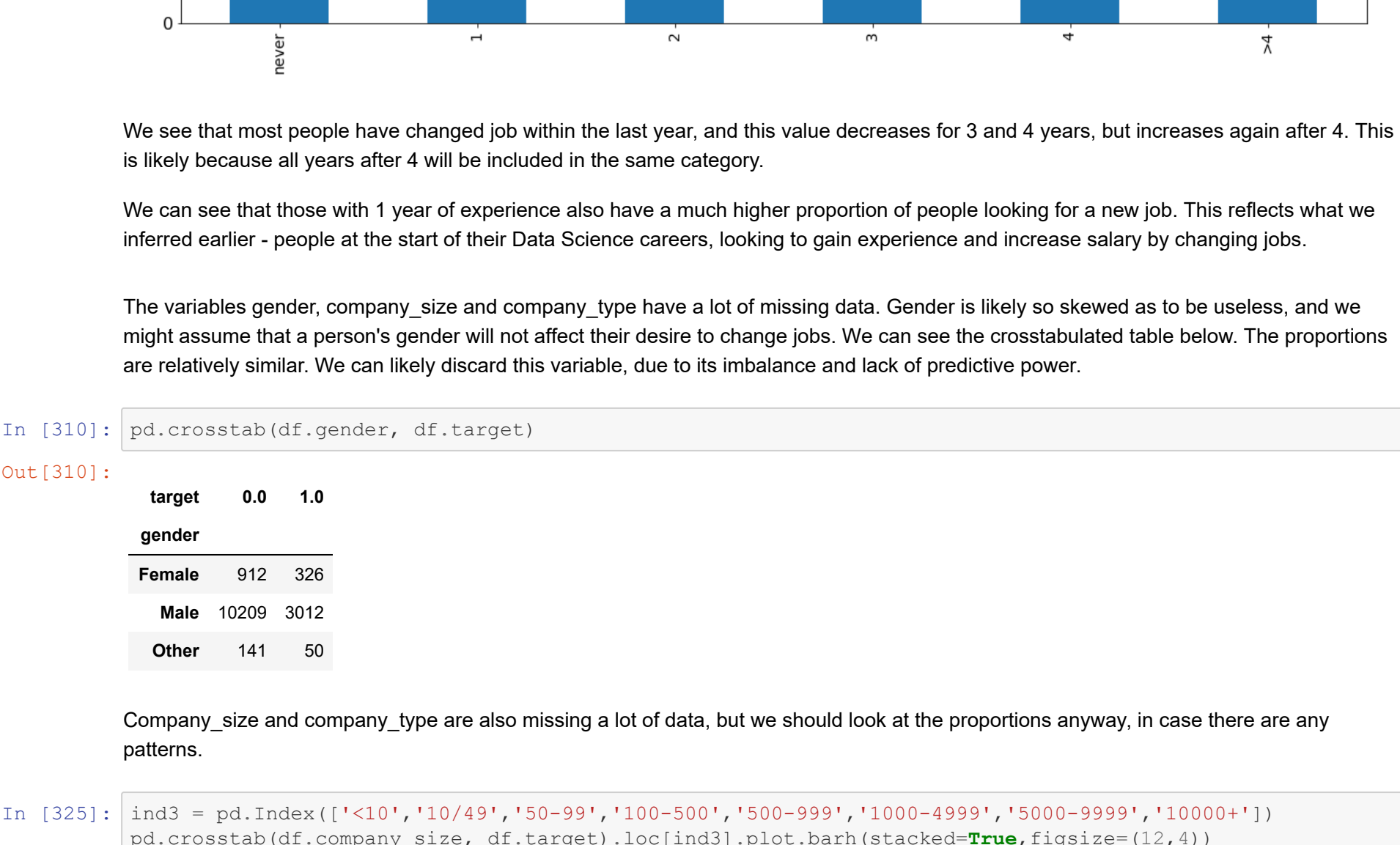


In [309]: # set index for target readability
ind2 = pd.Index(['never', '1', '2', '3', '4', '>4'])
years_since_job_change = df.last_new_job.value_counts()
put index in order
years_since_job_change = years_since_job_change.loc[ind2]
print out table of values
print(years_since_job_change)

plot distribution
years_since_job_change.plot.bar(figsize=(16,4))
plt.title('Years since last job change', fontsize=16)
cross tabulate experience with target variable
job_change_by_target = pd.crosstab(df.last_new_job, df.target)
sort, as above
job_change_by_target = job_change_by_target.loc[ind2]
plot as stacked bar plot
job_change_by_target.plot.bar(stacked=True, figsize=(20,6))
set legend
plt.legend(labels=my_labels, fontsize=16)
set title
plt.title('Years Since Last Job Change - Coloured by Looking For versus Not Looking For Job', fontsize=16):



We see that most people have changed job within the last year, and this value decreases for 3 and 4 years, but increases again after 4. This is likely because all years after 4 will be included in the same category.

We can see that those with 1 year of experience also have a much higher proportion of people looking for a new job. This reflects what we inferred earlier - people at the start of their Data Science careers, looking to gain experience and increase salary by changing jobs.

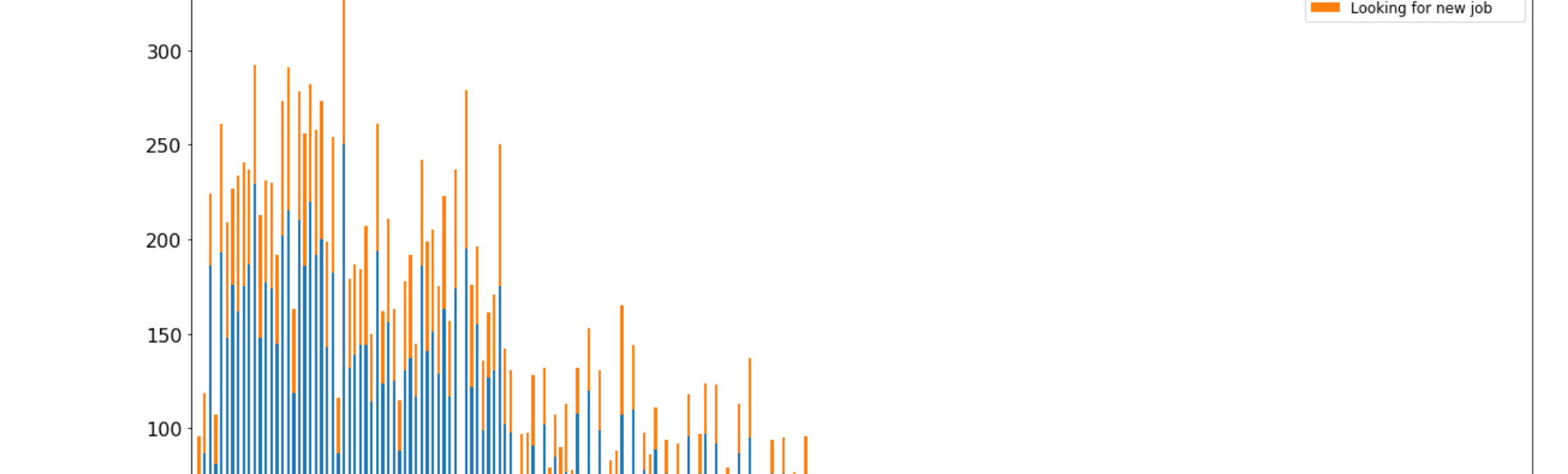
The variables gender, company_size and company_type have a lot of missing data. Gender is likely so skewed as to be useless, and we might assume that a person's gender will not affect their desire to change jobs. We can see the crosstab table below. The proportions are relatively similar. We can likely discard this variable, due to its imbalance and lack of predictive power.

In [310]: pd.crosstab(df.gender, df.target)
Out[310]:

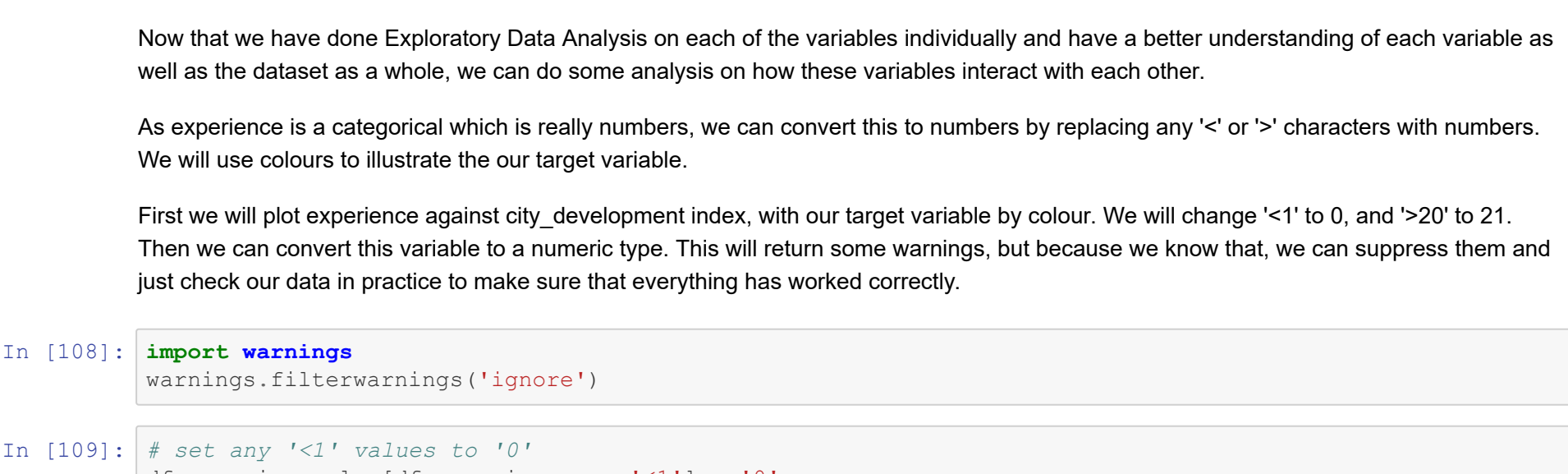
gender	0	1
Female	912	326
Male	10209	3012
Other	141	0

Company_size and company_type are also missing a lot of data, but we should look at the proportions anyway, in case there are any patterns:

In [325]: ind3 = pd.Index(['<10', '10/49', '50-99', '100-500', '500-999', '1000-4999', '5000-9999', '10000+'])
pd.crosstab(df.company_size, df.target).loc[ind3].plot.bar(stacked=True, figsize=(12,4))
plt.title('Company Size, Coloured by Target Variable')
plt.legend(labels = my_labels):



In [324]: pd.crosstab(df.company_type, df.target).plot.barh(stacked=True, figsize=(12,4))
plt.legend(labels=my_labels):



It seems that there is a larger proportion of people looking for a new job in companies of between 50 and 500 employees, with a similar proportion reflected in those companies over 10000 people. The relatively small proportion in companies of between 500 and 9999 people might be random chance, or due to missing values (i.e. a data entry error or some sort of error with the data file).

As we are aiming to predict retention for our company, and since we would know how many employees we have, we could perhaps use this data on an ad-hoc basis, but we shouldn't be certain of its predictive power just yet.

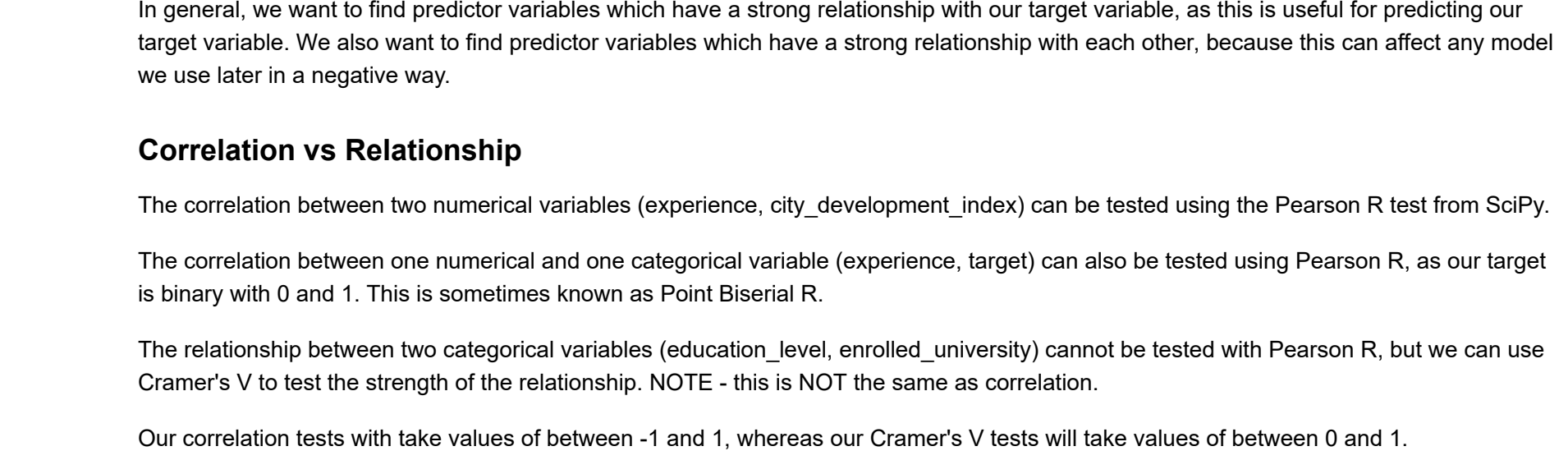
As for company type, we have the same issue of missing data. We might assume that those people working for NGOs (Non-Governmental Organisations) or Startups might also have reasons for not wanting to change job, or indeed, to admit to it. Therefore, with all these things in mind, we should be careful of making any assumptions or drawing any inferences on this variable alone.

Our final potential predictor variable is training_hours. We will look at its distribution, and its proportions for people looking for a new job. It is a numerical variable, and so has a lot of values. We could put these into categories, but it is not necessary at this point, and might be better left as a numerical variable for our predictive model.

In [362]: df.training_hours.hist(bins=100, figsize=(12,4))
plt.title('Histogram of Training Hours'):



In [366]: plt.corr('twick', labelsize=10)
pd.crosstab(df.training_hours, df.target).plot.bar(stacked=True, figsize=(20,10))
plt.title('Training Hours, Coloured by Target Variable', fontsize=16)
plt.xticks([1], fontsize=16)
plt.legend(labels=my_labels, fontsize=12):



We can see that the proportion of people looking for a new job decreases as training_hours increases. As this is a numerical variable with a large range of observed values, it could be useful to include in our predictive model.

Finally, a word on our response variable, df.target. We see the value counts for that variable below. This corresponds to about 75% of people in our dataset not looking for a job, and about 25% of people looking for a job.

In [368]: df.target.value_counts()
Out[368]:

0.0	14381
1.0	4777

Name: target, dtype: int64

In []:

Correlations and Scatterplots

Now that we have done Exploratory Data Analysis on each of the variables individually and have a better understanding of each variable as well as the dataset as a whole, we can do some analysis on how these variables interact with each other.

As experience is a categorical which is really numbers, we can convert this to numbers by replacing any '<' or '>' characters with numbers. We will use colours to illustrate the our target variable.

First we will plot experience against city_development_index, with our target variable by colour. We will change '<' to 0, and '>20' to 21. Then we can convert this variable to a numeric type. This will return some warnings, but because we know that, we can suppress them and just check our data in practice to make sure that everything has worked correctly.

In [108]: import warnings
warnings.filterwarnings('ignore')

In [109]: # set any '<' values to '0'
df.experience.loc[df.experience == '<20'] = '0'
set any '>20' values to '21'
df.experience.loc[df.experience == '>20'] = '21'
df.experience = df.experience.astype('int32')

Out[109]:

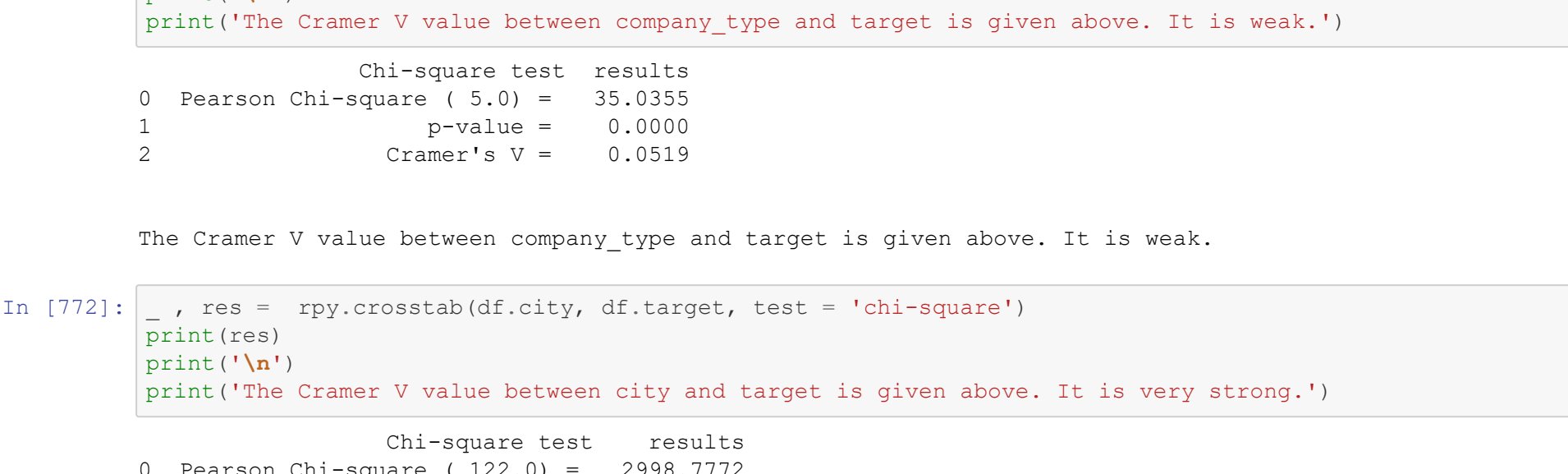
0	21
1	13
2	5
3	0
4	21
...	...
19153	14
19154	14
19155	21
19156	0
19157	2

Name: experience, Length: 19158, dtype: object

In [137]: # convert experience to numbers, making sure to drop NaN values, as these can't be converted
df.experience = df.experience.dropna().astype('int32')

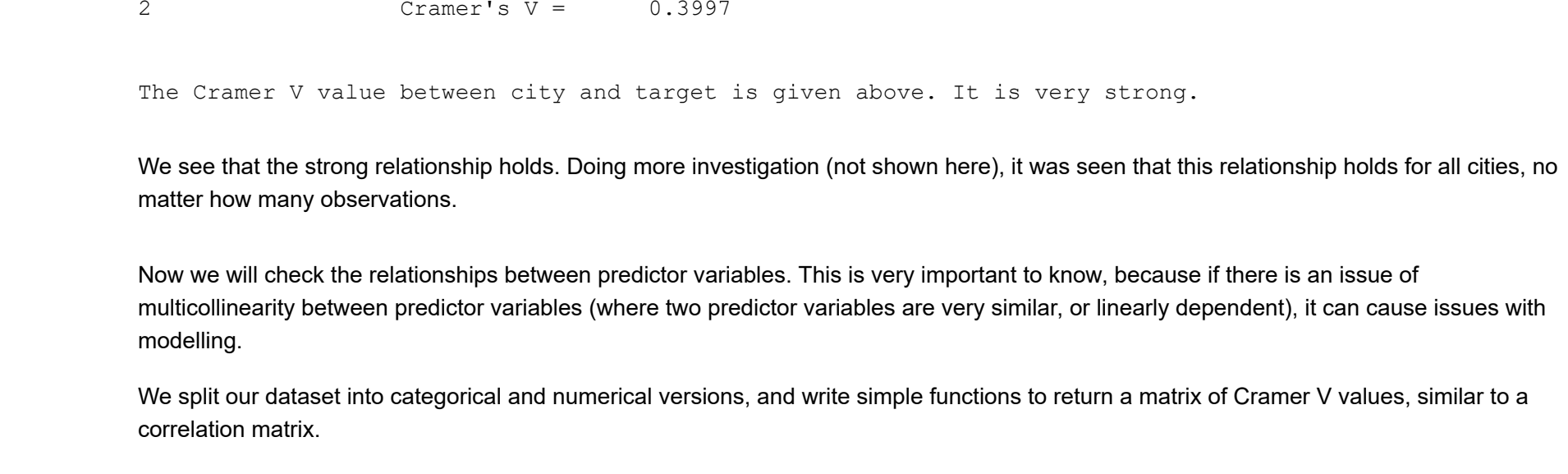
We're aiming to create plots where we can see some correlation pattern in the interaction of variables, and how our target variable changes depending on these. However, our two plots below are not so illustrative. In the first plot, we can see that there are more people looking for a new job (magenta dots) in cities with a lower development index (bottom left corner), especially with less years of experience. Moving towards the right, we see many more people who are not looking for a new job (black dots), as their experience increases, and especially as city_development_index increases.

In [742]: # plot experience against city_development_index, colour by target variable
plt.figure(figsize=(20,6))
sns.set_style('darkgrid')
sns.stripplot('city_development_index', data=df, jitter=0.4, hue='target', palette=['black', 'magenta'], alpha=0.5)
plt.xlabel('Experience', fontsize=16)
plt.ylabel('City Development Index', fontsize=16)
plt.title('Stripplot - Experience vs City Development Index', fontsize=16):



Below we have a similar plot of training_hours against last_new_job. This plot gives us some insight, but like with the plot above, there are so many data points that it is difficult to gain any real insight into correlation.

In [741]: plt.figure(figsize=(20,6))
sns.set_style('darkgrid')
sns.stripplot('training_hours', 'last_new_job', data=df, jitter=0.4, hue='target', palette=['black', 'magenta'], alpha=0.5)
plt.xlabel('Training hours', fontsize=16)
plt.ylabel('Years since last new job', fontsize=16)
plt.title('Stripplot - Training Hours versus Last New Job', fontsize=16):



These correlation-style scatterplots are not giving us much more insight. Instead, we can use some statistical tests to check the correlation and interaction between variables.

Statistical Tests

When checking interactions between variables, we have a few things to consider. We want to know:

- what predictor variables have a strong relationship with our target variable?
- what predictor variables have a strong relationship with each other?
- how can we quantify this relationship numerically?

In general, we want to find predictor variables which have a strong relationship with our target variable, as this is useful for predicting our target variable. We also want to find predictor variables which have a strong relationship with each other, because this can affect any model we use later in a negative way.

Correlation vs Relationship

The correlation between two numerical variables (experience, city_development_index) can be tested using Pearson R from SciPy.

The relationship between one numerical and one categorical variable (experience, target) can also be tested using Pearson R, as our target is binary with 0 and 1. This is sometimes known as Point Biserial R.

The relationship between two categorical variables (education_level, enrolled_university) cannot be tested using Pearson R, but we can use Cramer's V to test the strength of the relationship. NOTE - this is NOT the same as correlation.

Our correlation tests with take values of between -1 and 1, whereas our Cramer's V tests will take values of between 0 and 1.

For correlation, we have the following (generalised) interpretations:

- -1 is perfect negative correlation, 1 is perfect positive correlation (only happens when variables are identical)
- 0 is no correlation
- Between -0.3 and 0 or 0 and 0.3 is weak correlation (negative or positive)
- Between -0.6 and -0.3, or 0.3 and 0.6 is moderate correlation (negative or positive)
- Between -0.8 and -0.6, or 0.6 and 0.8, is strong correlation (negative or positive)
- Below -0.8 or above 0.8 is very strong correlation (negative or positive)

For Cramer's V, we have different values for our strength relationships:

- >0.25 is a very strong relationship
- >0.2 is a strong relationship
- >0.15 is a moderate relationship
- >0.10 is a weak-moderate relationship
- >0.05 is a weak relationship
- 0 is no relationship

In [117]: # import SciPy for statistical testing
import scipy.stats as stats

First we perform Pearson R tests on our target variable against numerical and categorical variables.

In [745]: corr, pval = stats.pearsonr(df.target, df.city_development_index)
print('The correlation coefficient for city_development_index vs target is: ', corr)
print('With a p-value of ', pval)
print('This means that there is a moderate negative correlation between these two variables, and it is statistically significant.')

The correlation coefficient for city_development_index vs target is: -0.3416505854522934
With a p-value of: 0.0
This means that there is a moderate negative correlation between these two variables, and it is statistically significant

In [748]: corr, pval = stats.pearsonr(df.target, df.training_hours)
print('The correlation coefficient for training_hours vs target is: ', corr)
print('With a p-value of ', pval)
print('This means that there is a very weak negative correlation between these two variables, and it is statistically significant.')

The correlation coefficient for training_hours vs target is: -0.02157724971361325
With a p-value of: 0.002819949452669254
This means that there is a very weak negative correlation between these two variables, and it is statistically significant.

In [749]: # experience is numerical now - converted it above
corr, pval = stats.pearsonr(df.target.loc[df.experience.dropna().index], df.experience.dropna())
print('The correlation coefficient for experience vs target is: ', corr)
print('With a p-value of ', pval)
print('This means that there is a weak negative correlation between these two variables, and it is statistically significant.')

The correlation coefficient for experience vs target is: -0.17672375550228786
With a p-value of: 0.49588976828782-134
This means that there is a weak negative correlation between these two variables, and it is statistically significant

So we see that the relationship of all three numerical variables to our target variable is negative. This means that, as the values of these variables increase, the chance of the person looking for a job decreases.

Now we will move on to testing the relationship between target and our categorical variables.

In [757]: _, res = rpy.crosstab(df.education_level, df.target, test = 'chi-square')
print(res)
print('\n')
print('The Cramer V value between education_level and target is given above. It is weak-moderate.')

Chi-square test results
0 Pearson Chi-square (4.0) = 165.6554
1 p-value = 0.0000
2 Cramer's V = 0.0941

The Cramer V value between education_level and target is given above. It is weak-moderate.

In [759]: _, res = rpy.crosstab(df.major_discipline, df.target, test = 'chi-square')
print(res)
print('\n')
print('The Cramer V value between major_discipline and target is given above. It is very weak.')

Chi-square test results
0 Pearson Chi-square (5.0) = 12.2071
1 p-value = 0.0321
2 Cramer's V = 0.0273

The Cramer V value between major_discipline and target is given above. It is very weak.

In [761]: _, res = rpy.crosstab(df.enrolled_university, df.target, test = 'chi-square')
print(res)
print('\n')
print('The Cramer V value between enrolled_university and target is given above. It is strong.')

Chi-square test results
0 Pearson Chi-square (2.0) = 455.1668
1 p-value = 0.0000
2 Cramer's V = 0.1357

The Cramer V value between enrolled_university and target is given above. It is strong.

In [763]: _, res = rpy.crosstab(df.relevant_experience, df.target, test = 'chi-square')
print(res)
print('\n')
print('The Cramer V value between relevant_experience and target is given above. It is moderate.')

Chi-square test results
0 Pearson Chi-square (1.0) = 315.9933
1 p-value = 0.0000
2 Cramer's phi = 0.1284

The Cramer V value between relevant_experience and target is given above. It is moderate.

In [764]: _, res = rpy.crosstab(df.gender, df.target, test = 'chi-square')
print(res)
print('\n')
print('The Cramer V value between gender and target is given above. It is very weak.')

Chi-square test results
0 Pearson Chi-square (2.0) = 9.0422
1 p-value = 0.0109
2 Cramer's V = 0.0248

The Cramer V value between gender and target is given above. It is very weak.

In [766]: _, res = rpy.crosstab(df.last_new_job, df.target, test = 'chi-square')
print(res)
print('\n')
print('The Cramer V value between last_new_job and target is given above. It is weak.')

Chi-square test results
0 Pearson Chi-square (5.0) = 132.4955
1 p-value = 0.0000
2 Cramer's V = 0.0841

The Cramer V value between last_new_job and target is given above. It is weak.

Now we test those variables which are missing a lot of data, or have very skewed distributions. We must bear that in mind when testing the relationship, and when using them in any models.

In [768]: _, res = rpy.crosstab(df.company_size, df.target, test = 'chi-square')
print(res)
print('\n')
print('The Cramer V value between company_size and target is given above. It is weak.')

Chi-square test results
0 Pearson Chi-square (7.0) = 45.5318
1 p-value = 0.0000
2 Cramer's V = 0.0587

The Cramer V value between company_size and target is given above. It is weak.

In [770]: _, res = rpy.crosstab(df.company_type, df.target, test = 'chi-square')
print(res)
print('\n')
print('The Cramer V value between company_type and target is given above. It is weak.')

Chi-square test results
0 Pearson Chi-square (5.0) = 35.0355
1 p-value = 0.0000
2 Cramer's V = 0.0519

The Cramer V value between company_type and target is given above. It is weak.

In [772]: _, res = rpy.crosstab(df.city, df.target, test = 'chi-square')
print(res)
print('\n')
print('The Cramer V value between city and target is given above. It is very strong.')

Chi-square test results
0 Pearson Chi-square (122.0) = 2998.7752
1 p-value = 0.0000
2 Cramer's V = 0.3956

The Cramer V value between city and target is given above. It is very strong.

We get a very strong relationship between city and target. Remember from above that the distribution of city is very imbalanced. We should restrict ourselves to cities for which we have a good number of observations, and test this relationship again. We will select our cutoff to be 60 observations.

In [773]: # get value counts of cities, and get index, to get those cities for which we have > 60 observations
df.index(keys=['enrolled_university'])
subset our dataframe by taking only observations whose cities are in our index list
df_cities_over_60 = df[df.index.isin(cities_over_60s)]

In [775]: _, res = rpy.crosstab(df_cities_over_60.city, df_cities_over_60.target, test = 'chi-square')
print(res)
print('\n')
print('The Cramer V value between city and target is given above. It is very strong.')

Chi-square test results
0 Pearson Chi-square (48.0) = 2834.8022
1 p-value = 0.0000
2 Cramer's V = 0.3997

The Cramer V value between city and target is given above. It is very strong.

We see that the strong relationship holds. Doing more investigation (not shown here), it was seen that this relationship holds for all cities, not matter how many observations.

Now we will check the relationships between predictor variables. This is very important to know, because if there is an issue of multicollinearity between predictor variables (where two predictor variables are very similar, or linearly dependent), it can cause issues with modelling.

We split our dataset into categorical and numerical versions, and write simple functions to return a matrix of Cramer V values, similar to a correlation matrix.

In [778]: # drop numerical variables, and target, as we just want categorical predictor variables
df_numerical = df.drop(['enrolled_university', 'city_development_index', 'experience', 'training_hours', 'target'], axis=1)

In [777]: # instantiate matrix as NxN table of zeros
cramer_matrix = np.zeros((9,9))
for i in range(0,9):
 for j in range(0,9):
 col1 = df_numerical.columns[i]
 col2 = df_numerical.columns[j]
 _, cramer_v = rpy.crosstab(df[col1], df[col2], test='chi-square')
 cramer_matrix[i,j] = cramer_v.loc[0].exp
matrix to dataframe, with variable as column and row names
cramer_dataframe = pd.DataFrame(cramer_matrix)
cramer_dataframe.columns = df_numerical.columns
cramer_dataframe.index = df_numerical.columns

Out[777]:

	city	gender	relevant_experience	enrolled_university	education_level	major_discipline	company_size	company_type
city	1.0000	1.0009	0.1469	0.2055	0.0579	0.1137	0.1314	0.1
gender	0.1200	1.0000	0.0584	0.0263	0.0850	0.0764	0.0270	0.0
relevant_experience	0.1469	0.0584	1.0000	0.3879	0.3175	0.0731	0.0603	0.2
enrolled_university	0.2055	0.0263	0.3879	1.0000	0.1659	0.0655	0.0481	0.0
education_level	0.0579	0.0850	0.3175	0.1659	1.0000	0.0465	0.0552	0.0
major_discipline	0.1137	0.0764	0.0731	0.0655	0.0465	1.0000	0.0322	0.0
company_size	0.1314	0.0270	0.0603	0.0481	0.0552	0.0322	1.0000	0.1
company_type	0.1342	0.0523	0.2058	0.1003	0.0810	0.0384	0.1928	1.0
last_new_job	0.1291	0.0345	0.3907	0.1982	0.1749	0.0247	0.0815	0.0

The diagonal of ones corresponds to each variable's relationship with itself. We see that there are moderate-strong relationships between

- city & relevant_experience
- city & enrolled_university
- city & education_level
- relevant_experience & enrolled_university
- relevant_experience & company_type
- company_type & company_size
- last_new_job & education_level

These are interesting from an analytical perspective, as well as from a modelling perspective.

In [780]: df_numerical = df[['city_development_index', 'experience', 'training_hours']]
df_numerical = df[['city_development_index', 'experience', 'training_hours']]

In [781]: # instantiate matrix as 3x3 table of zeros
pearson_matrix = np.zeros((3,3))
fill na of experience with the median value
only have 65 NaN values, so this shouldn't cause calculation problems
df_numerical.experience.fillna(df_numerical.experience.median(), inplace=True)
calculate Pearson R Correlation for each variable, update matrix entries
for i in range(0,3):
 for j in range(0,3):
 col1 = df_numerical.columns[i]
 col2 = df_numerical.columns[j]
 pearson_r = stats.pearsonr(df_numerical[col1], df_numerical[col2])
 pearson_matrix[i,j] = pearson_r
matrix to dataframe with variable as column and row names
pearson_dataframe = pd.DataFrame(pearson_matrix)
pearson_dataframe.columns = df_numerical.columns
pearson_dataframe.index = df_numerical.columns

Out[781]:

	city_development_index	experience	training_hours
city_development_index	1.000000	0.333368	0.001820
experience	0.333368	1.000000	0.000488
training_hours	0.001820	0.000488	1.000000

We see that there is only moderate correlation between experience and city_development_index (remember that Cramer's V and Pearson's R have different interpretations - 0.33 is very high for Pearson's R, but only moderate for Cramer's R).

Conclusions

We now have an idea of what variables are likely to be influential in whether a person is looking for a new job or not. The variables which seem to have the strongest relationship with our target variable are:

- city_development_index (-0.346, Pearson)
- experience (-0.176, Pearson)
- city (0.39, Cramer)
- enrolled_university (0.15, Cramer)
- relevant_experience (0.12, Cramer)

We also have some variables which might be helpful in our models, but don't have as strong of a relationship with target:

- education_level (0.09, Cramer)
- last_new_job (0.08, Cramer)
- company_size (0.08, Cramer)
- company_type (0.08, Cramer)
- training_hours (-0.02, Pearson)

We will take a selection of these variables, and see which combination gives us the best model for predicting our target variable. We must be careful about the variables city, company_size and company_type; city is very imbalanced, and both company_size and company_type are missing large amounts of data.

Further Analysis (can be skipped)

In a real-world situation, we could perform more analysis here on the interactions of all variables with each other. For example, we could analyse the impact that city has on education_level, or whether certain major_discipline or education_levels require more training_hours.

This kind of analysis that gives us huge insights into why we can improve hiring or training for employees, especially when combined with real-world domain knowledge gained by talking to teams in HR and other areas of the organisation.

Because it is so in-depth, we will not perform all that analysis in this notebook. However, below is an example of using the Kruskal-Wallis test to check for differences between groups (enrolled_university, education_level) with regard to variables such as experience, training_hours or city_development_index.

For example, we see that there is a significant difference between groups of education_level, and city_development_index, and we might hypothesise that the most highly educated workers can find jobs in the best cities, or rather, that they won't be stuck in less developed cities. This then requires further analysis, in order to get the most insights out of the data.

In [835]: enrun1 = df[['enrolled_university', 'city_development_index', 'training_hours', 'experience']].dropna().set_index(keys=['company_type'])
enrun1 = enrun1.rename(columns={'city_development_index': 'cdi', 'training_hours': 'trh', 'experience': 'exp'})
no_enrl = enrun1.loc['no_enrollment']
part_time = enrun1.loc['Part time course']
full_time = enrun1.loc['Full time course']
print('cdi', stats.kruskal(no_enrl.cdi, part_time.cdi, full_time.cdi))
print('trh', stats.kruskal(no_enrl.trh, part_time.trh, full_time.trh))
print('exp', stats.kruskal(no_enrl.exp, part_time.exp, full_time.exp))
cdi = KruskalResult(statistic=441.727444134247, pvalue=1.756212375961848e-96)
trh = KruskalResult(statistic=1.604553372708646, pvalue=0.4483071486022086)
exp = KruskalResult(statistic=2350.796072508217, pvalue=0.0)

Now that we have analysed the relationships between our variables, we can begin to model. It is important to check the relationships of these variables first, as one of the assumptions of a Logistic Regression model is that the predictor variables are fairly independent (i.e. don't have a strong relationship with each other). Therefore, we can use the information above to inform our model building.

We will try building a few different Logistic Regression models, using different variable combinations, and check their scores and prediction accuracy. We will then tune our model parameters, to see if we can improve performance. Finally, we will check other types of models, such as Support Vector Classifiers, Random Forests, and Neural Networks.

It is important to remember that the process of trial and error, based on exploratory data analysis which we did above. In building our models here, we must keep in mind what the purpose is. We are aiming to predict the probability of Data Science employees seeking another job, based on the information that we know about their background. We would then tailor this to our specific situation - for example, if we are based in Milan, we could disregard the cities which are not Milan. If we know the size and type of our company, we could disregard entries for company_size and company_type which don't fit our company.

We saw above that the variables which most correlate to our target variable are:

- city_development_index, experience, city, enrolled_university, relevant_experience (strongly).
- education_level, last_new_job, company_size, company_type, training_hours (moderately/weakly).

We will start with the following models:

- just city_development_index, as a baseline.
- city_development_index and all other strong variables (excluding city).
- city_development_index and all other strong and weak variables (excluding city, company_size and company_type).

As we are first trying to build a general model, we are excluding city because it is so unbalanced, and company_size and company_type because they have so many missing observations.

Then we will include city, using a subset of the data where we have over 60 observations for each city. This gives us some certainty that our results are accurate - if we include cities with only 1 or 2 observations, we can't be certain that these are representative of those cities as a whole.

Then finally, we will include company_size and company_type. We must restrict ourselves to an even smaller subset here, as these variables have so many missing observations.

```
In [927]: # import models and accuracy measures
import sklearn.metrics
import sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix, classification_report
from sklearn.model_selection import train_test_split
```

We will want to set aside some data which we won't train our model on. We can then use this data to validate how each model is performing. This is standard practice.

We will use a random number generator to split our dataset into training and test data. We will use about 80% of the data to train our models, and then can use the remaining unused 20% to test them.

```
In [940]: # generate random integers between 0 and maximum entries number of df
training_index = np.random.randint(0, len(df), round(0.8 * len(df)))
```

```
In [955]: # subset data - observations with index in training index going into training set
training_data = df.iloc[training_index]
# and then we drop those observations, and keep the rest as a validation set
validation_data = df.drop(training_index, axis=0)
```

```
In [957]: # test the shapes of the sets
print(training_data.shape)
print(validation_data.shape)

(15326, 14)
(8640, 14)
```

This looks correct. Now we can move on to building models.

Models 1-3

```
In [1024]: train1 = training_data[['city_development_index', 'experience', 'enrolled_university', 'relevant_experience', 'education_level', 'last_new_job', 'training_hours', 'target']].dropna()
val1 = validation_data[['city_development_index', 'experience', 'enrolled_university', 'relevant_experience', 'education_level', 'last_new_job', 'training_hours', 'target']].dropna()
```

```
In [1025]: model1 = smf.logit(formula = 'target ~ city_development_index', data = train1).fit()
model2 = smf.logit(formula = 'target ~ city_development_index + experience + enrolled_university + relevant_experience', data = train1).fit()
model3 = smf.logit(formula = 'target ~ city_development_index + experience + enrolled_university + relevant_experience + education_level + last_new_job + training_hours', data = train1).fit()
```

Optimization terminated successfully.
Current function value: 0.499787
Iterations 6
Optimization terminated successfully.
Current function value: 0.489555
Iterations 6
Optimization terminated successfully.
Current function value: 0.481681
Iterations 6

```
In [1026]: # test our predictions against our validation set
y_val = train1['target']
predictions1 = model1.predict(val1.drop(['target'], axis=1)).round()
predictions2 = model2.predict(val1.drop(['target'], axis=1)).round()
predictions3 = model3.predict(val1.drop(['target'], axis=1)).round()
```

We want to know how accurate we are, compared to the actual data. We see below that of our target variable for the whole dataset, about 75% of observations are not looking for a job, and 25% are. Therefore, we could randomly say that 3/4 people are not looking for work, and probably be right. Hence, we want our model to beat this 75%.

```
In [1034]: # not looking vs looking
df.target.value_counts()
```

```
Out[1034]: 0.0    14381
          1.0    4777
          Name: target, dtype: int64
```

```
In [1033]: # percentage of people not looking
1 - df.target.sum()/len(df.target)
```

```
Out[1033]: 0.7506524689424783
```

```
In [1029]: print('Results for Model 1: target ~ city_development_index')
print(accuracy_score(y, predictions1))
print(confusion_matrix(y, predictions1))
print('\n')
print('Results for Model 2: target ~ city_development_index + experience + enrolled_university + relevant_experience')
print(accuracy_score(y, predictions2))
print(confusion_matrix(y, predictions2))
print('\n')
print('Results for Model 3: target ~ city_development_index + experience + enrolled_university + relevant_experience + education_level + last_new_job + training_hours')
print(accuracy_score(y, predictions3))
print(confusion_matrix(y, predictions3))
```

```
Results for Model 1: target ~ city_development_index
0.75968
[[6000 116]
 [1840 312]]

Results for Model 2: target ~ city_development_index + experience + enrolled_university + relevant_experience
0.7657
[[5804 312]
 [1595 428]]
```

We see that our first model correctly predicts 6000 people who are not looking for a job, but falsely predicts 1840 who are looking. This is a problem, as this is exactly what we want to predict. Its accuracy is not much better than our 75% above.

Models 2 and 3 make small improvements, but Model 3 has specific multicollinearity issues because relevant_experience has a high Correlation V with both enrolled_university and education_level. The extra increase in accuracy is not worth this multicollinearity, so we should choose Model 2 to move forward with.

Model 4

Now we will include city. We must reduce our dataset to a subset which has over 60 observations of each city.

```
In [1101]: # get list of cities with >60 observations per city
cities_60_index = df.city.value_counts()[df.city.value_counts()>60].index
# subset dataframe
cities_60 = df.loc[df.city.isin(cities_60_index)]
# reset index in order to randomly subset for training and validation
cities_60.reset_index(drop=True, inplace=True)
# get random index
cities_60_random_index = np.random.randint(0, len(cities_60), int(0.8 * len(cities_60)))
```

```
In [1111]: # training set
cities_train = cities_60.iloc[cities_60_random_index]
# validation set
cities_val = cities_60.drop(cities_60_random_index, axis=0)
```

```
In [1094]: # selecting variables, based on Model 2
train2 = cities_train[['city', 'city_development_index', 'experience', 'enrolled_university', 'relevant_experience', 'education_level', 'last_new_job', 'training_hours', 'target']].dropna()
val2 = cities_val[['city', 'city_development_index', 'experience', 'enrolled_university', 'relevant_experience', 'education_level', 'last_new_job', 'training_hours', 'target']].dropna()
```

```
In [1096]: # fit model
model4 = smf.logit(formula = 'target ~ city + city_development_index + experience + enrolled_university + relevant_experience', data = train2).fit(maxiter=100)
Warning: Maximum number of iterations has been exceeded.
Current function value: 0.476459
Iterations: 100
```

This warning is likely because of our city variable. In practice we would investigate this further and maybe change the number of city observations.

```
In [1101]: # predictions vs validation target
y_val = model4.predict(val2).round()
y2 = val2.target
```

```
In [1102]: print(accuracy_score(y2, y_val).round(5))
print(confusion_matrix(y2, y_val))

0.79255
[[5254 510]
 [3160 744]]
```

By including city, we have improved our accuracy slightly from above, but at the expense of being able to use those cities for which we had 60 or less observations.

Models 5-7

Now we will incorporate company_size and company_type. These variables are likely strongly related, so we might want to use just one of these variables - including both might not give us any extra information, but might cost us in terms of multicollinearity.

We must also consider our dataset even further, due to missing values which we cannot reasonably replace. In practice, we would contact the data provider or data office, in order to see if we can get the information for these missing observations, instead of dropping them.

```
In [1103]: # subset further
train3 = training_data[['city_development_index', 'experience', 'enrolled_university', 'relevant_experience', 'education_level', 'last_new_job', 'training_hours', 'target', 'company_size', 'company_type']].dropna()
val3 = validation_data[['city_development_index', 'experience', 'enrolled_university', 'relevant_experience', 'education_level', 'last_new_job', 'training_hours', 'target', 'company_size', 'company_type']].dropna()
```

```
In [1112]: # including company_size
model5 = smf.logit(formula = 'target ~ city_development_index + experience + enrolled_university + relevant_experience + company_size', data = train3).fit()
# including company_type
model6 = smf.logit(formula = 'target ~ city_development_index + experience + enrolled_university + relevant_experience + company_type', data = train3).fit()
```

```
In [1113]: # including both
model7 = smf.logit(formula = 'target ~ city_development_index + experience + enrolled_university + relevant_experience + company_size + company_type', data = train3).fit()
Optimization terminated successfully.
Current function value: 0.373413
Iterations 7
Optimization terminated successfully.
Current function value: 0.373217
Iterations 7
Optimization terminated successfully.
Current function value: 0.371820
Iterations 7
```

```
In [1105]: # predictions versus validation target
y = val3.target
p5 = model5.predict(val3.drop(['target'], axis=1)).round()
p6 = model6.predict(val3.drop(['target'], axis=1)).round()
p7 = model7.predict(val3.drop(['target'], axis=1)).round()

print(accuracy_score(y, p5).round(5))
print(confusion_matrix(y, p5))
print('\n')
print(accuracy_score(y, p6).round(5))
print(confusion_matrix(y, p6))
print('\n')
print(accuracy_score(y, p7).round(5))
print(confusion_matrix(y, p7))

0.83531
[[4278 174]
 [ 716 236]]

0.83142
[[4316 136]
 [ 775 177]]

0.83235
[[4257 195]
 [ 711 240]]
```

We have increased our accuracy quite a bit, compared to our Model 2 above. The model with just company_size seems to be the best one - it has the highest accuracy, and avoids the multicollinearity issues of including both company_size and company_type.

Models 8-10

Finally, we can try incorporating city.

```
In [1106]: # subset cities over 60 observations with company size, remove missing values
train_f = cities_train[['city', 'city_development_index', 'experience', 'enrolled_university', 'relevant_experience', 'education_level', 'last_new_job', 'training_hours', 'target', 'company_size', 'company_type']].dropna()
val_f = cities_val[['city', 'city_development_index', 'experience', 'enrolled_university', 'relevant_experience', 'education_level', 'last_new_job', 'training_hours', 'target', 'company_size', 'company_type']].dropna()
```

```
In [1113]: model8 = smf.logit(formula = 'target ~ city + city_development_index + experience + enrolled_university + relevant_experience + company_size', data = train_f).fit()
model9 = smf.logit(formula = 'target ~ city + city_development_index + experience + enrolled_university + relevant_experience + company_type', data = train_f).fit()
model10 = smf.logit(formula = 'target ~ city + city_development_index + experience + enrolled_university + relevant_experience + company_size + company_type', data = train_f).fit()
```

Warning: Maximum number of iterations has been exceeded.
Current function value: 0.358863
Iterations: 35
Warning: Maximum number of iterations has been exceeded.
Current function value: 0.358578
Iterations: 35
Warning: Maximum number of iterations has been exceeded.
Current function value: 0.357240
Iterations: 35

```
In [1109]: # predictions versus validation target
y = val_f.target
p11 = model11.predict(val_f.drop(['target'], axis=1)).round()
p12 = model12.predict(val_f.drop(['target'], axis=1)).round()
p13 = model13.predict(val_f.drop(['target'], axis=1)).round()

print(accuracy_score(y, p11).round(5))
print(confusion_matrix(y, p11))
print('\n')
print(accuracy_score(y, p12).round(5))
print(confusion_matrix(y, p12))
print('\n')
print(accuracy_score(y, p13).round(5))
print(confusion_matrix(y, p13))

0.85703
[[5381 335]
 [ 394 489]]

0.85566
[[3865 351]
 [ 385 498]]
```

0.85664
[[5374 332]
 [399 484]]

We see that our most accurate model seems to have the formula as follows:

- target ~ city_development_index + experience + enrolled_university + education_level + city + company_size

We must then reach the decision of whether this increased accuracy is warranted, as it is less applicable than our previous models which didn't include city and company_size. This is a well-known issue in modelling, known as the Bias-Variance Tradeoff.

Further Modelling and Tuning

Now that we know a good combination of variables, we can try other model types, as well as tuning our Logistic Regression model parameters. This is given below, without much explanation, but more so as an illustration of the next steps.

Random Forest

```
In [1114]: train_extra = train_f.drop(['company_type'], axis=1)
val_extra = val_f.drop(['company_type'], axis=1)
X_train = train_extra.drop(['target'], axis=1)
y_train = train_extra.target
X_val = val_extra.drop(['target'], axis=1)
y_val = val_extra.target

print(X_train.shape)
print(y_train.shape)
print(X_val.shape)
print(y_val.shape)

(8953, 9)
(8953, 1)
(5099, 9)
(5099, 1)
```

```
In [1119]: X_train = pd.get_dummies(X_train, columns = ['city', 'enrolled_university', 'relevant_experience', 'education_level', 'last_new_job', 'company_size'])
X_val = pd.get_dummies(X_val, columns = ['city', 'enrolled_university', 'relevant_experience', 'education_level', 'last_new_job', 'company_size'])
```

```
In [1120]: from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier(n_estimators=100)
rfc.fit(X_train, y_train)
```

```
Out[1120]: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None, criterion='gini', max_depth=None, max_features='auto', max_leaf_nodes=None, max_samples=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=None, oob_score=False, random_state=None, verbose=0, warm_start=False)
```

```
In [1121]: rfpreds = rfc.predict(X_val)
print(accuracy_score(rfpreds, y_val))
print('\n')
print(confusion_matrix(rfpreds, y_val))

0.8442831927829887

[[3956 534]
 [ 260 349]]
```

Comparable to our best logistic regression above. Predicts people who are looking for a job more accurately, but also predicts more people who are not looking for a job, inaccurately.

Support Vector Classifier

```
In [462]: from sklearn.svm import SVC
svc = SVC()
```

```
In [1124]: svc.fit(X_train, y_train)
```

```
Out[1124]: SVC(C=1.0, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma='scale', kernel='rbf', max_iter=1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False)
```

```
In [1125]: spreds = svc.predict(X_val)
print(accuracy_score(spreds, y_val))
print('\n')
print(confusion_matrix(spreds, y_val))

0.826828789598155

[[4216 883]
 [  0  0]]
```

This is a big problem, it doesn't predict anyone to be looking for a job! We can do a grid search for the best parameters.

```
In [1126]: param_grid = {'C': [0.1, 1, 10, 100, 1000], 'gamma': [1.0, 1.01, 0.001, 0.0001], 'kernel': ['rbf']}
grid = GridSearchCV(SVC(), param_grid, refit=True, verbose=0)
```

```
In [ ]: grid.fit(X_train, y_train)
```

```
In [469]: grid.best_params_
```

```
Out[469]: {'C': 1000, 'gamma': 0.1, 'kernel': 'rbf'}
```

```
Out[1129]: Grid
          estimator=SVC(C=1000, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma='scale', kernel='rbf', max_iter=1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False),
          id='deprecated',
          n_jobs=None,
          param_grid={'C': [0.1, 1, 10, 100, 1000], 'gamma': [1.0, 1.01, 0.001, 0.0001], 'kernel': ['rbf']},
          pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
          scoring=None, verbose=0)
```

```
In [1132]: svc2 = SVC(C=1, gamma=0.1, kernel='rbf')
svc2.fit(X_train, y_train)
spreds2 = svc2.predict(X_val)
```

```
In [1133]: print(accuracy_score(spreds2, y_val))
print('\n')
print(confusion_matrix(spreds2, y_val))

0.8407530888409492

[[4116 712]
 [ 100 171]]
```

This is an improvement, and we could work further with this model, but is still predicting too few people who are looking for jobs.

Neural Network

A Neural Network is likely unsuitable for the type of data, but we can try it as an illustration. We first scale our data, then use a fully-connected network, followed by one with dropout layers.

```
In [1136]: from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Input, Activation, Dropout
from tensorflow.keras.callbacks import EarlyStopping
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
```

```
In [1140]: scaler.fit(X_train)
scaled_train = scaler.transform(X_train)
scaled_val = scaler.transform(X_val)
```

```
In [1142]: early_stop = EarlyStopping(monitor='val_loss', mode='min', verbose=1, patience=25)
```

```
In [1145]: model = Sequential()
model.add(Dense(units=76, activation='relu'))
model.add(Dense(units=34, activation='relu'))
model.add(Dense(units=19, activation='relu'))
model.add(Dense(units=10, activation='relu'))
model.add(Dense(units=1, activation='sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='adam')
```

```
In [1150]: model.fit(scaled_train,
                    y_train.values,
                    epochs=60,
                    validation_data=(scaled_val, y_val.values), verbose=1,
                    callbacks=[early_stop])
```

Train on 8953 samples, validate on 5099 samples
Epoch 1/600
8953/8953 [=====] - 22s 2ms/sample - loss: 0.3988 - val_loss: 0.3601
Epoch 2/600
8953/8953 [=====] - 1s 93us/sample - loss: 0.3560 - val_loss: 0.3614
Epoch 3/600
8953/8953 [=====] - 1s 93us/sample - loss: 0.3490 - val_loss: 0.3647
Epoch 4/600
8953/8953 [=====] - 1s 96us/sample - loss: 0.3425 - val_loss: 0.3669
Epoch 5/600
8953/8953 [=====] - 1s 98us/sample - loss: 0.3316 - val_loss: 0.3741
Epoch 6/600
8953/8953 [=====] - 1s 93us/sample - loss: 0.3220 - val_loss: 0.3793
Epoch 7/600
8953/8953 [=====] - 1s 96us/sample - loss: 0.3122 - val_loss: 0.3864
Epoch 8/600
8953/8953 [=====] - 1s 94us/sample - loss: 0.2998 - val_loss: 0.4093
Epoch 9/600
8953/8953 [=====] - 1s 98us/sample - loss: 0.2880 - val_loss: 0.4131
Epoch 10/600
8953/8953 [=====] - 1s 98us/sample - loss: 0.2823 - val_loss: 0.4261
Epoch 11/600
8953/8953 [=====] - 1s 103us/sample - loss: 0.2691 - val_loss: 0.4454
Epoch 12/600
8953/8953 [=====] - 1s 116us/sample - loss: 0.2619 - val_loss: 0.4461
Epoch 13/600
8953/8953 [=====] - 1s 113us/sample - loss: 0.2449 - val_loss: 0.4640
Epoch 14/600
8953/8953 [=====] - 1s 104us/sample - loss: 0.2335 - val_loss: 0.5111
Epoch 15/600
8953/8953 [=====] - 1s 91us/sample - loss: 0.2253 - val_loss: 0.5324
Epoch 16/600
8953/8953 [=====] - 1s 91us/sample - loss: 0.2193 - val_loss: 0.5629
Epoch 17/600
8953/8953 [=====] - 1s 98us/sample - loss: 0.1959 - val_loss: 0.6261
Epoch 18/600
8953/8953 [=====] - 1s 94us/sample - loss: 0.2151 - val_loss: 0.5554
Epoch 19/600
8953/8953 [=====] - 1s 92us/sample - loss: 0.2089 - val_loss: 0.5504
Epoch 20/600
8953/8953 [=====] - 1s 92us/sample - loss: 0.2033 - val_loss: 0.5705
Epoch 21/600
8953/8953 [=====] - 1s 92us/sample - loss: 0.1997 - val_loss: 0.5885
Epoch 22/600
8953/8953 [=====] - 1s 93us/sample - loss: 0.1901 - val_loss: 0.6261
Epoch 23/600
8953/8953 [=====] - 1s 98us/sample - loss: 0.1959 - val_loss: 0.6261
Epoch 24/600
8953/8953 [=====] - 1s 97us/sample - loss: 0.1902 - val_loss: 0.6374
Epoch 25/600
8953/8953 [=====] - 1s 94us/sample - loss: 0.1888 - val_loss: 0.6459
Epoch 26/600
8953/8953 [=====] - 1s 93us/sample - loss: 0.1847 - val_loss: 0.6777
Epoch 00026: early stopping

```
Out[1150]: <tensorflow.python.keras.callbacks.History at 0x19fc3b8c98>
```

```
In [1151]: model_loss = pd.DataFrame(model.history.history)
model_loss.plot()
```

```
Out[1151]: <matplotlib.axes._subplots.AxesSubplot at 0x19fc3b8f648>
```



This is a bad sign - our model is overfitting on our training data. We can try now with some dropout layers.

```
In [1153]: model = Sequential()
model.add(Dense(units=69, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(units=34, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(units=17, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(units=7, activation='relu'))
model.add(Dropout(0.5))
model.compile(loss='binary_crossentropy', optimizer='adam')
```

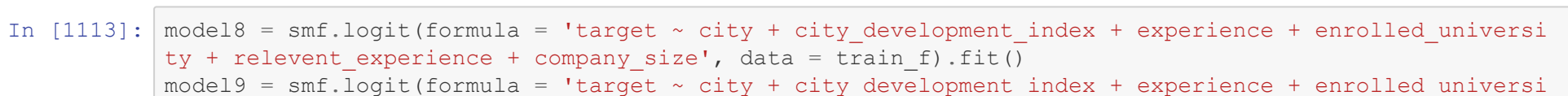
```
In [1154]: model.fit(scaled_train,
                    y_train.values,
                    epochs=60,
                    validation_data=(scaled_val, y_val.values), verbose=1,
                    callbacks=[early_stop])
```

Train on 8953 samples, validate on 5099 samples
Epoch 1/600
8953/8953 [=====] - 58s 6ms/sample - loss: 0.5771 - val_loss: 0.4521
Epoch 2/600
8953/8953 [=====] - 1s 117us/sample - loss: 0.4773 - val_loss: 0.4194
Epoch 3/600
8953/8953 [=====] - 1s 118us/sample - loss: 0.4536 - val_loss: 0.3751
Epoch 4/600
8953/8953 [=====] - 1s 120us/sample - loss: 0.4320 - val_loss: 0.3749
Epoch 5/600
8953/8953 [=====] - 1s 138us/sample - loss: 0.4226 - val_loss: 0.3814
Epoch 6/600
8953/8953 [=====] - 1s 118us/sample - loss: 0.4123 - val_loss: 0.3751
Epoch 7/600
8953/8953 [=====] - 1s 118us/sample - loss: 0.4084 - val_loss: 0.3732
Epoch 8/600
8953/8953 [=====] - 1s 118us/sample - loss: 0.4064 - val_loss: 0.3751
Epoch 9/600
8953/8953 [=====] - 1s 123us/sample - loss: 0.3993 - val_loss: 0.3711
Epoch 10/600
8953/8953 [=====] - 1s 117us/sample - loss: 0.3980 - val_loss: 0.3688
Epoch 11/600
8953/8953 [=====] - 1s 119us/sample - loss: 0.3993 - val_loss: 0.3675
Epoch 12/600
8953/8953 [=====] - 1s 118us/sample - loss: 0.3849 - val_loss: 0.3683
Epoch 13/600
8953/8953 [=====] - 1s 118us/sample - loss: 0.3761 - val_loss: 0.3699
Epoch 14/600
8953/8953 [=====] - 1s 118us/sample - loss: 0.3748 - val_loss: 0.3696
Epoch 15/600
8953/8953 [=====] - 1s 123us/sample - loss: 0.3762 - val_loss: 0.3699
Epoch 16/600
8953/8953 [=====] - 1s 119us/sample - loss: 0.3717 - val_loss: 0.3704
Epoch 17/600
8953/8953 [=====] - 1s 125us/sample - loss: 0.3703 - val_loss: 0.3690
Epoch 18/600
8953/8953 [=====] - 1s 125us/sample - loss: 0.3708 - val_loss: 0.3767
Epoch 19/600
8953/8953 [=====] - 1s 134us/sample - loss: 0.3660 - val_loss: 0.3706
Epoch 20/600
8953/8953 [=====] - 1s 118us/sample - loss: 0.3695 - val_loss: 0.3736
Epoch 21/600
8953/8953 [=====] - 1s 118us/sample - loss: 0.3630 - val_loss: 0.3748
Epoch 22/600
8953/8953 [=====] - 1s 122us/sample - loss: 0.3595 - val_loss: 0.3770
Epoch 23/600
8953/8953 [=====] - 1s 115us/sample - loss: 0.3562 - val_loss: 0.3751
Epoch 24/600
8953/8953 [=====] - 1s 126us/sample - loss: 0.3547 - val_loss: 0.3831
Epoch 25/600
8953/8953 [=====] - 1s 121us/sample - loss: 0.3507 - val_loss: 0.3757
Epoch 26/600
8953/8953 [=====] - 1s 117us/sample - loss: 0.3515 - val_loss: 0.3767
Epoch 27/600
8953/8953 [=====] - 1s 118us/sample - loss: 0.3538 - val_loss: 0.3795
Epoch 28/600
8953/8953 [=====] - 1s 123us/sample - loss: 0.3501 - val_loss: 0.3826
Epoch 29/600
8953/8953 [=====] - 1s 118us/sample - loss: 0.3478 - val_loss: 0.3807
Epoch 30/600
8953/8953 [=====] - 1s 119us/sample - loss: 0.3521 - val_loss: 0.3830
Epoch 00031: early stopping

```
Out[1154]: <tensorflow.python.keras.callbacks.History at 0x19fc7f64b88>
```

```
In [1155]: model_loss = pd.DataFrame(model.history.history)
model_loss.plot()
```

```
Out[1155]: <matplotlib.axes._subplots.AxesSubplot at 0x19fc7f64b88>
```



Unlike our first Neural Network, this one has not converged to the training data. However, it has made very little improvement in the validation loss, and so we conclude that a model type like this is not suitable for our dataset.

Grid Search on Logistic Regression

We find that our logistic regression model had the best results. We can try to tune its parameters also.

```
In [678]: # import Logistic Regression model which can be tuned via grid search
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
```