  
 [Parallel(n\_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.  
 [Parallel(n\_jobs=-1)]: Done 6 out of 10 | elapsed: 0.0s remaining: 0.0s

[Parallel(n\_jobs=-1)]: Done 10 out of 10 | elapsed: 0.0s finished  
 [Parallel(n\_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.

Linear Regression Accuracy: 0.97 (+/- 0.00)  
 Lasso Accuracy: 0.95 (+/- 0.00)  
 Ridge Accuracy: 0.97 (+/- 0.00)

[Parallel(n\_jobs=-1)]: Done 6 out of 10 | elapsed: 1.8s remaining: 1.2s  
 [Parallel(n\_jobs=-1)]: Done 10 out of 10 | elapsed: 2.6s finished  
 [Parallel(n\_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.  
 [Parallel(n\_jobs=-1)]: Done 6 out of 10 | elapsed: 0.0s remaining: 0.0s

[Parallel(n\_jobs=-1)]: Done 10 out of 10 | elapsed: 0.0s finished  
 [Parallel(n\_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.

SVR Accuracy: 0.84 (+/- 0.11)  
 Decision Tree Accuracy: 0.96 (+/- 0.01)

[Parallel(n\_jobs=-1)]: Done 6 out of 10 | elapsed: 3.5s remaining: 2.3s

Random Forest Accuracy: 0.98 (+/- 0.00)

[Parallel(n\_jobs=-1)]: Done 10 out of 10 | elapsed: 6.0s finished

In [10]:

```
1 cv_results = pd.DataFrame({"model_name": model_names, "mean_score": cv_mean_scores, "std_score": cv_std_scores})
2 cv_results.sort_values("mean_score", ascending=False, inplace=True)
3 cv_results.to_csv("cross_validation_results.csv", index=True)
4 cv_results
```

Out[10]:

	model_name	mean_score	std_score
5	Random Forest	0.978132	0.001437
2	Ridge	0.971263	0.001156
0	Linear Regression	0.971263	0.001156
4	Decision Tree	0.964552	0.002598
1	Lasso	0.946265	0.001854
3	SVR	0.843361	0.052938

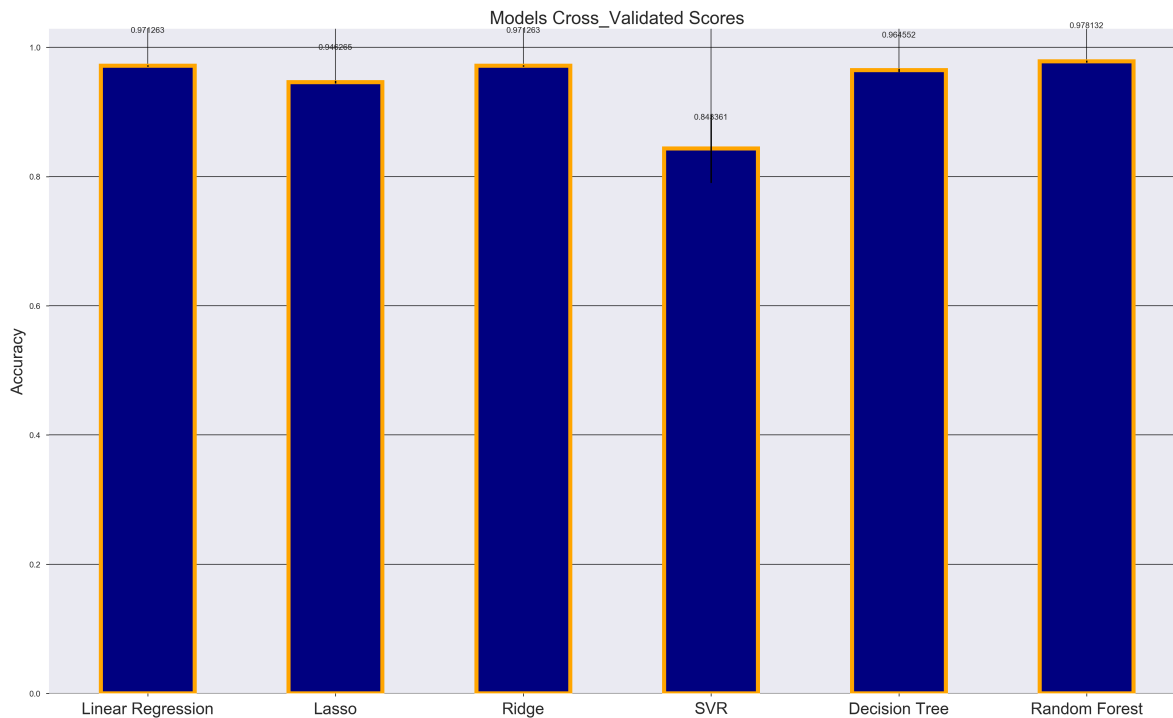
## Visualizing Cross Validated Models with Bar Charts

In [11]:

```

1 plt.figure(figsize=(25, 15), dpi=200)
2
3 x_loc = np.arange(len(models))
4 width = 0.5
5
6 models_graph = plt.bar(
7     x_loc, cv_mean_scores, width, yerr=cv_std_scores, color="navy", edgecolor="orange",
8 )
9 plt.title("Models Cross_Validated Scores", fontsize=22)
10 plt.xticks(x_loc, model_names, fontsize=20)
11 plt.ylabel("Accuracy", fontsize=20)
12 plt.grid(b=True, which="both", axis="both", color="black", linewidth=0.8)
13
14 addLabel(models_graph)
15
16 plt.savefig('Bar_Charts_of_Cross_Validated_Models_and_their_Accuracy.png', dpi=300, tra
17
18 plt.show()

```



## Training the Model with the Highest Score with Default Hyperparameters

In [12]:

```

1  # Instantiating model object
2  high_score_model = ensemble.RandomForestRegressor()
3
4  # Fitting the model on Train set
5  high_score_model.fit(X_train, y_train)
6
7  # Scoring the model on Test set
8  high_score_model_accuracy = high_score_model.score(X_test, y_test)
9
10 print(
11     f"Model without tuned hyperparameters has an accuracy of {high_score_model_accuracy}
12 )

```

Model without tuned hyperparameters has an accuracy of 0.9763555356961289

## Predicting with the Trained Model and Saving Predicted Results as csv

In [13]:

```

1  y_pred = high_score_model.predict(X_test)
2
3  data = pd.DataFrame({"actual_generation": list(y_test), "predicted_generation": list(y_
4  data.to_csv("model_predicted_values.csv", index=True)

```

In [14]:

```
1 data.head(10)
```

Out[14]:

	actual_generation	predicted_generation
0	17.596	17.875130
1	15.630	15.209890
2	10.850	11.478600
3	14.520	14.756200
4	8.420	8.510460
5	13.640	14.302863
6	14.010	13.693800
7	9.699	9.551900
8	13.900	12.544700
9	11.150	10.977800

## Optimizing the Hyperparameter of the Best Model with GridSearchCV

In [16]:

```

1  # Kfold with with n_splits = 5 to split the Dataset into 5-folds
2  kfold = model_selection.KFold(n_splits=5, shuffle=True, random_state=23)
3
4  # Dictionary of parameters to tune
5  parameters = {
6      "n_estimators" : [120, 500, 800, 1200],
7      "max_depth" : [15, 25, 30, None],
8      "min_samples_split" : [5, 10, 15, 100],
9      "min_samples_leaf" : [1, 2, 5, 10],
10     "max_features" : ["log2", "sqrt", None],
11 }
12
13 scorer = "r2"
14
15 # Instantiating Search object
16 grid = model_selection.RandomizedSearchCV(
17     estimator=high_score_model,
18     param_distributions=parameters,
19     scoring=scorer,
20     cv=kfold,
21     n_jobs=-1,
22     verbose=1,
23 )
24
25 # Fit the grid object on Training Dataset
26 grid.fit(X_train, y_train)

```

Fitting 5 folds for each of 10 candidates, totalling 50 fits

```

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done 34 tasks      | elapsed: 47.0s
[Parallel(n_jobs=-1)]: Done 50 out of 50 | elapsed: 1.4min finished

```

Out[16]:

```

RandomizedSearchCV(cv=KFold(n_splits=5, random_state=23, shuffle=True),
                  error_score=nan,
                  estimator=RandomForestRegressor(bootstrap=True,
                                                    ccp_alpha=0.0,
                                                    criterion='mse',
                                                    max_depth=None,
                                                    max_features='auto',
                                                    max_leaf_nodes=None,
                                                    max_samples=None,
                                                    min_impurity_decrease=0.
0,
                                                    min_impurity_split=None,
                                                    min_samples_leaf=1,
                                                    min_samples_split=2,
                                                    min_weight_fraction_leaf=
0....
                                                    random_state=None, verbose
e=0,
                                                    warm_start=False),
                  iid='deprecated', n_iter=10, n_jobs=-1,
                  param_distributions={'max_depth': [15, 25, 30, None],
                                      'max_features': ['log2', 'sqrt', Non
e],
                                      'min_samples_leaf': [1, 2, 5, 10],

```

```

    'min_samples_split': [5, 10, 15, 10
0],
    'n_estimators': [120, 500, 800, 120
0]}},
pre_dispatch='2*n_jobs', random_state=None, refit=True,
return_train_score=False, scoring='r2', verbose=1)

```

In [17]:

```

1 results = pd.DataFrame(grid.cv_results_[["params", "mean_test_score", "std_test_score",
2 results.sort_values("rank_test_score", inplace=True)
3 results.to_csv("hyperparameter_tuning_results.csv", index=True)
4 results

```

Out[17]:

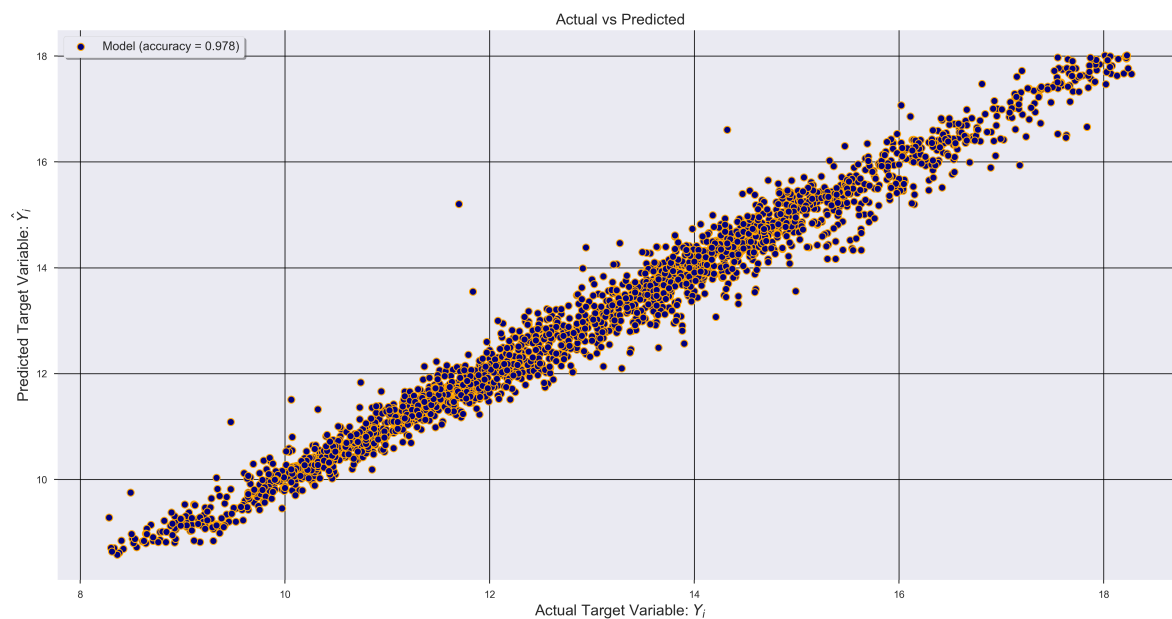
	params	mean_test_score	std_test_score	rank_test_score
7	{'n_estimators': 1200, 'min_samples_split': 10...	0.978835	0.000849	1
8	{'n_estimators': 1200, 'min_samples_split': 15...	0.978787	0.000782	2
4	{'n_estimators': 800, 'min_samples_split': 15,...	0.978703	0.000840	3
0	{'n_estimators': 1200, 'min_samples_split': 5,...	0.978700	0.000827	4
9	{'n_estimators': 1200, 'min_samples_split': 15...	0.978694	0.000824	5
2	{'n_estimators': 800, 'min_samples_split': 10,...	0.978583	0.000835	6
1	{'n_estimators': 1200, 'min_samples_split': 10...	0.978556	0.000854	7
6	{'n_estimators': 120, 'min_samples_split': 5, ...	0.978369	0.000830	8
3	{'n_estimators': 500, 'min_samples_split': 10,...	0.977691	0.001001	9
5	{'n_estimators': 800, 'min_samples_split': 100...	0.971914	0.001033	10

## Evaluating the Best Estimator from the GridSearch

In [18]:

```
1 best_estimator = grid.best_estimator_  
2 helper.evaluate(X_test, y_test, "Model", best_estimator)
```

Model Mean Squared Error: 0.1083768208580254  
Model Root Mean Squared Error: 0.32920634996613507  
Model R2 Score: 0.9775835310556501  
Model Explained Variance Score: 0.9775842437878877  
Model Mean Absolute Error: 0.23278235862104127  
Model Meadian Abosulte Error: 0.1657286307725485  
Model Mean Squared Log Error: 0.0005422388176224308



## Predicting with the Best Estimator and Saving Predicted Results as csv

In [19]:

```

1 tune_y_pred = best_estimator.predict(X_test)
2
3 hyp_tune_data = pd.DataFrame(
4     {"generation": list(y_test), "predicted_generation": list(tune_y_pred),}
5 )
6
7 hyp_tune_data.to_csv("tune_model_predicted_values.csv", index=True)
8
9 hyp_tune_data.head(10)

```

Out[19]:

	generation	predicted_generation
0	17.596	17.742404
1	15.630	15.212326
2	10.850	11.358161
3	14.520	14.700633
4	8.420	8.687071
5	13.640	14.385984
6	14.010	13.830590
7	9.699	9.565270
8	13.900	12.565679
9	11.150	11.021706

## Saving the Best Estimator with joblib

In [20]:

```

1 import joblib
2
3 joblib.dump(best_estimator, "Random_Forest_Regressor.joblib")

```

Out[20]:

```
['Random_Forest_Regressor.joblib']
```

## Feature Importance



In [21]:

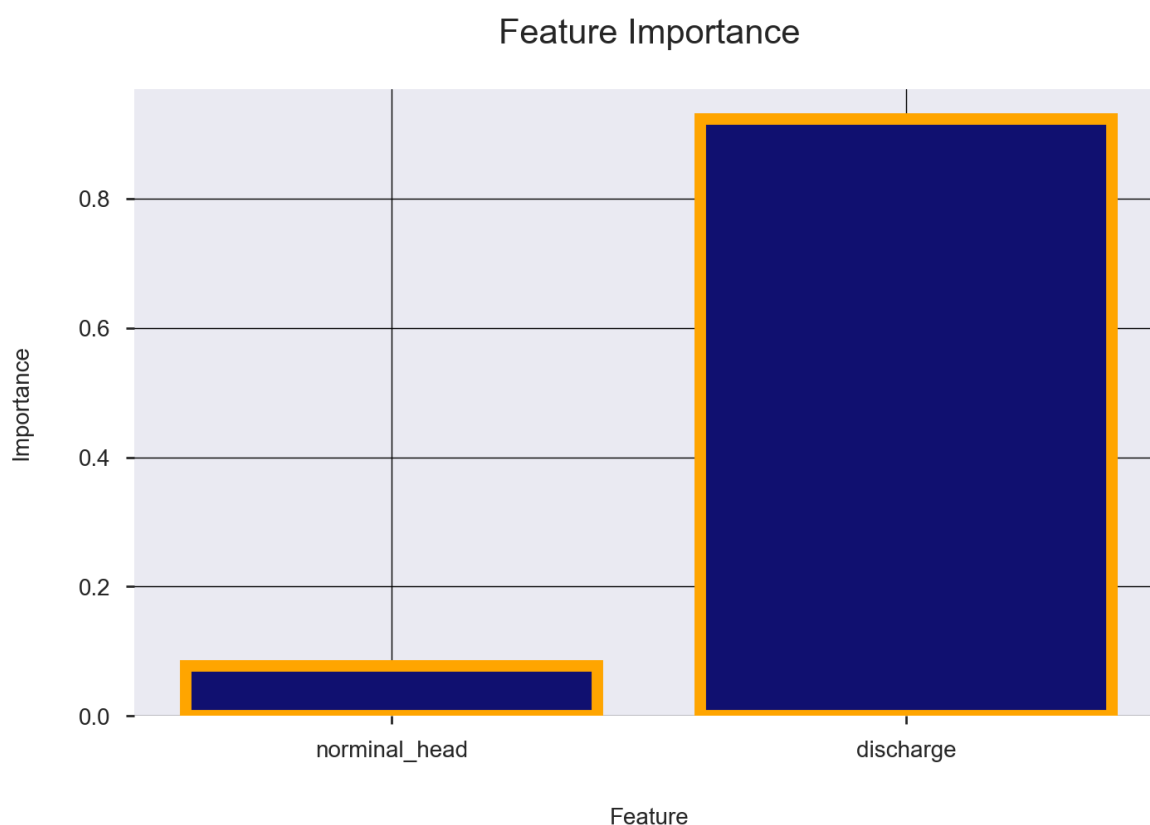
```
1 features = ['nominal_head', 'discharge']
2 importances = best_estimator.feature_importances_
3
4 feature_importance = pd.DataFrame({
5     'feature': features,
6     'importance': importances,
7 })
8
9 feature_importance.sort_values("importance", inplace=True, ascending=False)
10 feature_importance.to_csv("feature_importance.csv", index=True)
11 feature_importance
```

Out[21]:

	feature	importance
1	discharge	0.922701
0	nominal_head	0.077299

In [22]:

```
1 import seaborn as sns
2
3 plt.figure(figsize=(8, 5), dpi=200)
4
5 sns.barplot(x=features, y=importances, color="navy", edgecolor="orange", linewidth=5)
6 plt.title("Feature Importance", size=15, pad=20)
7 plt.xlabel("Feature", fontsize=10, labelpad=20)
8 plt.ylabel("Importance", fontsize=10, labelpad=20)
9
10 plt.grid(b=True, which="both", axis="both", color="black", linewidth=0.5)
11
12 plt.savefig('feature_importance.png', dpi=300, transparent=True)
13
14 plt.show()
```



In [ ]:

1