

Tipologia i cicle de vida de les dades

XAVIER JORDÀ MURRIA

ANNA SERENA LATRE

Màster Data Science UOC

MAIG 2019

Índex

	Pàg.
PRESENTACIÓ	1
BASE DE DADES	1
CARACTERÍSTIQUES DELS DATASETS	1
INTEGRACIÓ DE DADES	2
NETEJA DE DADES	2
VALIDESA DE LES DADES	2
ANÀLISI DESCRIPTIVA I INFERENCIAL	2
MODELS DE REGRESSIÓ, CLASSIFICACIÓ I PREDICCIÓ	3
PRESENTACIÓ DELS RESULTATS EN TAULES I GRÀFIQUES	3
RESOLUCIÓ DEL PROBLEMA	9
CODI	9
DOCUMENTS ANNEXOS	9
DATASETS TRACTATS	9
LINKS REPOSITORIS GITHUB	10
	10
REFERÈNCIES	10

PRESENTACIÓ

‘TITANIC: Machine Learning from Disaster.’

Disponible a: <https://www.kaggle.com/c/titanic>

En aquesta pràctica procedim a una descàrrega de bases de dades en relació a ‘TITANIC: Machine Learning from Disaster’ en el repositori kaggle.com en una proposta analítica a concurs orientada a la recerca del model predictiu que maximitzi la capacitat predictiva.

Perquè és important i a quina pregunta es pretén respondre?

La pregunta a la qual es pretén respondre en aquest exercici d’integració, neteja, validació i anàlisi de les dades és l’adequació del dataset per a la seva disposició a algorismes de classificació, regressió i predicció per tal que sigui el més eficient possible en favor d’un rendiment òptim i que posarem a prova en les tècniques d’aprenentatge supervisat train & test en una avaluació dels resultats. En aquest cas tenim en consideració que el conjunt de dades en el fitxer train ve etiquetat en la variable classe en la classificació de supervivents i ofegats en l’esdeveniment tràgic de l’enfonsament del Titànic el qual disposarem a tasques de neteja, integració i validació, procedirem a una analítica descriptiva i finalment amb el qual crearem els models predictius els resultats dels quals posarem a prova.

BASE DE DADES

Les bases de dades descarregades en el repositori ‘TITANIC: Machine Learning from Disaster’ del repositori kaggle.com són :

- Train.csv 60 KB
- Test.csv 28 KB
- gendersubmission 4 B

CARACTERÍSTIQUES DELS DATASETS

Train.csv 60KB

Dataset dimension: 891 rows x 12 columns

Test.csv 28 KB

Dataset dimension: 418 rows x 11 columns

gendersubmission.csv 4 B

Dataset dimension: 418 rows x 2 columns

Variables en les columnes:

- PassengerID : numer assignat al passatger. Tipus: numèric. [0-891]
- Survived : supervivent =1, no supervivent = 0. Tipus: numèric. [1,0]
- Pclass : Tipus de viatge, estatus socio-econòmic. Tipus: numèric. [1,2,3]

- Name : Nom del passatger. Tipus char.
- Sex : sexe dels passatgers; Masculí 'male' i Femení 'female'. Tipus: char.
- Age : edat dels passatgers. Tipus : float. [0-80]
- Sibsp : Nº de relacions de cònjuges i germans abord. Tipus : numèric. [0,1,2,3,4,5,8]
- Parch : Nº de progenitors o descendència abord. Tipus : numèric. [0,1,2,3,4,5,6]
- Ticket : Número del bitllet. Tipus char.
- Fare : Tarifa del bitllet. Tipus: float. [0-513]
- Cabin : cabina assignada al passatger. Tipus : char. [A, B,C,D,E,F] codi alfanumèric.
- Embarked : port d'embarcació. Tipus : char. ['C','Q','S']

En el dataset test la variable Survived està omesa.

INTEGRACIÓ DE DADES

Proposta d'integració de dades que s'aplica als fitxers train.csv i test.csv:

1. Fusió dels fitxers test.csv i gendersubmission.csv.
2. Supressió de les variables nom del passatger, número del bitllet i del passengeID.
3. Conversió de les variables categòriques a factorial.

NETEJA DE DADES

Tasques de neteja de dades que s'aplica als fitxers train.csv i test.csv:

1. Duplicitat en les dades.
2. Validesa de les variables: gestió dels valors NA.

Imputació de valors per KNN usant la mitja amb la informació de totes les variables numèriques, categòriques i semi-continues.

3. Valoració i tractament de la inconsistència de les dades.
4. Valors atípics: valoració d'outliers en variables quantitatives.
 - Boxplots.
 - Taula d'estimació de tendències centrals de dispersió.

VALIDESA DE LES DADES

1. Valoració de la validesa dels outliers en els intervals.
2. Valoració de la validesa de les categories.
3. Identificació de les categories en la factorització.

ANÀLISI DESCRIPTIVA I INFERENCIAL:

Arbres de classificació: Els arbres de classificació són models algorísmics de classificació que tenen un fort potencial visual per a la descripció de les tendències en les variables contínues, aquest cas Age i Fare en relació als casos d'èxit, és a dir procedim a una valoració de les tendències en l'edat i les tarifes de tiquet dels supervivents.

ANOVA: En aquesta analítica volem constatar que les ordres que es van donar en el moment de l'enfonsament del vaixell en última instància, que eren que pugen a bord dels bots d'emergència les dones i els nens en primer lloc, si es tracta d'una ordre real o bé tan sols va ser una consigna.

- ANOVA : Sex vs. Age.
- ANOVA: Sex vs. Pclass.

MODELS DE REGRESSIÓ, CLASSIFICACIÓ I PREDICCIÓ:

Machine Learning. El models proposats per a les proves Train&Test són:

Model de regressió logística:

- Entrenament del model de regressió logística.
- Valoració del model.
- Proves de predicció Train&Test per al model.

Model SVM: Machine Learning.

- Entrenament del model Suported Vectorial Machine.
- Valoració del model.
- Proves de predicció Train&Test per al model.

Una de les tècniques d'aplicació en aquest cas és la visualització del conjunt de dades train basat en els mètodes probabilístics en arbres de probabilitats o arbres de regressió.

PRESENTACIÓ DELS RESULTATS EN TAULES I GRÀFIQUES

Analítica descriptiva i inferencial: resultats en els arbres de classificació.

Resultats de l'arbre de classificació en les variables Age i Fare.

Conditional inference tree with 2 terminal nodes

Response: Age

Input: Fare

Number of observations: 889

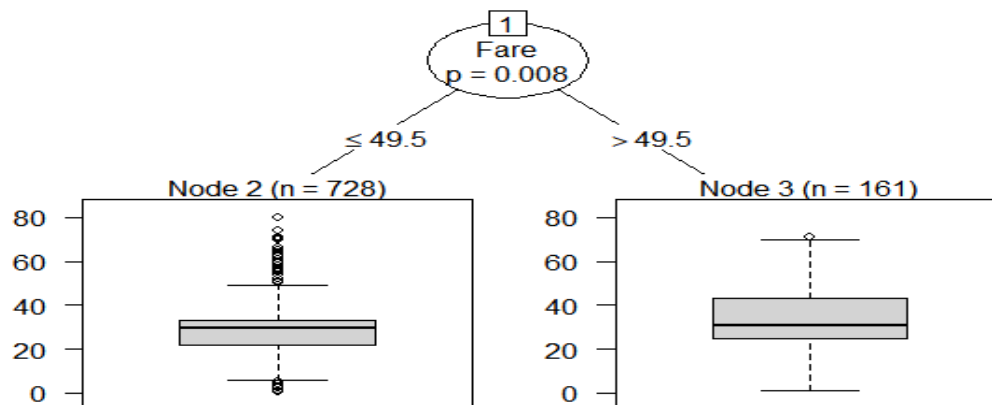
1) Fare <= 49.5; criterion = 0.992, statistic = 6.971

2)* weights = 728

1) Fare > 49.5

3)* weights = 161

Arbre de classificació en les variables Age i Fare:



Resultats de l'arbre de classificació en les variables Survived i Age.

Conditional inference tree with 2 terminal nodes

Response: Survived

Input: Age

Number of observations: 889

- 1) Age ≤ 6 ; criterion = 0.974, statistic = 4.952
- 2)* weights = 47
- 1) Age > 6
- 3)* weights = 842

Resultats de l'arbre de classificació en les variables Survived i Fare.

Conditional inference tree with 3 terminal nodes

Response: Survived

Input: Fare

Number of observations: 889

- 1) Fare ≤ 10.4625 ; criterion = 1, statistic = 57.874
- 2)* weights = 339
- 1) Fare > 10.4625
- 3) Fare ≤ 73.5 ; criterion = 1, statistic = 18.982
- 4)* weights = 455
- 3) Fare > 73.5
- 5)* weights = 95

Resultats de l'arbre de classificació en les variables Survived, Fare i Age.

Conditional inference tree with 3 terminal nodes

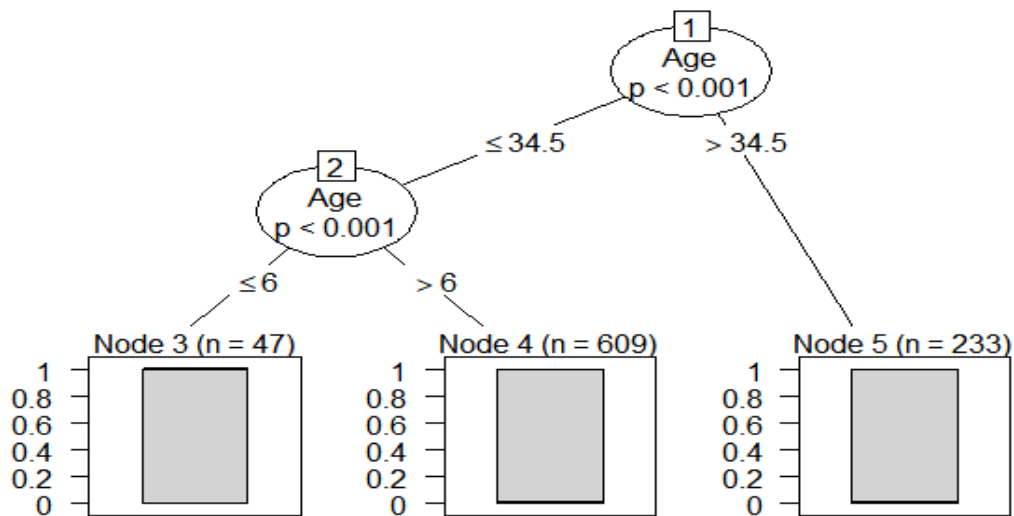
Responses: Survived, Fare

Input: Age

Number of observations: 889

- 1) Age ≤ 34.5 ; criterion = 1, statistic = 15.963
- 2) Age ≤ 6 ; criterion = 0.999, statistic = 14.589
- 3)* weights = 47
- 2) Age > 6
- 4)* weights = 609
- 1) Age > 34.5
- 5)* weights = 233

Arbre de classificació en les variables Survived, Fare i Age:



Resultats en l'arbre de classificació de les variables Survived, Age i Fare.

Conditional inference tree with 3 terminal nodes

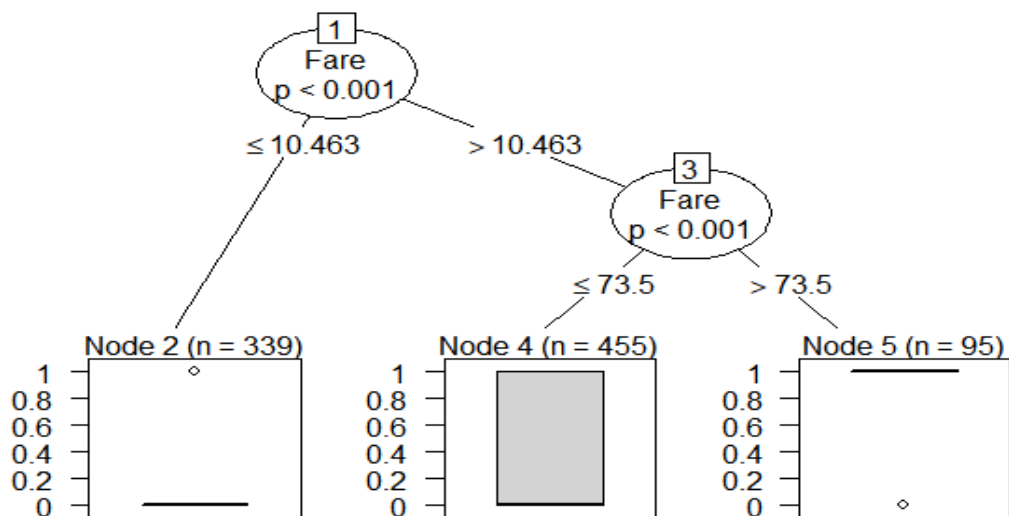
Responses: Survived, Age

Input: Fare

Number of observations: 889

- 1) Fare ≤ 10.4625 ; criterion = 1, statistic = 68.226
- 2)* weights = 339
- 1) Fare > 10.4625
- 3) Fare ≤ 73.5 ; criterion = 1, statistic = 23.771
- 4)* weights = 455
- 3) Fare > 73.5
- 5)* weights = 95

Arbre de classificació per a les variables Survived, Age i Fare.



Analítica descriptiva i inferencial: model ANOVA.

Resultats del model ANOVA en les variables Survived, Sex, Pclass i Age.

```
Call:
aov(formula = Survived ~ Sex * Pclass * Age, data = Train_T)

Terms:
              Sex      Pclass      Age Sex:Pclass  Sex:Age
Sum of Squares  61.58609  15.14444   3.34922   2.55158   0.49493
Deg. of Freedom      1         1         1         1         1

              Pclass:Age Sex:Pclass:Age Residuals
Sum of Squares    0.00631      0.06636  126.76732
Deg. of Freedom      1         1         881

Residual standard error: 0.3793287
Estimated effects may be unbalanced
```

Models de regressió, classificació i predicció: resultats en la regressió logística.

Resultats del model de regressió logística:

```
Call:
glm(formula = Survived ~ ., family = binomial(link = "logit"),
    data = Train_T)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.6446  -0.5907  -0.4230   0.6220   2.4431

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)  5.285188   0.564778   9.358 < 2e-16 ***
Pclass      -1.100058   0.143529  -7.664 1.80e-14 ***
Sexmale     -2.718695   0.200783 -13.540 < 2e-16 ***
Age         -0.039901   0.007854  -5.080 3.77e-07 ***
SibSp       -0.325777   0.109384  -2.978  0.0029 **
Parch       -0.092602   0.118708  -0.780  0.4353
Fare         0.001918   0.002376   0.807  0.4194
EmbarkedQ   -0.034076   0.381936  -0.089  0.9289
EmbarkedS   -0.418817   0.236794  -1.769  0.0769 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

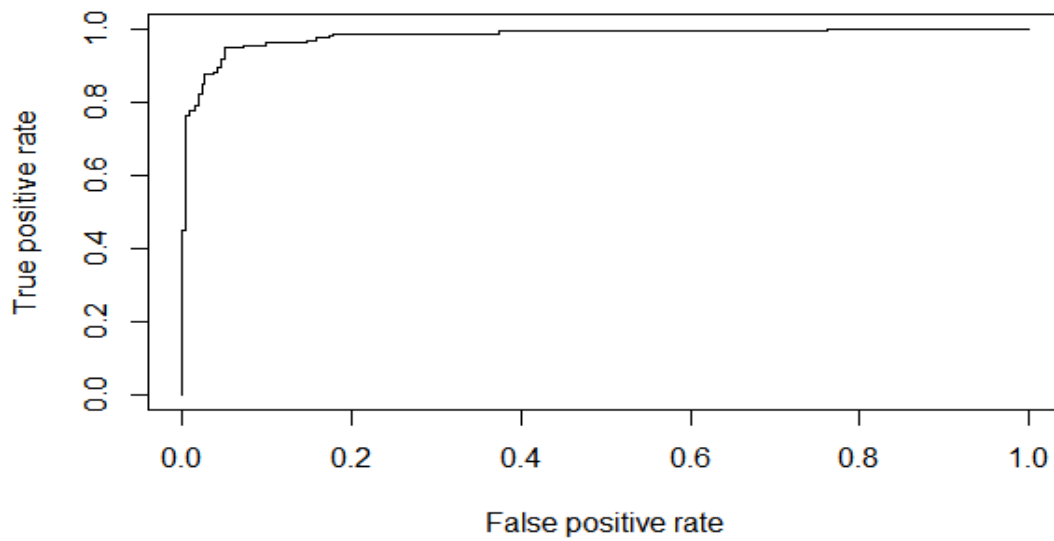
    Null deviance: 1182.82  on 888  degrees of freedom
Residual deviance:  784.19  on 880  degrees of freedom
AIC: 802.19

Number of Fisher Scoring iterations: 5
```

Precisió del model de regressió logística:

```
[1] "Accuracy 0.913669064748201"
```


Corba ROC del model de regressió logística:



Valor AUC:

```
Formal class 'performance' [package "ROCR"] with 6 slots
 ..@ x.name      : chr "None"
 ..@ y.name      : chr "Area under the ROC curve"
 ..@ alpha.name  : chr "none"
 ..@ x.values    : list()
 ..@ y.values    : List of 1
 .. ..$ : num 1
 ..@ alpha.values: list()
```

```
[1] 1
```

Models de regressió, classificació i predicció: resultats en el model SVM.

Resultats del model SVM:

Support Vector Machine object of class "ksvm"

SV type: eps-svr (regression)
parameter : epsilon = 0.1 cost C = 1

Linear (vanilla) kernel function.

Number of Support Vectors : 414

Objective Function Value : -353.5994

Training error : 0.82558

Resultats en el model predictiu SVM:

```
[,1]
1 0.04865698
2 0.95087662
3 0.04853553
4 0.04860953
5 0.95094333
6 0.04870747
```

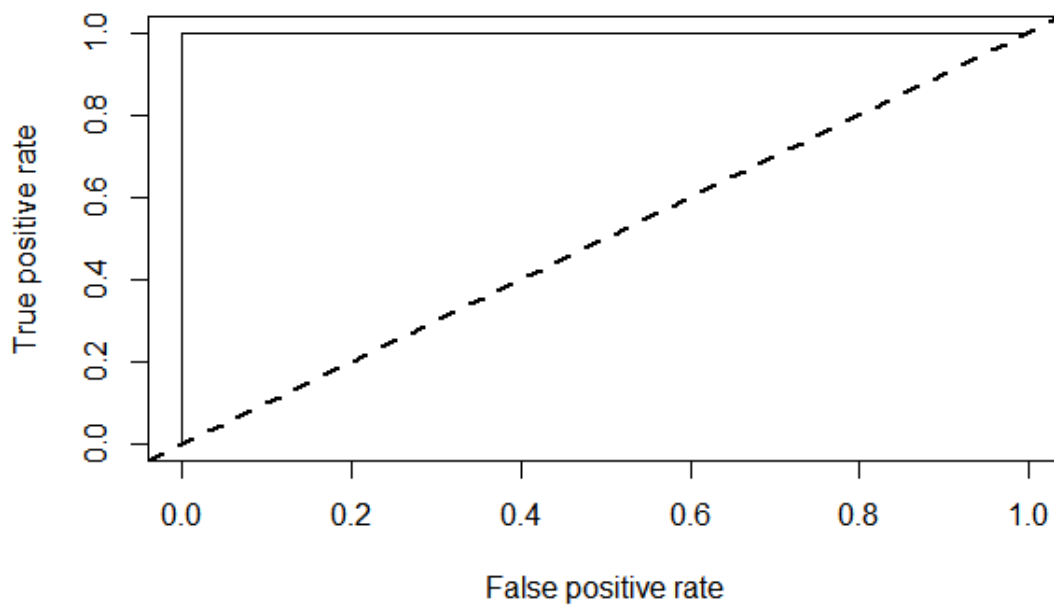
Matriu de confusió del model SVM:

```
svm.predict  0    1
            0 265    0
            1   0 152
```

Precisió del model SVM:

"Accuracy 1"

Corba ROC per a la classificació de vertaders i falsos positius:



Valor AUC del model SVM:

```
Formal class 'performance' [package "ROCR"] with 6 slots
 ..@ x.name      : chr "None"
 ..@ y.name      : chr "Area under the ROC curve"
 ..@ alpha.name  : chr "none"
 ..@ x.values    : list()
 ..@ y.values    : List of 1
 .. ..$ : num 1
 ..@ alpha.values: list()

[1] 1
```

RESOLUCIÓ DEL PROBLEMA

Conclusions.

En la fase analítica hem procedit a unes proves ANOVA per tal de valorar el succés èxit en la variable Sex, Age i Pclass.

El model de regressió logística ofereix un bon rendiment encara que no òptim i es tracta d'una modelització que pot servir de gran utilitat per a entendre i valorar la naturalesa de les variables que entren en joc i per a disposar-les per a altres algorismes com pot ser les xarxes neuronals.

El model que presenta millor rendiment és el SVM que junt amb les xarxes neuronals són els més utilitzats per el bon rendiment que ofereixen en quant a eficiència perquè són optimitzadors, si bé el SVM és un algorisme més simple però que funciona amb molta rapidesa en la realització i processament de còmputos en a màquina.

Els arbres de probabilitats són models que no aprenen en el sentit que no optimitzen el rendiment en la creació de l'algorisme però altrament resulten de gran utilitat en les tasques de classificació i regressió per a una comprensió visual de la segmentació i ordenació de les dades.

Quina és la finalitat del nostre estudi ?

La finalitat del nostre estudi es centra en la cerca d'un model òptim en quant a capacitat de representació de la variabilitat de les dades i que pugui funcionar amb eficiència en les tasques de classificació, regressió o bé predicció, en un model de regressió o en algorismes més complexos com pot ser SVM o com a funció logística activadora en un model de xarxa neuronal, ... intentarem buscar les tècniques i mecanismes que millor puguin funcionar per a la optimització del rendiment de l'algorisme en qüestió.

Podem afrontar la resolució del problema ?

La resolució del problema passa per unes tasques d'integració i neteja adequades per tal de poder procedir amb la resolució algorísmica que demanda de precisió.

CODI

DOCUMENTS ANNEXOS

Trobareu els datasets descarregats en la plataforma kaggle.com en el directori input_data.

- | | |
|------------------------------|--------|
| - Train_T.Rmd | 8 K |
| - Test_T.Rmd | 8 K |
| - Pràctiques_analítiques.Rmd | 6 KB |
| - Train_T.html | 861 KB |
| - Test_T.html | 857 KB |
| - Pràctiques analítiques | 945 KB |
| - PAC2_TCVD.pdf | 267 KB |

DATASETS TRACTATS

Trobareu els datasets resultants en el directori output_data.

- Test_T.csv 17 KB
- Train_T.csv 41 KB

LINKS REPOSITORI GITHUB

Els arxius, fitxers i documentació són disponibles en el repositori github.com:

<https://github.com/XavierJordaMurria/TipologiaPac2/>

<https://github.com/AnnaSerenaLatre/TipologiaPac2/>

REFERÈNCIES

https://rstudio-pubs-static.s3.amazonaws.com/283447_fd922429e1f0415c89b93b6da6dc1ccc.html

Taula de contribucions al treball:

Contribucions	Signa
Recerca prèvia	X.J.M, A.S.L.
Redacció de resposta	X.J.M, A.S.L.
Desenvolupament del codi	X.J.M, A.S.L.