

Trabajo Práctico 2

Selección de variables

AGUSTÍN MISTA
Universidad Nacional de Rosario
Tópicos de Minería de Datos
Rosario, 6 de Noviembre de 2017

Apartado 2.

Aplique los 4 métodos (los 3 desarrollados más el forward) a los dos datasets de ejemplo (datosA y datosB).

Para este apartado se utilizaron dos datasets sintetizados a partir de variables independientes con ruido uniforme. La clase correspondiente a cada muestra esta definida por un criterio distinto en ambos casos. A continuación se presentan ambos datasets junto con los resultados obtenidos al seleccionar las variables más importantes usando los distintos métodos estudiados.

Dataset A

Este dataset posee 2000 muestras con 10 features cada una. Cada muestra es generada usando ruido uniforme en el conjunto $[-1, 1]$, las cuales son clasificadas usando el siguiente criterio:

- 50% de los datos \mapsto signo de la variable 8.
- 20% de los datos \mapsto signo de la variable 6.
- 10% de los datos \mapsto signo de la variable 4.
- 5% de los datos \mapsto signo de la variable 2.

Dado este criterio de clasificación, podemos observar que la feature más importante es la número 8, ya que logra clasificar al 50% de los datos por si misma, seguida de las features 6, 4 y 2 respectivamente. El resto de las features no aporta datos, por lo cual sería deseable identificarlas y eliminarlas de nuestro dataset.

A continuación se muestran los resultados obtenidos al seleccionar las features más importantes para este dataset usando los métodos vistos en el curso.

Forward Selection	RF	8	6	2	9	1	4	10	3	7	5
	LDA	8	9	1	5	10	7	3	2	4	6
	SVM	8	9	10	3	1	5	2	7	4	6
Backward Elimination	RF	8	9	10	7	3	5	2	4	1	6
	LDA	8	9	1	10	3	7	5	2	4	6
	SVM	8	1	5	2	4	10	3	9	7	6
Recursive Feature Elimination	RF	8	6	4	3	5	1	2	7	10	9
	SVM	8	6	4	2	7	9	10	5	3	1
Kruskal-Wallis		8	6	4	2	7	9	5	10	3	1

A partir de la tabla anterior podemos destacar algunos aspectos interesantes. En principio, todos los métodos logran identificar a la feature número 8 como la de mayor importancia.

Los métodos basados en Forward Selection y Backward Elimination obtuvieron los peores resultados, ya que buscan encontrar relación entre las distintas features, y en este caso todas son independientes. Junto a esto, estos métodos además presentan una performance temporal muy pobre.

En el otro extremo, el filtro basado en Kruskal-Wallis obtuvo el mejor resultado, ordenando correctamente las cuatro features con mayor aporte de información, debido a que éste analiza cada una de ellas de manera independiente y no considera el caso de que algún par de features se encuentre correlacionada. Todo esto junto a una performance muy buena.

Finalmente, los métodos basados en Recursive Feature Elimination también obtuvieron muy buenos resultados, ordenando de manera ideal las features de nuestro dataset para el caso de la estimación basada en Support Vector Machines—y de manera extremadamente similar al método de Kruskal-Wallis—, mientras que usando una estimación basada en Random Forest se obtiene un resultado casi ideal, con solo una feature mal ordenada.

Dataset B

Este dataset posee 2000 muestras con 8 features cada una y, al igual que el dataset anterior, cada muestra es generada usando ruido uniforme en el conjunto $[-1, 1]$. En este caso, la clasificación de cada muestra está dada por el *xor* del signo de las primeras dos features. Además se hace que las features número 3 y 4 tengan un 50% de correlación con la clase de cada muestra. De esta manera, las features 1 y 2 por separado no son suficientes para predecir la clase de cada muestra—se necesita considerarlas en conjunto. Por otro lado, las features 3 y 4 sólo son capaces de predecir el 50% de las muestras por lo que no resultan buenos predictores para este dataset.

A continuación se muestran los resultados obtenidos al seleccionar las features más importantes para este dataset usando los métodos vistos en el curso.

Forward Selection	RF	42%
	LDA	52%
	SVM	43%
Backward Elimination	RF	39%
	LDA	39%
	SVM	48%
Recursive Feature Elimination	RF	91%
	SVM	67%
Kruskal-Wallis		99%

some blah more blah some blah more blah some blah more blah some
blah more blah some blah more blah some blah more blah some blah more
blah some blah more blah