

# Hackthon 6

Hugo Castillo

2018-01-06

## Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

## 1 Exploratory Data Analysis

The provided dataset contains personal information of several people (8164 samples), such as date of birth and current employment and maps it to the observation if the person became unemployed in the next 12 months.

The dataset was made available without any information about its contents except for the property names. Each one was analysed for insights and anomalies.

The following list details what could be obtained exploring the data (column name, value types and values).

**id** Integer, identifier for each subject, all values distinct no problems detected.

**target** Integer, did the subject become unemployed in the following 12 months.

**birth date** Date (YYYY-MM-DD), date of birth of the subject. The youngest being 2016-02-10 and the oldest 1928-01-09. Frequency plot (see fig. 1 on page 9) shows nothing surprising.

**country of origin** Categorical, country names come in a variety of inconsistent formats.

**domestic relationship type** Categorical, (see table 3 on page 6) who the subject lives with. Categories are unclear and ill defined. Moreover it is inconsistent with the *domestic status*, there are a considerable number of entries classified as *domestic relation type-never married* and *domestic status-d* (presumably divorced) (see table 5 on page 6).

**domestic status** Categorical, marital status or if has married several times (see table 4 on page 6). By elimination category *d* is supposedly divorced.

**earned dividends** Numerical, monetary amount (currency not specified). Return from distribution of corporate earnings, it's 0 for all samples.

**ethnicity** Categorical, categories have funny names (see table 6 on page 7).

**gender** Categorical, all female dataset (see table 7 on page 7).

**job type** Categorical, current job type of the subject (see table 8 on page 7), government, self employed, item etc...

**interest earned** Numerical, monetary amount (currency not specified). Returns from loaning money (see fig. 2 on page 10).

**monthly work** Numerical, number of hours of work per month (see fig. 3 on page 11)

**profession** Categorical, type of profession (see table 9 on page 7).

**school level** Categorical, subject level of schooling (see table 10 on page 8).

## 2 Data pre-processing

This section details any and all data pre-processing before modelling. The first section section 2.1 explains what was done to convert the dataset to usable, unambiguous types. Afterwards section 2.2 details the what was done to select and improve the features feed to the model. Finally section 2.3 on the next page breaks down the train/test split.

### 2.1 Data cleaning

On data import the *birth dates* where converted to naive dates as no timezone information was provided. It's doubtful the timezone would provide any useful information.

The *countries of origin* where converted to the corresponding *ISO 3166-1 alpha-2* representation. Some inputs required special rules. Especially ambiguous was *dr* which represents no country code this was converted to Dominican Republic (DO) even if the race for these inputs suggests it's not (mostly white).

All categorical data columns was kept as is, there was not enough information to reach any conclusion.

### 2.2 Feature engineering

*birth dates* where converted to timestamps. In the dataset earned dividends and gender do not change, these properties where dropped since they convey no useful information. If new samples include this value this decision will be reconsidered. All categorical data was turned into dummy class variables.

### 2.3 Train/test splitting

The train/test split was done holding out .4 of the data for final validation. Furthermore training was done using a shuffle split with .3 for testing.

## 3 Modelling

## 4 Initial tests

To get a feeling of the baseline performance of the models available in the *scikit* package several were tried with the default parameters (except where not possible), results in table 1.

Table 1:	
Model	AUC ROC Score
GradientBoostingClassifier	0.9212
AdaBoostClassifier	0.9008
BaggingClassifier	0.8428
VotingClassifier	0.8233
RandomForestClassifier	0.8160
ExtraTreesClassifier	0.8154
QuadraticDiscriminantAnalysis	0.8053
KNeighborsClassifier	0.7267
GaussianProcessClassifier	0.7022
DecisionTreeClassifier	0.6729
SVC	0.6698
SGDClassifier	0.6377
GaussianNB	0.6718

With the default parameters there is, as expected, a clear dominance of ensemble models. The top 5 were selected for further parameter tuning.

## 5 Model tuning

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

## 6 Model selection

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

## 7 Evaluation

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

## A EDA Tables

Table 2: Country of origin value counts.

u.s.	7330
unknown	126
mexico	111
philippines	60
de	50
puerto rico	39
jamaica	34
cuba	34
el-salvador	30
canada	28
dr	27
gb	22
south	20
italy	18
columbia	17
haiti	17
china	17
vietnam	17
guatemala	16
japan	15
poland	14
peru	11
taiwan	11
thailand	11
fr	9
trinadad/tobago	8
india	7
nicaragua	7
portugal	6
honduras	6
laos	6
ecuador	6
iran	6
ireland	5
us territory	4
hong	4
scotland	3
hungary	3
greece	3
yugoslavia	3
cambodia	2
netherlands	1

Table 3: Domestic relationship type value counts.

not living with family	2919
never married	2063
living with child	1750
has husband	1106
living with extende family	325
has wife	1

Table 4: Domestic status value counts.

single	3662
d	2073
married 2	1170
spouse passed	599
divorce pending	486
married not together	163
married 1	11

Table 5: Domestic relationship type grouped by domestic status counts.

domestic status	domestic relationship type	count
d	living with child	116
	living with extende family	47
	never married	1006
	not living with family	904
divorce pending	living with child	40
	living with extende family	25
	never married	293
	not living with family	128
married 1	has husband	10
	living with child	1
	has husband	1096
	has wife	1
married 2	living with child	28
	living with extende family	42
	not living with family	3
	living with child	25
married not together	living with extende family	7
	never married	73
	not living with family	58
	living with child	1530
single	living with extende family	174
	never married	449
	not living with family	1509
	living with child	10
spouse passed	living with extende family	30
	never married	242
	not living with family	317

Table 6: Ethnicity value counts.

white and privileged	6523
afro american	1210
asian	262
american indian	88
other	81

Table 7: Gender value counts.

Female	8164
--------	------

Table 8: Job type value counts.

private	5919
unknown	620
local-gov	618
state-gov	368
self-emp-not-inc	303
federal-gov	236
self-emp-inc	94
without-pay	4
never-worked	2

Table 9: Profession value counts.

secretarial	1949
other	1423
specialist technician	1096
sales	978
C-level	842
unknown	622
mechanic	420
technology support	247
vocational	184
household labor	131
estate employee	108
defense contractor	58
trucking	58
agriculture	48

Table 10: School level value counts.

secondary	2594
entry level college	2165
college graduate	1188
basic vocational	373
some post graduate	355
secondary 11	341
advanced vocational	326
10th	248
secondary-7 through 8	123
secondary 12	106
secondary-9	104
secondary-5 through 6	72
advanced post graduate	61
primary school	58
primary 1 through 4	37
kindergarten	13



## B EDA Figures

Figure 1: Date of birth frequency.

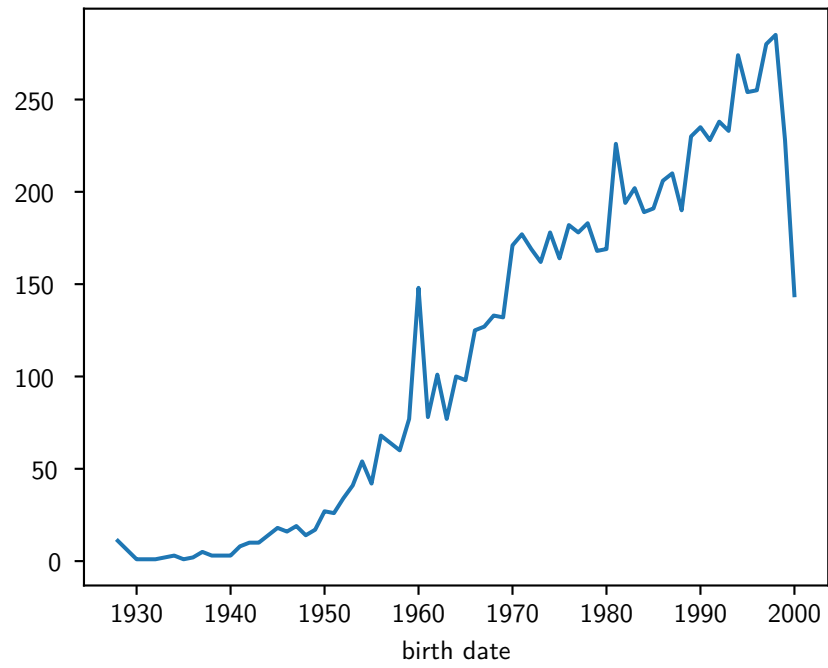


Figure 2: Interest earned frequency (logarithmic scale).

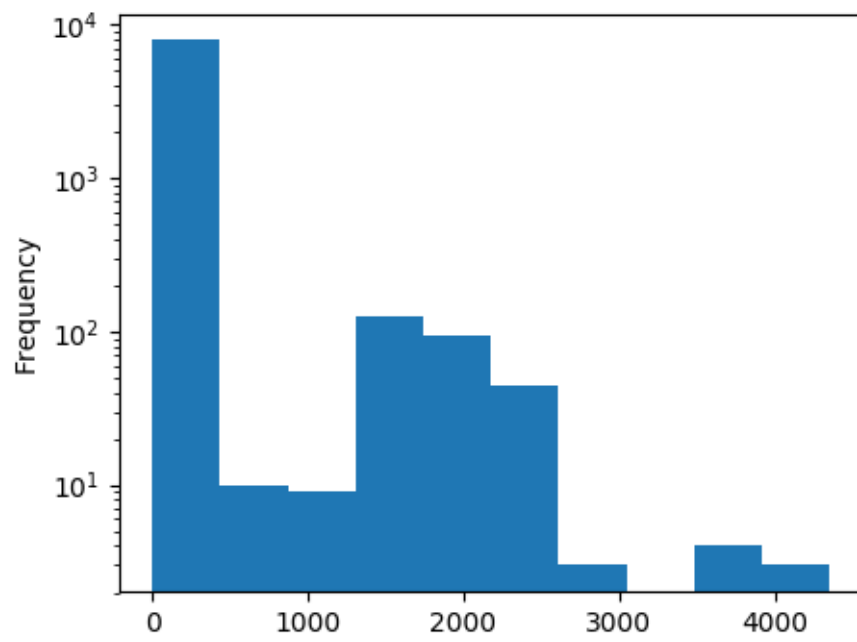


Figure 3: Monthly work frequency.

