

stove - regression

2022.12.05.

Table of contents

1	Introduction	1
2	Import sample data	2
3	Data Setup Tab	2
4	Modeling Tab	3
4.1	Linear Regression	3
4.2	K Nearest Neighbor	4
4.3	Decision Tree	5
4.4	Random Forest	6
4.5	XGBoost	7
4.6	lightGBM	7
4.7	MLP	8
5	Sources for report	9
5.1	Regression plot (actual vs predicted)	9
5.2	Evaluation metrics	10

1 Introduction

- 1) 본 문서는 stove 패키지를 Shiny app에서 사용하는 것을 상정해 작성했습니다.
- 2) 본 문서의 케이스 스타일은 Camel case와 Snake case가 혼용되어 있습니다.
 - Camel case : stove의 함수명 및 파라미터명
 - Snake case: 유저로부터 받는 입력, Shiny app의 server에서 사용(될 것이라고 예상)하는 Object명, snake case로 작성된 dependencies의 함수명 등

2 Import sample data

1) 전처리가 완료된 샘플데이터를 불러옵니다.

- NA가 없어야 함
- string value가 있는 열은 factor로 변환
- 한 열이 모두 같은 값으로 채워져 있을 경우 제외해야 함
- Date type column이 없어야 함
- Outcome 변수는 classification의 경우 factor, regression의 경우 numeric이어야 함 (clustering은 outcome변수를 사용하지 않음)

```
# remotes::install_github("statgarten/datatoys")
library(stove)
library(datatoys)
library(dplyr)
set.seed(1234)
cleaned_data <- datatoys::bloodTest
cleaned_data <- cleaned_data %>%
  mutate_at(vars(SEX, ANE, IHD, STK), factor) %>%
  sample_n(1000)
```

3 Data Setup Tab

User Input		description
target_var		목적 변수
train_set_ratio	전체 데이터 중 train set의 비율 (range: 0.0 - 1.0)	

1) User input을 다음과 같이 받습니다.

- formula는 user가 target_var를 입력할 때 함께 생성되도록 함

```
target_var <- "TG"
train_set_ratio <- 0.7
seed <- 1234
formula <- paste0(target_var, " ~ .")
```

2) Train-test split 작업이 완료된 Object를 저장하고, Train set을 보여줍니다.

```
split_tmp <- stove::trainTestSplit(data = cleaned_data,
                                   target = target_var,
                                   prop = train_set_ratio,
                                   seed = seed
                                   )
data_train <- split_tmp[[1]] # train data
data_test <- split_tmp[[2]] # test data
data_split <- split_tmp[[3]] # whole data with split information
```

3) train set에 적용할 전처리 정보를 담은 recipe를 생성합니다

```
rec <- stove::prepForCV(data = data_train,
                        formula = formula,
                        seed = seed
                        )
```

4 Modeling Tab

User Input	description
mode	mode 선택(분류/회귀)
algo	사용자정의 알고리즘명
engine	알고리즘 구현 engine 선택
v	교차검증시 훈련셋을 몇 번 분할할 것인지 입력
metric	Best performance에 대한 평가지표 선택
gridNum	각 하이퍼파라미터 별로 몇 개의 그리드를 할당해 베이지안 최적화를할지 설정 (ex. 모델의 하이퍼파라미터가 3개, gridNum이 5일 때, 하이퍼파라미터 최적화를 위한 그리드는 3*5=15개)
iter	베이지안 최적화 시 반복 횟수
seed	결과 재현을 위한 시드값 설정

모델 object를 저장할 빈 리스트를 생성합니다.

```
models_list <- list()
```

4.1 Linear Regression

```

# User input
mode <- "regression"
algo <- "linearRegression"
engine <- "glmnet" # glmnet (default), liblinear, brulee
v <- 2
metric <- "rmse" # rmse (default), rsq
gridNum <- 5
iter <- 10
seed <- 1234

# Modeling
finalized <- stove::linearRegression(
  algo = algo,
  engine = engine,
  mode = mode,
  trainingData = data_train,
  splitedData = data_split,
  formula = formula,
  rec = rec,
  v = v,
  gridNum = gridNum,
  iter = iter,
  metric = metric,
  seed = seed
)
# Add the model to models_list
models_list[[paste0(algo, "_", engine)]] <- finalized$finalFittedModel

```

4.2 K Nearest Neighbor

```

# User input
mode <- "regression"
algo <- "KNN"
engine <- "kkn" # kkn (default)
v <- 2
metric <- "rmse" # rmse (default), rsq
gridNum <- 5
iter <- 10
seed <- 1234

```

```

# Modeling
finalized <- stove::KNN(
  algo = algo,
  engine = engine,
  mode = mode,
  trainingData = data_train,
  splitedData = data_split,
  formula = formula,
  rec = rec,
  v = v,
  gridNum = gridNum,
  iter = iter,
  metric = metric,
  seed = seed
)
# Add the model to models_list
models_list[[paste0(algo, "_", engine)]] <- finalized$finalFittedModel

```

4.3 Decision Tree

```

# User input
mode <- "regression"
algo <- "decisionTree"
engine <- "rpart" # rpart (default), partykit
v <- 2
metric <- "rmse" # rmse (default), rsq
gridNum <- 5
iter <- 10
seed <- 1234

# Modeling
finalized <- stove::decisionTree(
  algo = algo,
  engine = engine,
  mode = mode,
  trainingData = data_train,
  splitedData = data_split,
  formula = formula,
  rec = rec,

```

```

    v = v,
    gridNum = gridNum,
    iter = iter,
    metric = metric,
    seed = seed
)
# Add the model to models_list
models_list[[paste0(algo, "_", engine)]] <- finalized$finalFittedModel

```

4.4 Random Forest

```

# User input
mode <- "regression"
algo <- "randomForest"
engine <- "ranger" # ranger (default), randomForest, partykit
v <- 2
metric <- "rmse" # rmse (default), rsq
gridNum <- 5
iter <- 10
seed <- 1234

# Modeling
finalized <- stove::randomForest(
  algo = algo,
  engine = engine,
  mode = mode,
  trainingData = data_train,
  splitedData = data_split,
  formula = formula,
  rec = rec,
  v = v,
  gridNum = gridNum,
  iter = iter,
  metric = metric,
  seed = seed
)
# Add the model to models_list
models_list[[paste0(algo, "_", engine)]] <- finalized$finalFittedModel

```

4.5 XGBoost

```
# User input
mode <- "regression"
algo <- "XGBoost"
engine <- "xgboost" # xgboost
v <- 2
metric <- "rmse" # rmse (default), rsq
gridNum <- 5
iter <- 10
seed <- 1234

# Modeling
finalized <- stove::xgBoost(
  algo = algo,
  engine = engine,
  mode = mode,
  trainingData = data_train,
  splitedData = data_split,
  formula = formula,
  rec = rec,
  v = v,
  gridNum = gridNum,
  iter = iter,
  metric = metric,
  seed = seed
)
# Add the model to models_list
models_list[[paste0(algo, "_", engine)]] <- finalized$finalFittedModel
```

4.6 lightGBM

```
# User input
mode <- "regression"
algo <- "lightGBM"
engine <- "lightgbm" # lightgbm
v <- 2
metric <- "rmse" # rmse (default), rsq
gridNum <- 5
```

```

iter <- 10
seed <- 1234

# Modeling
finalized <- stove::lightGbm(
  algo = algo,
  engine = engine,
  mode = mode,
  trainingData = data_train,
  splitedData = data_split,
  formula = formula,
  rec = rec,
  v = v,
  gridNum = gridNum,
  iter = iter,
  metric = metric,
  seed = seed
)
# Add the model to models_list
models_list[[paste0(algo, "_", engine)]] <- finalized$finalFittedModel

```

4.7 MLP

```

# User input
mode <- "regression"
algo <- "MLP"
engine <- "nnet" # nnet
v <- 2
metric <- "rmse" # rmse (default), rsq
gridNum <- 5
iter <- 10
seed <- 1234

# Modeling
finalized <- stove::MLP(
  algo = algo,
  engine = engine,
  mode = mode,
  trainingData = data_train,
  splitedData = data_split,

```



```

    formula = formula,
    rec = rec,
    v = v,
    gridNum = gridNum,
    iter = iter,
    metric = metric,
    seed = seed
)
# Add the model to models_list
models_list[[paste0(algo, "_", engine)]] <- finalized$finalFittedModel

```

5 Sources for report

5.1 Regression plot (actual vs predicted)

유저가 선택한 모델의 confusion matrix 출력 리스트 내 모델의 이름은 {algo}_{engine}의 형태로 저장되어 있음

```

# User input
names(models_list)

```

```

[1] "linearRegression_glmnet" "KNN_kknn"
[3] "decisionTree_rpart"     "randomForest_ranger"
[5] "XGBoost_xgboost"       "lightGBM_lightgbm"
[7] "MLP_nnet"

```

```

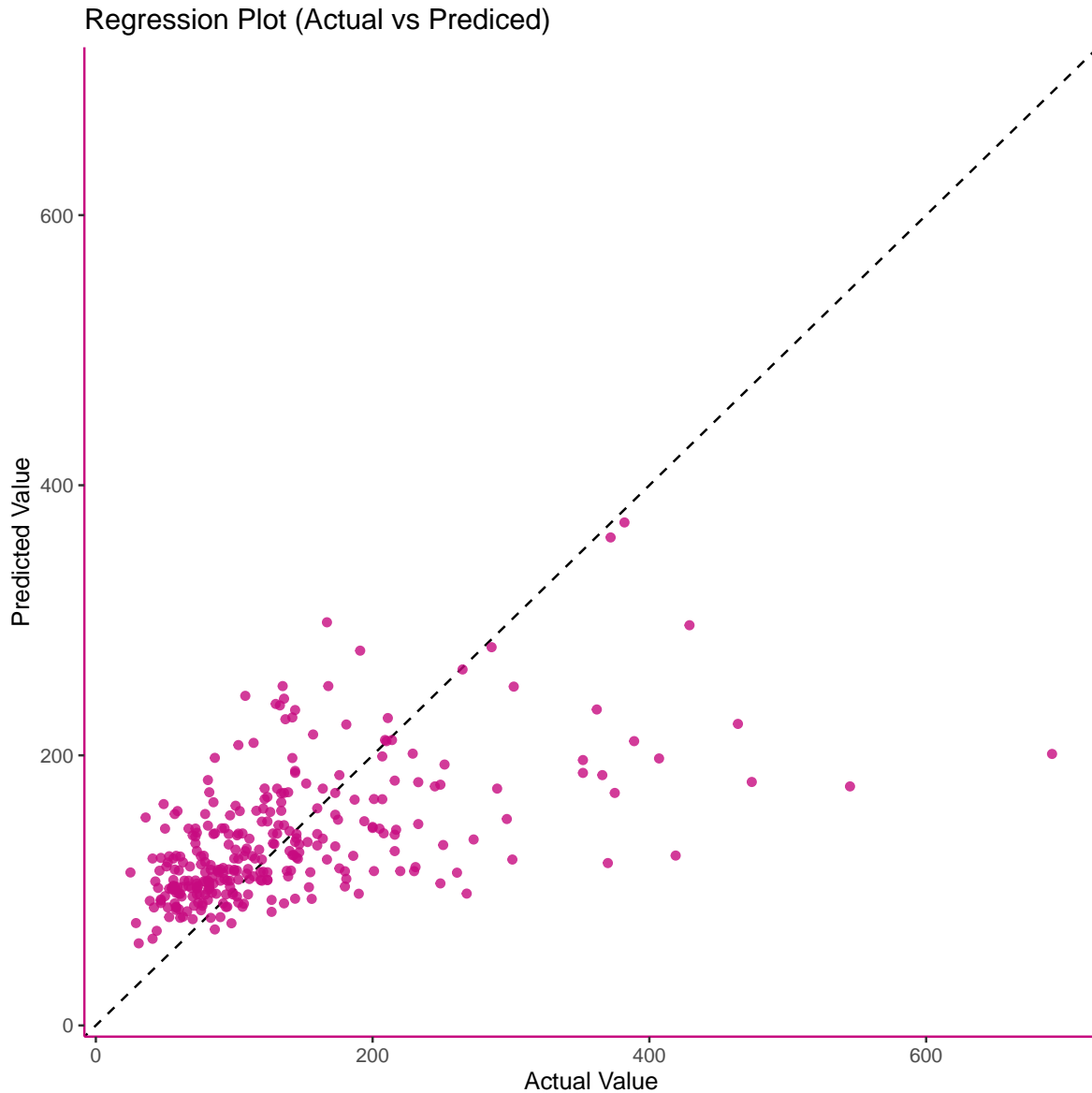
model_name <- "lightGBM_lightgbm"

```

```

rp <- stove::regressionPlot(modelName = model_name,
                             modelsList = models_list,
                             targetVar = target_var)
rp

```



5.2 Evaluation metrics

- 모델 성능 비교를 위한 표 출력

```
evalMet <- stove::evalMetricsR(models_list, target_var)
```

Warning: A correlation computation is required, but `estimate` is constant

and has 0 standard deviation, resulting in a divide by 0 error. `NA` will be returned.

```
knitr::kable(evalMet)
```

	RMSE	RSQ	MAE	MASE	RPD
linearRegression_glmnet	75.868	0.309	51.390	0.568	1.204
KNN_kknn	78.338	0.263	52.434	0.579	1.166
decisionTree_rpart	80.367	0.246	55.205	0.610	1.137
randomForest_ranger	76.063	0.311	52.144	0.576	1.201
XGBoost_xgboost	73.744	0.350	49.646	0.548	1.239
lightGBM_lightgbm	76.084	0.306	51.363	0.567	1.201
MLP_nnet	91.223	NA	62.891	0.695	1.001