

Predictive Analytics, Fall 2016
Professor Anasse Bari, PhD
Department of Computer Science

What will the real estate market look like after the L train closure in 2019?

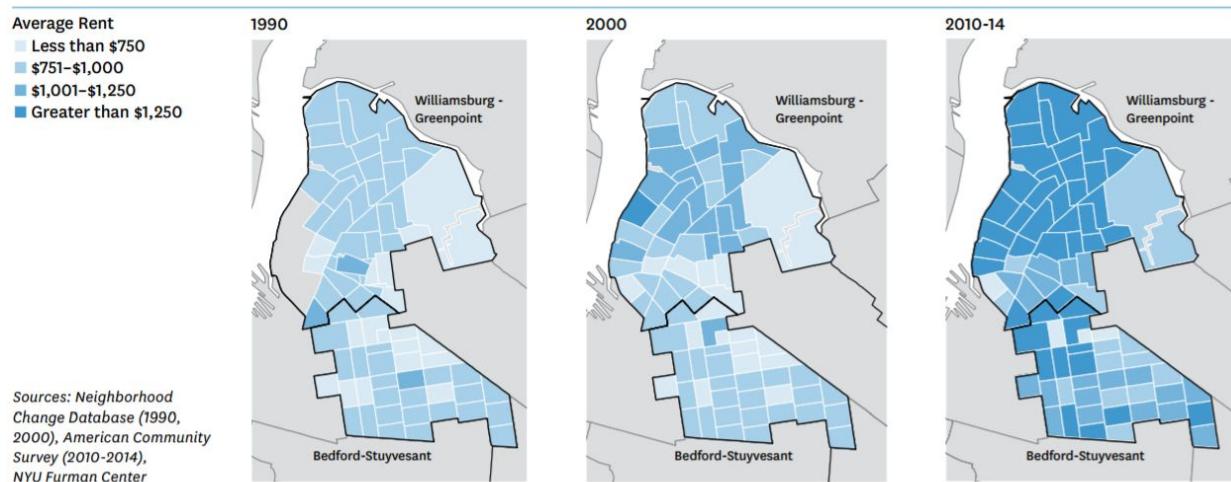
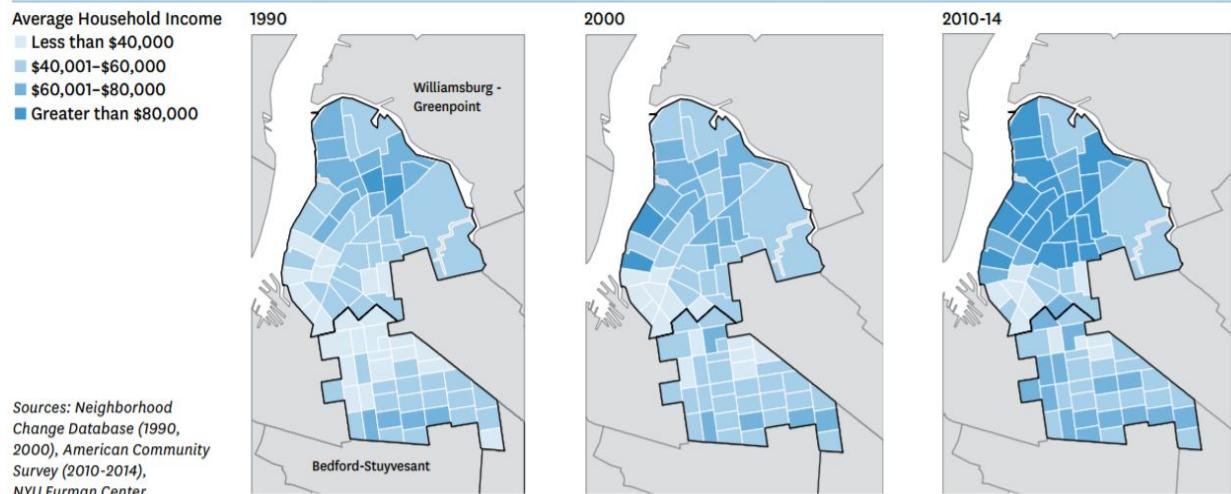
Problem Statement and Analytics Hypothesis

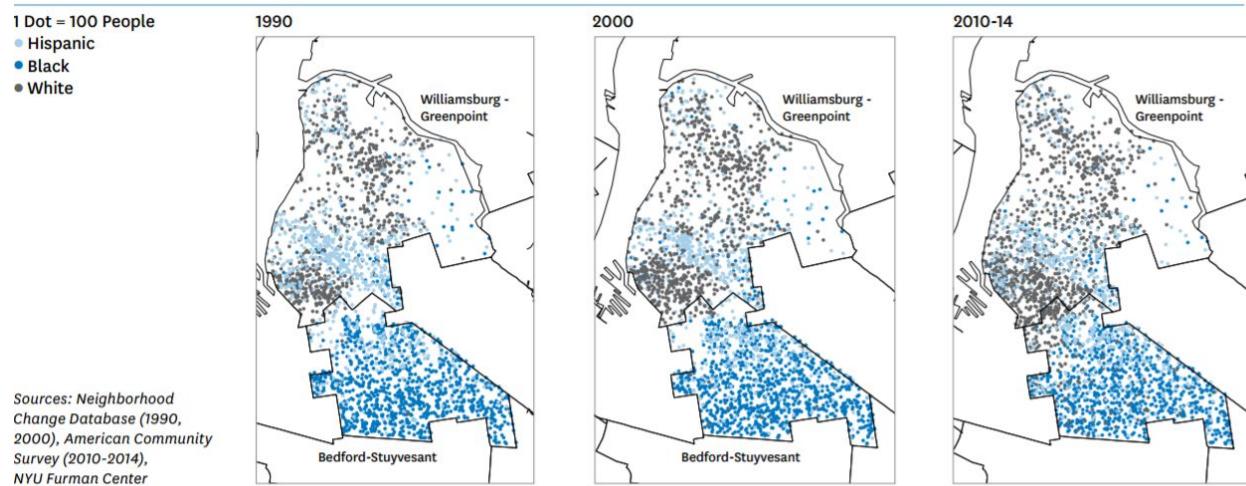
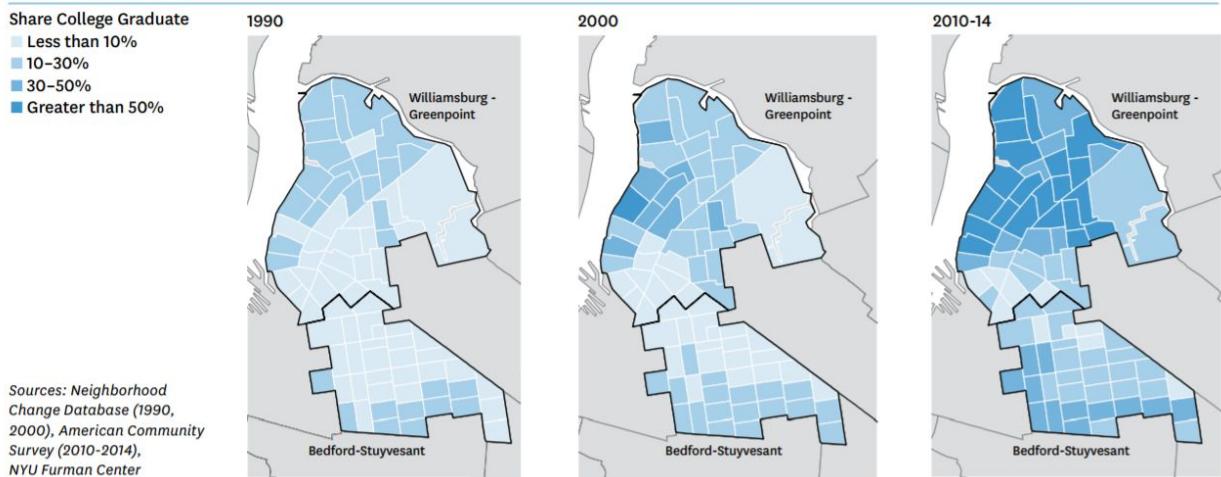
This study aims to characterize factors contributing to apartment rental prices in New York City, and to build a predictive model of rent prices following the closure of the L train line in early 2019. Suspended L train service from Brooklyn to Manhattan will likely demonstrate an increase of rental prices in areas off of other train lines that are proximal to Downtown Manhattan's commercial areas and easily accessed by public transportation. In contrast, we expect areas affected by the closure to experience a decrease in rental prices.

Business Understanding

In order to narrow down the necessary attributes used in the analysis, previous studies on New York City gentrification were examined. First, we confirmed that gentrification yields an increase in rental prices. Gentrification is a term describing cities or parts of cities, primarily in the US, that see a sharp increase in groups of stable socioeconomic status moving into areas that were not previously as expensive. Gentrification is of particular concern in New York City because it dramatically affects the cost of living. NYU Furman Center's 2016 study depicts a change in population demographics as the Williamsburg, Brooklyn area gentrifies between 1990 and 2014 the figures below. We chose 138 attributes based on their contributions to gentrification and increased rent prices: various crime indices, household income split by racial and ethnic group, culture index, transportation-related categories (average travel time to work, mode of transport, number of subway lines and stations, etc.), employment domains, educational attainment, and more.

The specific aim of this project is to build a predictive model that can accurately determine median rental values and subsequent percent changes by zip code after the 2019 L train closure.





Data Understanding

Initially, four datasets were examined for the right attributes: NYC Open Data, Simply Map 3.0, Zillow, and MTA and NYC.gov data. Simply Map 3.0 and MTA/NYC.gov were chosen because both contained necessary information in a readable format. Simply Map 3.0 provided data collected annually from 2000 and 2008-2016, and compiled from various established sources: New York City government data, United States census information (for demographics information), crime, culture, and restaurant indices, etc. MTA/NYC.gov provided relevant information about the location and names of subway lines. At first glance, the data from 2000, 2008, and 2009 contained more missing values than 2010-2016 (more detail provided later).

Missing values were noted at particular zip codes as well. These zip codes are non-residential zip codes that exist for United States Postal Service (USPS) purposes only, therefore data was not collected from these areas.

Data Preparation

Attributes were downloaded in .csv format from Simply Map 3.0:

Location	Food Services (\$000), 2016	Food Services (\$000), 2015	Food Services (\$000), 2014	Food Services (\$000), 2013	Food Services (\$000), 2012	Food Services (\$000), 2011	Food Services (\$000), 2010
11614, Hewlett Beach, NY	\$61,079.00	\$77,194.00	\$65,133.00	\$56,221.00	\$41,580.00	\$55,466.00	
11451, Kew Gardens, NY	\$6,974.00	\$7,893.00	\$8,333.00	\$20,109.00	\$7,193.00	\$5,591.00	
11456, Ozone Park, NY	\$12,036.00	\$14,482.00	\$24,133.00	\$14,657.00	\$25,598.00	\$12,773.00	
11457, Ozone Park, NY	\$59,584.00	\$62,104.00	\$21,000.00	\$0.00	\$13,439.00	\$17,205.00	
11618, Richmond Hill, NY	\$20,676.00	\$27,464.00	\$16,533.00	\$81,499.00	\$11,168.00	\$15,965.00	
11455, South Richmond Hill, NY	\$26,626.00	\$30,465.00	\$25,266.00	\$0.00	\$24,229.00	\$26,071.00	
11450, South Ozone Park, NY	\$12,653.00	\$11,612.00	\$11,600.00	\$0.00	\$13,285.00	\$8,457.00	
11421, Woodhaven, NY	\$21,417.00	\$16,374.00	\$15,667.00	\$6,868.00	\$22,021.00	\$12,831.00	
11422, Rosedale, NY	\$26,948.00	\$19,766.00	\$21,466.00	\$5,665.00	\$21,516.00	\$17,439.00	
11403, Hollis, NY	\$11,233.00	\$9,263.00	\$8,733.00	\$14,020.00	\$8,707.00	\$8,165.00	
11624, Jamaica, NY	\$0.00	\$0.00	\$0.00	\$57,637.00	\$0.00	\$0.00	
11425, Jamaica, NY	\$123.00	\$130.00	\$133.00	\$5,771.00	\$126.00	\$127.00	
11426, Bellerose, NY	\$10,369.00	\$9,981.00	\$8,800.00	\$26,057.00	\$8,139.00	\$8,749.00	
11427, Queens Village, NY	\$15,800.00	\$15,852.00	\$16,533.00	\$17,065.00	\$13,250.00	\$11,956.00	
11428, Queens Village, NY	\$7,036.00	\$8,742.00	\$7,267.00	\$15,082.00	\$5,742.00	\$2,566.00	
11429, Queens Village, NY	\$27,898.00	\$28,508.00	\$28,733.00	\$26,411.00	\$32,290.00	\$24,029.00	

There were 140 attributes and 2120 rows, including Year and Zip Code.:

Year,Zip,TotalPopulation,MedianAge,Population12to17Years,Population0to5Years,Population18to24Years,Population25to34Years,Population35to44Years,Population45to54Years,Population55to64Years,Population6to11Years,Population65to74Years,Population75to84Years,Population85YearsandOlder,EmploymentTravelTime3059Min,EmploymentTravelTimee90Min,EmploymentTravelTime6089Min,EmploymentTravelTime1529Min,EmploymentTravelTimeLessthan15Min,EmploymentWorkatHome,EmploymentWalkedtoWork,EmploymentSubwayorElevatedtoWork,EmploymentCarTruckVantoWork,EmploymentFerrytoWork,EmploymentPublicTransportationtoWork,EmploymentRailroadtoWork,EmploymentTaxitoWork,EmploymentOtherTransportationtoWork,EmploymentBicycletoWork,EmploymentBusorTrolleyBustoWork,EducationCollegeMastersorDoctorateDegree,@#EducationAttainmentSomeCollege,@#EducationAttainmentBachelorsDegree,@#EducationEnrolledPublicSchool,EducationHighSchool_A,EducationAttainmentMastersDegree,EducationHighSchool,EducationAttainmentDoctorateDegree,@#EducationAttainmentAssociatesDegree,@#EducationEnrolledPrivateSchool,EmploymentCivilianTotalEmployed,EmploymentBlueColl

ar, EmploymentWhiteCollar, EmploymentWholesaleTrade, EmploymentArtsEntertainmentAccommodationandFoodServicesetc, EmploymentOtherServices, EmploymentFinanceandInsurance, EmploymentInformation, EmploymentManufacturing, EmploymentEducationalServices, EmploymentAgricultureForestryFishingandHunting, EmploymentConstruction, EmploymentProfessionalScientificandTechnicalServices, EmploymentRetailTrade, EmploymentTransportationandWarehousing, EmploymentRealEstateandRentalandLeasing, EmploymentManagementofCompaniesandEnterprises, EmploymentHealthCareandSocialAssistance, EmploymentArtsEntertainmentandRecreation, EmploymentAdministrativeandSupportandWasteMgt.Services, EmploymentUtilities, EmploymentPublicAdministration, EmploymentMiningQuarryingandOilandGasExtraction, AsianHouseholdIncomeMedian\$, AsianHouseholdIncomeAverage\$, AsianHouseholdIncomePerCapita\$, AsianHouseholdIncomeTotal\$, AsianHouseholdIncomeHighIncomeAverage\$, BlackHouseholdIncomeMedian\$, BlackHouseholdIncomeTotal\$, BlackHouseholdIncomeAverage\$, BlackHouseholdIncomePerCapita\$, FamilyInc.PerCapita\$, FamilyIncomeTotal\$, FamilyIncomeAverage\$, FamilyIncomeHighIncomeAverage\$, FamilyIncomeMedian\$, HispanicHouseholdIncomeTotal\$, HispanicHouseholdIncomeAverage\$, HispanicHouseholdIncomeMedian\$, HispanicHouseholdIncomePerCapita\$, HispanicHouseholdIncomeHighIncomeAverage\$, WhiteHouseholdIncomePerCapita\$, OtherRaceHouseholdIncomePerCapita\$, NonFamilyInc.Median\$, NonFamilyInc.Average\$, NonFamilyInc.PerCapita\$, NonFamilyInc.HighIncomeAverage\$, NonFamilyInc.Total\$, OtherRaceHouseholdIncomeMedian\$, OtherRaceHouseholdIncomeTotal\$, OtherRaceHouseholdIncomeAverage\$, OtherRaceHouseholdIncomeHighIncomeAverage\$, WhiteHouseholdIncomeTotal\$, EducationIndex, CultureIndex, ReligiousIndex, @#PopulationPop, FamilyPopulation, NonFamilyPopulation, PopulationDensityperSq.mile, AggravatedAssaultIndex, EASIQualityofLifeIndex, MotorVehicleTheftIndex, Forcible RapeIndex, MurderIndex, LarcenyIndex, ForcibleRobberyIndex, BurglaryIndex, EASITotalCrimeIndex, AsianPopulation, AsianHouseholds, BlackPopulation, BlackHouseholds, HispanicPopulation, HispanicHouseholds, OtherRaceHouseholds, OtherRacePopulation, WhiteHouseholds, WhitePopulation, WhiteNonHispanicHouseholds, WhiteNonHispanicPopulation, HouseholdswRent\$1000\$1249, HouseholdswRent\$1250\$1499, HouseholdswRent\$1500\$1999, HouseholdswRent\$2000, HouseholdswRent\$250\$499, HouseholdswRent\$500\$749, HouseholdswRent\$750\$999, HouseholdswRentlessthan\$250, HousingMedianRent\$, VacantUnits, VacantUnitsForRent, VacantUnitsForSale, VacantUnitsSeasonal, VacantUnitsOther, @#Entries, @#Lines, @#Stations

Attributes highlighted in green are from the Metropolitan Transportation Authority (MTA) while the rest were downloaded from Simply Map 3.0.

Initial preprocessing was done using Microsoft Excel. Macros were formulated to eliminate columns that held redundant or unnecessary information, such as state ID:

The screenshot shows an Excel spreadsheet titled "Zip Codes in New York, NY". A macro dialog box is open over the spreadsheet. The macro code is:

```



```

Below the macro code, there are several buttons: Run, Cancel, Step Into, Edit, Create, Delete, Options..., and Description.

Afterward, each attribute from Simply Map 3.0 needed to be transformed from its current format to a working version (see figures below).

http://newyork.ny.gov/nygovdata.com

New Map New Tabular Report 123 New Ranking

Time TransportType Education CPS QualityOfLife_Crime Location Analysis Location Analysis 1

DATA FILTERS DISPLAY OPTIONS ACTIONS

Selected location: New York, NY Change Analyze data by: Zip Codes

Location	Food Services (\$'000), 2016	Food Services (\$'000), 2015	Food Services (\$'000), 2014	Food Services (\$'000), 2013	Food Services (\$'000), 2012	Food Services (\$'000), 2011	Food Services (\$'000), 2010 (in year 2010)
11414, Howard Beach, NY	\$61,079.00	\$73,194.00	\$61,133.00	\$56,221.00	\$41,380.00	\$15,466.00	
11415, Kew Gardens, NY	\$6,974.00	\$7,803.00	\$9,333.00	\$50,109.00	\$7,193.00	\$6,591.00	
11416, Ozone Park, NY	\$12,036.00	\$14,462.00	\$24,133.00	\$14,657.00	\$23,598.00	\$12,773.00	
11417, Ozone Park, NY	\$59,684.00	\$62,104.00	\$71,000.00	\$0.00	\$13,439.00	\$17,205.00	
11418, Richmond Hill, NY	\$30,676.00	\$27,460.00	\$16,533.00	\$91,499.00	\$11,168.00	\$10,965.00	
11419, South Richmond Hill, NY	\$29,626.00	\$30,462.00	\$25,266.00	\$0.00	\$74,229.00	\$26,071.00	
11420, South Ozone Park, NY	\$12,053.00	\$11,612.00	\$11,400.00	\$0.00	\$10,285.00	\$8,457.00	
11421, Woodhaven, NY	\$21,417.00	\$16,374.00	\$15,567.00	\$6,868.00	\$22,021.00	\$17,831.00	
11422, Rosedale, NY	\$31,948.00	\$19,766.00	\$27,466.00	\$5,661.00	\$71,516.00	\$17,439.00	
11423, Hollis, NY	\$11,233.00	\$9,263.00	\$8,733.00	\$14,020.00	\$8,707.00	\$8,165.00	
11424, Jamaica, NY	\$0.00	\$0.00	\$0.00	\$77,037.00	\$0.00	\$0.00	
11425, Jamaica, NY	\$123.00	\$136.00	\$133.00	\$9,771.00	\$126.00	\$117.00	
11426, Bellmore, NY	\$10,369.00	\$9,981.00	\$8,800.00	\$16,057.00	\$8,139.00	\$8,749.00	
11427, Queens Village, NY	\$15,800.00	\$15,952.00	\$16,533.00	\$17,061.00	\$13,250.00	\$11,956.00	
11428, Queens Village, NY	\$7,036.00	\$8,742.00	\$7,267.00	\$15,082.00	\$5,742.00	\$2,566.00	
11429, Queens Village, NY	\$27,998.00	\$28,359.00	\$28,733.00	\$26,411.00	\$37,260.00	\$24,029.00	

Year integer_id	Zip integer	TotalPopul... real	Population... real	Population... real	Population... real	Population... real
1 2010	10001	38.300	2.277	3.441	14.421	20.747
2 2010	10002	39.800	5.732	4.923	7.317	18.632
3 2010	10003	34.100	1.669	2.909	18.445	27.626
4 2010	10004	32.200	1.225	4.096	10.008	44.552
5 2010	10005	31.700	0.578	4.113	14.446	43.882
6 2010	10006 2010	31.900	1.226	4.847	11.387	43.562
7 2010	10007	33.900	3.673	10.415	7.528	27.535
8 2010	10009	37.900	4.295	4.265	9.470	23.699
9 2010	10010	35.300	1.681	4.012	10.492	31.320
10 2010	10011	40.200	2.573	3.356	5.565	24.267
11 2010	10012	35.700	2.680	3.179	7.818	32.318
12 2010	10013	39.600	3.680	6.760	9.483	17.216

This transformation was completed using several different methods. The quickest most visually informative process was through Python and the Jupyter Notebook:

The screenshot shows a Jupyter Notebook interface. The top bar includes the Jupyter logo, file menu, edit, view, insert, cell, kernel, widgets, and help options. A toolbar below has icons for file operations like new, open, save, and cell execution. The Python logo is in the top right corner.

```

In [2]: import os
import csv
import pandas as pd

In [4]: base_dir = '/Users/T/science/datahelp/ReshapeData/VacantUnit/'
# base_dir = '/Users/T/science/datahelp/data/'
filename = 'VacantUnitOther1.csv'
fname, ext = os.path.splitext(filename)
fullpath = base_dir + fname
infile = open(fullpath, 'rU')
data_list = list(csv.reader(infile, dialect=csv.excel_tab, delimiter=','))
headers = data_list[0]
n_columns = len(headers)

# column idx in the original data
zip_idx = 0
zip_to_rows = dict()

zipcodes = []
n_errors = 0
for r, row in enumerate(data_list):
    if row[zip_idx] != headers[0]:
        zip_to_rows[row[zip_idx]] = r
        zipcodes.append(row[zip_idx])

    if len(row) != n_columns:
        n_errors += 1

if n_errors > 0:
    print str(n_errors) + ' didn\'t match'

unique_zips = sorted(list(set(zipcodes)))
print headers

```

['Zip', '% Vacant Units Other, 2016', '% Vacant Units Other, 2015', '% Vacant Units Other, 2014', '% Vacant Units Other, 2013', '% Vacant Units Other, 2012', '% Vacant Units Other, 2011', '% Vacant Units Other, 2010', '% Vacant Units Other, 2009', '% Vacant Units Other, 2008']

Several other methods of data preparation were experimented with and reported in our final presentation, but Python and Jupyter was the best method overall.

MTA attributes such as number of lines and stations per zip code had to be calculated from converted from latitude/longitude pairs (reverse geocoding). We tried several methods using the Google API, Texas A&M University's Geoservices, and Doogal. All methods are in more detail within the presentation.

Data Cleaning

Since some records have too many missing values, we decided to drop rows that have missing values more than 20% of number of columns. In order to do that, we used the function dropna from dataframe in Python.

```

df = df.dropna(thresh = 0.80 * len(df.columns))
# thresh means it requires that many non-NA values in that row/column

```

Getting Zip Code by Google API

Since we can only get longitude and latitude of subway entrances from NYC government. We need to get the zip code by using these data. To do that, we used Google API.

Some of the code for getting zipcode:

```

googlemap =
https://maps.googleapis.com/maps/api/geocode/json?key=AIzaSyBRJoEcdgq0V_gRkJiWpVS
pJUD4JodOc-Q&latlng=

for address in addresslist:
    mapurl = ""
    mapresult = []
    getzip = False
    mapurl = googlemap + address[1] + ',' + address[0]
    maplink = urllib.urlopen(mapurl)
    mapdata = json.loads(maplink.read())

    mapresult = mapdata['results'][0]['address_components']
    for x in mapresult:
        if 'postal_code' in x['types']:
            az = address,'\t',x['long_name']
            zipcodelist.append(az)
            getzip = True
            break

```

During the data preparation phase, we are going to use StringIndexer, OneHotEncoder, VectorAssembler and ChiSqSelector from the extracting, transforming and selecting features in Spark

MLlib: Main Guide

- Pipelines
- **Extracting, transforming and selecting features**
- Classification and Regression
- Clustering
- Collaborative filtering
- Model selection and tuning
- Advanced topics

MLlib: RDD-based API Guide

- Data types
- Basic statistics
- Classification and regression
- Collaborative filtering
- Clustering
- Dimensionality reduction
- Feature extraction and transformation
- Frequent pattern mining
- Evaluation metrics
- PMML model export
- Optimization (developer)

Extracting, transforming and selecting features

This section covers algorithms for working with features, roughly divided into these groups:

- Extraction: Extracting features from "raw" data
- Transformation: Scaling, converting, or modifying features
- Selection: Selecting a subset from a larger set of features

Table of Contents

- Feature Extractors
 - TF-IDF
 - Word2Vec
 - CountVectorizer
- Feature Transformers
 - Tokenizer
 - StopWordsRemover
 - n-gram
 - Binarizer
 - PCA
 - PolynomialExpansion
 - Discrete Cosine Transform (DCT)
 - **StringIndexer**
 - **IndexToString**
 - **OneHotEncoder**
 - VectorIndexer
 - Normalizer
 - StandardScaler
 - MinMaxScaler
 - MaxAbsScaler
 - Bucketizer
 - ElementwiseProduct
 - SQLTransformer
 - **VectorAssembler**
 - QuantileDiscretizer
- Feature Selectors
 - VectorSlicer
 - RFormula
 - **ChiSqSelector**

One Hot Encoding

What is One Hot Encoding?

One Hot means there's only 1 "hot" or "on" value in this list, while the rest are "cold" or "off".

It can transform categorical features to a format that works better with classification and regression algorithms.

A common application of One Hot Encoding is in NLP, where words don't exactly turn into vectors so easily.

For example:

foo	color
1	green
2	red
3	blue

After transform, our data will be like:

foo	color
1	[1,0,0]
2	[0,1,0]
3	[0,0,1]

(vector)

or

foo	color=green	color=red	color=blue
1	1	0	0
2	0	1	0
3	0	0	1

(dummy)

Why do we need it?

Any learner based on standard distance metrics (such as k-nearest neighbors) between samples will get confused without one-hot encoding. With the naive encoding and Euclidean distance, the distance between Green and Red is 1. The distance between Green and Blue is 2. But with the one-hot encoding, the pairwise distances between [1, 0, 0], [0, 1, 0] and [0, 0, 1] are all equal to $\sqrt{2}$.

However, this is not true for all models. Decision trees and derived models such as random forests, if deep enough, can handle categorical variables without one-hot encoding.

In NLP, one hot encoding is very common, transform word to vector.

Reference: <https://achyutjoshi.github.io/datascience/one-hot-encoding>

How to do it?

For Python, we can use:

Pandas get_dummies

http://pandas.pydata.org/pandas-docs/stable/generated/pandas.get_dummies.html

Sklearn feature_extraction DictVectorizer

<http://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.OneHotEncoder.html>

Spark MLlib OneHotEncoder

<http://spark.apache.org/docs/latest/ml-features.html#onehotencoder>

Is Zip code a quantitative or categorical variable?

Although zip codes are written in numbers, the numbers are simply convenient labels and don't have numeric meaning (for example, we wouldn't add together two zip codes or calculate average of two zip codes)

<http://www.dummies.com/education/math/statistics/how-to-distinguish-quantitative-and-categorical-variables/>

How about year?

Same as zip code, it's categorical variable. Since we are not going to calculate anything by year, and in this project, year means interval of time.

<https://www.quora.com/Is-year-a-quantitative-or-categorical-variable>

Implementation in our project:

In our project, we decided to use Spark and since we have both Year and Zip code, we need to do encoding for twice.

First, we need to use StringIndexer to index both Year and Zip code attributes.

Second, we use OneHotEncoder to transform attributes to vector, since in MLlib, most of the models requires one vector as features.

```
import pandas as pd
from pyspark import SparkContext
from pyspark.sql import SQLContext
from pyspark.ml.feature import OneHotEncoder, StringIndexer

# Creates Spark Context
sc = SparkContext()
sql_sc = SQLContext(sc)
pandas_df = pd.read_csv('data.csv', na_values=[' '])
spark_df = sql_sc.createDataFrame(pandas_df)

# Step 1 - StringIndexer
firstIndexer = StringIndexer(inputCol="Year", outputCol="YearIndex")
firstModel = firstIndexer.fit(spark_df)
firstIndexed = firstModel.transform(spark_df)
secondIndexer = StringIndexer(inputCol="Zip", outputCol="ZipIndex")
secondModel = secondIndexer.fit(firstIndexed)
secondIndexed = secondModel.transform(firstIndexed)

# Step 2 - OneHotEncoder
yearEncoder = OneHotEncoder(dropLast=False, inputCol="YearIndex", outputCol="YearVec")
zipEncoder = OneHotEncoder(dropLast=False, inputCol="ZipIndex", outputCol="ZipVec")
yearEncoded = yearEncoder.transform(secondIndexed)
zipEncoded = zipEncoder.transform(yearEncoded)
final_df = zipEncoded
final_df.select("YearIndex","YearVec","ZipIndex","ZipVec").show()
```

Year and zip attributes after encoding:

YearIndex	YearVec	ZipIndex	ZipVec
4.0	(7,[4],[1.0])	47.0	(182,[47],[1.0])
4.0	(7,[4],[1.0])	20.0	(182,[20],[1.0])
4.0	(7,[4],[1.0])	153.0	(182,[153],[1.0])
4.0	(7,[4],[1.0])	115.0	(182,[115],[1.0])
4.0	(7,[4],[1.0])	87.0	(182,[87],[1.0])
4.0	(7,[4],[1.0])	30.0	(182,[30],[1.0])
4.0	(7,[4],[1.0])	163.0	(182,[163],[1.0])
4.0	(7,[4],[1.0])	21.0	(182,[21],[1.0])
4.0	(7,[4],[1.0])	25.0	(182,[25],[1.0])
4.0	(7,[4],[1.0])	102.0	(182,[102],[1.0])
4.0	(7,[4],[1.0])	53.0	(182,[53],[1.0])
4.0	(7,[4],[1.0])	8.0	(182,[8],[1.0])
4.0	(7,[4],[1.0])	146.0	(182,[146],[1.0])
4.0	(7,[4],[1.0])	66.0	(182,[66],[1.0])
4.0	(7,[4],[1.0])	44.0	(182,[44],[1.0])
4.0	(7,[4],[1.0])	134.0	(182,[134],[1.0])
4.0	(7,[4],[1.0])	170.0	(182,[170],[1.0])
4.0	(7,[4],[1.0])	26.0	(182,[26],[1.0])
4.0	(7,[4],[1.0])	107.0	(182,[107],[1.0])
4.0	(7,[4],[1.0])	70.0	(182,[70],[1.0])

only showing top 20 rows

As we can see, the first number of YearVec and ZipVec means the distinct values in this attribute. We can compare with the result from Weka.

Selected attribute			
Name: Year		Type: Numeric	
Selected attribute			
Name: Zip	Missing: 0 (0%)	Distinct: 7	Type: Numeric
Missing: 0 (0%)	Distinct: 182	Unique: 0 (0%)	Unique: 1 (0%)

The One Hot Encoding result from Spark is correct.

Feature Selection

The feature selection algorithms in Spark MLlib includes VectorSlicer, RFormula and ChiSqSelector. In our project, we are going to use ChiSqSelector.

What is ChiSqSelector?

ChiSqSelector implements Chi-Squared feature selection. It operates on labeled data with categorical features. ChiSqSelector orders features based on a Chi-Squared test of independence from the class, and then filters (selects) the top features which the class label depends on the most. This is akin to yielding the features with the most predictive power.

How to use it?

First, we need to assemble the features into one vector.

```

assembler = VectorAssembler(inputCols = df_columns, outputCol = "features")
output = assembler.transform(spark_df)
regModelData = output.select("features", "HousingMedianRent$")

```

Second, we apply the data into ChiSqSelector, in this case, we keep 50 features.

```

selector = ChiSqSelector(numTopFeatures = 50, featuresCol = "features", outputCol = "selectedFeatures", labelCol = "HousingMedianRent$")
result = selector.fit(output).transform(output)
model = selector.fit(output)

```

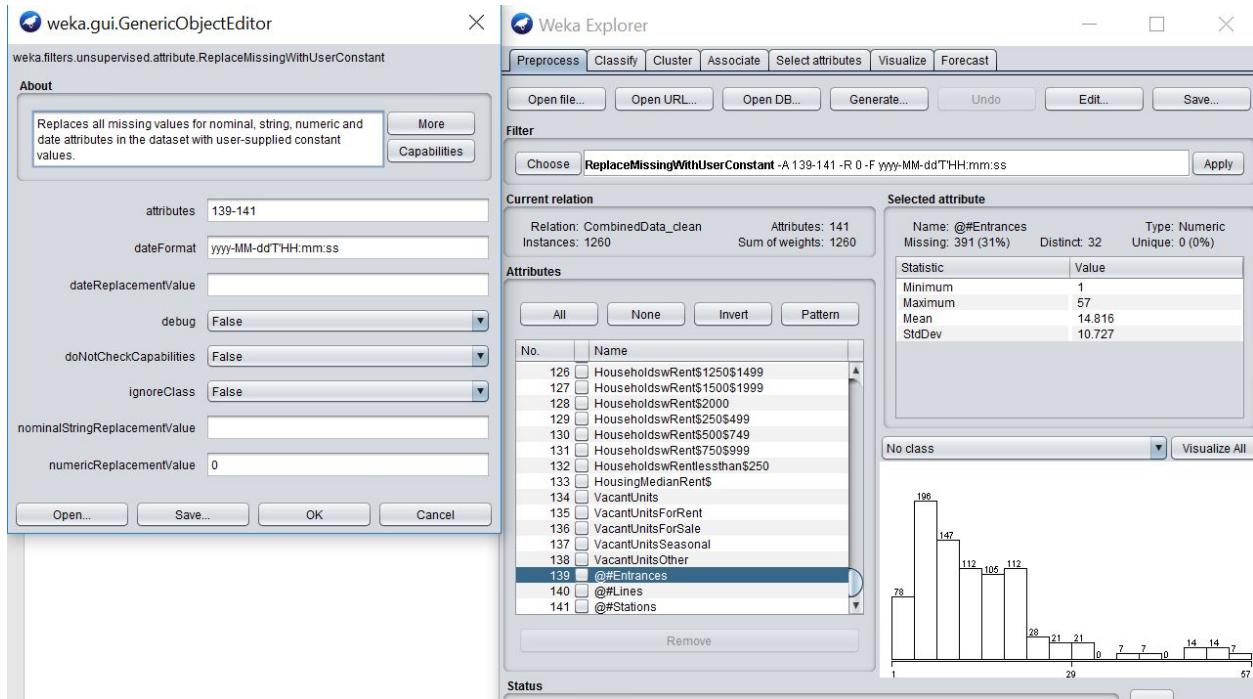
The result from ChiSqSelector

ChiSqSelector with top 50 features selected:

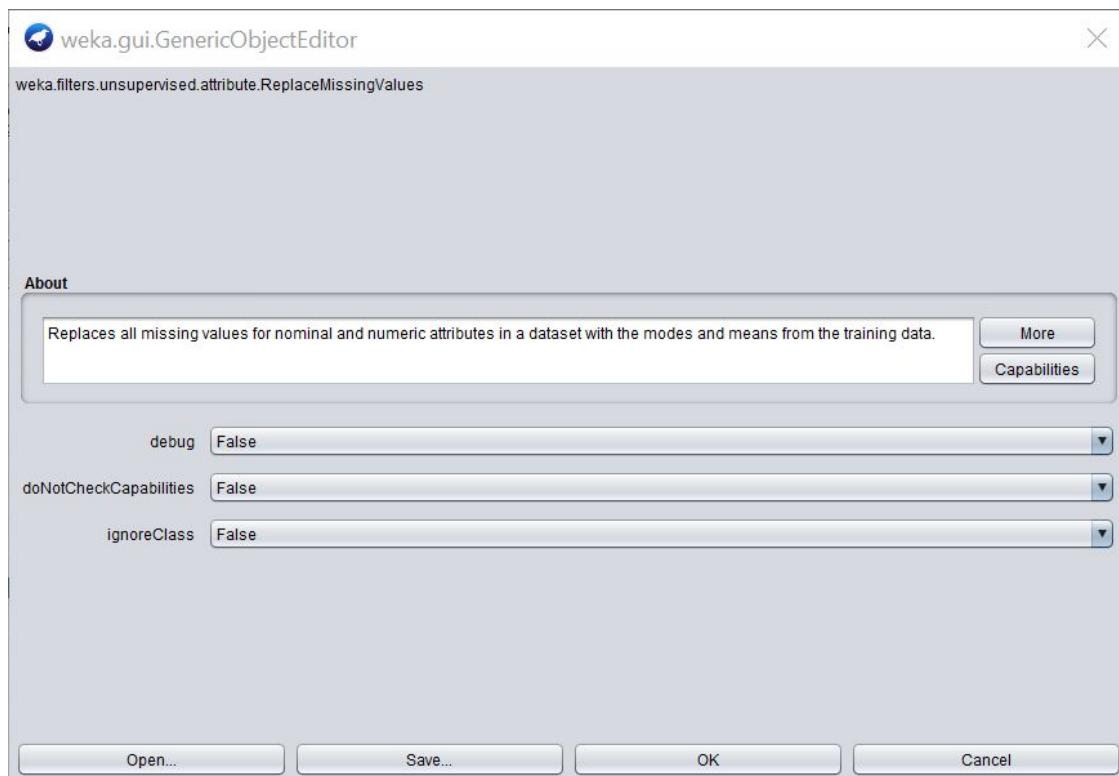
Population12to17Years	AsianHouseholdIncomeAverage\$
Population0to5Years	AsianHouseholdIncomeTotal\$
Population18to24Years	FamilyIncomeTotal\$
Population6to11Years	FamilyIncomeAverage\$
Population65to74Years	FamilyIncomeMedian\$
Population75to84Years	HispanicHouseholdIncomeTotal\$
EmploymentTravelTime3059Min	HispanicHouseholdIncomeAverage\$
EmploymentTravelTime90Min	NonFamilyIncAverage\$
EmploymentTravelTime6089Min	NonFamilyIncTotal\$
EmploymentTravelTime1529Min	WhiteHouseholdIncomeTotal\$
EmploymentTravelTimeLessthan15Min	FamilyPopulation
EmploymentWalkedtoWork	PopulationDensitypersq-mile
EmploymentSubwayorElevatedtoWork	AsianHouseholds
EmploymentCarTruckVantoWork	BlackPopulation
EducationCollegeMastersorDoctorateDegree	BlackHouseholds
EducationHighSchool_A	HispanicPopulation
EducationAttainmentMastersDegree	HispanicHouseholds
EducationHighSchool	OtherRaceHouseholds
EmploymentWhiteCollar	OtherRacePopulation
EmploymentArtsEntertainmentAccommodationandFoodServices	WhiteHouseholds
etc	WhitePopulation
EmploymentFinanceandInsurance	WhiteNonHispanicHouseholds
EmploymentEducationalServices	WhiteNonHispanicPopulation
EmploymentProfessionalScientificandTechnicalServices	HouseholdswRent\$1500\$1999
EmploymentHealthCareandSocialAssistance	VacantUnits
EmploymentPublicAdministration	

Handling missing values with Weka

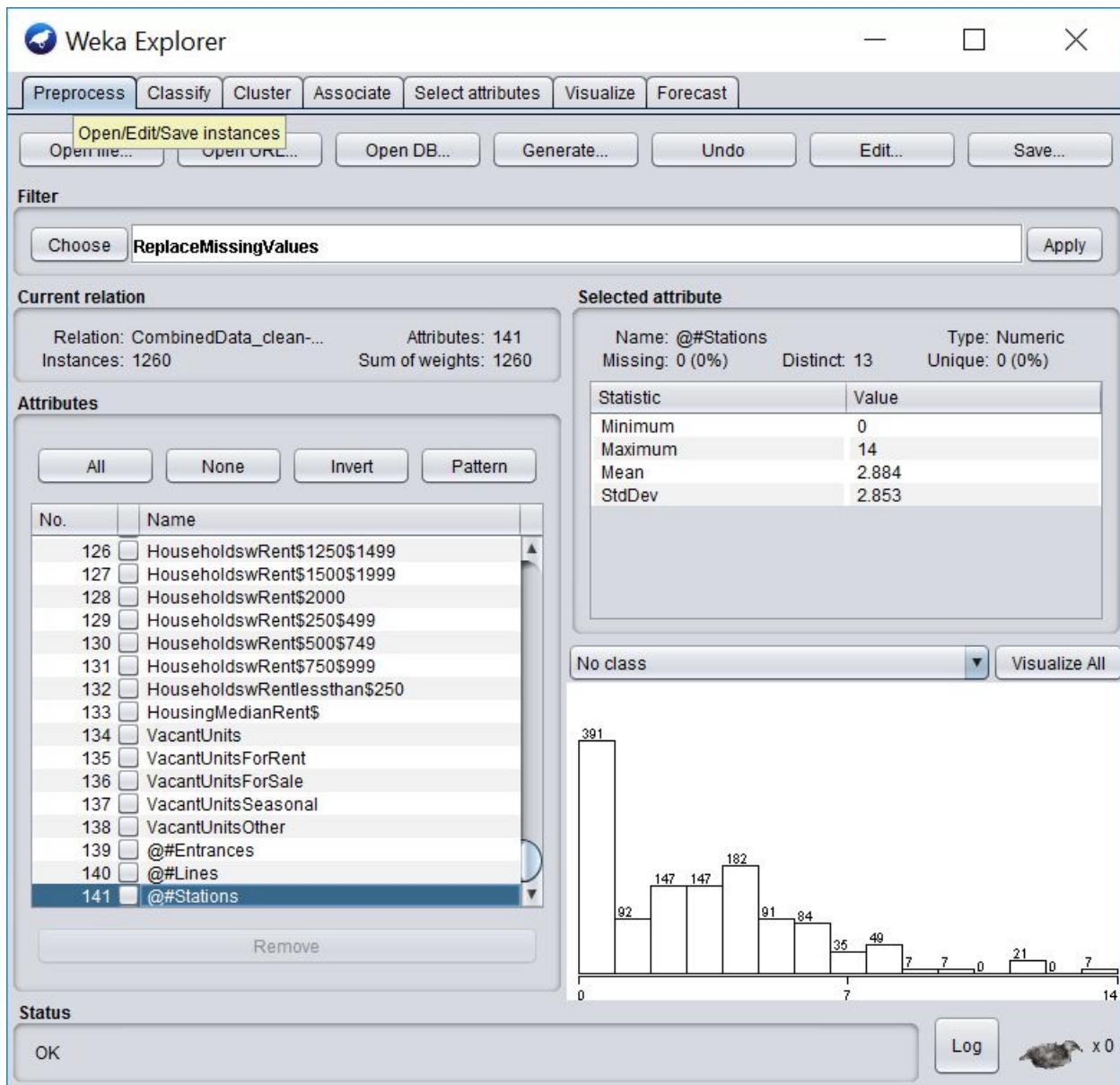
For attributes #Entrance, #Lines and #Stations we apply “ReplaceMissingWithUserConstant” filter and replace all the missing values with zero. Therefore, these three attributes have no longer missing values.



Then we apply “ReplaceMissingValues” to all attributes:

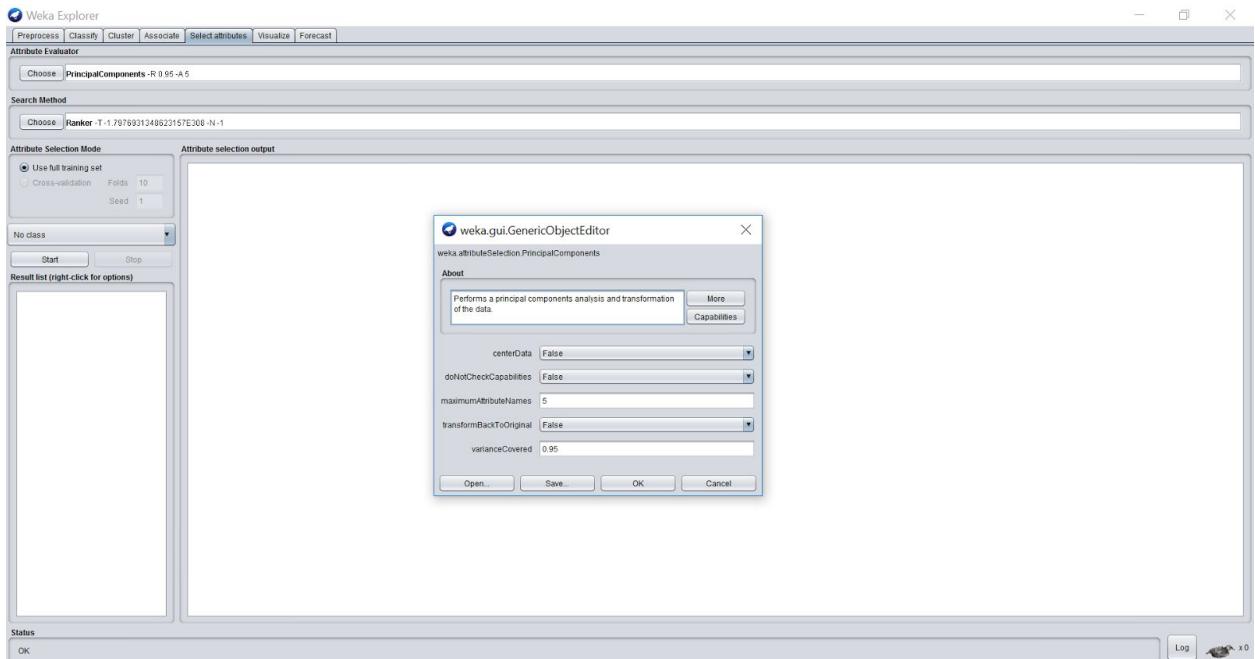


After replacing missing values:



Applying PCA on the data

Apply PCA by setting “VarianceCovered” to 0.95 (default):



Information

NAME
weka.attributeSelection.PrincipalComponents

SYNOPSIS
Performs a principal components analysis and transformation of the data. Use in conjunction with a Ranker search. Dimensionality reduction is accomplished by choosing enough eigenvectors to account for some percentage of the variance in the original data—default 0.95 (95%). Attribute noise can be filtered by transforming to the PC space, eliminating some of the worst eigenvectors, and then transforming back to the original space.

OPTIONS
centerData – Center (rather than standardize) the data. PCA will be computed from the covariance (rather than correlation) matrix

transformBackToOriginal – Transform through the PC space and back to the original space. If only the best n PCs are retained (by setting varianceCovered < 1) then this option will give a dataset in the original space but with less attribute noise.

varianceCovered – Retain enough PC attributes to account for this proportion of variance.

maximumAttributeNames – The maximum number of attributes to include in transformed attribute names.

doNotCheckCapabilities – If set, evaluator capabilities are not checked before evaluator is built (Use with caution to reduce runtime).

Correlation Matrix:

Eigen Values:

The screenshot shows the Weka Explorer interface with the following details:

- Top Bar:** Weka Explorer, Preprocess, Classify, Cluster, Associate, Selected attributes, Visualize.
- Attribute Evaluator:** Choose PrincipalComponents - R 0.65 - A 5
- Search Method:** Choose Ranker - T-1.797893134823157E300 - N-1
- Attribute Selection Mode:** Use full training set (radio button selected), Cross-validation Folds 10, Seed 1.
- Attribute selection output:** Shows the first 10 components with their eigenvalues, proportions, and cumulative values. The first component has an eigenvalue of 41.24247 and a proportion of 0.29675, contributing 0.151 to the total variance.
- Result list (right-click for options):** 0205-30 - Ranker + PrincipalComponent
- Table Data:** Displays the correlation matrix for the first 10 principal components. The matrix is as follows:

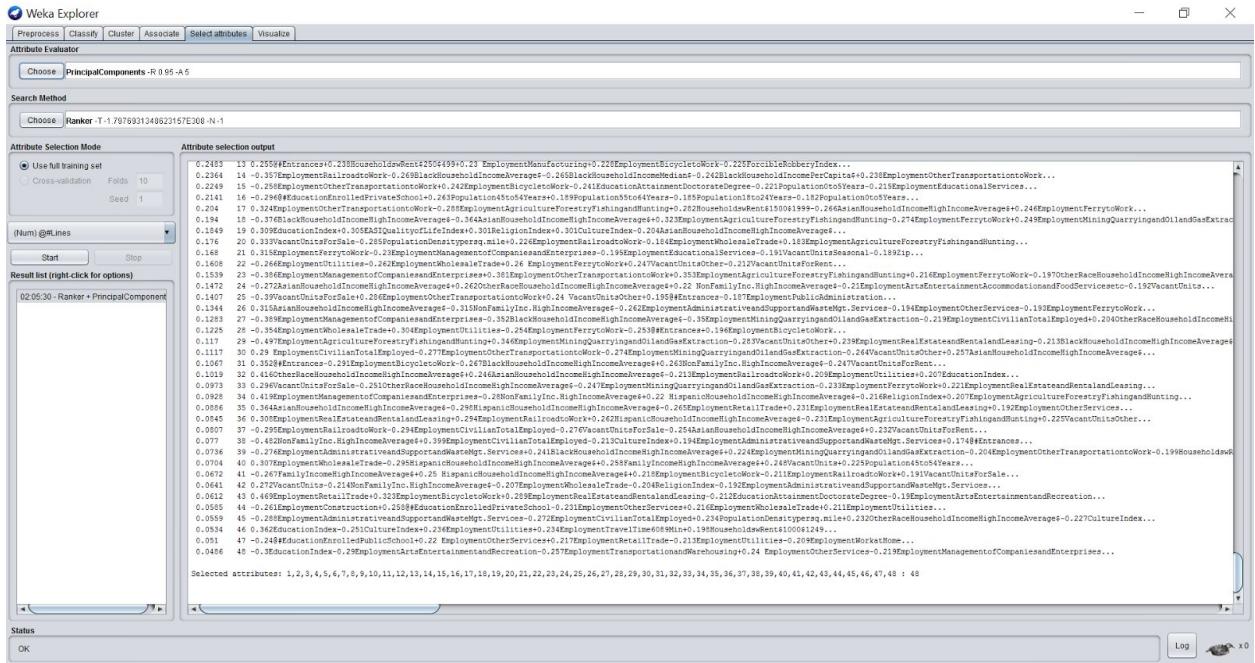
	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9	PC10
PC1	1.00000	-0.2088	-0.1990	-0.1940	-0.1920	-0.1900	-0.1880	-0.1860	-0.1840	-0.1820
PC2	-0.2088	1.00000	-0.1900	-0.1880	-0.1860	-0.1840	-0.1820	-0.1800	-0.1780	-0.1760
PC3	-0.1990	-0.1900	1.00000	-0.1860	-0.1840	-0.1820	-0.1800	-0.1780	-0.1760	-0.1740
PC4	-0.1940	-0.1880	-0.1860	1.00000	-0.1820	-0.1800	-0.1780	-0.1760	-0.1740	-0.1720
PC5	-0.1920	-0.1860	-0.1840	-0.1820	1.00000	-0.1780	-0.1760	-0.1740	-0.1720	-0.1700
PC6	-0.1900	-0.1840	-0.1820	-0.1800	-0.1780	1.00000	-0.1740	-0.1720	-0.1700	-0.1680
PC7	-0.1880	-0.1820	-0.1800	-0.1780	-0.1760	-0.1740	1.00000	-0.1700	-0.1680	-0.1660
PC8	-0.1860	-0.1800	-0.1780	-0.1760	-0.1740	-0.1720	-0.1700	1.00000	-0.1660	-0.1640
PC9	-0.1840	-0.1780	-0.1760	-0.1740	-0.1720	-0.1700	-0.1680	-0.1660	1.00000	-0.1620
PC10	-0.1820	-0.1760	-0.1740	-0.1720	-0.1700	-0.1680	-0.1660	-0.1640	-0.1620	1.00000

Eigen Vectors:

Weka Explorer																															
Preprocess Classify Cluster Associate Selected attributes Visualize																															
Attribute Evaluator																															
Choose	PrincipalComponents -R 0.95 -A 5																														
Search Method																															
Choose	Ranker -T 1.7976931248523157E308 -N 1																														
Attribute Selection Mode																															
<input checked="" type="radio"/> Use full training set																															
Cross-validation	Folds 10																														
Seed	1																														
(Num) @fLines																															
Start	Stop																														
Result list (right-click for options)																															
02:05:30 - Ranker + PrincipalComponent																															
Attribute selection output																															
Eigenvectors																															
V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	V14	V15	V16	V17	V18	V19	V20	V21	V22	V23	V24	V25	V26	V27	V28	V29			
0.0005	-0.0263	0.0712	-0.0746	-0.1961	0.2522	-0.0584	0.1463	-0.199	0.0258	0.0049	0.0804	0.0398	0.0866	-0.0408	0.0376	-0.0059	0.1336	0.0686	-0.0758	-0.0215	-0.0105	0.0338	-0.0307	0.051	0.0374	0.0304	-0.0246	0.004			
0.0124	-0.0084	0.0637	-0.0686	-0.2005	0.2506	-0.0377	0.1666	-0.1928	0.032	0.0141	0.09	0.0278	0.1013	-0.0507	0.0491	0.0047	0.1254	0.0581	-0.0767	0.0051	-0.0118	0.0368	-0.0427	0.0419	0.0436	0.0357	-0.0139	0.03			
0.0138	-0.0084	0.0637	-0.0686	-0.2005	0.2506	-0.0377	0.1666	-0.1928	0.032	0.0141	0.09	0.0278	0.1013	-0.0507	0.0491	0.0047	0.1254	0.0581	-0.0767	0.0051	-0.0118	0.0368	-0.0427	0.0419	0.0436	0.0357	-0.0139	0.03			
0.0423	-0.1524	0.0396	0.1455	0.048	0.173	0.1059	-0.0589	0.0623	0.155	-0.0284	-0.0048	0.015	0.0239	0.0531	0.1404	0.0104	-0.0389	0.0203	0.0577	0.0344	-0.0374	0.0033	0.0237	0.02	0.0553	-0.0245	-0.0033	0.005			
0.0138	-0.0312	-0.0427	0.0047	-0.0456	-0.0465	0.0469	0.0494	-0.0544	0.0218	0.0661	-0.109	0.0609	-0.1244	-0.1015	-0.0786	-0.2144	0.0346	0.0209	0.0241	-0.0012	-0.0147	0.0078	-0.0267	0.0194	0.0207	-0.0193	0.0419	0.019			
-0.0784	0.0304	-0.0265	-0.106	-0.0641	-0.1721	0.0289	-0.0664	-0.1341	0.0202	0.0231	-0.0378	0.1747	0.0069	-0.2209	-0.1824	-0.0526	0.1048	0.0350	0.0465	0.0538	-0.0273	0.1395	0.058	0.1067	-0.1586	0.0162	-0.0098	0.005			
-0.0416	0.1242	-0.0314	-0.0531	-0.0368	0.0602	-0.0781	0.2685	0.0006	-0.0217	0.039	0.1429	-0.1019	0.066	0.0384	-0.1854	0.0146	-0.0971	-0.0798	0.0658	-0.1495	-0.0394	-0.0111	0.0098	-0.0972	0.1562	0.0594	0.0693	0.03			
0.0491	0.0941	0.0186	-0.0282	0.0203	0.047	0.0943	0.0367	0.0179	0.0176	0.0099	0.1764	-0.232	-0.1081	0.0105	0.2626	0.0477	-0.0217	0.0409	-0.0255	-0.0239	0.0555	0.0522	0.0007	0.1183	-0.0525	-0.0232	0.0433	0.005			
0.0581	-0.1223	0.0194	0.0555	0.0117	-0.0083	-0.0499	0.0179	0.1767	0.0176	0.0099	0.1764	-0.232	-0.1081	0.0105	0.2626	0.0477	-0.0217	0.0409	-0.0255	-0.0239	0.0555	0.0522	0.0007	0.1183	-0.0525	-0.0232	0.0433	0.005			
0.0148	-0.1191	0.0796	0.1652	-0.0448	0.0973	0.0367	0.0187	0.0164	-0.0065	0.1184	-0.1332	-0.0681	-0.0127	-0.0587	0.1892	-0.219	0.0332	0.0616	0.1238	-0.1179	0.0454	0.0326	0.0755	0.0504	0.1064	-0.1832	0.0239	-0.0211	0.005		
-0.0126	-0.0196	-0.0282	-0.0203	-0.0447	-0.0943	0.0367	0.0399	-0.1053	0.04	0.0487	0.1068	0.102	-0.1111	-0.1205	-0.0236	0.0878	0.0221	0.0138	0.0383	0.0341	0.0455	0.0421	0.0191	0.0099	0.0284	0.01					
0.0211	0.0085	0.0154	0.0449	0.0246	0.0244	0.0141	-0.0448	0.0318	0.0449	0.0218	0.0099	0.1764	-0.232	-0.1081	0.0105	0.2626	0.0477	-0.0217	0.0409	-0.0255	-0.0239	0.0555	0.0522	0.0007	0.1183	-0.0525	-0.0232	0.0433	0.005		
0.0211	0.0085	0.0154	0.0449	0.0246	0.0244	0.0141	-0.0448	0.0318	0.0449	0.0218	0.0099	0.1764	-0.232	-0.1081	0.0105	0.2626	0.0477	-0.0217	0.0409	-0.0255	-0.0239	0.0555	0.0522	0.0007	0.1183	-0.0525	-0.0232	0.0433	0.005		
0.0239	-0.0749	0.0577	0.1098	0.0543	0.2576	0.167	-0.0957	-0.0009	-0.1749	-0.2019	0.0983	0.0748	-0.0229	0.0504	-0.0567	0.0776	-0.1016	0.0304	0.1374	0.1189	-0.0566	0.0563	0.1565	-0.0506	0.0312	0.0573	0.0632	0.005			
-0.0241	0.0581	0.0772	-0.0243	-0.0512	0.1045	-0.1245	-0.3276	-0.0662	0.1212	0.0822	-0.101	0.0301	0.0187	-0.0703	0.0703	0.0371	-0.0142	-0.0216	0.0985	0.0236	0.0649	0.1131	-0.0791	-0.1626	-0.0134	-0.1052	-0.0474	-0.0524	-0.01		
-0.0047	-0.1284	0.0968	0.0111	0.0717	-0.0396	-0.1818	0.1211	0.0139	0.0052	-0.0055	-0.0262	0.0415	0.0606	0.1098	0.0186	0.0306	0.0487	-0.1309	0.0669	0.0487	0.0391	0.1321	0.0184	0.0443	0.0765	0.02	0.0573	0.005			
0.0227	0.0145	-0.0434	-0.0466	-0.0208	0.0779	-0.0533	0.0394	0.1159	0.0618	0.0349	-0.0266	0.0499	0.0281	0.0687	0.0618	0.0281	0.0401	-0.0944	0.0442	0.0442	-0.1412	0.0348	0.1247	0.0247	0.0573	0.005	0.0573	0.005			
0.0252	0.0063	0.0446	-0.0109	0.1218	-0.0413	0.1192	0.2478	0.0116	-0.134	-0.1113	0.1228	-0.039	-0.0395	0.0242	-0.1027	0.0397	-0.0206	0.0796	0.0285	-0.0203	0.0303	0.0273	-0.028	0.0303	0.0273	-0.028	0.0303	0.0273	0.005		
0.0981	0.1118	0.1104	0.0165	-0.0058	0.0669	0.0463	0.0504	0.1485	0.0936	-0.1111	0.0722	0.0563	0.0465	0.0244	0.0423	-0.0054	0.0426	0.0328	0.0398	0.031	0.0442	-0.0173	0.1082	0.0102	0.019	0.0274	-0.0155	0.0421	0.11		
0.0398	0.1127	-0.0069	-0.0228	0.0299	0.0449	0.0403	0.1485	0.0936	-0.1111	0.0722	0.0563	0.0465	0.0244	0.0423	-0.0054	0.0426	0.0328	0.0398	0.031	0.0442	-0.0173	0.1082	0.0102	0.019	0.0274	-0.0155	0.0421	0.11			
0.0454	0.0118	0.0284	0.0285	0.0463	0.0463	0.0488	0.0835	-0.0749	0.0445	0.0504	0.0423	0.0338	0.0133	0.0123	0.0188	0.0445	0.0445	0.0504	0.0423	0.0338	0.0123	0.0124	0.0123	0.0124	0.0123	0.0124	0.0123	0.0124	0.0123		
-0.0144	-0.0475	0.0578	0.0041	0.0266	-0.0447	0.0721	0.1467	-0.1029	0.0124	0.046	0.1108	0.0139	0.0384	0.1678	-0.0618	0.0201	0.0465	-0.2743	0.0324	0.0264	0.0316	0.0257	0.0324	0.0257	0.0324	0.0257	0.0324	0.0257			
-0.0189	0.1189	0.0349	-0.0242	-0.0558	0.2269	-0.0532	-0.2256	0.0455	-0.0155	-0.0338	0.0542	0.0058	0.0788	-0.0751	0.0061	-0.0669	-0.0537	0.0788	0.0325	0.0988	-0.0985	-0.0278	0.0682	-0.0111	0.0045	0.0197	-0.0232	0.005			
-0.0339	-0.0871	-0.0184	-0.0127	0.0202	0.0732	-0.0538	0.0541	0.0395	0.0624	0.0711	-0.1485	0.0539	-0.0357	0.0668	0.058	0.0417	0.2336	-0.162	0.2259	0.0282	-0.1058	0.1707	0.1144	-0.1451	0.0809	-0.13	-0.11				
0.0170	0.0513	0.0513	0.0466	0.0466	0.0466	0.0466	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465	0.0465			
0.0202	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184	0.0184			
0.0321	0.0983	0.0922	0.0131	0.0398	0.0618	-0.0723	-0.0767	-0.1614	0.0124	0.1558	0.0668	0.228	-0.1203	0.2424	0.0579	-0.0073	0.1228	-0.0498	-0.1228	0.1386	0.0139	0.1134	0.0292	0.0193	0.0292	0.0193	0.0292	0.0193			
-0.0944	-0.0777	-0.0858	0.0543	-0.0777	-0.0598	0.1055	0.112	0.0713	0.0166	0.0331	-0.0423	0.0484	-0.0513	0.1498	-0.0373	-0.0253	0.1125	0.0832	-0.0051	0.0303	-0.1483	-0.0017	0.0707	-0.0577	-0.01	0.0045	0.0209	0.0045			
0.1508	0.0105	0.0098	0.0127	0.0076	0.0062	0.0039	0.0011	0.0076	0.0227	0.0118	0.0341	0.0041	0.1609	0.0484	-0.0513	0.1498	-0.0373	-0.0253	0.1125	0.0832	-0.0051	0.0303	-0.1483	-0.0017	0.0707	-0.0577	-0.01	0.0045	0.0209	0.0045	
0.0531	0.0205	0.0216	0.176	-0.0959	-0.0812	0.0041	-0.0342	-0.0118	-0.0234	-0.0098	-0.0099	0.0413	0.0238	-0.0323	-0.0595	0.034	0.0037	0.1247	0.0832	-0.0051	0.0303	-0.1483	-0.0017	0.0707	-0.0577	-0.01	0.0045	0.0209	0.0045		
0.0205	0.0567	0.2316	0.176	-0.0959	-0.0812	0.0041	-0.0342	-0.0118	-0.0234	-0.0098	-0.0099	0.0413	0.0238	-0.0323	-0.0595	0.034	0.0037	0.1247	0.0832	-0.0051	0.0303	-0.1483	-0.0017	0.0707	-0.0577	-0.01	0.0045	0.0209	0.0045		
-0.1	0.058	0.1091	0.0371	-0.1288	-0.1041	0.0031	-0.0093	-0.0189	-0.0457	-0.1291	0.0869	-0.0699	-0.0051	-0.075	0.0507	0.0406	0.0574	-0.1263	0.0487	-0.0394	-0.1149	0.0444	-0.0414	-0.045	-0.0221	0.005	0.0444	-0.0414	-0.045	-0.0221	0.005
-0.124	0.0869	0.0375	-0.0749	0.0373	0.0185	0.007	0.0354	0.0231	-0.0029	-0.0123	0.0943	0.0033	0.0014	-0.0441	-0.038	0.0184	0.0279	0.0303	0.0051	0.0303	-0.0143	0.0247	-0.0097	0.031	0.0427	0.0596	0.0382	-0.0454	0.0054		

Ranked attributes:

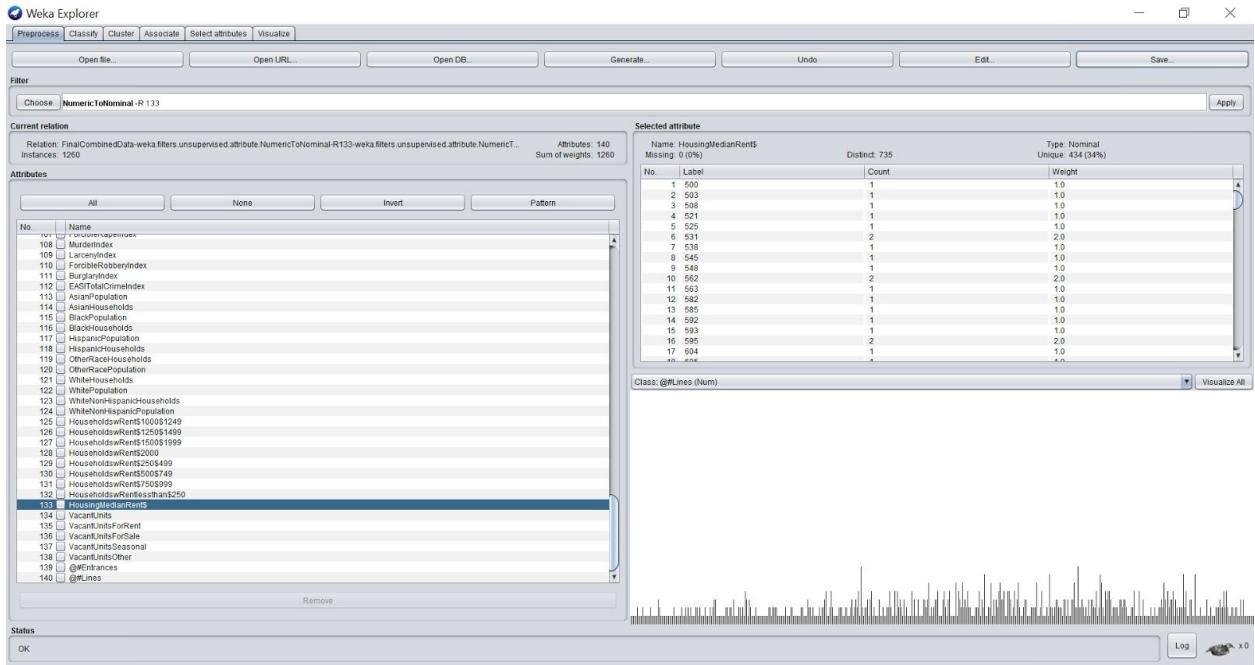
We can see only 48 attributes out of 140 attributes are selected:



After applying PCA, we lost interpretability of our data but since PCA returned 48 attributes, we decided to select “50” most important attributes in feature selection part.

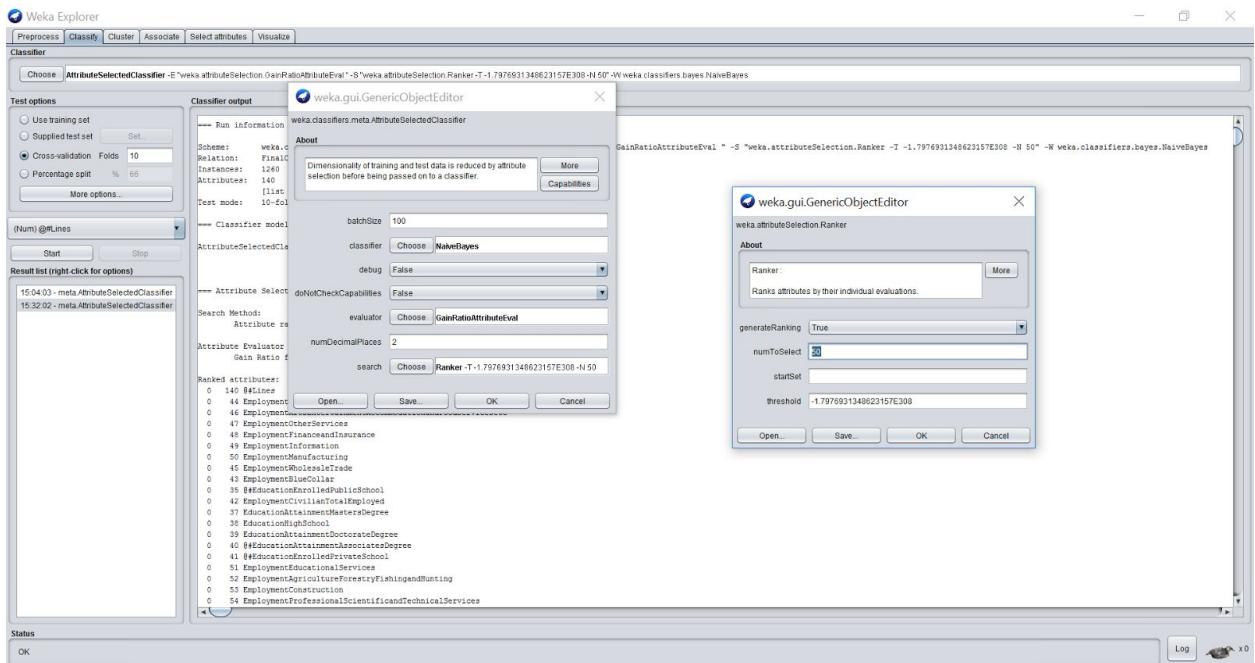
Attribute Selection with Weka

Attribute subset selection involves subset evaluation measure, and search method. But we have 140 attributes so searching is slow! An alternative could be using a single-attribute evaluator with ranking which evaluates each attribute individually instead of evaluating subsets of attributes. An advantage of this approach is that it can eliminate irrelevant attributes. Although, it can't eliminate redundant attributes. We choose “HousingMedianRent\$” attribute as the class. Since the data for this attribute is “Numeric”, we need to convert it to “Nominal”. We apply this transformation with “NumericToNominal” filter in Weka:

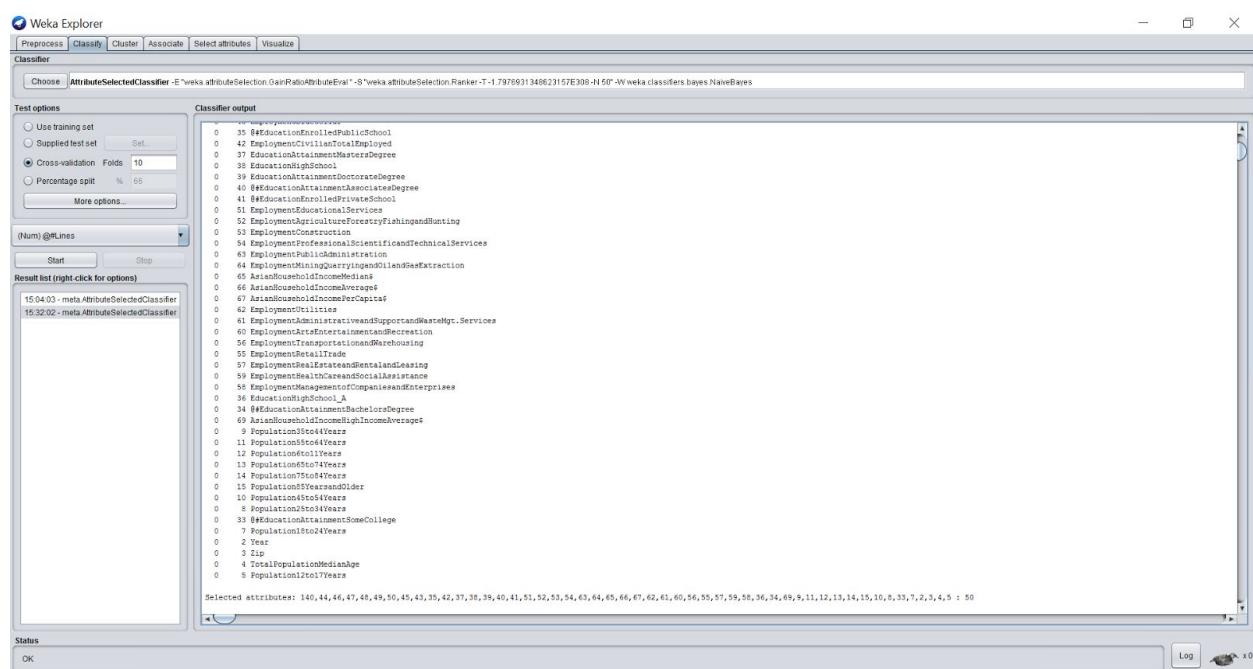
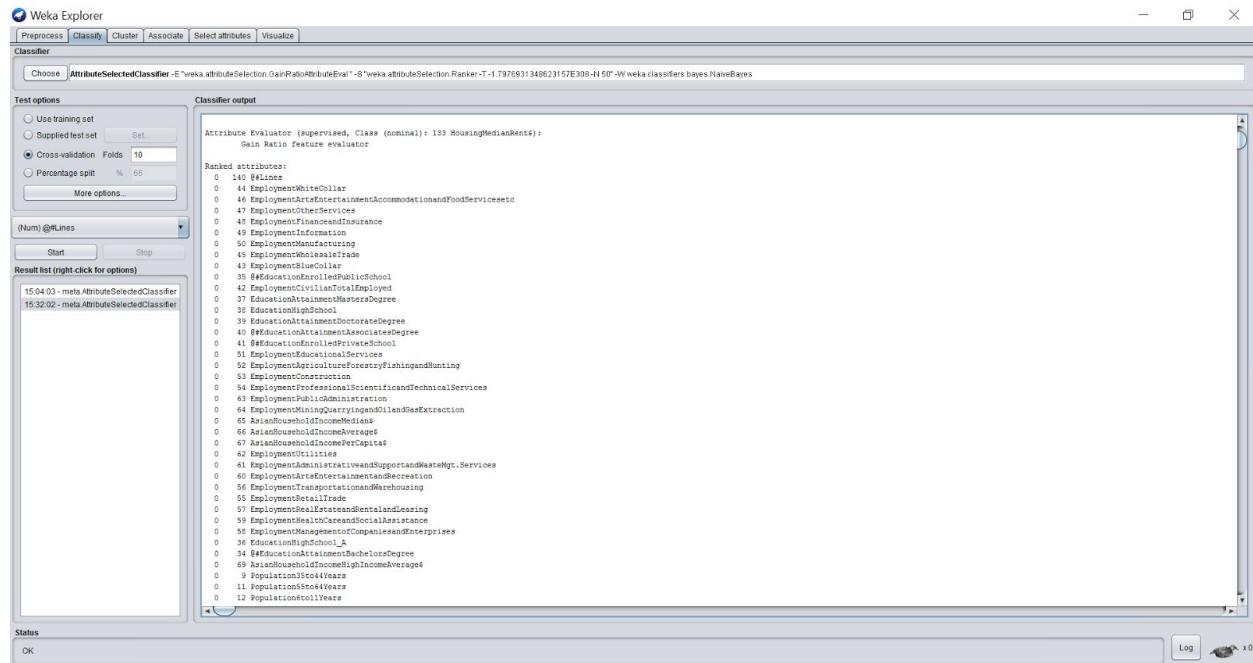


Gain Ratio Attribute Evaluator:

We use “AttributeSelectedClassifier” to get an evaluation of the attributes. We choose NaiveBayes as the classifier, Gain Ratio Attribute Evaluator for the evaluator with “Ranker” search method. For “Ranker”, we set the parameter “numToSelect” to 50 which outputs 50 best ranked attributes. We choose 50 since PCA returned back 50 attributes after applying data reduction. The default threshold parameter for “Ranker” is minus infinity (which is equal to the shown parameter in Java). We use “Cross-Validation” with 10 “Folds” for the evaluation.



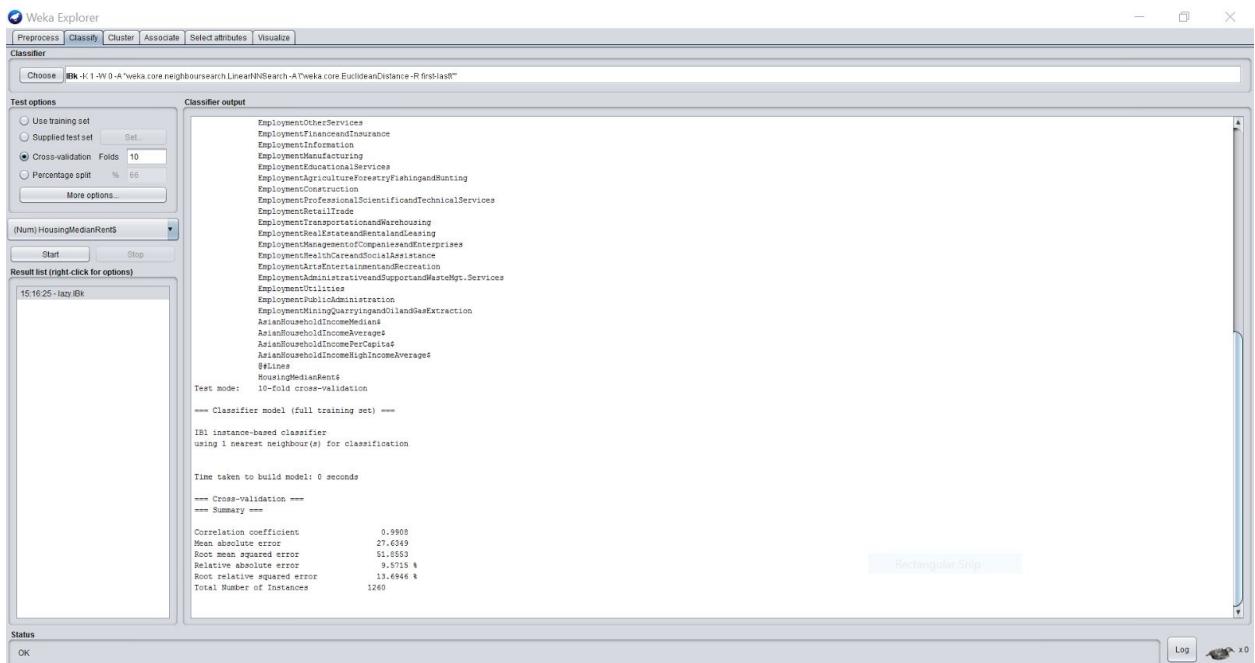
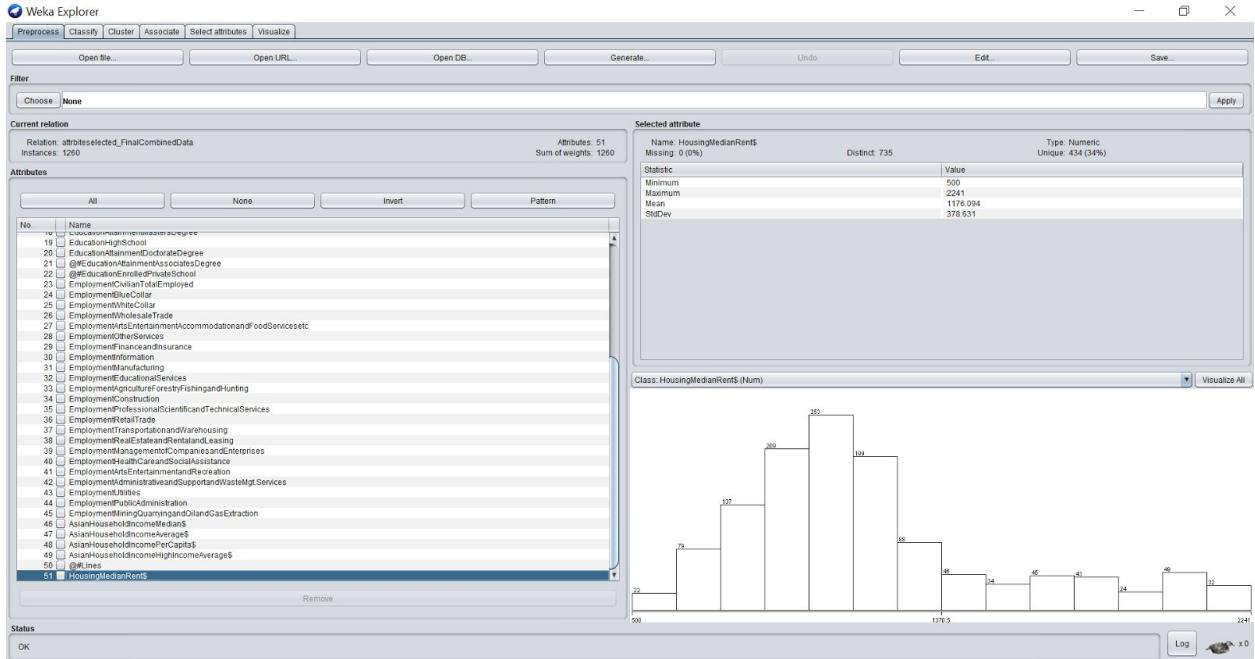
50 selected attributes are as following. As we see the most important selected attribute is “number of stations” which confirms our assumption that rental prices are directly dependent on access to the subway.



After selecting the 50 most important attributes, we create a data set with these attributes and convert the data of attribute “HousingMedianRent” to original (to avoid losing information). We choose “HousingMedianRent” as the class and then we apply IBK-lazy classifier, and as we see:

Relative absolute error: 9.5715 %

Root relative squared error: 13.6946 %



Modeling

Modeling with Spark

In this phase, we are going to use models from MLlib, there are Linear Regression, Decision Tree Regression and Random Forest Regression from its Classification and Regression.

MLlib: Main Guide

- Pipelines
- Extracting, transforming and selecting features
- **Classification and Regression**
- Clustering
- Collaborative filtering
- Model selection and tuning
- Advanced topics

MLlib: RDD-based API Guide

- Data types
- Basic statistics
- Classification and regression
- Collaborative filtering
- Clustering
- Dimensionality reduction
- Feature extraction and transformation
- Frequent pattern mining
- Evaluation metrics
- PMML model export
- Optimization (developer)

Classification and regression

This page covers algorithms for Classification and Regression. It also includes sections discussing specific classes of algorithms, such as linear methods, trees, and ensembles.

Table of Contents

- Classification
 - Logistic regression
 - Decision tree classifier
 - Random forest classifier
 - Gradient-boosted tree classifier
 - Multilayer perceptron classifier
 - One-vs-Rest classifier (a.k.a. One-vs-All)
 - Naive Bayes
- Regression
 - **Linear regression**
 - Generalized linear regression
 - Available families
 - **Decision tree regression**
 - **Random forest regression**
 - Gradient-boosted tree regression
 - Survival regression
 - Isotonic regression
 - Examples
- Linear methods
- Decision trees
 - Inputs and Outputs
 - Input Columns
 - Output Columns
- Tree Ensembles
 - Random Forests
 - Inputs and Outputs
 - Input Columns
 - Output Columns (Predictions)
 - Gradient-Boosted Trees (GBTs)

Linear Regression in MLlib

What is Linear Regression?

Linear regression attempts to model the relationship between two variables by fitting a linear equation to observed data. One variable is considered to be an explanatory variable, and the other is considered to be a dependent variable.

Some code for Linear Regression

```
sc = SparkContext()
sql_sc = SQLContext(sc)
pandas_df = pd.read_csv('data.csv', na_values=[' '])
spark_df = sql_sc.createDataFrame(pandas_df)

df_columns = spark_df.columns
df_columns = df_columns[3:-1]
assembler = VectorAssembler(inputCols = df_columns, outputCol = "features")
vector_df = assembler.transform(spark_df)
data = vector_df.select("features", "HousingMedianRent$")
(trainingData, testData) = data.randomSplit([0.7, 0.3])

lr = LinearRegression(labelCol = "HousingMedianRent$")
model = lr.fit(trainingData, {lr.regParam:0.3})
predictions = model.transform(testData)
```

Result

```
+-----+-----+-----+
|      features | HousingMedianRent$ |      prediction |
+-----+-----+-----+
|[29.6,6.82562,10....|          800|834.1419710054527|
|[29.8,9.00797,9.8...|          639|621.5723793103944|
|[30.0,9.33527,9.4...|          811|794.5769589562694|
|[30.0,10.4396,9.9...|          637|737.6132118800572|
|[30.2,9.37296,10....|          781|758.7192914819996|
+-----+-----+-----+
only showing top 5 rows
```

Linear Regression result:
Root Mean Squared Error (RMSE): 50.3514
Mean Squared Error (MSE): 2535.26
Mean Absolute Error (MAE): 36.049
R-squared: 0.981682

Linear Regression result:

Root Mean Squared Error (RMSE): 50.3514
Mean Squared Error (MSE): 2535.26
Mean Absolute Error (MAE): 36.049
R-squared: 0.981682

The result is very good. RMSE is not high and R-squared is close to 1.

What is RMSE?

RMSE measures how much error there is between two datasets. RMSE usually compares a predicted value and an observed value.

What is R squared?

R-squared is a statistical measure of how close the data are to the fitted regression line. It is also known as the coefficient of determination, or the coefficient of multiple determination for multiple regression. The definition of R-squared is the percentage of the response variable variation that is explained by a linear model, therefore R-squared is always between 0 and 100%.

R-squared = Explained variation/Total variation

0% indicates that the model explains none of the variability of the response data around its mean.

100% indicates that the model explains all the variability of the response data around its mean.

Decision Tree in MLlib

What is Decision Tree?

Decision trees and their ensembles are popular methods for the machine learning tasks of classification and regression. Decision trees are widely used since they are easy to interpret, handle categorical features, extend to the multiclass classification setting, do not require feature scaling, and can capture non-linearity and feature interactions. Tree ensemble algorithms such as random forests and boosting are among the top performers for classification and regression tasks.

The spark.ml implementation supports decision trees for binary and multiclass classification and for regression, using both continuous and categorical features.

Regression vs Classification

Regression: the output variable takes continuous or numeric values.

Classification: the output variable takes class labels.

In our case, we are going to use regression since we want to predict the rent (numeric) in the future.

Some code for Linear Regression

```
sc = SparkContext()
sql_sc = SQLContext(sc)
pandas_df = pd.read_csv('data.csv', na_values=[' '])
spark_df = sql_sc.createDataFrame(pandas_df)

df_columns = spark_df.columns
df_columns = df_columns[3:-1]
assembler = VectorAssembler(inputCols = df_columns, outputCol = "features")
vector_df = assembler.transform(spark_df)
data = vector_df.select("features", "HousingMedianRent$")
featureIndexer = VectorIndexer(inputCol="features", outputCol="indexedFeatures",
maxCategories=4).fit(data)
```

```
(trainingData, testData) = data.randomSplit([0.7, 0.3])

dt = DecisionTreeRegressor(featuresCol="indexedFeatures",
labelCol="HousingMedianRent$")
pipeline = Pipeline(stages=[featureIndexer, dt])
model = pipeline.fit(trainingData)
predictions = model.transform(testData)
```

Result

features	HousingMedianRent\$	prediction
[29.4, 9.95204, 14....]		891 879.2083333333334
[29.9, 8.8449, 9.77....]		644 783.9322033898305
[29.9, 9.36187, 9.5....]		804 879.2083333333334
[30.0, 9.33527, 9.4....]		811 879.2083333333334
[30.0, 10.4396, 9.9....]		637 536.6666666666666

only showing top 5 rows

Decision Tree Regression result:

Root Mean Squared Error (RMSE): 93.332
Mean Squared Error (MSE): 8710.86
Mean Absolute Error (MAE): 60.4492
R-squared: 0.944005
DecisionTreeRegressionModel (uid=DecisionTreeRegressor_4218b6c52f525a349c98) of depth 5 with 63 nodes

Decision Tree Regression result:

Root Mean Squared Error (RMSE): 93.332
Mean Squared Error (MSE): 8710.86
Mean Absolute Error (MAE): 60.4492
R-squared: 0.944005
DecisionTreeRegressionModel (uid=DecisionTreeRegressor_4218b6c52f525a349c98) of depth 5 with 63 nodes

The result is not as good as Linear Regression but still good.

How about trying Decision Tree Classifier?

In Decision Tree Classifier, we need to index label attribute (StringIndexer) into categorical variables.

```
# difference between Classifier and Regression
labelIndexer = StringIndexer(inputCol="HousingMedianRent$",
outputCol="indexedLabel").fit(data)
```

The result shows

```

+-----+-----+
|      features|indexedLabel|prediction|
+-----+-----+
|[29.4,10.0699,14....|      72.0|     2.0|
|[29.6,6.82562,10....|    278.0|     1.0|
|[30.0,10.4396,9.9...|    254.0|     1.0|
|[30.2,9.37296,10....|    452.0|     2.0|
|[30.4,9.55198,9.9...|    539.0|     1.0|
+-----+
only showing top 5 rows

```

Test Error = 0.97486

Test Error = 0.97486

Accuracy is very low, because we cannot change numeric into different classes to build the Decision Tree model.

Random Forest in MLlib

What is Random Forest?

A random forest is a meta estimator that fits several decision trees on various sub-samples of the dataset and use averaging to improve the predictive accuracy and control over-fitting.

In the code for Random Forest, we can just change the model to RandomForestRegressor in Decision Tree Regression.

```
dt = RandomForestRegressor(featuresCol="indexedFeatures", labelCol="HousingMedianRent$")
```

Result

```

+-----+-----+
|      features|HousingMedianRent$|      prediction|
+-----+-----+
|[29.6,6.94177,10....|      787| 765.367689330459|
|[29.8,9.00797,9.8...|      639| 727.6617914602946|
|[30.0,9.33527,9.4...|      811| 830.9129240992834|
|[30.0,10.4431,9.9...|      646| 684.7464470317793|
|[30.1,9.40755,10....|      773| 793.1010332199692|
+-----+
only showing top 5 rows

```

Random Forest Regression result:

Root Mean Squared Error (RMSE): 59.0543

Mean Squared Error (MSE): 3487.41

Mean Absolute Error (MAE): 39.9294

R-squared: 0.976438

RandomForestRegressionModel (uid=rfr_19429ebf3ba5) with 20 trees

Random Forest Regression result:

Root Mean Squared Error (RMSE): 59.0543
 Mean Squared Error (MSE): 3487.41
 Mean Absolute Error (MAE): 39.9294
 R-squared: 0.976438
 RandomForestRegressionModel (uid=rfr_19429ebf3ba5) with 20 trees

In the result, we can see when printing Random Forest, MLlib will show how many trees in this forest. As for the accuracy, Random Forest has better result than Decision Tree in every aspect.

Evaluation

Weka

	Correlation Coefficient	Mean Absolute Error	Root Mean Squared Error	Relative Absolute Error	Root Relative Squared Error
IBK-lazy	0.9908	28.0579	52.0729	9.718%	13.7521%

Figure (X) above shows the feature selection with the IBK-lazy (Weka's K-nearest neighbor) classifier. From the correlation coefficient, a sufficient positive correlation was shown between the top 50 attributes and the class label (Median Rent Price). Both RAE and RRSE were relatively low. However, MAE and RMSE are hard to evaluate without other classifiers to serve as a comparison. Future work will include at least one or two other classifiers in Weka's feature selection process.

Apache Spark

Compare results from different models

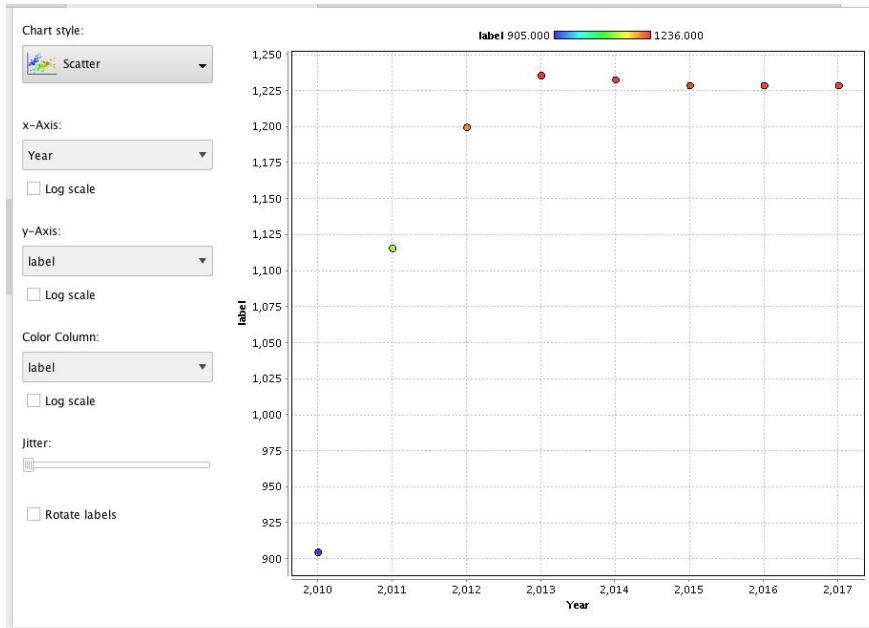
	RMSE	MSE	MAE	R-squared
Linear Regression	50.3514	2535.26	36.049	0.981682
Decision Tree	93.332	8710.86	60.4492	0.944005
Random Forest	59.0543	3487.41	39.9294	0.976438

As we can see, all three models have pretty good result, in the future, we are going to apply more models and keep testing our data.

The figure above depicts values from three different classifiers/regressions. Linear Regression outperformed both Decision Tree and Random Forest in all four measures of accuracy. Decision

Tree was notably the worst out of the three in all measures. Further investigation will be done in why a regression model outperformed a classifier.

RapidMiner



In Figure (X), the last data point is the predicted value of the Time Series analysis in RapidMiner. The analysis did not show a significant change in zip code 11211 (Williamsburg, Brooklyn) where we expected to see the most change in rent due to its reliance on the L train (Bedford Avenue is the first stop in Brooklyn on the L line). However, when we tried KNN, we believe there were not enough data points to properly evaluate or carry out the prediction. From these conclusions, we believe that Time Series, in theory, is the appropriate modeling approach but that it should be repeated using a software that can handle composite key values (year and zip code). By splitting up the zip codes, we lost a lot of critical information that rely on the interactions between multiple zip codes. Going forward, tools from the Apache ecosystem will be considered.

Deployment

Given our results, we believe our experiment should be repeated for more conclusive results using the proper tools that can handle a large dataset coupled with a predictive analytics question. Our modeling results should be considered pilot data with hypothesis generating ideas. Results from PCA and ChiSqSelector partially confirmed our hypothesis, in which we expected transportation to have a significant impact on rent decisions. Going forward, we would like to implement Time Series analysis in another software or tool, using all zip codes, years,

and attributes. These future results can be visualized as percent change and median rent prices overlaying a New York City map split by zip codes.