

Memòria pràctica 8

Josep Marc Mingot//, Gabriel Reines

30 de desembre de 2012

Resum

Índex

1	Introducció	3
2	Visualització de dades	3
3	Reducció de dimensionalitat	6
4	Classificadors	6
4.1	DDA, LDA, QDA	6
4.1.1	Selecció i projecció de variables	6
4.2	k-NN	7
4.3	SVM	7
4.4	Xarxes Neuronals	7
4.5	Arbres	11
4.6	Random Forest (opcional)	11
4.7	Unió de classificadors	11
5	Conclusions	11
A	Primer Apèndix	11

1 Introducció

Explicació del problema i llista de variables a usar. Especificació del marc de treball (train-test, evaluació de l'error).

2 Visualització de dades

El primer pas en qualsevol problema de classificació és entendre les dades que se'ns han donat: relació entre les diferents variables, rànkings de quines són potencialment més importants i distribucions d'aquestes. Per aquest motiu en aquesta secció pretenem donar a través de l'anàlisi visual, una serie de propietats sobre les variables que influeixen en el problema.

Comencem estudiant la correlació entre elles. Al tenir 116 variables l'eina més adequada per visualitzar l'autocorrelació entre elles és plasmar gràficament la matriu de correlació. En la figura 1 observem la matriu on cada casella representa un valor de la matriu i com més fosc més proper a 1. Hi ha variables que estan completament correlacionades (per exemple **neighborhood intensity feature 2** amb **neighborhood intensity feature 8**). Per eliminar-les, fixem un nivell a partir del qual considerem que les variables estan correlacionades i n'eliminem aquella que té una mitja de correlació més alta amb les altres. Després d'aquest procés i fixant un tall a 0.85 passem de 116 variables a 91 variables. Podem observar la matriu d'autocorrelació neta en la figura 2.

Després de netejar les variables autocorrelacionades, fem un anàlisi exploratòria d'algunes variables. De les variables ens interessa coneixer la seva distribució condicionada a la classe, boxplots segons classe, scatter plots 2 a 2 amb altres variables (per veure possibles parelles de variables que separin). Podem obtenir tots aquests descriptius de les variables en una sola imatge. En la figura ?? observem aquests descriptius per les variables de posició (coordenades x,y i z del PE) i la mida del PE candidat (amb el nom de V5). Observem com l'*scatter* de les coordenades de posició ens dibuixen els plans de tall d'uns pulmons com calia esperar. Tanmateix observem també que no és discriminador la posició del candidat per dir si és o no PE. Respecte a la variable "mida" si que veiem que és més discriminador: els scatters amb les coordenades ens coneixen separar prou bé els 1 dels 0.

Finalment en la figura ?? podem observar els mateixos descriptius comentats anteriorment per a les 7 variables més significatives del model segons el *CAT score* (ja en donarem més detalls en l'apartat LDA). Observem per anàlisi visual que cap de les distribucions condicionades (primera columna de la imatge) s'acosta a una distribució gaussiana, però les distribucions condicionades a les classes són diferents (no s'acaba de poder apreciar ja que motrem la freqüència enlloc de la densitat). Per altra banda, els boxplots condicionats (primera fila) en mostra com efectivament ens donen distribucions diferents segons la classe. Com més diferents siguin les distribucions condicionades a la classe més senzill en serà poder separar (no en va estem mostrant aquí les nostres top 7 variables).

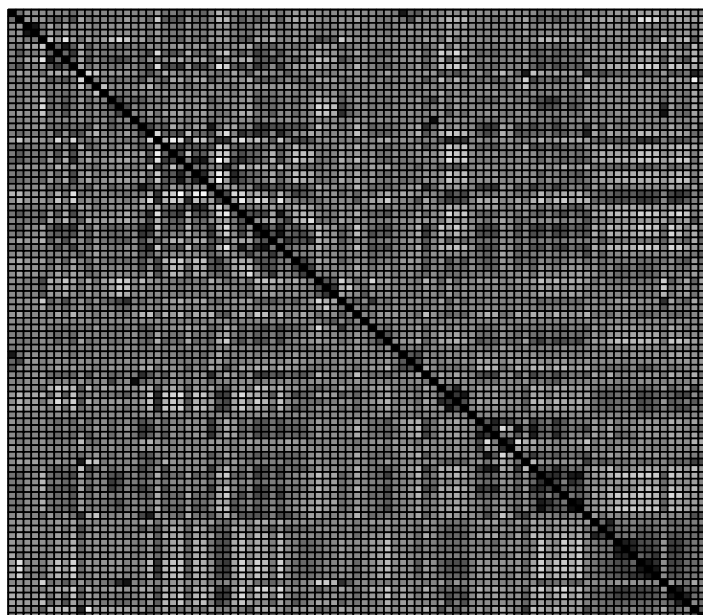


Figura 1: Autocorrelació de les 117 variables.

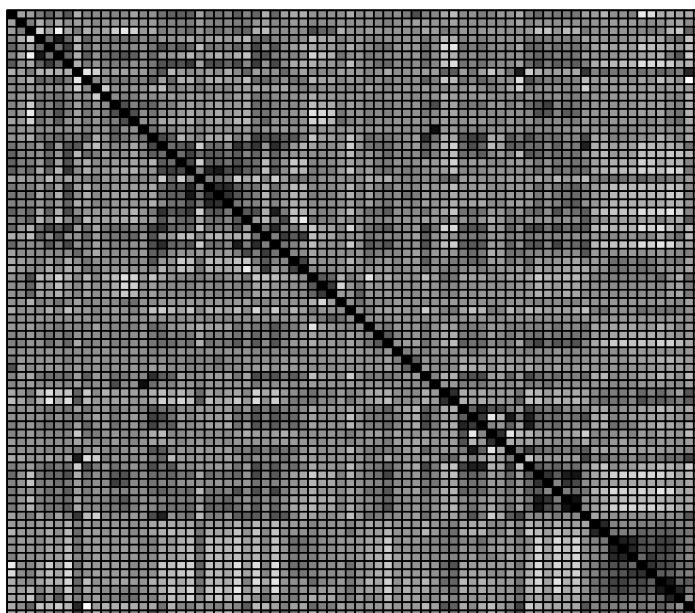


Figura 2: Autocorrelació de les 91 variables.

3 Reducció de dimensionalitat

4 Classificadors

4.1 DDA, LDA, QDA

En aquest apartat avaluarem els resultats dels classificadors bastats en discriminants. Concretament hem estudiat el DDA (*Diagonal Discriminant Analysis* o *Naive Bayes* on la matriu de covariàncies diagonals o, equivalentment, que les variables del model són independents), el LDA i QDA. Analitzem també diversos escenaris en funció de les variables d'entrada dels models. De tots ells en donem la probabilitat d'error en el conjunt de test i conjunt d'entrenament i quan és possible la curva ROC en el conjunt de test.

Comencem mostrant una comparativa entre els tres classificadors aplicats sobre les 91 variables (n'hem extret les no correlacionades tal com hem comentat en la secció de visualització de dades). En la taula 1 podem observar els errors de classificació dels tres algorismes en els conjunts de train i test. DDA, que fa la hipòtesi més forta de independència entre les variables obté el pitjor resultat en els dos conjunts. Tant LDA com QDA obtenen errors amb ordres de magnitud similars. Tanmateix QDA sembla produir una mica més de *overfitting* al tenir un error menor en el conjunt d'entrenament i un superior al de test respecte LDA. En sentit global podem dir que LDA obté millors resultats per aquest problema de QDA.

	error_train	error_test
DDA	0.18	0.22
LDA	0.09	0.10
QDA	0.07	0.11

Taula 1: Comparativa dels classificadors sobre les 91 variables no classificades

En la figura 3 hem inclòs també la curva ROC per als classificadors DDA i LDA per poder comparar-los millor.

4.1.1 Selecció i projecció de variables

En la segona part de l'anàlisi d'aquests classificadors, hem avaluat el LDA (per ser el que millors resultats ha obtingut) per diferents conjunts de variables. Pretenem aquí veure l'efecte dels classificadors en dos nous subconjunts de variables. Per una banda, usant les variables obtingudes per PCA i ICA (projecció). Per altra, usant només aquelles variables més significatives (selecció), on ja detallarem què entenem per més significatives.

PCA i ICA

La segona metodologia empleada per la selecció de variables és la de selecció de les més significatives. La selecció de variables més significatives sol ser una fase molt important en l'entrenament de classificadors per dos motius: per una banda redueix l'error de generalització i per altra disminueix el cost computacional de l'entrenament. Existeix una àmplia bibliografia sobre diferents metodologies de selecció de variables, tant seleccions genèriques (tipo filtres, tal com hem fet al eliminar les variables correlacionades) com seleccions específiques per cada classificador (*wrapper*). En el nostre cas hem empleat el marc teòric del classificador LDA per puntuar cada una de les variables mitjançant els *CAT scores* (una explicació detallada del mètode és pot trobar a [1]). En la imatge 4 podem observar les 20 variables més significatives així com la magnitud de la seva importància i el signe de les seva contribució (si el seu augment contribueix a una o una altra classe). Per tal de saber quin era el nombre de variables més adient per usar, hem fet una gràfica del error comés en el test set segons els nombre usat. Aquesta gràfica (figra 5) ens mostra que per 43 variables obtenim l'error en el test menor. A partir de llavors l'error incrementa de nou.

En la taula 2 es pot veure els resultat de tots els subconjunts de variables usats: LDA per al classificador amb les 91 variables, LDAPCA i LDAICA per als classificadors amb variables PCA i ICA respectivament i LDAselct per al classificador usant només les 43 variables més importants. [ANALISIS ICA I PCA]

	error_test
LDA	0.10
LDA_PCA	0.00
LDA_ICA	0.00
LDA_selec	0.10

Taula 2: Comparativa de LDA aplicat sobre diversos conjunts de variables

4.2 k-NN

4.3 SVM

4.4 Xarxes Neuronals

Presentem aquí els resultats de l'entrenament de les NN. En primer lloc, hem tingut en compte la dependència de les NN respecte el seus valors inicials per lo que hem canviat el marc de treball per fer-lo més robust. Hem també usat diferents mètodes per determinar els paràmetres òptims de les xarxes i presentem també els resultats de classificació i curves ROC per diferents conjunts de variables usats.

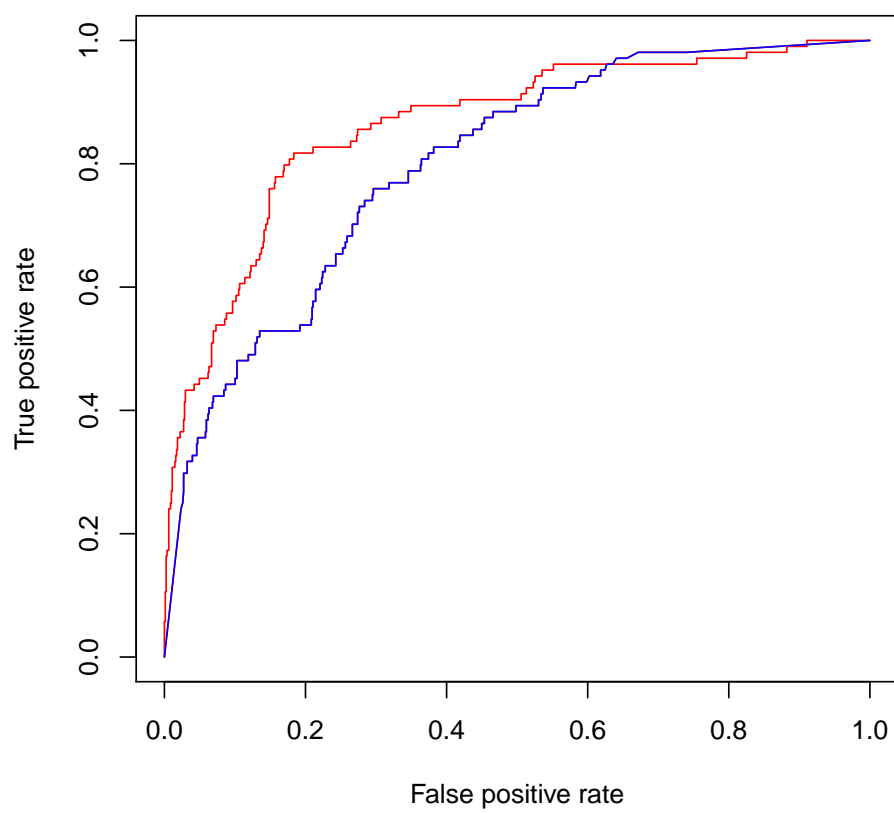


Figura 3: Curva ROC de LDA (en vermell) i DDA (en blau).

The 20 Top Ranking Features

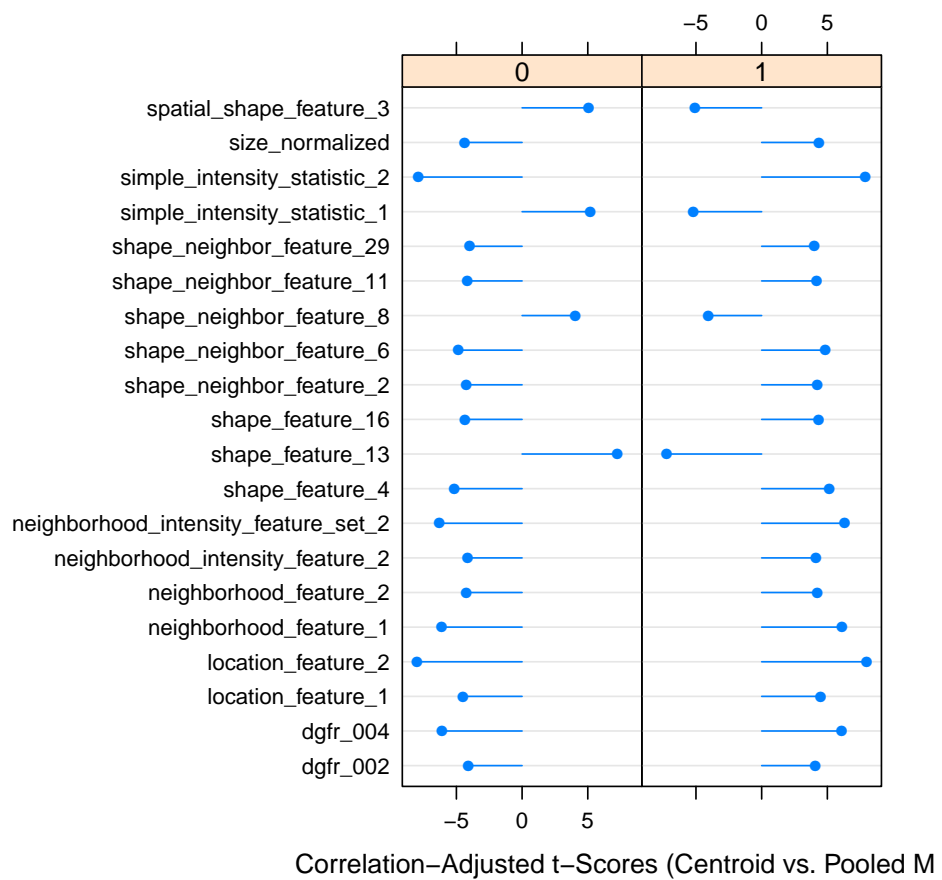


Figura 4: Ranking de les 20 variables més significatives i la seva contribució a cada una de les classes.

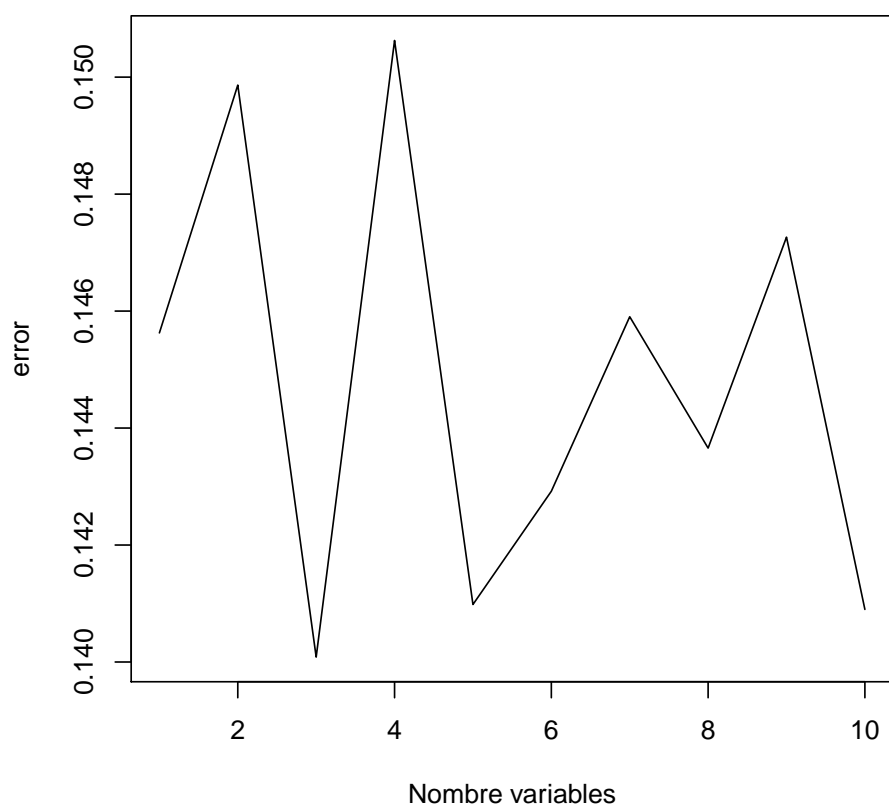


Figura 5: Error en el conjunt de test en funció del nombre de variables usades (ordenades segons importància de CAT score).

Per entrenar les xarxes neuronals hem canviat una mica el nostre marc de treball format, fins ara, per un conjunt d'entrenament i un de test. Degut a que les xarxes neuronals tenen una dependència elevada amb els seus valors inicials, hem optat per usar la tècnica de la validació creuada *k-fold*. Aquesta tècnica consisteix en partir el conjunt total de dades en k parts. Aleshores, s'entrena el classificador en $k - 1$ i s'evalua en la restant. D'aquesta manera és pot usar el mateix conjunt de dades per obtenir fins a k mesures de l'error sobre un conjunt de test. Respecte l'entrenament de les xarxes, hem seguit la filosofia general suggerida a [2]. Allí se'ns proposa entrenar les xarxes neuronals amb només una capa oculta. Per seleccionar el nombre de unitats d'aquesta capa, s'introdueix un terme de *shrinkage* que ens penalitza tenir molts paràmetres (regularització). Així doncs, la metodologia usada ha estat entrenar les xarxes amb moltes unitats imaginàries i determinar per validació creuada el valor òptim del terme de penalització.

En la figura 6 podem observar com varia l'error en el conjunt de test (mitja dels k entrenaments) en funció del terme de regularització. En totes elles estem usant 8. Després de diverses execucions probant diferents valors, obtenim que l'ordre de magnitud del valor òptim és sobre els 10^{-2} . Ens quedem doncs amb un valor de regularització de $5 \cdot 10^{-2}$. Per aquests paràmetres descrits podem observar en la figura 7 la curva ROC per les diferents realitzacions de *5-fold* validació creuada. Obtenim una mitja de 0.85 per l'AUC.

[Selecció i projecció de variables]

4.5 Arbres

4.6 Random Forest (opcional)

4.7 Unió de classificadors

5 Conclusions

A Primer Apèndix

Referències

- [1] Miika Ahdesmaki and Korbinian Strimmer, *Feature selection in omics prediction problems using CAT scores and False Nondiscovery Rate Control*, The Annals of Applied Statistics, Vol. 4, 2010.
- [2] Trevor Hastie, Robert Tibshirani, Jerome Friedman, *The elements of statistical learning*, Springer, 2nd ed., 2009.

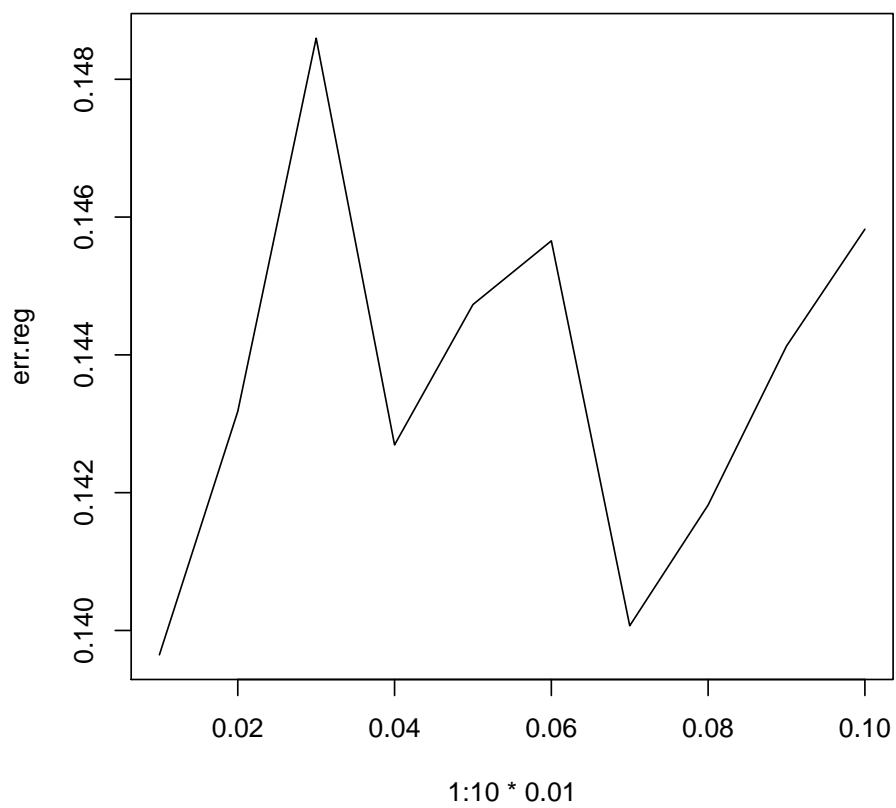


Figura 6: Error en el conjunt de test en funció del valor del paràmetre de regularització emprat en una nn d'una capa oculta i 8 unitats amagades.

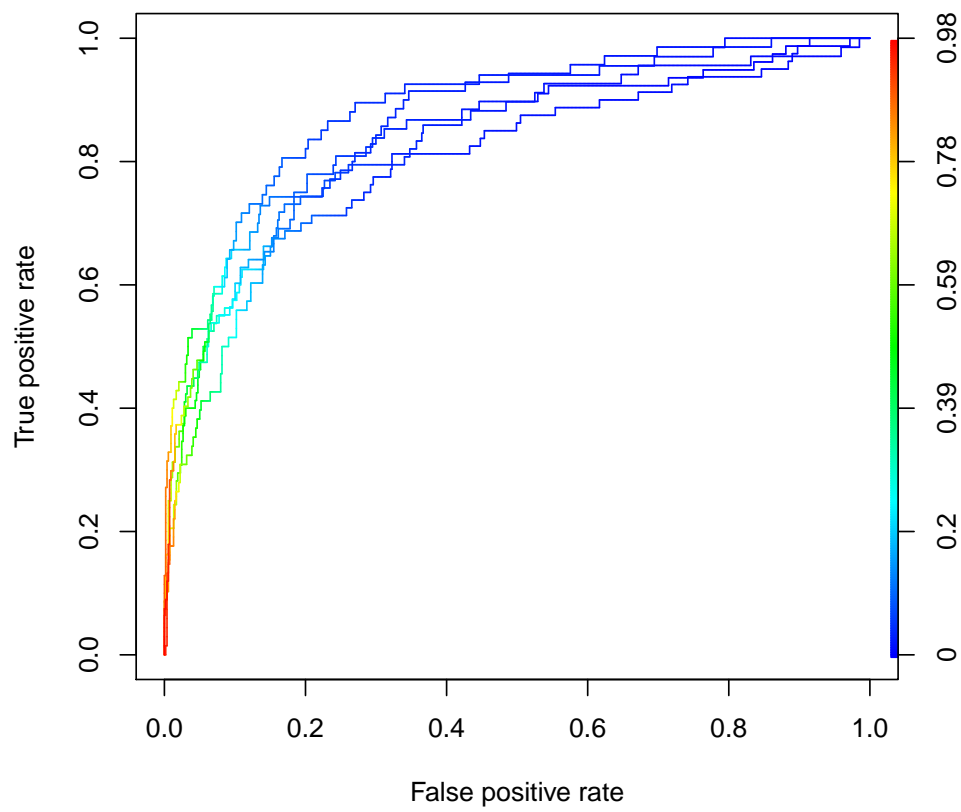


Figura 7: Curves ROC per la NN amb 8 neurones en capa ocultat, amb un terme de regularizació de 0.05. S'obté una AUC mitja de 0.85.