

```
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 data set-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.

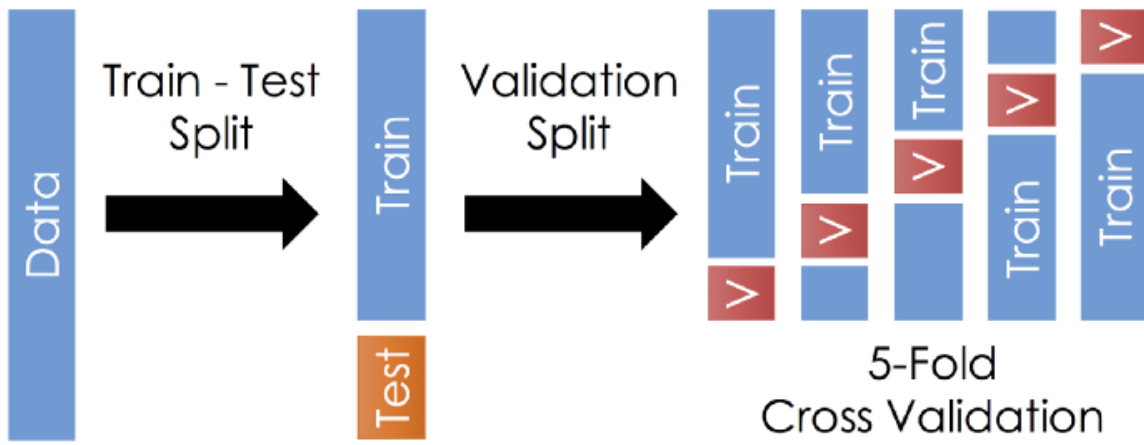
En el avance 1 se obtuvieron 4 datasets después del proceso de limpieza, los cuales se mencionan a continuación:

- 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\_fem\_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)
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 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 validación cruzada con "k-folds" (k=5), así, el valor de la exactitud presentado corresponde a la exactitud de validación 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 Last change: Fine Tree	Accuracy: 63.8% 28/28 features	<b>1.14</b> ☆ SVM Last change: Coarse Gaussian SVM	Accuracy: 6 28/28 features
<b>1.2</b> ☆ Tree Last change: Medium Tree	Accuracy: 65.1% 28/28 features	<b>1.15</b> ☆ KNN Last change: Fine KNN	Accuracy: 5 28/28 features
<b>1.3</b> ☆ Tree Last change: Coarse Tree	Accuracy: 66.7% 28/28 features	<b>1.16</b> ☆ KNN Last change: Medium KNN	Accuracy: 6 28/28 features
<b>1.4</b> ☆ Linear Discriminant Last change: Linear Discriminant	Accuracy: 66.9% 28/28 features	<b>1.17</b> ☆ KNN Last change: Coarse KNN	Accuracy: 6 28/28 features
<b>1.5</b> ☆ Quadratic Discriminant Last change: Quadratic Discriminant	<b>Failed</b> 28/28 features	<b>1.18</b> ☆ KNN Last change: Cosine KNN	Accuracy: 6 28/28 features
<b>1.6</b> ☆ Logistic Regression Last change: Logistic Regression	Accuracy: <b>67.4%</b> 28/28 features	<b>1.19</b> ☆ KNN Last change: Cubic KNN	Accuracy: 6 28/28 features
<b>1.7</b> ☆ Naive Bayes Last change: Gaussian Naive Bayes	Accuracy: 62.2% 28/28 features	<b>1.20</b> ☆ KNN Last change: Weighted KNN	Accuracy: 6 28/28 features
<b>1.8</b> ☆ Naive Bayes Last change: Kernel Naive Bayes	Accuracy: 62.0% 28/28 features	<b>1.21</b> ☆ Ensemble Last change: Boosted Trees	Accuracy: 6 28/28 features
<b>1.9</b> ☆ SVM Last change: Linear SVM	Accuracy: 66.0% 28/28 features	<b>1.22</b> ☆ Ensemble Last change: Bagged Trees	Accuracy: 6 28/28 features
<b>1.10</b> ☆ SVM Last change: Quadratic SVM	Accuracy: 65.3% 28/28 features	<b>1.23</b> ☆ Ensemble Last change: Subspace Discriminant	Accuracy: 6 28/28 features
<b>1.11</b> ☆ SVM Last change: Cubic SVM	Accuracy: 63.8% 28/28 features	<b>1.24</b> ☆ Ensemble Last change: Subspace KNN	Accuracy: 6 28/28 features
<b>1.12</b> ☆ SVM Last change: Fine Gaussian SVM	Accuracy: 61.8% 28/28 features	<b>1.25</b> ☆ Ensemble Last change: RUSBoosted Trees	Accuracy: 6 28/28 features
<b>1.13</b> ☆ SVM Last change: Medium Gaussian SVM	Accuracy: 65.8% 28/28 features	<b>2</b> ☆ Quadratic Discriminant Last change: 'Covariance structure' ...	Accuracy: 6 28/28 features

Resultados cds\_imputed

<b>1.1</b> ☆ Tree Last change: Fine Tree	Accuracy: 64.5% 33/33 features	<b>1.14</b> ☆ SVM Last change: Coarse Gaussian SVM	Accuracy: 6 33/33 fe
<b>1.2</b> ☆ Tree Last change: Medium Tree	Accuracy: 66.1% 33/33 features	<b>1.15</b> ☆ KNN Last change: Fine KNN	Accuracy: 6 33/33 fe
<b>1.3</b> ☆ Tree Last change: Coarse Tree	Accuracy: 66.8% 33/33 features	<b>1.16</b> ☆ KNN Last change: Medium KNN	Accuracy: 6 33/33 fe
<b>1.4</b> ☆ Linear Discriminant Last change: Linear Discriminant	Accuracy: <b>67.9%</b> 33/33 features	<b>1.17</b> ☆ KNN Last change: Coarse KNN	Accuracy: 6 33/33 fe
<b>1.5</b> ☆ Quadratic Discriminant Last change: Quadratic Discriminant	<b>Failed</b> 33/33 features	<b>1.18</b> ☆ KNN Last change: Cosine KNN	Accuracy: 6 33/33 fe
<b>1.6</b> ☆ Logistic Regression Last change: Logistic Regression	Accuracy: 67.8% 33/33 features	<b>1.19</b> ☆ KNN Last change: Cubic KNN	Accuracy: 6 33/33 fe
<b>1.7</b> ☆ Naive Bayes Last change: Gaussian Naive Bayes	Accuracy: 64.2% 33/33 features	<b>1.20</b> ☆ KNN Last change: Weighted KNN	Accuracy: 6 33/33 fe
<b>1.8</b> ☆ Naive Bayes Last change: Kernel Naive Bayes	Accuracy: 62.5% 33/33 features	<b>1.21</b> ☆ Ensemble Last change: Boosted Trees	Accuracy: 6 33/33 fe
<b>1.9</b> ☆ SVM Last change: Linear SVM	Accuracy: 66.8% 33/33 features	<b>1.22</b> ☆ Ensemble Last change: Bagged Trees	Accuracy: 6 33/33 fe
<b>1.10</b> ☆ SVM Last change: Quadratic SVM	Accuracy: 65.4% 33/33 features	<b>1.23</b> ☆ Ensemble Last change: Subspace Discriminant	Accuracy: 6 33/33 fe
<b>1.11</b> ☆ SVM Last change: Cubic SVM	Accuracy: 63.3% 33/33 features	<b>1.24</b> ☆ Ensemble Last change: Subspace KNN	Accuracy: 6 33/33 fe
<b>1.12</b> ☆ SVM Last change: Fine Gaussian SVM	Accuracy: 61.8% 33/33 features	<b>1.25</b> ☆ Ensemble Last change: RUSBoosted Trees	Accuracy: 6 33/33 fe
<b>1.13</b> ☆ SVM Last change: Medium Gaussian SVM	Accuracy: 65.5% 33/33 features	<b>2</b> ☆ Quadratic Discriminant Last change: 'Covariance structure' ...	Accuracy: 6 33/33 fe

### 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 Last change: Fine Tree	Accuracy: 53.9% 33/33 features	<b>1.15</b> ☆ KNN Last change: Fine KNN	Accuracy: 53.6% 33/33 features
<b>1.2</b> ☆ Tree Last change: Medium Tree	Accuracy: 60.5% 33/33 features	<b>1.16</b> ☆ KNN Last change: Medium KNN	Accuracy: 60.8% 33/33 features
<b>1.3</b> ☆ Tree Last change: Coarse Tree	Accuracy: <b>64.4%</b> 33/33 features	<b>1.17</b> ☆ KNN Last change: Coarse KNN	Accuracy: 60.8% 33/33 features
<b>1.4</b> ☆ Linear Discriminant Last change: Linear Discriminant	<b>Failed</b> 33/33 features	<b>1.18</b> ☆ KNN Last change: Cosine KNN	Accuracy: 61.2% 33/33 features
<b>1.5</b> ☆ Quadratic Discriminant Last change: Quadratic Discriminant	<b>Failed</b> 33/33 features	<b>1.19</b> ☆ KNN Last change: Cubic KNN	Accuracy: 60.2% 33/33 features
<b>1.6</b> ☆ Logistic Regression Last change: Logistic Regression	Accuracy: 61.8% 33/33 features	<b>1.20</b> ☆ KNN Last change: Weighted KNN	Accuracy: 57.7% 33/33 features
<b>1.7</b> ☆ Naive Bayes Last change: Gaussian Naive Bayes	<b>Failed</b> 33/33 features	<b>1.21</b> ☆ Ensemble Last change: Boosted Trees	Accuracy: 59.1% 33/33 features
<b>1.8</b> ☆ Naive Bayes Last change: Kernel Naive Bayes	Accuracy: 61.1% 33/33 features	<b>1.22</b> ☆ Ensemble Last change: Bagged Trees	Accuracy: 58.3% 33/33 features
<b>1.9</b> ☆ SVM Last change: Linear SVM	Accuracy: 61.2% 33/33 features	<b>1.23</b> ☆ Ensemble Last change: Subspace Discriminant	Accuracy: 63.4% 33/33 features
<b>1.10</b> ☆ SVM Last change: Quadratic SVM	Accuracy: 57.7% 33/33 features	<b>1.24</b> ☆ Ensemble Last change: Subspace KNN	Accuracy: 57.3% 33/33 features
<b>1.11</b> ☆ SVM Last change: Cubic SVM	Accuracy: 57.3% 33/33 features	<b>1.25</b> ☆ Ensemble Last change: RUSBoosted Trees	Accuracy: 58.3% 33/33 features
<b>1.12</b> ☆ SVM Last change: Fine Gaussian SVM	Accuracy: 58.9% 33/33 features	<b>2</b> ☆ Linear Discriminant Last change: 'Covariance structure' ...	Accuracy: 61.5% 33/33 features
<b>1.13</b> ☆ SVM Last change: Medium Gaussian SVM	Accuracy: 63.7% 33/33 features	<b>3</b> ☆ Quadratic Discriminant Last change: 'Covariance structure' ...	Accuracy: 60.6% 33/33 features
<b>1.14</b> ☆ SVM Last change: Coarse Gaussian SVM	Accuracy: 61.1% 33/33 features	<b>4</b> ☆ Naive Bayes Last change: Removed 3 features	Accuracy: 53.6% 29/33 features

## Resultados cds\_few\_minus\_alcohol



1.1 ☆ Tree Last change: Fine Tree	Accuracy: 54.2% 32/32 features	1.15 ☆ KNN Last change: Fine KNN	Accuracy: 53.7% 32/32 features
1.2 ☆ Tree Last change: Medium Tree	Accuracy: 55.6% 32/32 features	1.16 ☆ KNN Last change: Medium KNN	Accuracy: 55.1% 32/32 features
1.3 ☆ Tree Last change: Coarse Tree	Accuracy: 55.6% 32/32 features	1.17 ☆ KNN Last change: Coarse KNN	Accuracy: 56.2% 32/32 features
1.4 ☆ Linear Discriminant Last change: Linear Discriminant	<b>Failed</b> 32/32 features	1.18 ☆ KNN Last change: Cosine KNN	Accuracy: 55.4% 32/32 features
1.5 ☆ Quadratic Discriminant Last change: Quadratic Discriminant	<b>Failed</b> 32/32 features	1.19 ☆ KNN Last change: Cubic KNN	Accuracy: 54.7% 32/32 features
1.6 ☆ Logistic Regression Last change: Logistic Regression	Accuracy: 59.0% 32/32 features	1.20 ☆ KNN Last change: Weighted KNN	Accuracy: 55.0% 32/32 features
1.7 ☆ Naive Bayes Last change: Gaussian Naive Bayes	<b>Failed</b> 32/32 features	1.21 ☆ Ensemble Last change: Boosted Trees	Accuracy: 59.2% 32/32 features
1.8 ☆ Naive Bayes Last change: Kernel Naive Bayes	Accuracy: 55.0% 32/32 features	1.22 ☆ Ensemble Last change: Bagged Trees	Accuracy: 56.6% 32/32 features
1.9 ☆ SVM Last change: Linear SVM	Accuracy: 58.5% 32/32 features	1.23 ☆ Ensemble Last change: Subspace Discriminant	Accuracy: 58.1% 32/32 features
1.10 ☆ SVM Last change: Quadratic SVM	Accuracy: 55.3% 32/32 features	1.24 ☆ Ensemble Last change: Subspace KNN	Accuracy: 54.0% 32/32 features
1.11 ☆ SVM Last change: Cubic SVM	Accuracy: 53.3% 32/32 features	1.25 ☆ Ensemble Last change: RUSBoosted Trees	Accuracy: 57.3% 32/32 features
1.12 ☆ SVM Last change: Fine Gaussian SVM	Accuracy: 54.3% 32/32 features	2 ☆ Linear Discriminant Last change: 'Covariance structure' ...	Accuracy: <b>60.2%</b> 32/32 features
1.13 ☆ SVM Last change: Medium Gaussian SVM	Accuracy: 56.5% 32/32 features	3 ☆ Quadratic Discriminant Last change: 'Covariance structure' ...	Accuracy: 55.0% 32/32 features
1.14 ☆ SVM Last change: Coarse Gaussian SVM	Accuracy: 58.4% 32/32 features	4 ☆ Naive Bayes Last change: Removed 3 features	Accuracy: 54.4% 28/32 features

Por motivos exploratorios se realizaron pruebas aplicándole PCA a los datos, pero los resultados en general fueron inferiores a los obtenidos sin esta transformación, por lo que esta transformación de los datos no será utilizada. (*¿Uno si debería hacer PCA en datos categóricos?*)

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 regresion, SVM (linear y coarse) y Ensemble (BoostTrees)-

## 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 Maximun Relevance), del cual se puede obtener el "ranking" de importancia de los predictores teniendo en cuentas la respuesta.

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

```
idx = fscmrnr(cds_imputed, 'inten_prev');  
most_signif_features = cds_imputed.Properties.VariableNames(idx(1:7)).'  
less_signif_features = cds_imputed.Properties.VariableNames(idx(end-4:end)).'
```

## Optimización de hiperparámetros

Para la optimización de hiperparámetros serán utilizados dos enfoques:

Para los modelos simples (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.

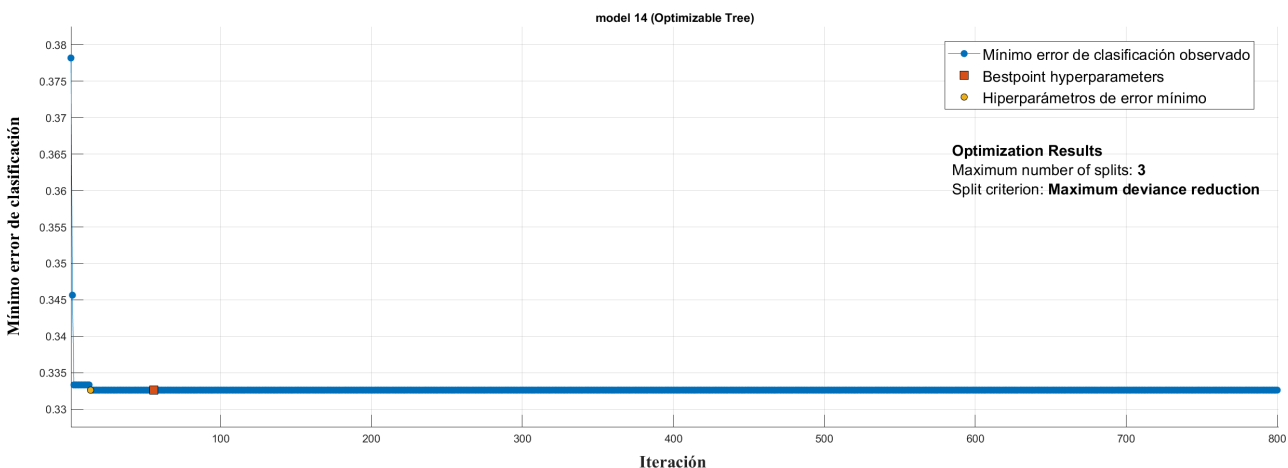
Como parte del proceso de optimización se obtiene el gráfico de error de clasificación mínimo en el cual se encuentran principalmente:

- 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 presentarán a continuación para cada modelo:

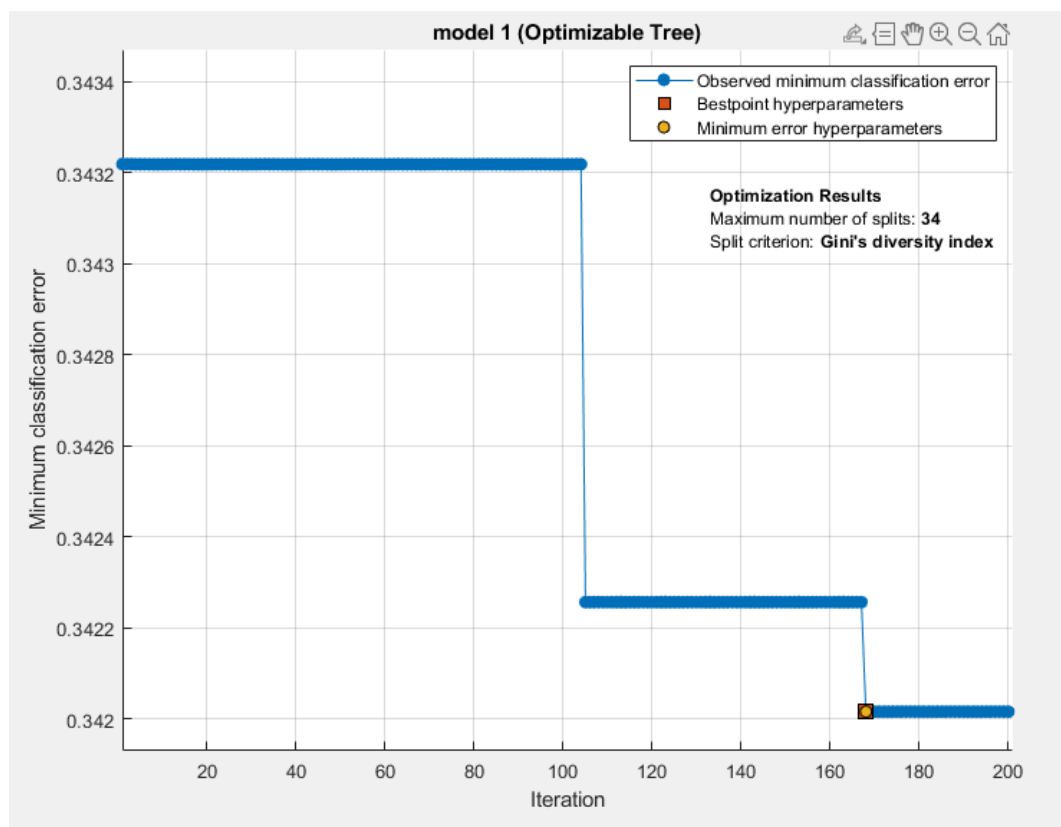
### Árboles de decisión:

En la anterior etapa mejores resultados fueron obtenidos con árboles de decisión gruesos (con poco número de splits), lo cual es confirmado con los resultados de este proceso



### Árboles de decisión (7 características):

Para este caso los mejores resultados fueron con árboles finos

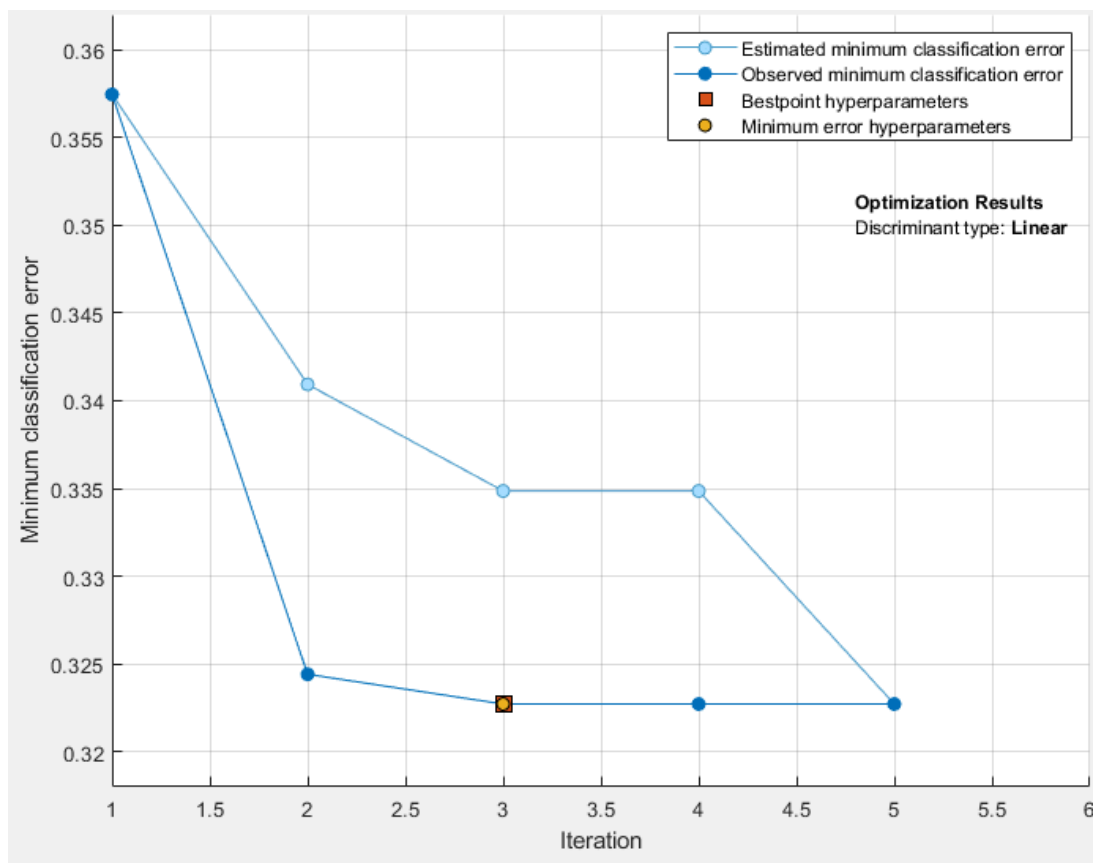


### Discriminante:

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

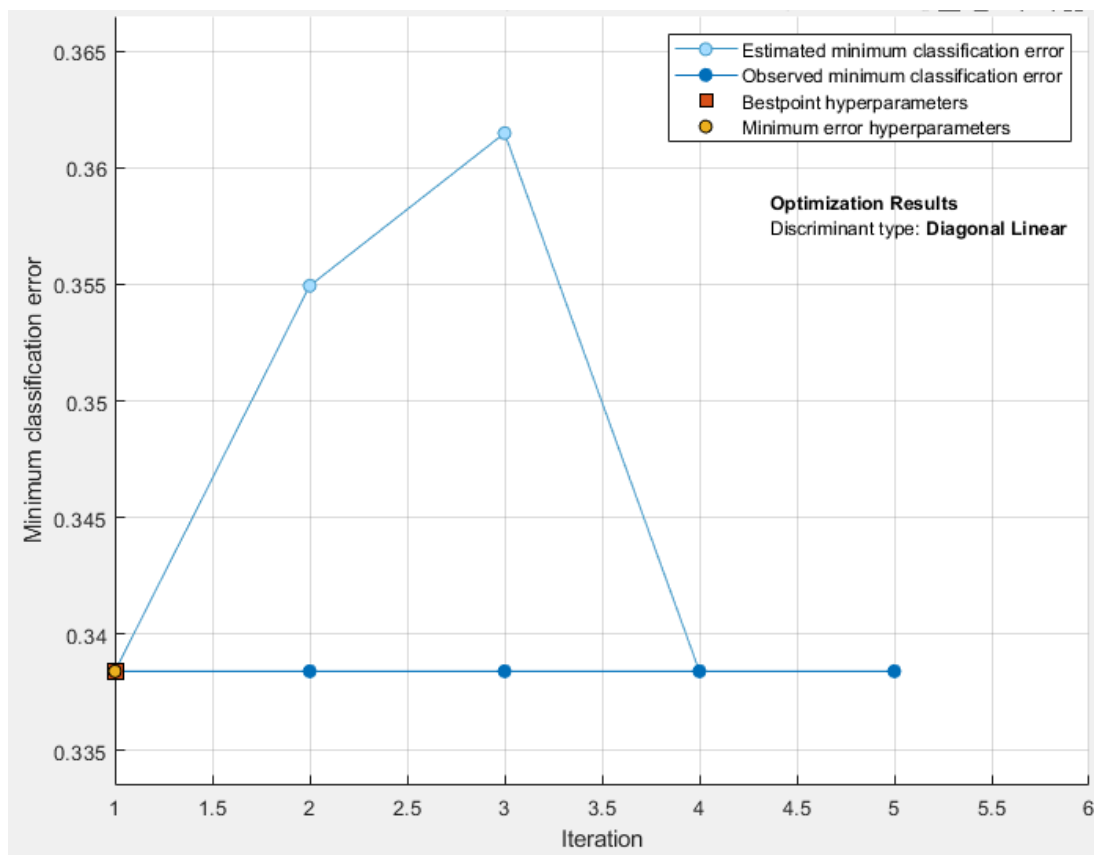
Las combinaciones disponibles para este tipo de modelos son pocas por lo que con pocas iteraciones es suficiente





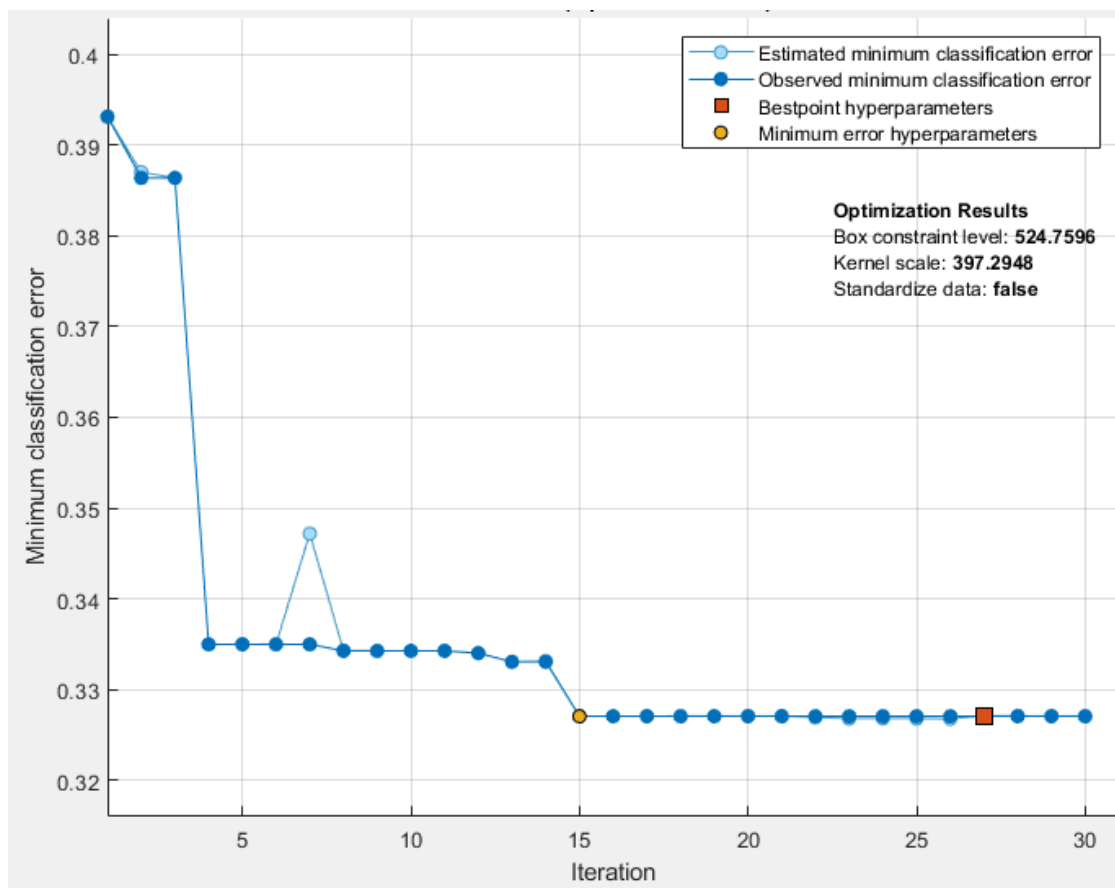
**Discriminante (7 características):**

Mejores resultados con diagonal linear

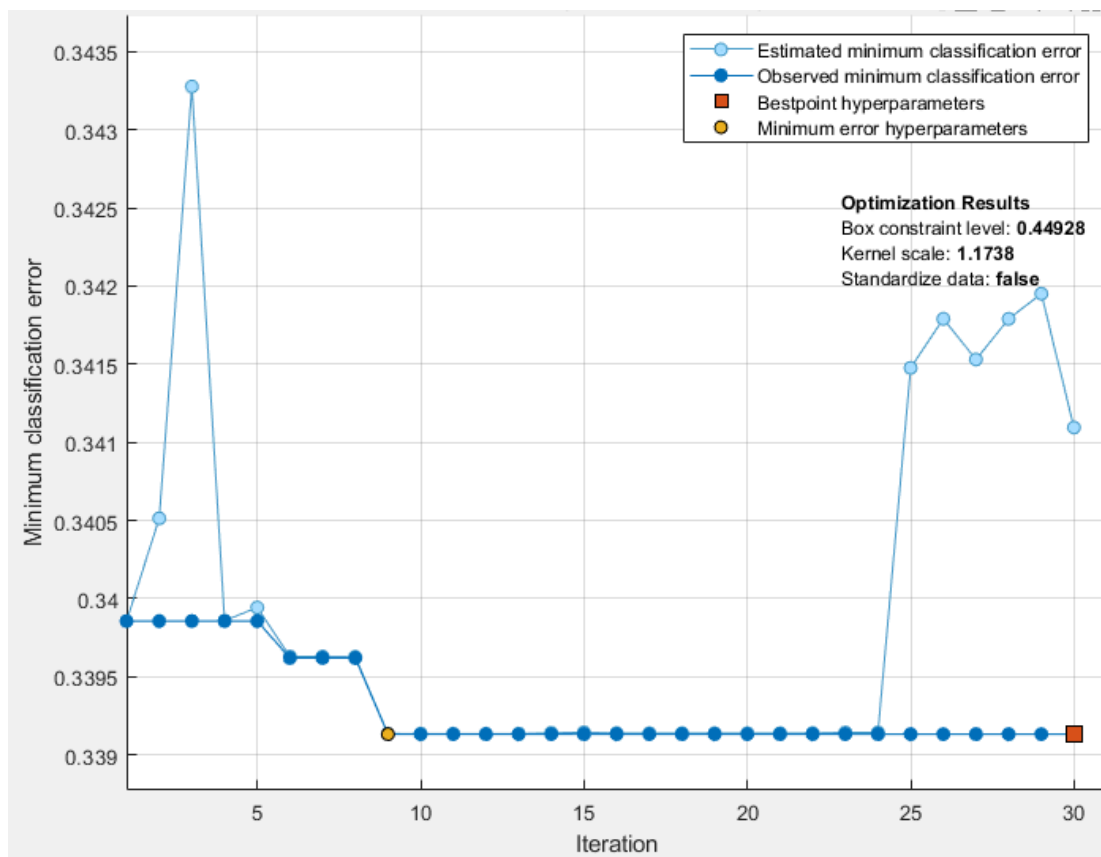


### SVM:

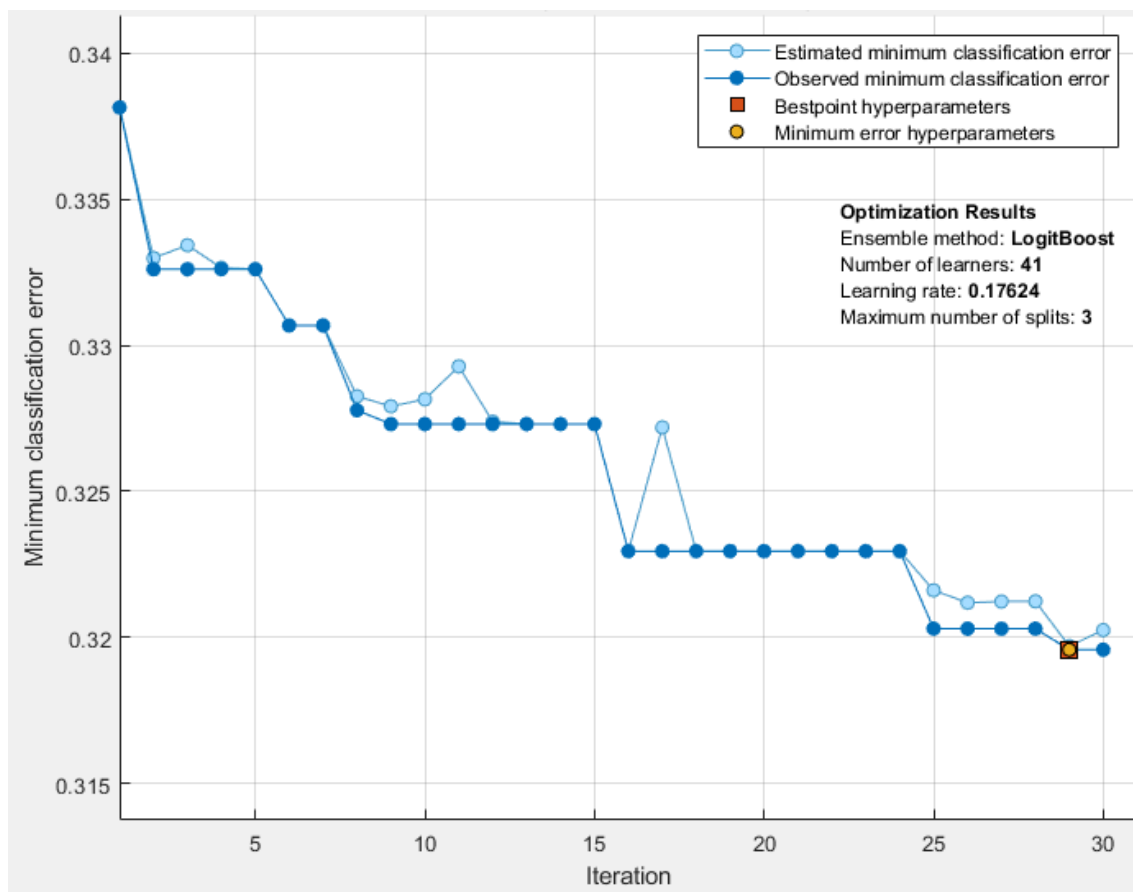
Para SVM los mejores resultados se obtuvieron para gaussiano "course" por lo que este será optimizado



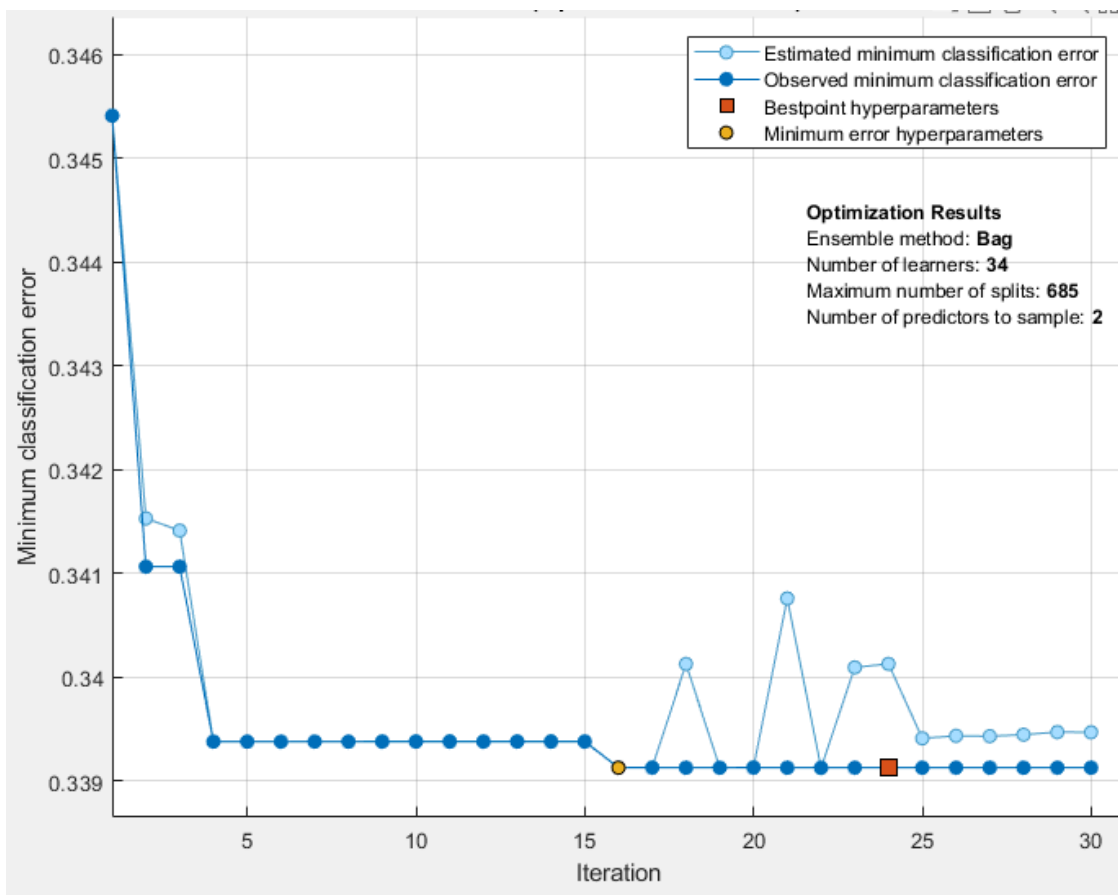
**SVM (7 características):**



Ensemble:



**Ensemble (7 características):**

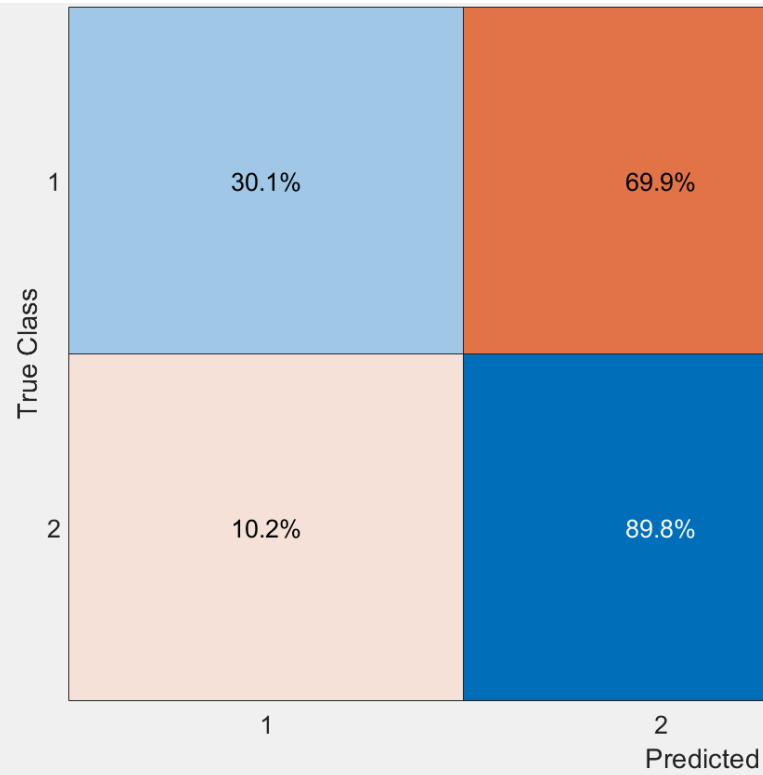
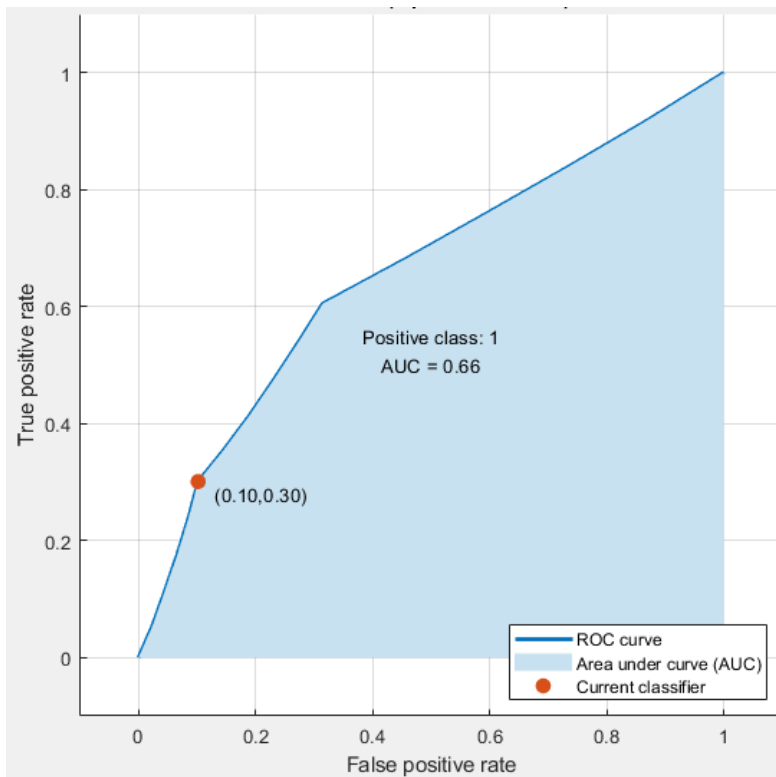


## Presentación de resultados de los modelos entrenados con el data set completo y el de características reducidas

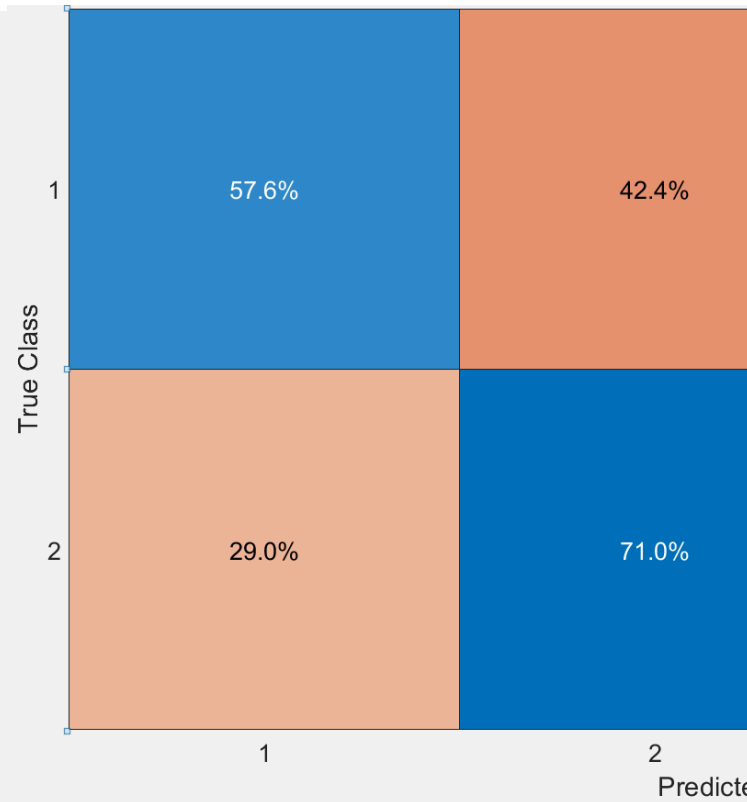
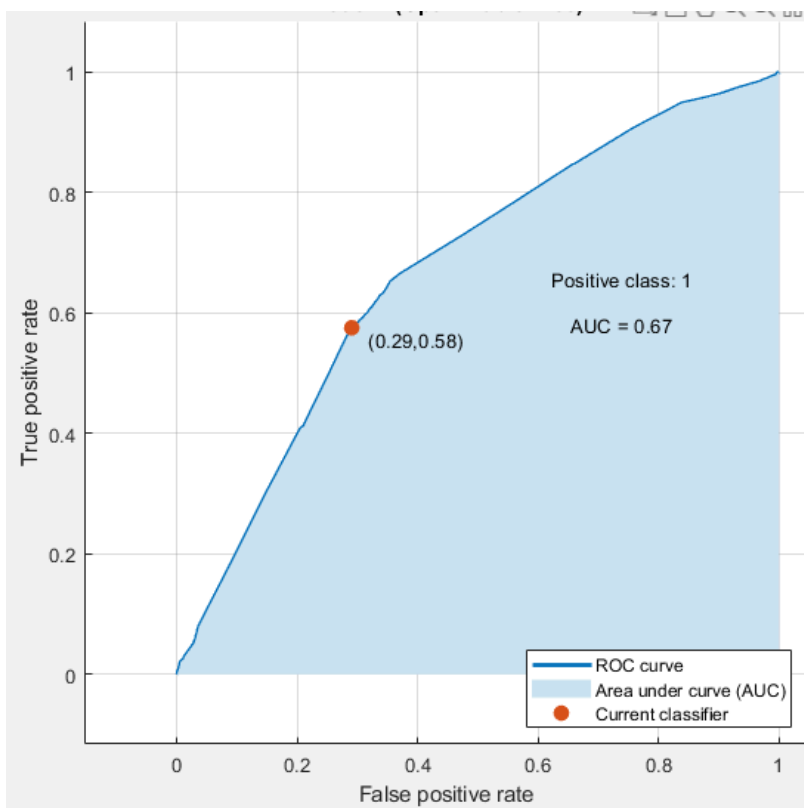
*El código para entrenar los modelos se encuentra en la carpeta "models".*

**Arboles de decisión:**

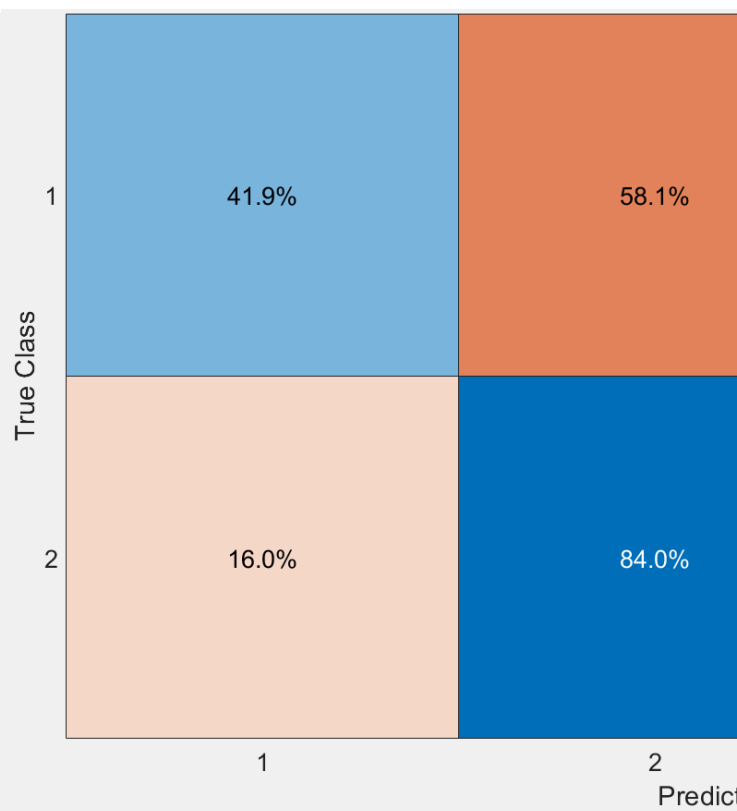
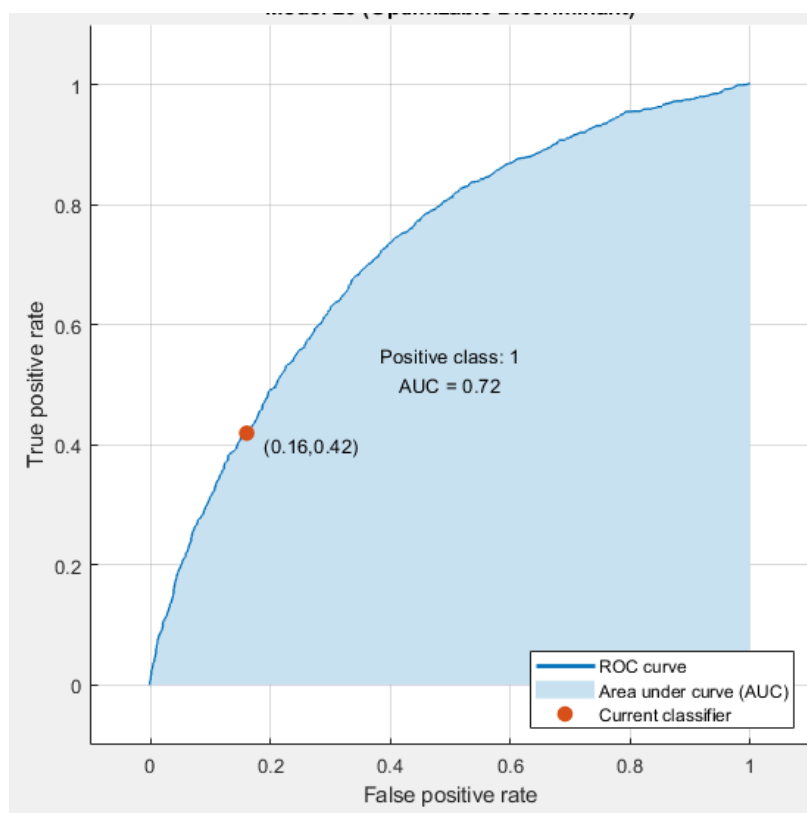




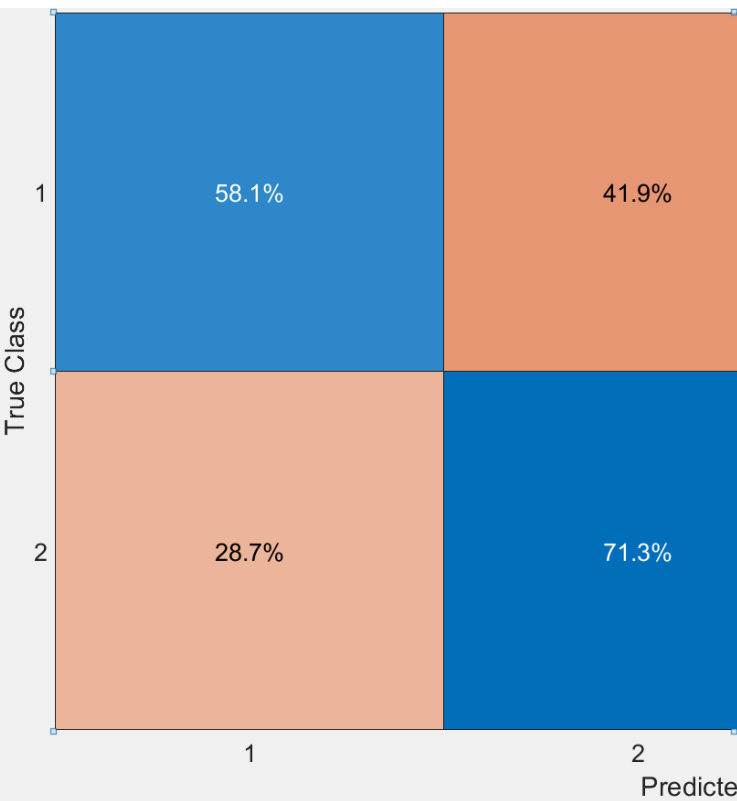
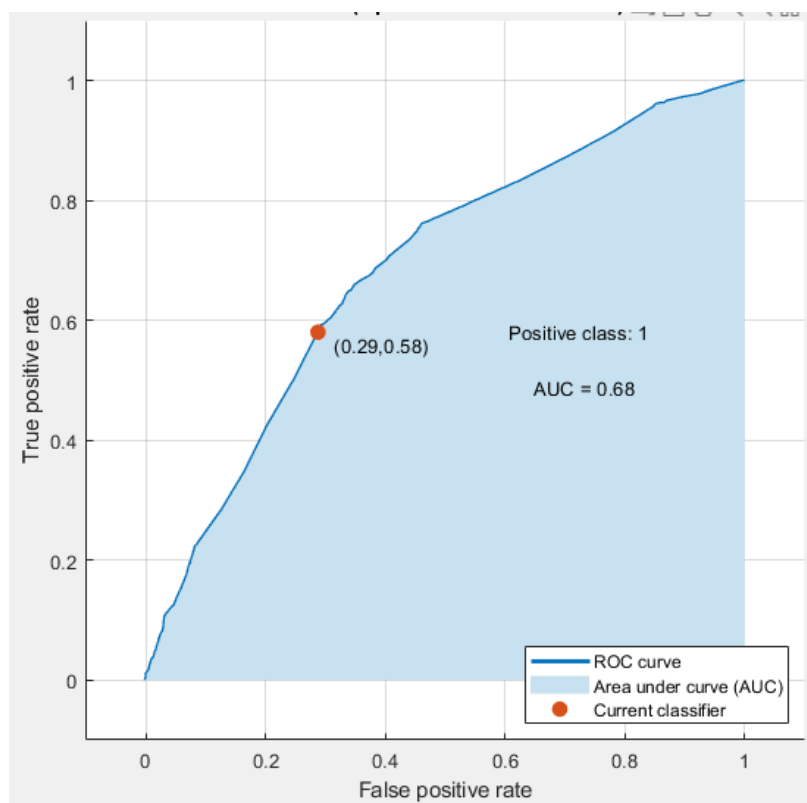
### Arboles de decisión (7 características):



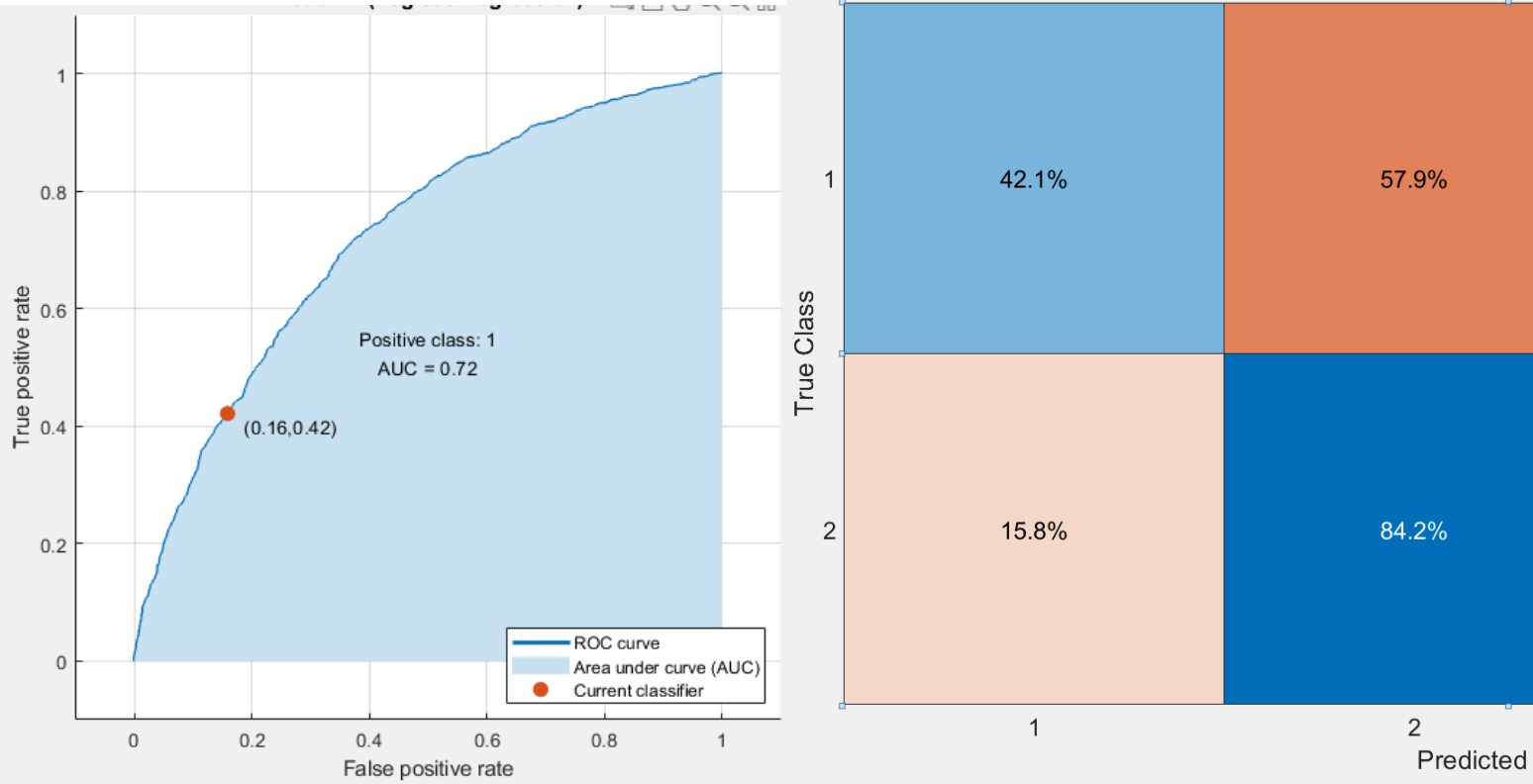
## Discriminante:



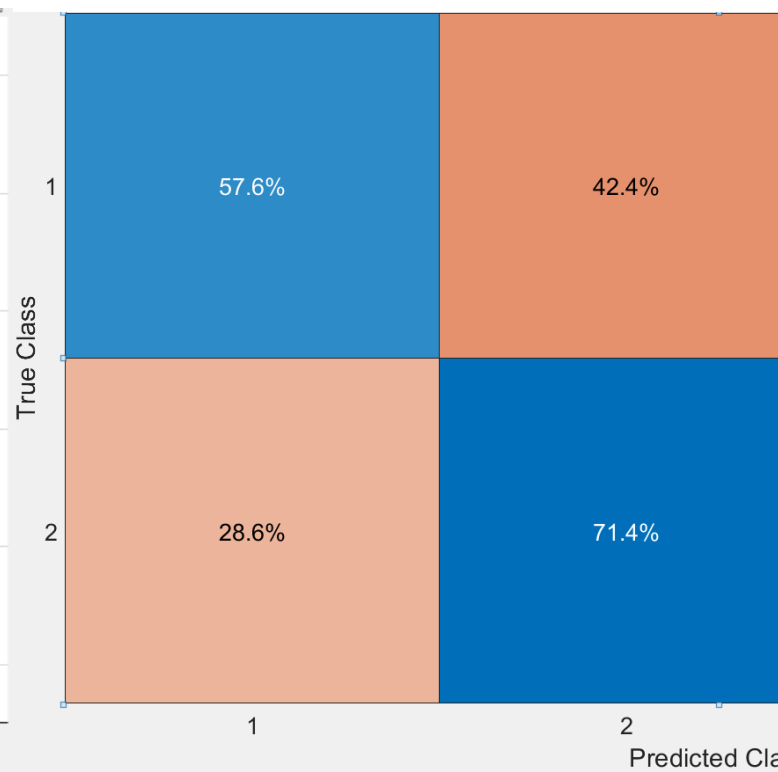
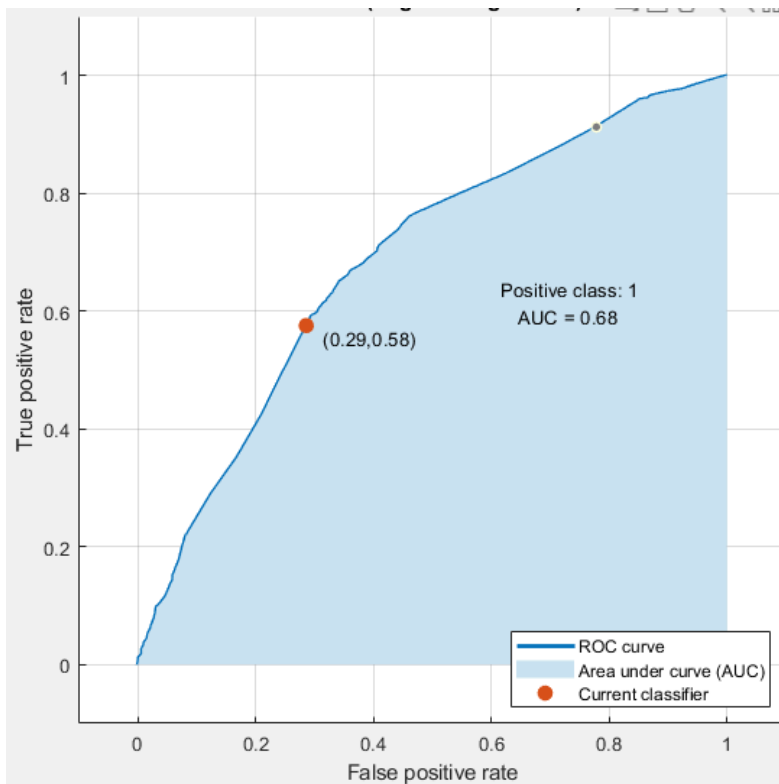
## Discriminante (7 características):



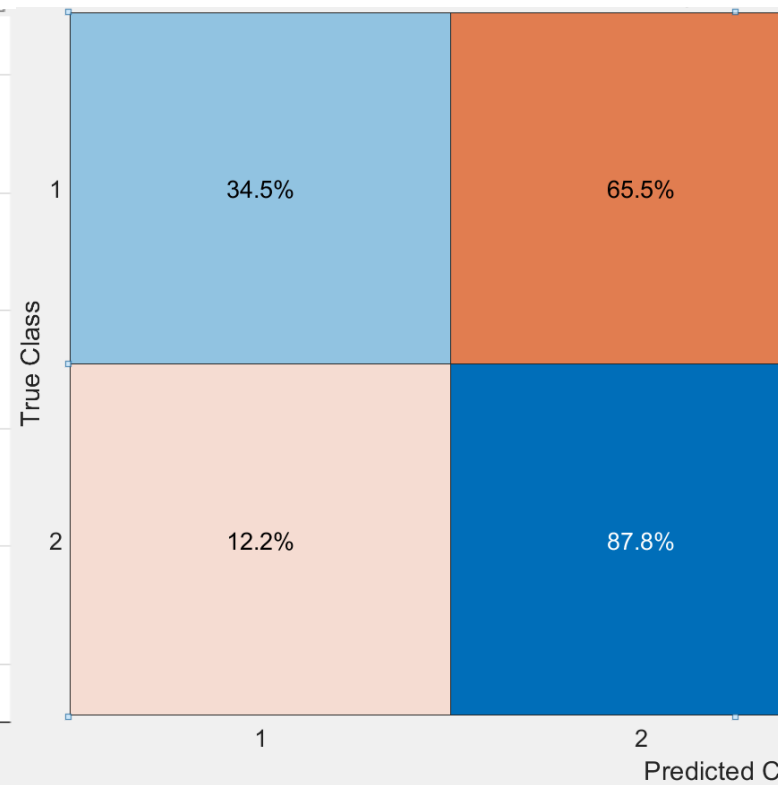
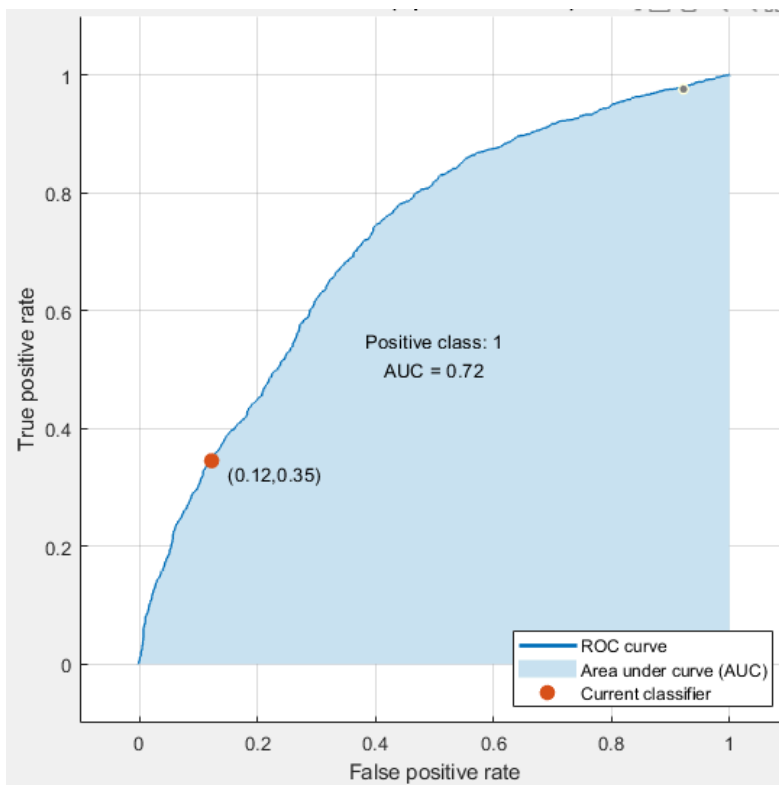
Logistic regression:



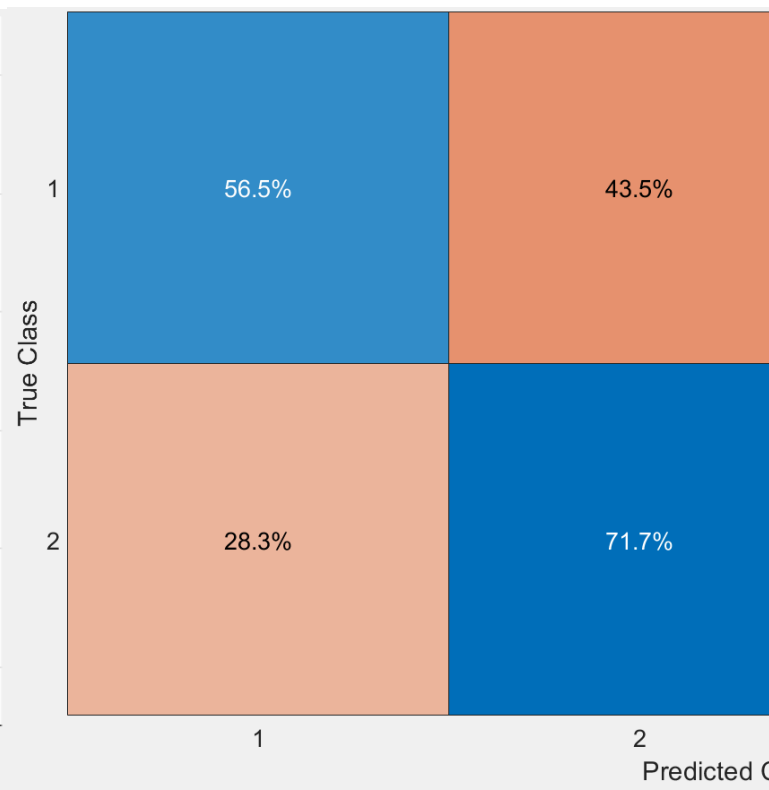
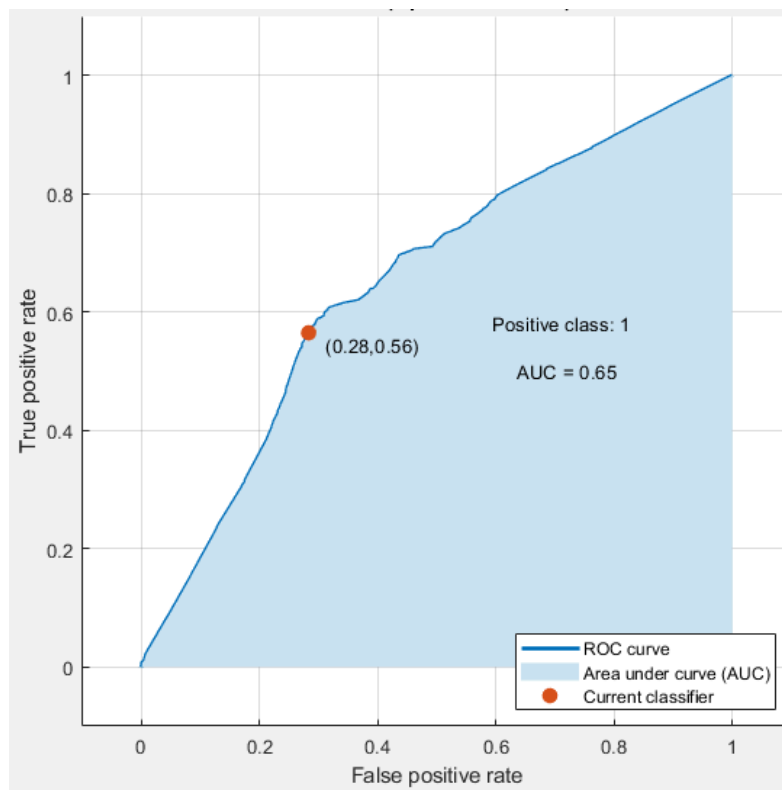
Logistic regression (7 características):



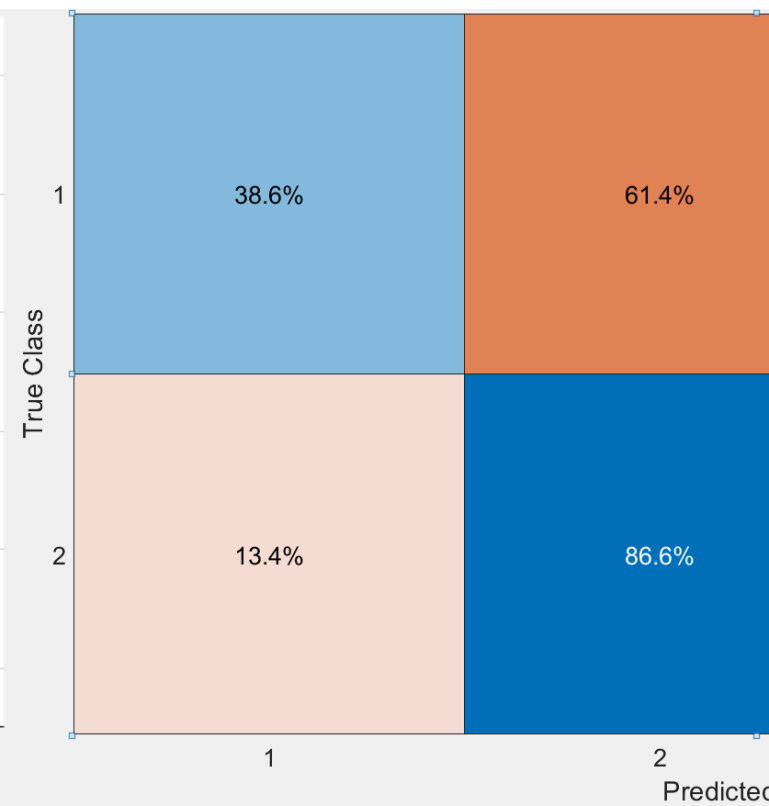
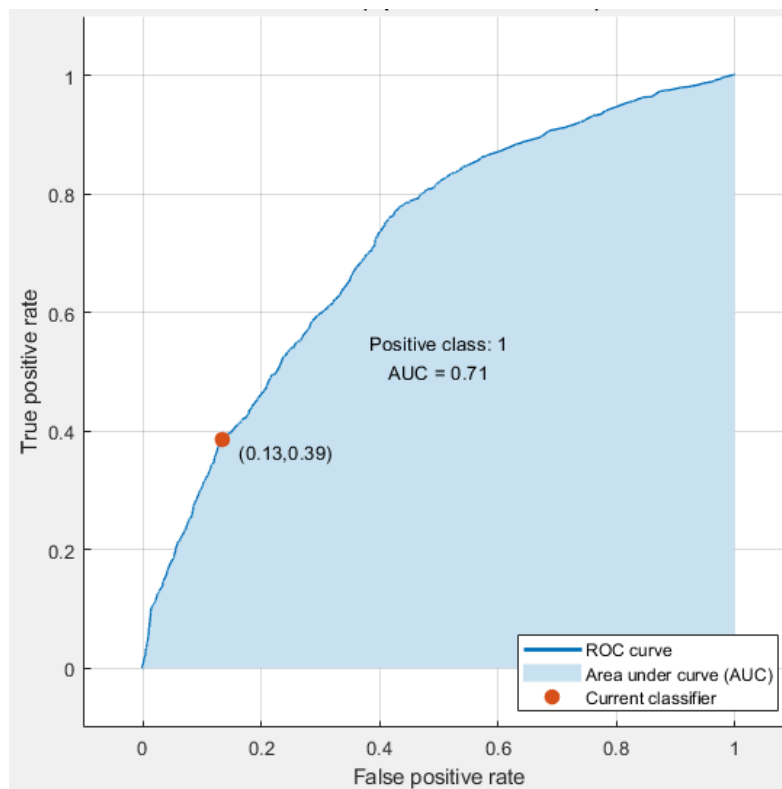
## SVM:



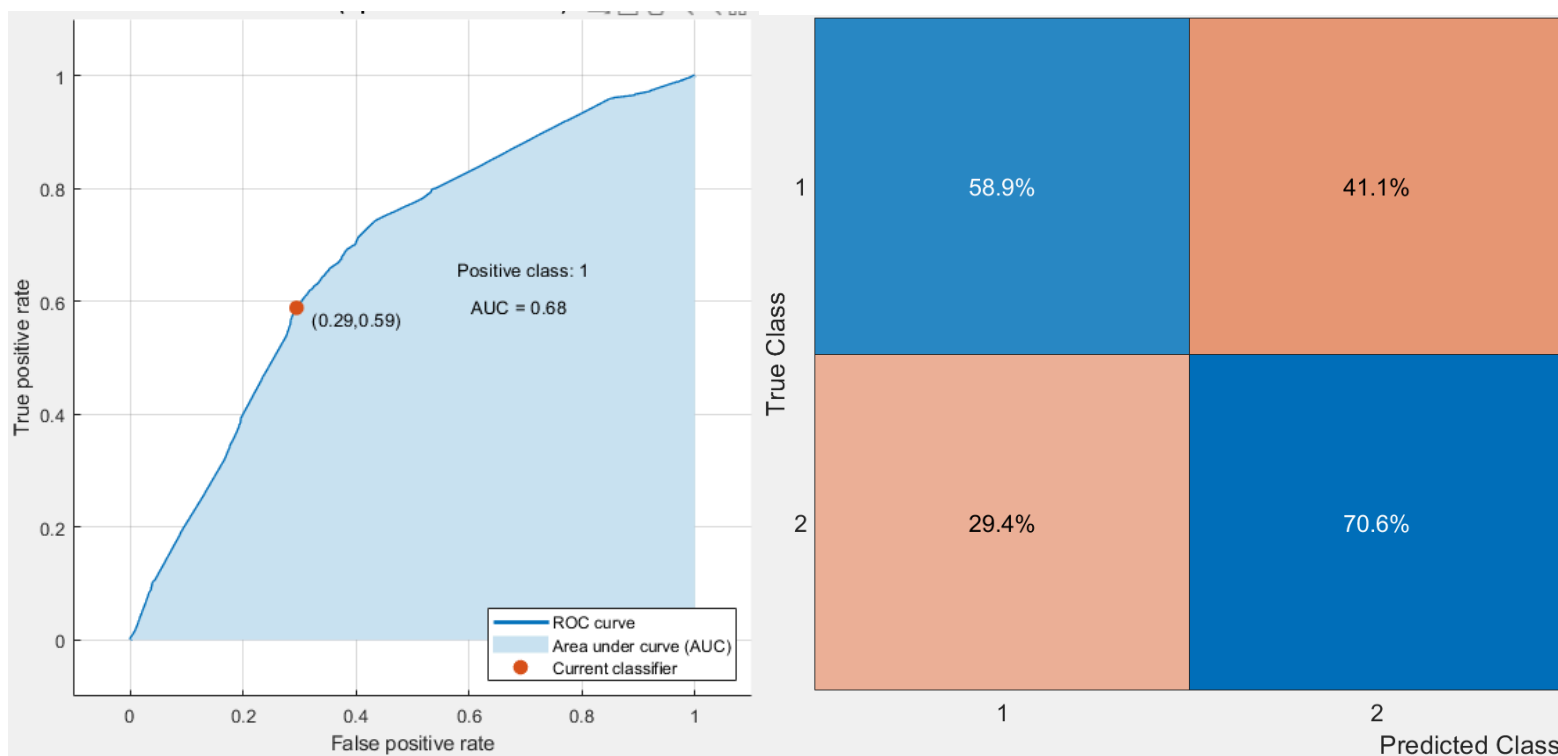
## SVM (7 características):



## Ensemble:



## Ensemble (7 características):



## Comparación de modelos

Para este punto se tendrán en cuenta varias cosas:

- La precisión y exactitud del modelo, mientras mayor mejor, sin llegar a un caso de sobreajuste.
- El número de parámetros en general es de interés obtener modelos que con un bajo número de parámetros sean capaces de cumplir con su objetivo a cabalidad, esto debido a que en un caso real es más difícil y costoso, en términos de dinero y tiempo obtener una cantidad grande de informaciones este caso no se les dará mayor importancia a unos parámetros sobre otros, solo será de interés el número de ellos.
- La complejidad del modelo, se preferirán modelos más simples
- Un último factor que se tendrá en cuenta para preferir un modelo sobre otro es la distribución de falsos negativos hacia cierta clase particular, i.e. en este contexto no sería nada bueno identificar erróneamente a aquellas personas con tendencia repetitiva al intento de suicidio, mientras que identificar erróneamente a aquellos que en realidad no (falso positivo), sería más aceptable.

### Provisional:

Teniendo en cuenta lo anterior lo que se busca es maximizar la predicción correcta de la etiqueta 1 {si a intentos previos de suicidio}, en este sentido la mayoría de los modelos son deficientes. Pero en general los de 7 características se comportan mejor que los completos. Así, uno de los mejores sería el de discriminante lineal con 7 características



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

**Pendiente Escobar:**

- **calcular score del modelo,**
- **poner análisis general de ROC y CM**
- **Hacer eso de la predicción probabilística y determinista**

## Conclusiones

**El chorro:**

- *Hablar de porque nuestros modelos son tan malos (posibles razones: ¿se hicieron las cosas mal? xd. ¿Características comunes para aquellos que son reincidentes por lo que es difícil separarlos de aquello que no? Preferiblemente una mayor cantidad de datos (y también mejor calidad ya que había muchos datos faltantes))*
- *Hay que comentar que podría significar según el contexto esas 7 características más importantes*
- *Como seria excelente de acuerdo con la propuesta inicial poder conseguir datasets con información de personas que han intentado 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>

```
%test = readtable('suicidio_null.csv'); size_test = size(test) %Tes con el  
%data set original sin procesar
```

## 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}