

# Capstone Project I

- **Problem statement and motivation:**

The client is Kiva.org. The goal is to understand what may cause a higher delay between the ‘posted time’ (the time at which the loan is posted on Kiva by the field agent), the ‘funded time’ (the time at which the loan posted to Kiva gets 100% funded by lenders), and the ‘disbursed time’ (the time at which the loan is disbursed by the field agent to the borrower). Subsequently, have a sense of its urgency and possibly discover hints on how they could resolve the issue. The problem is important since this could be a way of understanding a fundamental step in its business process, and take actions wherever needed. Lenders would get their borrowed money faster and it could possibly be an incentive to lend again more times, and/or with more money. Depending on what insights the analysis of the data provides, Kiva can then think about what they will do.

## Code for the project:

<https://github.com/MigBap/Springboard-Capstone-Project-I/blob/master/Capstone%20Project%20I%20--%20All.ipynb>

## Table of Contents:

### [1\) Data inspection and cleaning steps](#)

#### [1.1\) Main variables inspection and cleaning](#)

['Posted\\_time'](#)

['Funded\\_time'](#)

['Disbursed\\_time'](#)

['Posted\\_time' vs 'funded\\_time'](#)

#### [1.2\) Other variables](#)

#### [1.3\) Searching for outliers](#)

['loan\\_amount' and 'funded\\_amount'](#)

['posted\\_time' and 'funded\\_time'](#)

['posted\\_time' and 'disbursed\\_time'](#)

['funded\\_time' and 'disbursed\\_time'](#)

### [2\) Exploratory Data Analysis](#)

#### [2.1\) General EDA: the loan cycle](#)

#### [2.2\) Do countries have influence regarding delays?](#)

[2.3\) Do bigger loan amounts take more time to fund, on average?](#)

[2.4\) Do sectors of the loan influences delays?](#)

[2.5\) Field partners](#)

[2.6\) Number of lenders](#)

[2.7\) Does the repayment interval affect delays?](#)

[2.8\) When the loan is not completely funded](#)

[2.9\) Currency of the loan](#)

[2.10\) Currency policy](#)

[2.11\) Does lender term influence delays?](#)

### [3\) Modeling](#)

[3.1\) Initial Feature Selection & Extraction](#)

[3.2\) Problem statement](#)

[3.3\) Models initial results](#)

### [4\) Dimension reduction & Feature selection](#)

### [5\) Hyperparameter tuning](#)

### [6\) Conclusion](#)

## 1) Data inspection and cleaning steps

After getting an overview of the data, I **changed some column names, for clarity, and changed the dtype of the main variables to datetime64[ns]**.

### 1.1) Main variables inspection and cleaning

#### 'Posted\_time'

- There were no missing values and the filled ones seemed to be ok - Apparently, the 'posted\_time' column is clean.

#### 'Funded\_time'

There are 64282 missing values and 1355325 filled values.

- **Regarding the missing values:**

- Even though some null-values might correspond to loans which were indeed funded, there is not a way to know how to fill those values. Therefore, I **removed the rows where there is a missing value in funded\_time but (loan\_amount -**

- funded\_amount) <= 0**, so that there is no ambiguity. Note: I considered ok when the funded\_amount was larger than the loan\_amount, hence the < ).
- The remaining missing values were kept due to several reasons. Not every loan is funded or when there is a filled value in funded\_amount, because the loan may not be completely funded yet, that is, the funded\_amount is smaller than the loan\_amount.
  - **Regarding the filled values**, I compared with two columns to check if there were discrepancies: funded\_amount and status.
    - *funded\_time vs funded\_amount*:
      - Comparing with the funded\_amount column, there were no discrepancies.
      - However, when looking at the difference between loan\_amount and funded\_amount, there was inconsistency in the data. I then **removed the rows where the funded\_time was filled, but (loan\_amount - funded\_amount) > 0**. Note: the funded\_time is only filled when the loan is 100% funded by lenders.
    - *funded\_time vs status*:
      - The status column is composed of 4 hypotheses: funded, refunded, expired and fundraising.
      - Comparing with the *funded* status, everything was ok.
      - Comparing with the *refunded* status, everything was ok.
      - Comparing with the *expired* status, there was inconsistency in the data. There were two loans requests which supposedly expired but, at the same time, two values for the 'funded\_time' column were filled. In this case, however, the remaining variables seemed to indicate that the loan was indeed funded, so I considered them for analysis. For example, in the first case, what probably classified the status as expired was the fact that the funded\_time was slightly after (6 minutes) the planned\_expiration\_time. I **changed the status in these two cases to 'funded'**.
      - Comparing with the *fundraising* status, everything was ok.

## 'Disbursed\_time'

The "disbursed time" is the time at which the loan is disbursed by the field agent or group of lenders to the borrower.

- **Regarding the missing values:**
  - It may be the case that there was a posted\_time or funded\_time but the disbursal has not occurred yet. I will leave those values for further analysis since they could represent loans which were not delivered or just expected situations.

- **Regarding the filled values:**

- The timing of the disbursal can vary. For most Field Partner loans, the money is pre-disbursed, so the borrower can access the funds right away. Hence, it is not strange if the disbursed\_time is made even before the posted\_time. For direct loans, the money is disbursed only after the loan has been fully crowdfunded on the Kiva website. The disbursed\_time can then naturally occur after or before the posted\_time or the funded\_time. All of the values remained for analysis.

## 'Posted\_time' vs 'funded\_time'

The chronological order between the 'posted\_time' and the 'disbursal\_time' does not matter, as well as between the 'disbursal\_time' and the 'funded\_time'. What could bias the data is when the 'funded\_time' appears before the 'posted\_time'. There were indeed found cases where the funded\_time was way before the posted\_time, and they all corresponded to the first posting date. Therefore, I **removed these cases** for analysis.

## 1.2) Other variables

- The funded\_amount column was compared with the num\_lenders\_total column and everything was ok.
- If the funded\_amount is greater or equal to the loan\_amount, then the status should not appear as expired. No inconsistency was found.
- Checking for misspellings: No errors were found in the activity\_name and country\_name columns.
- Since there were only 23 missing values in the column 'lender\_term', I **did not consider them**.

## 1.3) Searching for outliers

### 'loan\_amount' and 'funded\_amount'

In the previous analysis, when looking at the filled values in the 'funded\_time' column, I did not include the cases where the funded\_amount was greater than the loan\_amount because it is acceptable. I checked, however, for outliers. Just looking at the general statistics of the differences, it was possible to see that they were not significant. Therefore, I **kept all of the data**.

### 'posted\_time' and 'funded\_time'

On one hand it was not necessary to look for outliers here since I had already excluded for analysis the cases where the funded\_time was before the posted\_time. Regarding the

remaining cases, I opted to maintain all of them, after making a comparison with the planned expiration time column.

### 'posted\_time' and 'disbursed\_time'

The chronological order of these two dates does not matter, because the disbursal can be made before the posted time for most of the entities (Field partners). It was necessary to check for outliers, though. I focused on the cases when the disbursal was made before the posted\_time.

- I defined a z-score function to detect outliers.
- With a threshold of 3, a few outliers were found. However, when examining them, in all of the cases the difference between the posted\_time and funded\_time was acceptable/minimal, which possibly means that although the disbursal was made long before the posted\_time, the loan got funded almost immediately. This could suggest that these cases were somehow managed by Kiva or the Field Partners. I then **opted for their maintenance** in the analysis. (A box plot was also presented).

### 'funded\_time' and 'disbursed\_time'

Here I focused on the cases on the cases where the disbursed\_time was filled after the loan was funded. I left apart from the cases where the disbursed\_time appeared before the funded\_time.

- After dividing the cases where there was a field partner involved and not, a few outliers were found, in both cases, using the previously defined z-score function with a threshold equal to 3. However, in both cases, when looking at the distribution of the outliers among countries, it was possible to see that they are somewhat representative, that is, they were concentrated in a small number of countries. I then **opted to keep the values** for further analysis.

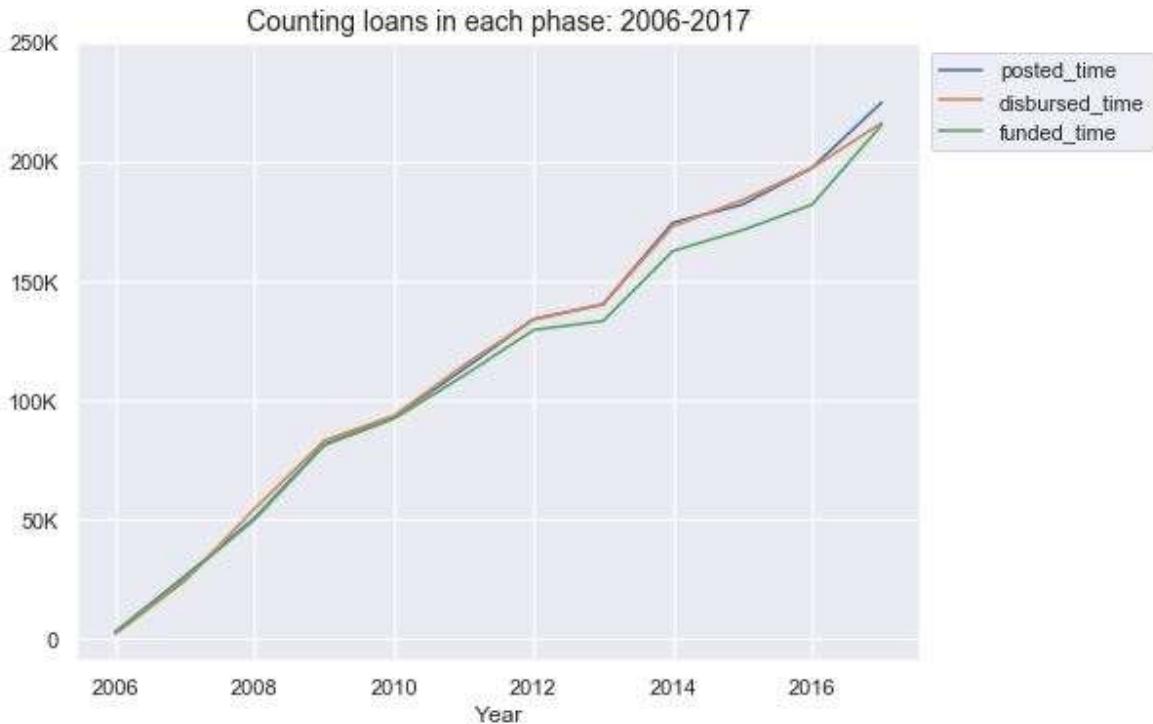
There was also 271 missing disbursed\_time values when the funded\_time was filled, all of them of USA and Kenya, and the values were off the scale. I **removed them**.

## 2) Exploratory Data Analysis

### 2.1) General EDA: the loan cycle

Now that the data is clean, it was time to investigate the data, first using visuals. To do that, I performed the following steps:

I began to study the characteristics of the main variables and how they interacted with each other over time. Looking at the evolution of the number of loans in each phase, during the period 2006-2017, it is clear that there are three major distinctive time periods:

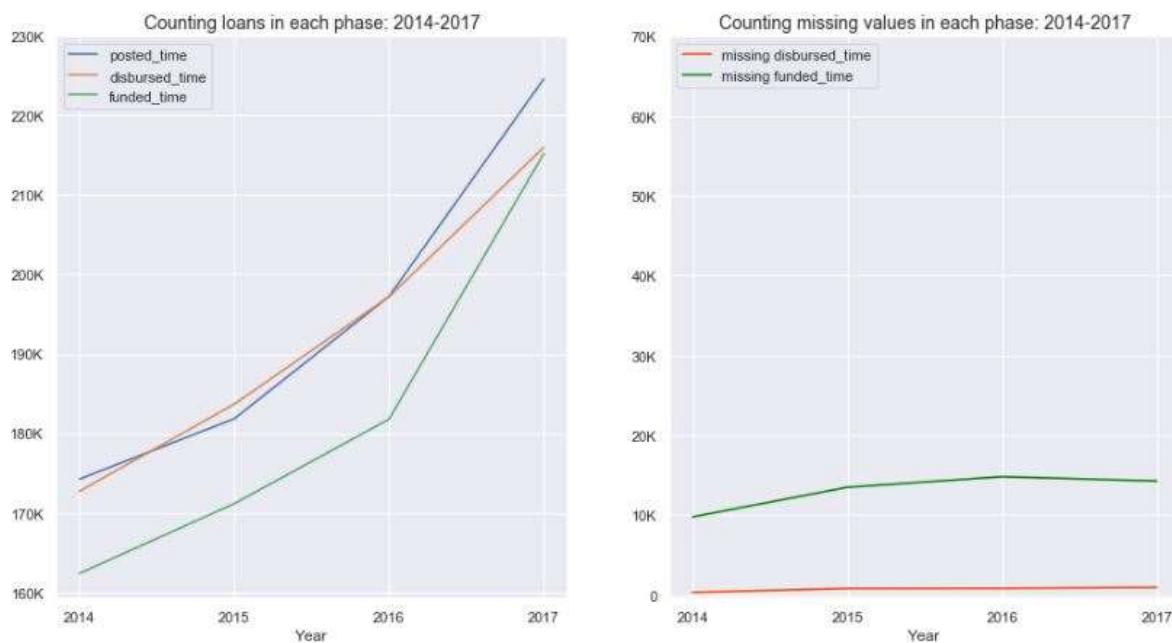


1. Until 2010, the amount of posted loans are closely matched in time by the number of fundings and disbursals.
2. After that, between 2010 and 2014, there is a clear separation between each of the main variables: it suggests that the loans, continuing with almost always a close match between the posting phase and disbursal, started to lose its efficiency regarding the funding of the projects.
3. After 2014 till the end, the differences continued to aggravate between the postings and fundings, and the number of disbursals and postings became just a bit volatile over time. Curiously, the number of fundings gained proximity with the number of disbursals, especially in the last year.

This could be still be related to missing values or other factors, so it deserved further exploration.

The process was smooth until 2010. Between 2010 and 2014, the funding started to lose track of the posted loans mainly between 2010 and 2011. This difference, however, could be due to the cyclical nature of the process, as it was later observed. Thereafter, the difference was mostly due to the funded missing values.

In the last phase, although the difference between the number of missing values of the disbursals and fundings remained relatively constant, the number of fundings augmented more than disbursals, equalizing them on the final stretch.

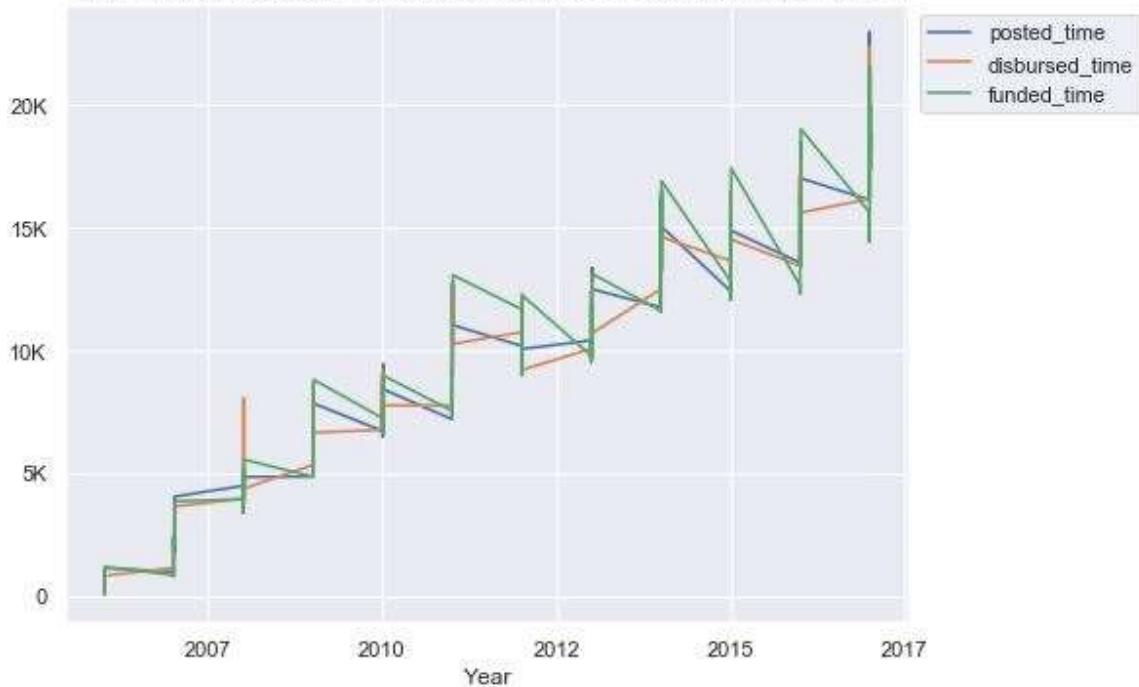


We can see that the deviation increase was mostly due to missing values until 2014. If we take that apart, the lines are somewhat at the same pace. This was not the case during the final 4 years, given the disbursals data we have. However, since it captures data at the end of the scale, and given the cyclical nature of the Kiva process, as shown below, it could still be related to the fact that the disbursal data was not yet inserted, lagging a few days or weeks.

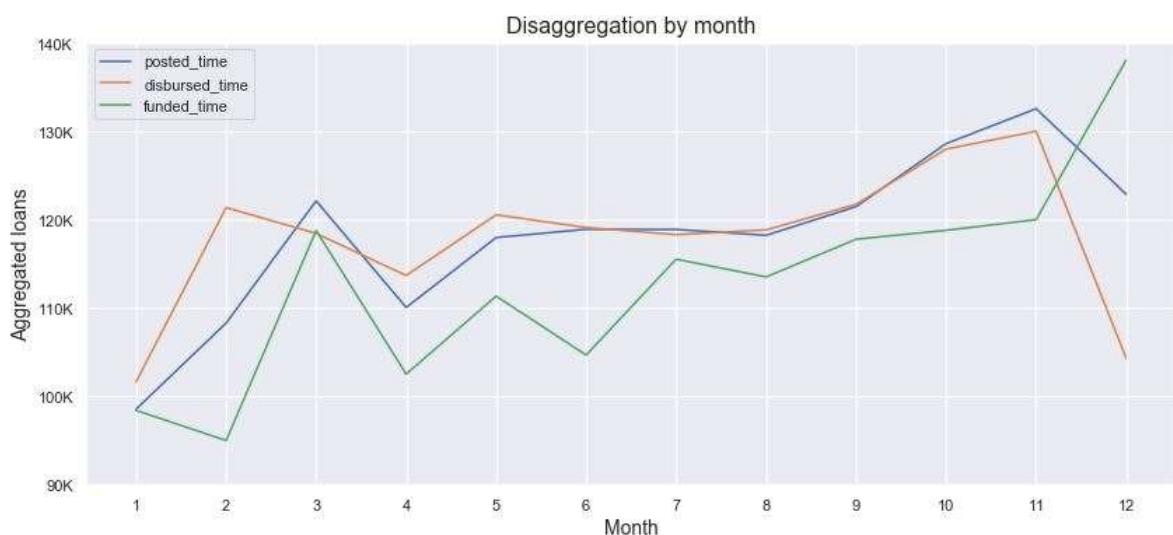
**In conclusion**, we see that from 2006 till 2010 the process was smooth and without significant delays. The number of loans not funded started then to increase year after year until 2017, ending with approximately 15K loans to fund. This does not mean delays got significantly worse year after year, but they're there.

When looking at the disaggregation by month, we see the **cyclical nature** of the journey of Kiva loans. These cycles suggest that Kiva gather many loan requests before they post on their website. The fundings and disbursals follow the announcements (or, since the disbursal is the first to occur, the posting and funding follow the disbursals).

Evolution of the number of loans in each phase: 2006-2017, by month

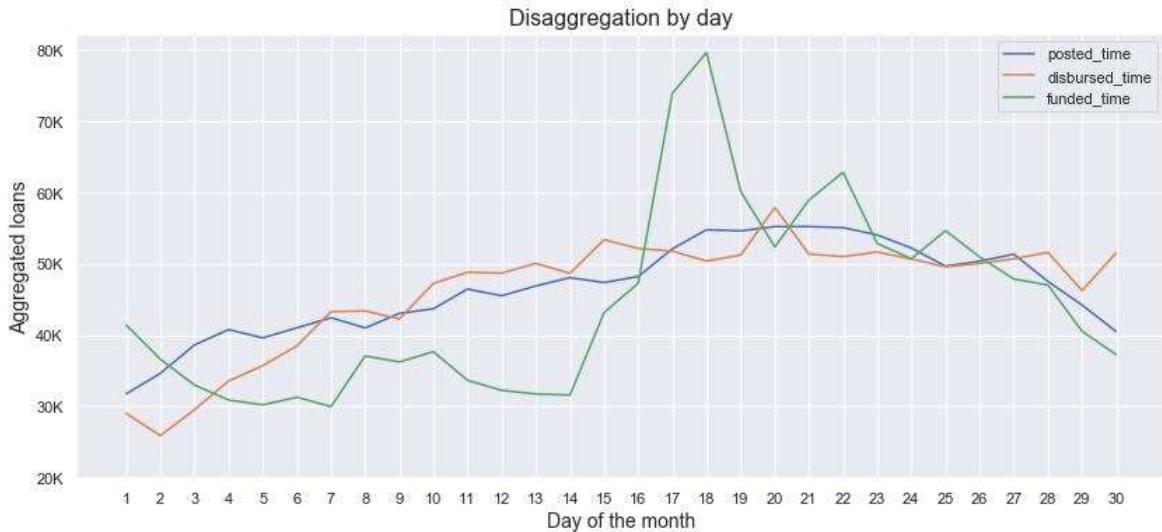


To further examine this cyclical nature, let us now display the aggregated values by month and finally, by day.



As we can see, there is a slight tendency to post, disburse and fund as the year comes to an end.

While the posting and disbursing periods relatively aggregate the same amount of loans throughout the month, there is a strong pressure to fund the loans during the **3rd week** of the month.

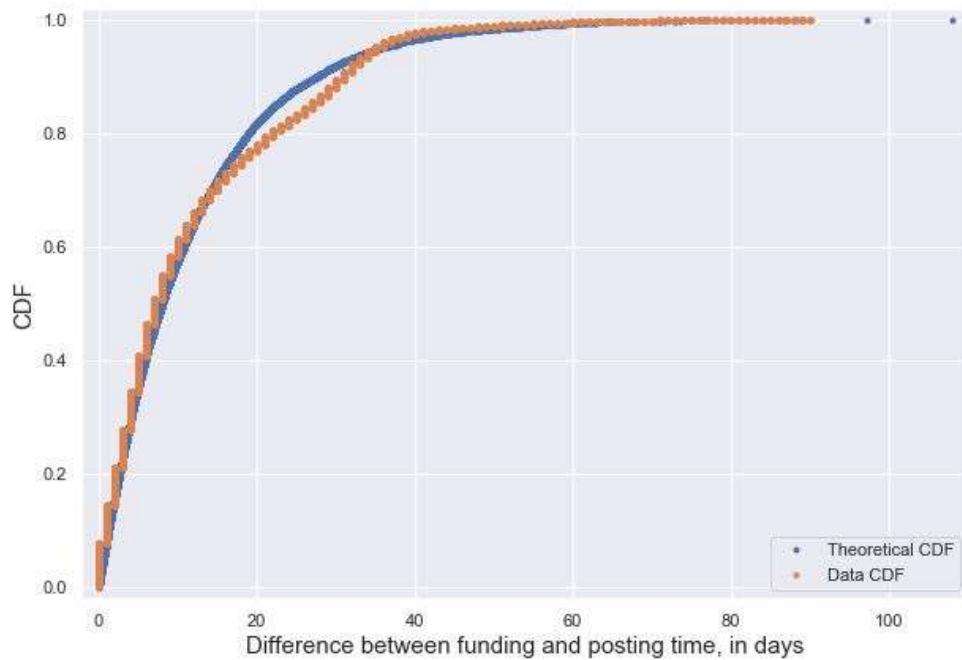


I then started to try to answer other questions:

I first focused on the time differences between the main variables:

98% of the differences between the time of disbursal and posting of loans occur on a range of -35 days and +32 days. Moreover, 90% of the differences are positive, that is, pre-disbursals. We could then immediately say that at least 90% of the disbursals are made by field partners, the only ones who can pre-disburse.

When we look at the time a loan takes to get funded since the time it is posted on Kiva website, we see that 99% of them are funded within 50 days, and almost 90% within a month. In fact, this behavior follows an exponential distribution, as presented in the CDF:

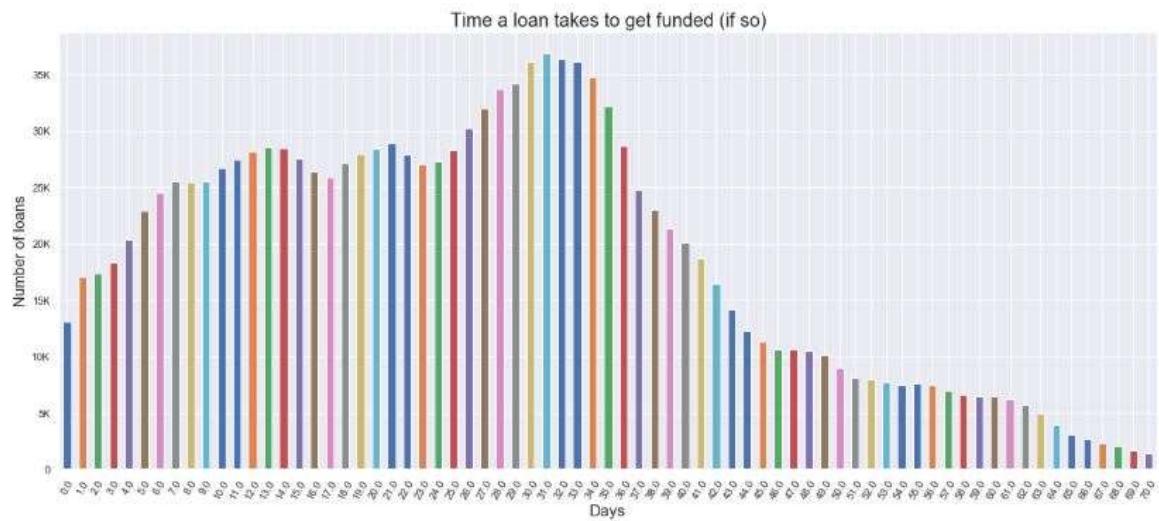


Close to 10% of the total loans are disbursed after the funding is complete, and 5% occur 14 days after the time of funding. This clear spike in the data suggests that 2 weeks could be a pre-staged arrangement. In most cases, a pre-disbursal occurs and only then the funding period commences. I also confirmed that the 5% cases where disbursal occurs 14 days after funding are distributed over time.

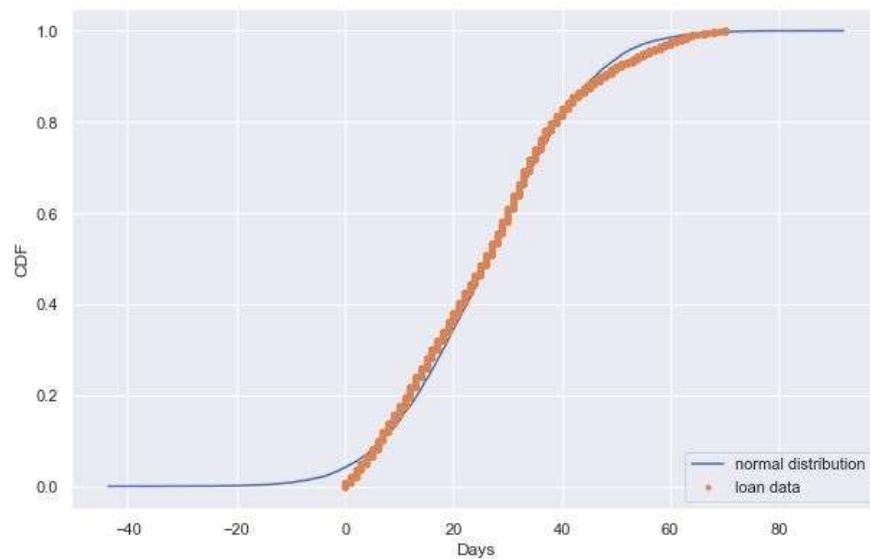
I then looked at the **real** time a loan takes to get funded, that is, looking at the difference between the funding time and the disbursal or posting time, depending on which of these occurred first. The following facts were checked:

- Close to 60% are funded within a month. Some irregularity (ups and downs) can be seen during this period. After that they start to get the funding in an exponential fashion way, as seen.
- 97% are funded within 2 months.

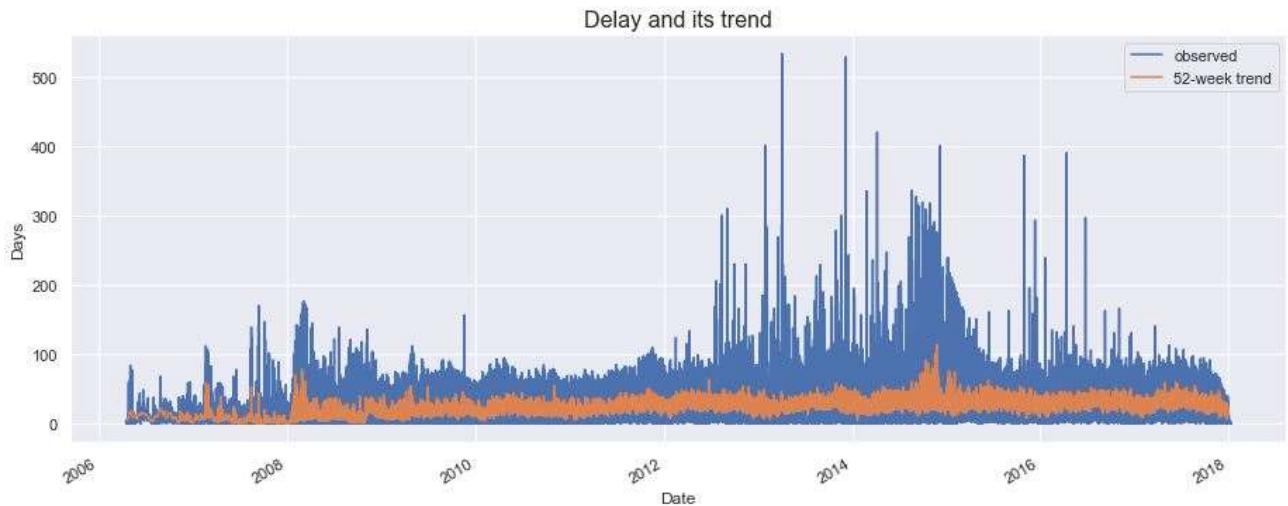
As a reminder, 8% of the total loans are only disbursed after they are funded (with 7% within two weeks). Considering these cases, this means that 99% of the loans are funded within 2 months and 2 weeks. Using D'Agostino's K-squared test, we reject the hypothesis that the funding period follows a normal distribution - with a null p-value.

**PDF:**

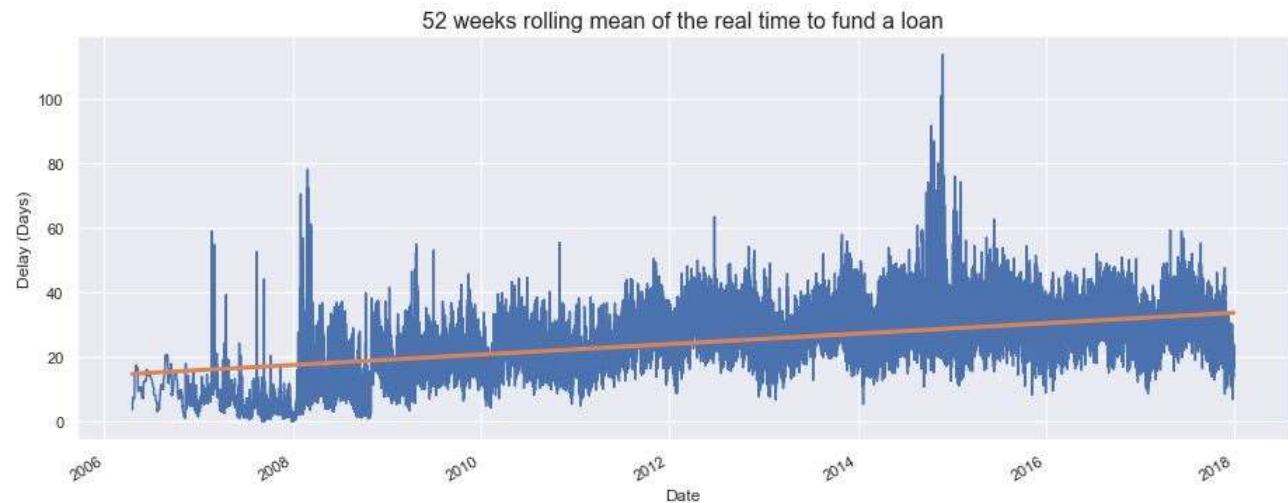
However, this test encompasses negative values, so when looking at the ECDF and comparing it with the theoretical CDF of a normal distribution, the resemblance is a bit more clear.

**ECDF:**

But even more important to this case is to look at the general representation of the real time to fund a loan over time:



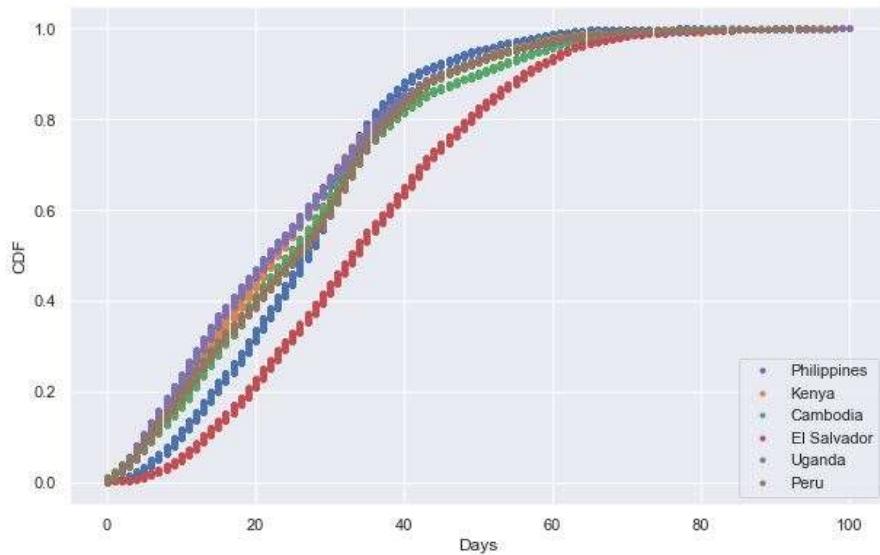
As we can see, there is a tendency to increase the fund delay over time. Looking more closely at the **52-week trend**, we notice that indeed the delays are increasing, on average, over time.



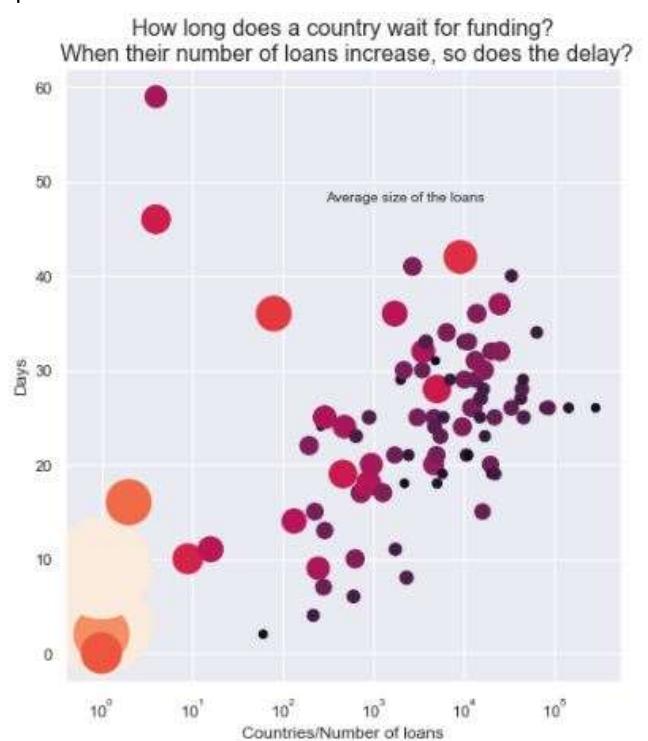
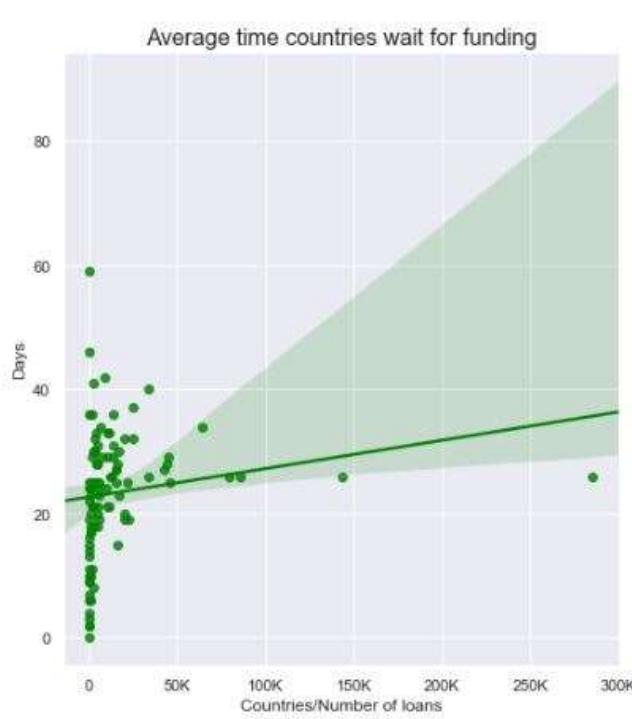
What may cause delays, then?

## 2.2) Do countries have influence regarding delays?

Half of the loans go to 6 countries (Philippines, Kenya, Peru, Cambodia, El Salvador, and Uganda), 77% to 20, 87% to 30. When looking specifically at the 6 countries PDFs and CDFs, we clearly distinguish El Salvador (takes 5% of the loans) from the other ones.



When a loan is funded (96% of the data), how long does it take to get that funding, on average, in each country? Does the number of loans have any impact?



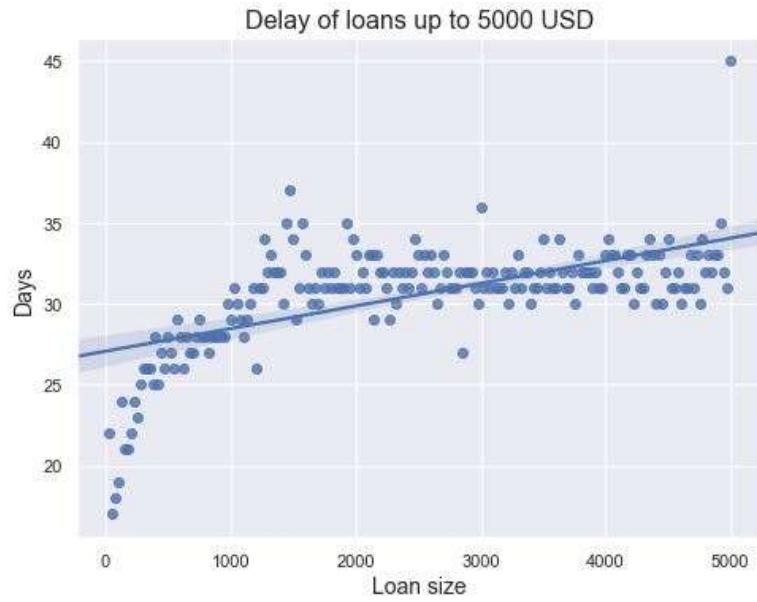
For countries with a few number of loans, the variation of average time to get the funding is significant. As soon as the number of loans starts to increase, the average number of days to get the funding rapidly get close to 30, where most of them reside. It seems to be that as the number of loans increases in a country, with all the remaining variables equal, so does the delay for funding a loan. There is a positive correlation between the two, but not significant when assuming a 5% significance level.

Given this result, however, other factors such as the characteristics of each country loan could be more influential, but we will keep this in mind. When filtering the countries that take a bit more than the average time to fund, USA, Colombia, Paraguay, Armenia, El Salvador, Lebanon, Bolivia, and Rwanda seem to have more impact.

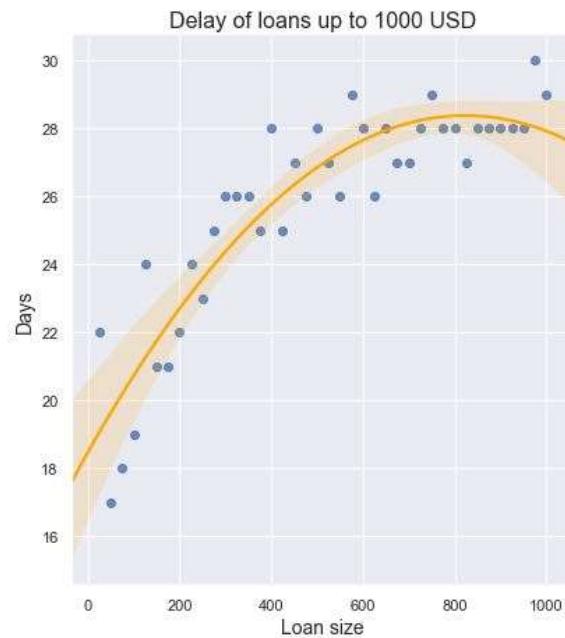
country	countries_nr_loans	avg_time_fund	size_loans
United States	9125	42.0	525.3
Colombia	33675	40.0	70.1
Paraguay	24787	37.0	217.7
Armenia	13951	36.0	164.5
El Salvador	64037	34.0	65.1
Lebanon	20083	32.0	136.4
Bolivia	25250	32.0	175.2
Rwanda	16773	30.0	173.8

## 2.3) Do bigger loan amounts take more time to fund, on average?

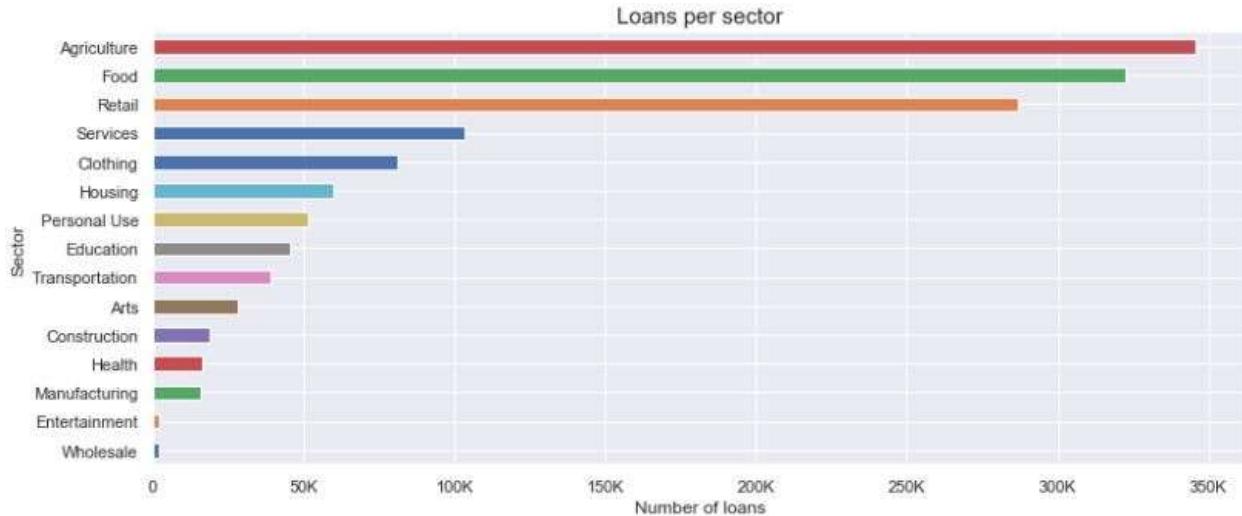
Globally, there is a significant negative correlation (-16%) between the size of the loan and the time it takes to get funded, contrary to intuition. But 99% of the data has loans inferior to 5000 USD, and in that case there is a significant positive correlation of 62%:



We notice that up until 1000USD there is an exponential increase in average delay. In fact, 77% of the loans are in this range. The correlation here is 84%.



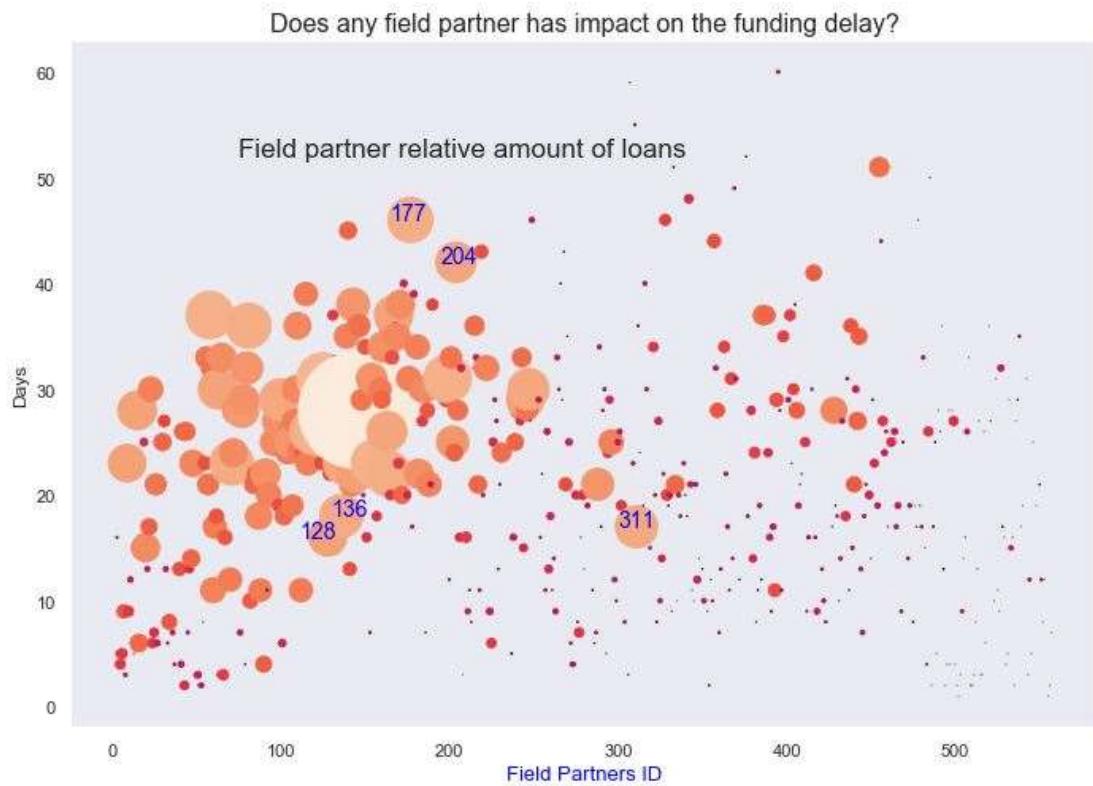
## 2.4) Do sectors of the loan influences delays?



67% of the loans go to 3 sectors: Agriculture (24%), Food (23%) and Retail (20%). With further analysis regarding each sector main statistics, we can notice that the main 6 sectors (which gather 85% of the loans) are the ones who also influence the most a higher delay in fundings, where Agriculture has the biggest impact.

sectors	sector_loans	perc_sector_loans	sector_mean	sector_std	sector_30	perc_30_total	perc_30_sector
Agriculture	329198	0.231973	28.277067	17.447268	149388	0.105268	0.453794
Food	311404	0.219434	25.949850	16.178154	123813	0.087246	0.397596
Retail	272475	0.192003	27.749236	16.898185	120001	0.084560	0.440411
Services	97819	0.068929	27.082254	18.387316	40988	0.028883	0.419019
Clothing	76856	0.054157	27.084248	18.371744	32693	0.023037	0.425380
Housing	54465	0.038379	31.714771	16.587375	29752	0.020965	0.546259
Education	44772	0.031549	25.440208	13.943170	18929	0.013339	0.422787
Transportation	36433	0.025673	28.307386	17.009264	16455	0.011595	0.451651
Personal Use	49595	0.034948	21.932150	14.215565	12993	0.009156	0.261982
Arts	27770	0.019568	21.314836	14.104928	7985	0.005627	0.287541
Construction	18379	0.012951	25.443169	16.998424	7122	0.005019	0.387507
Health	15812	0.011142	26.342525	16.110190	6710	0.004728	0.424361
Manufacturing	15840	0.011162	21.247096	12.713134	4279	0.003015	0.270139
Entertainment	1969	0.001387	25.169629	17.112305	695	0.000490	0.352971
Wholesale	2060	0.001452	22.014078	15.749968	640	0.000451	0.310680

## 2.5) Field partners

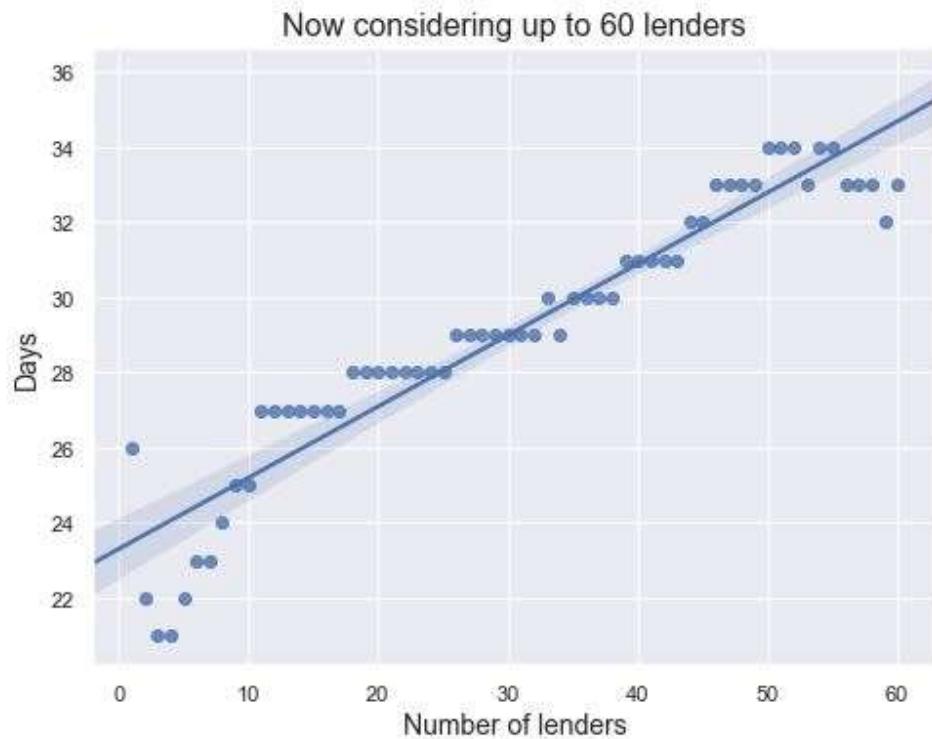


Of the 479 unique field partners, 11% of the loans are organized by field partner n° 145 and 90% by 120 field partners. There are not many field partners with a considerable size of loans that affect the overall average time loans take to get funded. Still, a few of them are relatively away from the concentrated 30-day delay. When filtering field partners that take more than 27 days on average to fund loans and with more than 1.5% of the loans, field partner with more loans, n° 145, is the one who causes more damage.

## 2.6) Number of lenders

Globally, as the number of lenders increased, the average number of days to fund a loan decreased.

However, 94% of the loans are composed by up to 60 lenders only. Here, there is a positive 95% correlation coefficient between the number of lenders and the time to fund a loan:

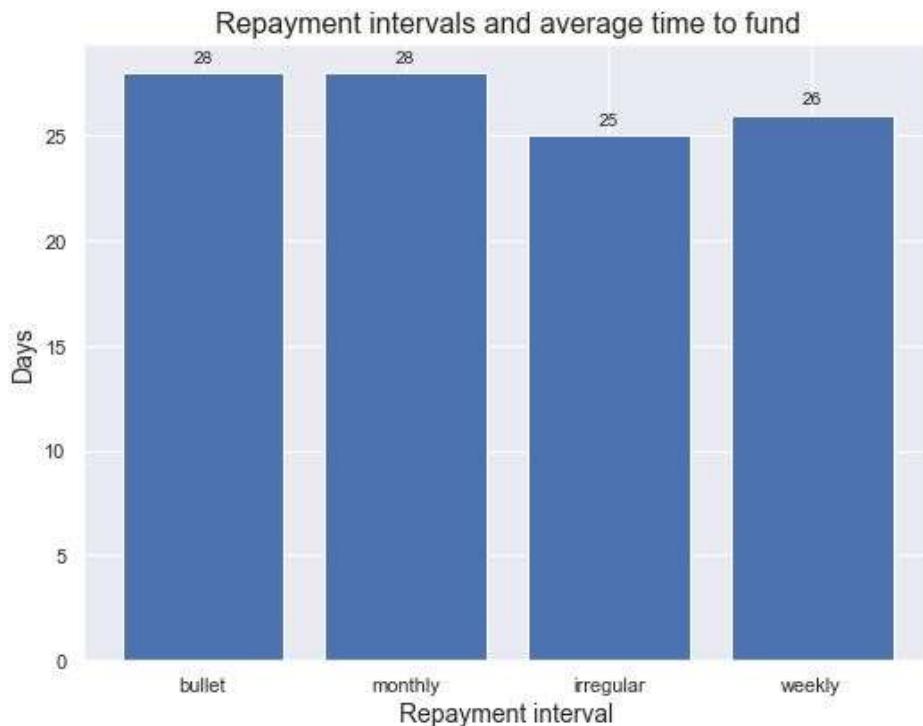


The data is presented in a seemingly organized disposition, which could suggest that the time to disburse and fund a loan is not subject to much uncertainty. In this sense, field partners may do pre-arranged commitments to fund the loans.

## 2.7) Does the repayment interval affect delays?

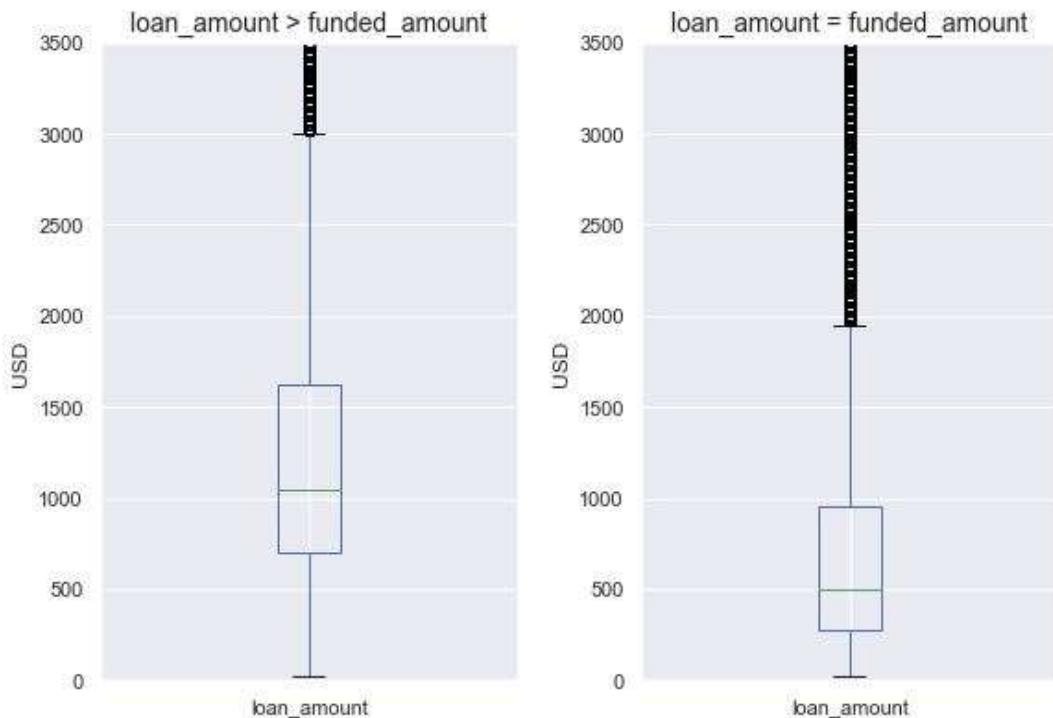
The main types of repayment intervals are “Monthly” (55%), “Irregular” (37%) and “Bullet”(8%):

There are slightly different statistics for each type. Dividing the data into two groups, one where the time to fund is below the median and the other above the median. Using a Chi-squared test, we conclude that there is a significant impact on the type of repayment interval. All in all, irregular repayments have a significantly better performance than monthly repayments, regarding delays.



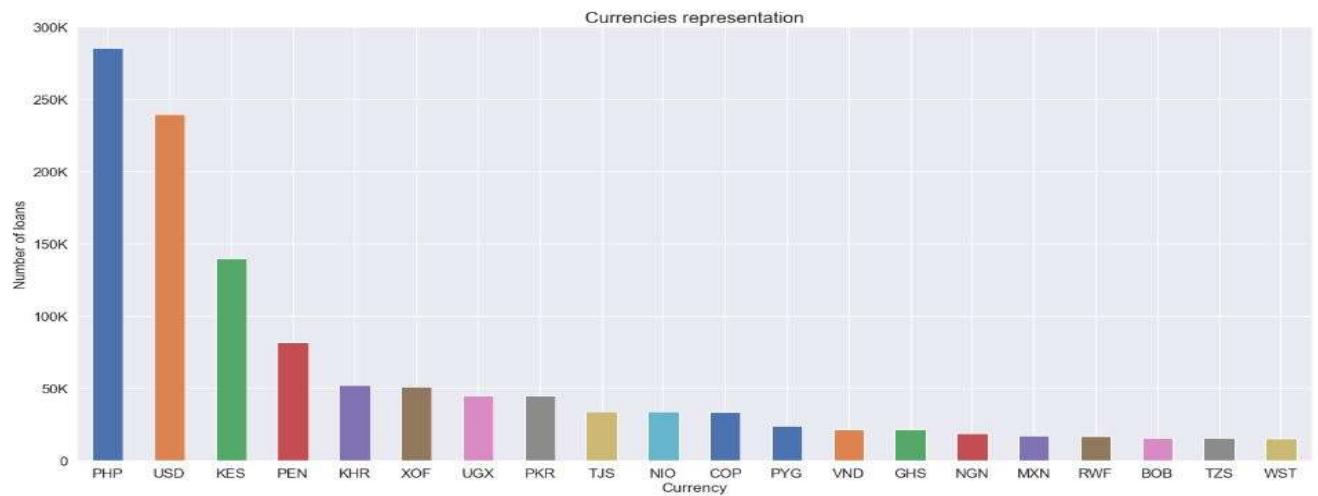
## 2.8) When the loan is not completely funded

They were bigger than usual, as expected.

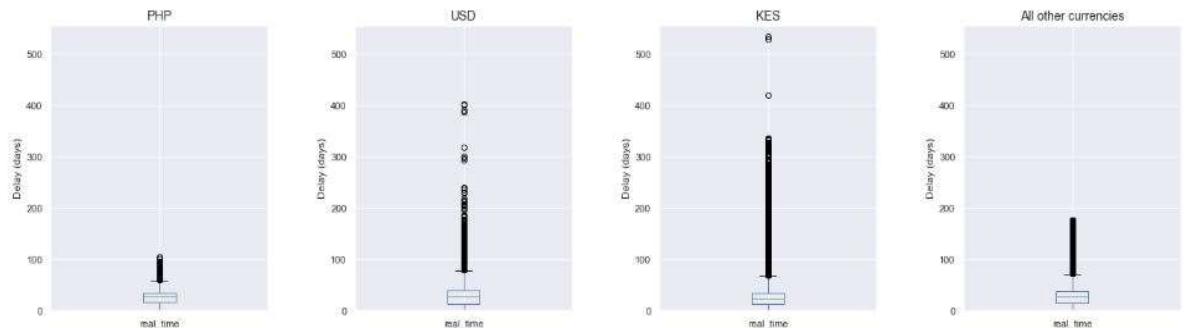


## 2.9) Currency of the loan

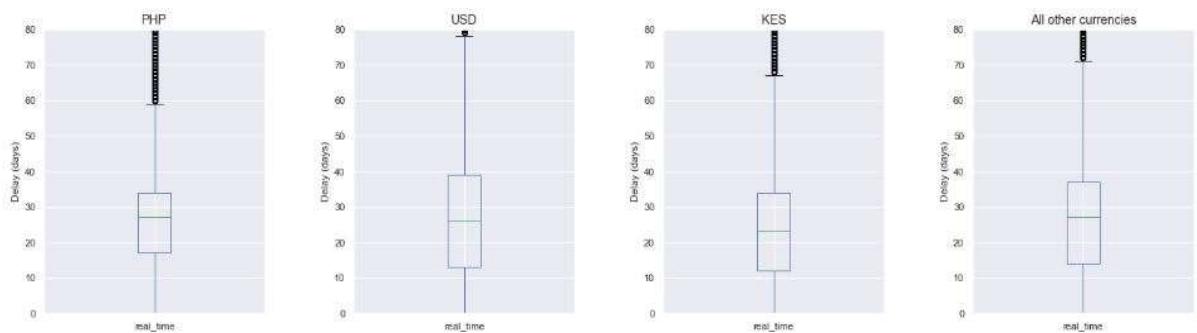
The 3 largest represented currencies are PHP, USD and KES, with almost half the loans. The 20 largest currencies represent 85% of the data:



Looking at the 3 major currencies, and comparing to the other ones, we see that USD is probably the currency that contributes more to delay in loans.

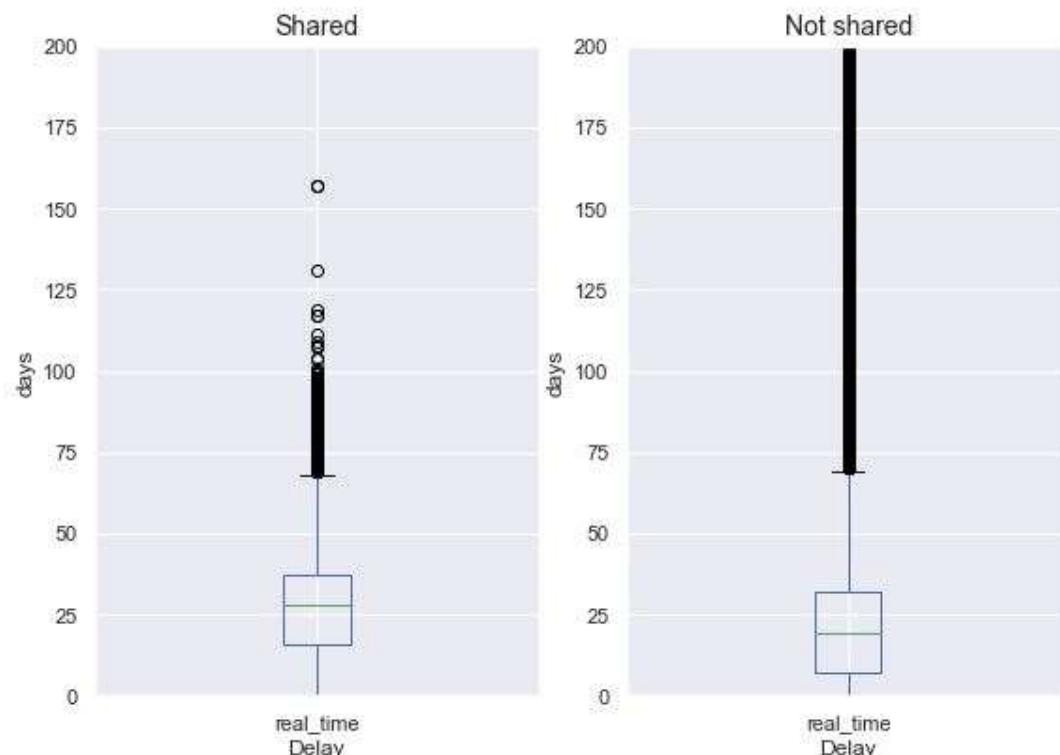


Zooming:



## 2.10) Currency policy

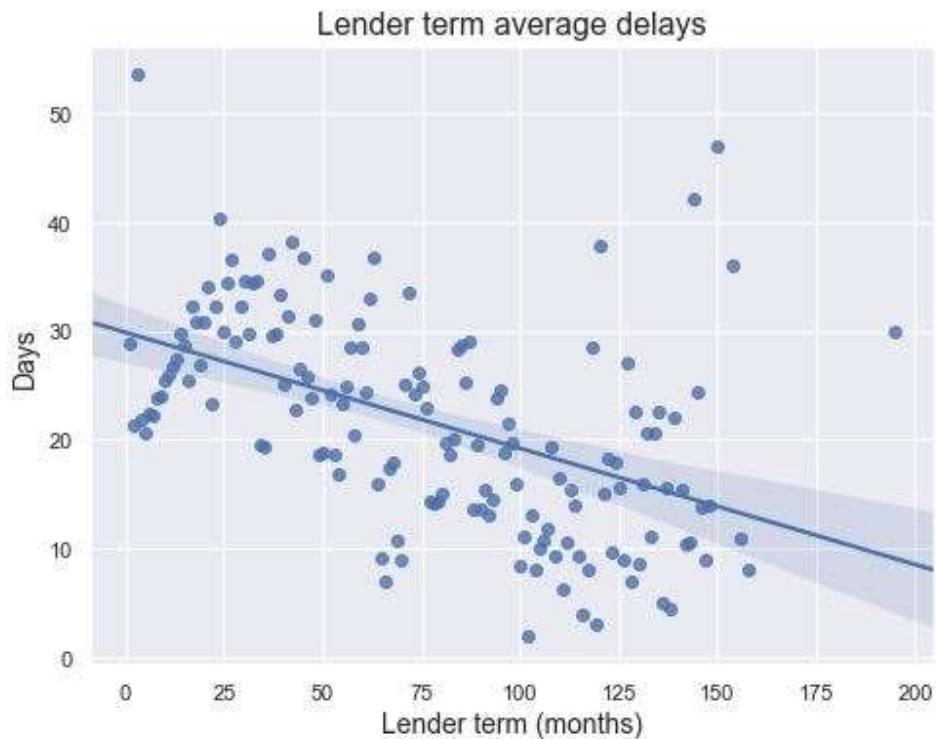
Currency policy can be of 2 kinds: shared or not shared. Kiva Field Partners have the choice to "opt-in" to Currency Risk Protection, which means they can choose to manage all Foreign Exchange Risk themselves, or alternatively they pass on the risk of currency devaluation over 20% to Kiva Lenders. 78% of the loans are shared and 22% not shared.



A not shared currency policy has the more problematic cases, as expected, as there are a lot more (long-time) outliers. However, overall, have a smaller median and general statistic properties than loans with a currency shared policy.

## 2.11) Does lender term influence delays?

As expected, the longer the lender term, the stronger the possibility that there is a smaller delay, on average. The correlation is -0.49.



## 3) Modeling

### 3.1) Initial Feature Selection & Extraction

Before running different models, I made the following adjustments: I selected only certain columns of interest, and excluded all the rows where a computation of the real time of delay was not possible. These don't affect the previous EDA made.

I did not consider the following variables:

- ['loan\_id', 'funded\_amount', 'status', 'activity\_name', 'loan\_use', 'country\_code', 'town\_name', 'posted\_time', 'planned\_expiration\_time', 'disbursed\_time', 'funded\_time', 'borrower\_genders']

They are either not necessary or do not possess sufficient quality data in order to be considered. I also disregarded 61732 rows where there was missing data relative to the computation of the real time delay, which also did not affect the results overall.

**Countries:** I kept the main 50 countries, which represent 97% of the data. I then turned the remaining countries into one variable named 'country\_name\_other'.

**Partners:** I kept the main 100 partners (includes partner\_id=1), which represents 85% of the data, and created 'partner\_id\_other' for any of the excluded partners.

**Currency:** I kept the main 25 currencies, which represent 90% of the data, and create a new feature named 'currency\_other' for any of the excluded currencies.

### 3.2) Problem statement

**Classification | "high" vs "low" delay:**

I treated the problem as a classification problem. The goal will be to classify a loan with either a **"high"** loan delay, that is, a loan that takes more than the median delay of all the loans in the data, or with a **"low"** loan delay, meaning that it will take less than the median, which in this case is equal to 26 days.

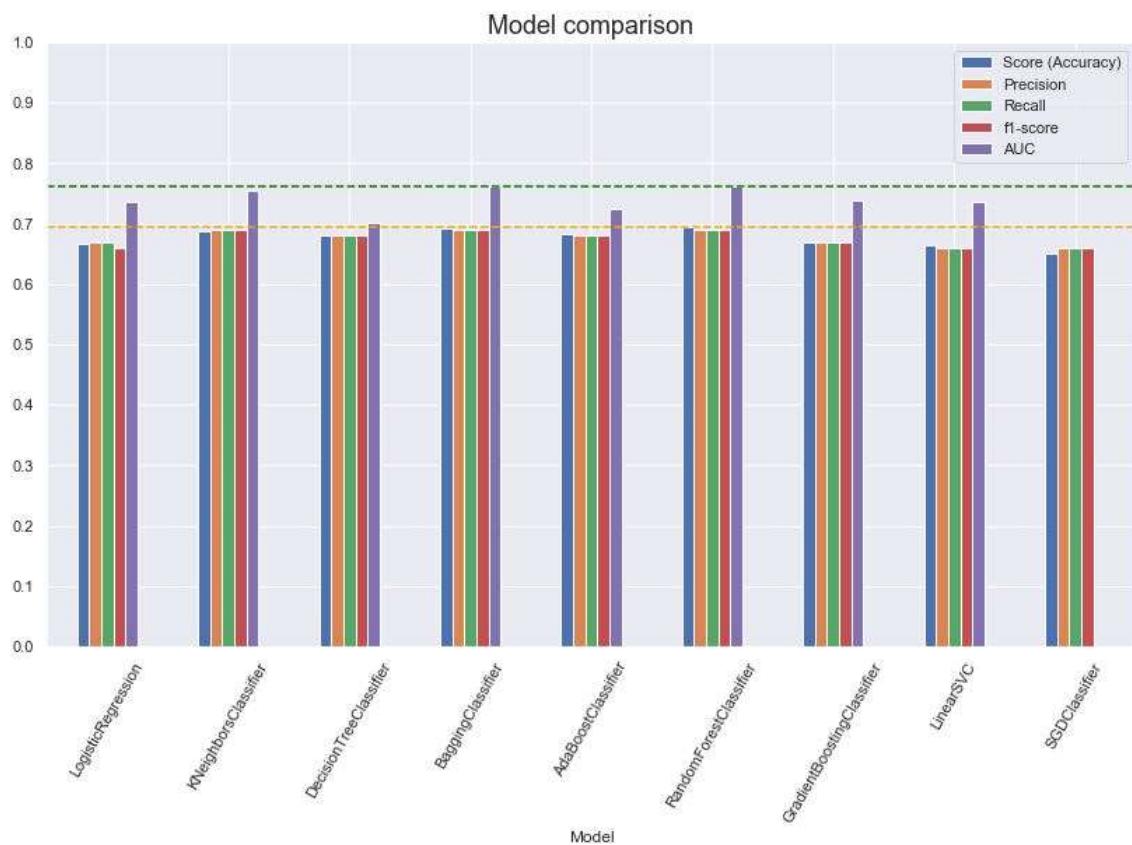
- "high" loan delay: a loan that takes more than 26 days to fund: 1
- "low" loan delay: a loan that takes less than 26 days to fund: 0

For linear models, I took the correlated variables out and scaled the numeric ones. I replaced the categorical variables with dummy variables both for linear and non-linear models.

### 3.3) Models initial results

Running 9 different models with their default parameters, I got the following results:

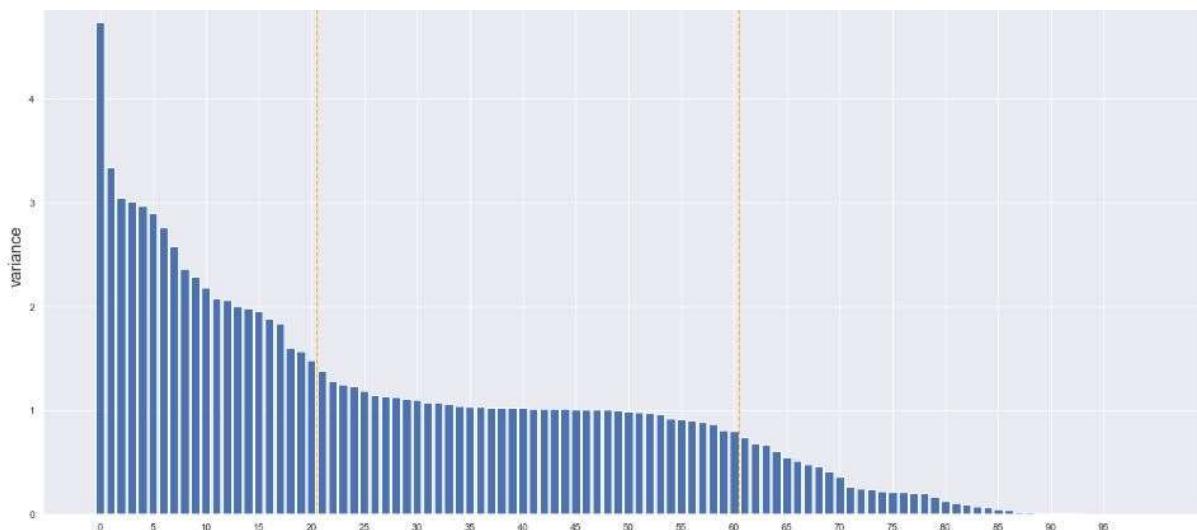
Model	Score (Accuracy)	Precision	Recall	f1-score	AUC score
LogisticRegression	0.666	0.67	0.67	0.66	0.737
KNeighborsClassifier	0.687	0.69	0.69	0.69	0.754
DecisionTreeClassifier	0.680	0.68	0.68	0.68	0.701
BaggingClassifier	0.693	0.69	0.69	0.69	0.761
AdaBoostClassifier	0.684	0.68	0.68	0.68	0.724
RandomForestClassifier	0.694	0.69	0.69	0.69	0.762
GradientBoostingClassifier	0.670	0.67	0.67	0.67	0.738
LinearSVC	0.664	0.66	0.66	0.66	0.737
SGDClassifier	0.650	0.66	0.66	0.66	0.



The **Random Forest classifier** got the highest scores on all categories.

## 4) Dimension reduction & Feature selection

Using Recursive Feature Elimination (RFE) with the Random Forest classifier as an estimator, I recursively fitted it up to until 85, 100, 110, 150 and the best performance was with 100 features, but still not as good as when we have all the 202 features. My computer memory could not handle EDA with all the variables, so I did it after this initial feature selection just to see the corresponding variances of each feature from that perspective. With 100 features, the explained variance started to decrease rapidly after approximately 60 features. However, we got better model performance with all the variables.



Hence, I opted to keep all the initial features for further tuning.

## 5) Hyperparameter tuning

The initial Random Forest classifier was defined as

- `rf = RandomForestClassifier(n_estimators=10, random_state=SEED)`

Running the following parameter grid:

- `params_rf = {'criterion': ['gini', 'entropy'], 'min_samples_leaf': [1, 2, 3, 4, 5, 10, 20], 'max_features': ['auto', 'log2'], 'max_depth': [None, 4], 'min_samples_split': [2, 3], n_estimators: [10, 100, 200, 400]}`

It only changed the default parameters in the case of the ‘min\_samples\_leaf’ to 4 and ‘n\_estimators’ to 400. The model performance improved:

Score: 0.7237761334489694

Confusion matrix:

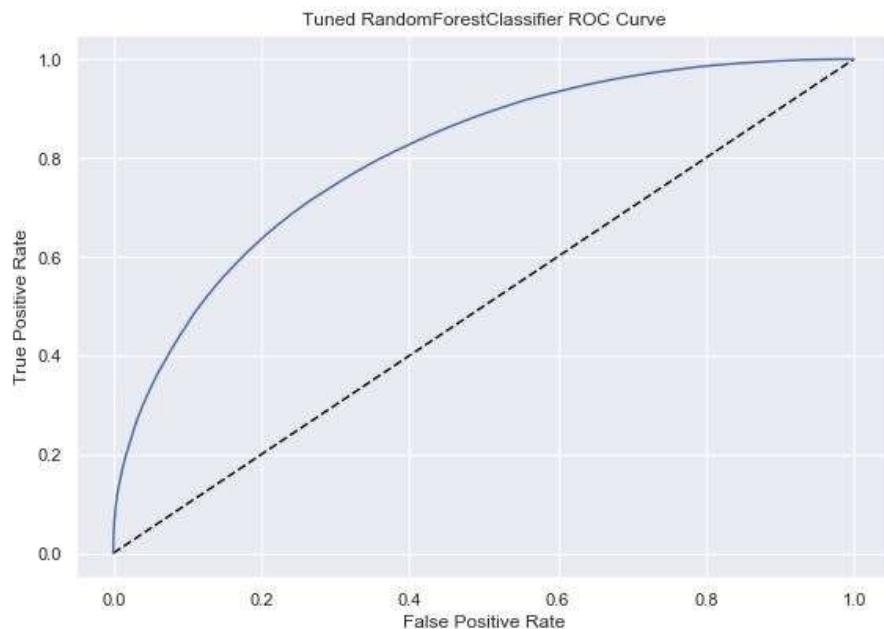
```
[[99106 37423]
 [37424 97012]]
```

Classification report:

	precision	recall	f1-score	support
0	0.73	0.73	0.73	136529
1	0.72	0.72	0.72	134436
micro avg	0.72	0.72	0.72	270965
macro avg	0.72	0.72	0.72	270965
weighted avg	0.72	0.72	0.72	270965

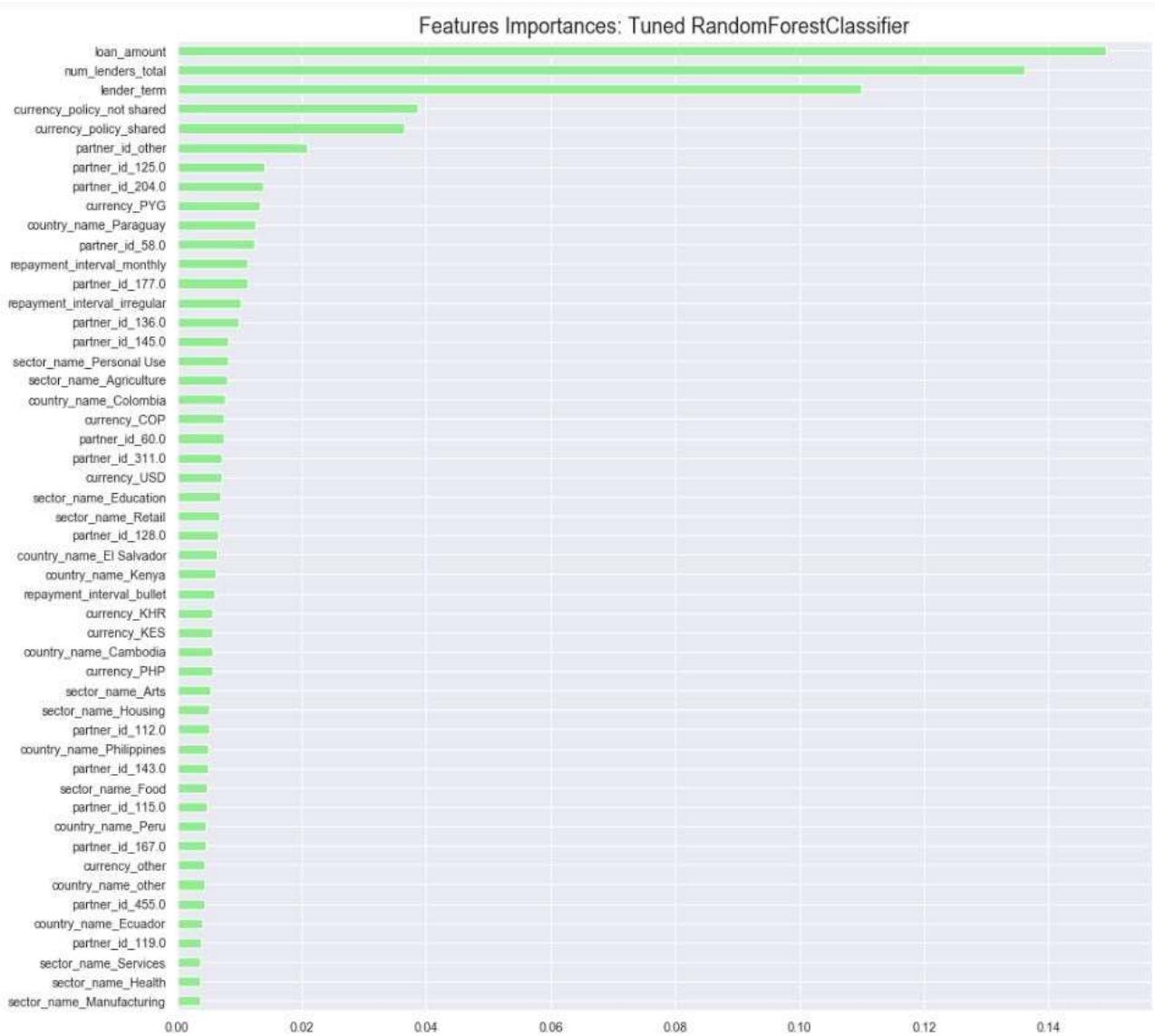
AUC: 0.8046943497441834

ROC curve:



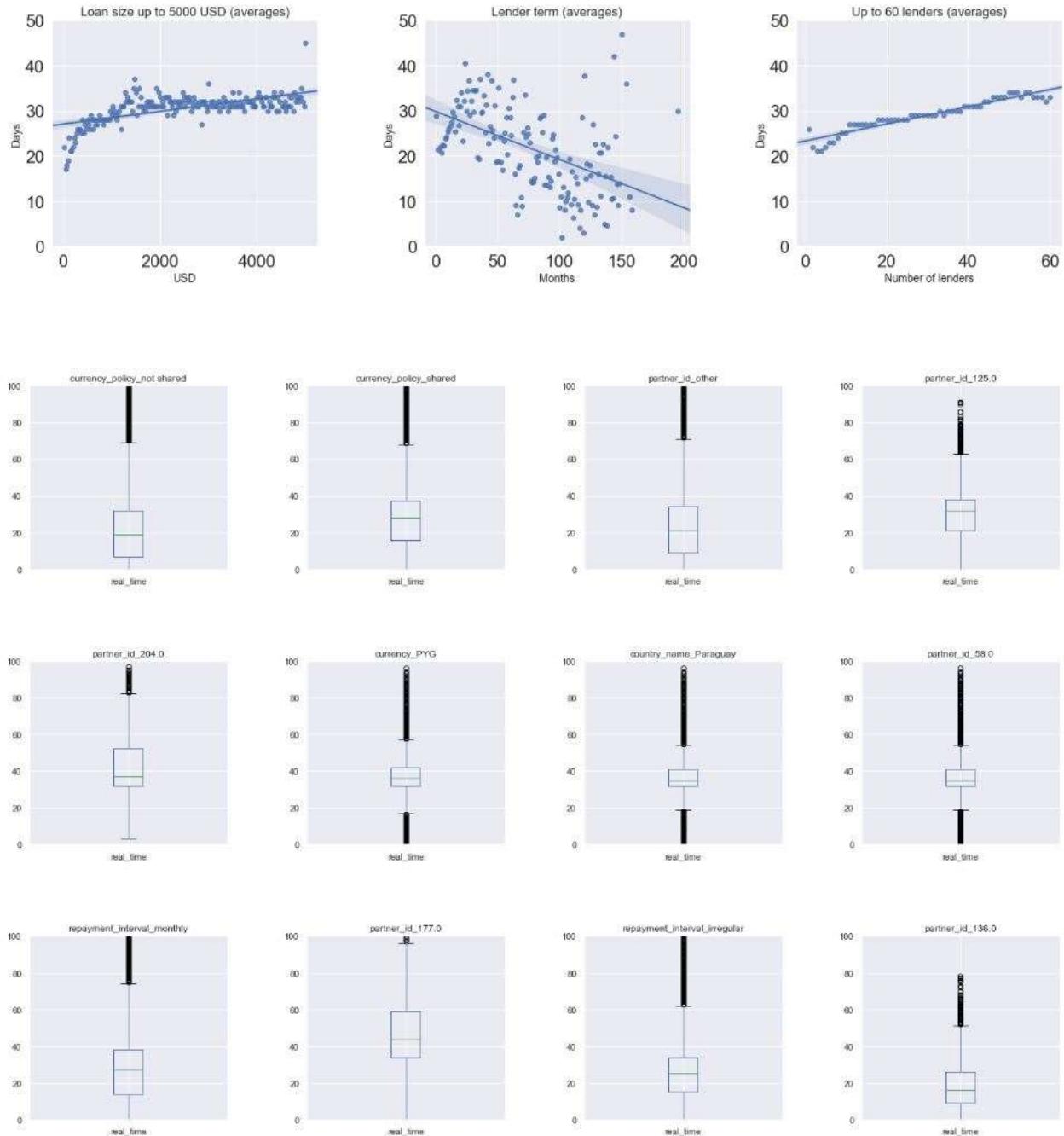
## 6) Conclusion

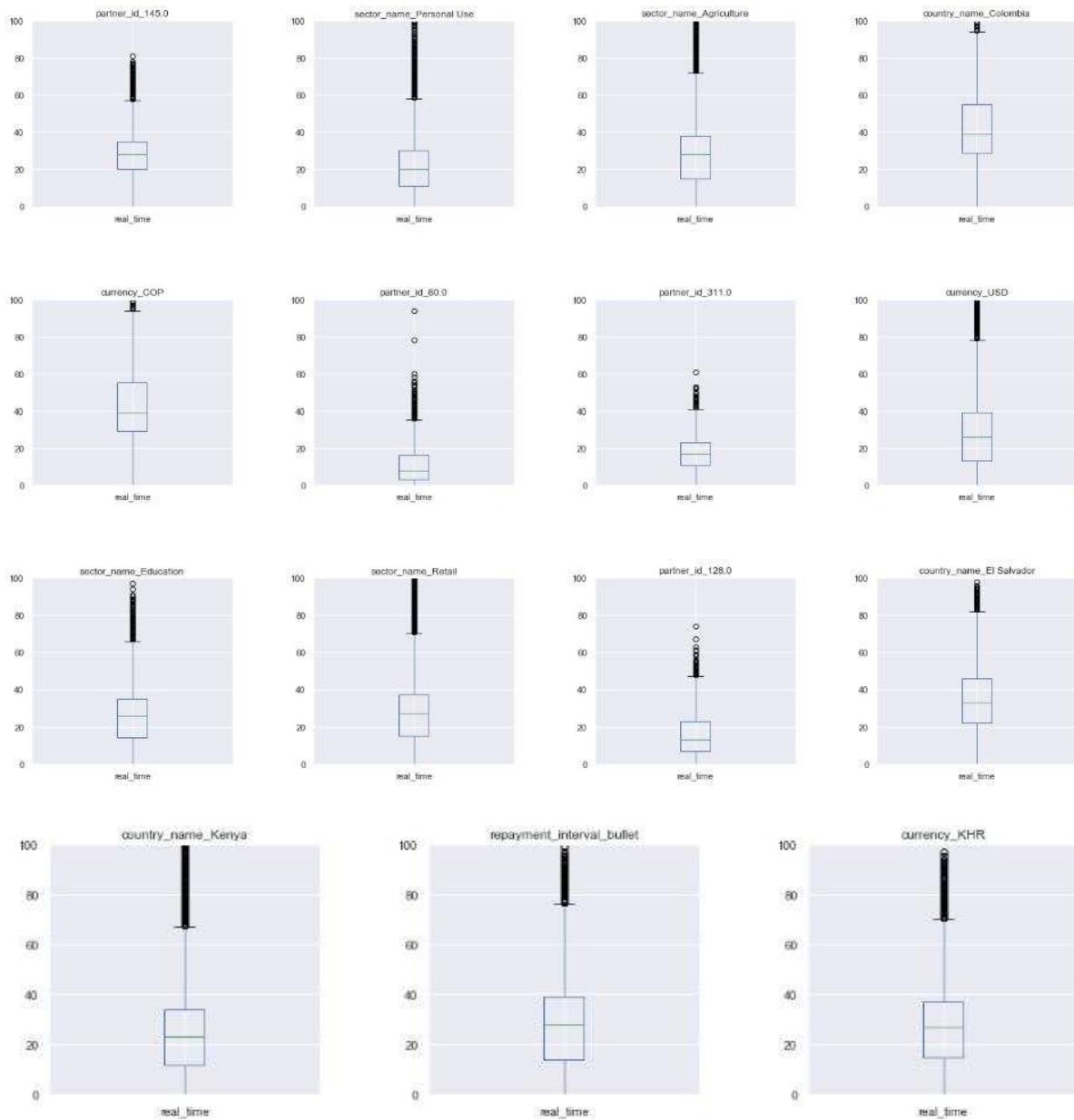
With the now tuned Random Forest classifier, I ran the feature importances of the model. Here are the 50 most important:



We see that the amount of the loan ('loan\_amount'), the total number of lenders ('num\_lenders\_total'), the lender term ('lender\_term') and the currency policy ('currency\_policy') are the most relevant features when predicting if a loan will be classified as one with a higher delay or not. Some of these are correlated, as seen before.

We can take a further look at the 30 most important features, according to our model. They represent a cumulative importance of 70%.

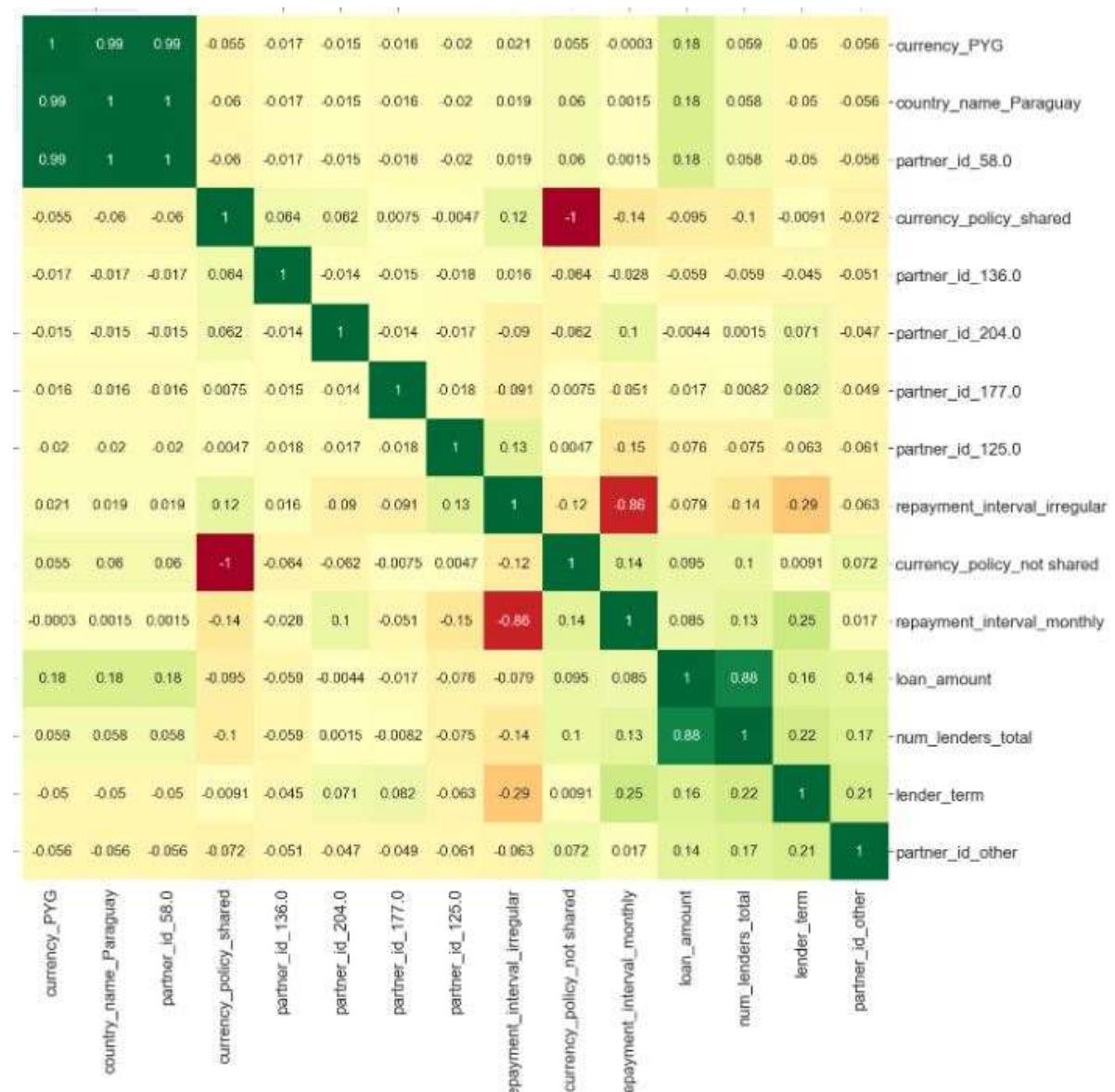




When relating to the previous EDA made, we see that some features are more helpful to predict when a loan will be delayed more than the desirable time, others to when it will have a smaller delay. For instance, we see that the the 2 most important ones (loan\_amount and num\_lenders\_total), which are positively correlated as we have also seen earlier, help the model in the sense that they contribute to higher delayed loans. Or, for instance, if it is a loan for Paraguay. Any loan with a high value in these features will be a motive of concern from the start, in the sense that they are good indicators of a higher delay. On the other hand, if, for instance, the loan has a shared currency policy, or it belongs to the 'Personal Use' sector, or the partner id is nº136, these are all good indicators that the model will be funded "on time".

We see that one of the most important features came up to be a product of our initial feature extraction, namely 'partner\_id\_other', which represents 15% of the data and 379 partners. Although each one has a small representation on the totality of loans, in this case it joined loans that, when seen as a whole, are not related with higher delays.

However, we should also look at the correlation between these important variables, before thinking if a loan with this or that feature may be problematic. These importances might hide the fact that a cause of a delay in a loan can have origin on one or few variables, but when other variables are highly correlated with those they will appear as "important" as well. Here is a clustermap with the relations between the 15 most important features:



For example, one of the most patent relations we see is that 99% of loans to Paraguay are in Guarani, their currency. In this case, the main cause is probably the sharp devaluation of the Guarani during the period of time of the dataset. As an example, the USD/PYG exchanged at 3972.000 in 2008, and by 2018 that rate was 5697.113. We can relate this to the fact that the feature 'currency\_policy\_not\_shared' contributed much more to delays than its counterpart, 'currency\_policy\_shared', which points to the fact that strong devaluations are difficult to deal with, specially after the 20% mark, as seen.

However, all in all most of the features are uncorrelated with each other, which is better in terms of interpretability ease.