

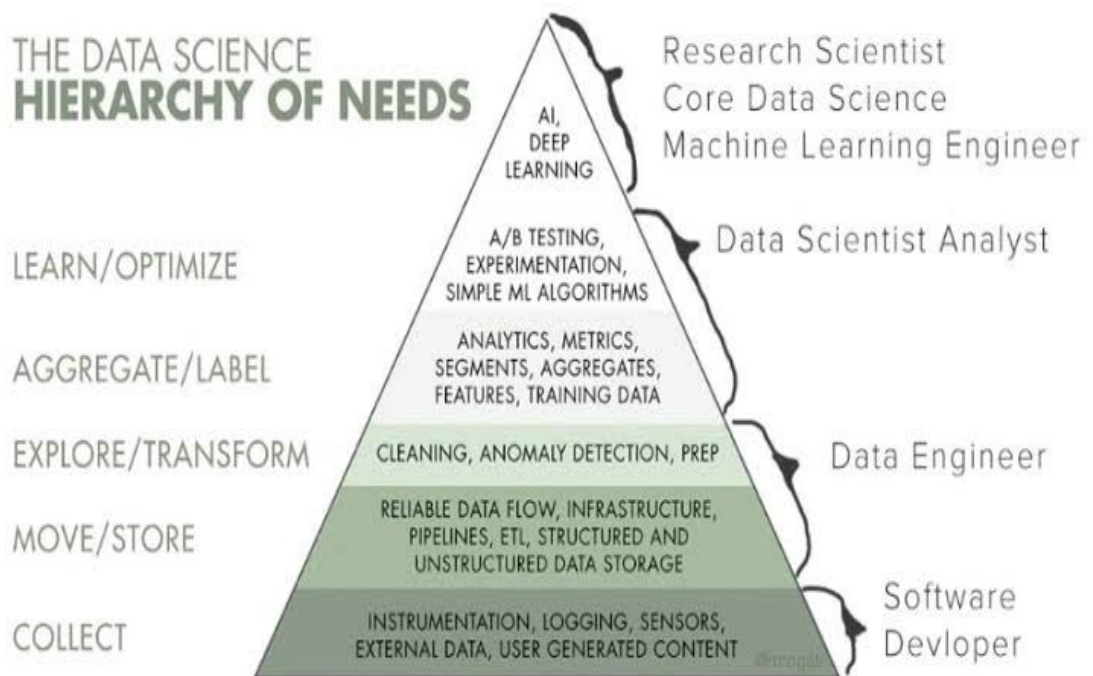
Loan Approval Prediction

Edmund Walsh - May 10th, 2020

Introduction

The project examines the data provided by the Home Mortgage Disclosure Act (HMDA) which requires mortgage lenders in the United States to disclose information about the mortgage lending decisions they have made. Specifically, we will be examining prediction of whether or not an application will be accepted or denied.

This notebook is part of a short 5-day project whose purpose is more about process than prediction of mortgages. The focus of this project will be setting up end-to-end engineering from raw data to results and presentation that will help us scale in the future. This notebook will focus on the data science approach and process and highlight a roadmap best illustrated in the image below:



A Ground Up Approach

The pyramid above is such a good illustration because the process truly is building in a step-by-step fashion towards the ultimate goal of finding useful results that are actionable and impactful in the real world. The end results get all of the attention, but a project is unlikely to be successful without these strong foundations.

Context

While this project request didn't specifically state 'why' we are looking into this data, I will work within the context of three important rationales.

1. Mortgage due diligence is expensive and time intensive. A process that can more reliably expedite the process will save lenders significant time and resources.
2. From a regulatory perspective and also importantly as Machine Intelligence become a larger and more common part of this process it is important for us to be aware of and highlight any bias.
3. Some financial institutions may be more or less likely to issue mortgages and this may reflect either an over or under utilization of the balance sheet or their risk appetite.

Preparation

Before digging in, let's install our python requirements and follow the instructions for setting up our docker environment and tools in the [README.md](https://github.com/ewalsh/hmda/blob/master/README.md) (<https://github.com/ewalsh/hmda/blob/master/README.md>).

```
In [1]: %%capture
        pip install --user -r ./misc/requirements.txt
```

Package Import

Now we can begin importing the required packages. Many of these are common python packages, the exception is seed which is a set of functions we will use to import our initial data and begin our collect step on our roadmap.

```
In [2]: import seed
        import pandas as pd
        import numpy as np
        import os
        import config
        import math
        import matplotlib.pyplot as plt
        import seaborn as sns
        sns.set_style('whitegrid')
        %matplotlib inline
```

A Discussion of the data

Our data sources will come from three major sources and we will use their APIs to download the data.

1. Our main dataset comes from the HMDA database and includes not only a large set of information about the individual loan approval decision but also information about the originating institutions.
 - a. A full list of available data on the mortgage approvals can be found [here](http://cfpb.github.io/api/hmda/fields.html) (<http://cfpb.github.io/api/hmda/fields.html>)
 - b. We also will look at the originating institutions and information about that dataset can be found [here](https://api.consumerfinance.gov/data/hmda/slice/institutions/metadata) (<https://api.consumerfinance.gov/data/hmda/slice/institutions/metadata>)
2. Our next and complementary set of information comes from the census bureau. We will use their county business patterns series which aggregate economic information at a county level. I hope that this information will provide some valuable economic insight into regional economics that may affect an mortgage approval decision.
 - a. Details about this dataset can be found [here](https://www.census.gov/programs-surveys/cbp.html) (<https://www.census.gov/programs-surveys/cbp.html>)
 - b. This API requires a key which you can sign up for [here](https://api.census.gov/data/key_signup.html) (https://api.census.gov/data/key_signup.html)
 - c. After you have signed up, please include this key in the `config.py` (https://github.com/ewalsh/hmda/blob/master/py_code/config.py) file.
3. Finally, we will also fill in some data from Federal Reserve Bank of St. Louis. This data will come into use towards the end of the project as we begin to look across time periods as it should give us an indication of both the financial conditions and sentiment of the originating institutions.
 - a. Details about this API can be found [here](https://fred.stlouisfed.org/docs/api/fred/) (<https://fred.stlouisfed.org/docs/api/fred/>)

```
In [ ]: # after configuring the census API re-import config and check ap
        i key
        import config
        print(config.api_key)
```

Data Collection

Luckily for us, the designers of the APIs have made this pretty easy. A big thank you to them!

I have selected a single year and a single state. Feel free to change to your preferences, data is available from 2007 - 2017. Some important caveats. As this is an illustrative project only, there are some important details about availability of data (i.e. when it was published) and data type issues that would require more attention in a production environment.

```
In [4]: # choose a first state by two letter code and year
        init_state = "OH"
        init_yr = 2016
```

```
In [5]: # pull data from the HMDA database on mortgage approvals -- this  
        # may take awhile  
        data_lar = seed.lar_pull(init_state, init_yr)
```

```

-----
-----
HTTPError                                     Traceback (most recent
call last)
<ipython-input-5-ed1e5a5b96af> in <module>
      1 # pull data from the HMDA database on mortgage approvals
-- this may take awhile
----> 2 data_lar = seed.lar_pull(init_state, init_yr)

~\Projects\hmda\py_code\seed.py in lar_pull(state, yr)
      15     url_end = "&$limit=0&$offset=0"
      16     url_full = url_base + state + url_middle + str(yr) +
url_end
----> 17     data_lar = pd.read_csv(url_full, dtype=object)
      18     # dropping rows numbers
      19     ids = {}

~\AppData\Roaming\Python\Python37\site-packages\pandas\io\parser
s.py in parser_f(filepath_or_buffer, sep, delimiter, header, nam
es, index_col, usecols, squeeze, prefix, mangle_dupe_cols, dtyp
e, engine, converters, true_values, false_values, skipinitialspa
ce, skiprows, skipfooter, nrows, na_values, keep_default_na, na_
filter, verbose, skip_blank_lines, parse_dates, infer_datetime_f
ormat, keep_date_col, date_parser, dayfirst, iterator, chunksiz
e, compression, thousands, decimal, lineterminator, quotechar, q
uoting, doublequote, escapechar, comment, encoding, dialect, tup
leize_cols, error_bad_lines, warn_bad_lines, delim_whitespace, l
ow_memory, memory_map, float_precision)
      700         skip_blank_lines=skip_blank_lines)
      701
--> 702         return _read(filepath_or_buffer, kwds)
      703
      704     parser_f.__name__ = name

~\AppData\Roaming\Python\Python37\site-packages\pandas\io\parser
s.py in _read(filepath_or_buffer, kwds)
      411     compression = _infer_compression(filepath_or_buffer,
compression)
      412     filepath_or_buffer, _, compression, should_close = g
et_filepath_or_buffer(
--> 413         filepath_or_buffer, encoding, compression)
      414     kwds['compression'] = compression
      415

~\AppData\Roaming\Python\Python37\site-packages\pandas\io\commo
n.py in get_filepath_or_buffer(filepath_or_buffer, encoding, com
pression, mode)
      200
      201     if _is_url(filepath_or_buffer):
--> 202         req = _urlopen(filepath_or_buffer)
      203         content_encoding = req.headers.get('Content-Enco
ding', None)
      204         if content_encoding == 'gzip':

c:\python37\lib\urllib\request.py in urlopen(url, data, timeout,
cafile, capath, cadefault, context)
      220     else:
      221         opener = _opener
--> 222     return opener.open(url, data, timeout)

```

```

223
224 def install_opener(opener):

c:\python37\lib\urllib\request.py in open(self, fullurl, data, timeout)
529         for processor in self.process_response.get(protocol, []):
530             meth = getattr(processor, meth_name)
--> 531             response = meth(req, response)
532
533         return response

c:\python37\lib\urllib\request.py in http_response(self, request, response)
639         if not (200 <= code < 300):
640             response = self.parent.error(
--> 641                 'http', request, response, code, msg, hdrs)
642
643         return response

c:\python37\lib\urllib\request.py in error(self, proto, *args)
567         if http_err:
568             args = (dict, 'default', 'http_error_default') + orig_args
--> 569         return self._call_chain(*args)
570
571 # XXX probably also want an abstract factory that knows
when it makes

c:\python37\lib\urllib\request.py in _call_chain(self, chain, kind, meth_name, *args)
501         for handler in handlers:
502             func = getattr(handler, meth_name)
--> 503             result = func(*args)
504             if result is not None:
505                 return result

c:\python37\lib\urllib\request.py in http_error_default(self, req, fp, code, msg, hdrs)
647 class HTTPDefaultErrorHandler(BaseHandler):
648     def http_error_default(self, req, fp, code, msg, hdrs):
--> 649         raise HTTPError(req.full_url, code, msg, hdrs, fp)
650
651 class HTTPRedirectHandler(BaseHandler):

```

HTTPError: HTTP Error 410: Gone

```
In [ ]: # a quick data snapshot
data_lar.head()
```

```
In [ ]: '{:,.0f}'.format(data_lar.shape[0]) + ' total rows for a total o
f ' + \
'{:,.0f}'.format(data_lar.shape[0]*data_lar.shape[1]) + ' data p
oints'
```

```

In [ ]: # Now pull county business patterns data from the census bureau
        census_df = seed.census_pull(init_state, init_yr, data_lar, config.api_key)

In [ ]: # A quick snapshot
        census_df.head()

In [ ]: '{:,.0f}'.format(census_df.shape[0]) + ' total rows for a total
        of ' + \
        '{:,.0f}'.format(census_df.shape[0]*census_df.shape[1]) + ' data
        points'

In [ ]: # Finally, let's pull data about the originating institutions
        data_inst = seed.inst_pull(init_state, init_yr)

In [ ]: # A quick snapshot
        data_inst.head()

In [ ]: '{:,.0f}'.format(data_inst.shape[0]) + ' total rows for a total
        of ' + \
        '{:,.0f}'.format(data_inst.shape[0]*data_inst.shape[1]) + ' data
        points'

```

Data Storage

We are very early on, but we have already pulled together a rather large dataset. So, after building on our foundation of data collection, where we now have functions that can pull and update the necessary data, we need to start thinking about data storage and the tools we are going to need to scale up analysis.

For data storage this project will rely on PostgreSQL. This is a powerful, open source object-relational database which is well suited for the type of data we are working with.

Thinking about how we can go from one small slice of data to the entire dataset over all available years ahead of time will save us a lot of headache going forward. To help us in this endeavor, I will be using Docker to run these services. This setup will be helpful when we scale up operations.

If you have not already, please follow the instructions in the [README.md](#) (file:///./README.md).

```

In [ ]: # after the setup we are ready to load the data we have download
        ed into our database
        import load
        load.load()

```

Exploratory Data Analysis, Transformation, & Feature Engineering

Now that we have our data loaded, we can begin to do some initial analysis. As we step through the available data we should take care to think about how it can be interpreted and to make sure we set ourselves up for success by looking at abnormal features like outliers or highly skewed distributions.

We have done our job in getting the data together and stored. But now we need to make sure we have informative features. In my humble opinion, this is where a lot of value can be added, or conversely, where things can go wrong. Due to its importance, we will spend a fair amount of time in this section. Before we can even truly to any exploratory analysis we will need to preprocess a lot of this data so we can make better sense of it.

First and most importantly, if we take a look at the variable we will be trying to predict/classify we can it can take a few different forms.

```
In [ ]: df = pd.DataFrame(set(data_lar.loc[:, 'action_taken_name'])).rename(
columns={0: 'action_taken_name'})
df.style.set_properties(**{'width': '100%', 'text-align': 'center'})
```

We need to define what meets our criteria. For this project, we are interested in approvals. If the loan is originated that is an approval, but so is if the application was approved by it was withdrawn by the applicant. On the other hand, we will define application denied by the financial institution as unapproved but also if the preapproval request was denied by the institution.

This leaves several actions that we will filter out of our dataset. These will be:

1. When loans are purchased by other institutions
2. For incompleteness
3. When the applicant withdraws the application

```
In [ ]: data_lar.loc[:, 'action_taken'] = pd.to_numeric(data_lar.loc[:, 'action_taken'])
lar = data_lar[(data_lar.loc[:, 'action_taken'] != 6)]
lar = lar[(lar.loc[:, 'action_taken'] != 5)]
lar = lar[(lar.loc[:, 'action_taken'] != 4)]
"""dropped {:.0f} observations""".format(data_lar.shape[0] - lar.shape[0])
```

Another important filter we should impose is when income is missing from the dataset. Income is likely to be one of the most important features here, so if that isn't included in the data we shouldn't include it for now.

```
In [ ]: start_obs = lar.shape[0]
lar = lar[(lar.loc[:, 'applicant_income_000s'].apply(lambda x: math.isnan(float(x))) == False)]
"""dropped {:.0f} observations""".format(start_obs - lar.shape[0])
```


A brief note on dropping observations. Although beyond the scope of this project, it is often worth examining patterns within data where it is missing, where it is an outlier or otherwise strange. In this project we will be removing this data, but in practice just as having data can be informative where data is missing or especially if the data is an outlier, that fact can be informative as well.

Moving on for now, we need to transform our data from text to our binary classification. 1 for approved and 0 for not approved.

It will be important to note how balanced/unbalanced our data is, since a highly unbalanced set will mean an model who selects the more likely outcome will appear to be more correct than it really may be.

```
In [ ]: approved = np.zeros(lar.shape[0])
        ids = (lar.loc[:, 'action_taken'] == 1) | (lar.loc[:, 'action_taken'] == 2)
        approved[ids] = 1
        lar.loc[:, 'action_taken'] = approved
```

```
In [ ]: lar[['action_taken', 'action_taken_name']].head()
```

```
In [ ]: plt.figure(figsize=(8,6))
        plt.hist(lar['action_taken'], facecolor='#3F4B8C', alpha=0.8)
        plt.ylabel('Count')
        plt.xlabel('Probability')
        plt.title('Unbalance Approval/Non-Approval');
```

While our metric of interest is unbalanced, it is not extremely so.

I have already mentioned that I expect income to be an important feature of our data. However, it comes in a few different forms. Let's examine our income data and look at how we might transform them below.

```
In [ ]: lar.loc[:, 'applicant_income_000s'] = pd.to_numeric(lar.loc[:, 'applicant_income_000s'])
        plt.figure(figsize=(8,6))
        plt.hist(lar.loc[:, 'applicant_income_000s'], 100, facecolor='#3F4B8C', alpha=0.8)
        plt.ylabel('Count')
        plt.xlabel('Income')
        plt.title('Applicant Income (in thousands USD)');
```

We can see here that this data has a very log-normal looking distribution (i.e. very long right tail). Most models that we will consider have an underlying assumption that the data is normally distributed. It will be important for us to keep the underlying model assumptions in mind through our analysis process.

To address this in the income data, we will perform a log transformation.

```
In [ ]: plt.figure(figsize=(8,6))
plt.hist(lar.loc[:, 'applicant_income_000s'].apply(lambda x: mat
h.log(x,10)),100, \
         facecolor='#3F4B8C',alpha=0.8)
plt.ylabel('Count')
plt.xlabel('Income -- Log 10 Scale')
plt.title('Applicant Income (in thousands USD)');
```

First, you may notice that I chose a base 10 for the log transform. In this case I prefer to work on base 10 for interpretation as it tends to be easier to work in units of 10. Also notice that all of our values are above zero. We could scale this distribution to be above/below average being above and below zero. This again comes down to distribution. I don't expect above or below average income to be especially important for acceptance, rather it will be proportional to the loan amount. I will always default to less transformation where possible with a keen focus on interpretation.

You will also notice now that the distribution is much more normal looking, but we can also see a very strange blip around zero. Under closer inspection, these are entries of \$1 for applicant income. For our purposes, we will filter these outliers from our dataset as they are likely an entry or other type of input error.

```
In [ ]: start_obs = lar.shape[0]
lar = lar[(lar.loc[:, 'applicant_income_000s'] != 1)]
        """dropped {:,.0f} observations""".format(start_obs - lar.shape
[0])
```

```
In [ ]: """The Range of our data is from {} up to {:,.0f}""".format(min
(lar.loc[:, 'applicant_income_000s']), \
                                                                    max(lar.lo
c[:, 'applicant_income_000s']))
```

Looking at the new range of our data, we can still see an arbitrarily low number of 2 and an equally suspicious high number of 9999. While these are also likely to be entries which are less than accurate, just as in the case of missing data, they can still be informative and I have kept these in the dataset. For instance, perhaps they represent sentiment of the loan officer. In any case, looking at our new distribution they don't cause the same type of outlier problem we witnessed before

```
In [ ]: plt.figure(figsize=(8,6))
plt.hist(lar.loc[:, 'applicant_income_000s'].apply(lambda x: mat
h.log(x,10)),100, \
         facecolor='#3F4B8C',alpha=0.8)
plt.ylabel('Count')
plt.xlabel('Income -- Log 10 scale')
plt.title('Applicant Income (in thousands USD)');
```

Taking a look at another example variable, let's look at the number of units in an area that are built to house less than 5 families. This is a good metric, beyond population, for density of an area.

```
In [ ]: lar['number_of_1_to_4_family_units'] = pd.to_numeric(lar.loc[:, '
number_of_1_to_4_family_units'])
plt.figure(figsize=(8,6))
plt.hist(lar.loc[:, 'number_of_1_to_4_family_units'], 100, faceco
lor='#3F4B8C', alpha=0.8)
plt.ylabel('Count')
plt.xlabel('Living Density')
plt.title('Number of Under 5 Family Homes');
```

Again here we have a very strong right tailed distribution. However, in this case moving to a log scale would only shift that right tail to a left tail. Furthermore, looking at the Shapiro-Wilk test of normality the transformation would move the metric from 0.92 to 0.99. I would not consider that to be a large enough change to be compelling. A large reason for this are the outliers to the right. However, these are likely meaningful observations within cities. As such, it wouldn't make much sense to remove these outliers.

In this case, I don't think removing outliers makes sense but I do think scaling the data to be above or below average has more interpretive power. So this is the route I will take when scaling.

```
In [ ]: from sklearn import preprocessing
plt.figure(figsize=(8,6))
plt.hist(preprocessing.scale(lar.loc[:, 'number_of_1_to_4_family
units']), 100, facecolor='#3F4B8C', alpha=0.8)
plt.ylabel('Count')
plt.xlabel('Living Density -- Normalized')
plt.title('Number of Under 5 Family Homes');
```

Up to this point, we have looked at numerical data. However, a large part of this dataset is categorical. So let us take a quick look at some of those examples.

The first of these examples is whether or not the applicant was the sole applicant (i.e. had a co-applicant or not).

```
In [ ]: pd.DataFrame(set(lar.loc[:, 'co_applicant_ethnicity_name']), colum
ns=['Co-Applicant Race']).head()
```

```
In [ ]: soleApplicant = np.zeros(lar.shape[0])
lar.loc[:, 'co_applicant_ethnicity'] = pd.to_numeric(lar.loc[:, 'c
o_applicant_ethnicity'])
ids = (lar.loc[:, 'co_applicant_ethnicity'] == 5)
soleApplicant[ids] = 1
```

```
In [ ]: plt.figure(figsize=(8,6))
plt.hist(soleApplicant, 3, facecolor='#3F4B8C', alpha=0.8)
plt.ylabel('Count')
plt.title('Sole Applicant');
```

Our data also has some detailed information about the race of the applicant and co-applicant. In this example I won't prepare a binary variable for each race. This is not because they are not all important to look at, but rather the focus of this document is on process. So we will focus on the most common categories and leave a more in-depth study for later.

The groups we look at are:

1. Black or African American
2. Asian
3. All other Non-White races

Setting up our variables this way should give us a clear indication of any bias.

```
In [ ]: blackApplicant = np.zeros(lar.shape[0])
ids = (lar.loc[:, 'applicant_race_name_1'] == 'Black or African A
merican')
blackApplicant[ids] = 1

asianApplicant = np.zeros(lar.shape[0])
ids = (lar.loc[:, 'applicant_race_name_1'] == 'Asian')
asianApplicant[ids] = 1

otherRaceApplicant = np.ones(lar.shape[0])
ids = (lar.loc[:, 'applicant_race_name_1'] == 'White')
otherRaceApplicant[ids] = 0
ids = (lar.loc[:, 'applicant_race_name_1'] == 'Black or African A
merican')
otherRaceApplicant[ids] = 0
ids = (lar.loc[:, 'applicant_race_name_1'] == 'Asian')
otherRaceApplicant[ids] = 0

white_pct = sum(lar.loc[:, 'applicant_race_name_1'] == 'White')/l
ar.shape[0]

pd.DataFrame([{'black': round(sum(blackApplicant)/lar.shape[0],
2),
'asian': round(sum(asianApplicant)/lar.shape[0], 2),
'other race': round(sum(otherRaceApplicant)/lar.shape[0], 2),
'white': round(white_pct,2)}])
```

Another variable we will consider around race is that sometimes a non-white applicant could have a co-applicant that is white. Given race is an important consideration here for our regulatory purposes, we want to control for the mitigating effect that having a white co-applicant may have. We will call this variable 'white friend'.

```
In [ ]: white_friend = np.zeros(lar.shape[0])
ids = (lar.loc[:, 'applicant_race_name_1'] != 'White') & (lar.loc
[:, 'co_applicant_race_name_1'] == 'White')
white_friend[ids] = 1
"""{:,.0f} non-white applicants have a white co-applicant""".for
mat(sum(white_friend))
```

We have a rather extensive dataset, so I won't go through all of my logic on choosing, normalizing, and transforming the data here. I hope the above section gives you some insight into my thought process and you will be able to see my other choices in the relevant code.

Recall, that we are combining different datasets and we would like our process to be scalable. In order to make that possible we will create a series of functions that will:

1. Extract and merge our data from the database.
2. Do the preprocessing and make the necessary transformations

The first is merging our HMDA data with the data we gathered from the census bureau.

```
In [ ]: merge_query = {0: 'select lar.action_taken, lar.applicant_income_000s, lar.applicant_race_name_1, \
lar.co_applicant_race_name_1, lar.applicant_sex, lar.lien_status, lar.loan_purpose, lar.loan_type, \
lar.agency_code, lar.owner_occupancy, lar.preapproval, lar.property_type, lar.purchaser_type_name, \
lar.hud_median_family_income, lar.loan_amount_000s, lar.number_of_1_to_4_family_units, \
lar.number_of_owner_occupied_units, \
lar.minority_population, lar.population, lar.tract_to_msamd_income, census.emp, census.estab, census.payann, \
census.pop, census.county_code, census.county_name, census.state_code from lar LEFT JOIN census ON \
lar.county_code = census.county_code AND lar.state_code = 39'}
```

```
In [ ]: from queries import get_query
from preproc import trans_actions
```

```
In [ ]: merged_data = get_query(merge_query[0])
```

```
In [ ]: merged_data.shape
```

```
In [ ]: # preprocessing the data
input_data = trans_actions(merged_data)
```

What we have just done above is combine two earlier steps in a scalable process. As we add data to the data base, either as it comes in or as we collect it, we are able to merge that data with supporting information, clean it, and also do transformations that get it ready to use. In terms of process I believe that this should be the type of goal folks should focus on.

Initial Modelling and Learning

Having all of our data in order, we can begin to look at our features and how they are related to our metric of interest (i.e. approval of a mortgage)

```
In [ ]: # some quick descriptors
input_data.describe()
```

```
In [ ]: # Correlations show where the strongest relationships exist
pd.options.display.float_format = '{:,.3F}'.format
input_data.corr()
```

```
In [ ]: # Easier to focus on a grouping, in this case income related met
rics
plt.figure(figsize=(16,10))
sns.set(font_scale=2)
sns.heatmap(input_data.loc[:,['approved','income_log','hud_sprea
d','inc_loan_ratio','local_income_ratio','pay_pop']].corr(), annot
ot = True, cmap=sns.color_palette("hls", 8))
plt.show()
```

In the graphic above, we can see some strong relationships. Both income and the spread of a person's income over the hud median family income are the same. In fact, on closer inspection these two metrics are 98% correlated. This could cause issues for us and our modelling. It would be safest for us to choose one or the other or dig deeper on how to get more information from our spread metric and remove the overall income effect from it (i.e. a form of dimension reduction perhaps?)

Below, we switch from income to ethnicity and gender. These are low metrics, but they do go along with pre-existing bias we may have had. It would be interesting to look at individual institutions to see if this is more prevalent in a few bad actors or a systemic issue.

```
In [ ]: plt.figure(figsize=(16,10))
sns.set(font_scale=2)
sns.heatmap(input_data.loc[:,['approved','sole_applicant','black
_applicant','asian_applicant','white_friend','is_female']].corr
(), annot = True, cmap=sns.color_palette("hls", 8))
plt.show()
```

Modelling

Our first model we will turn to the workhorse of classification, the Logistic Regression. This model often works extremely well for these types of problems. In essence, this model is link to a function which quickly goes from 0 to 1 which is well suited to the type of binary classification we are attempting here. I would normally use a standard model like this for benchmarking any further machine learning or other approaches.

Looking at our initial output we see a few interesting things. First, as confirmed in our earlier analysis looking at income only, you are more than likely to receive a mortgage. Given the output below a person in the bottom 25th percentil would break even with the intercept. So, while income is important it isn't an especially high bar. That said, the botton 25% does make up a substantial proportion of the population.

Moving on to our ethnicity and race metrics we do see being African American doesn't help. Unfortunately, with a p-value which is very strong. The effects are not as convincing for being Asian or female, but they are equally consistent in p-value terms with other race applicants with those of African Americans.

Lots of interesting results here to look through here.

```

call:
glm(formula = approved ~ incomeLog + soleApplicant + blackApplicant +
  asianApplicant + otherRaceApplicant + isFemale + firstLien +
  refinancing + homeImprovement + isHUD + isCreditUnion + isOwnerOccupied +
  isManufactured + hudsSpreadLogNormalized + incomeLoanRatioLog +
  lowDenseArea + selfOwnedArea + minorityPop + totalAreaPopLog +
  empPerPop + payPerPop + popLogNormalized, family = binomial(link = "logit
"),
  data = subData)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.8162   0.2560   0.4845   0.7483   2.0428

Coefficients:
                Estimate Std. Error z value Pr(>|z|)
(Intercept)      -3.70170    0.76640  -4.830 1.37e-06 ***
incomeLog         2.24331    0.28831   7.781 7.20e-15 ***
soleApplicant    -0.13754    0.03091  -4.450 8.60e-06 ***
blackApplicant   -0.70510    0.05853 -12.048 < 2e-16 ***
asianApplicant   -0.07319    0.10135  -0.722  0.47024
otherRaceApplicant -0.45625    0.05058  -9.021 < 2e-16 ***
isFemale         0.02997    0.02616   1.146  0.25197
firstLien        0.47198    0.08185   5.766 8.10e-09 ***
refinancing     -1.56475    0.03212 -48.715 < 2e-16 ***
homeImprovement -1.57521    0.06368 -24.736 < 2e-16 ***
isHUD           -0.40891    0.03042 -13.444 < 2e-16 ***
isCreditUnion   0.57124    0.06605   8.649 < 2e-16 ***
isOwnerOccupied  0.38625    0.06075   6.358 2.04e-10 ***
isManufacturedTRUE -0.93159    0.09887  -9.422 < 2e-16 ***
hudsSpreadLogNormalized -0.30589    0.10165  -3.009  0.00262 **
incomeLoanRatioLog 0.09153    0.06133   1.492  0.13561
lowDenseArea    -0.25449    0.04523  -5.626 1.84e-08 ***
selfOwnedArea    0.26670    0.04497   5.931 3.02e-09 ***
minorityPop     -0.02695    0.02022  -1.333  0.18254
totalAreaPopLog  0.35979    0.17355   2.073  0.03816 *
empPerPop       0.18467    0.07155   2.581  0.00985 **
payPerPop       -0.08865    0.05677  -1.562  0.11835
popLogNormalized -0.08496    0.02715  -3.130  0.00175 **
---
signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

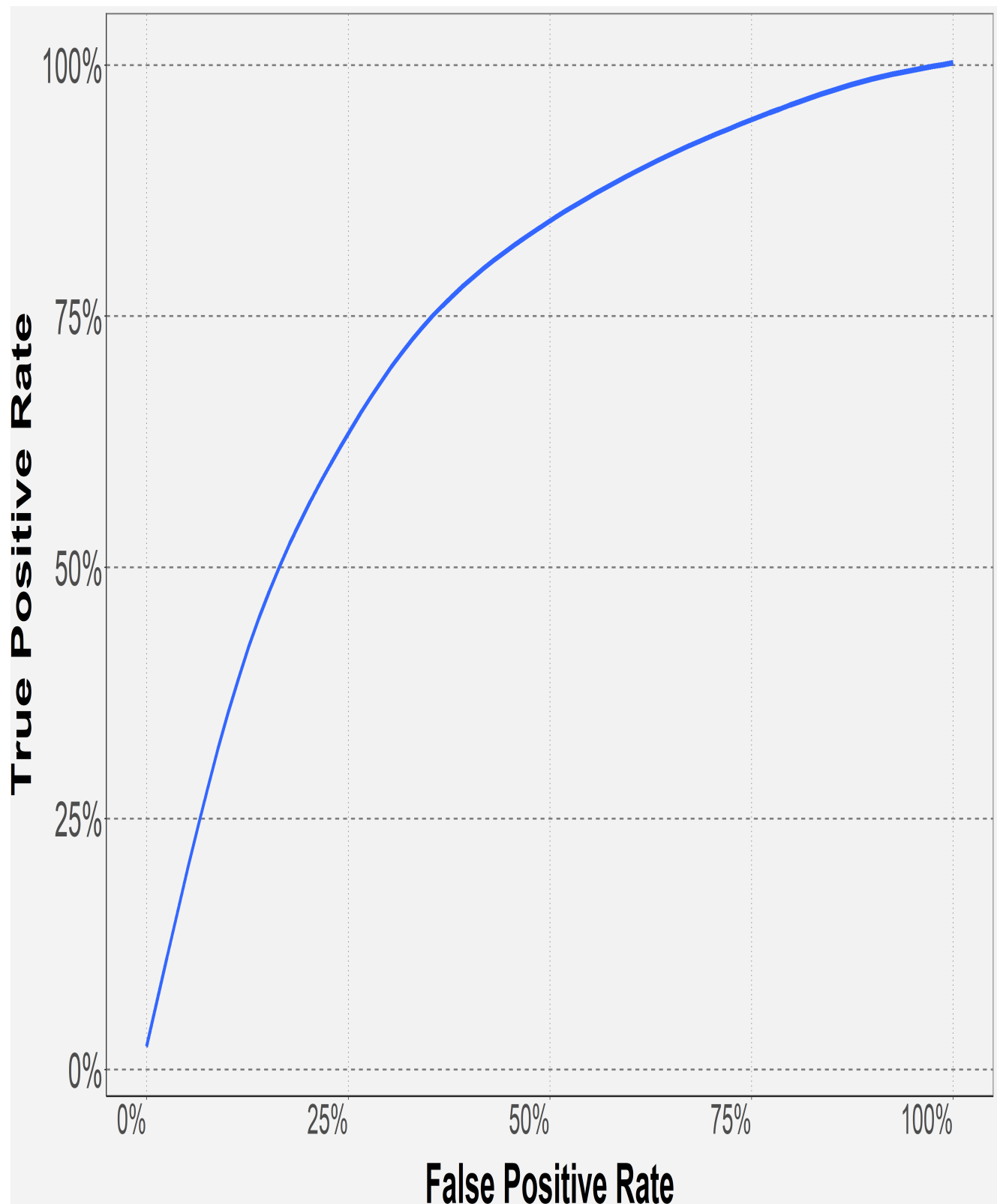
(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 37498  on 33890  degrees of freedom
Residual deviance: 32339  on 33868  degrees of freedom
AIC: 32385

Number of Fisher Scoring iterations: 5

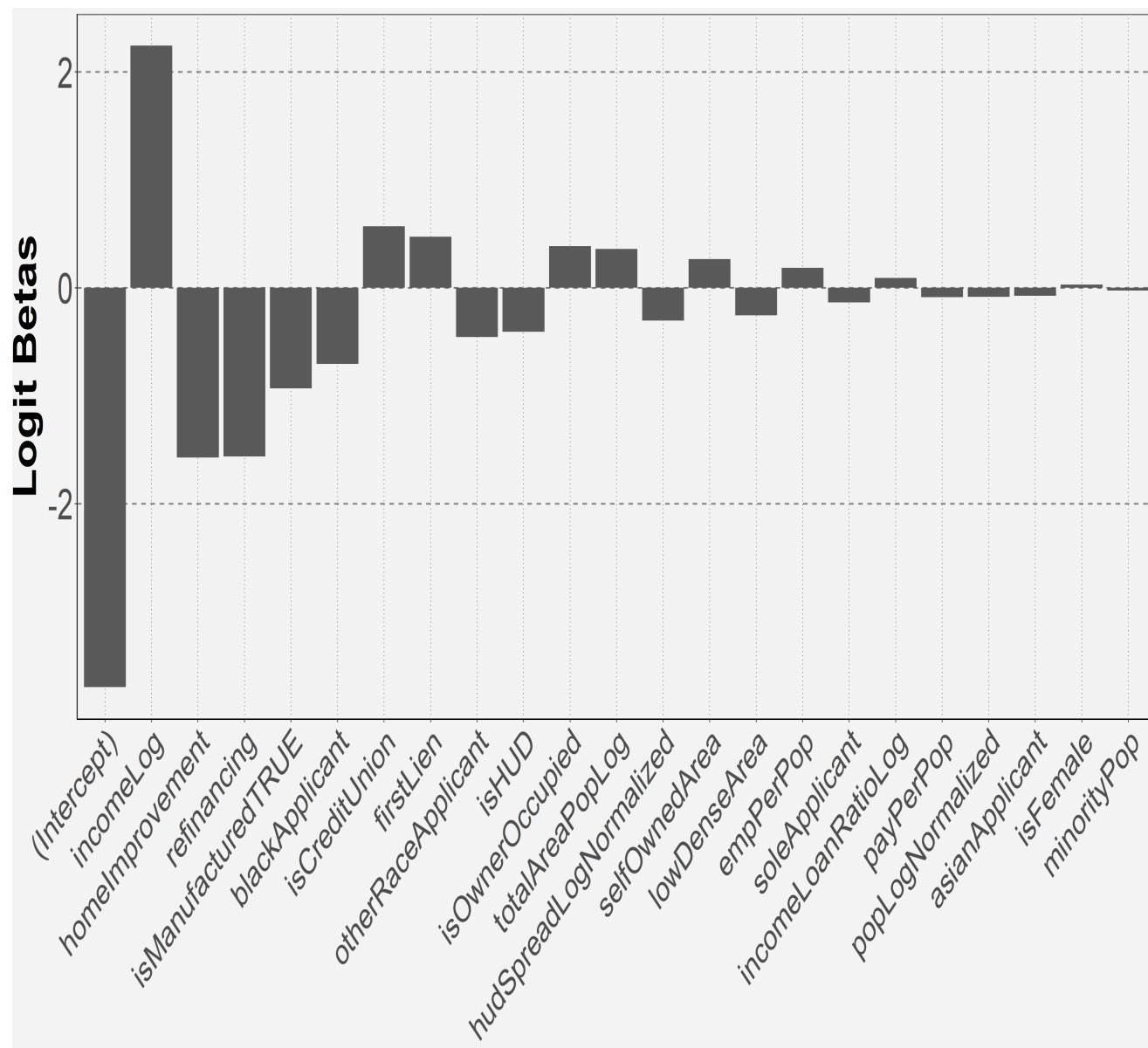
```

As we have an unbalanced dataset, it will be important to think about our false positives as well as true positives. Confusion matrices are a good way to look at this information, but being more graphical, I tend to prefer area under the curve graphics. We can see this for our Logistic Regression below.



This is a pretty good first run. The total Area Under the Curve metric for this model is roughly 75%.

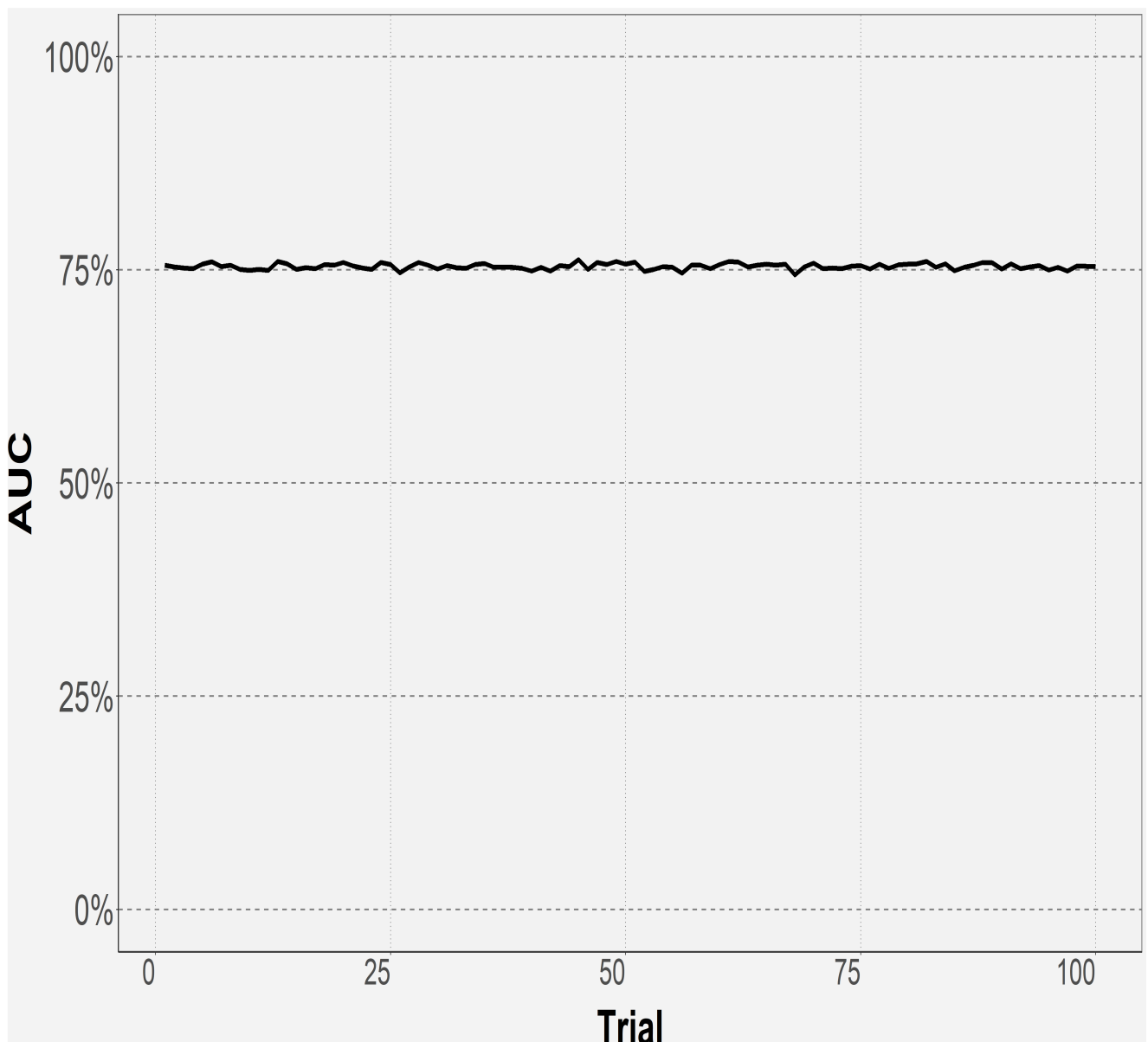
We have included a lot of variables in this first model. This leaves us open to bias. Normally, there is a trade-off between variance around the model and bias which lead to under or over fitting. There are many ways to think about combating this but a first look would just be looking at the loads of the model.



We have looked at this information from the model before, but this gives us a sense of feature importance (although I would adjust to variance of the underlying metrics) as well as a sense that no single feature dominates the model. Income, by far the most important isn't double the next few metrics and we did expect it to be strong to begin with

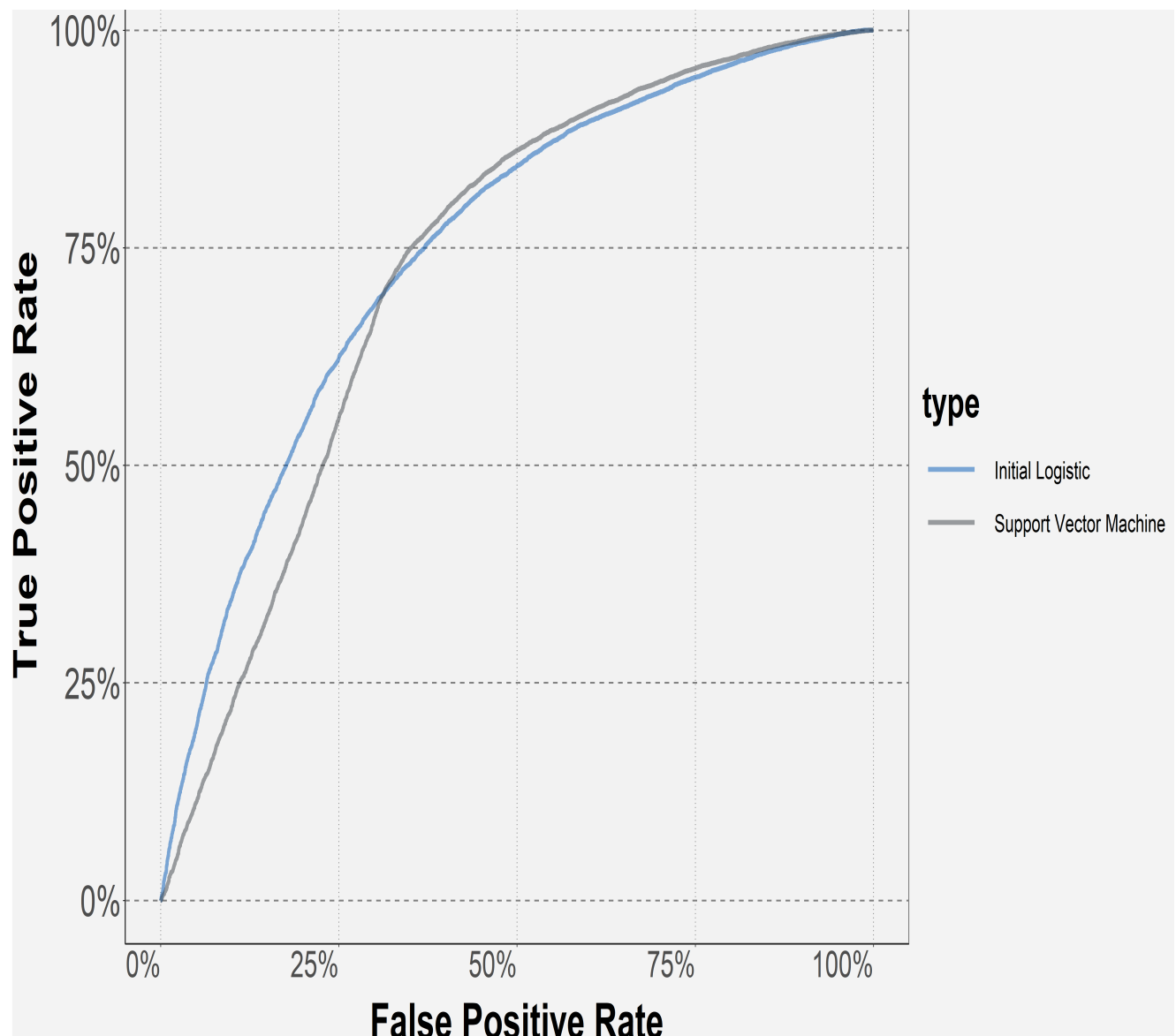
Another important technique we should exploit is cross-validation. This is closely tied to the idea of a learning curve where as we begin to feed more data into the model how much more accurate does it become, or perhaps, does the prediction deteriorate out of sample due to overconfidence.

We are only looking at a single state and single year for now and then randomizing the observations within a training and testing sample. I would expect this to become more difficult as regional difference become more important. For now, let's just see how the area under the curve metrics changes over different trials.



For the moment, our logistic regression seems to be very consistent. Again, a good start made possible by the foundations we set earlier.

The logistic regression is a workhorse function because it works very well. We are trying to fit a line to best predict a 1 or 0. However, what if we flip our approach and try to fit a line that splits our 1s from our 0s. This is the approach of a Support Vector Machine (SVM).



We can see here that the SVM has more false positive than the logistic regression.

Optimize

These models have shown themselves to be a good place to start, but as we get more into our data and modelling we may be able to exploit parameter tuning, optimization techniques, regularization, or deep learning. These models should be used with care, but can be very powerful. For an example, I will use a Elastic Net model below which is basically a step further on our Logistic Regression. Elastic Net is a form of Generalized Additive Model and can be useful for feature selection as well as stability through regularization of the L1 and L2 norms, which effectively combat outliers and bias. We will simply split our data to run our test.

```
In [ ]: xcols = input_data.columns
xcols = xcols[1:xcols.shape[0]]
X_train = input_data.loc[0:3000,xcols]
y_train = input_data.loc[0:3000,'approved']
X_test = input_data.loc[3001:5000,xcols]
y_test = input_data.loc[3001:5000,'approved']
```

```

In [ ]: from sklearn.linear_model import SGDClassifier
        from sklearn.model_selection import cross_val_predict
        from sklearn.metrics import confusion_matrix
        from sklearn.metrics import precision_score, recall_score
        from sklearn.metrics import roc_curve
        from sklearn.metrics import plot_roc_curve

In [ ]: clf = SGDClassifier(random_state = 0, loss='log', penalty='elasticnet')
        clf.fit(X_train, y_train)

In [ ]: y_train_pred = cross_val_predict(clf, X_train, y_train, cv=100)

In [ ]: tn, fp, fn, tp = confusion_matrix(y_train, y_train_pred).ravel()

In [ ]: # The Negative Class
        """"{} True Negatives and {} False Positives"".format(tn,fp)

In [ ]: # The Positive Class
        """"{} False negatives and {} True Positives"".format(fn,tp)

In [ ]: """"{:,.2F} Precision"".format(precision_score(y_train,y_train_pred))

In [ ]: """"{:,.2F} Recall"".format(recall_score(y_train,y_train_pred))

In [ ]: plt.figure(figsize=(50,50))
        plot_roc_curve(clf, X_test, y_test)
        plt.show()

```

We can see that this learning model provides substational improvement, moving our AUC measure from around 0.75 to 0.93. This is with minimal tuning, but it is still worth being wary as we have yet to scale up the process.

Conclusion

In this project we have focused on building a scalable process for analysis. I believe we have successfully met those goals as we have a data ingestion pipeline as well as cleaning and transformation that lead to an analytics process that while early stages, shows substation potential.

Keeping all of this in the context we set, with an AUC measure of 93%, this would be a high bar for any human to beat and so we could provide cost savings to mortgage lenders and/or help them make their lending deicions. We have also highlighted some suspect race and gender issues that could become regulatory issues in the future. It is worth further study here.

Finally, we have brought together many dataset relevant to the mortgage decision and financial instituions. This has set a strong foundation for us going forward and combining with our analytics we can see which institutions may need the most guidance/help or perhaps may just be unaware that they could lend more/less aggressively and with less bias.