

Metody Odkrywania Wiedzy

Dokumentacja końcowa projektu

„Predykcja zużycia energii na podstawie danych czujnikowych”

Krzysztof Belewicz
Paweł Pińczuk

25 stycznia 2020

1. Opis projektu

Celem projektu było wyznaczenie całkowitego zużycia energii dla zadanej chwili czasu, tzn. sumy poborów sprzętów AGD (kolumna ‘Appliances’) i oświetlenia (kolumna ‘Lights’). Zbiór danych został pozyskany z archiwum dostępnego na stronie: <https://archive.ics.uci.edu/ml/datasets/Appliances+energy+prediction>. Pojęciem docelowym jest wartość całkowitej pobieranej mocy przez gospodarstwo domowe.

TODO

2. Opis danych

2.1. Charakterystyka danych

Dane wykorzystywane do eksperymentów zostały zebrane za pomocą sieci czujników w niewielkim domu w czasie 4.5 miesiąca. Składają się z:

- daty i godziny pomiaru,
- poboru energii sprzętów domowych [Wh],
- poboru energii oświetlenia [Wh],
- pomiarów temperatury i wilgotności dla 8 różnych pomieszczeń ($[^{\circ}C]$, [%]),
- pomiarów temperatury i wilgotności dla zewnętrznej, północnej strony budynku ($[^{\circ}C]$, [%]),
- danych z pobliskiej stacji pogodowej:
 - temperatura powietrza $[^{\circ}C]$,
 - temperatura punktu rosy $[^{\circ}C]$,
 - ciśnienie atmosferyczne [$mm\ Hg$],
 - wilgotność [%],
 - prędkość wiatru [m/s],
 - widoczność [km].

2.2. Przygotowanie danych

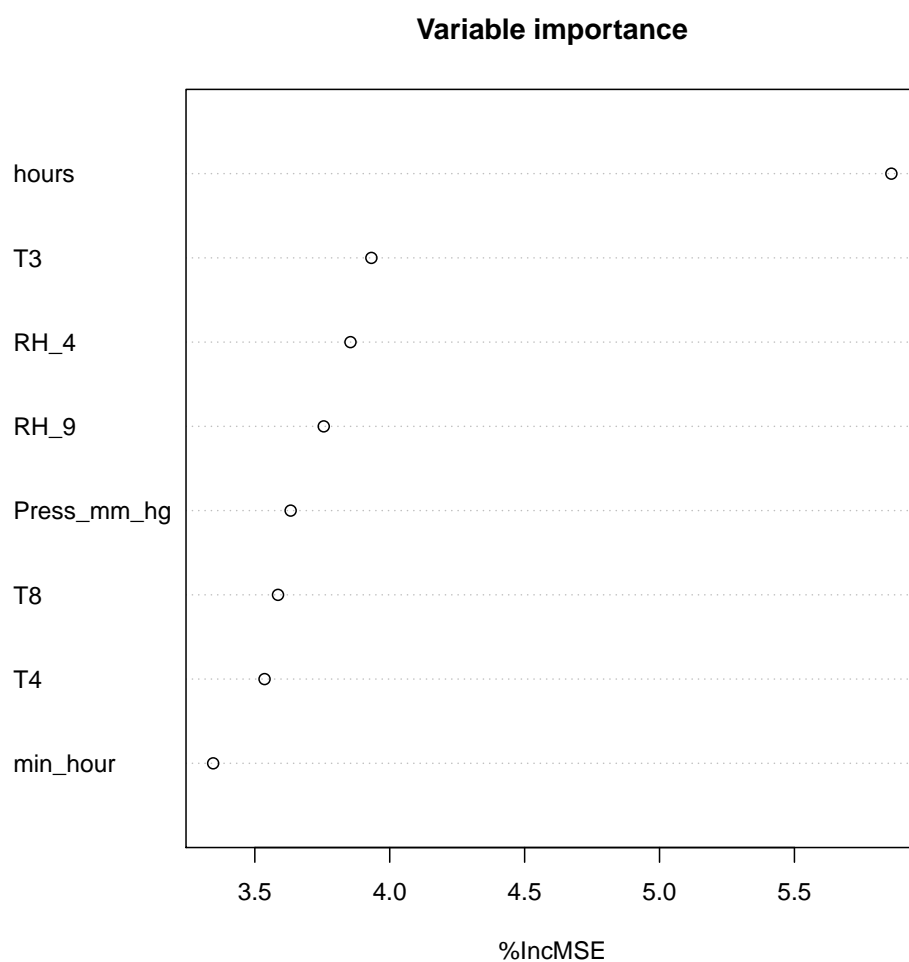
Każdy pomiar został uśredniony z 3 próbek wykonanych w równych odstępach co ok. 3,3 min. W ramach przygotowania danych, data i godzina pomiaru zostały rozdzielone na dwie oddzielne kolumny – ułatwi to późniejsze operacje na danych.

3. Opis algorytmów

TODO - ten chapter jest imo zawarty w 'konstrukcji i ocenie modeli'
ergo można usunąć, grafikę przenieść tam.

(chapter jest tutaj bo skopiowałem sprawko wstępne)

(grafika nie wiem czemu tutaj jest, tak wyszło)



4. Selekcja atrybutów

TODO jakoś ładnie opisać, to co niżej raczej dla inspiracji

Algorytm randomForest pozwala uszeregować atrybuty na dwa sposoby:

- przydatność atrybutu oceniania jest względem tego, jak bardzo model pogorszy się gdy tego atrybutu zabraknie,
- przydatność atrybutu określa zmniejszenie nieczystości po podziale w oparciu o ten atrybut.

Została wybrana pierwsza metoda ze względu na jej widoczny i bezpośredni związek z jakością całego modelu.

Najprostszym, zadowalającym rozwiązaniem jest zastosowanie pakietu `randomForest` w celu wyznaczenia predykcyjnej przydatności atrybutów, a następnie wybór (może w kilku wariantach) pewnej liczby najlepszych atrybutów (np. najlepsze 25%, 50% itp.). W tym przypadku zalecałbym użycie miary "mean decrease accuracy", co wymaga użycia w wywołaniu funkcji `randomForest` argumentu `importance=TRUE`, zaś w wywołaniu funkcji `importance` lub `varImpPlot` argumentu `type=1`. Inne bardziej zautomatyzowane podejścia można znaleźć w kilku pakietach wymienionych tutaj: <https://github.com/FrancisArgnR/R-FeatureSelection-Packages>

4.1. Wyniki selekcji atrybutów

TODO jakaś tabelka albo co + ew wnioski

5. Konstrukcja i ocena modeli

5.1. Metody konstrukcji modeli

Pojedyncze modele drzew regresji zazwyczaj cierpią z powodu wysokiej wariancji – jedną z metod jej redukcji jest tzw. Bagging (**B**ootstrap **a**ggregating). Metoda ta polega na łączeniu i uśrednianiu wielu modeli drzew, co zmniejsza wariancję i redukuje zbytne dopasowanie.

Bagging może zostać zrealizowany za pomocą pakietu *ipred* lub *caret*. *Ipred* jest z reguły prostszy w realizacji, jednakże stosowanie *caret* niesie za sobą kilka zalet. Znacznie prościej jest weryfikować krzyżowo wyniki – pomimo możliwości wykorzystania błędu OOB (Out-of-Bag) w *ipred*, weryfikacja krzyżowa daje dużo lepsze zrozumienie spodziewanego błędu. Dodatkowo możliwy jest dostęp do zmiennej odpowiedniości w wygenerowanych drzewach.

Do konstrukcji modeli klasyfikacji został wykorzystany pakiet *rpart* budujący drzewo klasyfikacji oraz pakiet *e1071*, który umożliwia waidacje skrosna oraz automatyczny dobór parametrów w celu zminimalizowania błędu.

5.2. Budowa modelu klasyfikacji z pakietu *ipred*

TODO - pakiety R, parametry (np kryteria stopu, etc)

5.3. Ocena modeli

TODO opisać że wybraliśmy np:
CC - współczynnik korelacji liniowej Pearsona

$$\frac{cov(P,A)}{var(P) \cdot var(A)}$$

TODO tutaj ładnie się wpasuje co po kolei w kodzie poszło

6. Wnioski

TODO

```
1 # main v1.1
2 #setwd("E:/Documents/Studies/MOW/MOW_project")
3 #setwd("Z:/Offtop/MOW/Project/MOW_project")
4 setwd("F:/GitHub/MOW_project")
5 setwd(".")
6
7 library(rsample)      # data splitting
8 library(dplyr)        # data wrangling
9 library(ipred)        # bagging
10 library(caret)       # bagging
11 library(dmr.regtree)
12
13 rm(list = ls())
14
15 # DATA COLLECTION AND ORGANIZATION GOES HERE
16
17 source('R/data_org.R', echo=TRUE)
18
19 # FEATURE SELECTION GOES HERE
20
21 source('R/feature_selection.R')
22 source('R/simple.filter.R')
23
24 res <- feature_selection(test_data = test_data,
25                           type = "rf",
26                           part = 0.25,
27                           trees_num=20)
28
29 # CREATING MODELS GOES HERE
30
31 # EVALUATION PROCEDURES GOES HERE
32
33 source("R/model_eval.R")
34
35 ctrl <- trainControl(method = "cv", number = 10)
36 args_t <- c(method="treebag",trControl=ctrl)
37
38 tempdataframe <- model_eval(test_data = test_data,
39                             fun = rpart,
40                             formula = as.character(res$attr_part[1]),
41                             crossval_number = 10
42                             )
43 tempdataframe2 <- model_eval(test_data = test_data,
44                             fun = lm,
45                             formula = as.character(res$attr_part[1]),
46                             crossval_number = 10
47                             )
48 tempdataframe3 <- model_eval(test_data = test_data,
49                             fun = bagging,
50                             formula = as.character(res$attr_part[1]),
51                             crossval_number = 10
52                             )
53 tempdataframe4 <- model_eval(test_data = test_data,
54                             fun = train,
55                             formula = as.character(res$attr_part[1]),
56                             crossval_number = 10,
57                             args=args_t
58                             )
```

```
1 #data_org.R
2 # data collecting and organization
3
4 training_data <- read.csv("training.csv")
5 complete_data <- read.csv("energydata_complete.csv")
6
7 month = as.numeric(substring(complete_data$date, 6, 7))
8 day = as.numeric(substring(complete_data$date, 9, 10))
9 hours = as.numeric(substring(complete_data$date, 12, 13))
10 minutes = as.numeric(substring(complete_data$date, 15, 16))
11 day_mon = day + 30 * month
12 min_hour = minutes + 60 * hours
13
14 test_data <- data.frame(
15   Appliances = complete_data$Appliances,
16   month = substring(complete_data$date, 6, 7),
17   day = substring(complete_data$date, 9, 10),
18   hours = substring(complete_data$date, 12, 13),
19   minutes = substring(complete_data$date, 15, 16),
20   day_mon = day_mon,
21   min_hour = min_hour,
22   T1 = complete_data$T1,
23   T2 = complete_data$T2,
24   T3 = complete_data$T3,
25   T4 = complete_data$T4,
26   T5 = complete_data$T5,
27   T6 = complete_data$T6,
28   T7 = complete_data$T7,
29   T8 = complete_data$T8,
30   T9 = complete_data$T9,
31   RH_1 = complete_data$RH_1,
32   RH_2 = complete_data$RH_2,
33   RH_3 = complete_data$RH_3,
34   RH_4 = complete_data$RH_4,
35   RH_5 = complete_data$RH_5,
36   RH_6 = complete_data$RH_6,
37   RH_7 = complete_data$RH_7,
38   RH_8 = complete_data$RH_8,
39   RH_9 = complete_data$RH_9,
40   T_out = complete_data$T_out,
41   RH_out = complete_data$RH_out,
42   Press_mm_hg = complete_data$Press_mm_hg,
43   Windspeed = complete_data$Windspeed,
44   Visibility = complete_data$Visibility,
45   Tdewpoint = complete_data$Tdewpoint,
46   rv1 = complete_data$rv1,
47   rv2 = complete_data$rv2
48 )
49
50 test_data <- test_data[1:1000,]
51
52 month = as.numeric(substring(training_data$date, 6, 7))
53 day = as.numeric(substring(training_data$date, 9, 10))
54 hours = as.numeric(substring(training_data$date, 12, 13))
55 minutes = as.numeric(substring(training_data$date, 15, 16))
56
57 day_mon = day + 30 * month
58 min_hour = minutes + 60 * hours
59
```

```

60 vars_train <- data.frame(
61   Appliances = training_data$Appliances,
62   month = substring(training_data$date, 6, 7),
63   day = substring(training_data$date, 9, 10),
64   hours = substring(training_data$date, 12, 13),
65   minutes = substring(training_data$date, 15, 16),
66   day_mon = day_mon,
67   min_hour = min_hour,
68   T1 = training_data$T1,
69   T2 = training_data$T2,
70   T3 = training_data$T3,
71   T4 = training_data$T4,
72   T5 = training_data$T5,
73   T6 = training_data$T6,
74   T7 = training_data$T7,
75   T8 = training_data$T8,
76   T9 = training_data$T9,
77   RH_1 = training_data$RH_1,
78   RH_2 = training_data$RH_2,
79   RH_3 = training_data$RH_3,
80   RH_4 = training_data$RH_4,
81   RH_5 = training_data$RH_5,
82   RH_6 = training_data$RH_6,
83   RH_7 = training_data$RH_7,
84   RH_8 = training_data$RH_8,
85   RH_9 = training_data$RH_9,
86   RH_out = training_data$RH_out,
87   T_out = training_data$T_out,
88   Press_mm_hg = training_data$Press_mm_hg,
89   Windspeed = training_data$Windspeed,
90   Visibility = training_data$Visibility,
91   Tdewpoint = training_data$Tdewpoint,
92   rv1 = training_data$rv1,
93   rv2 = training_data$rv2
94 )

```

Plik: feature_selection.R

```
1 # feature_selection.R
2
3 library(randomForest)
4 library(rpart)
5 library(rpart.plot) # plotting regression trees
6 library(Metrics)    # RMSE
7 library(dmr.disc)
8 library(dmr.stats)
9
10 #' @title Feature Selection
11 #'
12 #' @param test_data - data to perform feature ranking
13 #' @param type - "rf" - randomForest IMPORTANCE-based ranking (MSE) (default),
14 #' "bootstrap" - bootstrapped R2-based (coef. of determination) ranking
15 #' "simple" - simple filter based algorithms
16 #' @param part - part of attributes returned by feature selection, 0 < part <=1 (default=0)
17 #' @param trees_num - numbers of trees (randomForest) or bootstrap sets (bootstrap)
18 #'
19 #' @return formula with selected attributes
20 #'
21 #'
22 #'
23 feature_selection <- function(test_data, type="rf", part=1, trees_num=10) {
24
25   if ( !is.data.frame(test_data) ) {
26     message("Pass_data_frame_format")
27     stop()
28   }
29   if (trees_num<1) {
30     message("Trees_number_should_be_>=1")
31     stop()
32   }
33   if (part>1|part<=0) {
34     message("Parts_of_atrributes_must_be_>0_and_<=1")
35     stop()
36   }
37
38   print("FEATURE_SELECTION")
39   count <- ceiling((ncol(test_data)-1)*part)
40
41   if (type=="rf") {
42     full_model_RF <- randomForest(formula = Appliances ~ .,
43                                   data = test_data,
44                                   importance = TRUE,
45                                   ntree=trees_num)
46     importance_RF <- data.frame(importance(full_model_RF, type=1), "k"=1:(ncol(test_data)-1))
47     wyniki <- data.frame("importance"=importance_RF[order(importance_RF[,1],decreasing = TRUE)],
48                          "k"=importance_RF$k[order(importance_RF[,1],decreasing = TRUE)],
49                          "attr"=rownames(importance_RF)[order(importance_RF[,1],decreasing = TRUE)])
50
51     # plotting most important parameters
52     varImpPlot(x=full_model_RF,
53               n.var=count,
54               type=1,
55               main="Variable_importance")
56
57     # passing most important attributes from feature selection
58     count <- ceiling((ncol(test_data)-1)*part)
59     attr_part <- paste(wyniki$attr[2:count], collapse = "+")
60   }
```

```

60 attr_part <- paste("Appliances", "~", attr_part)
61 out <- data.frame(attr_part, wyniki)
62 return(out)
63 }
64 else if(type=="bootstrap"){
65
66   r2 <- function(pred.y, true.y)
67   { 1 - length(true.y)*mse(pred.y = pred.y, true.y=true.y) / ((length(true.y)-1)*var(true.y)
68
69   N=trees_num
70
71   wyniki <- data.frame(matrix(0, ncol=3, nrow=(ncol(test_data)-2)))
72
73   # bootstrap
74   for (i in 1:N)
75   {
76     wyniki_tmp <- data.frame("quality"=double(), "quality_perm"=double(), "quality_diff"
77
78     n <- sample(nrow(test_data), nrow(test_data), replace=TRUE)
79     data_train <- test_data[n,]
80     data_test <- test_data[-n,]
81
82     # tworzenie modelu na wszystkich atrybutach
83     model <- rpart(formula=Appliances~.,
84                   data=data_train,
85                   method='anova'
86
87     pred_full <- predict(model, data_test)
88
89     quality <- r2(pred.y = pred_full, true.y = data_test$Appliances)
90
91     # badanie wplywu permutacji
92     for (k in 2:(ncol(test_data)-1))
93     {
94       # permutacja k-tego atrybutu
95       data_test_perm <- data_test
96       data_test_perm[,k] <- sample(data_test_perm[,k])
97
98       # predykcja dla atrybutow o spermutowanych wartosciach
99       pred_perm <- predict(model, data_test_perm)
100
101       quality_perm <- r2(pred.y = pred_perm, true.y = data_test_perm$Appliances)
102
103       quality_diff <- quality - quality_perm
104       wyniki_tmp <- rbind(wyniki_tmp, c(quality, quality_perm, quality_diff))
105     }
106     # dodanie wynikow z i-tej iteracji
107     wyniki <- wyniki + wyniki_tmp
108   }
109
110   wyniki <- wyniki/N
111
112   wyniki <- data.frame(2:k, colnames(test_data)[2:k], wyniki[,3])
113   wyniki <- wyniki[order(wyniki[,3], decreasing = TRUE),]
114   colnames(wyniki) <- c("k", "attr", "R2_diff")
115
116   count <- ceiling((ncol(test_data)-1)*part)
117   message('Najwa?aniejsze_', part*100, '%_atrybut?w_to:\n', paste(wyniki$attr[1:count],
118   attr_part <- paste(wyniki$attr[2:count], collapse = "+")
119   attr_part <- paste("Appliances", "~", attr_part)
120

```

```

121     out <- data.frame(attr_part,wyniki)
122     return(out)
123 }
124
125 else if (type=="simple"){
126     # Symmetric uncertainty - symunc
127     # discnm.eqfreq - equal-frequency discretization
128
129     res <- simple.filter(formula = Appliances~.,
130                         data = discnm.eqfreq(~.,test_data,50),
131                         dd=symunc)
132
133     names_t <- names(res)
134     attr_part <- paste(names_t[1:count], collapse = "+")
135     attr_part <- paste("Appliances", "~", attr_part)
136     out <- data.frame(attr_part,res)
137     return(out)
138 }
139 else{
140     message("Wrong_type_of_feature_selection: \"bootstrap\" or \"rf\" or \"simple\"")
141     stop()
142 }
143 }

```

```
1
2 library(dmr.stats)
3 library(dmr.util)
4
5 dd.chi2 <- function(a1, a2) 1-chisq.test(a1, a2)$p.value
6 cd.kruskal <- function(a1, a2) 1-kruskal.test(a1, a2)$p.value
7 cc.spearman <- function(a1, a2) 1-cor.test(a1, a2, method="spearman")$p.value
8
9 #' @title Simple attribute filter
10 #'
11 #' @param formula
12 #' @param data
13 #' @param dd
14 #' @param cd
15 #' @param cc
16 #'
17 #' @return
18 #'
19 #'
20 #'
21 simple.filter <- function(formula, data, dd=dd.chi2, cd=cd.kruskal, cc=cc.spearman)
22 {
23   attributes <- x.vars(formula, data)
24   target <- y.var(formula)
25   utility <- function(a)
26   {
27     unname(switch(attr.type(data[[a]], data[[target]]),
28               dd = dd(data[[a]], data[[target]]),
29               cd = cd(data[[a]], data[[target]]),
30               dc = cd(data[[target]], data[[a]]),
31               cc = cc(data[[a]], data[[target]])))
32   }
33   sort(sapply(attributes, utility), decreasing=TRUE)
34 }
```

Plik: model_eval.R

```
1 # model evaluation
2
3 library(dmr.claseval)
4
5 #' @title Model Evaluation
6 #'
7 #' @param test_data - data to perform cross-validation
8 #' @param fun - regression algorithm (e.g. train,bagging,lm) default=rpart
9 #' @param formula - regression formula
10 #' @param crossval_number - number of cross-validations, default=10
11 #' @param args - passed arguments into regression algorithm, default=NULL
12 #'
13 #' @return data frame with measures in following order: MSE, RRSE, MAE, RMSE, RAE, COR, RSE
14 #'
15 #'
16 model_eval <-<= function(test_data, fun=rpart, formula, crossval_number=10, args=NULL) {
17
18   if ( !is.data.frame(test_data) ) {
19     message("Pass_data_frame_format")
20     stop()
21   }
22   if ( !is.character(formula) ) {
23     message("Pass_string_into_attr_parameter")
24     stop()
25   }
26   if (crossval_number<1){
27     message("Cross_validations_number_must_be_greater_than_0")
28     stop()
29   }
30   if (crossval_number==1){
31     warning("No_cross-validation")
32   }
33   temp_model <- crossval(fun,
34                           as.formula(formula),
35                           test_data,
36                           k=crossval_number,
37                           args = args)
38   MSE <- mse(temp_model$pred, temp_model$true)
39   RRSE <- rrse(temp_model$pred, temp_model$true)
40   MAE <- mae(temp_model$pred, temp_model$true)
41   RMSE <- rmse(temp_model$pred, temp_model$true)
42   RAE <- rae(temp_model$pred, temp_model$true)
43   COR <- cor(temp_model$pred, temp_model$true, method="pearson")
44   RSE <- rse(temp_model$pred, temp_model$true)
45
46   out <- data.frame(MSE, RRSE, MAE, RMSE, RAE, COR, RSE)
47   return(out)
48 }
```
