

PROBLEMAS NO ESTÁNDAR DE CLASIFICACIÓN SUPERVISADA

Exploración y análisis de datos

Jerónimo Hernández González

4 de octubre de 2017

Máster Universitario en Ingeniería
Computacional y Sistemas Inteligentes



1. Clasificación multi etiqueta

En este tutorial introduciremos el estudio de problemas de clasificación multi etiqueta. Para ello, usaremos el paquete **mlr** de R y el dataset Yeast:

```
1 library(mlr)
2
3 dataset <- getTaskData(yeast.task)
4 tarea <- makeMultilabelTask(id = "multi", data = dataset,
5   target = colnames(dataset)[1:14])
6 tarea
```

Código 1: Cargamos el dataset

El dominio “Yeast” está formado por datos de expresión de micro-arrays, hasta 103 variables para caracterizar un total de 2417 genes. Cada gen está asociado a un conjunto de clases funcionales (etiquetas). La función **makeMultilabelTask**, nativa del paquete **mlr**, transforma los datos facilitados a un formato propio del paquete.

Para caracterizar el dataset, atenderemos al conteo de etiquetas (número de apariciones de cada etiqueta), distribución del número de etiquetas por instancia o los cómo se agrupan habitualmente las etiquetas. Esta última medida da una impresión de la correlación entre etiquetas.

```
1 labelCount <- colSums(dataset[,1:14])
2 labelCount
3
4 nLabelsPerInst <- rowSums(dataset[, 1:14])
5 nLabelsPerInst
6 table(nLabelsPerInst)
7
8 labelsets <- do.call(paste, c(dataset[, 1:14]+0, sep = ""))
9 labelsets <- table(as.factor(labelsets))
10 labelsets
```

Código 2: Caracterizando los datos

Otras medidas tradicionales requieren cierta computación. Programaremos, por ejemplo, la cardinalidad, que sigue la siguiente fórmula:

$$card = \frac{1}{N} \sum_{i=1}^N |y_i|$$

donde N es el número de instancias (casos) del dataset, y_i el conjunto de etiquetas asociado a la instancia x_i y el operador $|\cdot|$ calcula el número de elementos en el conjunto.

```
1 # Label cardinality : Media de etiquetas por instancia
2 labelCardinality <- function(dataset) {
3   lc <- 0
```

```

-2
-3   return(lc)
-4 }
-5
-6 labelCardinality(dataset)

```

Código 3: Caracterizando los datos: cardinalidad

También podemos calcular la densidad, que sigue la siguiente fórmula:

$$card = \frac{1}{N} \sum_{i=1}^N \frac{|y_i|}{|L|}$$

donde L es el conjunto completo de posibles etiquetas.

```

1  # Label density : Media de la proporción de etiquetas por
2    instancia
3  labelDensity <- function(dataset) {
4    lc <- 0
5
6    return(lc)
7  }
8
9  labelDensity(dataset)

```

Código 4: Caracterizando los datos: densidad

A parte de los algoritmos expresamente adaptados al aprendizaje multi etiqueta, las mayorías de estrategias de aprendizaje están basadas en *ensembles* (conjuntos) de clasificadores que tratan sólo con un subproblema. Estos clasificadores, conocidos como clasificadores base, pueden ser de cualquier tipo. No obstante, suelen usarse clasificadores que devuelvan una probabilidad junto a la predicción para tener más información a la hora de combinar las diferentes predicciones. Entre los clasificadores implementados en este paquete encontramos:

- `classif.rpart` : Árboles de decision
- `classif.nnet` : Redes Neuronales
- `classif.logreg` : Regresion Logistica
- `classif.knn` : K vecinos mas cercanos

Se puede usar la función `listLearners()` para conocer todos los clasificadores base que se pueden usar. Seleccionamos así el tipo de clasificador base:

```

1  clasifBase <- makeLearner("classif.rpart", predict.type="prob")
2  clasifBase

```

Código 5: Seleccionar clasificador base

A continuación, debemos seleccionar la técnica de aprendizaje. Son varias las implementadas en este paquete. Las aquí seleccionadas son las principales:

- Binary relevance (**makeMultilabelBinaryRelevanceWrapper**): Convertir el problema multi-etiqueta (L etiquetas) en L subproblemas de clasificación binaria. Cada subproblema se resuelve independientemente con un clasificador estándar binario.
- Classifier chains (**makeMultilabelClassifierChainsWrapper**): Igualmente, tenemos L clasificadores. Difiere en que la predicción del clasificador 'l' se usa como variable predictora de los clasificadores subsiguientes ('l+1',...,L). Por lo tanto, es necesario establecer un orden entre las etiquetas y sus respectivos clasificadores. Esta técnica es capaz de modelar dependencias entre etiquetas.
- Stacking (**makeMultilabelStackingWrapper**): Igual que Dependent Binary Relevance, pero no usa tampoco las etiquetas reales en fase de aprendizaje. Usa un clasificador inicial basado en Binary Relevance para obtener una predicción de las etiquetas y con esas predicciones aprende el Dependent Binary Relevance.
- Nested stacking (**makeMultilabelNestedStackingWrapper**): Igual que Classifier Chains, pero no usa tampoco las etiquetas reales en fase de aprendizaje. Aprende el primer clasificador, obtiene una predicción de las etiquetas y con esas predicciones crea la columna de los subsiguientes clasificadores.

Una de las debilidades de este paquete es la falta de técnicas basadas en subconjuntos de etiquetas (ej., *Label Powerset* o *Random k-labelsets*). La selección de la técnica de aprendizaje se realiza de la siguiente manera:

```
1 tecnica <- makeMultilabelBinaryRelevanceWrapper(clasifBase)
2 tecnica
```

Código 6: Selección de la técnica de aprendizaje

Una vez están seleccionadas la técnica y el clasificador, podemos proceder a aprender el modelo y predecir la(s) etiqueta(s) de los casos que hemos dejado separados para test:

```
1 # Aprendizaje (con 1500 casos)
2 modelo = train(tecnica, tarea, subset = 1:1500)
3 modelo
4
5 # Predicción (con 917 casos)
6 prediccion <- predict(modelo, newdata = dataset[1501:nrow(
7   dataset),])
8 t(prediccion$data)
```

Código 7: Aprendizaje y clasificación

NOTA: Podemos intentar usar otros clasificadores base y técnicas de aprendizaje.

A la hora de evaluar un clasificador aprendido para un problema multi etiqueta, se pueden elegir entre dos estrategias: evaluar a nivel global o a nivel de etiqueta. Entre las medidas de evaluación específicas de clasificación multi etiqueta podemos elegir, usando la función **performance**, las siguientes:

- *multilabel.subset01*: Proporción de casos donde se acierta
- *multilabel.hamloss*: Hamming Loss - Proporción de etiquetas fuera de la intersección con respecto al total de etiquetas
- *multilabel.acc*: Media de la intersección entre la union de los sets de etiquetas reales y predicho (jaccard)
- *multilabel.f1*: Micro-F1 - Media armonica de precision y recall

La lista completa de medidas específicas de evaluación puede ser consultada usando la función **listMeasures("multilabel")**. Al evaluar a nivel de etiqueta podemos usar medidas “genéricas” de clasificación supervisada aplicadas al acierto en cada etiqueta individualmente. Entre este tipo de medidas de evaluación podemos elegir, usando la función **getMultilabelBinaryPerformances**, las siguientes:

- *acc*: Porcentaje de acierto
- *mmce*: Error de clasificación
- *tpr*: Recall
- *ppv*: Precisión
- *f1*: Media armónica de precisión y recall

La lista completa de medidas de evaluación puede ser consultada usando la función **listMeasures()**. La evaluación se realiza de la siguiente manera:

```
1 # Evaluación específica multi etiqueta
2 performance(prediccion, measures = list(multilabel.subset01,
3     multilabel.hamloss, multilabel.acc, multilabel.f1))
4
5 # Evaluación tradicional aplicada a los clasificadores base
6 # (independientemente por etiqueta)
7 getMultilabelBinaryPerformances(prediccion, measures=list(acc,
8     mmce, tpr, ppv, f1))
```

Código 8: Cálculo de varias distancia entre dos series

2. EM

Como práctica final, vamos a programar el algoritmo EM para aprender un sencillo Naive Bayes. Empezaremos, entonces, por programar este clasificador. Guardaremos en una lista el número de variables predictoras y los parámetros que usa: $p(c = 0)$, $p(c = 1)$, $p(x_i = 0|c = 0)$, $p(x_i = 1|c = 0)$, $p(x_i = 0|c = 1)$ y $p(x_i = 1|c = 1)$ (para todas las variables predictoras, i):

```
1 naiveBayes <- function(nVarPreds) {  
0   modelo <- list()  
-1  
-2   modelo$nPredictors <- nVarPreds  
-3  
-4   modelo$Pc <- vector("numeric",length=2)  
-5  
-6   modelo$Px_c <- matrix(0,nVarPreds,2*2)  
-7  
-8   return(modelo)  
-9 }
```

Código 9: Inicializar Naive Bayes

Una vez diseñado, vamos a lo importante. Un clasificador ha de poder ser aprendido y usado para predecir. Por ello, diseñaremos una función para cada procedimiento. La función de predicción sería así:

```
1 predecir <- function(instancia, modelo) {  
0   probs <- vector("numeric",length=2)  
-1  
-2   probs[1] <- modelo$Pc[1]  
-3   probs[2] <- modelo$Pc[2]  
-4  
-5   for (i in 1:modelo$nPredictors) {  
-6     probs[1] <- probs[1]*modelo$Px_c[i,1+instancia[i]]  
-7     probs[2] <- probs[2]*modelo$Px_c[i,3+instancia[i]]  
-8   }  
-9  
-10  return(probs/sum(probs))  
-11 }
```

Código 10: Inicializar Naive Bayes

Donde se calculan las siguiente probabilidades:

$$p(c = 0) \prod_{i=1}^v p(x_i|c = 0)$$

$$p(c = 1) \prod_{i=1}^v p(x_i|c = 1)$$

con v siendo el número de variables predictoras (2 en este caso). Finalmente, normalizamos para que la suma de $p(c = 0|x)$ y $p(c = 1|x)$ sea 1.

Considerando la estructura definida para guardar los parámetros y lo visto en clase sobre aprendizaje de parámetros, diseñaremos la siguiente función:

```

1  aprender <- function(dataset, modelo) {
2
3      modelo$Pc <- vector("numeric",length=2)
4      modelo$Px_c <- matrix(0,modelo$nPredictors,2*2)
5
6      #####
7      ## Tu código aquí ##
8      #####
9
10     return(modelo)
11 }

```

Código 11: Inicializar Naive Bayes

Por último, el esqueleto del algoritmo EM es el siguiente:

```

1  EM <- function(datos, epsilon) {
2      cDataset <- inicializar(datos)
3
4
5      modelo <- aprender(cDataset,modelo)
6
7      convergencia <- FALSE
8      while (!convergencia) {
9          cDataset <- EStep(cDataset,modelo)
10
11          antModelo <- modelo
12          modelo <- MStep(cDataset, antModelo)
13
14          convergencia <- testConvergencia(modelo,antModelo,epsilon)
15      }
16
17      return(list(modelo,cDataset))
18 }

```

Código 12: Esqueleto EM

Faltarían por programar las tres funciones fundamentales del EM: **EStep**, **MStep** y **testConvergencia**. Mediremos la convergencia atendiendo al siguiente criterio: la distancia euclídea entre los parámetros del paso anterior y del actual es menor que *epsilon* o no.

```

1  EStep <- function(dataset, modelo){
2
3      dataset$peso <- vector("numeric",length=length(dataset$peso))
4
5      #####
6      ## Tu código aquí ##
7      #####
8
9
10 }

```

```

-6      return(dataset)
-7    }
-8
-9
-10
-11 MStep <- function(dataset, modelo){
-12
-13     #####
-14     ## Tu código aquí ##
-15     #####
-16
-17     return(modelo)
-18 }
-19
-20 testConvergencia <- function(modeloA, modeloB, epsilon) {
-21
-22     resultado <- FALSE
-23
-24     #####
-25     ## Tu código aquí ##
-26     #####
-27
-28     return ( resultado )
-29 }

```

Código 13: Ejecución de la práctica

Con el objetivo de hacerlo comparativo, todos ejecutaremos el algoritmo realizando la siguiente llamada:

```

1  # Inicializamos el model con 2 variables predictoras
0  modelo <- naiveBayes(2)
-1  # Simulamos un dataset
-2  datos <- cbind(c(0,1,0,0,1,1,1,1,0,1),
-3                c(0,0,1,1,0,0,0,1,1,1),
-4                c(1,NA,NA,NA,1,0,NA,NA,0,NA))
-5
-6  EM(datos, 0.001)

```

Código 14: Ejecución de la práctica