

Tema 3 : Introducción a Python y a los paradigmas de datos

## IPPPD: Entrega ejercicios Kaggle (Cs570 y AllState)

Inmaculada Perea Fernández

## Índice

Introducción	2
CS570 - Descripción del problema	4
CS570 - Inspección de los datos	5
CS570 - Configuración del entorno	8
CS570 – Selección de modelos y entrenamiento	g
CS570 - Resultados	10
Submissions de Kaggle	
Ranking de Kaggle	11
AllState - Descripción del problema	
AllState - Inspección de los datos	14
AllState - Configuración del entorno	14
AllState – Selección de modelos y entrenamiento	
AllState - Resultados	16
Submission de Kaggle	17
Ranking de Kaggle	18
Referencias	19

### Introducción

En el siguiente documento se realizará un resumen del trabajo realizado para resolver los ejercicios propuestos para la evaluación del tema 3 de la asignatura IPPD.

Se detallan las acciones llevadas a cabo para encontrar la solución óptima a las competiciones de Kaggle propuestas. Algunas de las acciones más destacables son: la inspección de los datos, la configuración y preparación del entorno, las transformaciones de variables, la selección y entrenamiento de modelos y la optimización de parámetros.

Se ha utilizado el framework de machine learning estudiado en clase, configurándolo convenientemente e implementando las clases necesarias para resolver los ejercicios.

Por último se mostrarán los resultados obtenidos y la puntuación privada proporcionada por Kaggle.

El documento se divide en dos secciones distintas, una para la competición CS570 y otra para la competición AllState.

# **CS570**

(https://inclass.kaggle.com/c/newcs570midtermchallenge)

## CS570 - Descripción del problema

A continuación se muestra la descripción y los criterios de evaluación del problema que se plantea en la competición CS570 [1].

Completed • Knowledge • 43 teams

### New CS570 Midterm Challenge

Mon 10 Oct 2016 - Mon 5 Dec 2016 (31 days ago)

Competition Details » Get the Data » Make a submission

### Customer Retention Classifier

This will be CS 570's midterm assignment. For the midterm, we will be using everything you've learned so far to build a customer "churn" model.

Imagine you work for GreenPlate, the west coast's fastest growing new subscription box service. GreenPlate's customers get a monthly box, full of healthy and sometimes delicious snacks. GreenPlate wants to know which customers are most likely to cancel their subscription service so that they can use incentives to retain those customers.

You need to use the training dataset given to create a model that predicts if the customers in the testing dataset will cancel their subscription to GreenPlate.

The dataset contains multiple features, and you can use any of those features you want. The names of those features have been hidden from you.

The dependent variable is y. If y=0 that means the customer didn't cancel. If y=1 that means the customer DID cancel. You should predict the probability of cancellation because your model will be evaluated using AUC.

You have two separate deliverables, listed below:

- 1. Please submit your answer via Kaggle. You'll need to be on the leaderboard to receive credit for this assignment.
- 2. Please submit all the code you used for cleaning, prepping, and modeling your data via blackboard (NOT github).

### **Evaluation**

#### Evaluation

The evaluation metric that this competition will be using is Area Under the Curve (AUC).

#### **Submission Format**

For every entry in the test set, submission files should contain two columns: Id and Predicted Probability of cancellation.

The file should contain a header and have the following format:

id,y

1,0

2,.1

3,.9

4,.4

etc.

#### Grading

Your midterm will be scored by:

50% Your model implementation

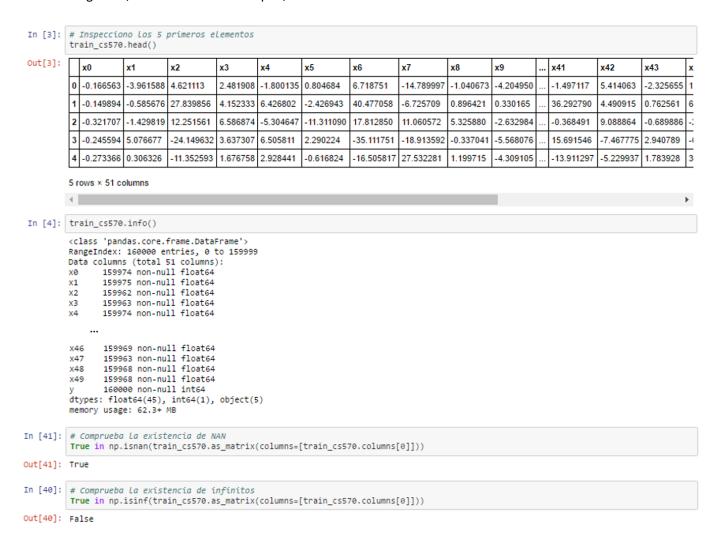
25% Your data cleaning/prep code

25% Your leaderboard score vs your classmates.

## CS570 - Inspección de los datos

El primer lugar inspeccionamos los datos con el propósito de localizar datos inexistentes, incompletos, erróneos, etc. Esto nos dará idea del pre-procesado necesario antes de comenzar con la selección de modelos y el entrenamiento de los mismos.

Para este cometido he creado un fichero notebook ipython [3] que utiliza las librerías pandas y numpy para explorar los ficheros csv de training y determinar el tipo de datos de las variables a analizar, la existencia de valores NaN, infinitos, variables categóricas, cadenas de texto a limpiar, etc.



Observamos que los datos contienen NaN, con lo cual en el pre-procesado a realizar antes de la selección de modelos será necesario limpiarlos. Para ello utilizaremos la clase *ImputerTransform* existente en *Feature\_transform*.py

```
class ImputerTransform(BaseEstimator, TransformerMixin):
64
65
            def init (self, missing values='NaN', strategy=-999999):
66
67
                self.missing_values = missing_values
68
                self.strategy = strategy
69
70
71
            def fit(self, X, y=None):
              return self
73
74
            def transform(self, X):
75
                if self.strategy in ['mean', 'median', 'most frequent']:
76
77
78
                    imp = Imputer(self.missing_values, self.strategy)
79
                    X = imp.fit_transform(X)
80
81
        # los algoritmos siempre trabajan con numeros, hay que sustituir los caracteres
82
                else:
83
                    for c in X.columns:
84
                        if X[c].dtype == np.object:
85
86
                            strategy = str(self.strategy)
87
                        else:
88
                            strategy = int(self.strategy)
89
                        if self.missing_values == "NaN":
90
91
                           X.loc[pd.isnull(X[c].values), c] = strategy
92
93
                            X.loc[X[c].isin([self.missing_values]), c] = strategy
94
                return X
95
```

A continuación inspeccionaremos las 5 variables tipo "object" para determinar qué procesado será necesario realizar. Encontramos que existen tres variables (x24, x29 y x30) categóricas que habrá que codificar correctamente para trabajar con los modelos.

```
In [18]: train_cs570["x24"].head()
Out[18]: 0
              euorpe
         3
                asia
                asia
         Name: x24, dtype: object
In [15]: train_cs570["x29"].head()
Out[15]: 0
              July
               Aug
              July
              July
              July
         Name: x29, dtype: object
In [17]: train_cs570["x30"].head()
Out[17]: 0
                tuesday
              wednesday
              wednesday
             wednesday
                tuesday
         Name: x30, dtype: object
```

Para el tratamiento de estas variables he implementado una nueva clase *Encoder* en *Feature\_transform.py*, esta clase está basada en un encoder utilizado en la versión 1 del framework, pero he tenido que adaptarla ya que en este caso el tipo de dato a reemplazar es de tipo 'object' y no 'str'. Esta clase utiliza el método LabelEncoder de sklearn para codificar las variables categóricas.

```
19
      class Encoder(BaseEstimator, TransformerMixin):
20
21
            def __init__(self, columns=None):
                self.columns = columns
22
23
24
            def fit(self, X, y):
25
                return self
26
            def transform(self, X):
27
                columns = self.columns
28
29
                if not columns:
30
                    columns = X.columns
31
32
                for c in columns:
                    if (X[c].dtype == 'object'):
33
                        X[c] = LabelEncoder().fit_transform(X[c])
34
35
                return X
```

Por último nos queda por inspeccionar las variables x32 y x37, en las que podemos observar que los datos incluyen el carácter % y \$ junto con datos de tipo numérico.

```
In [13]: train_cs570["x32"].head()
Out[13]: 0
                0.0%
              -0.02%
              -0.01%
               0.01%
               0.01%
         Name: x32, dtype: object
In [14]: train_cs570["x37"].head()
Out[14]: 0
               $1313.96
               $1962.78
                $430.47
              $-2366.29
               $-620.66
         Name: x37, dtype: object
```

En este caso utilizaremos la clase RemoveStringTransform existente en el framework, que reemplaza los caracteres que se le pasan como argumento por una cadena vacía.

```
37
38
       class RemoveStringsTransform(BaseEstimator, TransformerMixin):
39
            def __init__(self, strings, columns=None):
40
                self.strings = strings
41
                self.columns = columns
42
43
44
45
            def fit(self, X, y):
46
                return self
47
48
49
            def transform(self, X):
50
                # Si no se especifica nada se aplica a todas las columnas
51
                columns = self.columns
52
                if not columns:
53
                    columns = X.columns
54
55
                for c in columns:
56
                    for string in self.strings:
                        # Compruebo si tiene caracter, y si lo tiene lo sustituyo para cadena vacia
57
58
                        if X[c].dtype == np.object and X[c].str.contains(string).any():
59
                            X[c] = X[c].str.replace(string, '')
60
61
                return X
62
```

## CS570 - Configuración del entorno

En este apartado detallaré los principales parámetros configurados en el fichero CS570/config.py [5] utilizado para el preprocesado, selección de modelos y generación de predicciones con el framework.

He realizado varias pruebas de pre-procesado para este ejercicio, y el que se muestra a continuación es el que me ha ofrecido mejores resultados.

```
49
        PREPROCESSING = [
50
51
             ('scaler int', [
52
                 ('dropNaN', ImputerTransform()),
53
                 ('clean', RemoveStringsTransform(strings=['%', '$'])),
                 ('encode', Encoder()),
                # ('log', LogTransform()),
56
                 ('inter', InteractionTransform(
57
                     interactions=['add'],
                     columns=['x%d' % i for i in range(1, 10)])),
58
                 ('std', StandardScaler()),
59
60
                 ('fs', TreeBased('extra_trees_regressor', 20, 190)),
61
62
            ]),
63
```

He utilizado métrica log\_loss para calcular el score y comparar los distintos modelos generados

```
SCORER = make_scorer(log_loss,
needs_proba=True,
greater_is_better=False)
```

También he cambiado el parámetro TARGET porque en este ejercicio la variable objetivo es la "y".

```
TARGET = "y"

TARGET_TRANSFORM = LabelTransform
```

Para crear las predicciones tal y como las requiere kaggle tuve que utilizar la clase Dummyld para que generase la columna con el ld en el csv de submission.

```
PREDICTIONS_DIR = 'cs570/submissions'

PREDICTIONS_PREPROCESS = [DummyId()]

PREDICTION_COLUMNS = ['id']
```

A continuación muestro una captura de todos los modelos con los que he trabajado y que he configurado convenientemente en el fichero config. El modelo KNN no lo he generado para este ejercicio puesto que tarda mucho en procesar (más de 3 horas) y he abortado su generación para poder realizar más pruebas con otros modelos que no necesitan tanto tiempo de procesado.

```
MODELS = {
68
            ('lr', LogisticRegression(fit_intercept=True, solver='newton-cg',
                                    multi_class='multinomial')),
69
            ('rf', RandomForestClassifier(random state=RANDOM STATE)),
70
            ('xgbClasifier', XGBClassifier(base_score=0.5, colsample_bylevel=0.6, colsample_bytree=0.6,
72
              gamma=0, learning_rate=0.1, max_delta_step=0,
73
              max_depth=10, min_child_weight=1,
              n_estimators=1000, nthread=-1, objective='binary:logistic',
74
75
              reg alpha=0.6, reg lambda=1, scale pos weight=1, silent=True,
              subsample=1)),
78
           ('xqb', XGBRegressor()),
           ('lasso', Lasso(fit intercept=True)), # inma
           ('ridge', Ridge(fit intercept=True)),
           ('elastic net', ElasticNet(fit_intercept=True)),
81
           ('bayes_ridge', BayesianRidge(fit_intercept=True)),
82
83
            ('SGD', SGDRegressor()),
84
            ('gb', GradientBoostingRegressor()),
            #('KNN', KNeighborsRegressor(n_neighbors=10)),
```

## CS570 – Selección de modelos y entrenamiento

Una vez obtuve resultados con distintos modelos de clasificación y regresión logística probé a realizar transformaciones de variables y optimización de parámetros para mejorar los resultados:

```
Project: cs570
Project: cs570
                                         Date: Tue Jan 3 22:23:26 2017
Date: Sun Nov 27 17:29:45 2016
Random state: 2016
                                         Random state: 2016
Target transform: LabelTransform
                                         Target transform: LabelTransform
Parameter optimization: True
                                         Parameter optimization: True
Train data shape: 160000, 52
                                         Train data shape: 160000, 52
Preprocess
             | Model
                       | Score | STD
                                         | Preprocess
                                                                                | STD
                                                           | Model | Score
                      1 0.33567 | 0.00200
                                         | ----- | ----- | ------ | ------ | -----
                                          | scaler int
                                                            Project: cs570
Project: cs570
Date: Sun Nov 27 18:13:32 2016
                                          Date: Tue Jan 3 22:27:43 2017
Random state: 2016
                                         Random state: 2016
Target transform: LabelTransform
                                          Target transform: LabelTransform
Parameter optimization: True
                                          Parameter optimization: True
Train data shape: 160000, 52
                                          Train data shape: 160000, 52
| Preprocess | Model | Score | STD
                                         | Preprocess | Model
                                                                      | Score | STD
|-----|----|
| scaler_int
                                                           | elastic_net | 0.67351 | 0.00000
```

- 1) Probé a realizar el logaritmo usando la clase LogTransform, pero obtuve errores, creo que se debe a que hay números demasiado grandes, así que abandoné ese camino.
- 2) Probé a realizar operaciones entre columnas (suma y multiplicación) con la clase InteractionTransform, pero no obtuve mejora del score. Probé primero con las 10 primeras variables (x1,..., x10) y luego con 20 variables (x1,..., x20) sin resultados positivos.
- 3) Los que mejores resultados los obtuve con los modelos rf y gb, asi que me centré en ellos e hice varias pruebas variando algunos de los META\_PARAMETERS: n\_estimators, max\_features y max\_depth pero tampoco he conseguido mejoras significativas.

- 4) También he probado a variar los parámetros learning\_rate, max\_depth y n\_estimators del modelo xgboostClassifier, pero no he conseguido mejorar los resultados obtenidos con RandomForestClassifier
- 5) No he podido probar KNN porque el tiempo que necesita para finalizar es muy alto y me he visto obligada abortar la ejecución en 2 ocasiones.

## CS570 - Resultados

El modelo con el que mejores resultados he obtenido es con RandomForestClassifier

Project: cs570

Date: Sun Nov 27 17:29:45 2016

Random state: 2016

Target transform: LabelTransform Parameter optimization: True Train data shape: 160000, 52

### Con los siguientes parámetros:

- min\_samples\_split=1,
- n\_estimators=20,
- criterion=gini
- min samples leaf=3
- bootstrap=True
- min\_impurity\_split=1e-07
- max\_features=auto
- n\_jobs=1
- random\_state=2016

### Submissions de Kaggle

A continuación los 3 modelos con los que he obtenido los mejores resultados en kaggle.

## New CS570 Midterm Challenge

Mon 10 Oct 2016 - Mon 5 Dec 2016 (31 days ago)

You are submitting as part of team inm	a perea. Mal	ke a submi:	ssion »	
The competition deadline has already p While this competition was active, you c information is provided for historical pu	ould select up to	_	_	
Submission	Files	Public Score	Private Score	Selected?
Post-Deadline: Sun, 18 Dec 2016 12:38:49 Edit description	rf_0.33567.cs v	0.95448	0.95705	
Post-Deadline: Wed, 04 Jan 2017 20:35:10 Edit description	gb_0.44365.c	0.90490	0.90843	
Post-Deadline: Thu, 05 Jan 2017 18:43:50 Edit description	gb_0.44365.c	0.90490	0.90843	
Post-Deadline: Wed. 04 Jan 2017 20:38:40	bayes ridge	0.50173	0.50005	

## Ranking de Kaggle

23	_	AlvaroSanchez	0.96994	3	Mon, 05 Dec 2016 13:30:02
24	_	gamgee	0.96327	2	Sun, 30 Oct 2016 19:28:45
-		inma perea	0.95705		Fri, 06 Jan 2017 18:30:56 Post-Deadline
		dline Entry uld have submitted this entry during the competitio	on, you would l	nave be	een around here on the leaderboard.
			on, you would I	nave be	een around here on the leaderboard.  Mon, 24 Oct 2016 22:13:01

# **Allstate**

(https://www.kaggle.com/c/allstateclaimsseverity)

## AllState - Descripción del problema

A continuación se muestra la descripción y el modo de evaluación del problema que se plantea en la competición AllState [2].

Completed • Jobs • 3,055 teams

### Allstate Claims Severity

Mon 10 Oct 2016 - Mon 12 Dec 2016 (24 days ago)

Competition Details » Get the Data » Make a submission

### How severe is an insurance claim?

When you've been devastated by a serious car accident, your focus is on the things that matter the most: family, friends, and other loved ones. Pushing paper with your insurance agent is the last place you want your time or mental energy spent. This is why Allstate, a personal insurer in the United States, is continually seeking fresh ideas to improve their claims service for the over 16 million households they protect.











Allstate is currently developing automated methods of predicting the cost, and hence severity, of claims. In this recruitment challenge, Kagglers are invited to show off their creativity and flex their technical chops by creating an algorithm which accurately predicts claims severity. Aspiring competitors will demonstrate insight into better ways to predict claims severity for the chance to be part of Allstate's efforts to ensure a worry-free customer experience.

New to Kaggle? This competition is a recruiting competition, your chance to get a foot in the door with the hiring team at Allstate.

### **Evaluation**

Submissions are evaluated on the mean absolute error (MAE) between the predicted loss and the actual loss.

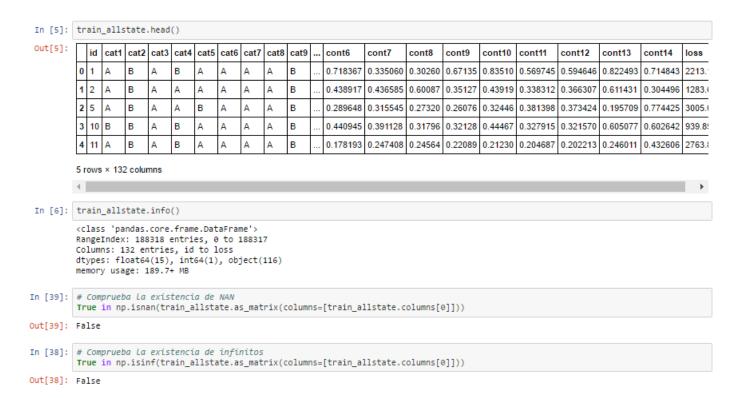
### Submission File

For every id in the test set, you should predict the loss value. The file should contain a header and have the following format:

id,loss 4,0 6,1 9,99.3 etc.

## AllState - Inspección de los datos

Del mismo modo que en el ejercicio anterior inspeccionamos los datos para determinar el pre-procesado antes de empezar con la selección de modelos.



Al contrario que en el ejercicio anterior en este ejercicio no es necesario realizar limpieza de datos ya que no existen datos NAN, o caracteres incorrectos. Será necesario utilizar el encoder implementado en el ejercicio anterior para codificar las variables categóricas.

## AllState - Configuración del entorno

En este caso la configuración utilizada es similar a la del ejercicio anterior, exceptuando algunos parámetros obvios como los de los directorios para almacenar las predicciones y los modelos generados.

La variable objetivo, que en este caso es "loss", y por tanto hay que modificar convenientemente el parámetro TARGET.

Otra diferencia está en el parámetro TARGET\_TRANSFORM en este caso indicamos que se realice la transformación logarítmica de la columna respuesta.

También hay diferencias en la métrica usada para calcular el error, en este caso se usa la métrica MAE (parámetro SCORER)

```
TARGET = "loss"

TARGET_TRANSFORM = LogTransform

SCORER = make_scorer(median_absolute_error, greater_is_better=False)
```

Con respecto a los modelos, la configuración utilizada es la siguiente, he ido descomentando en cada prueba el modelo que me interesaba y cambiando los META\_PARAMETERS.

```
MODELS = {
   #('lr', LinearRegression(fit intercept=True)),
    #('xgb', XGBRegressor(n_estimators=100, colsample_bytree=0.6, colsample_bylevel=0.6,
                               #subsample=0.5, learning_rate=0.1,
                               #max_depth=2, reg_alpha=0.6, min_child_veight=1)),
    ('xgb', XGBRegressor(n_estimators=1000, colsample_bytree=0.6, colsample_bylevel=0.6,
              subsample=0.5, learning_rate=0.1,
                               max_depth=4, reg_alpha=0.6, min_child_weight=1)),
   # ('lasso', Lasso(fit intercept=True)),
    #('ridge', Ridge(fit intercept=True)),
    #('elastic_net', ElasticNet(fit_intercept=True)),
    #('bayes_ridge', BayesianRidge(fit_intercept=True)),
    #('SGD', SGDRegressor()),
    #('random_forest', RandomForestRegressor()),
    #('gb', GradientBoostingRegressor()),
    #('KNN', KNeighborsRegressor(n_neighbors=10)),
}
```

## AllState – Selección de modelos y entrenamiento

Los primeros resultados que obtuve fueron los siguientes:

Project: allstate

Date: Sun Dec 18 17:35:38 2016

Random state: 2016

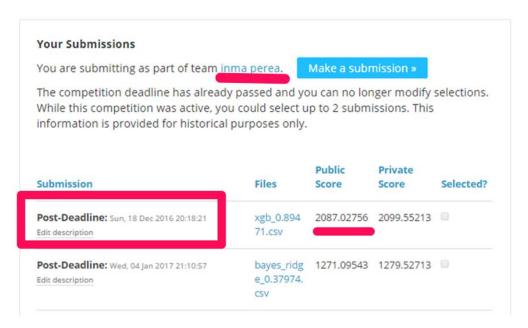
Target transform: LogTransform Parameter optimization: True Train data shape: 188318, 133

1	Preprocess	Ī	Model	Ī	Score	1	STD	Ī	OP	1	Parameters
-1		1		I		1		1		1	
-1	none	Ī	elastic_net	. 1	0.53795		0.00060		F		alpha=1.0, max_it
-1	none	1	lasso	I	0.56923	1	0.00031	1	F	1	alpha=1.0, max_ite
-1	none	I	xgb	I	0.89471	1	0.00081	1	T	1	n_estimators=20, s
-1	none	Ī	ridge	I	0.37982	1	0.00162	1	F	1	fit_intercept=True
-1	none	1	lr	I	0.38345	1	0.00124	1	T	1	n_jobs=1, copy_X=1
-1	none	I	random_fore	st	0.49204		0.00286		- 1	Т	n_estimators=1,
-1	none	ı	KNN	I	0.48267	1	0.00105	1	T	1	algorithm=ball_tre
-1	none	1	bayes_ridge		0.37974		0.00128		F		alpha_1=1e-06, a:
-1	none	I	SGD	I	10018498884	12	73.96094   3	32:	1593	332	294511.86719   F
-1	none	ī	gb	ī	0.56592	Ī	0.00012	ī	T	I	loss=ls, criterio:

Total execution time: 179.74367972215018 minutes

### Allstate Claims Severity

Mon 10 Oct 2016 - Mon 12 Dec 2016 (24 days ago)



Para mejorar los resultados anteriores realicé las siguientes pruebas:

- 1) En primer lugar hice pruebas para ver si realizando operaciones entre variables, sumas, multiplicaciones mejoraba los resultados pero apenas si obtuve mejora, y algunos modelos empeoraban.
- 2) Luego probé a realizar pre-procesado, añadiendo lo siguiente por separado en el fichero config.py y tampoco obtuve mejora

```
('std', StandardScaler()),
('fs', TreeBased('extra trees regressor', 20, 190)),
```

3) Luego probé a aumentar el parámetro *n\_estimators* del modelo *XGBRegressor*, y la mejora fue bastante significativa. Pasando de 0.89565 a 0.35377. De modo que continúe experimentando con los parámetros de este modelo. Cambié *n\_estimators* y *learning\_rate* de forma inversa, es decir, si aumento uno disminuyo el otro. Y una vez no obtuve mejora con los anteriores comencé a aumentar el parámetro *max\_depth*. Con esto he conseguido disminuir sustancialmente el score de 0.894 a 0.341.

## AllState - Resultados

A continuación los mejores resultados de score y STD obtenidos con el framework tras el entrenamiento y la optimización de parámetros que he realizado.

### Con los siguientes parámetros:

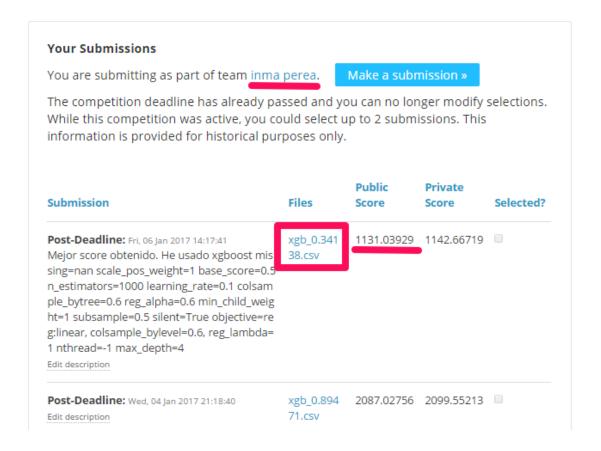
- missing=nan
- scale\_pos\_weight=1
- base\_score=0.5
- n estimators=1000
- learning\_rate=0.1
- colsample\_bytree=0.6
- reg\_alpha=0.6
- min\_child\_weight=1
- subsample=0.5
- silent=True
- objective=reg:linear
- colsample\_bylevel=0.6
- reg\_lambda=1
- nthread=-1
- max\_depth=4

## Submission de Kaggle

A continuación una captura de la mejor puntuación obtenida en kaggle para este ejercicio.

## **Allstate Claims Severity**

Mon 10 Oct 2016 - Mon 12 Dec 2016 (24 days ago)



### Ranking de Kaggle

A continuación una captura con el ranking privado mostrado por kaggle para las predicciones que he obtenido.

1648	<b>↓14</b>	Mougatine	1142.32241	22	Mon, 14 Nov 2016 17:28:02 (-31.4d)
1649	<b>↑18</b>	AchyutJoshi	1142.41956	4	Fri, 11 Nov 2016 11:32:37 (-5.4h)
1650	<b>↓12</b>	murphy	1142.44269	3	Wed, 30 Nov 2016 15:29:38 (-2d)
1651	↓6	PrashantMishra	1142.53022	7	Fri, 28 Oct 2016 19:07:25 (-2.3d)
			4440.66740		Fri, 06 Jan 2017 14:17:41
-		inma perea	1142.66719	-	Post-Deadline
		line Entry Id have submitted this entry during the competition, you		ound h	Post-Deadline
	ı woul	line Entry		- ound h	Post-Deadline
If you	ı woul	line Entry Id have submitted this entry during the competition, you	would have been ar		Post-Deadline nere on the leaderboard.
<b>If you</b>	ı woul	line Entry Id have submitted this entry during the competition, you MikeFrantz	would have been ar 1142.75908	1	Mon, 28 Nov 2016 19:41:44  Tue, 29 Nov 2016 03:12:08

## Referencias

- [1] https://inclass.kaggle.com/c/newcs570midtermchallenge
- [2] https://www.kaggle.com/c/allstateclaimsseverity
- [3] Fichero de inspección de datos: data\_inspection\_kaggle.ipynb
- [4] Fichero de configuración cs570/config.py: cs570\_config.py
- [5] Fichero de configuración allstate/config.py: allstate\_config.py