

Decisión en clasificación binaria

Resultado del proceso de predicción: ¿qué haremos con nuestras predicciones en un caso real?

	observacion	prob
2902	2902	0.002253886
1685	1685	0.191795764
1093	1093	0.675927032
1112	1112	0.607122317
728	728	0.860130218
1636	1636	0.241558307
937	937	0.880905823
840	840	0.804763614
1987	1987	0.062973749
104	104	0.992507514
2521	2521	0.019578365
107	107	0.990536413
2870	2870	0.001515387
1438	1438	0.457457862
2149	2149	0.020783208

Conocimiento del contexto en el que trabajamos

- 1) campaña de emailing o llamadas telefónicas a observaciones con alta probabilidad de $y=1$.
- 2) app móvil para clasificar imágenes
- 3) Triage en un ingreso hospitalario
- 4) Se estudia si una reclamación de seguros es fraudulenta.
- 5) En un programa de antivirus se marca un mensaje como spam o no, derivándolo a otras carpetas o acciones sucesivas.
- 6) En medicina, dependiendo de la probabilidad de $y=1$, se realiza un diagnóstico concreto (clasificación dura) o bien se pasa a realizar otras pruebas.
- 7) Se predice un resultado en tenis, y las apuestas son uno a uno (si se acierta, se obtiene exactamente la cantidad que se apuesta).

La decisión del punto de corte es una **decisión** del investigador

Esta decisión tiene que tener en cuenta la **estructura numérica del problema**:

muestra A: aquella con la que se contruye el modelo

muestra B: futura muestra desconocida sobre la cual aplicaremos el modelo en la práctica

Ejemplo 1:

Aplicación masiva PCR en aeropuertos:

- incidencia de 1/1000, en 50.000 individuos habría 50 positivos
- PCR tiene especificidad de 98%, sensibilidad 85%
- Se detectarían 43 positivos y se darían 999 falsos positivos y 7 falsos negativos.

Ejemplo 2:

Test antígenos aplicado en residencias:

- El test de antígenos tiene sensibilidad (detección de 1), de un 93% pero en pacientes que llevan más de 5 días de evolución desde el inicio de los síntomas produce demasiados falsos negativos (baja la sensibilidad).
- En estos últimos casos habría que hacer PCR, no test de antígenos SARS.
- Realmente al final se ha llegado a la conclusión de que en residencias, si el test de antígenos da negativo, es prudente hacer un test PCR posterior.

Aspectos Relativos a la muestra A:

- 1) El número de observaciones con $y=1$ en la muestra
- 2) El número de observaciones total en la muestra
- 3) Aspectos relativos al desempeño del algoritmo predictivo:
 - 3.1 Número de variables input
 - 3.2 Auc obtenido en validación cruzada; Accuracy o tasa de fallos obtenida en validación cruzada
 - 3.3 Número de observaciones clasificadas como $y=1$ con punto de corte básico 0.5
 - 3.4 Sensitividad, especificidad con punto de corte básico 0.5, FP, FN
- 4) Porcentaje de $y=1$ en la muestra

Aspectos Relativos a la futura muestra B:

5) Si se conoce, número de observaciones a predecir en una muestra B, en la aplicación práctica futura del algoritmo. Puede ser una cifra aproximada.

El número aproximado de observaciones $y=1$ que se espera en esta muestra B también es una cifra importante.

En algunos contextos las predicciones no se van a aplicar sobre una muestra B de casos, sino en observaciones aisladas una a una. Entonces no aplica lo dicho en este punto 5).

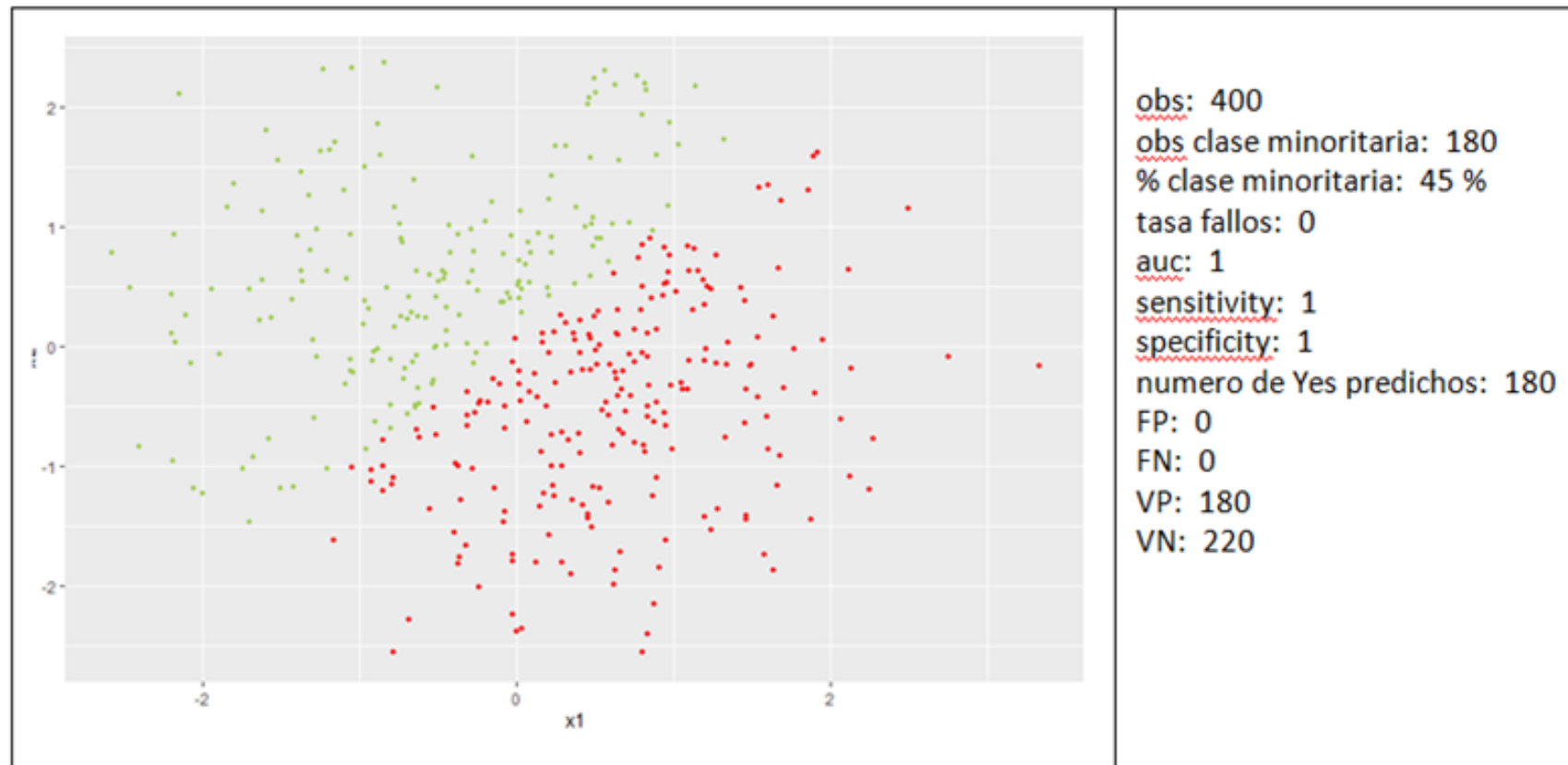
Ejemplos artificiales básicos

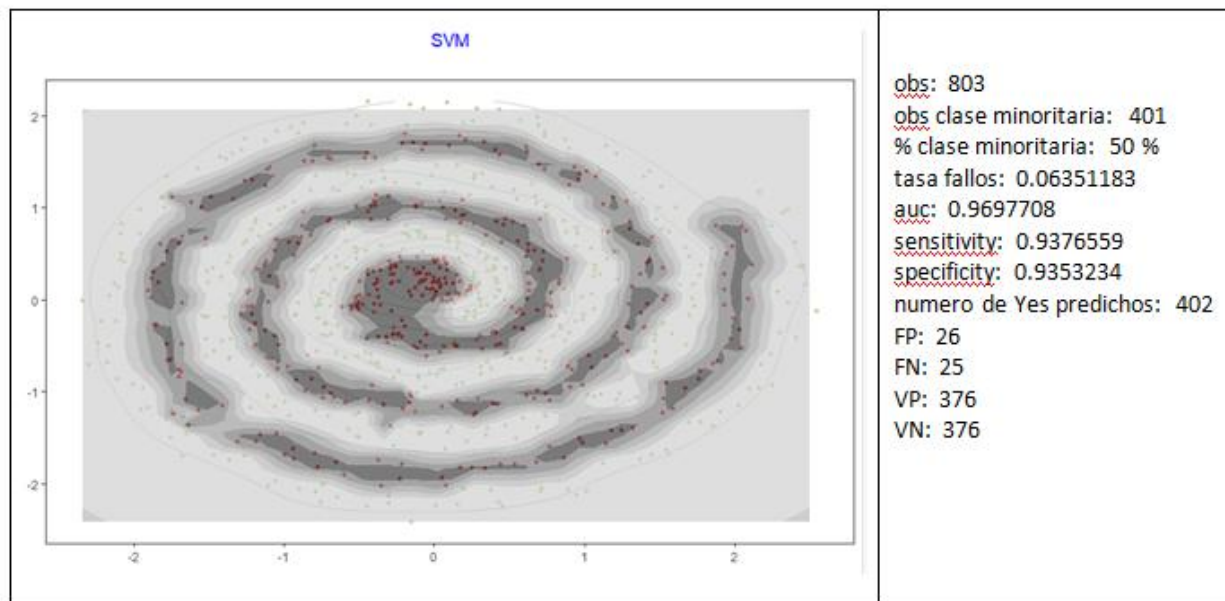
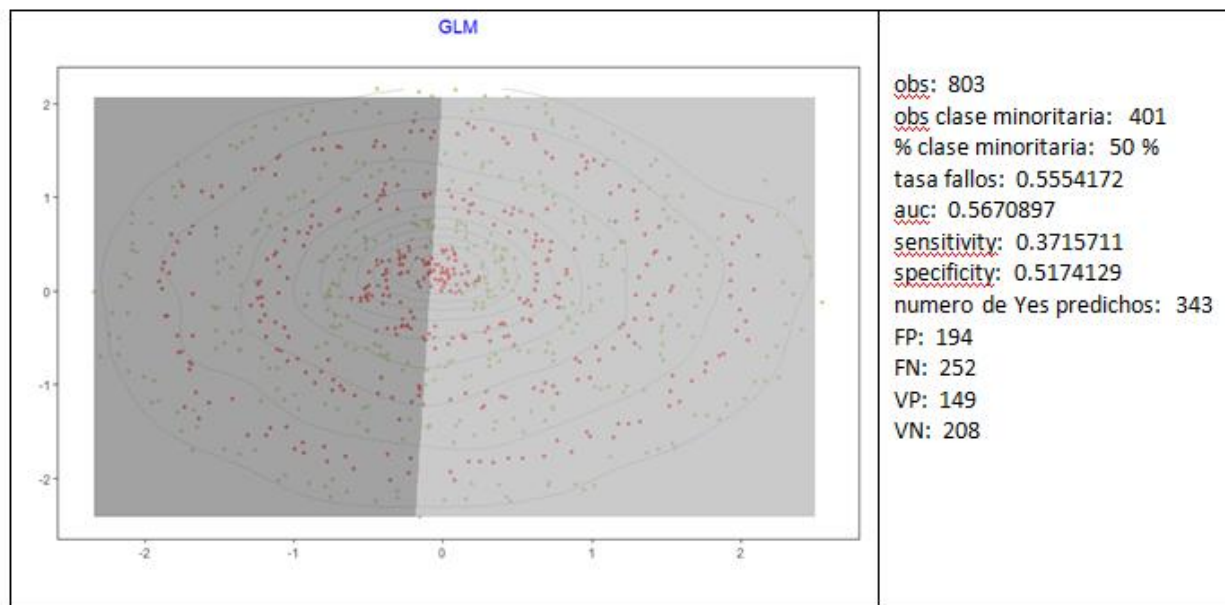
Archivo `graficos tema decision.R`

Sensitividad=Probabilidad de que la predicción sea 1, dado que la observación es realmente 1= capacidad de detectar positivos
 $=VP/(VP+FN)=0.65$. También se le puede llamar **recall**.

Especificidad=Probabilidad de que la predicción sea 0, dado que la observación es realmente 0=capacidad de detectar negativos
 $=VN/(VN+FP)=0.80$

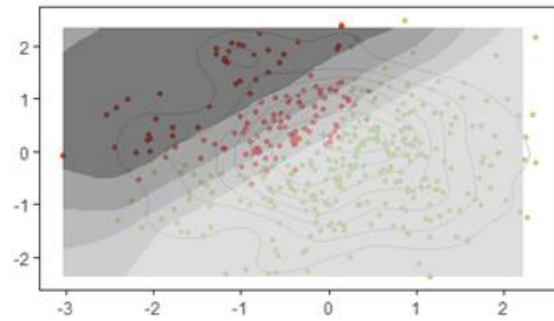
1



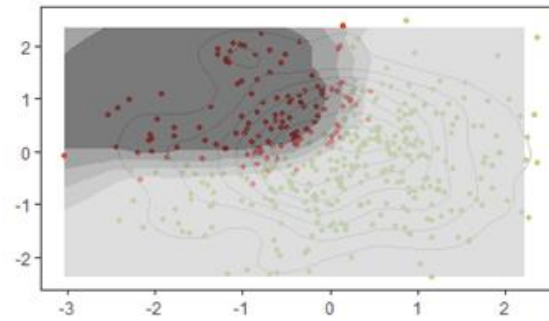


3

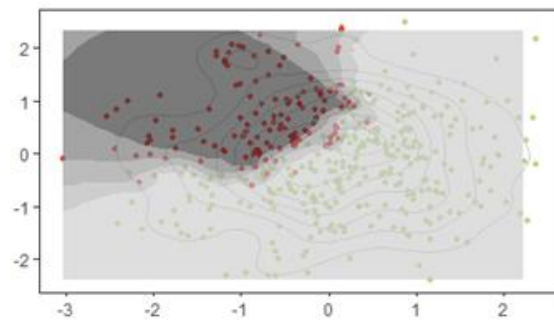
GLM



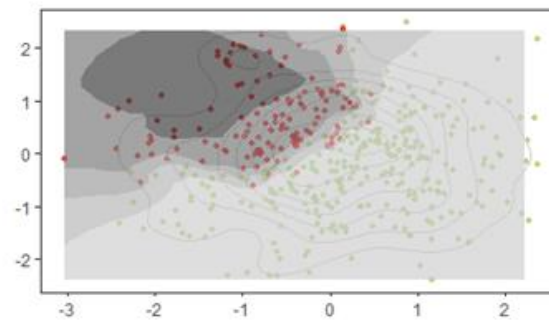
NNET



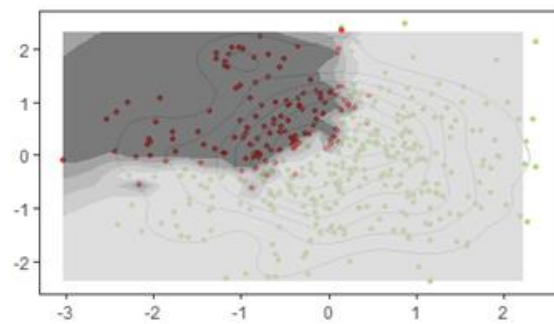
RF



GBM



SVM

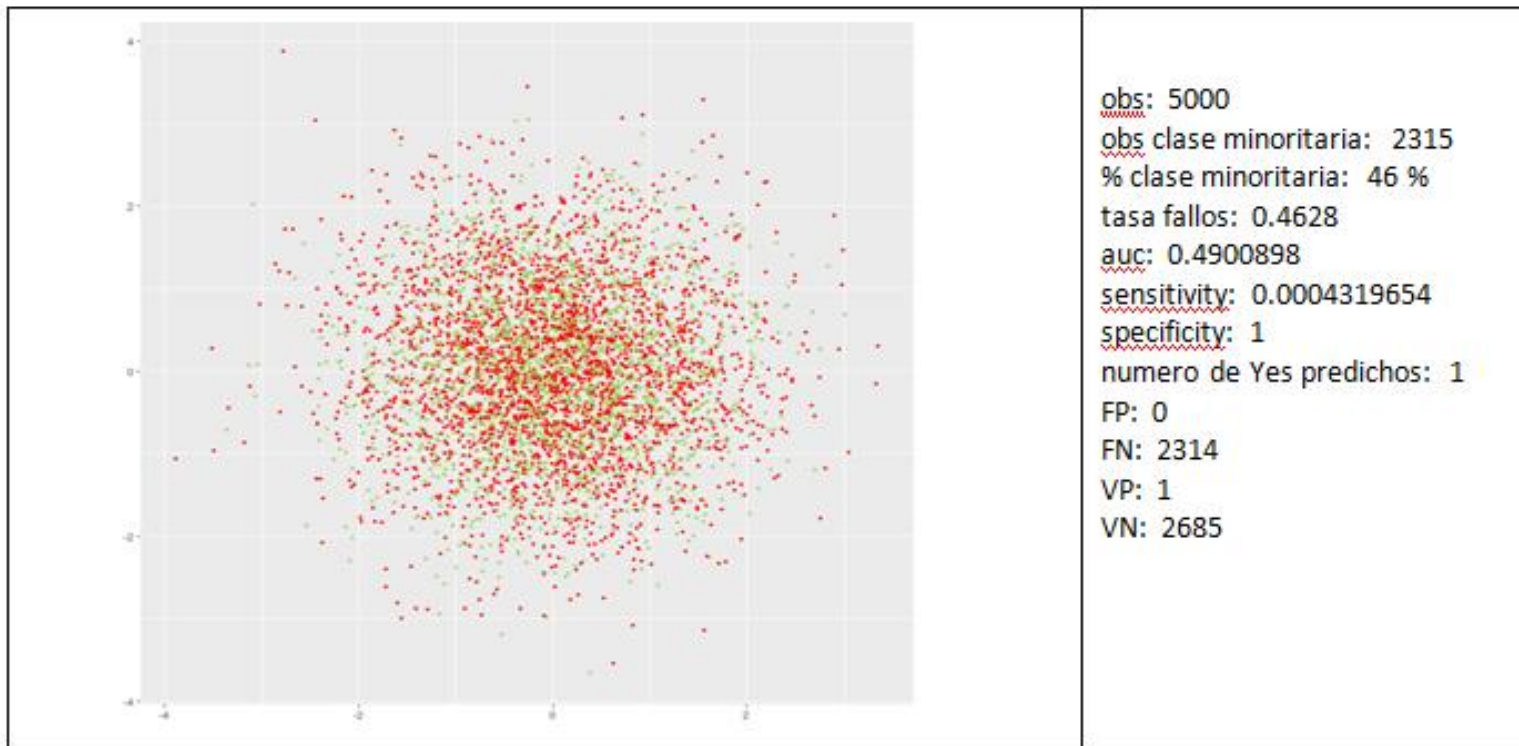


Recordatorio

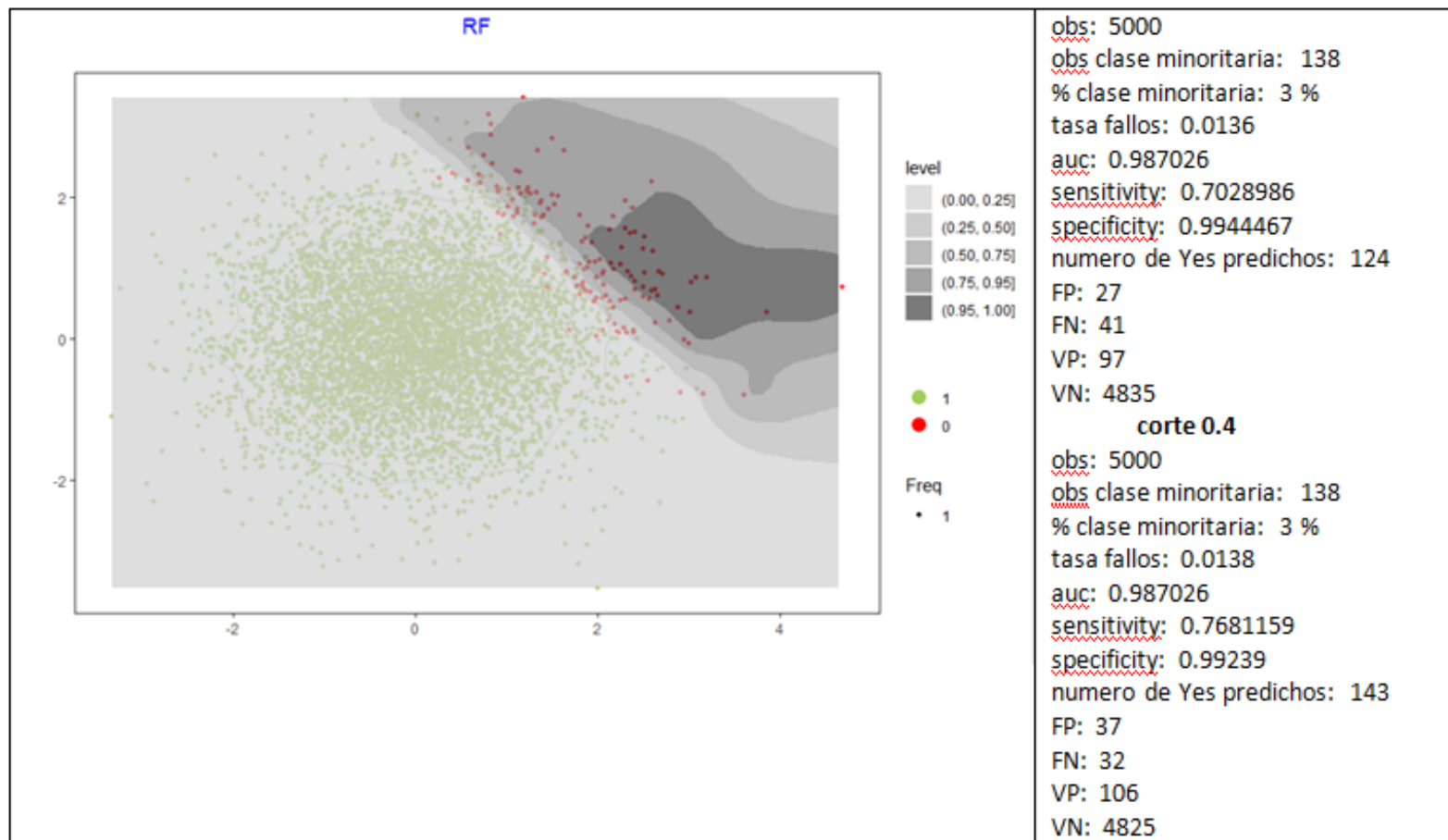
Sensitividad=Probabilidad de que la predicción sea 1, dado que la observación es realmente 1= capacidad de detectar positivos
 $=VP/(VP+FN)=0.65$. También se le puede llamar *recall*.

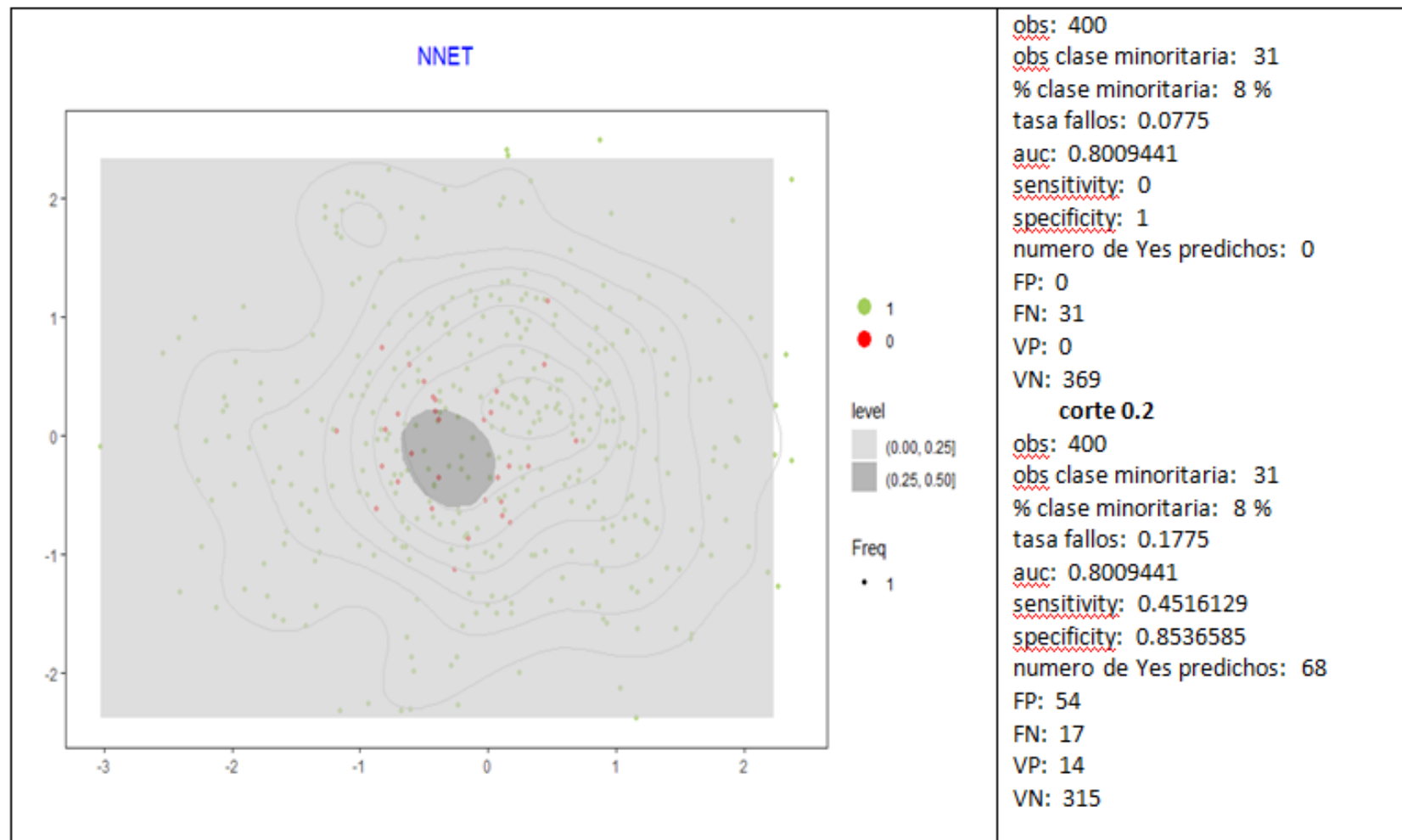
Especificidad=Probabilidad de que la predicción sea 0, dado que la observación es realmente 0=capacidad de detectar negativos
 $=VN/(VN+FP)=0.80$

4



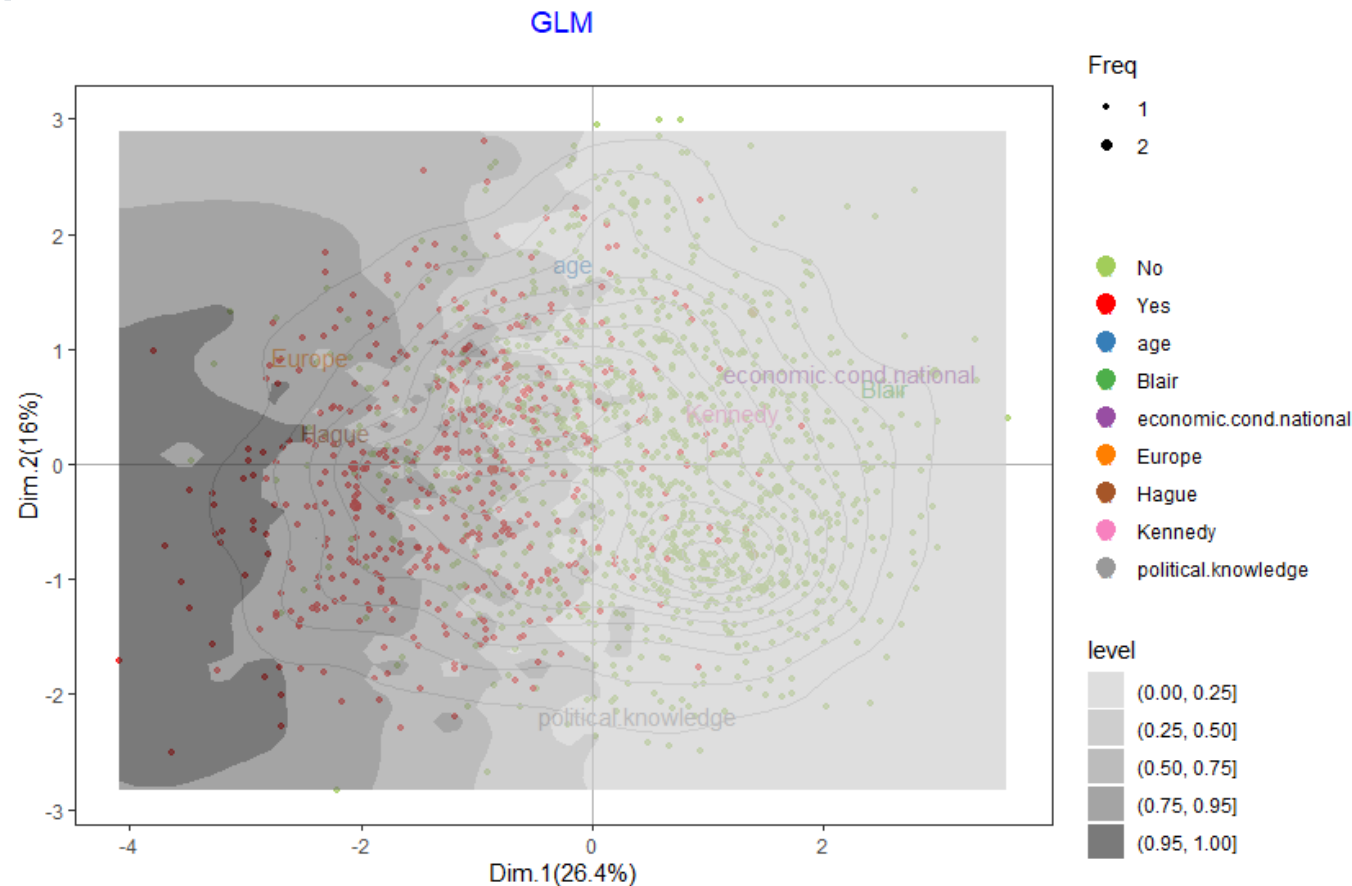
5





Paquete visualpred

- 1) MCA o FAMD para reducir el espacio de variables input en factores
- 2) Se presentan las observaciones sobre los dos primeros factores coloreadas según la variable objetivo
- 3) Se aplican algoritmos predictivos que dan lugar a una predicción de probabilidad para cada observación. Permite glm, rf,gbmy svm radial.
- 4) Se construye una rejilla de valores en el rango de los valores que toman los dos primeros factores, y se da valores a esa rejilla por interpolación, tomando como referencia las predicciones de probabilidad obtenidas en 3.
- 5) Con los valores interpolados de probabilidades obtenidas en 4 para esa rejilla, se superpone un gráfico de contour sobre los gráficos básicos de puntos.



Información y descarga:

<https://cran.r-project.org/web/packages/visualpred/index.html>

Vignettes=ejemplos

https://cran.r-project.org/web/packages/visualpred/vignettes/Basic_example.html

<https://cran.r-project.org/web/packages/visualpred/vignettes/Comparing.html>

<https://cran.r-project.org/web/packages/visualpred/vignettes/Outliers.html>

<https://cran.r-project.org/web/packages/visualpred/vignettes/Advanced.html>

Manual:

<https://cran.r-project.org/web/packages/visualpred/visualpred.pdf>

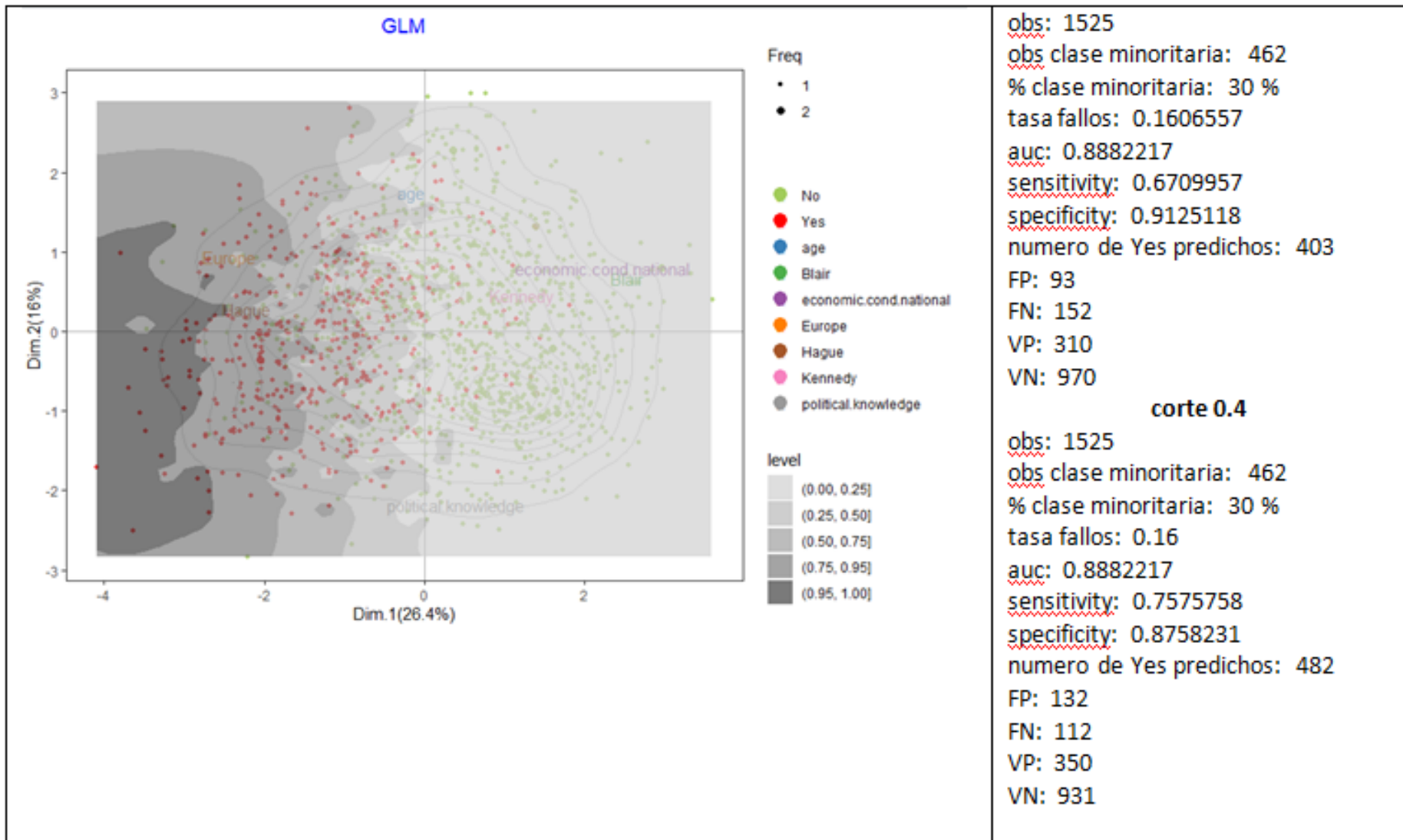
Ejemplo sintaxis

```
library(visualpred)
load("BEPS2.Rda")
archivo<-BEPS
listconti<-c("age", "economic.cond.national", "economic.cond.household",
             "Blair", "Hague", "Kennedy", "Europe", "political.knowledge")
listclass<-c("gender")
vardep<- "y"

result<-famdcontour(dataf=archivo,listconti=listconti,listclass=listclass,vardep=vardep,
                    title="rf",title2="
",Dime1="Dim.1",Dime2="Dim.2",selec=1,modelo="rf",classvar=0,mtry=5)
result[[2]]
result[[3]]
result[[4]]
result[[4]]
result[[5]]
result[[6]]
```

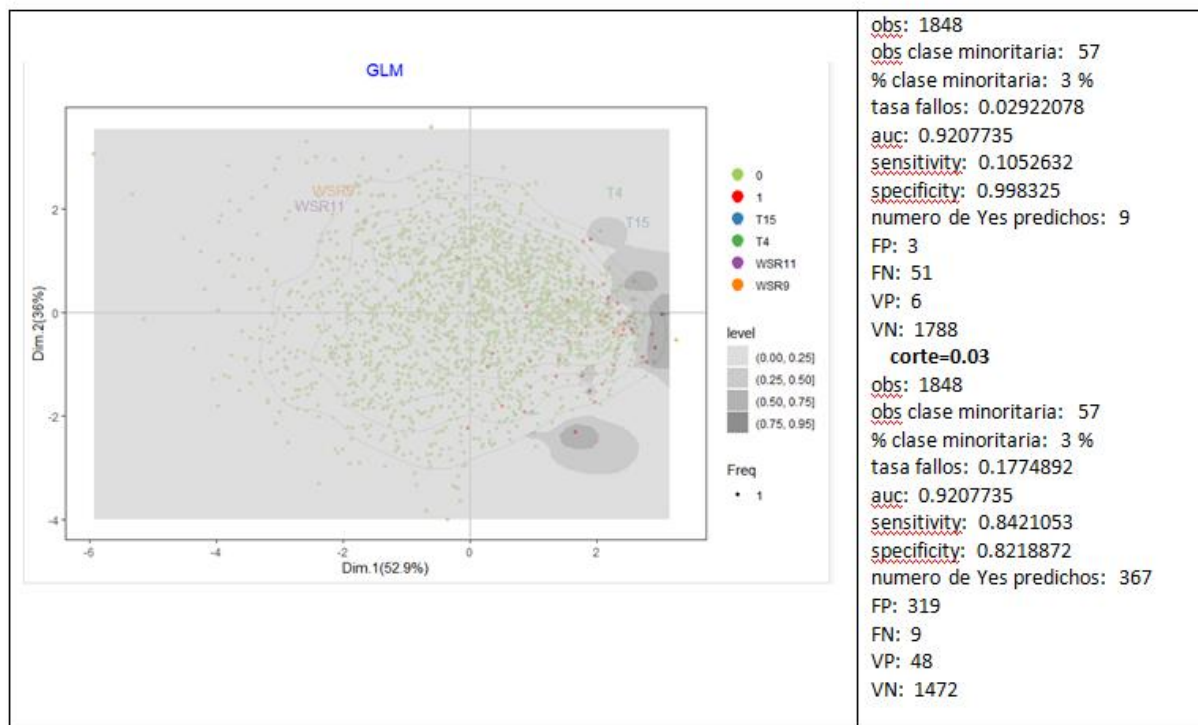
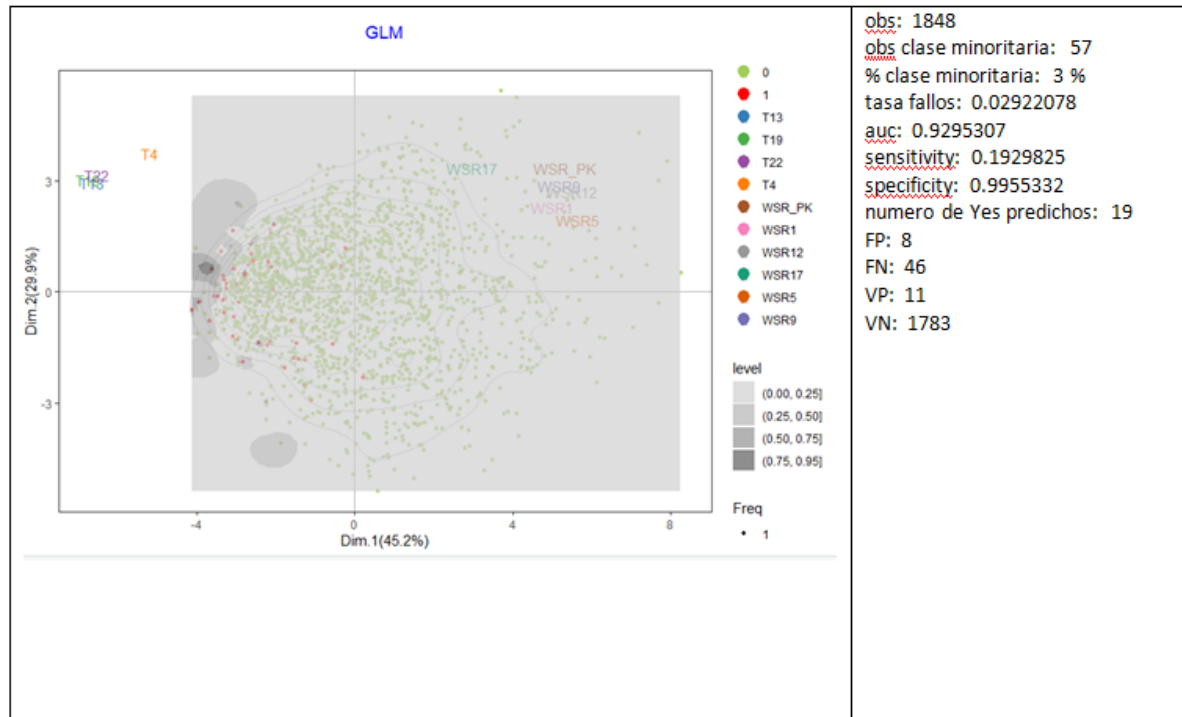
Ejemplos datos reales

BEPS



OZONO

10 variables

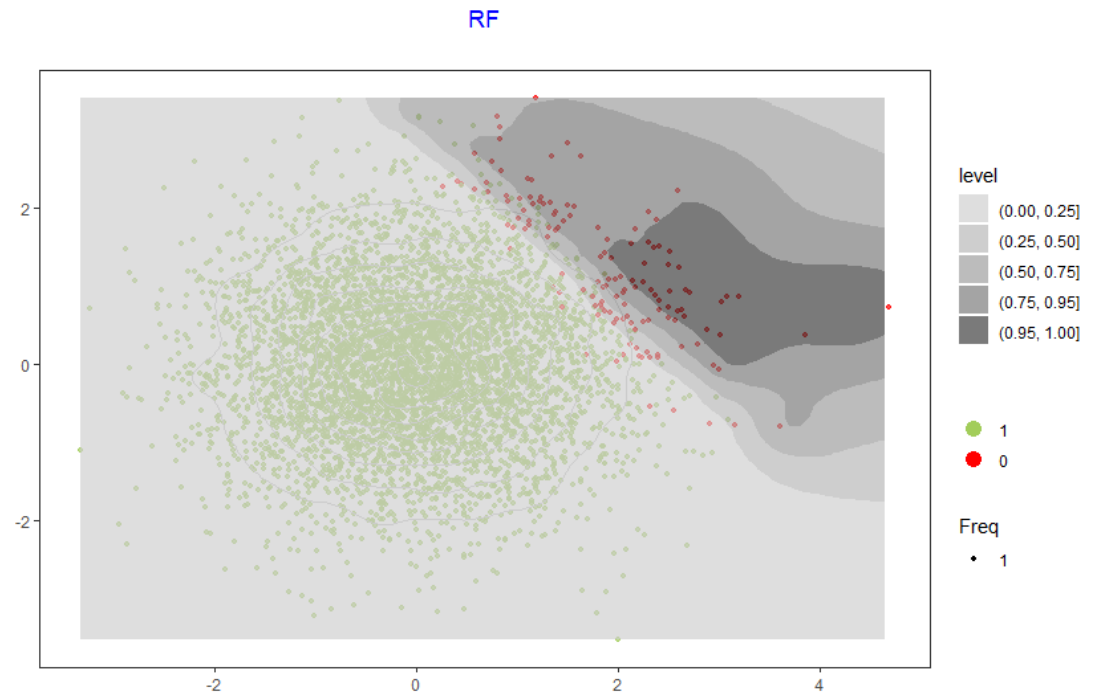
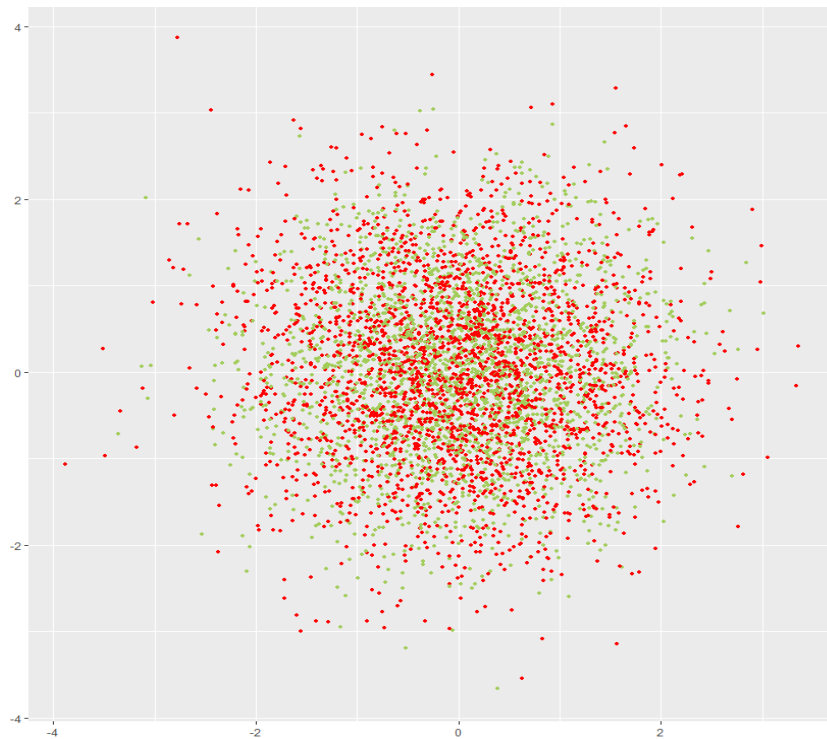


4 variables

La subcultura del “imbalanced dataset”

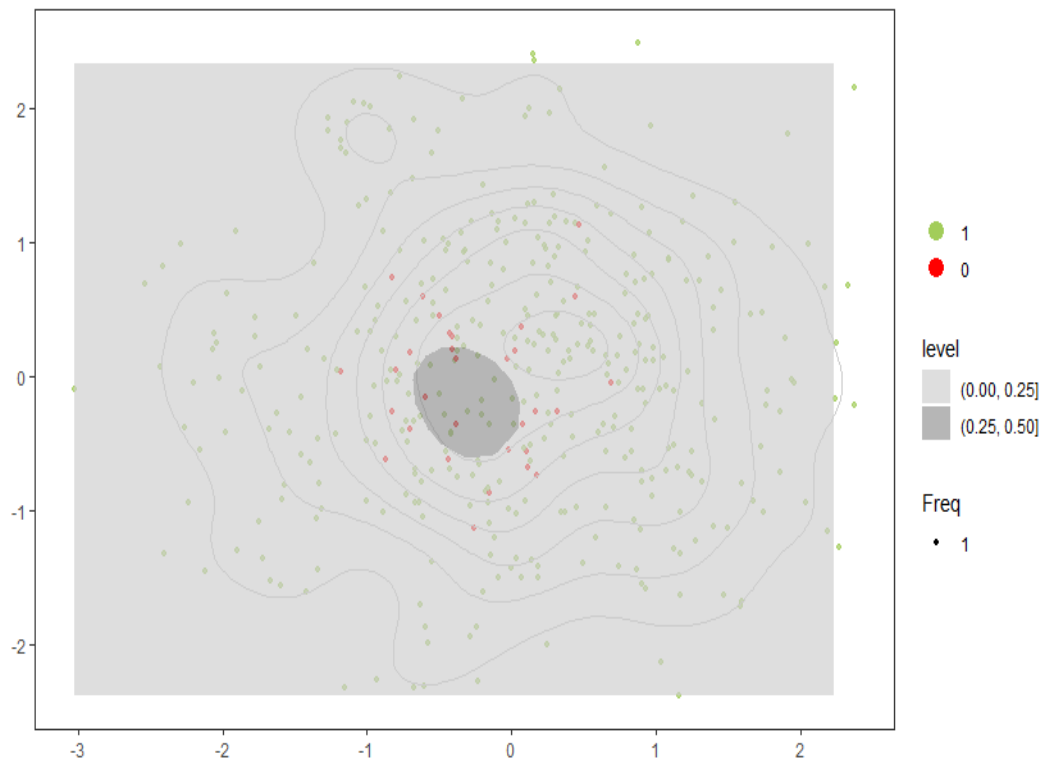
(sección de opinión)

- Izquierda: Clases equilibradas pero mala separabilidad
- Derecha: Clase desequilibradas pero buena separabilidad

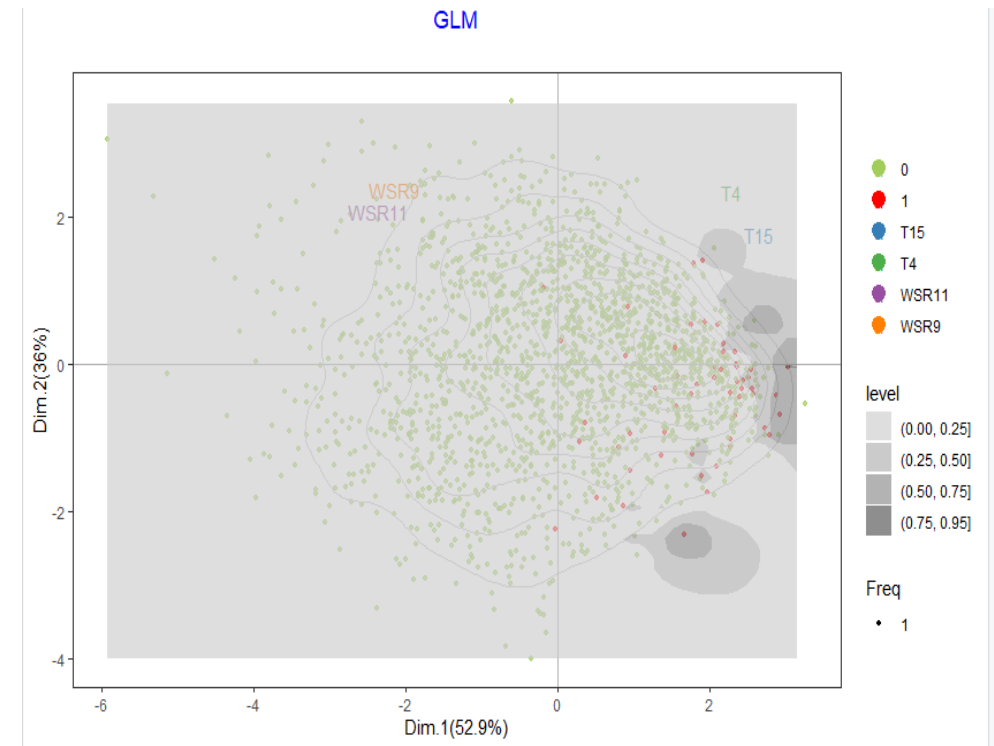


- Los algoritmos detectan bien las regiones y estiman correctamente las probabilidades, tanto en datos desequilibrados como que en todo tipo de datos.
- Es la **presencia o no de variables input útiles** (predictivas) la que determina la buena o mala separabilidad

NNET



GLM



Undersampling =quedarse con solo una porción de la clase mayoritaria con el objetivo de equilibrar las clases .

En mi opinión, es una estrategia equivocada:

- No se obtienen buenos estimadores de las medidas de diagnóstico y el error puede estar subestimado; las conclusiones no se pueden extrapolar directamente a la población al haber alterado la muestra
- Las regiones detectadas y las probabilidades estimadas no coinciden con la realidad
- Los algoritmos no funcionan “mejor”, funcionan sobre una muestra diferente

"Oversampling" (SMOTE y otros)= generar observaciones artificiales en la clase minoritaria para equilibrar las clases.

En mi opinión, también es una estrategia equivocada por las mismas razones anteriores:

- No se obtienen buenos estimadores de las medidas de diagnóstico y el error puede estar subestimado; las conclusiones no se pueden extrapolar directamente a la población al haber alterado la muestra
- Las regiones detectadas y las probabilidades estimadas no coinciden con la realidad
- Los algoritmos no funcionan “mejor”, funcionan sobre una muestra diferente

Comportamiento comparativo de algoritmos

Archivo `comparaciones básicas.R`

```
# listconti<-c("a1")
```

```
# listclass<-c("a2","a3","a4","a5","a6")
```

