

PROJECT: PREDICT ENERGY CONSUMPTION





Imagine living in a house where every single watt of electricity you use is meticulously recorded, each of which contributes to a vast pool of data. By analyzing this detailed household power consumption data recorded over nearly 4 years, an energy company can help customers achieve sustainable energy usage while balancing their energy generation. With predictive models, the company can optimize energy usage, forecast future consumption, and provide tailored recommendations. Your task is to use this dataset to build a model that predicts power consumption, benefiting both the energy provider and its customers.

The Data

Available in `df_train.csv` and `df_test.csv`:

Column	Type	Description
date	chr	Date of the measurement
power_consumption	dbl	Daily power consumption (in kilowatts)
year	int	Year of the measurement
semester	int	Semester of the measurement (1 for Jan-Jun, 2 for Jul-Dec)
quarter	int	Quarter of the measurement (1 for Q1, 2 for Q2, 3 for Q3, 4 for Q4)
day_in_week	chr	Day of the week of the measurement (e.g., Monday, Tuesday)
week_in_year	int	Week number in the year of the measurement

How likely are you to recommend DataLab to a friend or co-worker?

Not at all likely 0 1 2 3 4 5 6 7 8 9 10 Extremely likely

powered by InMoment

```
# Load necessary libraries
suppressPackageStartupMessages(library(dplyr))
library(lubridate)
library(ranger)
library(xgboost)
library(ggplot2)

# Load and inspect the training and testing datasets
df_train <- read.csv("df_train.csv")
df_test <- read.csv("df_test.csv")

## Explore the structure of the dataset
glimpse(df_train)

# Load necessary libraries
library(dplyr)
library(lubridate)
library(ranger)
library(xgboost)
library(ggplot2)

# Read training and testing data from CSV files
df_train <- read.csv("df_train.csv")
df_test <- read.csv("df_test.csv")

# Display structure of the training data
glimpse(df_train)

# Convert 'date' column to Date type and 'day_in_week' column to factor in both
datasets
df_train <- df_train %>%
  mutate(date = as.Date(date, format = "%m/%d/%Y"),
         day_in_week = factor(day_in_week))
df_test <- df_test %>%
  mutate(date = as.Date(date, format = "%m/%d/%Y"),
         day_in_week = factor(day_in_week))

# Convert categorical variable 'day_in_week' to indicator variables using one-hot
encoding in both datasets
df_onehot_train <- model.matrix(~ day_in_week - 1, data = df_train) %>%
```

How likely are you to recommend DataLab to a friend or co-worker?

Not at all likely

0

1

2

3

4

5

6

7

8

9

10

Extremely likely

powered by InMoment

```

df_test <- mutate(df_test, df_onehot_test) %>% select(-c(day_in_week))

# Separate features and target variable for both training and testing datasets
train_x <- df_train %>% select(-power_consumption, -date)
train_y <- df_train[["power_consumption"]]
test_x <- df_test %>% select(-power_consumption, -date)
test_y <- df_test[["power_consumption"]]

# Train models, predict on test dataset and calculate RMSE for each model.
## Linear regression
lm_model <- lm(train_y ~ ., data = train_x)
lm_pred <- predict(lm_model, newdata = test_x)
lm_rmse <- sqrt(mean((test_y - lm_pred)^2))

## Random forest
rf_model <- ranger(power_consumption ~., data = df_train %>% select(-date),
  num.trees = 1000)
rf_pred <- predict(rf_model, data = df_test %>% select(-date))$predictions
rf_rmse <- sqrt(mean((test_y - rf_pred)^2))

## XGBoost
xgb_model <- xgboost(
  data = as.matrix(train_x),
  label = train_y,
  nrounds = 500,
  objective = "reg:squarederror",
  eta = 0.1,
  max_depth = 1,
  verbose = FALSE
)
xgb_pred <- predict(xgb_model, newdata = as.matrix(test_x))
xgb_rmse <- sqrt(mean((test_y - xgb_pred)^2))

# RMSE scores
data.frame(
  Model = c("Linear Regression", "Random Forest", "XGBoost"),
  RMSE = c(lm_rmse, rf_rmse, xgb_rmse)
)

# Get the lowest RMSE and assign it to selected_rmse
selected_rmse <- min(lm_rmse, rf_rmse, xgb_rmse)

```

How likely are you to recommend DataLab to a friend or co-worker?

Not at all likely

0

1

2

3

4

5

6

7

8

9

10

Extremely likely

```
Warning message in predict.lm(lm_model, newdata = test_x):
"prediction from a rank-deficient fit may be misleading"
```

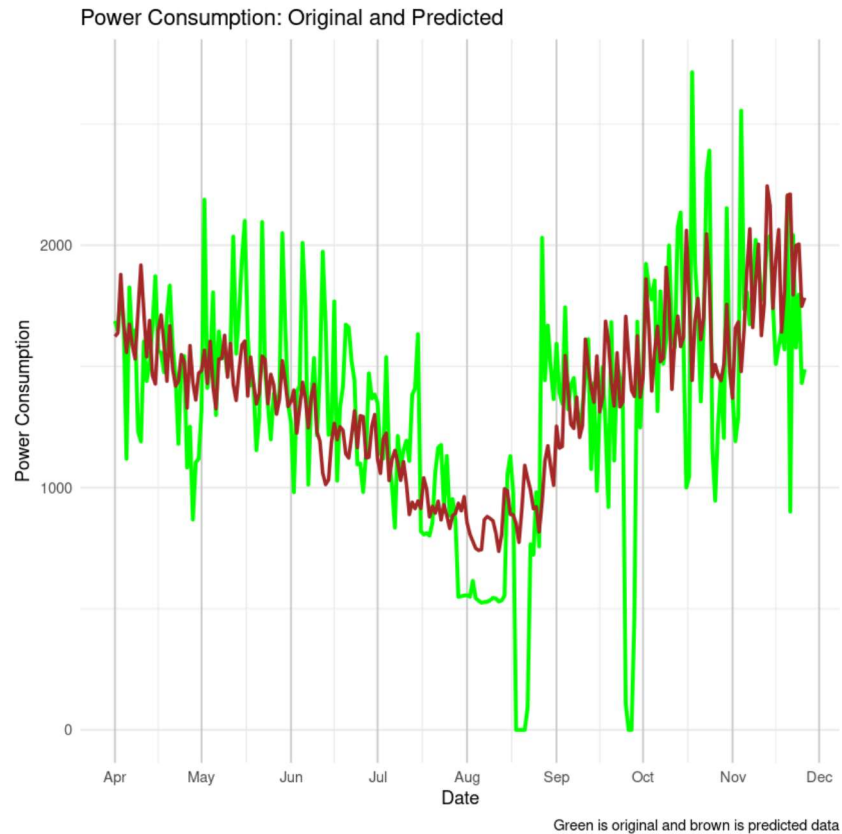
...	↑↓	Model	...	↑↓	...	↑↓	

Extremely likely

<https://www.datacamp.com/datalab/w/46d883bb-a224-40e2-907b-002ea58a2496/print-notebook/notebook.ipynb?emitCellOutputs=false&reducedMen...> 5/6

selected_rmse: 391.7094 kW

trend_similarity: Yes



How likely are you to recommend DataLab to a friend or co-worker?



Not at all likely

0

1

2

3

4

5

6

7

8

9

10

Extremely likely

powered by InMoment