

```
clc; clear all; clf;
```

## Avance 2: Entrenamiento, adecuación y evaluación de modelos

Para este avance se tomarán los diferentes datasets obtenidos en el avance 1 y se procederá a buscar el mejor/mejores modelos de clasificación para predecir la posibilidad de intento de suicidio recurrente.

En un primer momento se probarán diferentes modelos, haciendo uso de la herramienta "Classification Learner" de Matlab, debido a su facilidad y rapidez para probar múltiples modelos simultáneamente. Se tomarán los pares dataset-modelo que mejores resultados den (preferiblemente por encima de 70% de acierto) para seguirlos desarrollando, en términos de selección de parámetros y optimización de hiperparámetros.

Para la selección de parámetros y ajuste de hiperparámetros, en donde sea posible se usarán herramientas las interactivas o automáticas que provee Matlab.

Como métodos de validación y calificación de los modelos se pretenden usar los datos a continuación (**To Do: añadir breve descripción de cada uno**)

- Score?
- Matriz de confusión
- ROC curve

Al momento de realizar predicciones se generarán dos, una determinística y otra probabilística.

### Data sets de entrada.

En el avance 1 se obtuvieron 4 datasets después del proceso de limpieza, los cuales se mencionan a continuación:

- cds\_imputed : dataset con 33 características y 4146 registros
- cds : dataset con 28 características 4146 registros,
- cds\_few : dataset 33 características y 655 registros
- cds\_few\_minus\_alcohol: dataset 32 características 1690 registros.

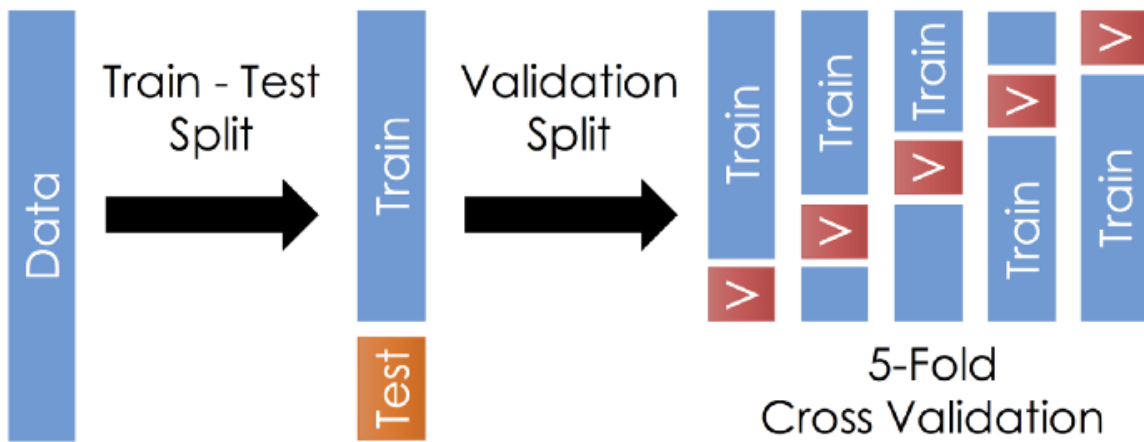
```
%cds = readtable('clean_datasets\cds.csv'); size_cds = size(cds)
cds_imputed = readtable('clean_datasets\cds_imputed.csv'); size_imputed = size(cds_imputed)
```

```
size_imputed = 1x2
              4146      34
```

```
cds_imputed = movevars(cds_imputed,'inten_prev','after','tipo_ss_S');
%cds_few = readtable('clean_datasets\cds_few.csv'); size_few = size (cds_few)
%cds_few_minus_alcohol = readtable('clean_datasets\cds_few_minus_alcohol.csv');
%      size_few_minus_alcohol = size (cds_few_minus_alcohol)
```

Con estos dataset se procede a realizar un entrenamiento exploratorio de modelos, para continuar con los más prometedores. Sin embargo, es necesario definir el concepto de "más prometedor". En este primer momento se tendrá en cuenta la exactitud de los modelos

Es de utilidad tener en cuenta que para el entrenamiento de los modelos fue usada validacion cruzaada con "k-folds" (  $k=5$  ),asi, el valor de la exactitud presentado corresponde a la exactitud de validacion y esta sirve como un estimado del desempeño del modelo en nuevos datos comparados con el conjunto de entrenamiento.



**Resultados cds**

<b>1.1</b>	☆ Tree	Accuracy: 63.8%
Last change: Fine Tree		28/28 features
<b>1.2</b>	☆ Tree	Accuracy: 65.1%
Last change: Medium Tree		28/28 features
<b>1.3</b>	☆ Tree	Accuracy: 66.7%
Last change: Coarse Tree		28/28 features
<b>1.4</b>	☆ Linear Discriminant	Accuracy: 66.9%
Last change: Linear Discriminant		28/28 features
<b>1.5</b>	☆ Quadratic Discriminant	<b>Failed</b>
Last change: Quadratic Discriminant		28/28 features
<b>1.6</b>	☆ Logistic Regression	Accuracy: <b>67.4%</b>
Last change: Logistic Regression		28/28 features
<b>1.7</b>	☆ Naive Bayes	Accuracy: 62.2%
Last change: Gaussian Naive Bayes		28/28 features
<b>1.8</b>	☆ Naive Bayes	Accuracy: 62.0%
Last change: Kernel Naive Bayes		28/28 features
<b>1.9</b>	☆ SVM	Accuracy: 66.0%
Last change: Linear SVM		28/28 features
<b>1.10</b>	☆ SVM	Accuracy: 65.3%
Last change: Quadratic SVM		28/28 features
<b>1.11</b>	☆ SVM	Accuracy: 63.8%
Last change: Cubic SVM		28/28 features
<b>1.12</b>	☆ SVM	Accuracy: 61.8%
Last change: Fine Gaussian SVM		28/28 features
<b>1.13</b>	☆ SVM	Accuracy: 65.8%
Last change: Medium Gaussian SVM		28/28 features

<b>1.14</b>	☆ SVM	Accuracy: 66.8%
Last change: Coarse Gaussian SVM 28/28 features		
<b>1.15</b>	☆ KNN	Accuracy: 59.5%
Last change: Fine KNN 28/28 features		
<b>1.16</b>	☆ KNN	Accuracy: 63.8%
Last change: Medium KNN 28/28 features		
<b>1.17</b>	☆ KNN	Accuracy: 65.6%
Last change: Coarse KNN 28/28 features		
<b>1.18</b>	☆ KNN	Accuracy: 63.7%
Last change: Cosine KNN 28/28 features		
<b>1.19</b>	☆ KNN	Accuracy: 63.3%
Last change: Cubic KNN 28/28 features		
<b>1.20</b>	☆ KNN	Accuracy: 62.4%
Last change: Weighted KNN 28/28 features		
<b>1.21</b>	☆ Ensemble	Accuracy: 66.6%
Last change: Boosted Trees 28/28 features		
<b>1.22</b>	☆ Ensemble	Accuracy: 65.4%
Last change: Bagged Trees 28/28 features		
<b>1.23</b>	☆ Ensemble	Accuracy: 66.8%
Last change: Subspace Discriminant 28/28 features		
<b>1.24</b>	☆ Ensemble	Accuracy: 63.6%
Last change: Subspace KNN 28/28 features		
<b>1.25</b>	☆ Ensemble	Accuracy: 64.7%
Last change: RUSBoosted Trees 28/28 features		
<b>2</b>	☆ Quadratic Discriminant	Accuracy: 62.2%
Last change: 'Covariance structure' ... 28/28 features		

Resultados cds\_imputed

<b>1.1</b>	☆ Tree	Accuracy: 64.5%
Last change: Fine Tree		33/33 features
<b>1.2</b>	☆ Tree	Accuracy: 66.1%
Last change: Medium Tree		33/33 features
<b>1.3</b>	☆ Tree	Accuracy: 66.8%
Last change: Coarse Tree		33/33 features
<b>1.4</b>	☆ Linear Discriminant	Accuracy: <b>67.9%</b>
Last change: Linear Discriminant		33/33 features
<b>1.5</b>	☆ Quadratic Discriminant	<b>Failed</b>
Last change: Quadratic Discriminant		33/33 features
<b>1.6</b>	☆ Logistic Regression	Accuracy: 67.8%
Last change: Logistic Regression		33/33 features
<b>1.7</b>	☆ Naive Bayes	Accuracy: 64.2%
Last change: Gaussian Naive Bayes		33/33 features
<b>1.8</b>	☆ Naive Bayes	Accuracy: 62.5%
Last change: Kernel Naive Bayes		33/33 features
<b>1.9</b>	☆ SVM	Accuracy: 66.8%
Last change: Linear SVM		33/33 features
<b>1.10</b>	☆ SVM	Accuracy: 65.4%
Last change: Quadratic SVM		33/33 features
<b>1.11</b>	☆ SVM	Accuracy: 63.3%
Last change: Cubic SVM		33/33 features
<b>1.12</b>	☆ SVM	Accuracy: 61.8%
Last change: Fine Gaussian SVM		33/33 features
<b>1.13</b>	☆ SVM	Accuracy: 65.5%
Last change: Medium Gaussian SVM		33/33 features

<b>1.14</b>	☆ SVM	Accuracy: 67.6%
Last change: Coarse Gaussian SVM 33/33 features		
<b>1.15</b>	☆ KNN	Accuracy: 61.3%
Last change: Fine KNN 33/33 features		
<b>1.16</b>	☆ KNN	Accuracy: 64.4%
Last change: Medium KNN 33/33 features		
<b>1.17</b>	☆ KNN	Accuracy: 64.8%
Last change: Coarse KNN 33/33 features		
<b>1.18</b>	☆ KNN	Accuracy: 64.6%
Last change: Cosine KNN 33/33 features		
<b>1.19</b>	☆ KNN	Accuracy: 64.5%
Last change: Cubic KNN 33/33 features		
<b>1.20</b>	☆ KNN	Accuracy: 63.7%
Last change: Weighted KNN 33/33 features		
<b>1.21</b>	☆ Ensemble	Accuracy: 67.6%
Last change: Boosted Trees 33/33 features		
<b>1.22</b>	☆ Ensemble	Accuracy: 65.6%
Last change: Bagged Trees 33/33 features		
<b>1.23</b>	☆ Ensemble	Accuracy: 67.7%
Last change: Subspace Discriminant 33/33 features		
<b>1.24</b>	☆ Ensemble	Accuracy: 63.8%
Last change: Subspace KNN 33/33 features		
<b>1.25</b>	☆ Ensemble	Accuracy: 64.7%
Last change: RUSBoosted Trees 33/33 features		
<b>2</b>	☆ Quadratic Discriminant	Accuracy: 64.2%
Last change: 'Covariance structure' ... 33/33 features		

## Resultados cds\_few

Para este dataset algunos modelos se hicieron individualmente, porque presentaban problemas con las características 'antec\_tran', 'tipo\_ss\_l', 'suici\_fm\_a' y 'tipo\_SS\_P' ya que la mayoría o casi todos sus valores son iguales por lo que no aportan información o no presentan variación con respecto a una de las clases por hallar.

<b>1.1</b> ☆ Tree	Accuracy: 53.9%
Last change: Fine Tree	33/33 features
<b>1.2</b> ☆ Tree	Accuracy: 60.5%
Last change: Medium Tree	33/33 features
<b>1.3</b> ☆ Tree	Accuracy: <b>64.4%</b>
Last change: Coarse Tree	33/33 features
<b>1.4</b> ☆ Linear Discriminant	<b>Failed</b>
Last change: Linear Discriminant	33/33 features
<b>1.5</b> ☆ Quadratic Discriminant	<b>Failed</b>
Last change: Quadratic Discriminant	33/33 features
<b>1.6</b> ☆ Logistic Regression	Accuracy: 61.8%
Last change: Logistic Regression	33/33 features
<b>1.7</b> ☆ Naive Bayes	<b>Failed</b>
Last change: Gaussian Naive Bayes	33/33 features
<b>1.8</b> ☆ Naive Bayes	Accuracy: 61.1%
Last change: Kernel Naive Bayes	33/33 features
<b>1.9</b> ☆ SVM	Accuracy: 61.2%
Last change: Linear SVM	33/33 features
<b>1.10</b> ☆ SVM	Accuracy: 57.7%
Last change: Quadratic SVM	33/33 features
<b>1.11</b> ☆ SVM	Accuracy: 57.3%
Last change: Cubic SVM	33/33 features
<b>1.12</b> ☆ SVM	Accuracy: 58.9%
Last change: Fine Gaussian SVM	33/33 features
<b>1.13</b> ☆ SVM	Accuracy: 63.7%
Last change: Medium Gaussian SVM	33/33 features
<b>1.14</b> ☆ SVM	Accuracy: 61.1%
Last change: Coarse Gaussian SVM	33/33 features



<b>1.15</b>	☆ KNN	Accuracy: 53.6%
Last change: Fine KNN		33/33 features
<b>1.16</b>	☆ KNN	Accuracy: 60.8%
Last change: Medium KNN		33/33 features
<b>1.17</b>	☆ KNN	Accuracy: 60.8%
Last change: Coarse KNN		33/33 features
<b>1.18</b>	☆ KNN	Accuracy: 61.2%
Last change: Cosine KNN		33/33 features
<b>1.19</b>	☆ KNN	Accuracy: 60.2%
Last change: Cubic KNN		33/33 features
<b>1.20</b>	☆ KNN	Accuracy: 57.7%
Last change: Weighted KNN		33/33 features
<b>1.21</b>	☆ Ensemble	Accuracy: 59.1%
Last change: Boosted Trees		33/33 features
<b>1.22</b>	☆ Ensemble	Accuracy: 58.3%
Last change: Bagged Trees		33/33 features
<b>1.23</b>	☆ Ensemble	Accuracy: 63.4%
Last change: Subspace Discriminant		33/33 features
<b>1.24</b>	☆ Ensemble	Accuracy: 57.3%
Last change: Subspace KNN		33/33 features
<b>1.25</b>	☆ Ensemble	Accuracy: 58.3%
Last change: RUSBoosted Trees		33/33 features
<b>2</b>	☆ Linear Discriminant	Accuracy: 61.5%
Last change: 'Covariance structure' ...		33/33 features
<b>3</b>	☆ Quadratic Discriminant	Accuracy: 60.6%
Last change: 'Covariance structure' ...		33/33 features
<b>4</b>	☆ Naive Bayes	Accuracy: 53.6%
Last change: Removed 3 features		29/33 features

Resultados cds\_few\_minus\_alcohol



<b>1.1</b>	☆ Tree	Accuracy: 54.2%
Last change: Fine Tree		32/32 features
<b>1.2</b>	☆ Tree	Accuracy: 55.6%
Last change: Medium Tree		32/32 features
<b>1.3</b>	☆ Tree	Accuracy: 55.6%
Last change: Coarse Tree		32/32 features
<b>1.4</b>	☆ Linear Discriminant	<u>Failed</u>
Last change: Linear Discriminant		32/32 features
<b>1.5</b>	☆ Quadratic Discriminant	<u>Failed</u>
Last change: Quadratic Discriminant		32/32 features
<b>1.6</b>	☆ Logistic Regression	Accuracy: 59.0%
Last change: Logistic Regression		32/32 features
<b>1.7</b>	☆ Naive Bayes	<u>Failed</u>
Last change: Gaussian Naive Bayes		32/32 features
<b>1.8</b>	☆ Naive Bayes	Accuracy: 55.0%
Last change: Kernel Naive Bayes		32/32 features
<b>1.9</b>	☆ SVM	Accuracy: 58.5%
Last change: Linear SVM		32/32 features
<b>1.10</b>	☆ SVM	Accuracy: 55.3%
Last change: Quadratic SVM		32/32 features
<b>1.11</b>	☆ SVM	Accuracy: 53.3%
Last change: Cubic SVM		32/32 features
<b>1.12</b>	☆ SVM	Accuracy: 54.3%
Last change: Fine Gaussian SVM		32/32 features
<b>1.13</b>	☆ SVM	Accuracy: 56.5%
Last change: Medium Gaussian SVM		32/32 features
<b>1.14</b>	☆ SVM	Accuracy: 58.4%
Last change: Coarse Gaussian SVM		32/32 features

<b>1.15</b>	☆ KNN	Accuracy: 53.7%
Last change: Fine KNN		32/32 features
<b>1.16</b>	☆ KNN	Accuracy: 55.1%
Last change: Medium KNN		32/32 features
<b>1.17</b>	☆ KNN	Accuracy: 56.2%
Last change: Coarse KNN		32/32 features
<b>1.18</b>	☆ KNN	Accuracy: 55.4%
Last change: Cosine KNN		32/32 features
<b>1.19</b>	☆ KNN	Accuracy: 54.7%
Last change: Cubic KNN		32/32 features
<b>1.20</b>	☆ KNN	Accuracy: 55.0%
Last change: Weighted KNN		32/32 features
<b>1.21</b>	☆ Ensemble	Accuracy: 59.2%
Last change: Boosted Trees		32/32 features
<b>1.22</b>	☆ Ensemble	Accuracy: 56.6%
Last change: Bagged Trees		32/32 features
<b>1.23</b>	☆ Ensemble	Accuracy: 58.1%
Last change: Subspace Discriminant		32/32 features
<b>1.24</b>	☆ Ensemble	Accuracy: 54.0%
Last change: Subspace KNN		32/32 features
<b>1.25</b>	☆ Ensemble	Accuracy: 57.3%
Last change: RUSBoosted Trees		32/32 features
<b>2</b>	☆ Linear Discriminant	Accuracy: <b>60.2%</b>
Last change: 'Covariance structure' ...		32/32 features
<b>3</b>	☆ Quadratic Discriminant	Accuracy: 55.0%
Last change: 'Covariance structure' ...		32/32 features
<b>4</b>	☆ Naive Bayes	Accuracy: 54.4%
Last change: Removed 3 features		28/32 features

Por motivos exploratorios se realizaron pruebas aplicandole PCA a los datos, pero los resultados en general fueron inferiores a los obtenidos sin esta transformacion, por lo que esta transformacion de los datos no sera utilizada. (*¿Uno si deberia hacer PCA en datos categoricos?*)

Como se puede notar, ningún par dataset-modelo obtuvo una precisión mayor al 70% tal y como se había definido inicialmente para su aceptación. Por este motivo se tomará aquel dataset que produjo el modelo con la mayor precisión (cds\_imputed) y los mejores modelos obtenidos a partir de este -Coarse Tree, Linear discriminant, Logistic regresión, SVM (lineal y coarse) y Ensemble(BoostTrees, SubsD)-

## Feature selection

Buscando reducir la dimensionalidad y explorar diferentes opciones se pretende realizar un proceso de selección de características. Esto se hará filtrando aquellas características menos importantes para la respuesta 'inten\_prev' mediante el algoritmo MRMR (Minimum Redundancy Maximum Relevance), del cual se puede obtener el "ranking" de importancia de los predictores teniendo en cuenta la respuesta.

Se entrenarán 2 modelos, uno con todas las características y adicionalmente otro con el conjunto de las 7 más importantes

```
idx = fscmr(r(cds_imputed, 'inten_prev'));  
most_signif_features = cds_imputed.Properties.VariableNames(idx(1:7)).'
```

```
most_signif_features = 7x1 cell  
'antec_tran'  
'hist_famil'  
'muerte_fam'  
'antec_v_a'  
'prob_consu'  
'plan_suici'  
'gp_psiquia'
```

```
less_signif_features = cds_imputed.Properties.VariableNames(idx(end-4:end)).'
```

```
less_signif_features = 5x1 cell  
'escolarid'  
'esco_educ'  
'tipo_ss_C'  
'trab_socia'  
'sexo_'
```

## Optimización de hiperparámetros

Para la optimización de hiperparámetros serán utilizados dos enfoques:

Para los modelos simples (e.g. árboles de decisión) se realizará mediante GridSearch, mientras que para los más complejos (Ensamble, SVM) será utilizado un método de optimización bayesiana, el cual, a través de 30 iteraciones se va redirigiendo hacia aquellos hiperparámetros del espacio de búsqueda que proveen mejores resultados para el modelo. Esta elección se hace debido a los costos computacionales elevados de realizar Grid o Random Search en modelos complejos.

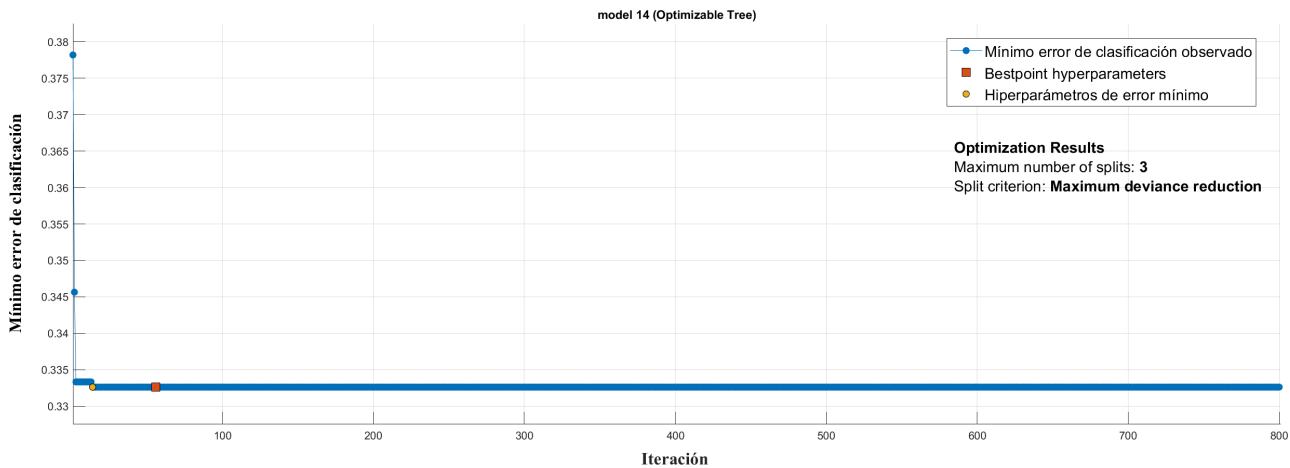
Como parte del proceso de optimización se obtiene el gráfico de error de clasificación mínimo en el cual se encuentran principalmente :

- Los resultados de la optimización
- El mínimo error de clasificación observado(puntos azules) hasta la iteración actual
- Bestpoint hyperparameters(cuadrado rojo), indica la iteración que corresponde a los valores de los hiperparámetros optimizados

Los resultados se presentaran a continuacion para cada modelo:

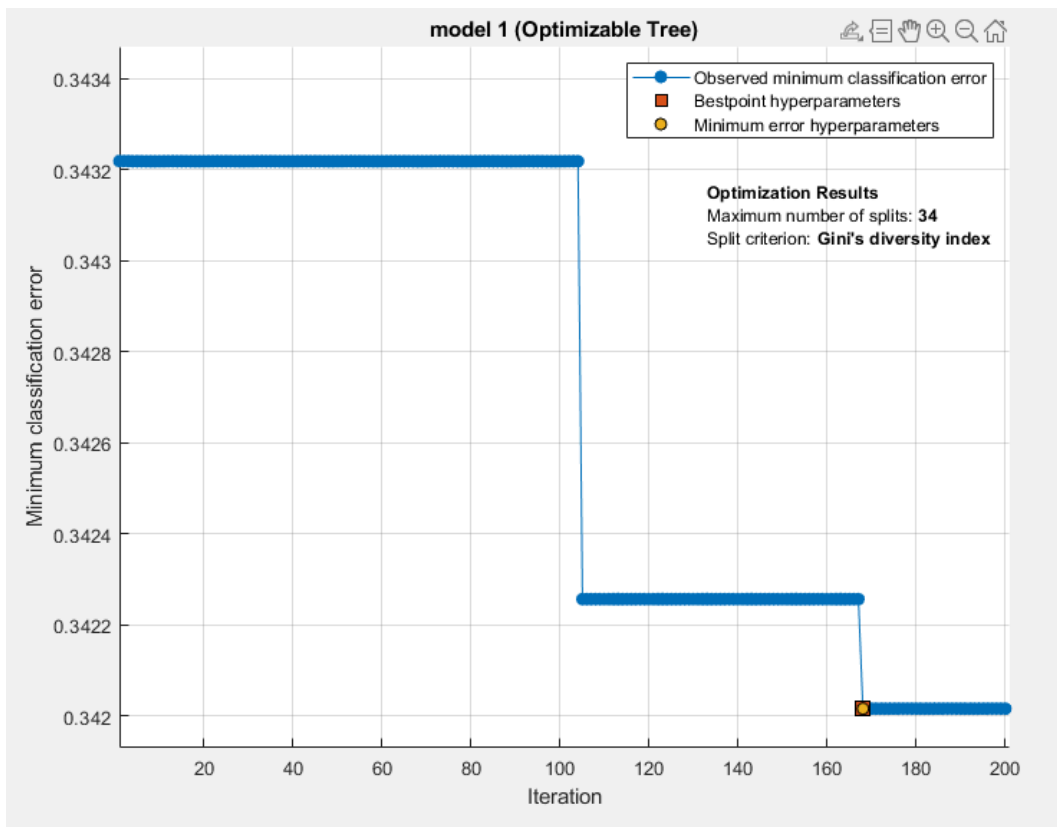
### Arboles de decision:

En la anterior etapa mejores resultados fueron obtenidos con arboles de decision gruesos(con poco numero de splits), lo cual es confirmado con los resultados de este proceso



### Arboles de decision( 7 características):

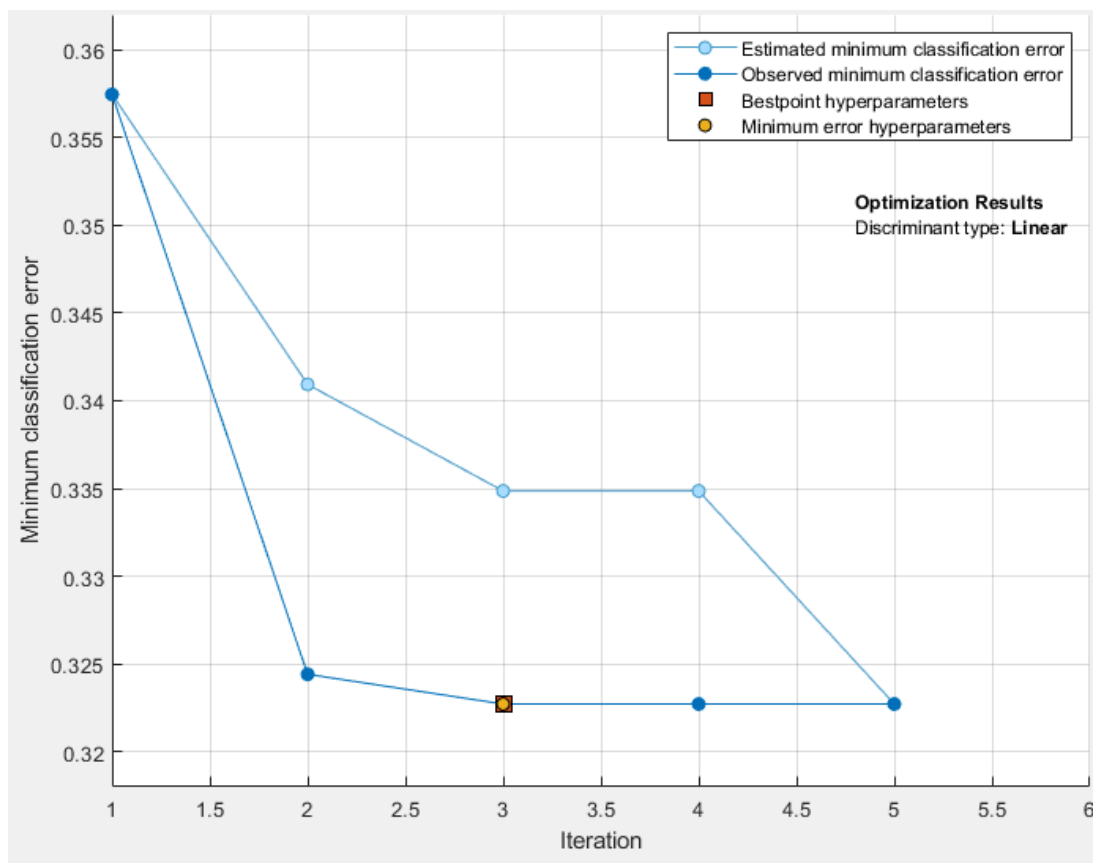
Para este caso los mejores resultados fueron con arboles finos



### Discriminante:

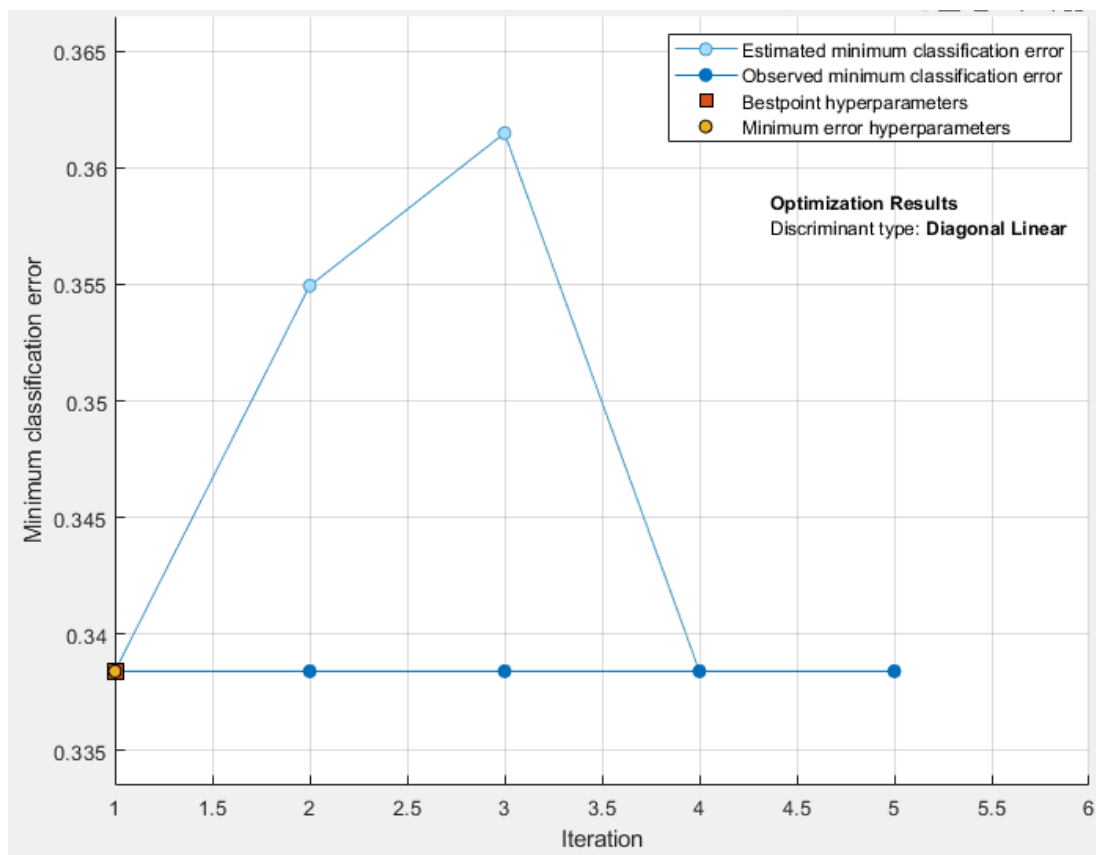
Anteriormente se habia hallado que el lineal era el que mejores resultados presentaba, esto se comprueba/ reafirma al raelizar este paso.

Las combinaciones disponibles para este tipo de modelos son pocas por lo que con pocas iteraciones es suficiente



**Discriminante(7 características):**

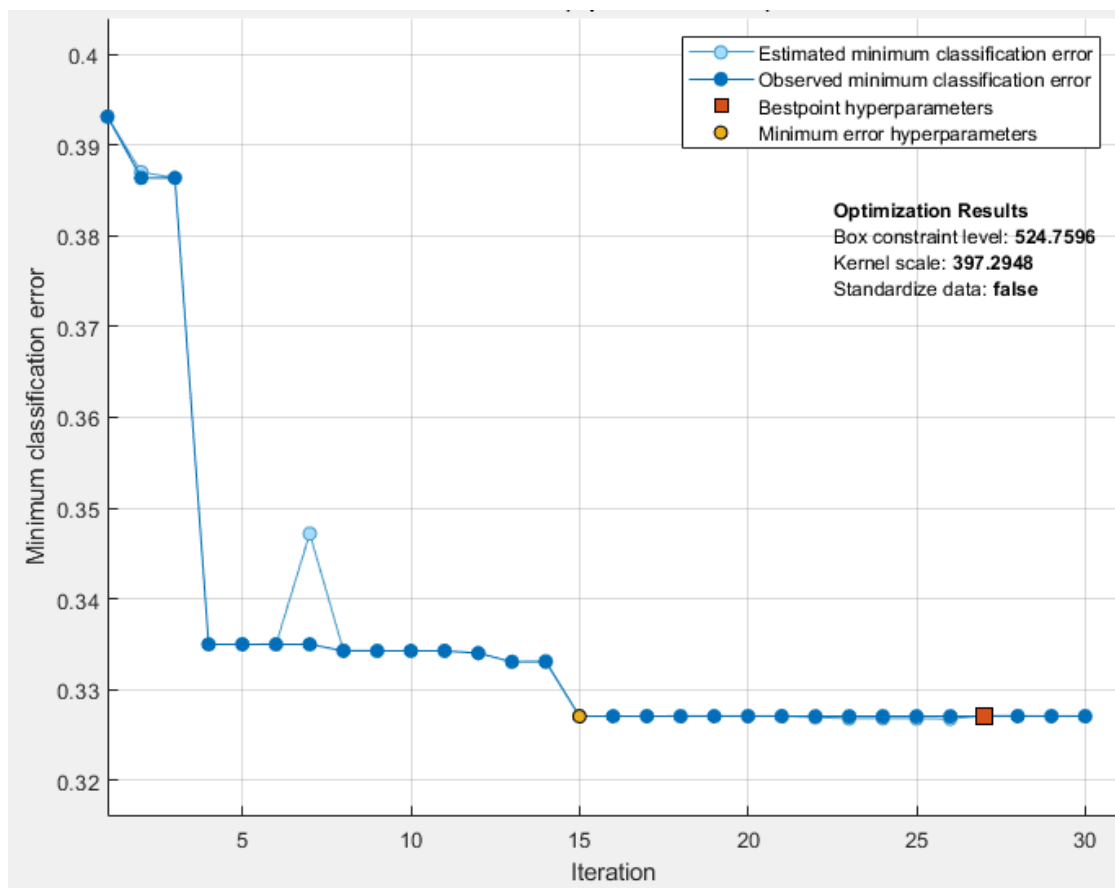
Mejores resultados con diagonal linear



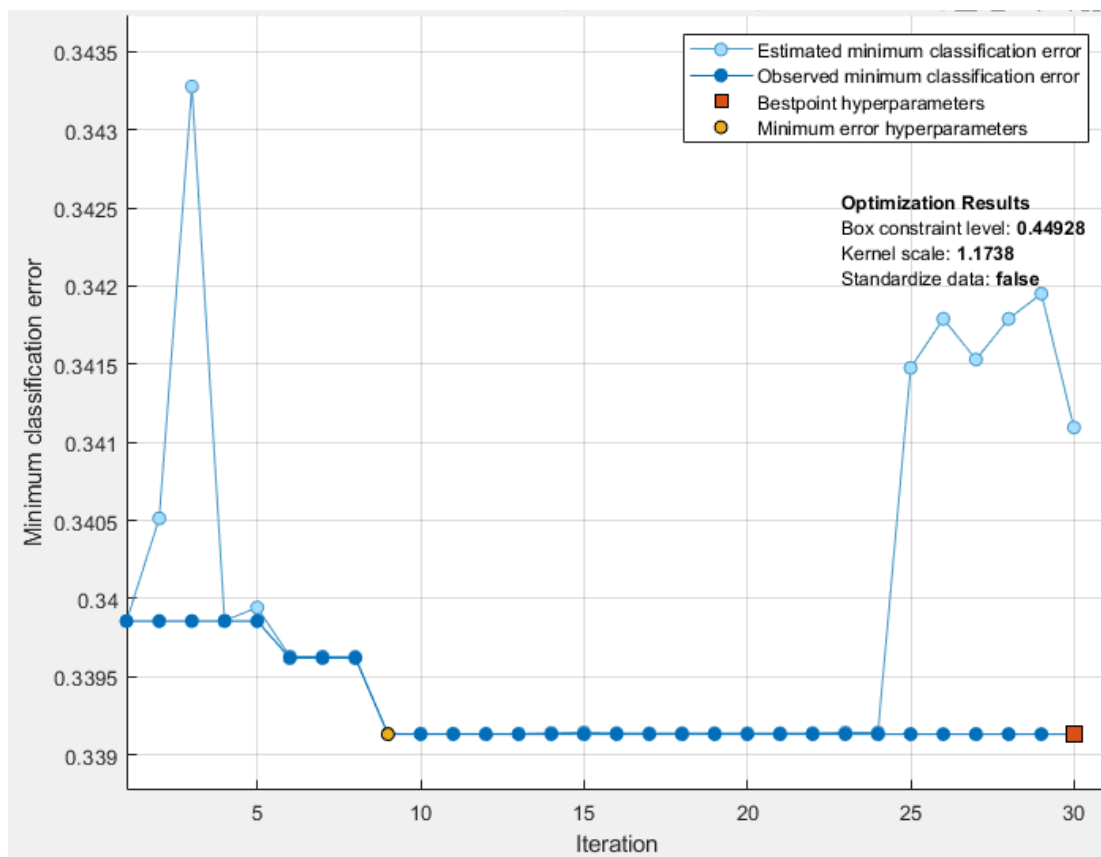
### SVM:

Para SVM los mejores resultados se obtuvieron para gaussiano "course" por lo que este sera optimizado

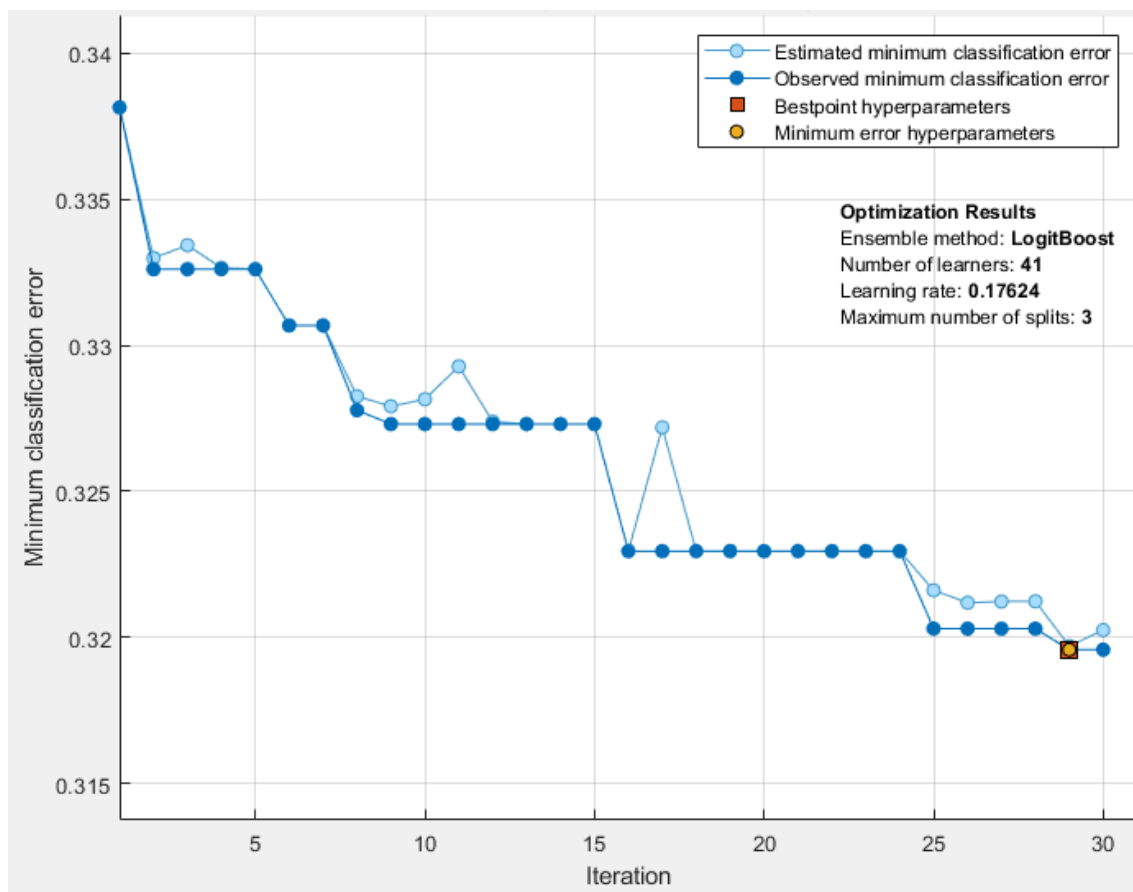




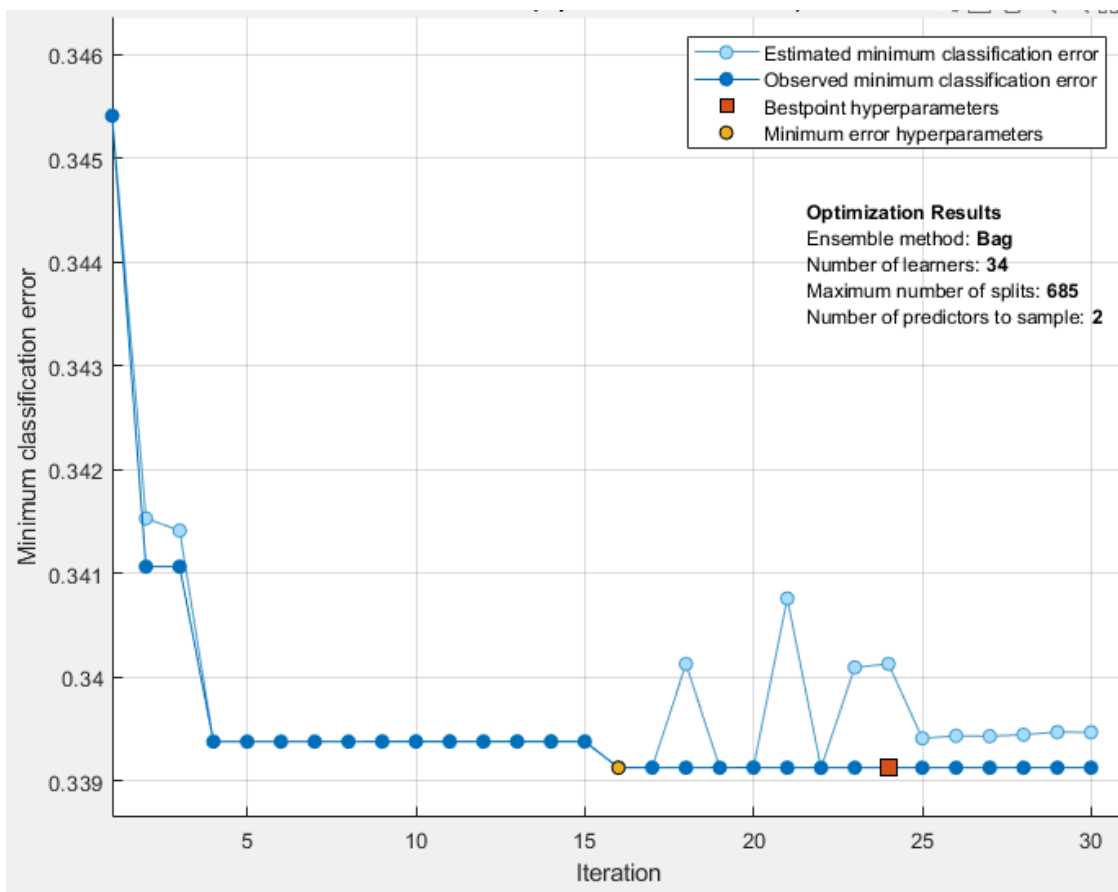
**SVM(7 características):**



Ensemble:



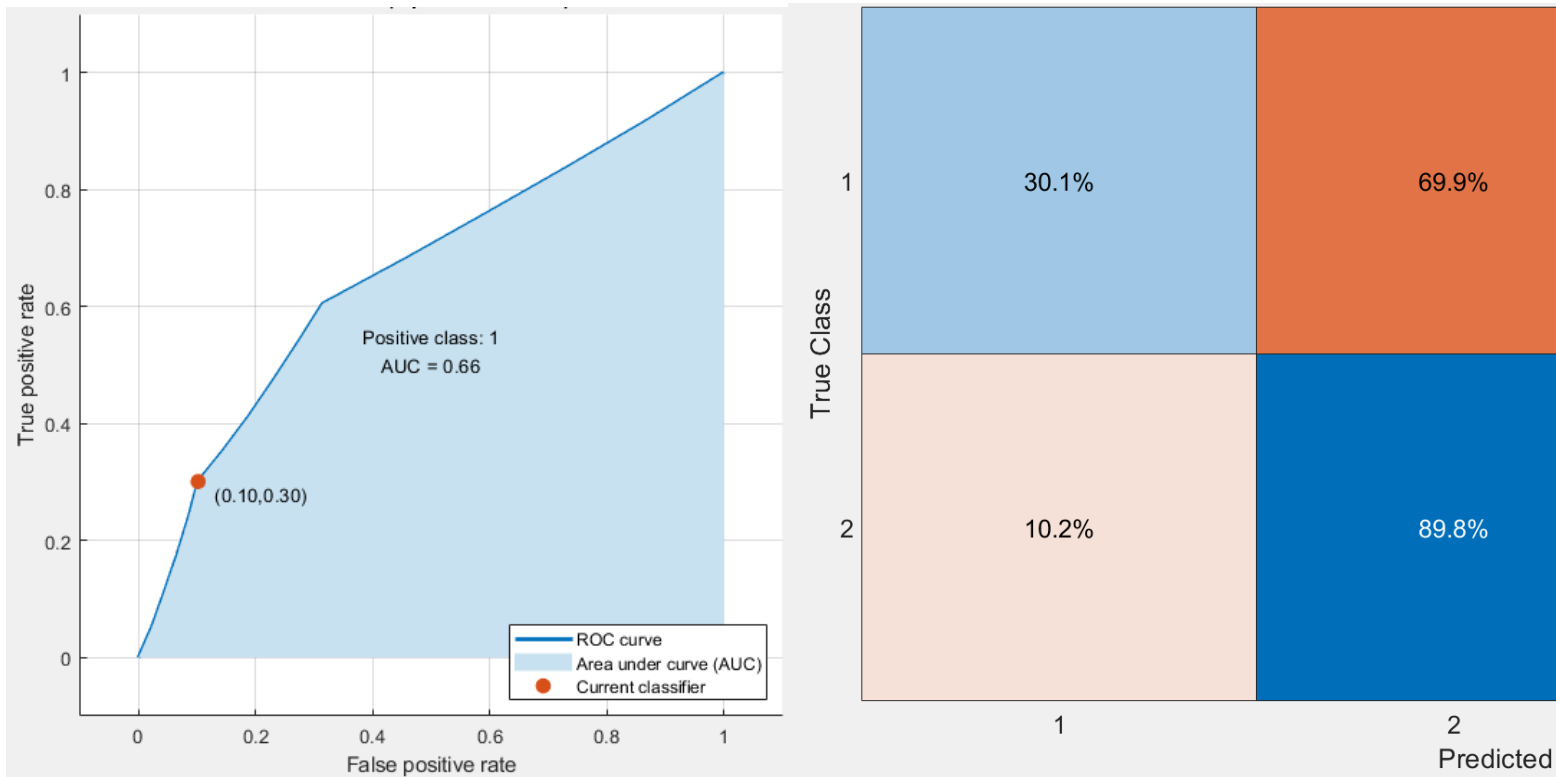
Ensemble (7 caracteristicas):



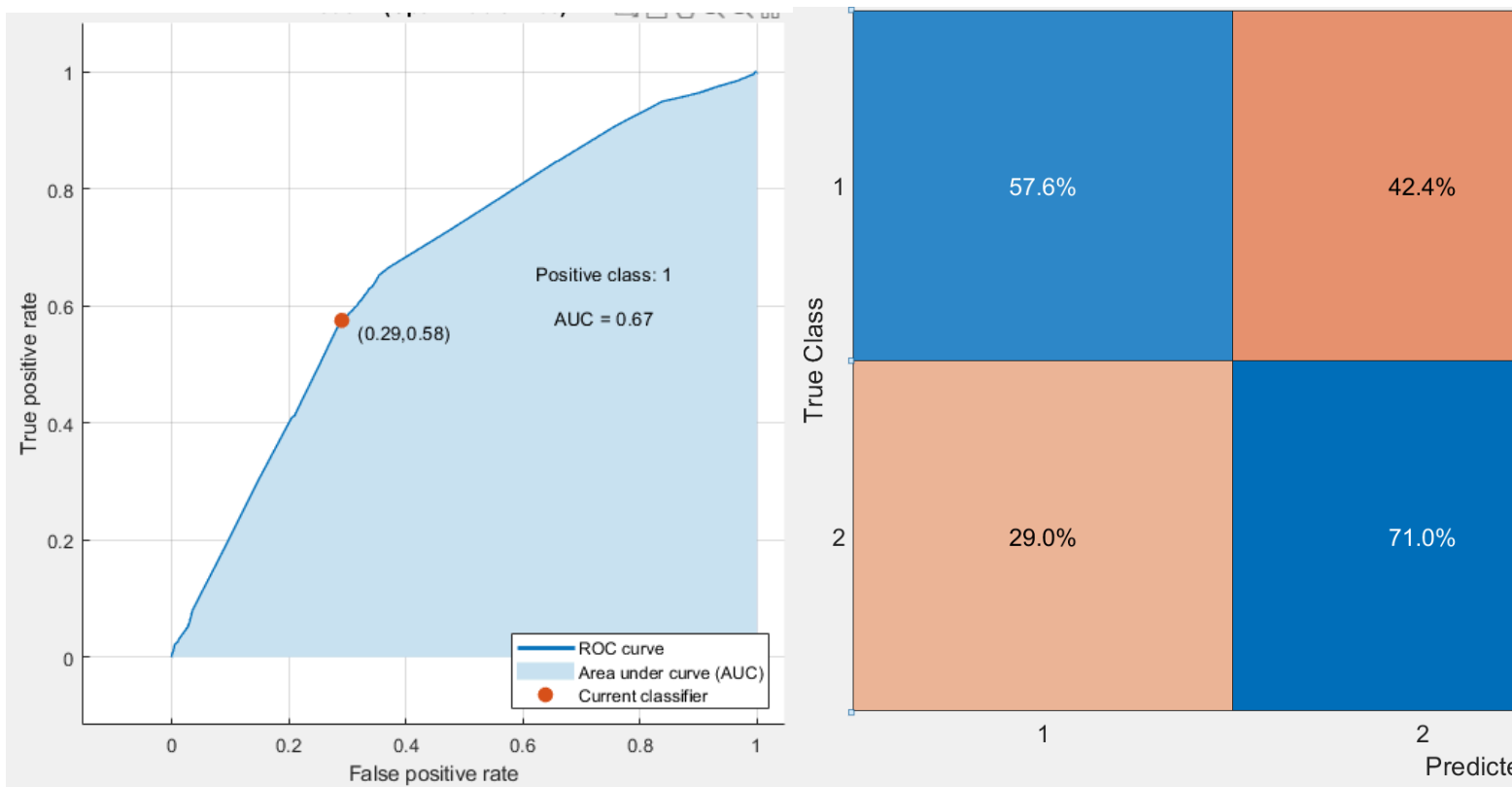
## Presentacion de resultados de los modelos entrenados con el data set completo y el de características reducidas

*El código para entrenar los modelos se encuentra en la carpeta "models".*

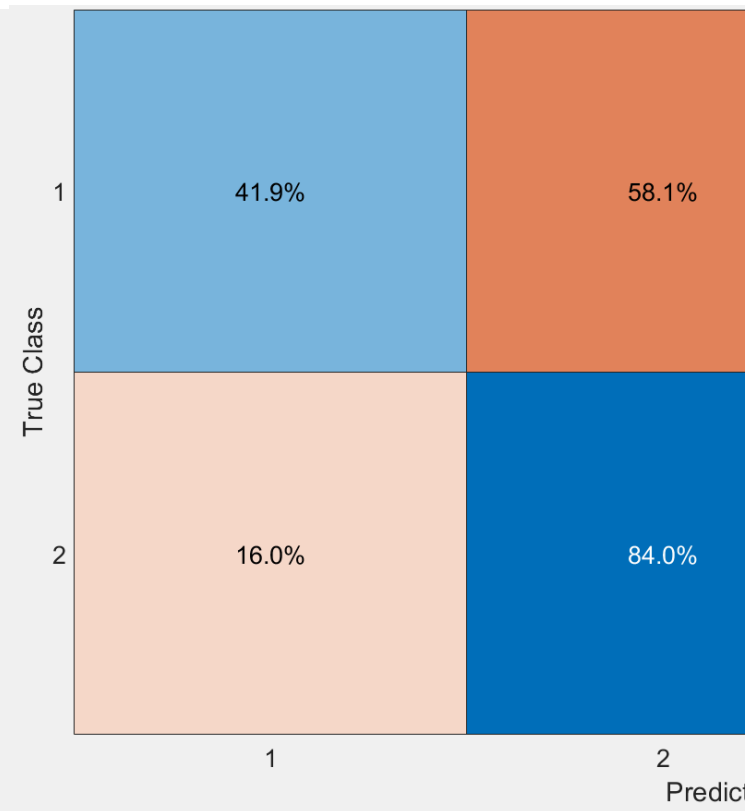
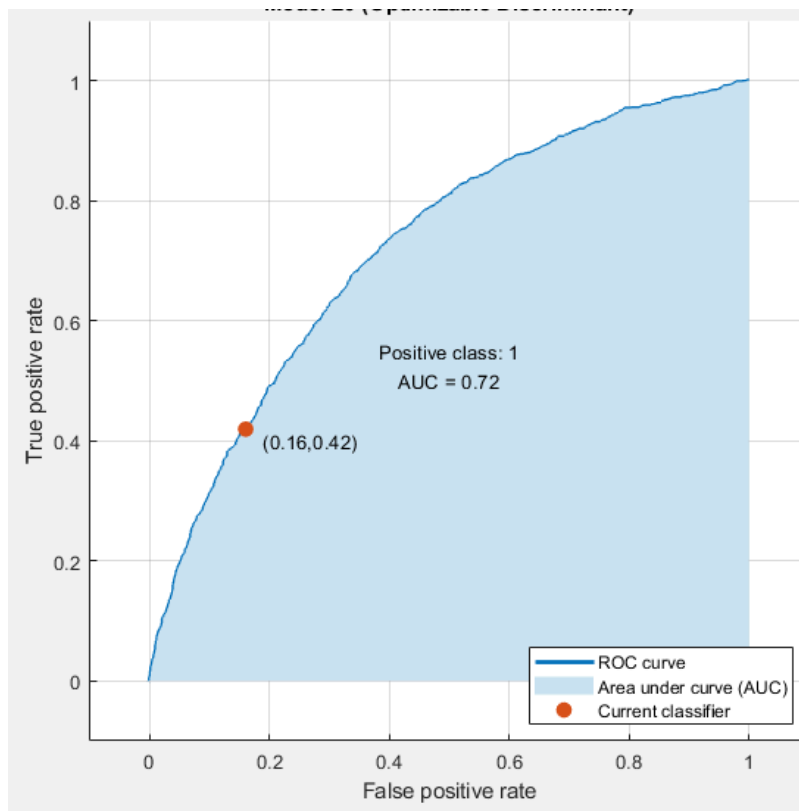
### Arboles de decision:



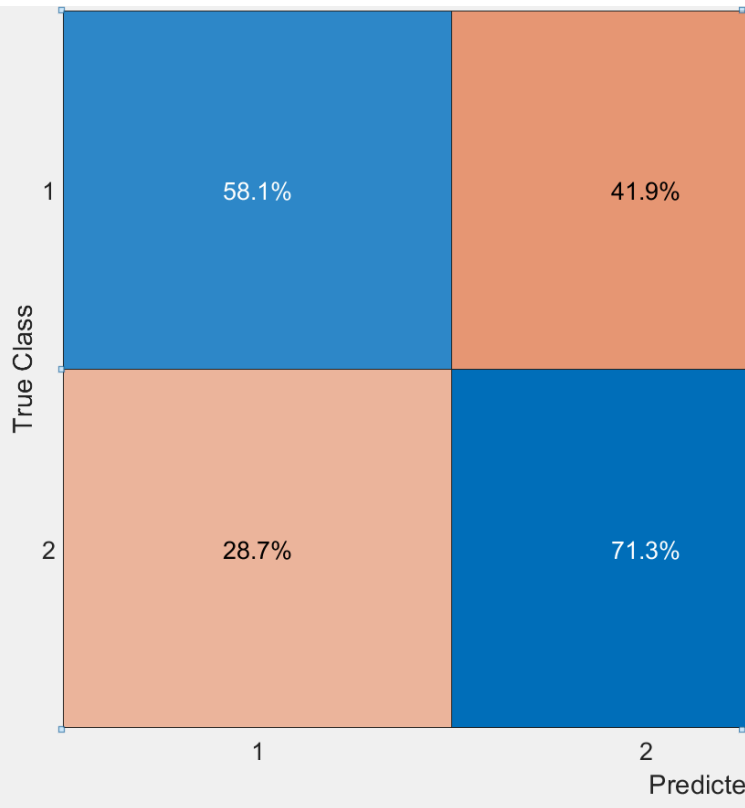
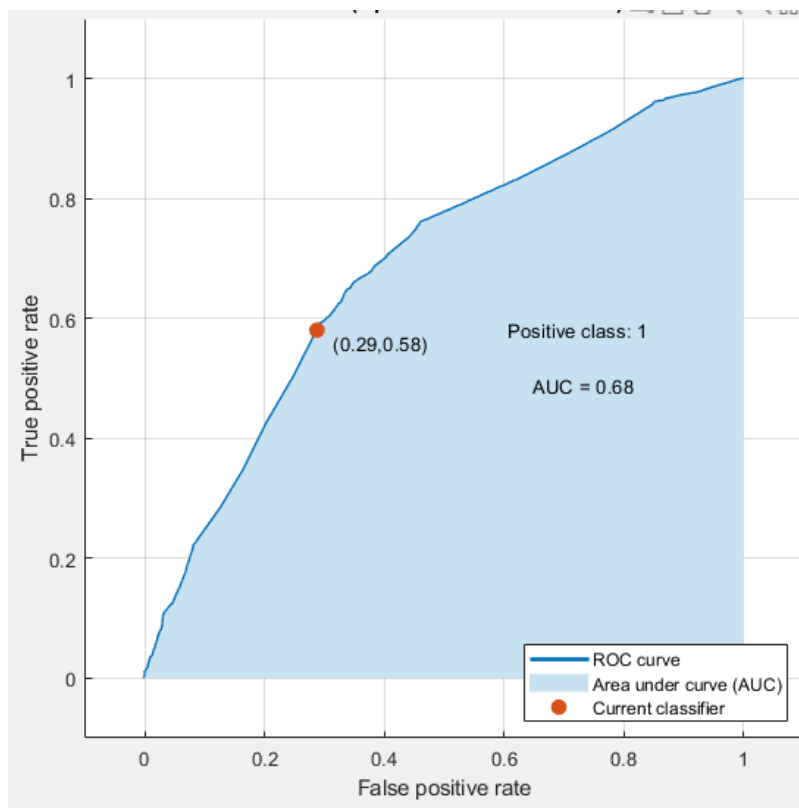
### Arboles de decision( 7 características):



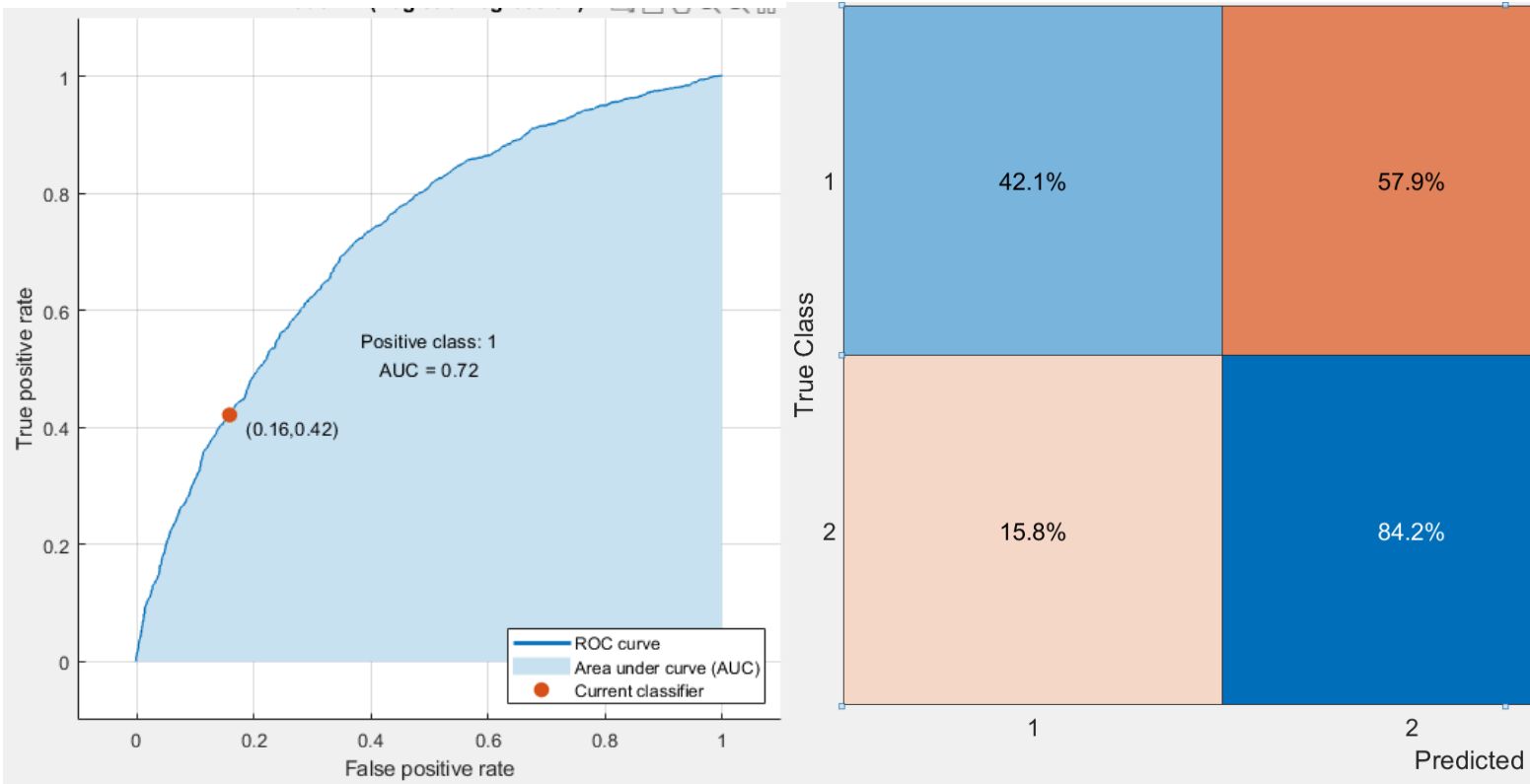
## Discriminante:



## Discriminante(7 características):

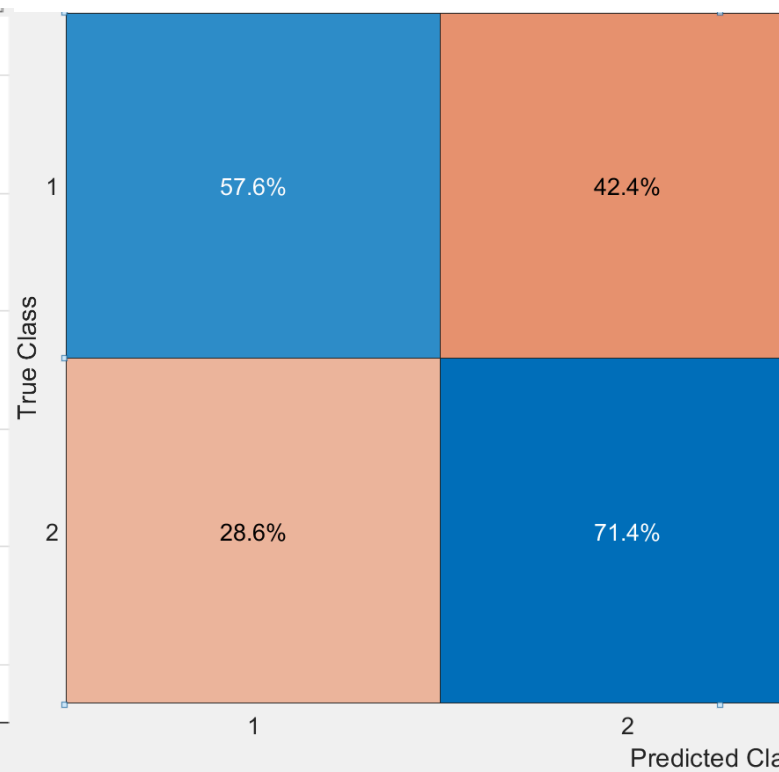
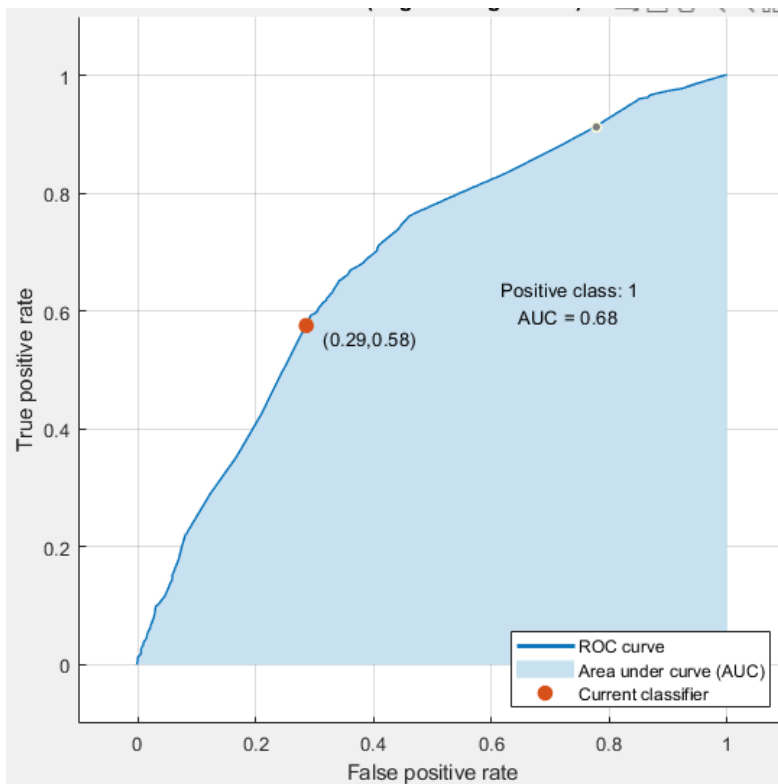


Logistic regression:

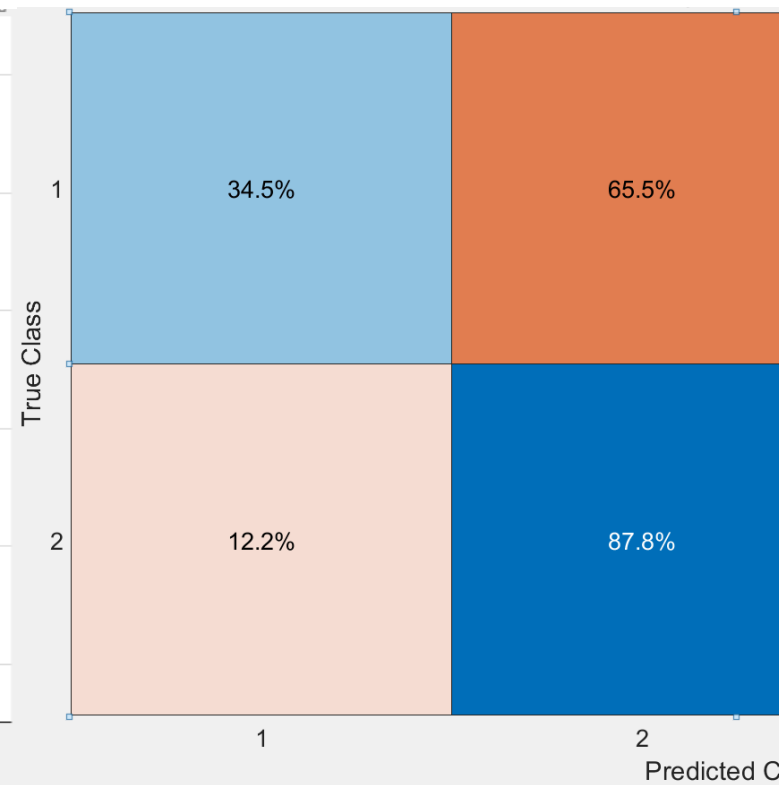
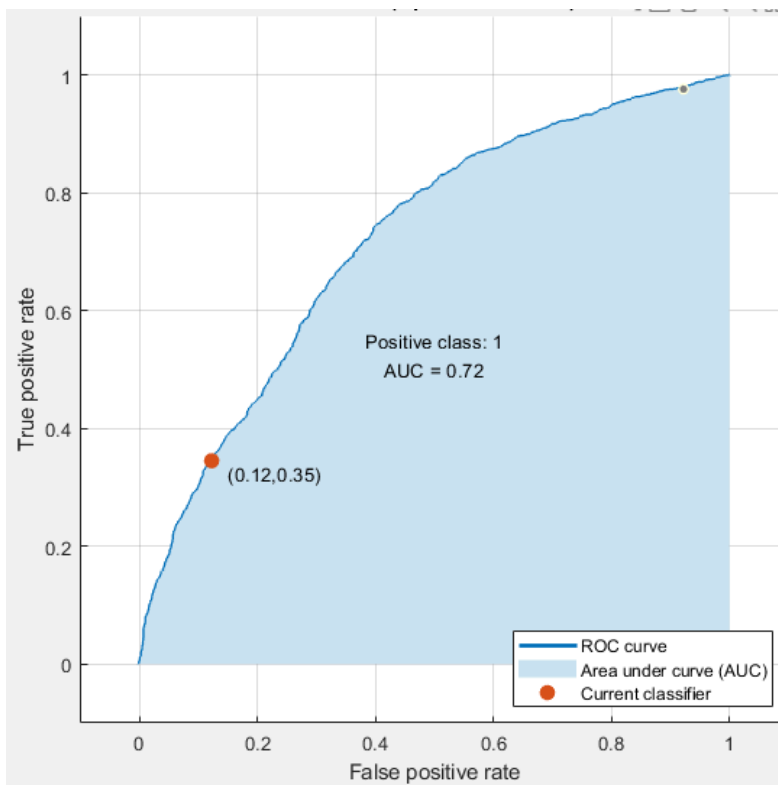


Logistic regression (7 caracteristicas):

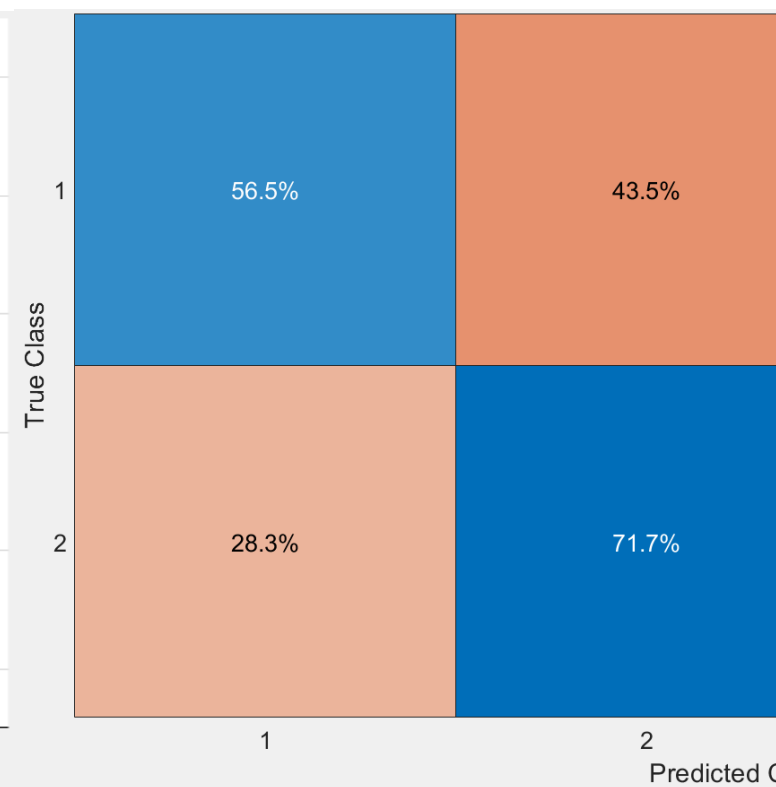
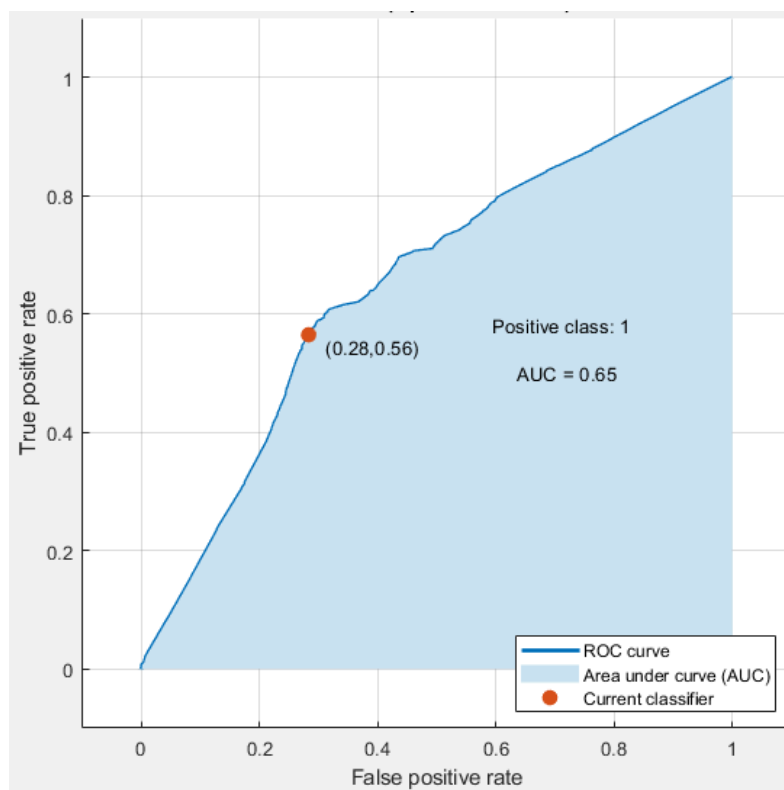




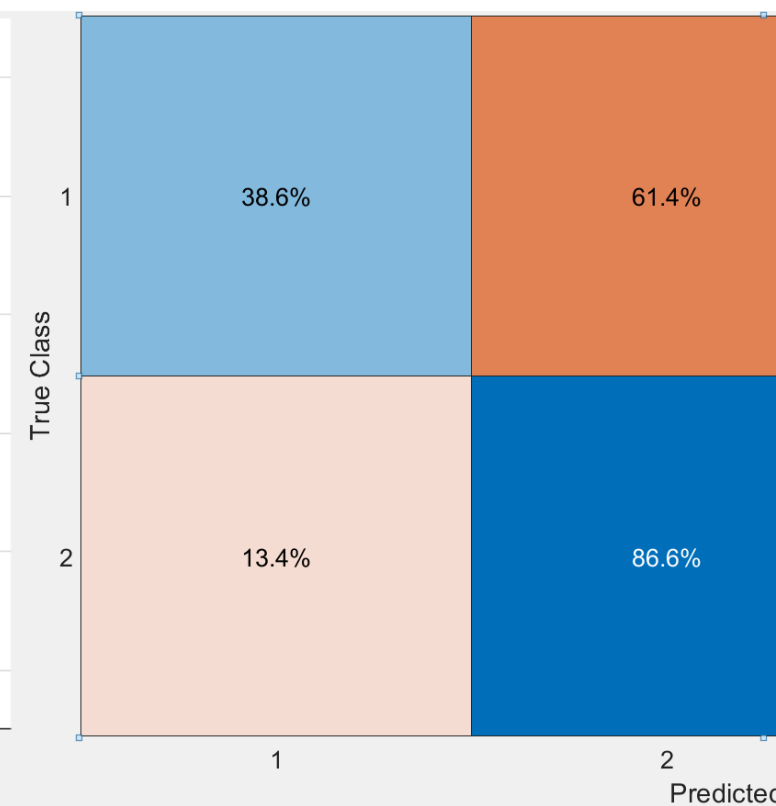
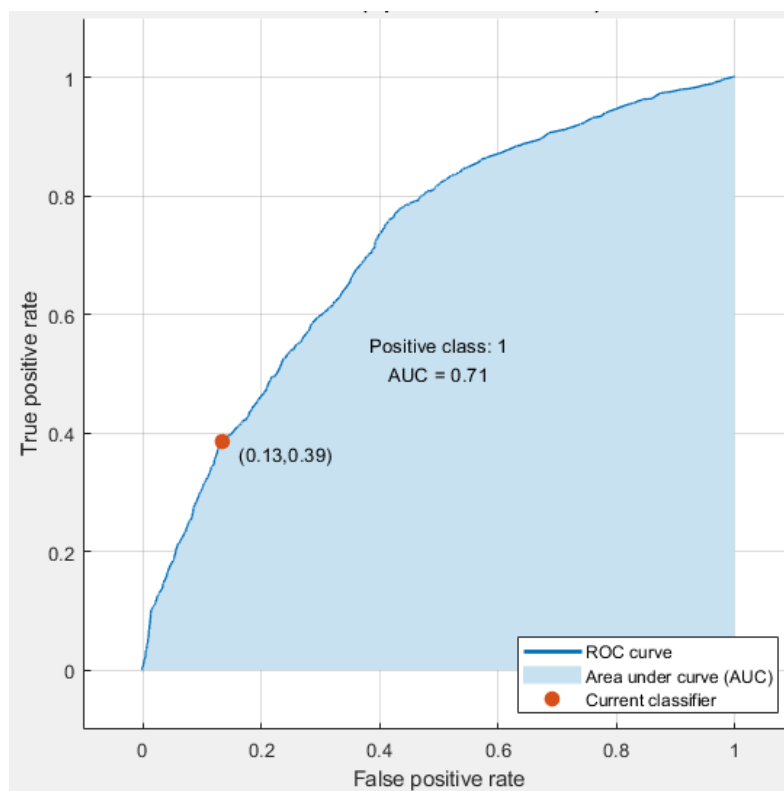
## SVM:



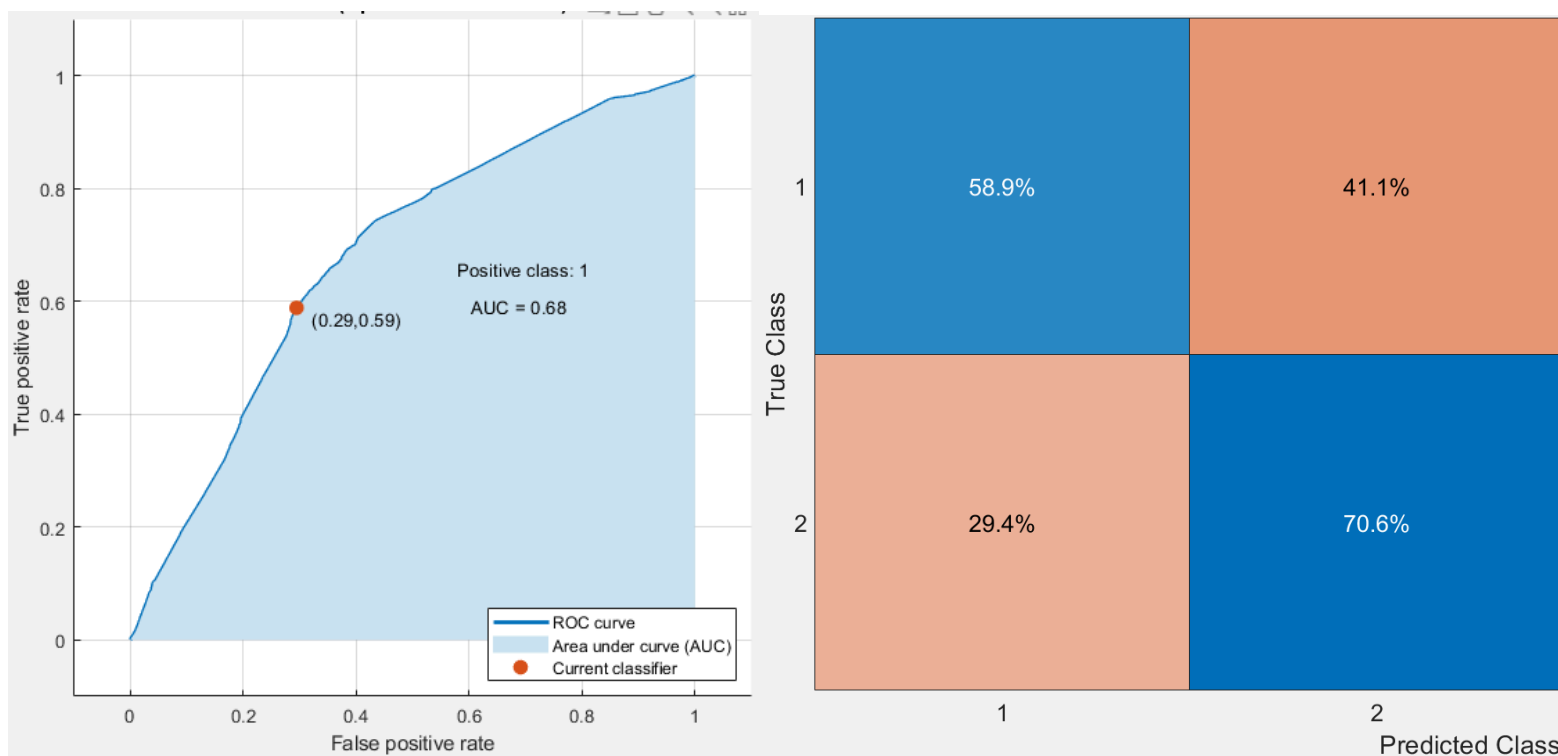
## SVM(7 características):



## Ensemble:



## Ensemble(7 características):



## Comparacion de modelos

Para este punto se tendran en cuenta varias cosas:

- La precision y exactitud del modelo, mientras mayor mejor, sin llegar a un caso de sobreajuste.
- El numero de parametros, en general es de interes obtener modelos que con un bajo numero de parametros sean capaces de cumplir con su objetivo a cabalidad, esto debido a que en un caso real es mas dificil y costoso, en terminos de dinero y tiempo obtener una cantidad grande de informacion. En este caso no se le dara mayor importancia a unos parametros sobre otros, solo sera de interes el numero de ellos.
- La complejidad del modelo, se preferiran modelos mas simples
- Un ultimo factor que se tendra en cuenta para preferir un modelo sobre otro, es la distribucion de falsos negativos hacia cierta clase particular, i.e. en este contexto no seria nada bueno identificar erroneamente a aquellas personas con tendencia repetitiva al intento de suicidio, mientras que identificar erroneamente a aquellos que en realidad no (falso positivo), seria mas aceptable.

## Provisional:

Teniendo en cuenta lo anterior lo que se busca es maximar la prediccion correcta de la etiqueta 1 {si a intentos previos de suicidio}, en este sentido la mayoria de modelos son deficientes. Pero en general los de 7 caracteristicas se comportan mejor que los completos. Asi, uno de los mejores seria el de discriminante lineal con 7 caracteristicas

ya que es el que mas se acerca a lo requerido, tienen poco numero de características y no hace parte de los modelos demasiado complejos

**Pendinete Escobar:**

- calcular score del modelos,
- poner analisis general de ROC y CM
- Hacer eso de la prediccion probabilistica y determinista

## Conclusiones

**El chorro:**

- *Hablar de porque nuestros modelos son tan malos(posibles razones:se hicieron las cosas mal? xd. Caracteristicas comunes para aquellos que son reincidentes por lo que es dificl separarlos de aquello que no?. Prefereiblemente una mayor cantidad de datos(y tambien mejor calidad ya que habia muchos datos faltants))*
- *Comentar que podria significar segun el contexto esas 7 caracteristicas mas importantes*
- *Como seria excelente de acuerdo a la propuesta inicial poder consefuir datasets con informacion de personas que han intentando previamente el suicidio como de aquella que no*

## Referencias (arreglar)

<https://www.mathworks.com/help/stats/feature-selection-and-feature-transformation.html>

<https://www.mathworks.com/help/stats/train-classification-models-in-classification-learner-app.html>

<https://www.mathworks.com/help/stats/assess-classifier-performance.html>

<https://towardsdatascience.com/intuitive-hyperparameter-optimization-grid-search-random-search-and-bayesian-search-2102dbfaf5b>

<https://towardsdatascience.com/automated-machine-learning-hyperparameter-tuning-in-python-dfda59b72f8a>

<https://towardsdatascience.com/workflow-of-a-machine-learning-project-ec1dba419b94>

<https://www.mathworks.com/help/stats/feature-selection.html>

<https://towardsdatascience.com/understanding-confusion-matrix-a9ad42dcfd62>

<https://towardsdatascience.com/understanding-auc-roc-curve-68b2303cc9c5>

*Notas*

Modelos sugeridos por el profesor: Regresion Logistica, SVM, Arboles de decision, Redes neuronales, LMP, Random Forest

*Intent\_prev{1 = SI; 2 = NO}*

*Para traducir graficas{Mínimo error de clasificación,Iteración,Mínimo error de clasificación observado,Hiperparámetros de error mínimo}*