

▼ Introduction: Home Credit Default Risk Competition

This notebook is intended for those who are new to machine learning competitions or want a gentle introduction to the problem. I purposely avoid jumping into complicated models or joining together lots of data in order to show the basics of how to get started in machine learning! Any comments or suggestions are much appreciated.

In this notebook, we will take an initial look at the **Home Credit default risk** machine learning competition currently hosted on Kaggle. The objective of this competition is to use historical loan application data to predict whether or not an applicant will be able to repay a loan. This is a standard supervised classification task:

- **Supervised:** The labels are included in the training data and the goal is to train a model to learn to predict the labels from the features
- **Classification:** The label is a binary variable, 0 (will repay loan on time), 1 (will have difficulty repaying loan)

Data

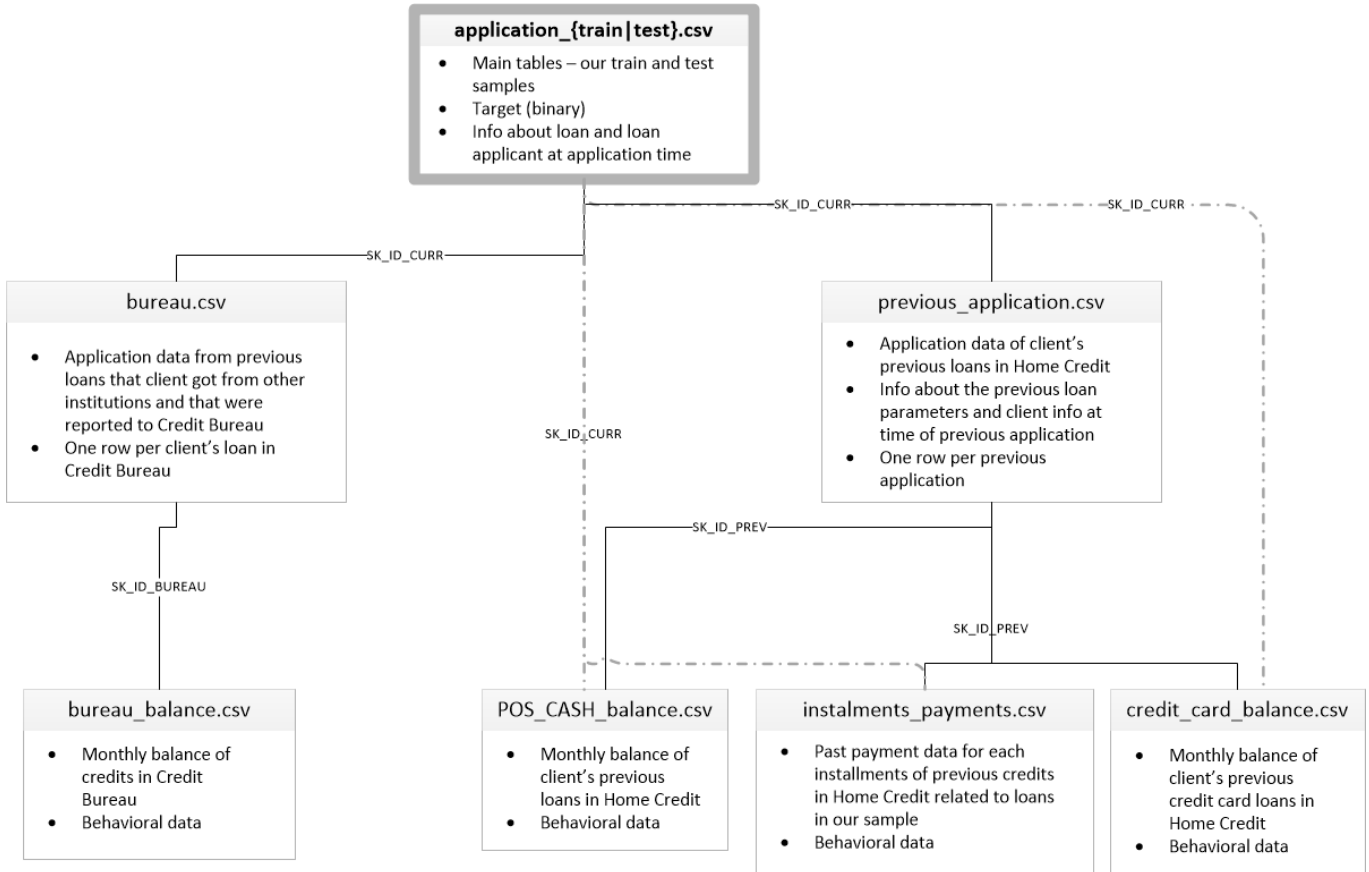
The data is provided by [Home Credit](#), a service dedicated to provided lines of credit (loans) to the unbanked population. Predicting whether or not a client will repay a loan or have difficulty is a critical business need, and Home Credit is hosting this competition on Kaggle to see what sort of models the machine learning community can develop to help them in this task.

There are 7 different sources of data:

- **application_train/application_test:** the main training and testing data with information about each loan application at Home Credit. Every loan has its own row and is identified by the feature `SK_ID_CURR`. The training application data comes with the `TARGET` indicating 0: the loan was repaid or 1: the loan was not repaid.
- **bureau:** data concerning client's previous credits from other financial institutions. Each previous credit has its own row in bureau, but one loan in the application data can have multiple previous credits.
- **bureau_balance:** monthly data about the previous credits in bureau. Each row is one month of a previous credit, and a single previous credit can have multiple rows, one for each month of the credit length.
- **previous_application:** previous applications for loans at Home Credit of clients who have loans in the application data. Each current loan in the application data can have multiple previous loans. Each previous application has one row and is identified by the feature `SK_ID_PREV`.
- **POS_CASH_BALANCE:** monthly data about previous point of sale or cash loans clients have had with Home Credit. Each row is one month of a previous point of sale or cash loan, and a single previous loan can have many rows.

- **credit_card_balance**: monthly data about previous credit cards clients have had with Home Credit. Each row is one month of a credit card balance, and a single credit card can have many rows.
- **installments_payment**: payment history for previous loans at Home Credit. There is one row for every made payment and one row for every missed payment.

This diagram shows how all of the data is related:



Moreover, we are provided with the definitions of all the columns (in `HomeCredit_columns_description.csv`) and an example of the expected submission file.

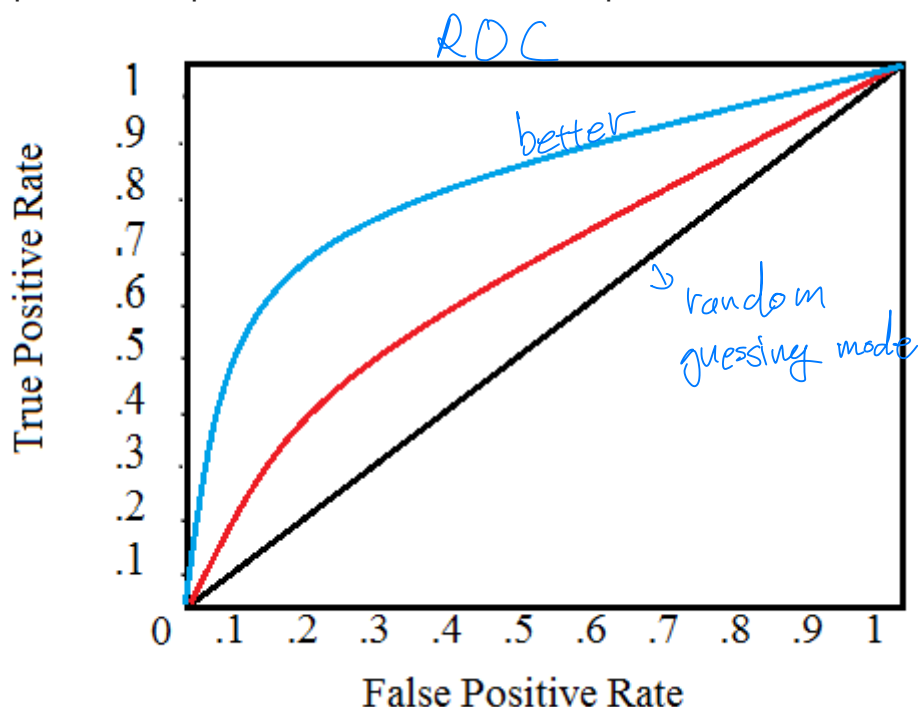
In this notebook, we will stick to using only the main application training and testing data. Although if we want to have any hope of seriously competing, we need to use all the data, for now we will stick to one file which should be more manageable. This will let us establish a baseline that we can then improve upon. With these projects, it's best to build up an understanding of the problem a little at a time rather than diving all the way in and getting completely lost!

Metric: ROC AUC

Once we have a grasp of the data (reading through the [column descriptions](#) helps immensely), we need to understand the metric by which our submission is judged. In this case, it is a common classification metric known as the [Receiver Operating Characteristic Area Under the Curve \(ROC AUC, also sometimes called AUROC\)](#).

The ROC AUC may sound intimidating, but it is relatively straightforward once you can get your head around the two individual concepts. The [Receiver Operating Characteristic \(ROC\) curve](#)

graphs the true positive rate versus the false positive rate:



A single line on the graph indicates the curve for a single model, and movement along a line indicates changing the threshold used for classifying a positive instance. The threshold starts at 0 in the upper right to and goes to 1 in the lower left. A curve that is to the left and above another curve indicates a better model. For example, the blue model is better than the red model, which is better than the black diagonal line which indicates a naive random guessing model.

The Area Under the Curve (AUC) explains itself by its name! It is simply the area under the ROC curve. (This is the integral of the curve.) This metric is between 0 and 1 with a better model scoring higher. A model that simply guesses at random will have an ROC AUC of 0.5.

When we measure a classifier according to the ROC AUC, we do not generation 0 or 1 predictions, but rather a probability between 0 and 1. This may be confusing because we usually like to think in terms of accuracy, but when we get into problems with imbalanced classes (we will see this is the case), accuracy is not the best metric. For example, if I wanted to build a model that could detect terrorists with 99.9999% accuracy, I would simply make a model that predicted every single person was not a terrorist. Clearly, this would not be effective (the recall would be zero) and we use more advanced metrics such as ROC AUC or the F1 score to more accurately reflect the performance of a classifier. A model with a high ROC AUC will also have a high accuracy, but the ROC AUC is a better representation of model performance.

Not that we know the background of the data we are using and the metric to maximize, let's get into exploring the data. In this notebook, as mentioned previously, we will stick to the main data sources and simple models which we can build upon in future work.

Follow-up Notebooks

For those looking to keep working on this problem, I have a series of follow-up notebooks:

ROC AUC

start-here-a-gentle-introduction

January 21, 2021

0.1 Imports

We are using a typical data science stack: numpy, pandas, sklearn, matplotlib.

```
[8]: # numpy and pandas for data manipulation
import numpy as np
import pandas as pd

# sklearn preprocessing for dealing with categorical variables
from sklearn.preprocessing import LabelEncoder

# File system manangement
import os

# Suppress warnings
import warnings
warnings.filterwarnings('ignore')

# matplotlib and seaborn for plotting
import matplotlib.pyplot as plt
import seaborn as sns
```

```
[9]: from google.colab import drive
drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call `drive.mount("/content/drive", force_remount=True)`.

0.2 Read in Data

First, we can list all the available data files. There are a total of 9 files: 1 main file for training (with target) 1 main file for testing (without the target), 1 example submission file, and 6 other files containing additional information about each loan.

```
[10]: # List files available
print(os.listdir("/content/drive/MyDrive/kaggle/ Start Here: A Gentle_
→Introduction/home-credit-default-risk.zip (Unzipped Files)"))
```

```
['HomeCredit_columns_description.csv', 'application_test.csv',
'application_train.csv', 'bureau.csv', 'POS_CASH_balance.csv',
'bureau_balance.csv', 'sample_submission.csv', 'credit_card_balance.csv',
'previous_application.csv', 'installments_payments.csv']
```

```
[11]: # Training data
app_train = pd.read_csv('/content/drive/MyDrive/kaggle/ Start Here: A Gentle_
→Introduction/home-credit-default-risk.zip (Unzipped Files)/application_train.
→csv')
print('Training data shape: ', app_train.shape)
app_train.head()
```

Training data shape: (307511, 122)

```
[11]: SK_ID_CURR  TARGET  ... AMT_REQ_CREDIT_BUREAU_QRT  AMT_REQ_CREDIT_BUREAU_YEAR
0      100002      1    ...                      0.0                      1.0
1      100003      0    ...                      0.0                      0.0
2      100004      0    ...                      0.0                      0.0
3      100006      0    ...                      NaN                      NaN
4      100007      0    ...                      0.0                      0.0
```

[5 rows x 122 columns]

The training data has 307511 observations (each one a separate loan) and 122 features (variables) including the TARGET (the label we want to predict).

```
[12]: # Testing data features
app_test = pd.read_csv('/content/drive/MyDrive/kaggle/ Start Here: A Gentle_
→Introduction/home-credit-default-risk.zip (Unzipped Files)/application_test.
→csv')
print('Testing data shape: ', app_test.shape)
app_test.head()
```

Testing data shape: (48744, 121)

```
[12]: SK_ID_CURR  ... AMT_REQ_CREDIT_BUREAU_YEAR
0      100001  ...                      0.0
1      100005  ...                      3.0
2      100013  ...                      4.0
3      100028  ...                      3.0
4      100038  ...                      NaN
```

[5 rows x 121 columns]

The test set is considerably smaller and lacks a TARGET column.

1 Exploratory Data Analysis

Exploratory Data Analysis (EDA) is an open-ended process where we calculate statistics and make figures to find trends, anomalies, patterns, or relationships within the data. The goal of EDA is to

learn what our data can tell us. It generally starts out with a high level overview, then narrows in to specific areas as we find intriguing areas of the data. The findings may be interesting in their own right, or they can be used to inform our modeling choices, such as by helping us decide which features to use.

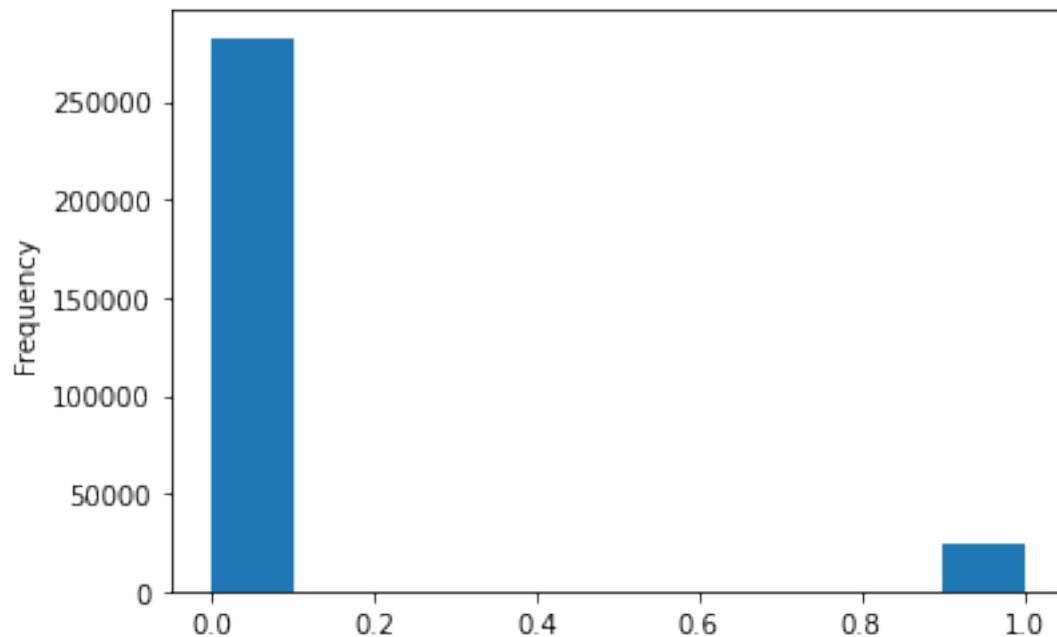
1.1 Examine the Distribution of the Target Column

The target is what we are asked to predict: either a 0 for the loan was repaid on time, or a 1 indicating the client had payment difficulties. We can first examine the number of loans falling into each category.

```
[13]: app_train['TARGET'].value_counts()
```

```
[13]: 0    282686  
      1     24825  
      Name: TARGET, dtype: int64
```

```
[14]: app_train['TARGET'].astype(int).plot.hist();
```



From this information, we see this is an *imbalanced class problem*. There are far more loans that were repaid on time than loans that were not repaid. Once we get into more sophisticated machine learning models, we can *weight the classes* by their representation in the data to reflect this imbalance.

1.2 Examine Missing Values

Next we can look at the number and percentage of missing values in each column.

```
[15]: # Function to calculate missing values by column# Funct
def missing_values_table(df):
    # Total missing values
    mis_val = df.isnull().sum()

    # Percentage of missing values
    mis_val_percent = 100 * df.isnull().sum() / len(df)

    # Make a table with the results
    mis_val_table = pd.concat([mis_val, mis_val_percent], axis=1)

    # Rename the columns
    mis_val_table_ren_columns = mis_val_table.rename(
        columns = {0 : 'Missing Values', 1 : '% of Total Values'})

    # Sort the table by percentage of missing descending
    mis_val_table_ren_columns = mis_val_table_ren_columns[
        mis_val_table_ren_columns.iloc[:,1] != 0].sort_values(
        '% of Total Values', ascending=False).round(1)

    # Print some summary information
    print ("Your selected dataframe has " + str(df.shape[1]) + " columns.\n"
          "There are " + str(mis_val_table_ren_columns.shape[0]) +
          " columns that have missing values.")

    # Return the dataframe with missing information
    return mis_val_table_ren_columns

[16]: # Missing values statistics
missing_values = missing_values_table(app_train)
missing_values.head(20)
```

Your selected dataframe has 122 columns.
There are 67 columns that have missing values.

```
[16]:
```

	Missing Values	% of Total Values
COMMONAREA_MEDI	214865	69.9
COMMONAREA_AVG	214865	69.9
COMMONAREA_MODE	214865	69.9
NONLIVINGAPARTMENTS_MEDI	213514	69.4
NONLIVINGAPARTMENTS_MODE	213514	69.4
NONLIVINGAPARTMENTS_AVG	213514	69.4
FONDKAPREMONT_MODE	210295	68.4
LIVINGAPARTMENTS_MODE	210199	68.4
LIVINGAPARTMENTS_MEDI	210199	68.4
LIVINGAPARTMENTS_AVG	210199	68.4
FLOORSMIN_MODE	208642	67.8
FLOORSMIN_MEDI	208642	67.8

FLOORSMIN_AVG	208642	67.8
YEARS_BUILD_MODE	204488	66.5
YEARS_BUILD_MEDI	204488	66.5
YEARS_BUILD_AVG	204488	66.5
OWN_CAR_AGE	202929	66.0
LANDAREA_AVG	182590	59.4
LANDAREA_MEDI	182590	59.4
LANDAREA_MODE	182590	59.4

When it comes time to build our machine learning models, we will have to fill in these missing values (known as imputation). In later work, we will use models such as XGBoost that can [handle missing values with no need for imputation](#). Another option would be to drop columns with a high percentage of missing values, although it is impossible to know ahead of time if these columns will be helpful to our model. Therefore, we will keep all of the columns for now.

1.3 Column Types

Let's look at the number of columns of each data type. `int64` and `float64` are numeric variables ([which can be either discrete or continuous](#)). `object` columns contain strings and are [categorical features](#).

```
[17]: # Number of each type of column
app_train.dtypes.value_counts()
```

```
[17]: float64    65
      int64     41
      object    16
      dtype: int64
```

Let's now look at the number of unique entries in each of the `object` (categorical) columns.

```
[18]: # Number of unique classes in each object column
app_train.select_dtypes('object').apply(pd.Series.nunique, axis = 0)
```

```
[18]: NAME_CONTRACT_TYPE    2
      CODE_GENDER        3
      FLAG_OWN_CAR       2
      FLAG_OWN_REALTY    2
      NAME_TYPE_SUITE    7
      NAME_INCOME_TYPE   8
      NAME_EDUCATION_TYPE 5
      NAME_FAMILY_STATUS 6
      NAME_HOUSING_TYPE   6
      OCCUPATION_TYPE    18
      WEEKDAY_APPR_PROCESS_START 7
      ORGANIZATION_TYPE  58
      FONDKAPREMONT_MODE  4
      HOUSETYPE_MODE     3
      WALLSMATERIAL_MODE  7
      EMERGENCYSTATE_MODE 2
      dtype: int64
```

```

ORGANIZATION_TYPE      58
FONDKAPREMONT_MODE      4
HOUSETYPE_MODE          3
WALLSMATERIAL_MODE      7
EMERGENCYSTATE_MODE     2
dtype: int64

```

Most of the categorical variables have a relatively small number of unique entries. We will need to find a way to deal with these categorical variables!

▼ Encoding Categorical Variables

Before we go any further, we need to deal with pesky categorical variables. A machine learning model unfortunately cannot deal with categorical variables (except for some models such as LightGBM). Therefore, we have to find a way to encode (represent) these variables as numbers before handing them off to the model. There are two main ways to carry out this process:

- Label encoding: assign each unique category in a categorical variable with an integer. No new columns are created. An example is shown below

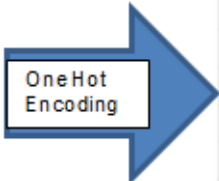
	occupation
0	programmer
1	data scientist
2	engineer
3	manager
4	ceo



	occupation
0	4
1	1
2	2
3	3
4	0

- One-hot encoding: create a new column for each unique category in a categorical variable. Each observation receives a 1 in the column for its corresponding category and a 0 in all other new columns.

	occupation
0	programmer
1	data scientist
2	engineer
3	manager
4	ceo



	occupation_ceo	occupation_data scientist	occupation_engineer	occupation_manager	occupation_programmer
0	0	0	0	0	1
1	0	1	0	0	0
2	0	0	1	0	0
3	0	0	0	1	0
4	1	0	0	0	0

The problem with label encoding is that it gives the categories an arbitrary ordering. The value assigned to each of the categories is random and does not reflect any inherent aspect of the category. In the example above, programmer receives a 4 and data scientist a 1, but if we did the same process again, the labels could be reversed or completely different. The actual assignment of the integers is arbitrary. Therefore, when we perform label encoding, the model might use the relative value of the feature (for example programmer = 4 and data scientist = 1) to assign weights which is not what we want. If we only have two unique values for a

categorical variable (such as Male/Female), then label encoding is fine, but for more than 2 unique categories, one-hot encoding is the safe option.

There is some debate about the relative merits of these approaches, and some models can deal with label encoded categorical variables with no issues. [Here is a good Stack Overflow discussion](#). I think (and this is just a personal opinion) for categorical variables with many classes, one-hot encoding is the safest approach because it does not impose arbitrary values to categories. The only downside to one-hot encoding is that the number of features (dimensions of the data) can explode with categorical variables with many categories. To deal with this, we can perform one-hot encoding followed by PCA or other [dimensionality reduction methods](#) to reduce the number of dimensions (while still trying to preserve information).

In this notebook, we will use Label Encoding for any categorical variables with only 2 categories and One-Hot Encoding for any categorical variables with more than 2 categories. This process may need to change as we get further into the project, but for now, we will see where this gets us. (We will also not use any dimensionality reduction in this notebook but will explore in future iterations).

▼ Label Encoding and One-Hot Encoding

Let's implement the policy described above: for any categorical variable (`dtype == object`) with 2 unique categories, we will use label encoding, and for any categorical variable with more than 2 unique categories, we will use one-hot encoding.

For label encoding, we use the Scikit-Learn `LabelEncoder` and for one-hot encoding, the pandas `get_dummies(df)` function.

```
# Create a label encoder object
le = LabelEncoder()
le_count = 0

# Iterate through the columns
for col in app_train:
    if app_train[col].dtype == 'object':
        # If 2 or fewer unique categories
        if len(list(app_train[col].unique())) <= 2:
            # Train on the training data
            le.fit(app_train[col])
            # Transform both training and testing data
            app_train[col] = le.transform(app_train[col])
            app_test[col] = le.transform(app_test[col])

            # Keep track of how many columns were label encoded
            le_count += 1

print('%d columns were label encoded.' % le_count)
```

3 columns were label encoded.

Most of the categorical variables have a relatively small number of unique entries. We will need to find a way to deal with these categorical variables!

1.3.1 Label Encoding and One-Hot Encoding

Let's implement the policy described above: for any categorical variable (`dtype == object`) with 2 unique categories, we will use label encoding, and for any categorical variable with more than 2 unique categories, we will use one-hot encoding.

For label encoding, we use the Scikit-Learn `LabelEncoder` and for one-hot encoding, the pandas `get_dummies(df)` function.

```
[19]: # Create a label encoder object
le = LabelEncoder()
le_count = 0

# Iterate through the columns
for col in app_train:
    if app_train[col].dtype == 'object':
        # If 2 or fewer unique categories
        if len(list(app_train[col].unique())) <= 2:
            # Train on the training data
            le.fit(app_train[col])
            # Transform both training and testing data
            app_train[col] = le.transform(app_train[col])
            app_test[col] = le.transform(app_test[col])

            # Keep track of how many columns were label encoded
            le_count += 1

print('%d columns were label encoded.' % le_count)
```

unique value ≤ 2

3 columns were label encoded.

```
[20]: # one-hot encoding of categorical variables
app_train = pd.get_dummies(app_train)
app_test = pd.get_dummies(app_test)

print('Training Features shape: ', app_train.shape)
print('Testing Features shape: ', app_test.shape)
```

Training Features shape: (307511, 243)) 문제
Testing Features shape: (48744, 239)

1.3.2 Aligning Training and Testing Data

There need to be the same features (columns) in both the training and testing data. One-hot encoding has created more columns in the training data because there were some categorical variables with categories not represented in the testing data. To remove the columns in the training data

that are not in the testing data, we need to align the dataframes. First we extract the target column from the training data (because this is not in the testing data but we need to keep this information). When we do the align, we must make sure to set `axis = 1` to align the dataframes based on the columns and not on the rows!

```
[21]: train_labels = app_train['TARGET']

# Align the training and testing data, keep only columns present in both
→ dataframes
app_train, app_test = app_train.align(app_test, join = 'inner', axis = 1)

# Add the target back in
app_train['TARGET'] = train_labels

print('Training Features shape: ', app_train.shape)
print('Testing Features shape: ', app_test.shape)
```

Training Features shape: (307511, 240)

Testing Features shape: (48744, 239)

The training and testing datasets now have the same features which is required for machine learning. The number of features has grown significantly due to one-hot encoding. At some point we probably will want to try dimensionality reduction (removing features that are not relevant) to reduce the size of the datasets.

ex) PCA

1.4 Back to Exploratory Data Analysis

1.4.1 Anomalies

One problem we always want to be on the lookout for when doing EDA is anomalies within the data. These may be due to mis-typed numbers, errors in measuring equipment, or they could be valid but extreme measurements. One way to support anomalies quantitatively is by looking at the statistics of a column using the `describe` method. The numbers in the `DAYS_BIRTH` column are negative because they are recorded relative to the current loan application. To see these stats in years, we can multiply by -1 and divide by the number of days in a year:

```
[22]: (app_train['DAYS_BIRTH'] / -365).describe()
```

```
[22]: count    307511.000000
      mean      43.936973
      std       11.956133
      min       20.517808
      25%       34.008219
      50%       43.150685
      75%       53.923288
      max       69.120548
      Name: DAYS_BIRTH, dtype: float64
```

Those ages look reasonable. There are no outliers for the age on either the high or low end. How about the days of employment?

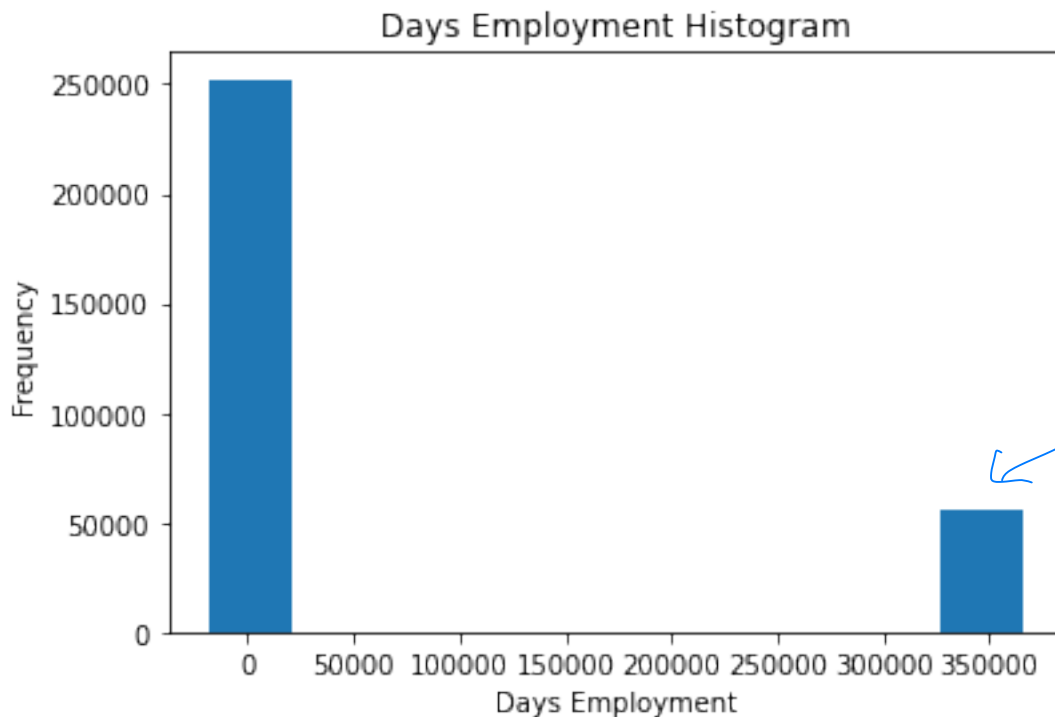
```
[23]: app_train['DAYS_EMPLOYED'].describe()
```

```
[23]: count    307511.000000
      mean      63815.045904
      std      141275.766519
      min      -17912.000000
      25%      -2760.000000
      50%      -1213.000000
      75%      -289.000000
      max      365243.000000
      Name: DAYS_EMPLOYED, dtype: float64
```

→ 이상치

That doesn't look right! The maximum value (besides being positive) is about 1000 years!

```
[24]: app_train['DAYS_EMPLOYED'].plot.hist(title = 'Days Employment Histogram');
      plt.xlabel('Days Employment');
```



Just out of curiosity, let's subset the anomalous clients and see if they tend to have higher or low rates of default than the rest of the clients.

```
[25]: anom = app_train[app_train['DAYS_EMPLOYED'] == 365243]
      non_anom = app_train[app_train['DAYS_EMPLOYED'] != 365243]
      print('The non-anomalies default on %0.2f%% of loans' % (100 *
        ↳non_anom['TARGET'].mean()))
      print('The anomalies default on %0.2f%% of loans' % (100 * anom['TARGET'].
        ↳mean()))
      print('There are %d anomalous days of employment' % len(anom))
```

The non-anomalies default on 8.66% of loans
The anomalies default on 5.40% of loans
There are 55374 anomalous days of employment

Well that is extremely interesting! It turns out that the anomalies have a lower rate of default.

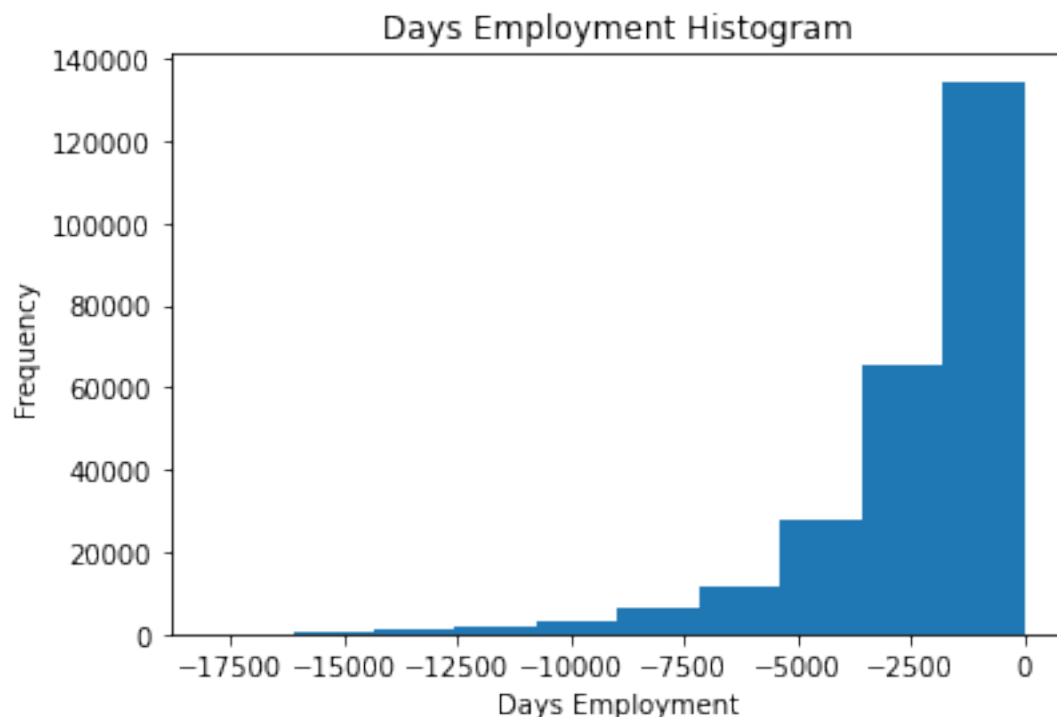
Handling the anomalies depends on the exact situation, with no set rules. One of the **safest approaches** is just to **set** the anomalies **to a missing value** and **then** have them **filled in** (using Imputation) before machine learning. In this case, since all the anomalies have the exact same value, we want to fill them in with the same value in case all of these loans share something in common. The anomalous values seem to have some importance, so we want to tell the machine learning model if we did in fact fill in these values. As a solution, we will fill in the anomalous values with not a number (`np.nan`) and then create a new boolean column indicating whether or not the value was anomalous.

```
[26]: # Create an anomalous flag column
app_train['DAYS_EMPLOYED_ANOM'] = app_train["DAYS_EMPLOYED"] == 365243

# Replace the anomalous values with nan
app_train['DAYS_EMPLOYED'].replace({365243: np.nan}, inplace = True)

app_train['DAYS_EMPLOYED'].plot.hist(title = 'Days Employment Histogram');
plt.xlabel('Days Employment');
```

365243 을 갖는 자료가 의미를 갖을 수 있으므로
column 제거 X



The distribution looks to be much more in line with what we would expect, and we also have created a new column to tell the model that these values were originally anomalous (because we

will have to fill in the nans with some value, probably the median of the column). The other columns with DAYS in the dataframe look to be about what we expect with no obvious outliers.

As an extremely important note, anything we do to the training data we also have to do to the testing data. Let's make sure to create the new column and fill in the existing column with np.nan in the testing data.

```
[27]: app_test['DAYS_EMPLOYED_ANOM'] = app_test["DAYS_EMPLOYED"] == 365243
      app_test["DAYS_EMPLOYED"].replace({365243: np.nan}, inplace = True)

      print('There are %d anomalies in the test data out of %d entries' % (
        →(app_test["DAYS_EMPLOYED_ANOM"].sum(), len(app_test)))
```

There are 9274 anomalies in the test data out of 48744 entries

1.4.2 Correlations

Now that we have dealt with the categorical variables and the outliers, let's continue with the EDA. One way to try and understand the data is by looking for correlations between the features and the target. We can calculate the Pearson correlation coefficient between every variable and the target using the .corr dataframe method.

The correlation coefficient is not the greatest method to represent "relevance" of a feature, but it does give us an idea of possible relationships within the data.

- .00-.19 "very weak"
- .20-.39 "weak"
- .40-.59 "moderate"
- .60-.79 "strong"
- .80-1.0 "very strong"

```
[28]: # Find correlations with the target and sort
      correlations = app_train.corr()['TARGET'].sort_values()

      # Display correlations
      print('Most Positive Correlations:\n', correlations.tail(15))
      print('\nMost Negative Correlations:\n', correlations.head(15))
```

Most Positive Correlations:

OCCUPATION_TYPE_Laborers	0.043019
FLAG_DOCUMENT_3	0.044346
REG_CITY_NOT_LIVE_CITY	0.044395
FLAG_EMP_PHONE	0.045982
NAME_EDUCATION_TYPE_Secondary / secondary special	0.049824
REG_CITY_NOT_WORK_CITY	0.050994
DAYS_ID_PUBLISH	0.051457
CODE_GENDER_M	0.054713
DAYS_LAST_PHONE_CHANGE	0.055218
NAME_INCOME_TYPE_Working	0.057481
REGION_RATING_CLIENT	0.058899


```

REGION_RATING_CLIENT_W_CITY      0.060893
DAYS_EMPLOYED                    0.074958
DAYS_BIRTH                       0.078239
TARGET                          1.000000
Name: TARGET, dtype: float64

```

Most Negative Correlations:

```

EXT_SOURCE_3                    -0.178919
EXT_SOURCE_2                    -0.160472
EXT_SOURCE_1                    -0.155317
NAME_EDUCATION_TYPE_Higher education -0.056593
CODE_GENDER_F                   -0.054704
NAME_INCOME_TYPE_Pensioner      -0.046209
DAYS_EMPLOYED_ANOM              -0.045987
ORGANIZATION_TYPE_XNA           -0.045987
FLOORSMAX_AVG                   -0.044003
FLOORSMAX_MEDI                  -0.043768
FLOORSMAX_MODE                  -0.043226
EMERGENCYSTATE_MODE_No          -0.042201
HOUSETYPE_MODE_block of flats   -0.040594
AMT_GOODS_PRICE                 -0.039645
REGION_POPULATION_RELATIVE      -0.037227
Name: TARGET, dtype: float64

```

Let's take a look at some of more significant correlations: the DAYS_BIRTH is the most positive correlation. (except for TARGET because the correlation of a variable with itself is always 1!) Looking at the documentation, DAYS_BIRTH is the age in days of the client at the time of the loan in negative days (for whatever reason!). The correlation is positive, but the value of this feature is actually negative, meaning that as the client gets older, they are less likely to default on their loan (ie the target == 0). That's a little confusing, so we will take the absolute value of the feature and then the correlation will be negative.

*Days-birth + minus
2TH plus corr $\approx 4\frac{1}{2}$*

1.4.3 Effect of Age on Repayment

```

[29]: # Find the correlation of the positive days since birth and target
app_train['DAYS_BIRTH'] = abs(app_train['DAYS_BIRTH'])
app_train['DAYS_BIRTH'].corr(app_train['TARGET'])

```

```
[29]: -0.07823930830982694
```

As the client gets older, there is a negative linear relationship with the target meaning that as clients get older, they tend to repay their loans on time more often.

Let's start looking at this variable. First, we can make a histogram of the age. We will put the x axis in years to make the plot a little more understandable.

```

[30]: # Set the style of plots
plt.style.use('fivethirtyeight')

# Plot the distribution of ages in years

```

```
plt.hist(app_train['DAYS_BIRTH'] / 365, edgecolor = 'k', bins = 25)
plt.title('Age of Client'); plt.xlabel('Age (years)'); plt.ylabel('Count');
```



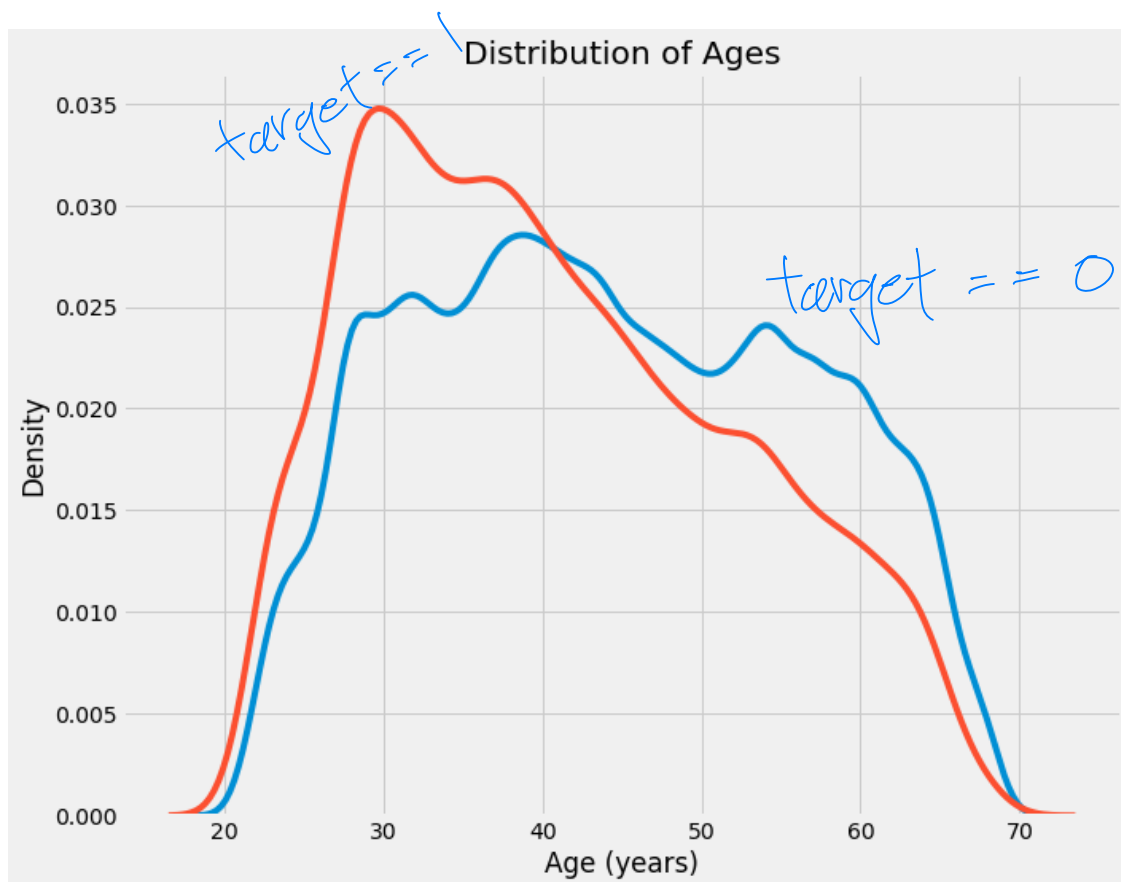
By itself, the distribution of age does not tell us much other than that there are no outliers as all the ages are reasonable. To visualize the effect of the age on the target, we will next make a **kernel density estimation plot** (KDE) colored by the value of the target. A **kernel density estimate plot shows the distribution of a single variable** and can be thought of as a smoothed histogram (it is created by computing a kernel, usually a Gaussian, at each data point and then averaging all the individual kernels to develop a single smooth curve). We will use the seaborn kdeplot for this graph.

```
[31]: plt.figure(figsize = (10, 8))

# KDE plot of loans that were repaid on time
sns.kdeplot(app_train.loc[app_train['TARGET'] == 0, 'DAYS_BIRTH'] / 365, label='target == 0')

# KDE plot of loans which were not repaid on time
sns.kdeplot(app_train.loc[app_train['TARGET'] == 1, 'DAYS_BIRTH'] / 365, label='target == 1')

# Labeling of plot
plt.xlabel('Age (years)'); plt.ylabel('Density'); plt.title('Distribution of Ages');
```



The target == 1 curve skews towards the younger end of the range. Although this is not a significant correlation (-0.07 correlation coefficient), this variable is likely going to be useful in a machine learning model because it does affect the target. Let's look at this relationship in another way: average failure to repay loans by age bracket.

To make this graph, first we cut the age category into bins of 5 years each. Then, for each bin, we calculate the average value of the target, which tells us the ratio of loans that were not repaid in each age category.

```
[32]: # Age information into a separate dataframe
age_data = app_train[['TARGET', 'DAYS_BIRTH']]
age_data['YEARS_BIRTH'] = age_data['DAYS_BIRTH'] / 365

# Bin the age data
age_data['YEARS_BINNED'] = pd.cut(age_data['YEARS_BIRTH'], bins = np.
    ↳ linspace(20, 70, num = 11))
age_data.head(10)
```

```
[32]:  TARGET  DAYS_BIRTH  YEARS_BIRTH  YEARS_BINNED
0         1      9461      25.920548  (25.0, 30.0]
1         0     16765      45.931507  (45.0, 50.0]
2         0     19046      52.180822  (50.0, 55.0]
3         0     19005      52.068493  (50.0, 55.0]
```

4	0	19932	54.608219	(50.0, 55.0]
5	0	16941	46.413699	(45.0, 50.0]
6	0	13778	37.747945	(35.0, 40.0]
7	0	18850	51.643836	(50.0, 55.0]
8	0	20099	55.065753	(55.0, 60.0]
9	0	14469	39.641096	(35.0, 40.0]

```
[33]: # Group by the bin and calculate averages
age_groups = age_data.groupby('YEARS_BINNED').mean()
age_groups
```

```
[33]:
```

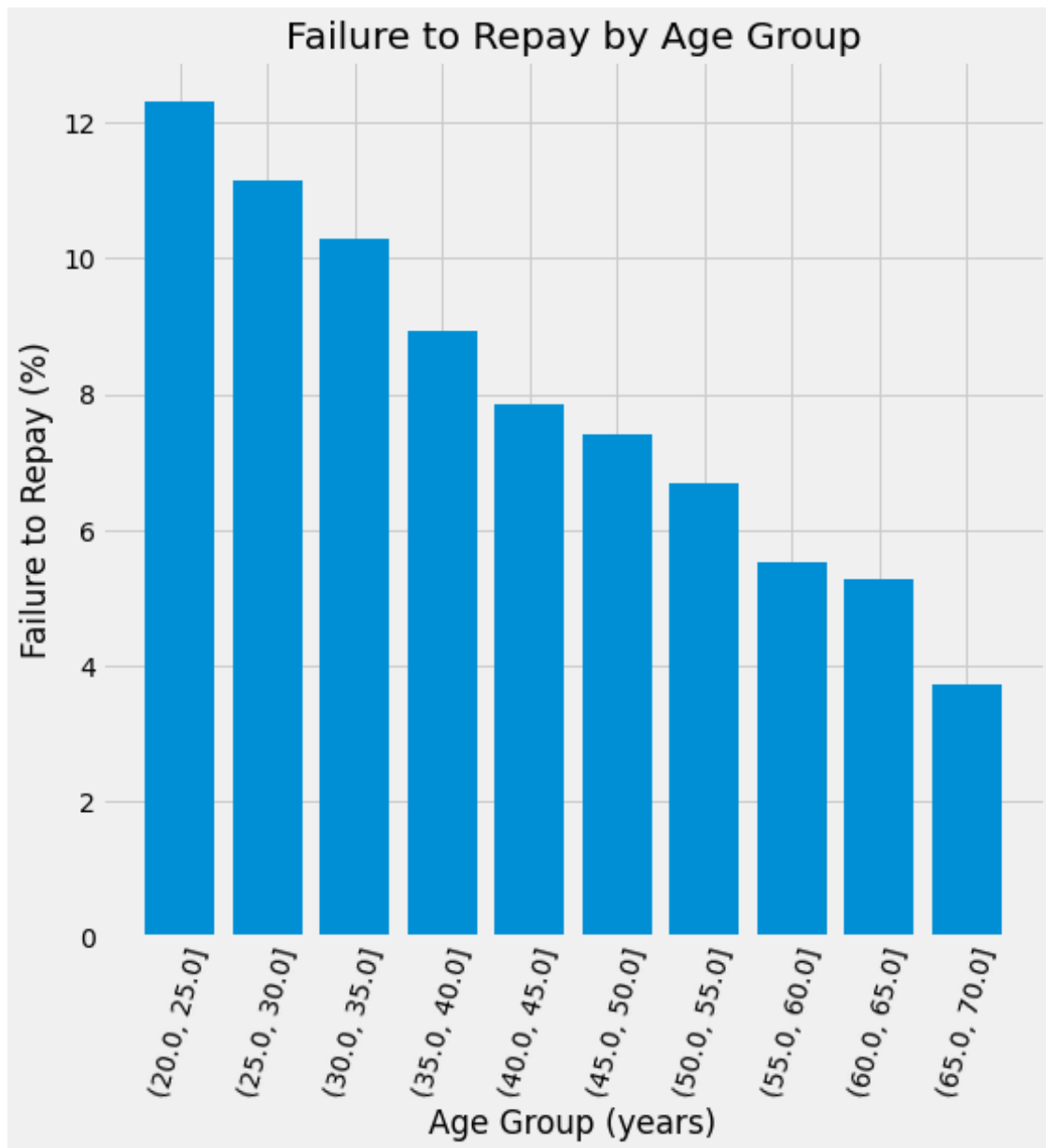
	TARGET	DAYS_BIRTH	YEARS_BIRTH
YEARS_BINNED			
(20.0, 25.0]	0.123036	8532.795625	23.377522
(25.0, 30.0]	0.111436	10155.219250	27.822518
(30.0, 35.0]	0.102814	11854.848377	32.479037
(35.0, 40.0]	0.089414	13707.908253	37.555913
(40.0, 45.0]	0.078491	15497.661233	42.459346
(45.0, 50.0]	0.074171	17323.900441	47.462741
(50.0, 55.0]	0.066968	19196.494791	52.593136
(55.0, 60.0]	0.055314	20984.262742	57.491131
(60.0, 65.0]	0.052737	22780.547460	62.412459
(65.0, 70.0]	0.037270	24292.614340	66.555108

Age ↗ Default ↘

```
[34]: plt.figure(figsize = (8, 8))

# Graph the age bins and the average of the target as a bar plot
plt.bar(age_groups.index.astype(str), 100 * age_groups['TARGET'])

# Plot labeling
plt.xticks(rotation = 75); plt.xlabel('Age Group (years)'); plt.ylabel('Failure_
→to Repay (%)')
plt.title('Failure to Repay by Age Group');
```



There is a clear trend: younger applicants are more likely to not repay the loan! The rate of failure to repay is above 10% for the youngest three age groups and below 5% for the oldest age group.

This is information that could be directly used by the bank: because younger clients are less likely to repay the loan, maybe they should be provided with more guidance or financial planning tips. This does not mean the bank should discriminate against younger clients, but it would be smart to take precautionary measures to help younger clients pay on time.

1.4.4 Exterior Sources

The 3 variables with the strongest negative correlations with the target are EXT_SOURCE_1, EXT_SOURCE_2, and EXT_SOURCE_3. According to the documentation, these features represent a "normalized score from external data source". I'm not sure what this exactly means, but it may be a cumulative sort of credit rating made using numerous sources of data.

Let's take a look at these variables.

First, we can show the correlations of the EXT_SOURCE features with the target and with each other.

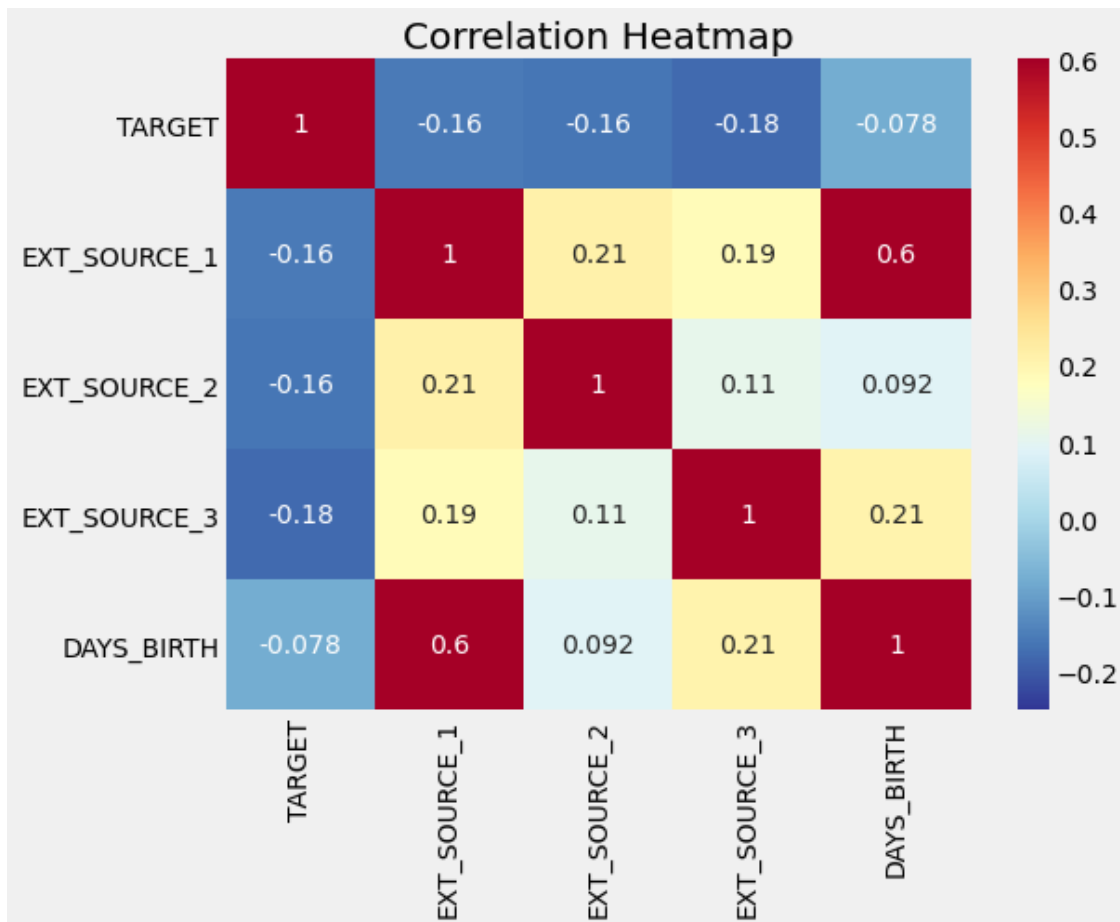
```
[35]: # Extract the EXT_SOURCE variables and show correlations
ext_data = app_train[['TARGET', 'EXT_SOURCE_1', 'EXT_SOURCE_2', 'EXT_SOURCE_3',
    → 'DAYS_BIRTH']]
ext_data_corrs = ext_data.corr()
ext_data_corrs
```

```
[35]:
```

	TARGET	EXT_SOURCE_1	EXT_SOURCE_2	EXT_SOURCE_3	DAYS_BIRTH
TARGET	1.000000	-0.155317	-0.160472	-0.178919	-0.078239
EXT_SOURCE_1	-0.155317	1.000000	0.213982	0.186846	0.600610
EXT_SOURCE_2	-0.160472	0.213982	1.000000	0.109167	0.091996
EXT_SOURCE_3	-0.178919	0.186846	0.109167	1.000000	0.205478
DAYS_BIRTH	-0.078239	0.600610	0.091996	0.205478	1.000000

```
[36]: plt.figure(figsize = (8, 6))

# Heatmap of correlations
sns.heatmap(ext_data_corrs, cmap = plt.cm.RdYlBu_r, vmin = -0.25, annot = True,
    → vmax = 0.6)
plt.title('Correlation Heatmap');
```



All three EXT_SOURCE features have negative correlations with the target, indicating that as the value of the EXT_SOURCE increases, the client is more likely to repay the loan. We can also see that DAYS_BIRTH is positively correlated with EXT_SOURCE_1 indicating that maybe one of the factors in this score is the client age.

Next we can look at the distribution of each of these features colored by the value of the target. This will let us visualize the effect of this variable on the target.

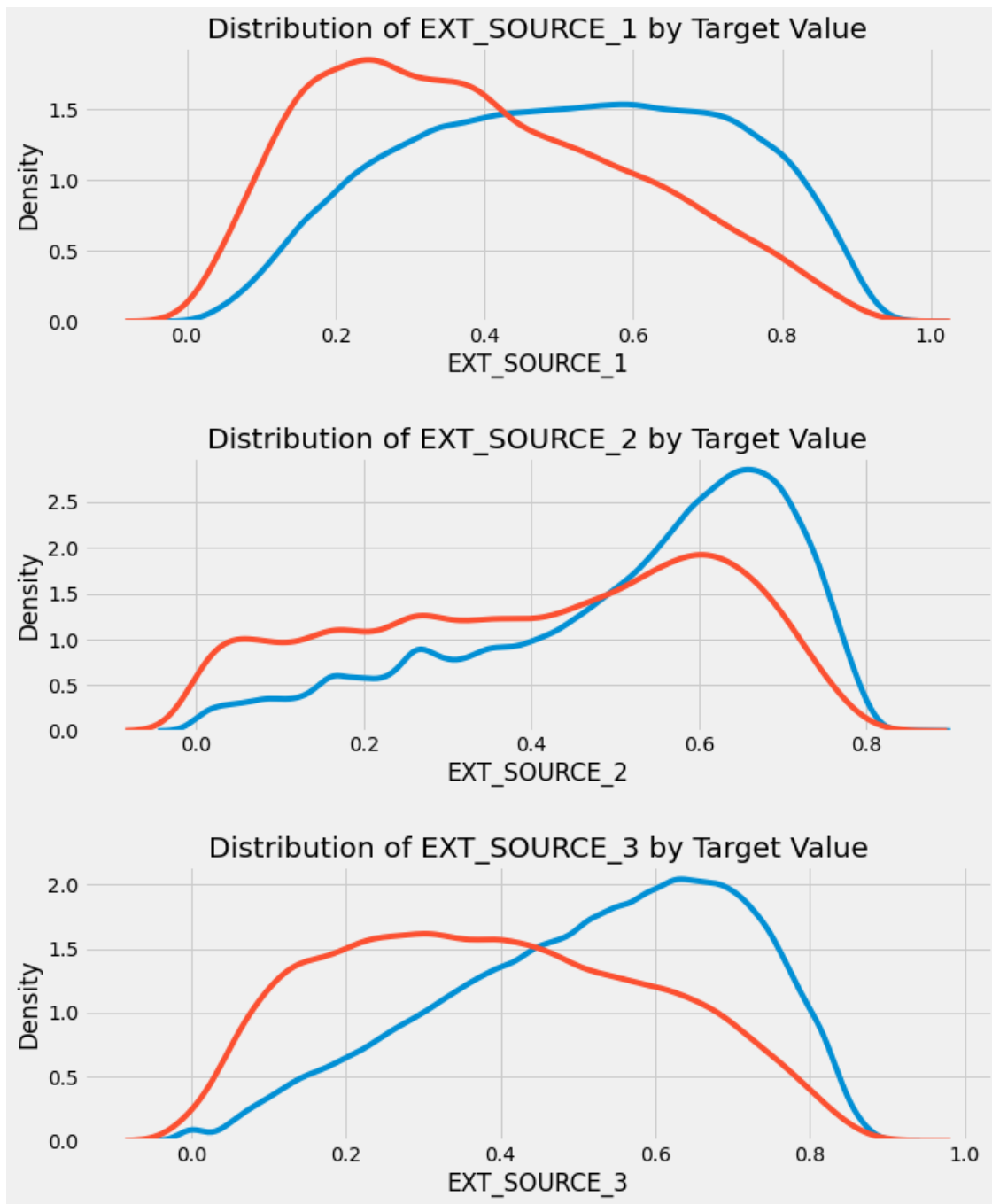
```
[37]: plt.figure(figsize = (10, 12))

# iterate through the sources
for i, source in enumerate(['EXT_SOURCE_1', 'EXT_SOURCE_2', 'EXT_SOURCE_3']):

    # create a new subplot for each source
    plt.subplot(3, 1, i + 1)
    # plot repaid loans
    sns.kdeplot(app_train.loc[app_train['TARGET'] == 0, source], label = 'target == 0')
    # plot loans that were not repaid
    sns.kdeplot(app_train.loc[app_train['TARGET'] == 1, source], label = 'target == 1')
```

```
# Label the plots
plt.title('Distribution of %s by Target Value' % source)
plt.xlabel('%s' % source); plt.ylabel('Density');

plt.tight_layout(h_pad = 2.5)
```



EXT_SOURCE_3 displays the greatest difference between the values of the target. We can clearly see that this feature has some relationship to the likelihood of an applicant to repay a loan. The relationship is not very strong (in fact they are all **considered very weak**, but these variables will still be useful for a machine learning model to predict whether or not an applicant will repay a loan on time.

1.5 Pairs Plot

As a final exploratory plot, we can make a pairs plot of the EXT_SOURCE variables and the DAYS_BIRTH variable. The **Pairs Plot** is a great exploration tool because it lets us see relationships between multiple pairs of variables as well as distributions of single variables. Here we are using the seaborn visualization library and the PairGrid function to create a Pairs Plot with scatterplots on the upper triangle, histograms on the diagonal, and 2D kernel density plots and correlation coefficients on the lower triangle.

If you don't understand this code, that's all right! Plotting in Python can be overly complex, and for anything beyond the simplest graphs, I usually find an existing implementation and adapt the code (don't repeat yourself)!

```
[38]: # Copy the data for plotting
plot_data = ext_data.drop(columns = ['DAYS_BIRTH']).copy()

# Add in the age of the client in years
plot_data['YEARS_BIRTH'] = age_data['YEARS_BIRTH']

# Drop na values and limit to first 100000 rows
plot_data = plot_data.dropna().loc[:100000, :]

# Function to calculate correlation coefficient between two columns
def corr_func(x, y, **kwargs):
    r = np.corrcoef(x, y)[0][1]
    ax = plt.gca()
    ax.annotate("r = {:.2f}".format(r),
                xy=(.2, .8), xycoords=ax.transAxes,
                size = 20)

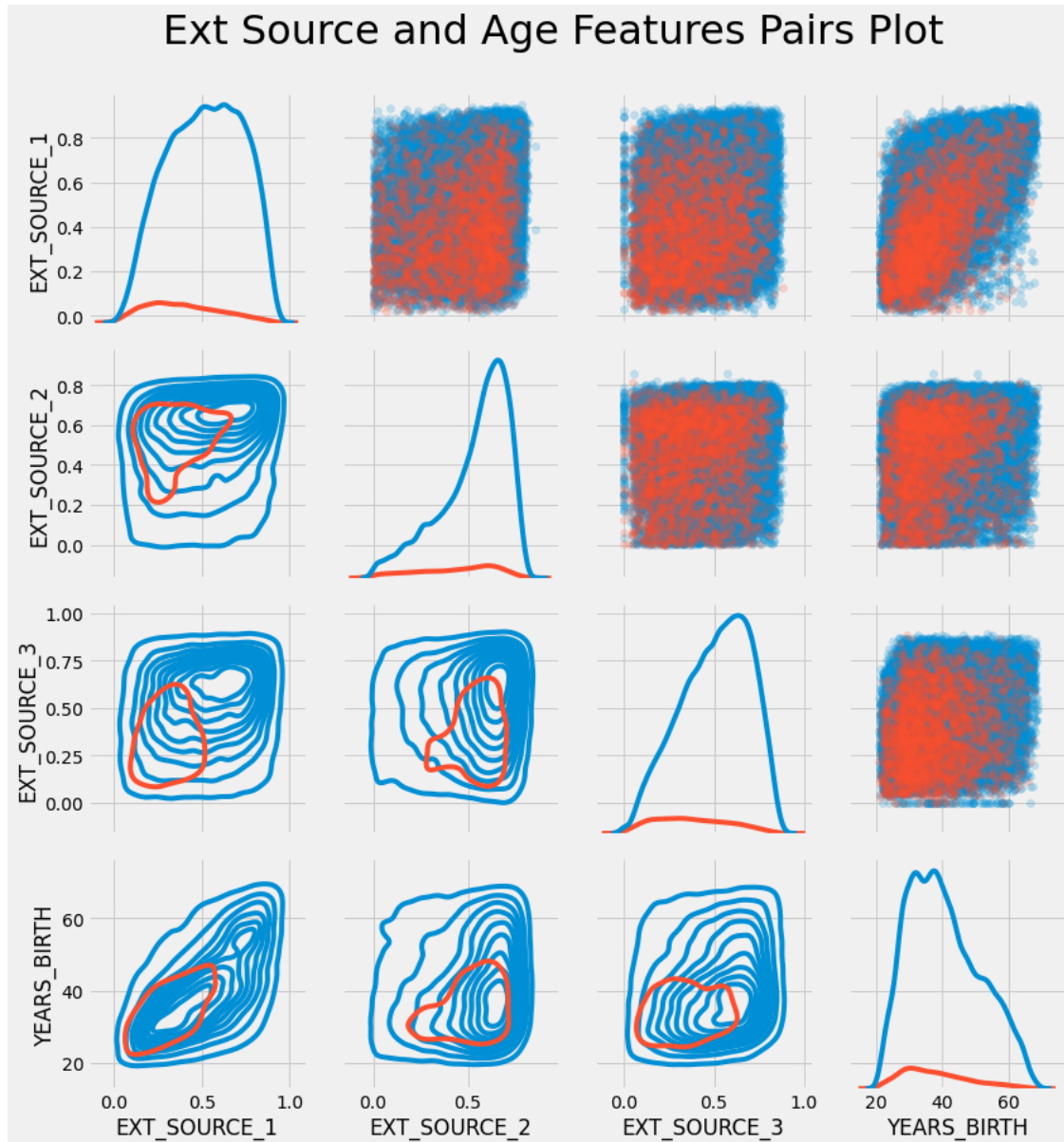
# Create the pairgrid object
grid = sns.PairGrid(data = plot_data, size = 3, diag_sharey=False,
                    hue = 'TARGET',
                    vars = [x for x in list(plot_data.columns) if x != 'TARGET'])

# Upper is a scatter plot
grid.map_upper(plt.scatter, alpha = 0.2)

# Diagonal is a histogram
grid.map_diag(sns.kdeplot)
```

```
# Bottom is density plot
grid.map_lower(sns.kdeplot, cmap = plt.cm.OrRd_r);

plt.suptitle('Ext Source and Age Features Pairs Plot', size = 32, y = 1.05);
```



In this plot, the red indicates loans that were not repaid and the blue are loans that are paid. We can see the different relationships within the data. There does appear to be a moderate positive linear relationship between the EXT_SOURCE_1 and the DAYS_BIRTH (or equivalently YEARS_BIRTH), indicating that this feature may take into account the age of the client.

2 Feature Engineering

Kaggle competitions are won by feature engineering: those who win are those who can create the most useful features out of the data. (This is true for the most part as the winning models, at least for structured data, all tend to be variants on [gradient boosting](#)). This represents one of the patterns in machine learning: [feature engineering has a greater return](#) on investment than model building and hyperparameter tuning. [This is a great article on the subject](#)). As Andrew Ng is fond of saying: "applied machine learning is basically feature engineering."

While choosing the right model and optimal settings are important, the model can only learn from the data it is given. Making sure this data is as relevant to the task as possible is the job of the data scientist (and maybe some [automated tools](#) to help us out).

Feature engineering refers to a general process and can involve both feature construction: [adding new features from the existing data](#), and [feature selection](#): choosing only the most important features or other methods of dimensionality reduction. There are many techniques we can use to both create features and select features.

We will do a lot of feature engineering when we start using the other data sources, but in this notebook we will try only two simple feature construction methods:

- Polynomial features
- Domain knowledge features

2.1 Polynomial Features

One simple feature construction method is called [polynomial features](#). In this method, we [make features that are powers of existing features as well as interaction terms between existing features](#). For example, we can create variables $EXT_SOURCE_1^2$ and $EXT_SOURCE_2^2$ and also variables such as $EXT_SOURCE_1 \times EXT_SOURCE_2$, $EXT_SOURCE_1 \times EXT_SOURCE_2^2$, $EXT_SOURCE_1^2 \times EXT_SOURCE_2^2$, and so on. These features that are a combination of multiple individual variables are called [interaction terms](#) ([https://en.wikipedia.org/wiki/Interaction_\(statistics\)](https://en.wikipedia.org/wiki/Interaction_(statistics))) because they capture the interactions between variables. In other words, while two variables by themselves may not have a strong influence on the target, [combining them together into a single interaction variable might show a relationship with the target](#). [Interaction terms are commonly used in statistical models](#) to capture the effects of multiple variables, but I do [not see them used as often](#) in machine learning. Nonetheless, we can try out a few to see if they might help our model to predict whether or not a client will repay a loan.

Jake VanderPlas writes about [polynomial features in his excellent book Python for Data Science](#) for those who want more information.

In the following code, we create polynomial features using the `EXT_SOURCE` variables and the `DAYS_BIRTH` variable. [Scikit-Learn has a useful class called PolynomialFeatures](#) that creates the polynomials and the interaction terms up to a specified degree. We can use a degree of 3 to see the results (when we are creating polynomial features, we want to [avoid using too high of a degree](#), both because the number of features scales exponentially with the degree, and because we can run into [problems with overfitting](#)).

```
[40]: # Make a new dataframe for polynomial features
poly_features = app_train[['EXT_SOURCE_1', 'EXT_SOURCE_2', 'EXT_SOURCE_3',
    → 'DAYS_BIRTH', 'TARGET']]
poly_features_test = app_test[['EXT_SOURCE_1', 'EXT_SOURCE_2', 'EXT_SOURCE_3',
    → 'DAYS_BIRTH']]
```

```

# imputer for handling missing values
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy = 'median')

poly_target = poly_features['TARGET']

poly_features = poly_features.drop(columns = ['TARGET'])

# Need to impute missing values
poly_features = imputer.fit_transform(poly_features)
poly_features_test = imputer.transform(poly_features_test)

from sklearn.preprocessing import PolynomialFeatures

# Create the polynomial object with specified degree
poly_transformer = PolynomialFeatures(degree = 3)

```

```

[41]: # Train the polynomial features
poly_transformer.fit(poly_features)

# Transform the features
poly_features = poly_transformer.transform(poly_features)
poly_features_test = poly_transformer.transform(poly_features_test)
print('Polynomial Features shape: ', poly_features.shape)

```

Polynomial Features shape: (307511, 35)

This creates a considerable number of new features. To get the names we have to use the polynomial features `get_feature_names` method.

```

[42]: poly_transformer.get_feature_names(input_features = ['EXT_SOURCE_1',
↳ 'EXT_SOURCE_2', 'EXT_SOURCE_3', 'DAYS_BIRTH'])[:15]

```

```

[42]: ['1',
      'EXT_SOURCE_1',
      'EXT_SOURCE_2',
      'EXT_SOURCE_3',
      'DAYS_BIRTH',
      'EXT_SOURCE_1^2',
      'EXT_SOURCE_1 EXT_SOURCE_2',
      'EXT_SOURCE_1 EXT_SOURCE_3',
      'EXT_SOURCE_1 DAYS_BIRTH',
      'EXT_SOURCE_2^2',
      'EXT_SOURCE_2 EXT_SOURCE_3',
      'EXT_SOURCE_2 DAYS_BIRTH',
      'EXT_SOURCE_3^2',
      'EXT_SOURCE_3 DAYS_BIRTH',
      'DAYS_BIRTH^2']

```

There are 35 features with individual features raised to powers up to degree 3 and interaction terms. Now, we can see whether any of these new features are correlated with the target.

```
[43]: # Create a dataframe of the features
poly_features = pd.DataFrame(poly_features,
                             columns = poly_transformer.
                             →get_feature_names(['EXT_SOURCE_1', 'EXT_SOURCE_2',
                             →'EXT_SOURCE_3', 'DAYS_BIRTH']))

# Add in the target
poly_features['TARGET'] = poly_target

# Find the correlations with the target
poly_corrs = poly_features.corr()['TARGET'].sort_values()

# Display most negative and most positive
print(poly_corrs.head(10))
print(poly_corrs.tail(5))
```

```
EXT_SOURCE_2 EXT_SOURCE_3          -0.193939
EXT_SOURCE_1 EXT_SOURCE_2 EXT_SOURCE_3 -0.189605
EXT_SOURCE_2 EXT_SOURCE_3 DAYS_BIRTH  -0.181283
EXT_SOURCE_2^2 EXT_SOURCE_3          -0.176428
EXT_SOURCE_2 EXT_SOURCE_3^2          -0.172282
EXT_SOURCE_1 EXT_SOURCE_2            -0.166625
EXT_SOURCE_1 EXT_SOURCE_3            -0.164065
EXT_SOURCE_2                        -0.160295
EXT_SOURCE_2 DAYS_BIRTH              -0.156873
EXT_SOURCE_1 EXT_SOURCE_2^2          -0.156867
Name: TARGET, dtype: float64
DAYS_BIRTH      -0.078239
DAYS_BIRTH^2    -0.076672
DAYS_BIRTH^3    -0.074273
TARGET          1.000000
1              NaN
Name: TARGET, dtype: float64
```

Several of the new variables have a greater (in terms of absolute magnitude) correlation with the target than the original features. When we build machine learning models, we can **try with and without these features** to determine if they actually help the model learn.

We will add these features to a copy of the training and testing data and then evaluate models with and without the features. Many times in machine learning, the only way to know if an approach will work is to try it out!

```
[44]: # Put test features into dataframe
poly_features_test = pd.DataFrame(poly_features_test,
                                   columns = poly_transformer.
                                   →get_feature_names(['EXT_SOURCE_1', 'EXT_SOURCE_2',
```

```

→'EXT_SOURCE_3', 'DAYS_BIRTH'])))

# Merge polynomial features into training dataframe
poly_features['SK_ID_CURR'] = app_train['SK_ID_CURR']
app_train_poly = app_train.merge(poly_features, on = 'SK_ID_CURR', how = 'left')

# Merge polynomial features into testing dataframe
poly_features_test['SK_ID_CURR'] = app_test['SK_ID_CURR']
app_test_poly = app_test.merge(poly_features_test, on = 'SK_ID_CURR', how =
→'left')

# Align the dataframes
app_train_poly, app_test_poly = app_train_poly.align(app_test_poly, join =
→'inner', axis = 1)

# Print out the new shapes
print('Training data with polynomial features shape: ', app_train_poly.shape)
print('Testing data with polynomial features shape: ', app_test_poly.shape)

```

Training data with polynomial features shape: (307511, 275)

Testing data with polynomial features shape: (48744, 275)

2.2 Domain Knowledge Features

Maybe it's not entirely correct to call this "domain knowledge" because I'm not a credit expert, but perhaps we could call this "attempts at applying limited financial knowledge". In this frame of mind, we can make a couple features that attempt to capture what we think may be important for telling whether a client will default on a loan. Here I'm going to use five features that were inspired by this script by Aguiar:

- CREDIT_INCOME_PERCENT: the percentage of the credit amount relative to a client's income
- ANNUITY_INCOME_PERCENT: the percentage of the loan annuity relative to a client's income
- CREDIT_TERM: the length of the payment in months (since the annuity is the monthly amount due)
- DAYS_EMPLOYED_PERCENT: the percentage of the days employed relative to the client's age

Again, thanks to Aguiar and [his great script](#) for exploring these features.

```

[45]: app_train_domain = app_train.copy()
      app_test_domain = app_test.copy()

      app_train_domain['CREDIT_INCOME_PERCENT'] = app_train_domain['AMT_CREDIT'] /
→app_train_domain['AMT_INCOME_TOTAL']
      app_train_domain['ANNUITY_INCOME_PERCENT'] = app_train_domain['AMT_ANNUITY'] /
→app_train_domain['AMT_INCOME_TOTAL']
      app_train_domain['CREDIT_TERM'] = app_train_domain['AMT_ANNUITY'] /
→app_train_domain['AMT_CREDIT']

```

```
app_train_domain['DAYS_EMPLOYED_PERCENT'] = app_train_domain['DAYS_EMPLOYED'] /
    ↳app_train_domain['DAYS_BIRTH']
```

```
[46]: app_test_domain['CREDIT_INCOME_PERCENT'] = app_test_domain['AMT_CREDIT'] /
    ↳app_test_domain['AMT_INCOME_TOTAL']
app_test_domain['ANNUITY_INCOME_PERCENT'] = app_test_domain['AMT_ANNUITY'] /
    ↳app_test_domain['AMT_INCOME_TOTAL']
app_test_domain['CREDIT_TERM'] = app_test_domain['AMT_ANNUITY'] /
    ↳app_test_domain['AMT_CREDIT']
app_test_domain['DAYS_EMPLOYED_PERCENT'] = app_test_domain['DAYS_EMPLOYED'] /
    ↳app_test_domain['DAYS_BIRTH']
```

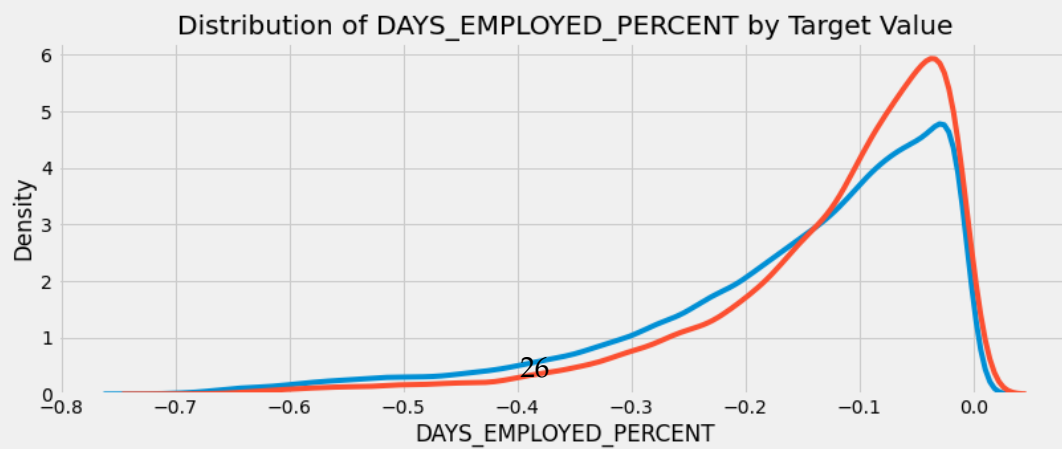
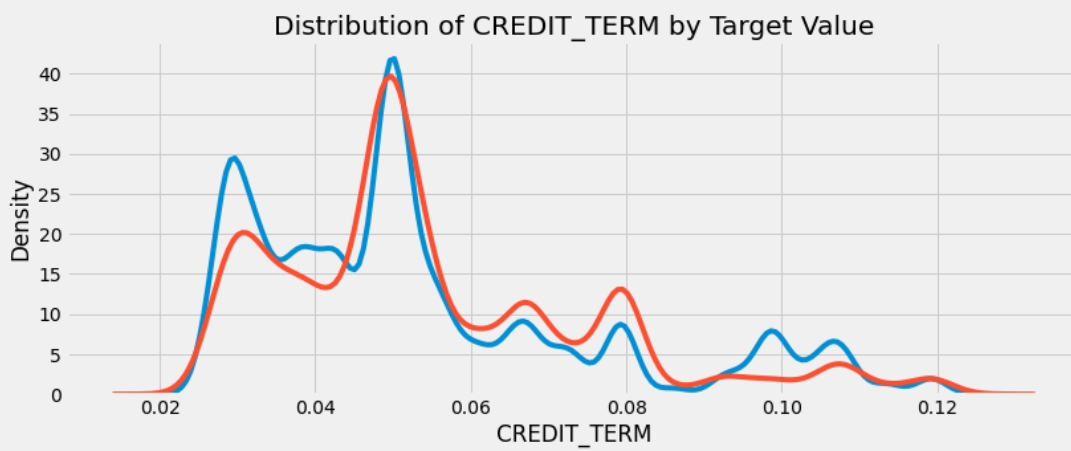
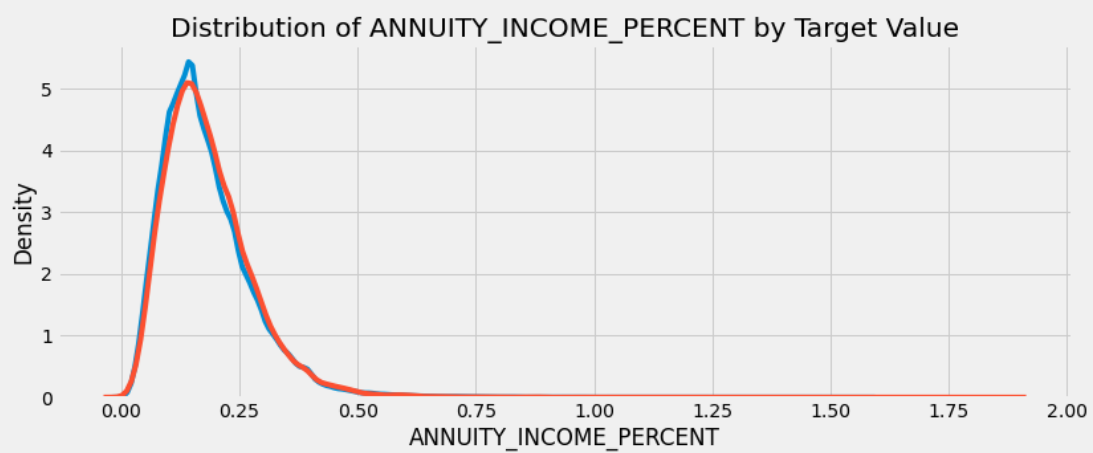
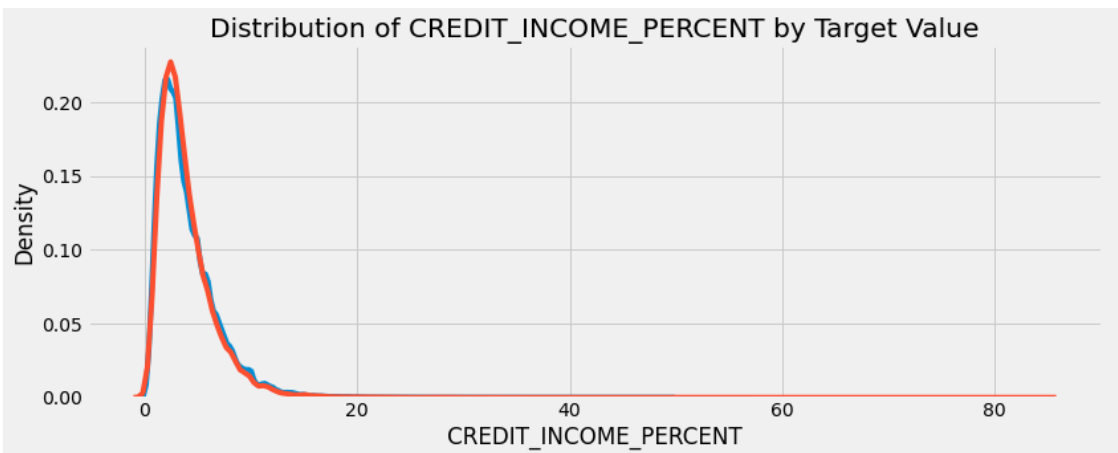
Visualize New Variables We should explore these **domain knowledge** variables visually in a graph. For all of these, we will make the same KDE plot colored by the value of the TARGET.

```
[47]: plt.figure(figsize = (12, 20))
# iterate through the new features
for i, feature in enumerate(['CREDIT_INCOME_PERCENT', 'ANNUITY_INCOME_PERCENT',
    ↳'CREDIT_TERM', 'DAYS_EMPLOYED_PERCENT']):

    # create a new subplot for each source
    plt.subplot(4, 1, i + 1)
    # plot repaid loans
    sns.kdeplot(app_train_domain.loc[app_train_domain['TARGET'] == 0, feature],
    ↳label = 'target == 0')
    # plot loans that were not repaid
    sns.kdeplot(app_train_domain.loc[app_train_domain['TARGET'] == 1, feature],
    ↳label = 'target == 1')

    # Label the plots
    plt.title('Distribution of %s by Target Value' % feature)
    plt.xlabel('%s' % feature); plt.ylabel('Density');

plt.tight_layout(h_pad = 2.5)
```



It's hard to say ahead of time if these new features will be useful. The only way to tell for sure is to try them out!

3 Baseline

For a naive baseline, we could guess the same value for all examples on the testing set. We are asked to predict the probability of not repaying the loan, so if we are entirely unsure, we would guess 0.5 for all observations on the test set. This will get us a **Reciever Operating Characteristic Area Under the Curve** (AUC ROC) of 0.5 in the competition ([random guessing on a classification task will score a 0.5](#)).

Since we already know what score we are going to get, we don't really need to make a naive baseline guess. Let's use a slightly more sophisticated model for our actual baseline: Logistic Regression.

3.1 Logistic Regression Implementation

Here I will focus on implementing the model rather than explaining the details, but for those who want to learn more about the theory of machine learning algorithms, I recommend both [An Introduction to Statistical Learning](#) and [Hands-On Machine Learning with Scikit-Learn and TensorFlow](#). Both of these books present the theory and also the code needed to make the models (in R and Python respectively). They both teach with the mindset that the best way to learn is by doing, and they are very effective!

To get a baseline, we will use all of the features after encoding the categorical variables. We will preprocess the data by filling in the missing values (imputation) and normalizing the range of the features (feature scaling). The following code performs both of these preprocessing steps.

```
[ ]: [50]: from sklearn.preprocessing import MinMaxScaler
      from sklearn.impute import SimpleImputer

      # Drop the target from the training data
      if 'TARGET' in app_train:
          train = app_train.drop(columns = ['TARGET'])
      else:
          train = app_train.copy()

      # Feature names
      features = list(train.columns)

      # Copy of the testing data
      test = app_test.copy()

      # Median imputation of missing values
      imputer = SimpleImputer(strategy = 'median')
```

```

# Scale each feature to 0-1
scaler = MinMaxScaler(feature_range = (0, 1))

# Fit on the training data
imputer.fit(train)

# Transform both training and testing data
train = imputer.transform(train)
test = imputer.transform(app_test)

# Repeat with the scaler
scaler.fit(train)
train = scaler.transform(train)
test = scaler.transform(test)

print('Training data shape: ', train.shape)
print('Testing data shape: ', test.shape)

```

Training data shape: (307511, 240)

Testing data shape: (48744, 240)

We will use [LogisticRegression](#) from [Scikit-Learn](#) for our first model. The only change we will make from the default model settings is to lower the [regularization parameter](#), *C*, which controls the amount of overfitting (a lower value should decrease overfitting). This will get us slightly better results than the default `LogisticRegression`, but it still will set a low bar for any future models.

Here we use the familiar Scikit-Learn modeling syntax: we first create the model, then we train the model using `.fit` and then we make predictions on the testing data using `.predict_proba` (remember that we want probabilities and not a 0 or 1).

```

[51]: from sklearn.linear_model import LogisticRegression

# Make the model with the specified regularization parameter
log_reg = LogisticRegression(C = 0.0001)

# Train on the training data
log_reg.fit(train, train_labels)

```

```

[51]: LogisticRegression(C=0.0001, class_weight=None, dual=False, fit_intercept=True,
                        intercept_scaling=1, l1_ratio=None, max_iter=100,
                        multi_class='auto', n_jobs=None, penalty='l2',
                        random_state=None, solver='lbfgs', tol=0.0001, verbose=0,
                        warm_start=False)

```

Now that the model has been trained, we can use it to make predictions. We want to predict the probabilities of not paying a loan, so we use the model `predict_proba` method. This returns an $m \times 2$ array where m is the number of observations. The first column is the probability of the target being 0 and the second column is the probability of the target being 1 (so for a single row, the two columns must sum to 1). We want the probability the loan is not repaid, so we will select the second column.

The following code makes the predictions and selects the correct column.

```
[52]: # Make predictions
# Make sure to select the second column only
log_reg_pred = log_reg.predict_proba(test)[: , 1]
```

The predictions must be in the format shown in the `sample_submission.csv` file, where there are only two columns: `SK_ID_CURR` and `TARGET`. We will create a dataframe in this format from the test set and the predictions called `submit`.

```
[53]: # Submission dataframe
submit = app_test[['SK_ID_CURR']]
submit['TARGET'] = log_reg_pred

submit.head()
```

```
[53]: SK_ID_CURR    TARGET
0      100001    0.078515
1      100005    0.137926
2      100013    0.082194
3      100028    0.080921
4      100038    0.132618
```

The predictions represent a probability between 0 and 1 that the loan will not be repaid. If we were using these predictions to classify applicants, we could set a probability threshold for determining that a loan is risky.

```
[54]: # Save the submission to a csv file
submit.to_csv('log_reg_baseline.csv', index = False)
```

The submission has now been saved to the virtual environment in which our notebook is running. To access the submission, at the end of the notebook, we will hit the blue Commit & Run button at the upper right of the kernel. This runs the entire notebook and then lets us download any files that are created during the run.

Once we run the notebook, the files created are available in the Versions tab under the Output sub-tab. From here, the submission files can be submitted to the competition or downloaded. Since there are several models in this notebook, there will be multiple output files.

The logistic regression baseline should score around 0.671 when submitted.

3.2 Improved Model: Random Forest

To try and beat the poor performance of our baseline, we can update the algorithm. Let's try using a Random Forest on the same training data to see how that affects performance. The Random Forest is a much more powerful model especially when we use hundreds of trees. We will use 100 trees in the random forest.

```
[55]: from sklearn.ensemble import RandomForestClassifier

# Make the random forest classifier
random_forest = RandomForestClassifier(n_estimators = 100, random_state = 50,
→ verbose = 1, n_jobs = -1)
```

```
[56]: # Train on the training data
random_forest.fit(train, train_labels)

# Extract feature importances
feature_importance_values = random_forest.feature_importances_
feature_importances = pd.DataFrame({'feature': features, 'importance':
    ↳feature_importance_values})

# Make predictions on the test data
predictions = random_forest.predict_proba(test)[: , 1]
```

```
[Parallel(n_jobs=-1)]: Using backend ThreadingBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done 46 tasks      | elapsed: 1.1min
[Parallel(n_jobs=-1)]: Done 100 out of 100 | elapsed: 2.4min finished
[Parallel(n_jobs=2)]: Using backend ThreadingBackend with 2 concurrent workers.
[Parallel(n_jobs=2)]: Done 46 tasks      | elapsed: 0.7s
[Parallel(n_jobs=2)]: Done 100 out of 100 | elapsed: 1.7s finished
```

```
[57]: # Make a submission dataframe
submit = app_test[['SK_ID_CURR']]
submit['TARGET'] = predictions

# Save the submission dataframe
submit.to_csv('random_forest_baseline.csv', index = False)
```

These predictions will also be available when we run the entire notebook.
This model should score around 0.678 when submitted.

3.2.1 Make Predictions using Engineered Features

The only way to see if the Polynomial Features and Domain knowledge improved the model is to train a test a model on these features! We can then compare the submission performance to that for the model without these features to gauge the effect of our feature engineering.

```
[59]: poly_features_names = list(app_train_poly.columns)

# Impute the polynomial features
imputer = SimpleImputer(strategy = 'median')

poly_features = imputer.fit_transform(app_train_poly)
poly_features_test = imputer.transform(app_test_poly)

# Scale the polynomial features
scaler = MinMaxScaler(feature_range = (0, 1))

poly_features = scaler.fit_transform(poly_features)
poly_features_test = scaler.transform(poly_features_test)
```

```
random_forest_poly = RandomForestClassifier(n_estimators = 100, random_state = 50, verbose = 1, n_jobs = -1)
```

```
[60]: # Train on the training data
random_forest_poly.fit(poly_features, train_labels)

# Make predictions on the test data
predictions = random_forest_poly.predict_proba(poly_features_test)[: , 1]
```

```
[Parallel(n_jobs=-1)]: Using backend ThreadingBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done 46 tasks      | elapsed: 2.0min
[Parallel(n_jobs=-1)]: Done 100 out of 100 | elapsed: 4.3min finished
[Parallel(n_jobs=2)]: Using backend ThreadingBackend with 2 concurrent workers.
[Parallel(n_jobs=2)]: Done 46 tasks      | elapsed: 0.7s
[Parallel(n_jobs=2)]: Done 100 out of 100 | elapsed: 1.4s finished
```

```
[61]: # Make a submission dataframe
submit = app_test[['SK_ID_CURR']]
submit['TARGET'] = predictions

# Save the submission dataframe
submit.to_csv('random_forest_baseline_engineered.csv', index = False)
```

This model scored 0.678 when submitted to the competition, exactly the same as that without the engineered features. Given these results, it does not appear that our feature construction helped in this case.

Testing Domain Features Now we can test the domain features we made by hand.

```
[62]: app_train_domain = app_train_domain.drop(columns = 'TARGET')

domain_features_names = list(app_train_domain.columns)

# Impute the domainnomial features
imputer = SimpleImputer(strategy = 'median')

domain_features = imputer.fit_transform(app_train_domain)
domain_features_test = imputer.transform(app_test_domain)

# Scale the domainnomial features
scaler = MinMaxScaler(feature_range = (0, 1))

domain_features = scaler.fit_transform(domain_features)
domain_features_test = scaler.transform(domain_features_test)

random_forest_domain = RandomForestClassifier(n_estimators = 100, random_state = 50, verbose = 1, n_jobs = -1)

# Train on the training data
```

```

random_forest_domain.fit(domain_features, train_labels)

# Extract feature importances
feature_importance_values_domain = random_forest_domain.feature_importances_
feature_importances_domain = pd.DataFrame({'feature': domain_features_names,
→ 'importance': feature_importance_values_domain})

# Make predictions on the test data
predictions = random_forest_domain.predict_proba(domain_features_test)[: , 1]

```

```

[Parallel(n_jobs=-1)]: Using backend ThreadingBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done 46 tasks      | elapsed: 1.4min
[Parallel(n_jobs=-1)]: Done 100 out of 100 | elapsed: 3.0min finished
[Parallel(n_jobs=2)]: Using backend ThreadingBackend with 2 concurrent workers.
[Parallel(n_jobs=2)]: Done 46 tasks      | elapsed: 0.9s
[Parallel(n_jobs=2)]: Done 100 out of 100 | elapsed: 1.9s finished

```

```

[63]: # Make a submission dataframe
submit = app_test[['SK_ID_CURR']]
submit['TARGET'] = predictions

# Save the submission dataframe
submit.to_csv('random_forest_baseline_domain.csv', index = False)

```

This scores 0.679 when submitted which probably shows that the engineered features do not help in this model (however they do help in the Gradient Boosting Model at the end of the notebook).

In later notebooks, we will do more **feature engineering** by using the information from the other data sources. From experience, this will definitely help our model!

3.3 Model Interpretation: Feature Importances

As a simple method to see which variables are the most relevant, we can look at the feature importances of the random forest. Given the correlations we saw in the exploratory data analysis, we should expect that the most important features are the **EXT_SOURCE** and the **DAYS_BIRTH**. We may use these feature importances as a method of **dimensionality reduction** in future work.

```

[64]: def plot_feature_importances(df):
    """
    Plot importances returned by a model. This can work with any measure of
    feature importance provided that higher importance is better.

    Args:
        df (dataframe): feature importances. Must have the features in a column
        called `features` and the importances in a column called `importance`

    Returns:
        shows a plot of the 15 most importance features
    """

```

```

    df (dataframe): feature importances sorted by importance (highest to
    →lowest)
    with a column for normalized importance
    """

    # Sort features according to importance
    df = df.sort_values('importance', ascending = False).reset_index()

    # Normalize the feature importances to add up to one
    df['importance_normalized'] = df['importance'] / df['importance'].sum()

    # Make a horizontal bar chart of feature importances
    plt.figure(figsize = (10, 6))
    ax = plt.subplot()

    # Need to reverse the index to plot most important on top
    ax.barh(list(reversed(list(df.index[:15]))),
            df['importance_normalized'].head(15),
            align = 'center', edgecolor = 'k')

    # Set the yticks and labels
    ax.set_yticks(list(reversed(list(df.index[:15]))))
    ax.set_yticklabels(df['feature'].head(15))

    # Plot labeling
    plt.xlabel('Normalized Importance'); plt.title('Feature Importances')
    plt.show()

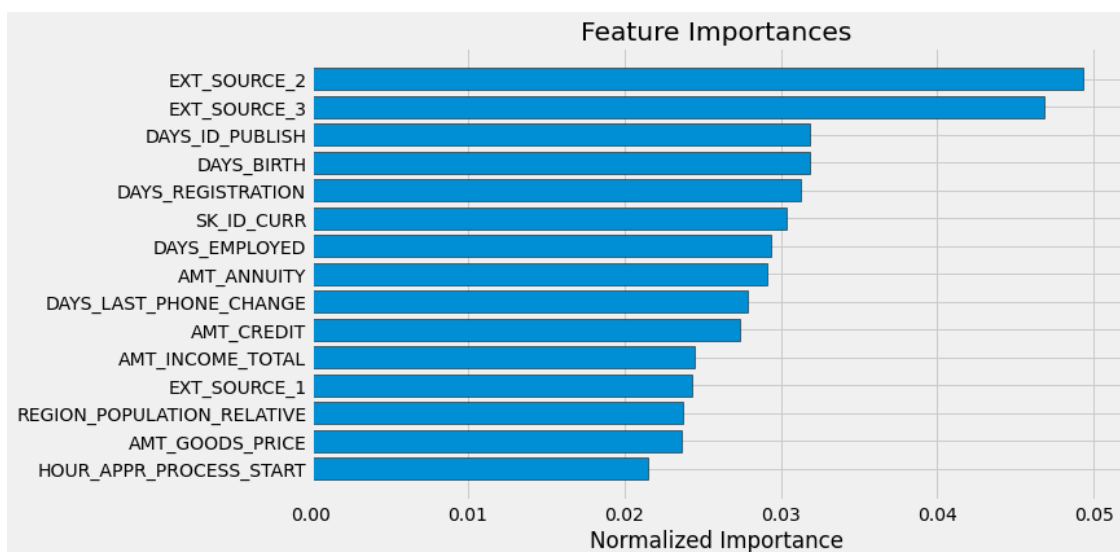
    return df

```

```

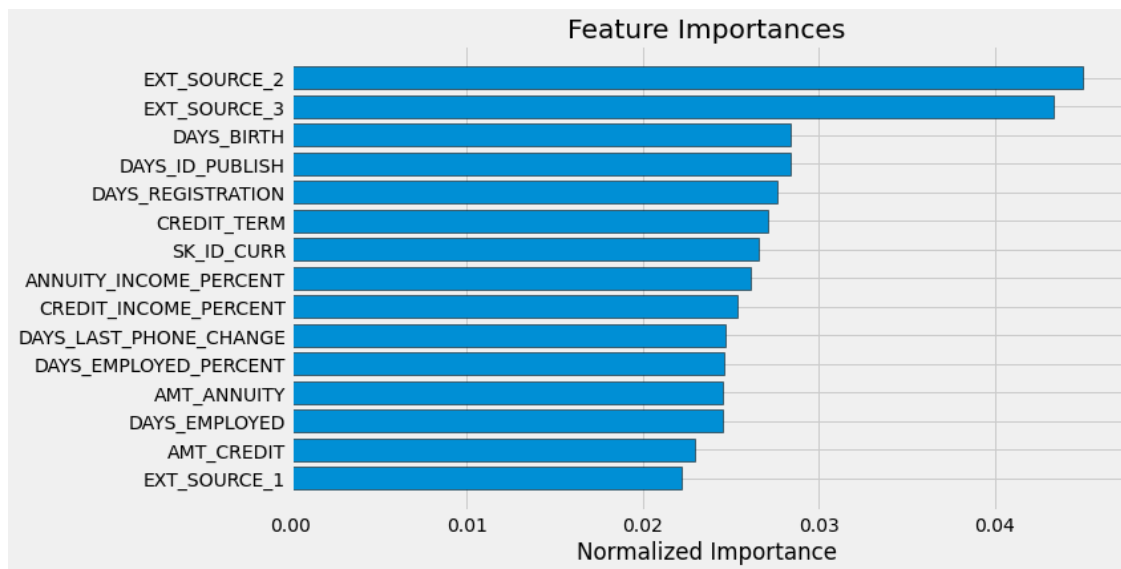
[65]: # Show the feature importances for the default features
feature_importances_sorted = plot_feature_importances(feature_importances)

```



As expected, the most important features are those dealing with EXT_SOURCE and DAYS_BIRTH. We see that there are only a handful of features with a significant importance to the model, which suggests we may be able to drop many of the features without a decrease in performance (and we may even see an increase in performance.) Feature importances are not the most sophisticated method to interpret a model or perform dimensionality reduction, but they let us start to understand what factors our model takes into account when it makes predictions.

```
[66]: feature_importances_domain_sorted =  
      → plot_feature_importances(feature_importances_domain)
```



We see that all four of our hand-engineered features made it into the top 15 most important! This should give us confidence that our domain knowledge was at least partially on track.

4 Conclusions

In this notebook, we saw how to get started with a Kaggle machine learning competition. We first made sure to understand the data, our task, and the metric by which our submissions will be judged. Then, we performed a fairly simple EDA to try and identify relationships, trends, or anomalies that may help our modeling. Along the way, we performed necessary preprocessing steps such as encoding categorical variables, imputing missing values, and scaling features to a range. Then, we constructed new features out of the existing data to see if doing so could help our model.

Once the data exploration, data preparation, and feature engineering was complete, we implemented a baseline model upon which we hope to improve. Then we built a second slightly more complicated model to beat our first score. We also carried out an experiment to determine the effect of adding the engineering variables.

We followed the general outline of a [machine learning project](#):

1. Understand the problem and the data
2. Data cleaning and formatting (this was mostly done for us)
3. Exploratory Data Analysis
4. Baseline model
5. Improved model
6. Model interpretation (just a little)

Machine learning competitions do differ slightly from typical data science problems in that we are concerned only with achieving the best performance on a single metric and do not care about the interpretation. However, by attempting to understand how our models make decisions, we can try to improve them or examine the mistakes in order to correct the errors. In future notebooks we will look at incorporating more sources of data, building more complex models (by following the code of others), and improving our scores.

I hope this notebook was able to get you up and running in this machine learning competition and that you are now ready to go out on your own - with help from the community - and start working on some great problems!

Running the notebook: now that we are at the end of the notebook, you can hit the blue Commit & Run button to execute all the code at once. After the run is complete (this should take about 10 minutes), you can then access the files that were created by going to the versions tab and then the output sub-tab. The submission files can be directly submitted to the competition from this tab or they can be downloaded to a local machine and saved. The final part is to share the notebook: go to the settings tab and change the visibility to Public. This allows the entire world to see your work!

4.0.1 Follow-up Notebooks

For those looking to keep working on this problem, I have a series of follow-up notebooks:

- [Manual Feature Engineering Part One](#)
- [Manual Feature Engineering Part Two](#)
- [Introduction to Automated Feature Engineering](#)
- [Advanced Automated Feature Engineering](#)
- [Feature Selection](#)
- [Intro to Model Tuning: Grid and Random Search](#)

As always, I welcome feedback and constructive criticism. I write for Towards Data Science at <https://medium.com/@williamkoehrsen/> and can be reached on Twitter at https://twitter.com/koehrsen_will
Will

5 Just for Fun: Light Gradient Boosting Machine

Now (if you want, this part is entirely optional) we can step off the deep end and use a real machine learning model: the [gradient boosting machine](#) using the [LightGBM library](#)! The Gradient Boosting Machine is currently the leading model for learning on structured datasets (especially on Kaggle) and we will probably need some form of this model to do well in the competition. Don't worry, even if this code looks intimidating, it's just a series of small steps that build up to a complete model. I added this code just to show what may be in store for this project, and because

it gets us a slightly better score on the leaderboard. In future notebooks we will see how to work with more advanced models (which mostly means adapting existing code to make it work better), feature engineering, and feature selection. See you in the next notebook!

```
[67]: from sklearn.model_selection import KFold
from sklearn.metrics import roc_auc_score
import lightgbm as lgb
import gc

def model(features, test_features, encoding = 'ohe', n_folds = 5):

    """Train and test a light gradient boosting model using
    cross validation.

    Parameters
    -----
        features (pd.DataFrame):
            dataframe of training features to use
            for training a model. Must include the TARGET column.
        test_features (pd.DataFrame):
            dataframe of testing features to use
            for making predictions with the model.
        encoding (str, default = 'ohe'):
            method for encoding categorical variables. Either 'ohe' for one-hot_
            →encoding or 'le' for integer label encoding
        n_folds (int, default = 5): number of folds to use for cross_
            →validation

    Return
    -----
        submission (pd.DataFrame):
            dataframe with `SK_ID_CURR` and `TARGET` probabilities
            predicted by the model.
        feature_importances (pd.DataFrame):
            dataframe with the feature importances from the model.
        valid_metrics (pd.DataFrame):
            dataframe with training and validation metrics (ROC AUC) for each_
            →fold and overall.

    """

    # Extract the ids
    train_ids = features['SK_ID_CURR']
    test_ids = test_features['SK_ID_CURR']

    # Extract the labels for training
    labels = features['TARGET']
```

```

# Remove the ids and target
features = features.drop(columns = ['SK_ID_CURR', 'TARGET'])
test_features = test_features.drop(columns = ['SK_ID_CURR'])

# One Hot Encoding
if encoding == 'ohe':
    features = pd.get_dummies(features)
    test_features = pd.get_dummies(test_features)

    # Align the dataframes by the columns
    features, test_features = features.align(test_features, join = 'inner',
→axis = 1)

    # No categorical indices to record
    cat_indices = 'auto'

# Integer label encoding
elif encoding == 'le':

    # Create a label encoder
    label_encoder = LabelEncoder()

    # List for storing categorical indices
    cat_indices = []

    # Iterate through each column
    for i, col in enumerate(features):
        if features[col].dtype == 'object':
            # Map the categorical features to integers
            features[col] = label_encoder.fit_transform(np.
→array(features[col].astype(str)).reshape((-1,)))
            test_features[col] = label_encoder.transform(np.
→array(test_features[col].astype(str)).reshape((-1,)))

            # Record the categorical indices
            cat_indices.append(i)

# Catch error if label encoding scheme is not valid
else:
    raise ValueError("Encoding must be either 'ohe' or 'le'")

print('Training Data Shape: ', features.shape)
print('Testing Data Shape: ', test_features.shape)

# Extract feature names
feature_names = list(features.columns)

```

```

# Convert to np arrays
features = np.array(features)
test_features = np.array(test_features)

# Create the kfold object
k_fold = KFold(n_splits = n_folds, shuffle = True, random_state = 50)

# Empty array for feature importances
feature_importance_values = np.zeros(len(feature_names))

# Empty array for test predictions
test_predictions = np.zeros(test_features.shape[0])

# Empty array for out of fold validation predictions
out_of_fold = np.zeros(features.shape[0])

# Lists for recording validation and training scores
valid_scores = []
train_scores = []

# Iterate through each fold
for train_indices, valid_indices in k_fold.split(features):

    # Training data for the fold
    train_features, train_labels = features[train_indices],
→labels[train_indices]
    # Validation data for the fold
    valid_features, valid_labels = features[valid_indices],
→labels[valid_indices]

    # Create the model
    model = lgb.LGBMClassifier(n_estimators=10000, objective = 'binary',
                                class_weight = 'balanced', learning_rate = 0.
→05,
                                reg_alpha = 0.1, reg_lambda = 0.1,
                                subsample = 0.8, n_jobs = -1, random_state =
→50)

    # Train the model
    model.fit(train_features, train_labels, eval_metric = 'auc',
              eval_set = [(valid_features, valid_labels), (train_features,
→train_labels)],
              eval_names = ['valid', 'train'], categorical_feature =
→cat_indices,
              early_stopping_rounds = 100, verbose = 200)

```

```

    # Record the best iteration
    best_iteration = model.best_iteration_

    # Record the feature importances
    feature_importance_values += model.feature_importances_ / k_fold.
→n_splits

    # Make predictions
    test_predictions += model.predict_proba(test_features, num_iteration =
→best_iteration)[:, 1] / k_fold.n_splits

    # Record the out of fold predictions
    out_of_fold[valid_indices] = model.predict_proba(valid_features,
→num_iteration = best_iteration)[:, 1]

    # Record the best score
    valid_score = model.best_score_['valid']['auc']
    train_score = model.best_score_['train']['auc']

    valid_scores.append(valid_score)
    train_scores.append(train_score)

    # Clean up memory
    gc.enable()
    del model, train_features, valid_features
    gc.collect()

    # Make the submission dataframe
    submission = pd.DataFrame({'SK_ID_CURR': test_ids, 'TARGET':
→test_predictions})

    # Make the feature importance dataframe
    feature_importances = pd.DataFrame({'feature': feature_names, 'importance':
→feature_importance_values})

    # Overall validation score
    valid_auc = roc_auc_score(labels, out_of_fold)

    # Add the overall scores to the metrics
    valid_scores.append(valid_auc)
    train_scores.append(np.mean(train_scores))

    # Needed for creating dataframe of validation scores
    fold_names = list(range(n_folds))
    fold_names.append('overall')

```

```

# Dataframe of validation scores
metrics = pd.DataFrame({'fold': fold_names,
                        'train': train_scores,
                        'valid': valid_scores})

return submission, feature_importances, metrics

```

```

[68]: submission, fi, metrics = model(app_train, app_test)
print('Baseline metrics')
print(metrics)

```

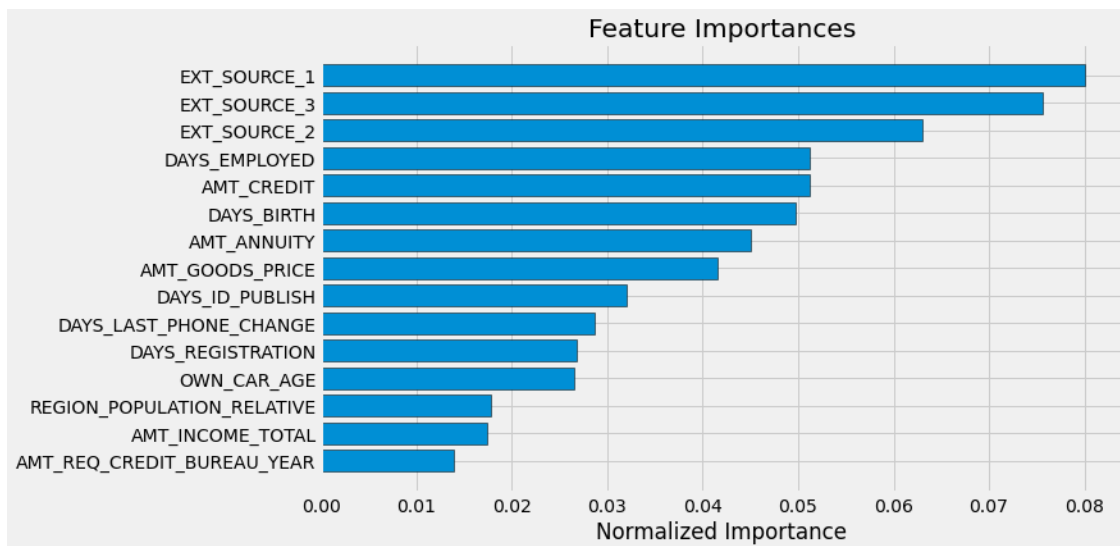
```

Training Data Shape: (307511, 239)
Testing Data Shape: (48744, 239)
Training until validation scores don't improve for 100 rounds.
[200] train's auc: 0.79887 train's binary_logloss: 0.547648 valid's
auc: 0.754949 valid's binary_logloss: 0.563125
Early stopping, best iteration is:
[208] train's auc: 0.80025 train's binary_logloss: 0.546264 valid's
auc: 0.755109 valid's binary_logloss: 0.562276
Training until validation scores don't improve for 100 rounds.
[200] train's auc: 0.798518 train's binary_logloss: 0.548144 valid's
auc: 0.758539 valid's binary_logloss: 0.563479
Early stopping, best iteration is:
[217] train's auc: 0.801374 train's binary_logloss: 0.545314 valid's
auc: 0.758619 valid's binary_logloss: 0.561732
Training until validation scores don't improve for 100 rounds.
[200] train's auc: 0.79774 train's binary_logloss: 0.54923 valid's auc:
0.762652 valid's binary_logloss: 0.564246
[400] train's auc: 0.827288 train's binary_logloss: 0.520152 valid's
auc: 0.762202 valid's binary_logloss: 0.546576
Early stopping, best iteration is:
[320] train's auc: 0.81638 train's binary_logloss: 0.531111 valid's
auc: 0.763103 valid's binary_logloss: 0.553039
Training until validation scores don't improve for 100 rounds.
[200] train's auc: 0.799107 train's binary_logloss: 0.547723 valid's
auc: 0.757496 valid's binary_logloss: 0.562014
Early stopping, best iteration is:
[183] train's auc: 0.796125 train's binary_logloss: 0.550639 valid's
auc: 0.75759 valid's binary_logloss: 0.563795
Training until validation scores don't improve for 100 rounds.
[200] train's auc: 0.798268 train's binary_logloss: 0.548197 valid's
auc: 0.758099 valid's binary_logloss: 0.564499
Early stopping, best iteration is:
[227] train's auc: 0.802746 train's binary_logloss: 0.543868 valid's
auc: 0.758251 valid's binary_logloss: 0.561904
Baseline metrics
      fold      train      valid
0         0  0.800250  0.755109

```

1	1	0.801374	0.758619
2	2	0.816380	0.763103
3	3	0.796125	0.757590
4	4	0.802746	0.758251
5	overall	0.803375	0.758537

```
[69]: fi_sorted = plot_feature_importances(fi)
```



```
[70]: submission.to_csv('baseline_lgb.csv', index = False)
```

This submission should score about 0.735 on the leaderboard. We will certainly best that in future work!

```
[71]: app_train_domain['TARGET'] = train_labels

# Test the domain knowledge features
submission_domain, fi_domain, metrics_domain = model(app_train_domain,
→app_test_domain)
print('Baseline with domain knowledge features metrics')
print(metrics_domain)
```

Training Data Shape: (307511, 243)

Testing Data Shape: (48744, 243)

Training until validation scores don't improve for 100 rounds.

[200]	train's auc: 0.804531	train's binary_logloss: 0.541661	valid's
	auc: 0.762577	valid's binary_logloss: 0.557281	

Early stopping, best iteration is:

[237]	train's auc: 0.810671	train's binary_logloss: 0.535426	valid's
	auc: 0.762858	valid's binary_logloss: 0.553438	

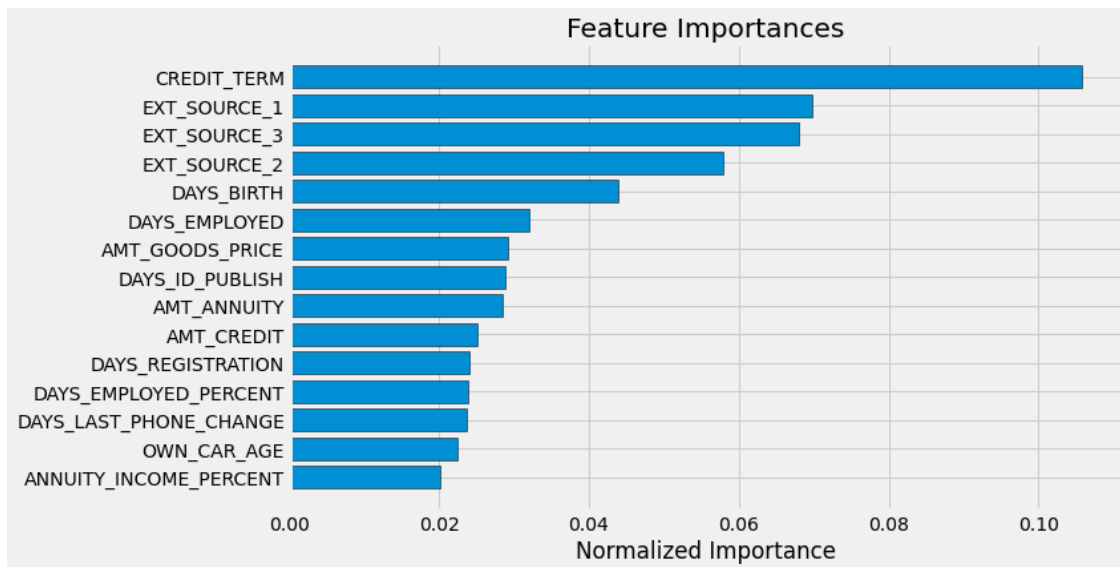
Training until validation scores don't improve for 100 rounds.

```

[200]   train's auc: 0.804304   train's binary_logloss: 0.542018   valid's
auc: 0.765594   valid's binary_logloss: 0.55808
Early stopping, best iteration is:
[227]   train's auc: 0.808665   train's binary_logloss: 0.537574   valid's
auc: 0.765861   valid's binary_logloss: 0.555268
Training until validation scores don't improve for 100 rounds.
[200]   train's auc: 0.803753   train's binary_logloss: 0.542936   valid's
auc: 0.770139   valid's binary_logloss: 0.557892
[400]   train's auc: 0.834338   train's binary_logloss: 0.511693   valid's
auc: 0.770328   valid's binary_logloss: 0.538395
Early stopping, best iteration is:
[302]   train's auc: 0.820401   train's binary_logloss: 0.526044   valid's
auc: 0.770629   valid's binary_logloss: 0.547303
Training until validation scores don't improve for 100 rounds.
[200]   train's auc: 0.804487   train's binary_logloss: 0.542071   valid's
auc: 0.765653   valid's binary_logloss: 0.556352
Early stopping, best iteration is:
[262]   train's auc: 0.815066   train's binary_logloss: 0.53137   valid's auc:
0.766318   valid's binary_logloss: 0.549785
Training until validation scores don't improve for 100 rounds.
[200]   train's auc: 0.804527   train's binary_logloss: 0.541724   valid's
auc: 0.764456   valid's binary_logloss: 0.55882
Early stopping, best iteration is:
[235]   train's auc: 0.810422   train's binary_logloss: 0.535826   valid's
auc: 0.764517   valid's binary_logloss: 0.55519
Baseline with domain knowledge features metrics
      fold    train    valid
0         0  0.810671  0.762858
1         1  0.808665  0.765861
2         2  0.820401  0.770629
3         3  0.815066  0.766318
4         4  0.810422  0.764517
5  overall  0.813045  0.766050

```

```
[72]: fi_sorted = plot_feature_importances(fi_domain)
```

Again, we see that some of our features made it into the most important. Going forward, we will need to think about whether domain knowledge features may be useful for this problem (or we should consult someone who knows more about the financial industry!)

```
[73]: submission_domain.to_csv('baseline_lgb_domain_features.csv', index = False)
```

This model scores about 0.754 when submitted to the public leaderboard indicating that the domain features do improve the performance! [Feature engineering](#) is going to be a critical part of this competition (as it is for all machine learning problems)!

```
[74]: !apt-get install texlive texlive-xetex texlive-latex-extra pandoc
      !pip install py pandoc
```

```
Reading package lists... Done
Building dependency tree
Reading state information... Done
pandoc is already the newest version (1.19.2.4~dfsg-1build4).
pandoc set to manually installed.
The following additional packages will be installed:
  fonts-droid-fallback fonts-lato fonts-lmodern fonts-noto-mono fonts-texgyre
  javascript-common libcupsfilters1 libcupsimage2 libgs9 libgs9-common
  libijs-0.35 libjbig2dec0 libjs-jquery libkpathsea6 libpotrace0 libptexenc1
  libruby2.5 libsyntaxtex1 libtexlua52 libtexluajit2 libzzip-0-13 lmodern
  poppler-data preview-latex-style rake ruby ruby-did-you-mean ruby-minitest
  ruby-net-telnet ruby-power-assert ruby-test-unit ruby2.5
  rubygems-integration tlutils tex-common tex-gyre texlive-base
  texlive-binaries texlive-fonts-recommended texlive-latex-base
  texlive-latex-recommended texlive-pictures texlive-plain-generic tipa
Suggested packages:
  fonts-noto apache2 | lighttpd | httpd poppler-utils ghostscript
  fonts-japanese-mincho | fonts-ipafont-mincho fonts-japanese-gothic
```

```

Setting up libruby2.5:amd64 (2.5.1-1ubuntu1.7) ...
Processing triggers for mime-support (3.60ubuntu1) ...
Processing triggers for libc-bin (2.27-3ubuntu1.3) ...
/sbin/ldconfig.real: /usr/local/lib/python3.6/dist-
packages/ideep4py/lib/libmkldnn.so.0 is not a symbolic link

Processing triggers for man-db (2.8.3-2ubuntu0.1) ...
Processing triggers for fontconfig (2.12.6-0ubuntu2) ...
Processing triggers for tex-common (6.09) ...
Running updpmap-sys. This may take some time... done.
Running mktexlsr /var/lib/texmf ... done.
Building format(s) --all.
    This may take some time... done.
Collecting py pandoc
  Downloading https://files.pythonhosted.org/packages/d6/b7/5050dc1769c8a93d3ec7
c4bd55be161991c94b8b235f88bf7c764449e708/py pandoc-1.5.tar.gz
Requirement already satisfied: setuptools in /usr/local/lib/python3.6/dist-
packages (from py pandoc) (51.1.2)
Requirement already satisfied: pip>=8.1.0 in /usr/local/lib/python3.6/dist-
packages (from py pandoc) (19.3.1)
Requirement already satisfied: wheel>=0.25.0 in /usr/local/lib/python3.6/dist-
packages (from py pandoc) (0.36.2)
Building wheels for collected packages: py pandoc
  Building wheel for py pandoc (setup.py) ... done
  Created wheel for py pandoc: filename=py pandoc-1.5-cp36-none-any.whl size=17037
sha256=de70df329bd0bf565532a4362523d18c8cb6e691977f92363cd6ca345e7d1b2a
  Stored in directory: /root/.cache/pip/wheels/bb/7d/d6/2f9af55e800d37e42e546106
bcbd36a86e24e725e303d17e04
Successfully built py pandoc
Installing collected packages: py pandoc
Successfully installed py pandoc-1.5

```

```
[91]: !cp /content/drive/MyDrive/kaggle/week1/start-here-a-gentle-introduction.ipynb .
      ↪ /
```

```
[92]: !jupyter nbconvert --to pdf './start-here-a-gentle-introduction.ipynb'
```

```

[NbConvertApp] Converting notebook ./start-here-a-gentle-introduction.ipynb to
pdf
[NbConvertApp] Support files will be in start-here-a-gentle-introduction_files/
[NbConvertApp] Making directory ./start-here-a-gentle-introduction_files
[NbConvertApp] Making directory ./start-here-a-gentle-introduction_files
[NbConvertApp] Making directory ./start-here-a-gentle-introduction_files
[NbConvertApp] Making directory ./start-here-a-gentle-introduction_files
[NbConvertApp] Making directory ./start-here-a-gentle-introduction_files
[NbConvertApp] Making directory ./start-here-a-gentle-introduction_files
[NbConvertApp] Making directory ./start-here-a-gentle-introduction_files
[NbConvertApp] Making directory ./start-here-a-gentle-introduction_files

```