

TRABALHO DE ESTATÍSTICA APLICADA II (IAA005)

EQUIPE:

Gabriel Rulka Tardoski

Jordhan Emmanuel Marciano da Silva

Raul Cordeiro Chiarella

Zaira Mendonça Amorim

1. Sobre o trabalho

"O enunciado pedia que fizéssemos as regressões Ridge, Lasso e ElasticNet com a variável dependente 'lwage' (salário-hora da esposa em logaritmo neperiano) e todas as demais variáveis explicativas na base de dados (todas aquelas que tentam explicar o salário-hora da esposa). Com esse intuito, fizemos o seguinte processo:"

Pré-processamento dos Dados:

```
1 load("trabalhosalarios.RData")
2 library(plyr)
3 library(readr)
4 library(dplyr)
5 library(caret)
6 library(ggplot2)
7 library(repr)
8 library(glmnet)
9
10 # Padronizar os resultados
11 set.seed(75)
12
```

Importamos a biblioteca glmnet (utilizada para as regressões Ridge, Lasso e Elastic Net) e configuramos a seed para o número fixo de '75', para garantir a replicabilidade e consistência da amostragem dos dados entre cada modelo. Ainda neste passo, a coluna earns também é removida do dataset, pois não será utilizada na análise.

```
1 # Passando o dataset em uma variavel para ficar mais facil a chamada
2 dat <- trabalhosalarios
3
4 #Tirando a coluna earns pois não será utilizada
5 dat$earn <- NULL
```

2. Rotina Genérica

Existem variáveis e dados em comum, acessados e utilizados pelos modelos em questão, que, portanto, devem ser normalizados (no caso dos valores não binários) e indexados em um dataset:

```
1 # Criar um indice e particionar o dataset em 80:20
2 index = sample(1:nrow(dat),0.8*nrow(dat))
3 train = dat[index,]
```

```

4 test = dat[-index,]
5
6 # Dimensão das bases, ambas com 17 colunas
7 dim(train) # 2059 amostras
8 dim(test) # 515 amostras
9
10 # Colunas que necessitam de padronização e não são binárias exceto lwage
11 cols = c('husage', 'husearns', 'huseduc', 'hushrs',
12          'age', 'educ', 'exper', 'lwage')
13
14 # Padronizando a base de treinamento e teste
15 pre_proc_val <- preProcess(train[,cols],
16                             method = c("center", "scale"))
17
18 train[,cols] = predict(pre_proc_val, train[,cols])
19 test[,cols] = predict(pre_proc_val, test[,cols])
20

```

```

✓ > dim(train) # 2059 amostras
  [1] 5059 16

> dim(test) # 515 amostras
  [1] 515 16

```

Também se instancia uma matriz, que logo adiante será usada para predição (com base nos valores do trabalho e aplicando logaritmo neperiano):

```

1 # Criando o dataframe para predição no futuro
2 husage_norm <- (40-pre_proc_val[["mean"]][["husage"]])/
3   pre_proc_val[["std"]][["husage"]]
4 husearns_norm <- (600-pre_proc_val[["mean"]][["husearns"]])/
5   pre_proc_val[["std"]][["husearns"]]
6 huseduc_norm <- (13-pre_proc_val[["mean"]][["huseduc"]])/
7   pre_proc_val[["std"]][["huseduc"]]
8 hushrs_norm <- (40-pre_proc_val[["mean"]][["hushrs"]])/
9   pre_proc_val[["std"]][["hushrs"]]
10 age_norm <- (38-pre_proc_val[["mean"]][["age"]])/
11   pre_proc_val[["std"]][["age"]]
12 educ_norm <- (13-pre_proc_val[["mean"]][["educ"]])/
13   pre_proc_val[["std"]][["educ"]]
14 exper_norm <- (18-pre_proc_val[["mean"]][["exper"]])/
15   pre_proc_val[["std"]][["exper"]]
16
17 our_pred <- as.matrix(data.frame(husage=husage_norm,
18                                 husunion=0,
19                                 husearns=husearns_norm,
20                                 huseduc=huseduc_norm,
21                                 husblk=1,
22                                 hushisp=0,
23                                 hushrs=hushrs_norm,
24                                 kidge6=1,
25                                 age=age_norm,
26                                 black=0,
27                                 educ=educ_norm,
28                                 hispanic=1,
29                                 union=0,

```

```

30         exper=exper_norm,
31         kidlt6=1))

```

Devido à baixa legibilidade da imagem fornecida, a tabela foi manualmente transcrita para garantir precisão e clareza dos dados. Abaixo está a tabela transcrita:

husage	husunion	husearns	huseduc	husblck	hushisp	hushrs	kidpre	age	black	educ	hispanic	union	exper	kidlt6
3.582579	0	3.940915	3.243491	0	0	3.252873	1	3.378238	0	3.243491	1	0	2.93072	1

Essa transcrição foi feita para assegurar que todas as informações da imagem original sejam facilmente acessíveis e compreensíveis.

```

1 # Objeto com os valores do modelo
2 cols_reg = c('husage', 'husunion', 'husearns', 'huseduc', 'husblck',
3             'hushisp', 'hushrs', 'kidge6', 'age', 'black', 'educ',
4             'hispanic', 'union', 'exper', 'kidlt6', 'lwage')
5
6 # Estamos interessado em estimar o
7 # salário-hora da esposa em logaritmo neperiano (lwage)
8 dummies <- dummyVars(lwage~husage+husunion+husearns+huseduc+
9                      husblck+hushisp+hushrs+kidge6+age+
10                     black+educ+hispanic+union+exper+kidlt6,
11                     data = dat[,cols_reg])
12 train_dummies = predict(dummies, newdata = train[,cols_reg])
13 test_dummies = predict(dummies, newdata = test[,cols_reg])
14 print(dim(train_dummies)); print(dim(test_dummies))

```

```

✔ > print(dim(train_dummies)); print(dim(test_dummies))

1] 2059    15

[1] 515     15

```

Feito isso, são gravados os valores para aplicação posterior em todos os modelos aqui sendo aplicados, com exceção de 'lambdas', que só é utilizada em Ridge e Lasso.

```

1 # Vamos guardar a matriz de dados de treinamento das
2 # variaveis explicativas para o modelo em um objeto
3 # chamado "x"
4 x = as.matrix(train_dummies)
5
6 # Vamos guardar o vetor de dados de treinamento da
7 # variavel dependente para o modelo em um objeto
8 # chamado "y_train"
9 y_train = train$lwage
10
11 # Vamos guardar a matriz de dados de teste das variaveis

```

```

12 # explicativas para o modelo em um objeto chamado
13 # "x_test"
14 x_test = as.matrix(test_dummies)
15
16 # Vamos guardar o vetor de dados de teste da variavel
17 # dependente para o modelo em um objeto chamado "y_test"
18 y_test = test$lwage
19
20 # Lambda padrao para Ridge e Lasso
21 lambdas <- 10^seq(2, -3, by = -.1)

```

Preview de como ficou os nossos dados genéricos:

Data		
dat	2574 obs. of 16 variables	
dummies	List of 9	
our_pred	num [1, 1:15] -0.0204 0 -0.0138 -0.189...	
pre_proc_val	List of 21	
test	515 obs. of 16 variables	
test_dummies	num [1:515, 1:15] -0.727 0.181 -1.837 ...	
trabalhosalarí...	2574 obs. of 17 variables	
train	2059 obs. of 16 variables	
train_dummies	num [1:2059, 1:15] 0.787 1.191 -1.03 -...	
x	num [1:2059, 1:15] 0.787 1.191 -1.03 -...	
x_test	num [1:515, 1:15] -0.727 0.181 -1.837 ...	
Values		
age_norm	0.0169964286412261	
cols	chr [1:8] "husage" "husearns" "huseduc" ...	
cols_reg	chr [1:16] "husage" "husunion" "husearns..."	
educ_norm	-0.182288131568206	
exper_norm	-0.0410975205763276	
husage_norm	-0.0204437481969515	
husearns_norm	-0.0138088728527147	
huseduc_norm	-0.189633936181637	
hushrs_norm	-0.195349202025089	
index	int [1:2059] 1728 1192 873 759 473 2512 ...	
lambdas	num [1:51] 100 79.4 63.1 50.1 39.8 ...	
y_test	num [1:515] 0.45 0.202 -0.358 0.815 -0.2...	
y_train	num [1:2059] 0.9669 -0.0332 1.8882 -0.86...	

3. Regressão Ridge

Calculando o lambda:

```

1 # Calculando o lambda:
2 ridge_lamb <- cv.glmnet(x, y_train, alpha = 0,
3                         lambda = lambdas)
4 # Vamos ver qual o lambda otimo
5 best_lambda_ride <- ridge_lamb$lambda.min
6 best_lambda_ride

```

```

✓ > best_lambda_ride
[1] 0.06309573

```

Estimativa de coeficientes:

```

1 # Estimando o modelo Ridge
2 ridge_reg = glmnet(x, y_train, nlambda = 25, alpha = 0,

```

```

3         family = 'gaussian',
4         lambda = best_lambda_ridge)
5
6 # Vamos ver o resultado (valores) da estimativa
7 # (coeficientes)
8 ridge_reg[["beta"]]

```

```

✓ > ridge_reg[["beta"]]

15 x 1 sparse Matrix of class "dgCMatrix"

              s0
husage      0.01545592
husunion   -0.01943655
husearns    0.10500981
huseduc     0.03877441
husblck     0.02475497
hushisp     0.04068344
hushrs     -0.03665247
kid6ge     -0.07602030
age         0.01699643
black      -0.02524068
educ       0.15938866
educ       0.15938866
hispanic   -0.02511423
union      0.21152823
exper      -0.01710070
kidlt6     -0.02114730

```

Calculando o RMSE e R^2 do modelo em Ridge:

```

1 # Vamos calcular o R^2 dos valores verdadeiros e
2 # preditos conforme a seguinte funcao:
3 ridge_eval_results <- function(true, predicted, df) {
4   SSE <- sum((predicted - true)^2)
5   SST <- sum((true - mean(true))^2)
6   R_square <- 1 - SSE / SST
7   RMSE = sqrt(SSE/nrow(df))
8
9   # As metricas de performace do modelo:
10  data.frame(
11    RMSE = RMSE,
12    Rsquare = R_square
13  )
14 }
15
16 # Predicao e avaliacao nos dados de treinamento:

```

```

17 ridge_predictions_train <- predict(ridge_reg,
18                                   s = best_lambda_ridge,
19                                   newx = x)
20
21 # As metricas da base de treinamento sao:
22 ridge_eval_results(y_train, ridge_predictions_train, train)
23
24 # Predicao e avaliacao nos dados de teste:
25 ridge_predictions_test <- predict(ridge_reg,
26                                   s = best_lambda_ridge,
27                                   newx = x_test)

```

✓ > ridge_eval_results(y_train, ridge_predictions_train, train)

```

      RMSE  Rsquare
1 0.8459544 0.2840134

```

> # Predicao e avaliacao nos dados de teste:

```

> ridge_predictions_test <- predict(ridge_reg,
+                                   s = best_lambda_ridge,
+                                   newx = x_test)

```

> # As metricas da base de teste sao:

```

      RMSE  Rsquare
1 0.8692358 0.2836009

```

Fazendo a predição em Ridge:

```

1 # Fazendo a predicao:
2 predict_our_ridge <- predict(ridge_reg,
3                               s = best_lambda_ridge,
4                               newx = our_pred)
5 # Normalizando o resultado:
6 lwage_pred_ridge=(predict_our_ridge*
7                   pre_proc_val[["std"]][["lwage"]])+
8   pre_proc_val[["mean"]][["lwage"]]
9
10 # Antilog no resultado:
11 ridge_lwage <- exp(lwage_pred_ridge)

```

✓ > predict_our_ridge

```

      s1
[1,] -0.2389061

```

> # O resultado normalizado é:

```

> lwage_pred_ridge

      s1
[1,] 2.073597

```

```
> # O salário hora da esposa é (antilog logaritmo neperiano)

> ridge_lwage

      s1

[1,] 7.953378
```

Valor original predito foi -0.238, após a normalização, voltando a ficar em logaritmo neperiano, foi de \$3.24, por fim, após a aplicação do antilog, se tornou \$7.95.

O intervalo de confiança para nosso modelo ficou em 7.93 o inferior e 7.98 o superior.

```
1 # O intervalo de confianca para o nosso exemplo é:
2 n_ridge <- nrow(train) # tamanho da amostra
3 m_ridge <- ridge_lwage # valor medio predito
4 s_ridge <- pre_proc_val[["std"]][["lwage"]] # desvio padrao
5 dam_ridge <- s_ridge/sqrt(n_ridge) # distribuicao da amostragem da media
6 CIlwr_ridge <- m_ridge + (qnorm(0.025))*dam_ridge # intervalo inferior
7 CIupr_ridge <- m_ridge - (qnorm(0.025))*dam_ridge # intervalo superior
```

```
✓ > CIlwr_ridge

      s1

[1,] 7.930951

> CIupr_ridge

      s1

[1,] 7.975806
```

O que significa que o salário pode variar entre 7.93 e 7.97.

4. Regressão Lasso

Como feito na Regressão Ridge, não será necessário recriar variáveis para dataset ou para predição, pois são os mesmos valores, além do próprio lambda, porém será feito um lambda otimizado para essa regressão.

```
1 # Vamos atribuir alpha = 1 para implementar a regressao
2 # lasso
3 lasso_lamb <- cv.glmnet(x, y_train, alpha = 1,
4                       lambda = lambdas,
5                       standardize = TRUE, nfolds = 5)
6
7 # Vamos guardar o lambda "otimo" em um objeto chamado
8 # best_lambda_lasso
9 best_lambda_lasso <- lasso_lamb$lambda.min
10 best_lambda_lasso
11
12 # Vamos estimar o modelo Lasso
13 lasso_model <- glmnet(x, y_train, alpha = 1,
14                     lambda = best_lambda_lasso,
15                     standardize = TRUE)
16
17 # Vamos visualizar os coeficientes estimados
```

```
18 lasso_model[["beta"]]
```

```
✓ > best_lambda_lasso

[1] 0.01258925

> # Vamos estimar o modelo Lasso

> lasso_model <- glmnet(x, y_train, alpha = 1,
+                       lambda = best_lambda_lasso,
+                       standardize = TRUE)

> # Vamos visualizar os coeficientes estimados

> lasso_model[["beta"]]

15 x 1 sparse Matrix of class "dgCMatrix"

              s0
husage      0.022316956
husunion   -0.008868825
husearns    0.205639157
huseduc     0.052309556
husblck     .
hushisp     .
hushrs     -0.064254431
kid6ge     -0.113650167
age         0.016579810
black       .
educ        0.332126998
hispanic    .
union       0.388032081
exper       .
kidlt6     -0.003836047
```

Agora será realizado a predição do modelo Lasso e encontrar o RMSE e R^2 :

```
1 # Vamos fazer as predicoes na base de treinamento e
2 # avaliar a regressao Lasso
3 lasso_predictions_train <- predict(lasso_model,
4                                   s = best_lambda_lasso,
5                                   newx = x)
6
7 # Vamos calcular o R^2 dos valores verdadeiros e
8 # preditos conforme a seguinte funcao:
9 lasso_eval_results <- function(true, predicted, df) {
10   SSE <- sum((predicted - true)^2)
11   SST <- sum((true - mean(true))^2)
```



```

12 R_square <- 1 - SSE / SST
13 RMSE <- sqrt(SSE / nrow(df))
14 data.frame(
15     RMSE = RMSE,
16     Rsquare = R_square
17 )
18 }
19
20 # As metricas da base de treinamento sao:
21 lasso_eval_results(y_train, lasso_predictions_train, train)

```

```

✓ # As metricas da base de treinamento sao:

> lasso_eval_results(y_train, lasso_predictions_train, train)

      RMSE      Rsquare
1 0.8464903 0.2831059

# Vamos fazer as predicoes na base de teste

lasso_predictions_test <- predict(lasso_model,

                                   s = best_lambda_lasso,

                                   newx = x_test)

# As metricas da base de teste sao:

lasso_eval_results(y_test, lasso_predictions_test, test)

# As metricas da base de teste sao:

> lasso_eval_results(y_test, lasso_predictions_test, test)

      RMSE      Rsquare
1 0.870381 0.281217

```

Em comparação a Regressão Ridge, o RMSE e R^2 também apresentou resultados baixos, explicando que o modelo não está tão bem otimizado.

Vamos realizar a predição com os dados montados:

```

1 # Vamos para a predicao
2 predict_our_lasso <- predict(lasso_model,
3                             s = best_lambda_lasso,
4                             newx = our_pred)
5
6 # Novamente, o resultado esta padronizado, nos temos de
7 # converte-lo para valor compativel com o dataset original
8 lwage_pred_lasso=(predict_our_lasso*
9                  pre_proc_val[["std"]][["lwage"]])+
10 pre_proc_val[["mean"]][["lwage"]]
11
12 # Antilog no resultado:
13 lasso_lwage <- exp(lwage_pred_lasso)

```

```

✓ > # O resultado da predicao é:

> predict_our_lasso

      s1
[1,] -0.1947567

> # O salário hora da esposa é (antilog logaritmo neperiano)

> lasso_lwage

      s1
[1,] 8.137804

```

Valor original predito foi -0.19, após a normalização, voltando a ficar em logaritmo neperiano, foi de \$2.10, por fim, após a aplicação do antilog, se tornou \$8.14.

Agora calculando o intervalo de confiança:

```

1 # Vamos criar o intervalo de confianca para o nosso
2 # exemplo
3 n_lasso <- nrow(train)
4 m_lasso <- lasso_lwage
5 s_lasso <- pre_proc_val[["std"]][["lwage"]]
6 dam_lasso <- s_lasso/sqrt(n_lasso)
7 CIlwr_lasso <- m_lasso + (qnorm(0.025))*dam_lasso
8 CIupr_lasso <- m_lasso - (qnorm(0.025))*dam_lasso

```

```

✓ > CIlwr_lasso

      s1
[1,] 8.115377

> CIupr_lasso

      s1
[1,] 8.160231

```

O que significa que o salário pode variar entre 8.12 e 8.16 na Regressão Lasso.

5. Regressão ElasticNet [↗](#)

```

1 #####
2 #                REGRESSAO ELASTICNET                #
3 #####
4
5 train_cont <- trainControl(method = "repeatedcv",
6                             number = 10,
7                             repeats = 5,
8                             search = "random",
9                             verboseIter = TRUE)
10
11 # Vamos treinar o modelo
12 elastic_reg <- train(lwage~usage+union+earnings+educ+

```

```

13         husblk+hushisp+hushrs+kidge6+age+
14         black+educ+hispanic+union+exper+kidlt6,
15     data = train,
16     method = "glmnet",
17     tuneLength = 10,
18     trControl = train_cont)
19
20 # O melhor parametro alpha escolhido eh:
21 elastic_reg$bestTune
22
23 # E os parametros sao:
24 elastic_reg[["finalModel"]][["beta"]]

```

Melhor alpha escolhido:

```

✓ > # O melhor parametro alpha escolhido eh:

> elastic_reg$bestTune

   alpha   lambda

6 0.5594109 0.02112781

> # E os parametros sao:

> elastic_reg[["finalModel"]][["beta"]]

15 x 75 sparse Matrix of class "dgCMatrix"

[ suppressing 75 column names 's0', 's1', 's2', 's3', 's4', 's5', 's6', 's7', 's8', 's9', 's10', 's11', 's12',
's13', 's14', 's15', 's16', 's17', 's18', 's19', 's20', 's21', 's22', 's23', 's24',
's25', 's26', 's27', 's28', 's29', 's30', 's31', 's32', 's33', 's34', 's35', 's36', 's37', 's38', 's39',
's40', 's41', 's42', 's43', 's44', 's45', 's46', 's47', 's48', 's49', 's50', 's51', 's52',
's53', 's54', 's55', 's56', 's57', 's58', 's59', 's60', 's61', 's62', 's63', 's64', 's65', 's66', 's67',
's68', 's69', 's70', 's71', 's72', 's73', 's74']

```

Fazendo a predição e adquirindo os valores de RMSE e R²:

```

✓ > en_eval_results(y_train, en_predictions_train, train)

   RMSE  Rsquare

1 0.8464912 0.2831044

> # Vamos fazer as predicoes na base de teste

> en_predictions_test <- predict(elastic_reg, x_test)

> # As metricas de performance na base de teste sao:

> en_eval_results(y_test, en_predictions_test, test)

   RMSE  Rsquare

1 0.8704495 0.281599

```

Novamente não saíram tão bons.

```

1 # Vamos fazer a predicao com base nos parametros que
2 # selecionamos
3 predict_our_elastic <- predict(elastic_reg,our_pred)
4
5 # Novamente, o resultado eh padronizado, nos temos que
6 # reverte-lo para o nivel dos valores originais do
7 # dataset, vamos fazer isso:
8 lwage_pred_elastic=(predict_our_elastic*
9   pre_proc_val[["std"]][["lwage"]])+
10   pre_proc_val[["mean"]][["lwage"]]
11
12 elastic_lwage <- exp(lwage_pred_elastic)
13
14 # O resultado da predicao eh:
15 predict_our_elastic
16
17 # O resultado normalizado eh:
18 lwage_pred_elastic
19
20 # O salário hora da esposa é (antilog logaritmo neperiano)
21 elastic_lwage

```

Valor original predito foi -0.19, após a normalização, voltando a ficar em logaritmo neperiano, foi de \$2.10, por fim, após a aplicação do antilog, se tornou \$8.13, bem semelhante ao Lasso.

Para os intervalos de confiança:

```

1 # Vamos criar o intervalo de confianca para o nosso
2 # exemplo
3 n_elastic <- nrow(train)
4 m_elastic <- elastic_lwage
5 s_elastic <- pre_proc_val[["std"]][["lwage"]]
6 dam_elastic <- s_elastic/sqrt(n_elastic)
7 CIlwr_elastic <- m_elastic + (qnorm(0.025))*dam_elastic
8 CIupr_elastic <- m_elastic - (qnorm(0.025))*dam_elastic

```

✓ `> CIlwr_elastic`

```
[1] 8.107046
```

`> CIupr_elastic`

```
[1] 8.151901
```

```

1 # O melhor parametro alpha escolhido eh:
2 elastic_reg$bestTune
3
4 # E os parametros sao:
5 elastic_reg[["finalModel"]][["beta"]]

```

✓ `> # O melhor parametro alpha escolhido eh:`

```
> elastic_reg$bestTune
```

```
alpha    lambda
```

```
6 0.5594109 0.02112781

> # E os parametros sao:

> elastic_reg[["finalModel"]][["beta"]]

15 x 75 sparse Matrix of class "dgCMatrix"

[ suppressing 75 column names 's0', 's1', 's2', 's3', 's4', 's5', 's6', 's7', 's8', 's9', 's10', 's11', 's12',
's13', 's14', 's15', 's16', 's17', 's18', 's19', 's20', 's21', 's22', 's23', 's24',
's25', 's26', 's27', 's28', 's29', 's30', 's31', 's32', 's33', 's34', 's35', 's36', 's37', 's38', 's39',
's40', 's41', 's42', 's43', 's44', 's45', 's46', 's47', 's48', 's49', 's50', 's51', 's52',
's53', 's54', 's55', 's56', 's57', 's58', 's59', 's60', 's61', 's62', 's63', 's64', 's65', 's66', 's67',
's68', 's69', 's70', 's71', 's72', 's73', 's74']
```

Fazendo a predição e adquirindo os valores de RMSE e R²:

```
✓ > en_eval_results(y_train, en_predictions_train, train)

      RMSE      Rsquare
1 0.8464912
```

Valores entre 8.11 e 8.15.

6. Conclusão

Em resumo, os modelos ficaram da seguinte forma:

Modelo	RMSE Treino	R ² Treino	RMSE Teste	R ² Teste
Ridge Regression	0.8459544	0.2840134	0.8692358	0.2836009
Lasso Regression	0.8464903	0.2831059	0.870381	0.281712
ElasticNet	0.8464912	0.2831044	0.8704495	0.281599

De forma geral, os modelos sem o atributo 'earn' performaram bem mal, porém entre eles, o melhor foi o que usou regressão Ridge.

Outra coisa é que ele também foi o que mostrou o menor salário por hora (\$7.94/h) , dando a ideia que talvez o salário por hora correto deve ser muito menor do que o apresentado pelos modelos.