

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

## Avance 2: Entrenamiento, adecuacion y evaluacion de modelos

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

En un primer momento se probaran diferentes modelos , haciendo uso de la herramienta "Classification Learner" de Matlab, debido a su facilidad y rapidez para probar multiples modelos simultaneamente. Se tomaran los pares dataset-modelo que mejores resultados den (preferiblemente por encima de 70% de acierto) para seguirlos trabajando, en cuanto a seleccion de parametros y ajuste de hiperparametros.

Para la seleccion de parametros y ajuste de hiperparametros en donde sea posible se usaran herramientas interactivas o automaticas que provee MATLAB.

Como metodos de validacion y calificacion de los modelos se pretenden usar los dados a continuacion (To Do: añadir breve descripcion de cada uno)

- Score
- Matriz de confusion
- ROC curve

Debido al problema a resolver, al momento de realizar predicciones se generaran dos. Un deterministica y otra probabilistica.

### Descripcion de los data set de entrada.

En el avance 1 se obtuvieron 4 datasets despues del proceso de limpieza, los cuales se describiran a continuacion:

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

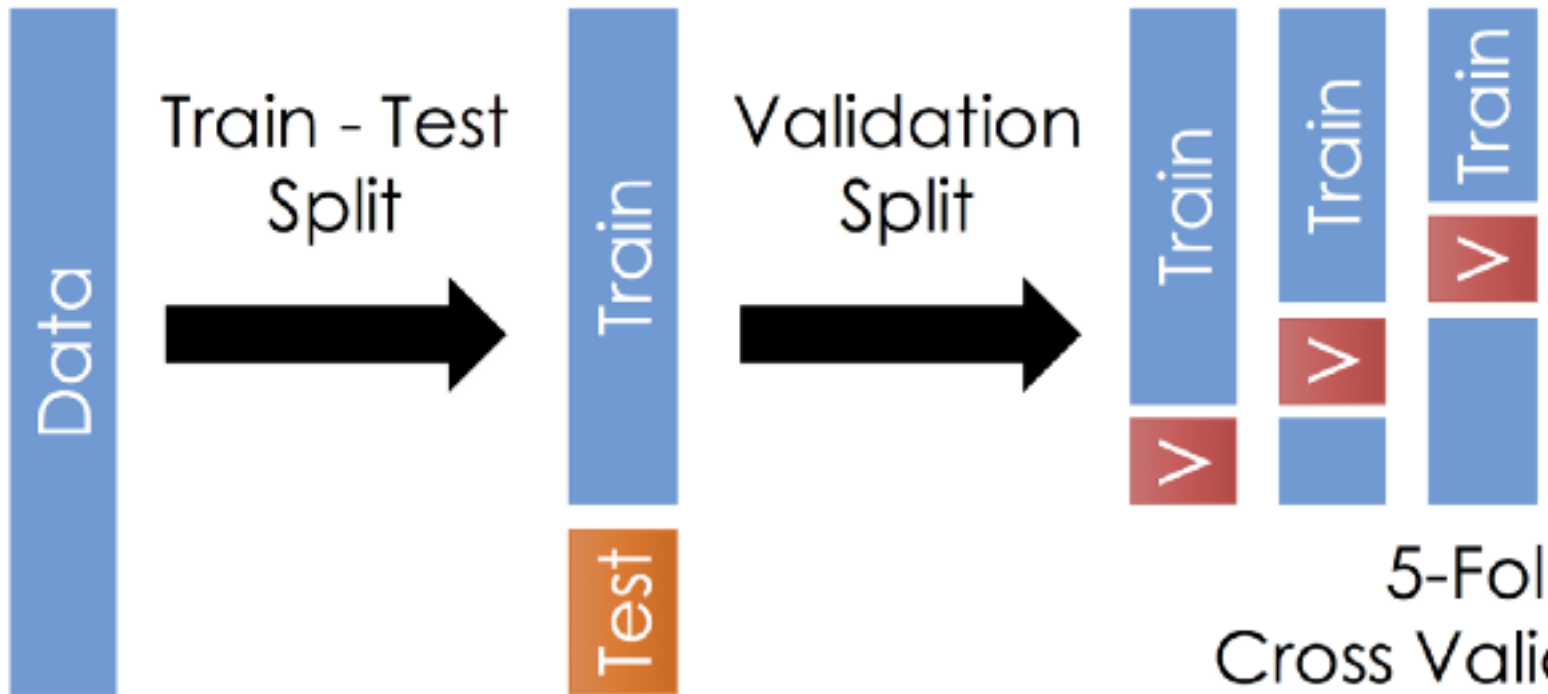
```
%cds = readtable('cds.csv'); size_cds = size(cds)
cds_imputed = readtable('cds_imputed.csv'); size_imputed = size(cds_imputed)
```

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

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

Con estos dataset se procede realizar un enntrenamiento exploratorio de modelos , para continuar con los mas prometedores. Sin embargo, es necesarian definir el concepto de "mas prometedor". En este primer momeno se tendra en cuenta la exactitud del modelos

Es de utilidad tener en cuenta que para lo expuesto a continuacion fue usada validacion cruzaada con "k-folds"( k=5) ,y el valor de precision presentado es correspondiente a esto.



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 fea
<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 fea
<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 fea
<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 fea
<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 fea
<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 fea
<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 fea
<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 fea
<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 fea
<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 fea
<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 fea
<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 fea
<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 fea

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 o no presentan variación con respecto a una de las clases por hallar, por lo que no aportan información

<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: 5
Last change: Fine KNN	33/33 fe
<b>1.16</b> ☆ KNN	Accuracy: 6
Last change: Medium KNN	33/33 fe
<b>1.17</b> ☆ KNN	Accuracy: 6
Last change: Coarse KNN	33/33 fe
<b>1.18</b> ☆ KNN	Accuracy: 6
Last change: Cosine KNN	33/33 fe
<b>1.19</b> ☆ KNN	Accuracy: 6
Last change: Cubic KNN	33/33 fe
<b>1.20</b> ☆ KNN	Accuracy: 5
Last change: Weighted KNN	33/33 fe
<b>1.21</b> ☆ Ensemble	Accuracy: 5
Last change: Boosted Trees	33/33 fe
<b>1.22</b> ☆ Ensemble	Accuracy: 5
Last change: Bagged Trees	33/33 fe
<b>1.23</b> ☆ Ensemble	Accuracy: 6
Last change: Subspace Discriminant	33/33 fe
<b>1.24</b> ☆ Ensemble	Accuracy: 5
Last change: Subspace KNN	33/33 fe
<b>1.25</b> ☆ Ensemble	Accuracy: 5
Last change: RUSBoosted Trees	33/33 fe
<b>2</b> ☆ Linear Discriminant	Accuracy: 6
Last change: 'Covariance structure' ...	33/33 fe
<b>3</b> ☆ Quadratic Discriminant	Accuracy: 6
Last change: 'Covariance structure' ...	33/33 fe
<b>4</b> ☆ Naive Bayes	Accuracy: 5
Last change: Removed 3 features	29/33 fe

Resultados cds\_few\_minus\_alcohol



<b>1.1</b> ☆ Tree	Accuracy: 54.2%	<b>1.15</b> ☆ KNN	Accuracy: 53.3%
Last change: Fine Tree	32/32 features	Last change: Fine KNN	32/32 features
<b>1.2</b> ☆ Tree	Accuracy: 55.6%	<b>1.16</b> ☆ KNN	Accuracy: 55.6%
Last change: Medium Tree	32/32 features	Last change: Medium KNN	32/32 features
<b>1.3</b> ☆ Tree	Accuracy: 55.6%	<b>1.17</b> ☆ KNN	Accuracy: 56.5%
Last change: Coarse Tree	32/32 features	Last change: Coarse KNN	32/32 features
<b>1.4</b> ☆ Linear Discriminant	<b>Failed</b>	<b>1.18</b> ☆ KNN	Accuracy: 55.6%
Last change: Linear Discriminant	32/32 features	Last change: Cosine KNN	32/32 features
<b>1.5</b> ☆ Quadratic Discriminant	<b>Failed</b>	<b>1.19</b> ☆ KNN	Accuracy: 54.3%
Last change: Quadratic Discriminant	32/32 features	Last change: Cubic KNN	32/32 features
<b>1.6</b> ☆ Logistic Regression	Accuracy: 59.0%	<b>1.20</b> ☆ KNN	Accuracy: 55.6%
Last change: Logistic Regression	32/32 features	Last change: Weighted KNN	32/32 features
<b>1.7</b> ☆ Naive Bayes	<b>Failed</b>	<b>1.21</b> ☆ Ensemble	Accuracy: 59.0%
Last change: Gaussian Naive Bayes	32/32 features	Last change: Boosted Trees	32/32 features
<b>1.8</b> ☆ Naive Bayes	Accuracy: 55.0%	<b>1.22</b> ☆ Ensemble	Accuracy: 56.5%
Last change: Kernel Naive Bayes	32/32 features	Last change: Bagged Trees	32/32 features
<b>1.9</b> ☆ SVM	Accuracy: 58.5%	<b>1.23</b> ☆ Ensemble	Accuracy: 58.5%
Last change: Linear SVM	32/32 features	Last change: Subspace Discriminant	32/32 features
<b>1.10</b> ☆ SVM	Accuracy: 55.3%	<b>1.24</b> ☆ Ensemble	Accuracy: 54.3%
Last change: Quadratic SVM	32/32 features	Last change: Subspace KNN	32/32 features
<b>1.11</b> ☆ SVM	Accuracy: 53.3%	<b>1.25</b> ☆ Ensemble	Accuracy: 57.1%
Last change: Cubic SVM	32/32 features	Last change: RUSBoosted Trees	32/32 features
<b>1.12</b> ☆ SVM	Accuracy: 54.3%	<b>2</b> ☆ Linear Discriminant	Accuracy: 60.0%
Last change: Fine Gaussian SVM	32/32 features	Last change: 'Covariance structure' ...	32/32 features
<b>1.13</b> ☆ SVM	Accuracy: 56.5%	<b>3</b> ☆ Quadratic Discriminant	Accuracy: 55.6%
Last change: Medium Gaussian SVM	32/32 features	Last change: 'Covariance structure' ...	32/32 features
<b>1.14</b> ☆ SVM	Accuracy: 58.4%	<b>4</b> ☆ Naive Bayes	Accuracy: 54.3%
Last change: Coarse Gaussian SVM	32/32 features	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 (en términos de precisión) a los obtenidos sin esta transformación. (*¿Uno si debería hacer PCA en datos categoricos?*)

Como se puede notar, ningun par dataset-modelo obtuvo una precision mayore al 70% tal y como se habia definido inicialmente para su aceptacion. Por este motivo se tomara aquel dataset que produjo el modelo con la mayor precision(cds\_imputed) y los mejores modelos obtenidos a partir de este.

Coarse Tree, Linear discriminant, Logistic regresion , SVM coarse emse

## Feature selection

Bucando reducir la dimensionalidad y explorar direferente opciones se pretende realizar un proceso de seleccion de caracteriscas. Esto se hara filtrando aquellas caracteristicas 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 entrenaran 2 modelos, uno con todas las caracteristicas y adicionalmente otro con el conjunto de las 7 mas importantes

```
idx = fscmrnr(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_'
```

## Optimizacion de hiperparametros

Se presentara el proceso para cada uno de los modelos reaizados

Arboles de decision:

**Presentacion de resultados usando el data set completo y el reducido despues de la seelccion de caracteristicas**

caso se tendran en cuenta varias cosas:

- La precision 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 mas informacion. En este caso no le daremos mayor importancia a unos parametros sobre otros, nos concentraremos principalmente en el numero de ellos.
- Un ultimo factor que se tendra en cuenta para preferir un modelo sobre otro es la predisposicion de falsos positivos 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 seria mas aceptable. Esto se tendra en cuenta al observar matrices de confusion despues de realizar la optimizacion de hiperparametros.

When you open the plot, the rows show the true class, and the columns show the predicted class. If you are using holdout or cross-validation, then the confusion matrix is calculated using the predictions on the held-out observations. The diagonal cells show where the true class and predicted class match. If these diagonal cells are blue, the classifier has classified observations of this true class are classified correctly.

## Conclusiones

### Select Features to Include

In Classification Learner, you can specify different features (or predictors) to include in the model. See if you can improve models by removing features with low predictive power. If data collection is expensive or difficult, you might prefer a model that performs satisfactorily without some predictors.

dando prioridad a los sugeridos por el profesor: Regresion logistica, SVM, Arboles de decision, Redes neuronales, LMP, Random Fores

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

Refs:

<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>

```
%Feature selection tests
%X = table2array([cds(:,1:2), cds(:, 4:end)])
%y = table2array(cds(:,3))
%ncaMdl = fscnca(X,y,'FitMethod','exact','Verbose',1,'Solver','lbfgs');
%Model1 = stepwiselm(X,Y,'constant','ResponseVar','ResResistant')
%Model2 = stepwiselm(X,Y,'linear','ResponseVar','ResResistant')
```

```
% figure
% semilogx(ncaMdl.FeatureWeights,'ro')
% xlabel('Feature index')
% ylabel('Feature weight')
% grid on
```

```
% cds_new = movevars(cds,'inten_prev','after','tipo_ss_S');
```

```
% [idx,scores] = fscmrnr(cds_new,"antec_tran");
% bar(scores(idx))
% xlabel('Predictor rank')
% ylabel('Predictor importance score')
```

```
% idx(1:5)
% cds_new.Properties.VariableNames()
% cds_new.Properties.VariableNames(idx(1:5))
```

```
% cds_new.Properties.VariableNames(idx(end-4:end))
```