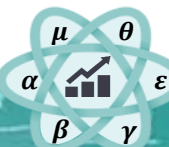
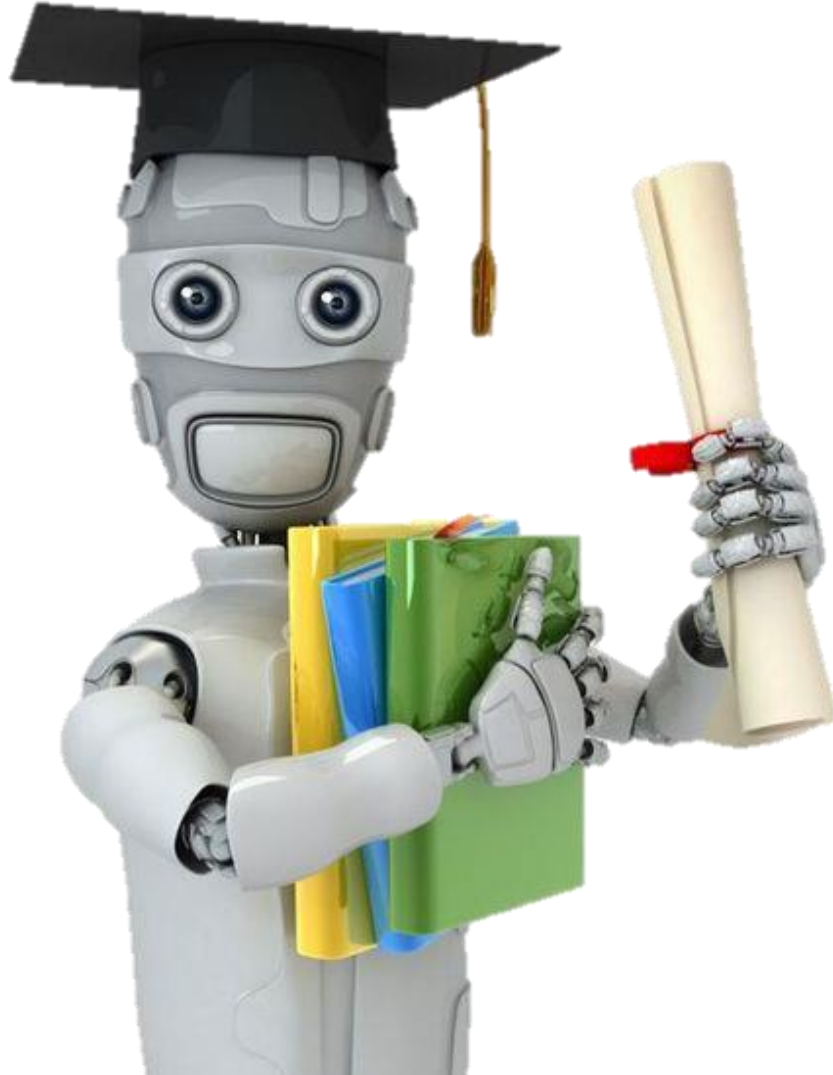


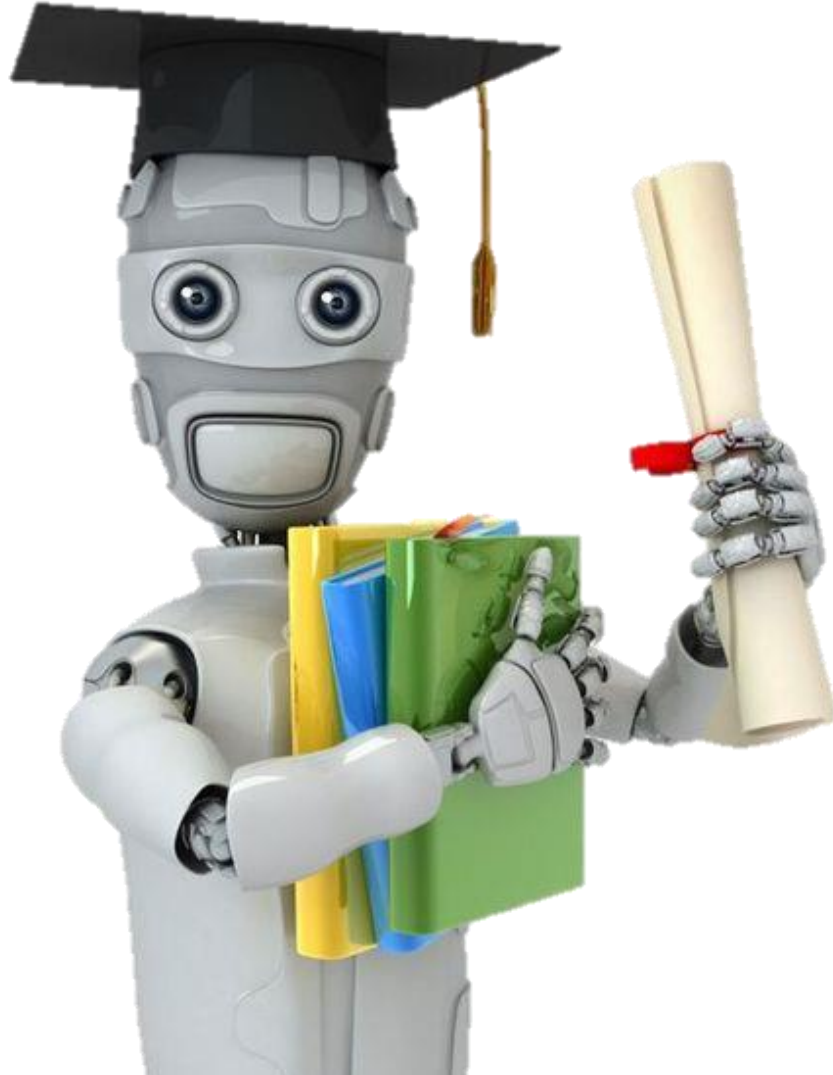
# Charla 1:

## CARACTERÍSTICAS DE SERIES DE TIEMPO PARA UN MODELO DE CLASIFICACIÓN



# Charla 1:

## CARACTERÍSTICAS DE SERIES DE TIEMPO PARA UN MODELO DE CLASIFICACIÓN



Alex Castaño Ballesteros



[github.com/acastanob](https://github.com/acastanob)



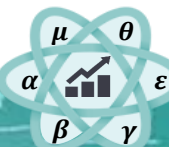
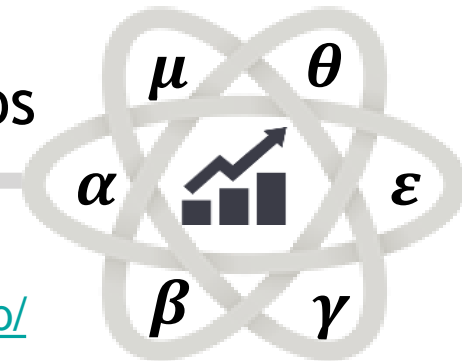
[linkedin.com/in/acastanob/](https://linkedin.com/in/acastanob/)



[acastanob@gmail.com](mailto:acastanob@gmail.com)



[312 840 88 87](tel:3128408887)



# P r e d i c c i ó n



¿Se Fugará el cliente?

¿Está Lavando Activos?

¿ Es Fraude?

¿Tengo el virus xyz?

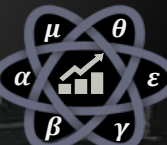
¿Es un cliente potencial?

¿Comprará otro  
producto?

¿Me pagará la deuda?

¿Eso será Pitufeo?

¿Colombia ganará el  
Mundial?





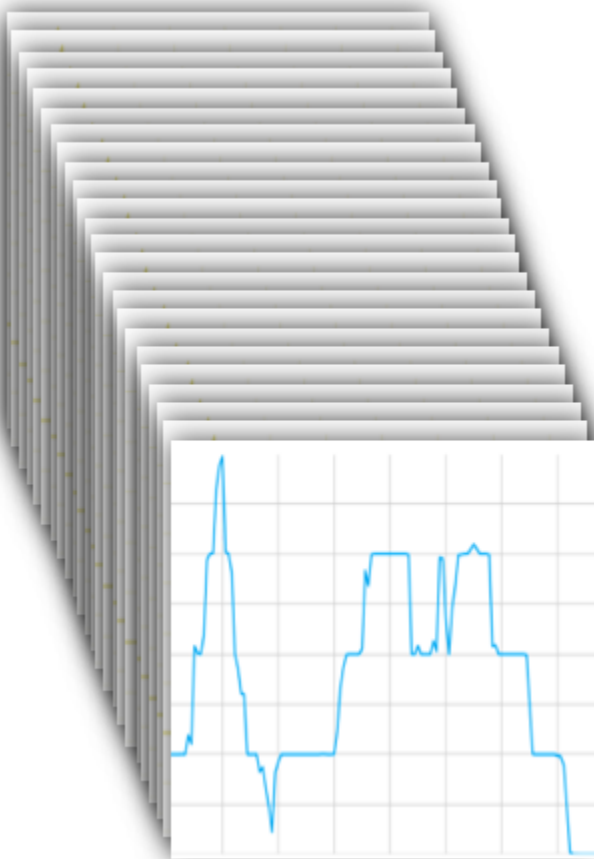
## DATOS

- Sociodemográficos
- Transaccionales

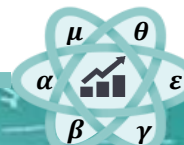
¿Qué variables o características tengo?



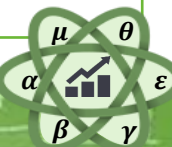
# Proceso



12	ST	SEX	AGE	STU	PRO	NBV	NBI	NBPS	NBPR	MB	ENF	L.T	w.t	we.t	IN.T	K.A	D	C	D	TA	PO	PMAX	PMIN	
03:16:00	1	1	21	1	1	10	10	1	1	2	0	1	3	4	7	14	3	4	1	180	48.1	94.5	45.4	
22:44:00	1	1	34	1	1	10	30	30	1	2	0	4	6	6	6	14	6	1	4	165	94	75	54	
20:16:00	1	1	22	1	1	30	20	50	5	1	0	3	9	8	8	18	6	4	5	170	98	92	61.5	
19:02:00	1	1	27	0	0	4	5	3	3	5	0	3	3	5	9	18	9	1	7	175	154	165	104	
22:41:00	1	0	24	1	0	2	18	10	2	3	0	1	6	6	5	17	3	1	2	165	90	90	52	
23:54:00	1	1	24	1	0	1	20	3	3	2	0	1	4	5	6	18	6	1	3	165	90.4	94.4	42.2	
23:47:00	1	1	16	1	0	35	100	10	2	2	0	2	3	2	4	4	1	3	187	104	138	108		
14:10:00	1	1	37	0	1	15	100	10	4	3	1	4	7	10	5	17	8	2	5	170	100	118	63	
22:54:00	1	1	17	0	0	2	5	1	3	0	0	2	1	1	3	15	5	1	3	152	58	65	58	
23:46:00	1	1	21	1	1	2	3	2	3	0	4	8	4	4	14	6	1	4	182	48.5	90	46		
23:49:00	1	0	28	1	0	3	30	1	6	0	3	6	7	9	21	2	4	2	165	65	115	55		
06:41:00	1	1	17	1	0	3	3	3	3	4	0	5	2	4	7	9	1	3	162	91	53	44		
01:16:00	1	1	17	1	0	15	2	10	3	0	2	3	3	3	4	4	2	3	157	92	90	46		
01:28:00	1	1	17	1	0	5	7	2	2	4	0	4	6	9	8	16	5	2	3	185	63.5	66.6	56.8	
01:53:00	1	1	19	1	0	50	99	20	1	1	0	3	6	2	9	13	4	1	1	180	40	92	38	
02:15:00	1	1	18	1	0	1	1	1	1	0	2	0	2	4	4	7	17	4	1	2	187	90	64	52
03:34:00	1	1	24	0	1	5	25	7	3	3	0	3	6	4	7	21	4	2	3	170	85	87	54	
03:54:00	1	1	21	1	0	6	10	2	1	1	0	2	4	5	8	18	4	2	2	180	90	96	51	
08:19:00	1	1	16	1	1	45	25	2	2	3	0	5	2	3	5	5	1	3	158	51	57	46.5		
12:27:00	1	1	18	1	0	15	5	8	2	3	0	1	1	1	2	5	1	1	140	25	30	23		
12:46:00	1	1	20	1	0	4	12	6	3	3	0	5	9	9	9	19	9	1	3	165	61	86	58	
14:41:00	1	1	19	1	0	3	70	6	3	6	0	3	7	11	9	6	1	8	161	77	101	44		
18:39:00	1	1	20	1	0	1	3	2	2	2	0	2	4	4	4	5	16	4	1	137	47.6	52.2	43.1	
19:55:00	1	1	17	1	0	9	3	1	1	4	0	1	3	3	7	7	1	4	153	51	63	51		
22:13:00	1	1	32	0	1	8	5	4	1	1	0	2	3	5	5	16	4	1	2	182	90	94	48	
22:23:00	1	1	18	1	0	4	7	2	1	4	0	1	4	5	4	6	2	6	163	63	65	52		
22:41:00	1	0	22	1	0	4	8	3	1	9	0	3	5	4	9	15	6	1	3	172	85	75	60	
22:55:00	1	1	18	0	1	15	100	5	3	3	0	3	2	1	8	14	8	1	7	177	71.3	82.7	43	
23:41:00	1	1	16	1	0	350	700	10	3	2	0	1	5	6	9	6	1	3	159	40	52	44		
04:09:00	1	1	21	1	0	5	7	3	2	4	0	1	8	7	7	18	5	3	5	155	90	98	30	
05:02:00	1	1	17	1	1	3	7	5	3	2	0	4	2	8	8	5	2	2	152	48.5	94	46		
12:50:00	1	1	21	0	0	6	2	3	2	2	0	4	2	2	5	18	8	3	9	150	63	66	63	
12:55:00	1	1	19	1	1	8	20	5	5	3	0	4	1	1	6	16	4	1	2	175	85	81	60	
16:15:00	1	1	23	1	1	yes	yes	yes	yes	1	1	2	4	1	7	16	6	2	4	185	57.2	114.3	49.3	
16:34:00	1	1	26	1	0	2	2	5	3	1	0	1	6	1	8	16	2	2	1	164	44.4	50.8	36	
16:40:00	1	1	26	1	0	30	3	2	2	5	0	3	2	3	8	6	2	3	175	72	97	57		
16:37:00	1	1	17	1	0	20	10	10	5	1	0	3	2	2	8	14	2	1	1	167	44	94	37	
21:24:00	1	0	24	0	1	50	15	20	10	4	0	2	3	1	6	11	4	1	9	40	90	30	30	
22:35:00	1	1	17	1	1	300	800	5	3	4	0	3	5	5	5	16	6	1	2	172	58	64	47	
07:01:00	1	1	17	1	0	40	250	6	6	4	0	1	4	4	5	2	1	2	180	50.7	58.96	53.1		
09:35:00	1	1	27	0	0	yes	yes	yes	yes	1	0	3	8	7	5	5	6	1	3	152	88	74	30	
10:52:00	1	1	42	0	1	3	20	5	3	1	0	2	2	3	6	5	2	3	180	87	43	30		
13:21:00	1	1	22	1	1	20	25	5	3	4	0	1	1	1	3	18	2	1	2	159	55	80	50	
13:37:00	1	1	27	0	1	20	10	5	2	2	0	2	5	4	6	5	2	4	170	87	70	64		
03:33:00	1	1	24	1	0	4	15	4	2	1	0	5	8	7	9	21	5	2	4	163	65	78	60	
07:11:00	1	1	23	0	0	2	1	1	2	0	1	6	4	7	15	5	3	5	180	65	80	47		
09:10:00	1	1	20	1	0	10	8	2	1	3	0	2	4	4	8	6	2	5	165	70.3	95.2	68.2		
16:22:00	1	1	26	0	0	8	20	10	4	1	0	2	5	5	6	18	2	1	2	175	49	82	44	
16:33:00	1	1	18	1	1	7	5	5	3	4	0	1	1	2	3	13	4	1	2	168	55.8	65.77	51.3	
17:45:00	1	1	18	1	1	12	4	4	1	4	0	4	4	4	5	4	1	3	167	75	87	72		
02:08:00	1	1	17	1	1	7	5	1	1	5	0	1	4	6	5	6	1	3	167	92	57.2	50		
12:13:00	1	1	24	1	0	5	5	3	1	3	0	2	5	4	6	18	7	1	2	149	102	129	88	
21:42:00	1	0	36	0	1	20	40	10	5	2	0	2	4	2	5	16	4	3	3	190	70	75	62	
22:30:00	1	1	16	1	0	50	10	3	1	2	0	2	5	8	7	4	1	2	158	44	90	38.5		
17:42:00	1	1	16	1	0	10	5	4	3	3	0	1	4	4	8	4	2	3	173	69	70	64		
16:57:00	1	1	22	1	1	7	100	4	2	2	0	3	5	4	6	17	4	1	3	165	87	70	57	
16:19:00	1	1	17	1	0	yes	yes	4	3	6	0	1	4	7	9	16	7	1	4	155	50.3	63.5	48	
23:21:00	1	1	16	1	0	1	10	1	1	5	0	2	1	5	5	15	6	1	4	168	94	94	16	
17:28:00	1	1	16	1	0	5	20	5	1	6	0	2	3	2	9	15	7	1	4	164	63.6	70.5	56.1	
17:33:00	1	1	18	1	0	3	6	5	3	4	0	2	2	5	6	3	2	1	161	45	90	42		
19:41:00	1	1	21	1	1	Yes	Yes	No	No	5	1	2	6	10	8	16	5	1	3	175	72	116	57	
00:44:00	1	1	24	1	1	20	50	5	3	1	0	5	8	9	7	16	4	1	2	161	50	94	46	
10:50:00	1	1	28	0	1	6	4	4	4	2	0	1	2	3	7	4	2	6	2	151	70	77	4	
18:52:00	1	1	17	1	0	30	40	10	5	3	0	4	4	6	5	15	4	1	2	165	48.9	61.7	43.1	
00:23:00	1	1	19	1	0	9	2	5	3	4	0	3	1	1	4	13	5	1	3	141	94	55	50	
01:13:00	1	1	19	0	0	3	4	5	6	0	1	1	1	2	4	8	2	4	150	60	80	50		
01:22:00	1	0	23	1	0	8	6	8	2	2	0	2	1	1	5	18	7	2	1	151	60	55	50	
01:39:00	1	1	16	0	0	2	1	2	1	2	1	0	1	1	3	9	4	4	140	100	120	80		
15:04:00	1	1	35	0	0	yes	yes	yes	yes	4	1	4	1	4	3	6	14	5	3	5	163	77	91	61
15:50:00	1	1	22	1	1	4	10	5	1	6	0	4	5	4	7	16	5	2	4	170	76	85	56	
16:19:00	1	1	19	1	0	4	2	2	1	1	0	4	4	1	6	16	5	2	9	158	90	90	40	



<a href="#"><u>maximum(x)</u></a>	highest value of the time series x.
<a href="#"><u>mean(x)</u></a>	mean of x
<a href="#"><u>mean_abs_change(x)</u></a>	mean over the absolute differences between subsequent time series values which is
<a href="#"><u>mean_change(x)</u></a>	mean over the absolute differences between subsequent time series values which is
<a href="#"><u>median(x)</u></a>	median of x
<a href="#"><u>minimum(x)</u></a>	lowest value of the time series x.
<a href="#"><u>standard_deviation(x)</u></a>	standard deviation of x
<a href="#"><u>variance(x)</u></a>	variance of x
<a href="#"><u>kurtosis(x)</u></a>	kurtosis of x (calculated with the adjusted Fisher-Pearson standardized moment coefficient G2).
<a href="#"><u>skewness(x)</u></a>	sample skewness of x (calculated with the adjusted Fisher-Pearson standardized moment coefficient G1).
<a href="#"><u>sum_values(x)</u></a>	sum over the time series values
<a href="#"><u>length(x)</u></a>	length of x
<a href="#"><u>first_location_of_maximum(x)</u></a>	first location of the maximum value of x.
<a href="#"><u>first_location_of_minimum(x)</u></a>	first location of the minimal value of x.
<a href="#"><u>last_location_of_maximum(x)</u></a>	relative last location of the maximum value of x.
<a href="#"><u>last_location_of_minimum(x)</u></a>	last location of the minimal value of x.
<a href="#"><u>count_above_mean(x)</u></a>	number of values in x that are higher than the mean of x
<a href="#"><u>count_below_mean(x)</u></a>	number of values in x that are lower than the mean of x
<a href="#"><u>longest_strike_above_mean(x)</u></a>	length of the longest consecutive subsequence in x that is bigger than the mean of x
<a href="#"><u>longest_strike_below_mean(x)</u></a>	length of the longest consecutive subsequence in x that is smaller than the mean of x

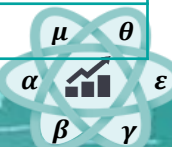




<a href="#"><u>linear_trend(x, param)</u></a>	linear least-squares regression for the values of the time series
<a href="#"><u>agg_linear_trend(x, param)</u></a>	linear least-squares regression for values of the time series that were aggregated over chunks.
<a href="#"><u>quantile(x, q)</u></a>	q quantile of x.
<a href="#"><u>number_peaks(x, n)</u></a>	number of peaks of at least support n in the time series x.
<a href="#"><u>abs_energy(x)</u></a>	absolute energy of the time series which is the sum over the squared values
<a href="#"><u>energy_ratio_by_chunks(x, param)</u></a>	sum of squares of chunk (10_segments)
<a href="#"><u>absolute_sum_of_changes(x)</u></a>	sum over the absolute value of consecutive changes in the series x
<a href="#"><u>change_quantiles(x, ql, qh, isabs, f_agg)</u></a>	First fixes a corridor given by the quantiles ql and qh of the distribution of x.
<a href="#"><u>percentage_of_reoccurring_datapoints_to_all_data_points(x)</u></a>	percentage of unique values, that are present in the time series more than once.
<a href="#"><u>percentage_of_reoccurring_values_to_all_values(x)</u></a>	ratio of unique values, that are present in the time series more than once.
<a href="#"><u>range_count(x, min, max)</u></a>	Count observed values within the interval [min, max).
<a href="#"><u>value_count(x, value)</u></a>	Count occurrences of value in time series x.
<a href="#"><u>sum_of_reoccurring_data_points(x)</u></a>	sum of all data points, that are present in the time series more than once.
<a href="#"><u>sum_of_reoccurring_values(x)</u></a>	sum of all values, that are present in the time series more than once.
<a href="#"><u>has_duplicate(x)</u></a>	Checks if any value in x occurs more than once
<a href="#"><u>has_duplicate_max(x)</u></a>	Checks if the maximum value of x is observed more than once
<a href="#"><u>has_duplicate_min(x)</u></a>	Checks if the minimal value of x is observed more than once
<a href="#"><u>approximate_entropy(x, m, r)</u></a>	Implements a vectorized Approximate entropy algorithm.
<a href="#"><u>sample_entropy(x)</u></a>	Calculate and return sample entropy of x.
<a href="#"><u>binned_entropy(x, max_bins)</u></a>	First bins the values of x into max_bins equidistant bins.

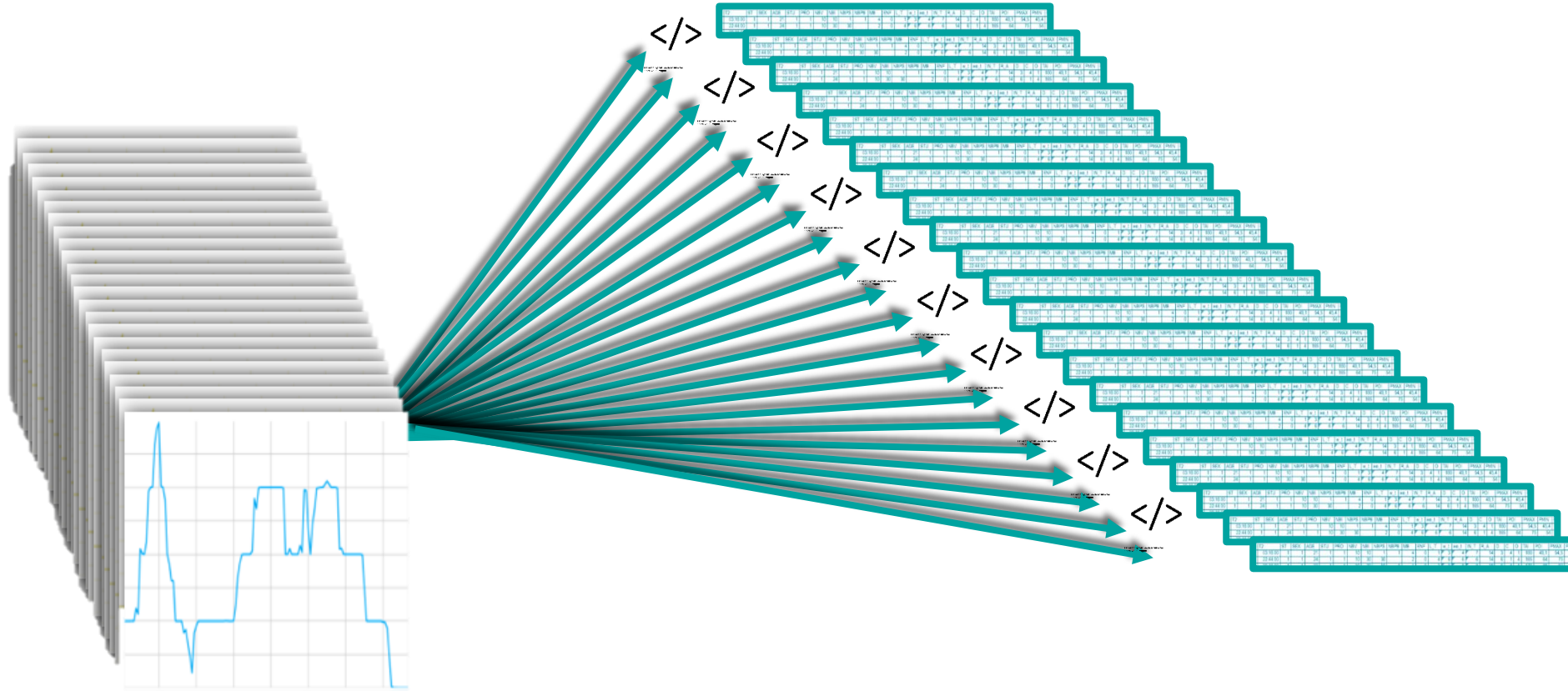
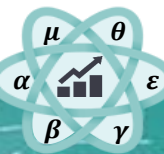


<a href="#"><u>ar_coefficient(x, param)</u></a>	This feature calculator fits the unconditional maximum likelihood of an autoregressive AR(k) process.
<a href="#"><u>agg_autocorrelation(x, param)</u></a>	value of an aggregation function f_agg (e.g.
<a href="#"><u>autocorrelation(x, lag)</u></a>	autocorrelation of the specified lag, according to the formula [1]
<a href="#"><u>partial_autocorrelation(x, param)</u></a>	value of the partial autocorrelation function at the given lag.
<a href="#"><u>augmented_dickey_fuller(x, param)</u></a>	The Augmented Dickey-Fuller test is a hypothesis test which checks whether a unit root is present.
<a href="#"><u>c3(x, lag)</u></a>	Measure of non linearity in the time series.
<a href="#"><u>cwt_coefficients(x, param)</u></a>	Continuous wavelet transform for the Ricker wavelet, also known as the “Mexican hat wavelet” which is
<a href="#"><u>fft_coefficient(x, param)</u></a>	fourier coefficients of the one-dimensional discrete Fourier Transform for real input by fast
<a href="#"><u>friedrich_coefficients(x, param)</u></a>	Coefficients of polynomial , which has been fitted to
<a href="#"><u>index_mass_quantile(x, param)</u></a>	Those apply features calculate the relative index i where q% of the mass of the time series x lie left of i.
<a href="#"><u>large_standard_deviation(x, r)</u></a>	Boolean variable denoting if the standard dev of x is higher than 'r'*(max - min of x.) "r=0.05:1"
<a href="#"><u>variance_larger_than_standard_deviation(x)</u></a>	Boolean variable denoting if the variance of x is greater than its standard deviation.
<a href="#"><u>symmetry_looking(x, param)</u></a>	Boolean variable denoting if the distribution of x looks symmetric.
<a href="#"><u>max_langevin_fixed_point(x, r, m)</u></a>	Largest fixed point of dynamics : $\text{argmax}_x \{h(x)=0\}$ estimated from polynomial ,
<a href="#"><u>number_crossing_m(x, m)</u></a>	number of crossings of x on m.
<a href="#"><u>number_cwt_peaks(x, n)</u></a>	This feature calculator searches for different peaks in x.
<a href="#"><u>ratio_beyond_r_sigma(x, r)</u></a>	Ratio of values that are more than r*std(x) (so r sigma) away from the mean of x.
<a href="#"><u>ratio_value_number_to_time_series_length(x)</u></a>	factor which is 1 if all values in the time series occur only once, and below one if this is not the case.
<a href="#"><u>set_property(key, value)</u></a>	This method decorator that sets the property key of the function to value
<a href="#"><u>spkt_welch_density(x, param)</u></a>	Cross power spectral density of the time series x at different frequencies.
<a href="#"><u>time_reversal_asymmetry_statistic(x, lag)</u></a>	Proposed by Fulcher and Jones as a promising feature to extract from time series.





# Construcción de características de la ts

[illegible]

# Construcción de características de la ts

## Algoritmo resumen mediante código Spark

```
DB = spark.read(...)
#[Row(Hash='id1', valor=100', periodo='201208'), Row(...), ... ]

DB = DB.groupBy(Hash).agg(F.collect_list(valor), F.collect_list(periodo))
#[Row(Hash='id1', valor=[100,105,110,...]', periodo=['201208', '201209', '201210',...], Row(...),...]
```

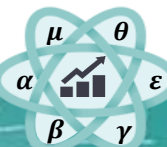
```
def Create_Features(pyspark_row):
    ts = create_pandas(pyspark_row)
    ts = fix_time_series(ts)

    f_fuga = tsfresh.extract_features(ts, column_id='id', column_sort='periodo', n_jobs=0)\
        .T.id.to_dict()

    r = {'Hash':pyspark_row.Hash, 'features_fuga':f_fuga}

    return Row(**r)
```

```
caracteristicas = DB.rdd.map(lambda x: all_features_entrenamiento(x))
```





Project Jupyter exists to develop open-source software, open-standards, and services for interactive computing across dozens of programming languages.

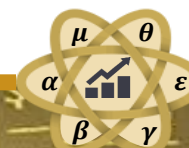
This screenshot shows the Jupyter web interface. At the top, there's a "jupyter" logo and a "Logout" button. Below that are tabs for "Files", "Running", "Clusters", "Conda", and "Nbextensions". The "Files" tab is active, showing a file browser for the path "/ Riesgos / Segmentacion SARLAFT FONDO". A filter input is present. A list of files and folders is shown, including "..", "data", and "Desarrollo". A context menu is open over the "data" folder, listing options: "Text File", "Folder", "Terminals Unavailable", "Notebooks", "Python [conda root]", "Python [default]", and "R". There are also "Upload", "New", and "Refresh" buttons at the top of the file list.

This screenshot shows a Jupyter notebook titled "Untitled". The top bar indicates "Last Checkpoint: 7 minutes ago (unsaved changes)". The menu bar includes "File", "Edit", "View", "Insert", "Cell", "Kernel", "Widgets", "Help", and "Snippets". The toolbar contains various icons for file operations and cell management. The notebook content includes a "Table of Contents" section with links to "1 Lectura de datos" and "2 Preparacion de datos". Below this is a section titled "1 Lectura de datos". The first code cell contains the following Python code:

```
InfoCliente= read.table("InfoCliente.txt", header=TRUE)
head(datos)
```

The second code cell contains the following R code:

```
library(plyr)
dt <- data.frame(age=rchisq(20,10),group=sample(1:2,20,rep=T))
ddply(dt,~group,summarise,mean=mean(age),sd=sd(age))
```

The third code cell is empty.



## IBM Data Science Experience

Watson Data Platform

Welcome Alex!

IBM Data Science Experience is part of Watson Data Platform.

[Try out other Watson Data Platform apps.](#)

Get started with key tasks



New project



New notebook



New model



New streams flow

New in the community

Explore

**ARTICLE**  
Neural networks for beginners: popular types...

**AUTHOR**  
Stats and Bots

**DATE**  
Nov 16, 2017

**TOPIC**  
Neural Networks

**FORMAT**  
Web page



**DATA SET**  
Health insurance (2015): United States...

**AUTHOR**  
IBM

**DATE**  
Oct 31, 2017

**TOPIC**  
Society



**NOTEBOOK**  
Social media insights with Watson Developer...

**AUTHOR**  
IBM

**DATE**  
Nov 02, 2017

**TOPIC**  
Society



**TUTORIAL**  
Probabilistic Graphical Models...

**AUTHOR**  
Stats and Bots

**DATE**  
Nov 02, 2017

**LEVEL**  
Intermediate

**TOPIC**  
Data Science +1



IBM Data Science Experience | Projects Tools Data Services Community US South

My Projects > Churn > AC-02-Features\_Time\_Series ...

File Edit View Insert Cell Kernel Help Not Trusted | Python 3.5 (Experimental) with Spark 2.0

Format

Code

### 3.3. Comparar Modelos

[Top](#)

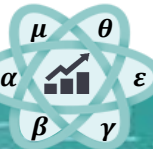
EN esta sección se compararán los diferentes modelos entrenados basado en el accuracy, AUC y recall

```
In [35]: var = 'accuracy'
for m in modelos:
    sns.distplot(scores[m][var], label=var+' '+m)
sns.plt.legend()
```

/usr/local/src/conda3\_runtime.v16/4.1.1/lib/python3.5/site-packages/statsmodels/nonparametric/kdetools.py:20: VisibleDeprecationWarning:


using a non-integer number instead of an integer will result in an error in the future

```
Out[35]: <matplotlib.legend.Legend at 0x7fc3280449b0>
```







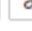













FLOW

Flow
Cell
Data
Model
Score
Admin













Untitled Flow

assist

?

Assistance

	Routine	Description
	<a href="#">importFiles</a>	Import file(s) into H <sub>2</sub> O
	<a href="#">getFrames</a>	Get a list of frames in H <sub>2</sub> O
	<a href="#">splitFrame</a>	Split a frame into two or more frames
	<a href="#">mergeFrames</a>	Merge two frames into one
	<a href="#">getModels</a>	Get a list of models in H <sub>2</sub> O
	<a href="#">getGrids</a>	Get a list of grid search results in H <sub>2</sub> O
	<a href="#">getPredictions</a>	Get a list of predictions in H <sub>2</sub> O
	<a href="#">getJobs</a>	Get a list of jobs running in H <sub>2</sub> O
	<a href="#">buildModel</a>	Build a model
	<a href="#">runAutoML</a>	Automatically train and tune many models
	<a href="#">importModel</a>	Import a saved model
	<a href="#">predict</a>	Make a prediction

getModel "model\_glm"

372ms

Model

Model ID: model\_glm

Algorithm: Generalized Linear Modeling

Actions:

Refresh

Predict...

Download POJO

Download Model Deployment Package (MOJO)

Delete

Export

Inspect

Download Gen Model

MODEL PARAMETERS

SCORING HISTORY

STANDARDIZED COEFFICIENT MAGNITUDES

poblacion

outlier

snnlv

tasa\_empleo\_13\_ciudades

OUTLINE

FLOW

CLIPS

HELP

Help

Using Flow for the first time?

Quickstart Videos

Or, [view example Flows](#) to explore and learn H<sub>2</sub>O.

STAR H2O ON GITHUB!

Star

2,575

GENERAL

Flow Web UI ...

... Importing Data

... Building Models

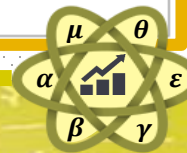
... Making Predictions

jupyter AC-02-Clasification Model-AF Last Checkpoint: 10/04/2017 (autosaved) ?

File Edit View Insert Cell Kernel Navigate Widgets Help Snippets Python [default]

Disabled Snippets Code CellToolbar Autosave interval (min): 2 (default) Notify:

In [ ]: # initialize grid search  
 gsearch = H2OGradientBoostingEstimator,  
   hyper\_params=hyper\_parameters,  
   search\_criteria=search\_criteria)  
  
 print(' - Obteniendo HyperParametros del: ',clf\_type)  
 gsearch.train(x=X,  
                   y=y,  
                   training\_frame=trainh2o,  
                   validation\_frame=validh2o)  
  
 - Obteniendo HyperParametros del: xgboost  
 gbm Grid Build progress: 100%



# Construcción de características

## Algoritmo resumen H2O mediante código python

```
baseh2o = h2o.H2OFrame(base.values.tolist(), column_names=base.columns.values.tolist())
```

```
y = 'cliente_fuga'
```

```
X = selected_features
```

```
#Grid search
```

```
hyper_parameters = {'ntrees':[100,500,...], 'max_depth':[4,5,...], 'sample_rate':np.arange(0.3,0.6,...),  
'col_sample_rate':np.arange(0.3,0.6,...), ...}
```

```
search_criteria = {'strategy':'RandomDiscrete', 'max_runtime_secs':10800, 'seed': 1234}
```

```
gsearch = H2OGridSearch(H2OGradientBoostingEstimator, hyper_params=hyper_parameters, search_criteria=search_criteria)
```

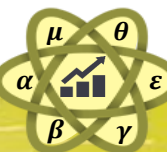
```
gsearch.train(x=X, y=y, training_frame=trainh2o, validation_frame=validh2o)
```

```
#Train Model
```

```
model_h2o = H2OGradientBoostingEstimator(ntrees=700, max_depth=5, col_sample_rate=0.7, col_sample_rate_per_tree=0.9, \  
                                          sample_rate=0.9, learn_rate=0.01, min_rows=50, seed=1234)
```

```
model_h2o.train(X, y, training_frame=train_completoh2o, validation_frame=testh2o)
```

```
model_h2o.model_id='Modelo_Fuga_h2o'
```





Alex Castaño Ballesteros



[github.com/acastanob](https://github.com/acastanob)



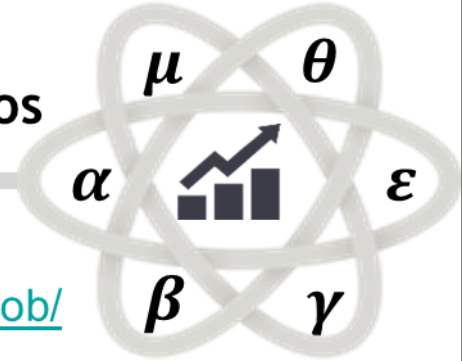
[linkedin.com/in/acastanob/](https://www.linkedin.com/in/acastanob/)



[acastanob@gmail.com](mailto:acastanob@gmail.com)



[312 840 88 87](tel:3128408887)



312 840 88 87



acastanob@gmail.com

# GRACIAS

