

Bike sharing

Load default libraries

```
library(data.table)
library(dplyr)
library(lubridate)
library(ggplot2)
library(caret)
library(stringi)
library(xgboost)
library(reshape2)
theme_set(theme_bw())
set.seed(42)
```

Read the data

```
full_train = fread("./train.csv", header = T, sep = ",", integer64 = "numeric")
test = fread("./test.csv", header = T, sep = ",", integer64 = "numeric")
```

Change type to categorical variables

```
fix_types = function(data) {
  data %>%
    mutate(
      datetime = ymd_hms(datetime),
      workingday = 1 - workingday,
      weather = weather) %>%
    mutate(
      year = year(datetime),
      month = month(datetime),
      day = as.character(mday(datetime)),
      wday = wday(datetime),
      hour = hour(datetime),
      hour4 = as.character(round(hour(datetime) / 4)),
      hour3 = as.character(round(hour(datetime) / 3)),
      hour6 = as.character(round(hour(datetime) / 6)),
      hour8 = as.character(round((hour(datetime) - 4) / 8)),
      hour12 = as.character(round((hour(datetime)) / 12)),
      month2 = as.character(round((month(datetime)) / 2)))
}
full_train = fix_types(full_train)
test = fix_types(test)

train = full_train[as.integer(full_train$day) < 13]
validate = full_train[as.integer(full_train$day) >= 13]
```

Prepare features

```
dummy = dummyVars(~ workingday + season + weather + wday + hour +
                  year + hour4 + hour8 + hour12 + hour3 + hour6 + holiday +
                  workingday * hour4 + workingday * hour8 + workingday * hour12 +
                  workingday * hour3 + workingday * hour6 +
                  holiday * hour4 + holiday * hour8 + holiday * hour12 +
                  holiday * hour3 + holiday * hour6,
                  data = train,
                  fullRank = T,
                  levelsOnly = T)

summarise_stats = function(groups) {
  groups %>% summarise(
    avg_count = median(as.double(count)),
    var_count = var(as.double(count)),
    var_registered = var(as.double(registered)),
    casual_ratio1 = sum(casual) / sum(count),
    casual_ratio2 = mean(casual / count))
}

train_hour_statistics = train %>% group_by(hour, wday) %>% summarise_stats()
full_train_hour_statistics = full_train %>% group_by(hour, wday) %>% summarise_stats()
train_month_statistics = train %>% group_by(month, year) %>% summarise_stats()
full_train_month_statistics = full_train %>% group_by(month, year) %>% summarise_stats()

select_statistics = function(data, suffix) {
  stats = data %>% select(avg_count:casual_ratio2)
  names(stats) = paste0(names(stats), suffix)
  stats
}

names = names(predict(dummy, train) %>% as.data.frame())
prepare_dataset_features = function(data, hour_statistics, month_statistics) {
  categorical_data = predict(dummy, data) %>% as.data.frame()
  Missing = setdiff(names, names(categorical_data)) # Find names of missing columns
  categorical_data[Missing] = 0 # Add them, filled with '0's
  categorical_data = categorical_data[names]
  numeric_data = data %>% transmute(temp, atemp - temp, humidity, windspeed)
  hour_stats_data = inner_join(data, hour_statistics, by=c("hour", "wday"), p) %>%
    select_statistics(".hour")
  month_stats_data = inner_join(data, month_statistics, by=c("month", "year")) %>%
    select_statistics(".month")

  prepared_data = cbind(
    categorical_data,
    numeric_data,
    hour_stats_data,
    month_stats_data)
  as.matrix(prepared_data)
}
```

```

train_matrix = prepare_dataset_features(train,
                                       train_hour_statistics,
                                       train_month_statistics)
validate_matrix = prepare_dataset_features(validate,
                                           train_hour_statistics,
                                           train_month_statistics)

test_matrix = prepare_dataset_features(test,
                                       full_train_hour_statistics,
                                       full_train_month_statistics)
full_train_matrix = prepare_dataset_features(full_train,
                                             full_train_hour_statistics,
                                             full_train_month_statistics)

```

Prepare target

We going to predict $\log(Y+1)$ to optimize the target cost function

```

preparedTrainTarget = log(train$count + 1)

preparedValidateTarget = log(validate$count + 1)
preparedFullTarget = log(full_train$count + 1)

```

Xgboost train and cross-validate

Interesting link: [how to tune hyperparameters](#)

```

dtrain <- xgb.DMatrix(train_matrix, label = preparedTrainTarget)

xgbControl = list(
  subsample=0.8, colsample_bytree = 0.8, metrics=list("rmse"), gamma = 0.9,
  max.depth = 6, eta = 0.11, alpha = 1, lambda = 1, objective = "reg:linear"
)
model = xgboost(params = xgbControl, data = dtrain,
               nround=2000, nthread = 4, print.every.n = 500)

```

```

## [0] train-rmse:3.848741
## [500] train-rmse:0.233766
## [1000] train-rmse:0.229112
## [1500] train-rmse:0.226564

```

```

history <- xgb.cv(params = xgbControl, data = dtrain,
                 nround=2000, nthread = 4, nfold = 10, print.every.n = 500)

```

```

## [0] train-rmse:3.849032+0.002681 test-rmse:3.849260+0.025239
## [500] train-rmse:0.234618+0.001421 test-rmse:0.286852+0.015078
## [1000] train-rmse:0.230440+0.001408 test-rmse:0.284658+0.014445
## [1500] train-rmse:0.228490+0.001286 test-rmse:0.283699+0.014335

```

```
print(tail(history))
```

```
##      train.rmse.mean train.rmse.std test.rmse.mean test.rmse.std
## 1:      0.227079      0.001799      0.283079      0.014174
## 2:      0.227066      0.001811      0.283069      0.014165
## 3:      0.227066      0.001811      0.283069      0.014167
## 4:      0.227066      0.001811      0.283068      0.014166
## 5:      0.227066      0.001811      0.283066      0.014171
## 6:      0.227066      0.001811      0.283063      0.014168
```

Validate Score

```
validate_predictions = predict(model, validate_matrix)
RMSE(validate_predictions, preparedValidateTarget)
```

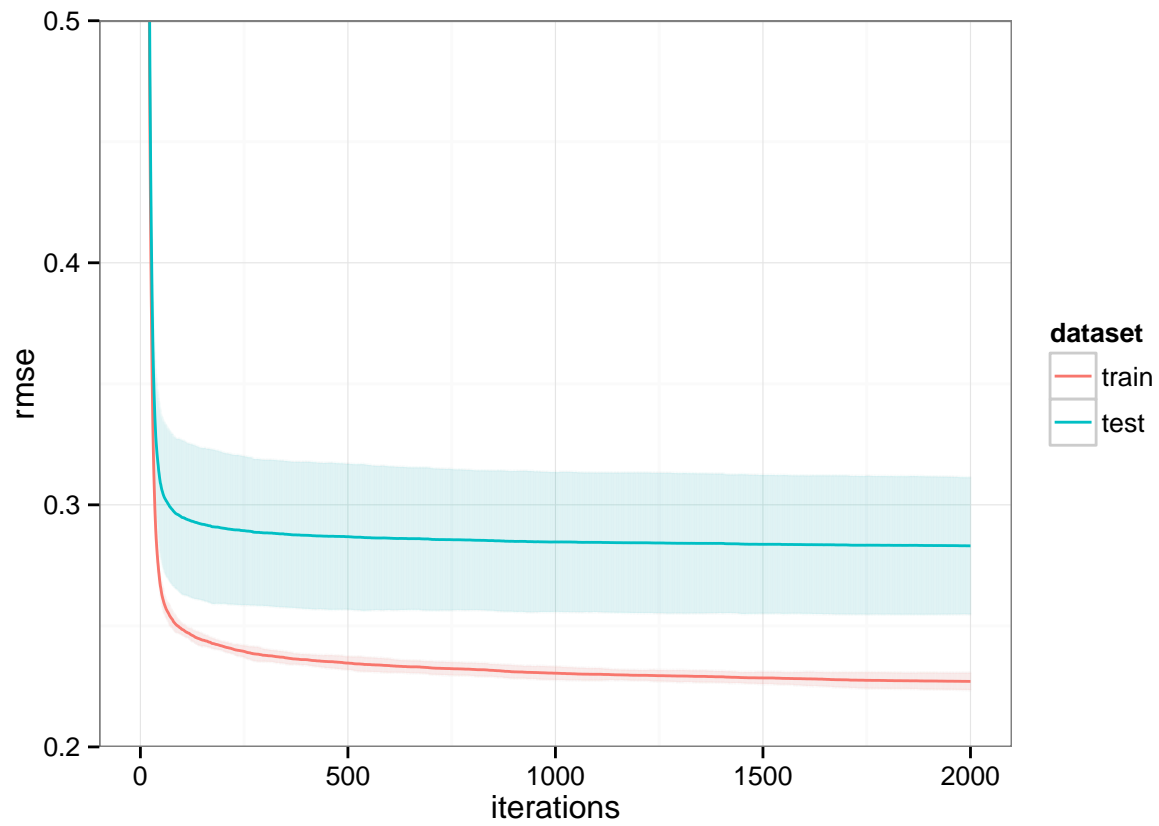
```
## [1] 0.3097785
```

Plot learning curve

```
plot_learning_curve = function(learning_curves) {
  learning_curves_train = NULL
  learning_curves_train$dataset = 'train'
  learning_curves_train$rmse = learning_curves$train.rmse.mean
  learning_curves_train$rmse.se = learning_curves$train.rmse.std
  learning_curves_train$iterations = 1:nrow(learning_curves)
  learning_curves_test = NULL
  learning_curves_test$dataset = 'test'
  learning_curves_test$rmse = learning_curves$test.rmse.mean
  learning_curves_test$rmse.se = learning_curves$test.rmse.std
  learning_curves_test$iterations = 1:nrow(learning_curves)
  learning_curves_prepared = rbind(as.data.frame(learning_curves_train),
                                   as.data.frame(learning_curves_test))

  ggplot(data = learning_curves_prepared,
        mapping = aes(x=iterations, y=rmse, group = dataset, colour=dataset)) +
    geom_errorbar(aes(ymin=rmse-2*rmse.se, ymax=rmse+2*rmse.se), width=.01, alpha=0.02) +
    geom_line() + coord_cartesian(ylim = c(0.2, 0.5))
}

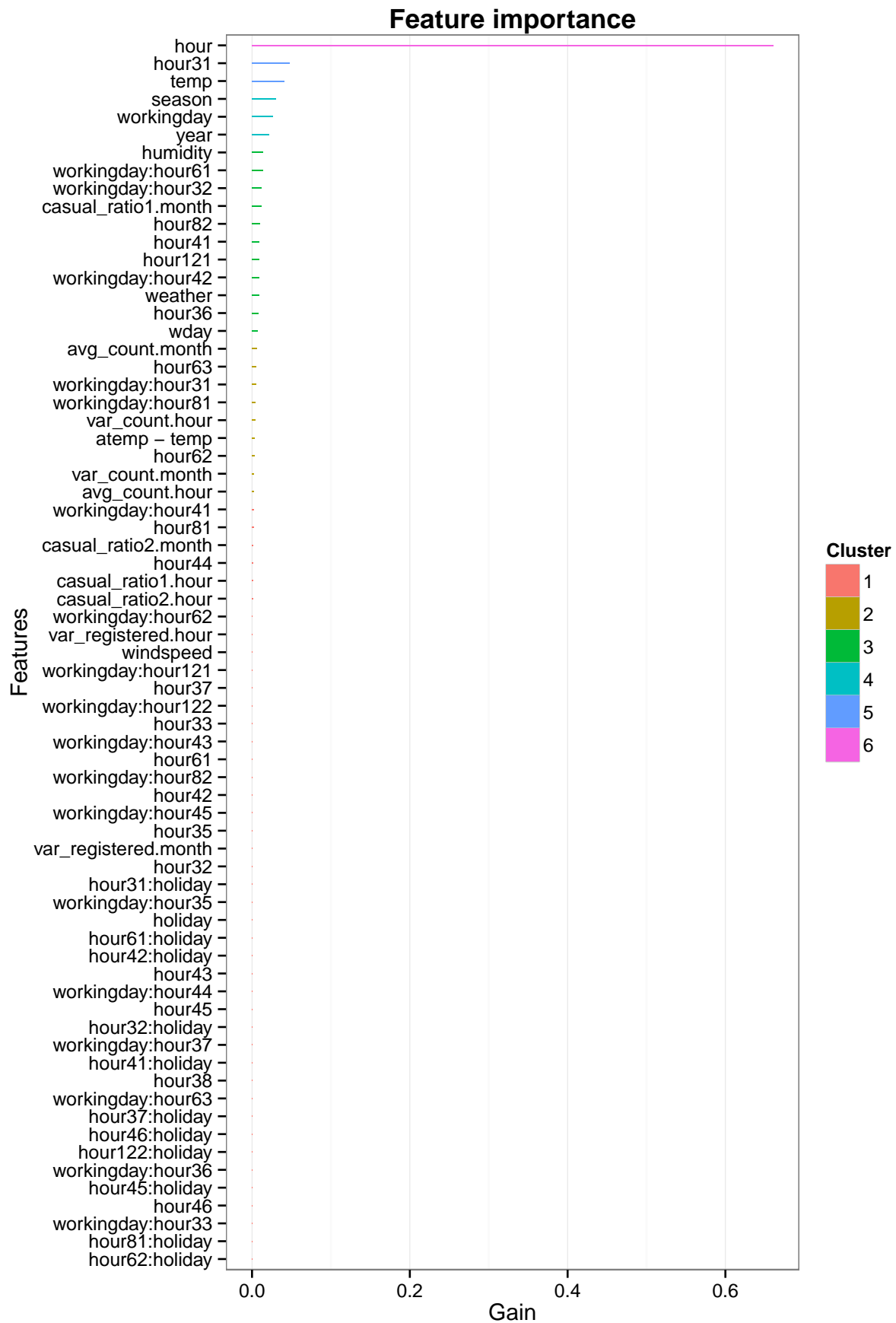
plot_learning_curve(history)
```



Feature importance

```
imp = xgb.importance(colnames(train_matrix), model = model)
```

```
xgb.plot.importance(imp)
```



Train on full dataset

```
dtrain_final <- xgb.DMatrix(full_train_matrix, label = preparedFullTarget)

model_final = xgboost(params = xgbControl, data = dtrain_final,
                      nround=2000, nthread = 4, print.every.n = 500)
```

```
## [0]  train-rmse:3.863858
## [500]  train-rmse:0.233306
## [1000]  train-rmse:0.228246
## [1500]  train-rmse:0.226401
```

```
history <- xgb.cv(params = xgbControl, data = dtrain_final,
                  nround=2000, nthread = 4, print.every.n = 500, nfold = 3)
```

```
## [0]  train-rmse:3.863931+0.004294  test-rmse:3.864195+0.009682
## [500]  train-rmse:0.238692+0.000752  test-rmse:0.287640+0.006449
## [1000]  train-rmse:0.234104+0.000233  test-rmse:0.285751+0.006612
## [1500]  train-rmse:0.231935+0.000594  test-rmse:0.285175+0.006573
```

Predict and un-log predictions

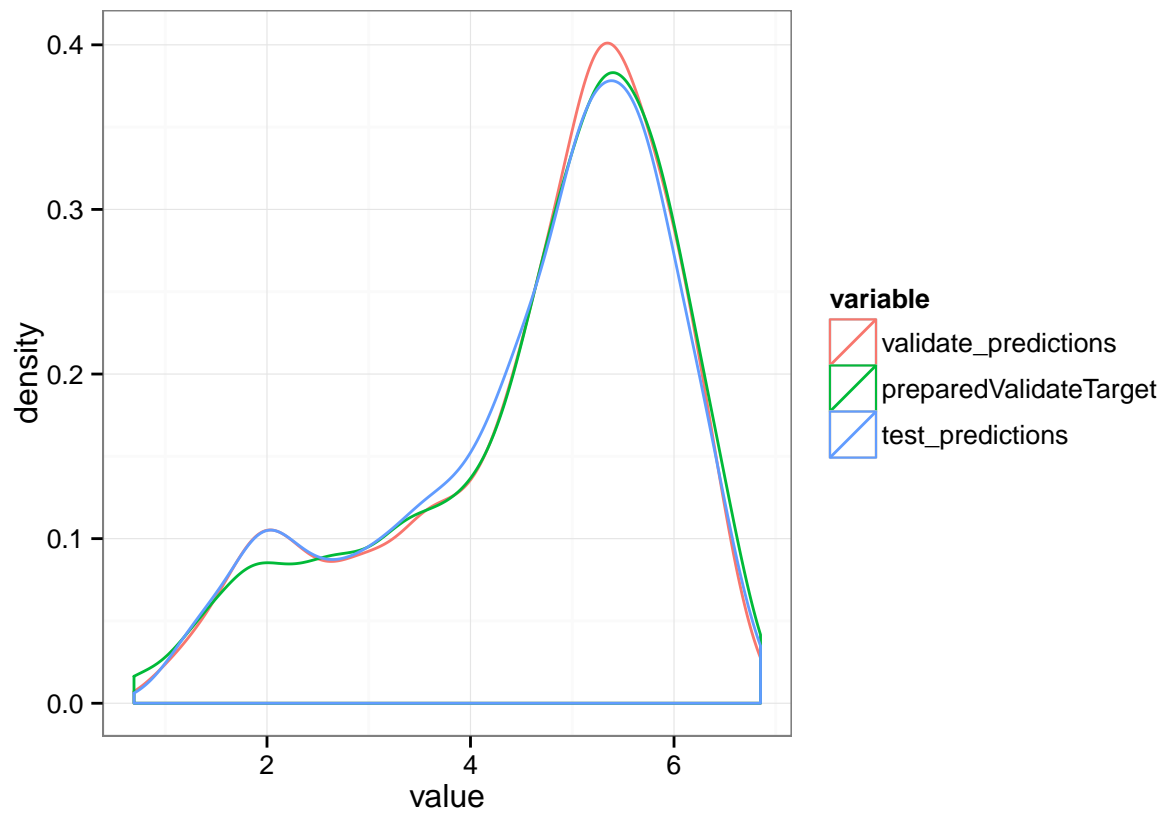
```
predictions = predict(model_final, test_matrix)
fixed_predictions = exp(predictions) - 1
```

Prediction vs Target density plot

```
test_predictions = sample(predictions, length(validate_predictions))
count_value = data.frame(validate_predictions, preparedValidateTarget, test_predictions)
count_value = melt(count_value)
```

```
## No id variables; using all as measure variables
```

```
ggplot(count_value, aes(group = variable, color = variable, x = value)) + geom_density()
```



Write the result

```
result = cbind(as.character(test$datetime), fixed_predictions) %>% as.data.frame()
names(result) = c('datetime', 'count')
write.csv(result, 'submission.csv', quote = F, row.names = F)
```