# ANALYSIS OF
# NYC SERVICE REQUEST CALLS TO 311

## *Outline*

- Data ETL
- MVA
  - CA, PCA, MCA
  - Clustering, Tree classification
- Conclusions

By: **Cedric Bhihe** <cedric.bhihe@gmail.com>  **Santi Calvo** <s.calvo93@gmail.com>

Data from:
- NYC OpenData
- NYPD DB
- IRS
- US Census Office

*All material and code available at:*
    https://www.github.com/Cbhihe/nyc311

**Legend:**
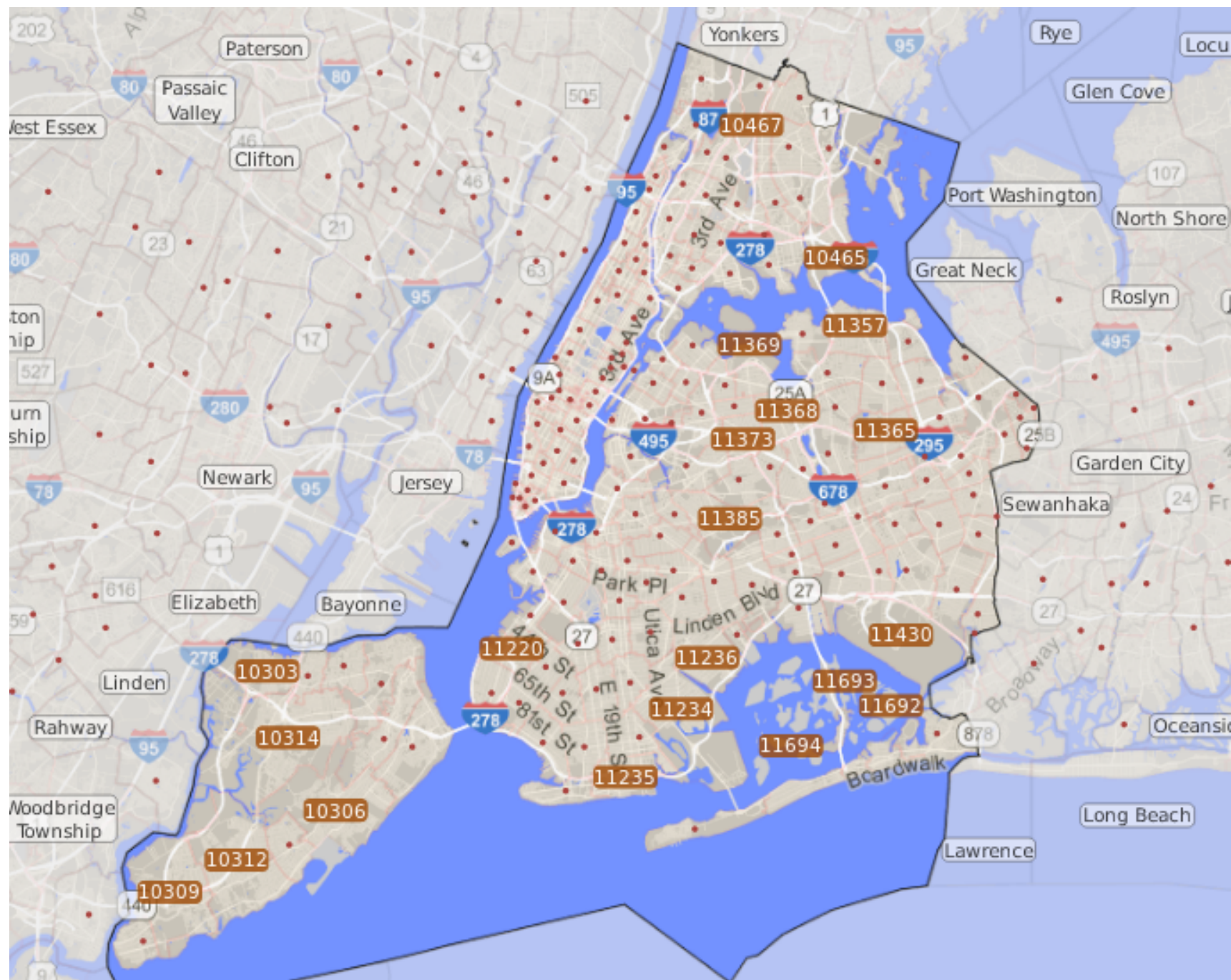1. Manhattan
2. Brooklyn
3. Queens
4. The Bronx
5. Staten Island

**5 boroughs:**

⇒ > 200 ZIPs

⇒ 8.7 M. people

⇒ 100,000 SRC
    to 311 /month

_Objective(s)_:

- ■ to explore data with MVA tools
- ■ to extract features and descriptive information, so we may:
    - detect trends
    - optimize urban resources

# Data ETL – extract

| ZIP | Date | Complaint | Descriptor | Address | planeX | planeY | GPS |
|---|---|---|---|---|---|---|---|
| 10463 | 2014-04-01 | PAINT/PLASTER | CEILING | 1 JACOBUS PLACE | 1008544 | 258129 | (40.875145400143914, -73.9121540... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | ENTIRE BUILDING | 3631 CORLEAR AVENUE | 1011190 | 261705 | (40.884952723690205, -73.9025718... |
| 10038 | 2014-04-01 | GENERAL | VENTILATION SYSTEM | 41 JOHN STREET | 982002 | 197700 | (40.70931729970831, -73.9128285342... |
| 10463 | 2014-04-01 | FLOORING/STAIRS | FLOOR | 110 TERRACE VIEW AVE | 1008357 | 258602 | (40.876444151331, -74.00810843... |
| 10306 | 2014-04-01 | Rodent | Condition Attracting Rodents | 1934 NORTH RAILROAD | 948358 | 144556 | (40.563376021503466, -74.0087726... |
| 10007 | 2014-04-01 | UNSANITARY CONDITION | GARBAGE/RECYCLING STORA | 23 PARK PLACE | 981818 | 199058 | (40.713044638283755, -73.8975242... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | ENTIRE BUILDING | 3810 BAILEY AVENUE | 1012586 | 261482 | (40.87422023300874, -73.9049668{... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | ENTIRE BUILDING | 2840 BAILEY AVENUE | 1010532 | 257794 | (40.884336286878096, -73.9043180... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | ENTIRE BUILDING | 3810 BAILEY AVENUE | 1012586 | 260028 | (40.880351339751236, -74.00810843... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | APARTMENT ONLY | 225 WEST 232 STREET | 1010709 | 197700 | (40.70931729970831, -73.9043180... |
| 10463 | 2014-04-01 | WATER LEAK | APARTMENT ONLY | 41 JOHN STREET | 982002 | 260028 | (40.875145400143914, -73.9121540... |
| 10038 | 2014-04-01 | PAINT/PLASTER | HEAVY FLOW | 225 WEST 232 STREET | 1010709 | 258129 | (40.884336286878096, -73.9121540... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | CEILING | 1 JACOBUS PLACE | 1008544 | 261482 | (40.875145400143914, -73.9128285342... |
| 10463 | 2014-04-01 | WATER LEAK | ENTIRE BUILDING | 3810 BAILEY AVENUE | 1012586 | 258129 | (40.876444151331, -73.9121540... |
| 10463 | 2014-04-01 | PLUMBING | SLOW LEAK | 1 JACOBUS PLACE | 1008544 | 258602 | (40.75145400143914, -73.9128285342... |
| 10463 | 2014-04-01 | PAINT/PLASTER | BASIN/SINK | 110 TERRACE VIEW AVE | 1008357 | 258602 | (40.875145400143914, -73.8943148... |
| 10463 | 2014-04-01 | WATER LEAK | WALL | 110 TERRACE VIEW AVE | 1008357 | 258129 | (40.876444151331, -73.8975242... |
| 10463 | 2014-04-01 | UNSANITARY CONDITION | HEAVY FLOW | 1 JACOBUS PLACE | 1008544 | 258602 | (40.885343438331205, -73.90886308... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | SLOW LEAK | 1 JACOBUS PLACE | 1008357 | 261850 | (40.884336286878096, -73.9017708... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | MOLD | 110 TERRACE VIEW AVE | 1013473 | 261482 | (40.88602554641576, -73.9034516... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | ENTIRE BUILDING | 140 VAN CORTLANDT A | 1012586 | 262094 | (40.878790215257844, -73.9128285342... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | ENTIRE BUILDING | 3810 BAILEY AVENUE | 1009450 | 259460 | (40.871850303847744, -73.8975242... |
| 10463 | 2014-04-01 | Rodent | ENTIRE BUILDING | 3555 OXFORD AVENUE | 1011414 | 256931 | (40.876444151331, -73.8975242... |
| 10463 | 2014-04-01 | PAINT/PLASTER | APARTMENT ONLY | 3150 BAILEY AVENUE | 1010952 | 258602 | (40.884336286878096, -73.9128285342... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | Rat Sighting | 2738 KINGSBRIDGE TER | 1008357 | 261482 | (40.884336286878096, -73.1291780... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | CEILING | 110 TERRACE VIEW AVE | 1012586 | 261482 | (40.876444151331, -73.93511991... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | ENTIRE BUILDING | 3810 BAILEY AVENUE | 1012586 | 258602 | (40.563376021503466, -73.8975242... |
| 10463 | 2014-04-01 | Rodent | ENTIRE BUILDING | 3810 BAILEY AVENUE | 1008357 | 144556 | (40.844453523297704, -74.1074178... |
| 10463 | 2014-04-01 | PAINT/PLASTER | ENTIRE BUILDING | 3810 BAILEY AVENUE | 948358 | 246971 | (40.884336286878096, -74.1177974... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | DOOR | 110 TERRACE VIEW AVE | 1002201 | 261482 | (40.558770684605705, -73.95728767... |
| 10463 | 2014-04-01 | HEAT/HOT WATER | Condition Attracting Rodents | 1934 NORTH RAILROAD | 1012586 | 142870 | (40.57748284313505, -73.97828340... |
| 10463 | 2014-04-01 | DOOR/WINDOW | Rat Sighting | 534 WEST 175TH STREE | 954402 | 149691 | (40.58351622059898, -73.95458695... |
| 10463 | 2014-04-01 | Rodent | ENTIRE BUILDING | 3810 BAILEY AVENUE | 951527 | 151870 | (40.77303218808042, -73.99406471... |
| 10463 | 2014-04-01 | Rodent | Rat Sighting | 2341 RICHMOND ROAD | 996114 | 220914 | (40.61617004369414, ... |
| 10306 | 2014-04-01 | HEAT/HOT WATER | Loud Music/Party | 2775 EAST 12 STREET | 990265 | 163767 | (40.72997723507379, ... |
| 10033 | 2014-04-01 | Rodent | Banging/Pounding | | 996858 | 205227 | ... |
| 10463 | 2014-04-01 | Noise - Residential | Noise: Construction Before/After Hours (NM1) | | 985895 | | ... |
| 10306 | 2014-04-01 | Noise - Residential | Noise: Alarms (NR3) | | | | |
| 10306 | 2014-04-01 | Noise | Noise: Construction Before/After Hours (NM1) | | | | |
| 10306 | 2014-04-01 | Noise | | | | | |
| 11235 | 2014-04-01 | Noise | | | | | |
| 10023 | 2014-04-01 | Noise | | | | | |
| 11230 | 2014-04-01 | Noise | | | | | |
| 10003 | 2014-04-01 | Noise | | | | | |

# Data ETL – reduce

| Period | Raw data's obs number | Obs # with missing ZIP | Obs # missing all location info | Service requests' modalities # | Unique ZIP |
|---|---|---|---|---|---|
| April 2014 | 81645 | 3206 | 2740 | 170 | 278 |
| April 2015 | 101890 | 4231 | 3069 | 178 | 260 |

- Between 150 and 200 different SRCs' raw features

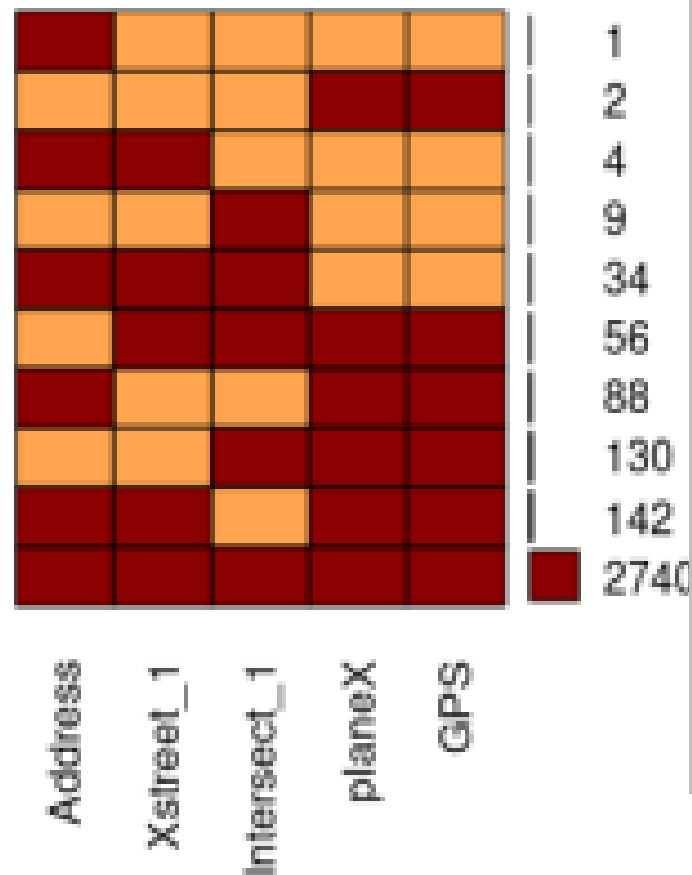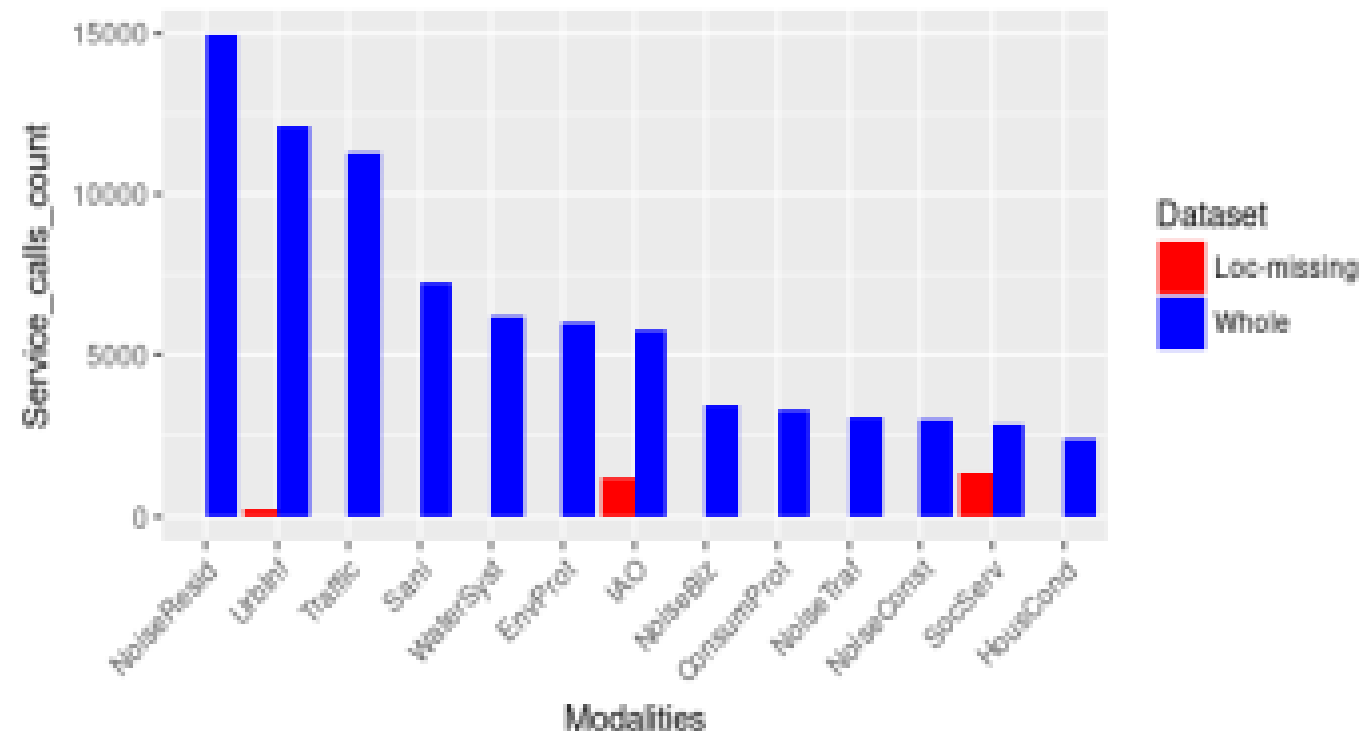- Reduce to 13 features by consolidating calls' objects

**Data ETL – reduce**

| Service request calls' modalities | Modality description | Service request call frequencies | | Change in rank from 2014 to 2015 |
|---|---|---|---|---|
| | | April 2014 | April 2015 | |
| *NoiseResid* | Residential Noise | 19.00% | 17.50% | – |
| UrbInf | Urban Infrastructure | 15.00% | 13.40% | ↘ |
| Traffic | Traffic related Issues | 14.30% | 17.20% | ↗ |
| Sani | Unsanitary Conditions | 9.20% | 10.50% | – |
| WaterSyst | Water Systems | 7.80% | 7.60% | – |
| EnvProt | Environmental Protection | 7.60% | 5.90% | – |
| IAO | Inspect, Audit, Order | 5.80% | 5.20% | ↘ |
| *NoiseBiz* | Commercial Noise | 4.40% | 4.90% | ↘ |
| ConsumProt | Comsumer Protection | 4.20% | 3.40% | ↘ |
| *NoiseTraf* | Traffic Noise | 3.90% | 5.40% | ↗↗ |
| *NoiseConst* | Construction Noise | 3.80% | 3.70% | ↗ |
| HousCond | Housing Conditions | 3.10% | 3.40% | – |
| SocServ | Social Services | 1.90% | 1.90% | – |
| | Total number of SRCs | 78825 | 98649 | ↗↗ |

# Data ETL – impute missings

| Period¶ | Raw data's obs number¶ | Obs # with missing ZIP¶ | Obs # missing all location info¶ | Service requests' modalities #¶ | Unique¶ ZIP¶ |
|---|---|---|---|---|---|
| April 2014¶ | 81645¶ | 3206¶ | 2740¶ | 170¶ | 278¶ |

# Data ETL – clean

# Data ETL – clean

NYC ZIP codes neighboring with "00083"

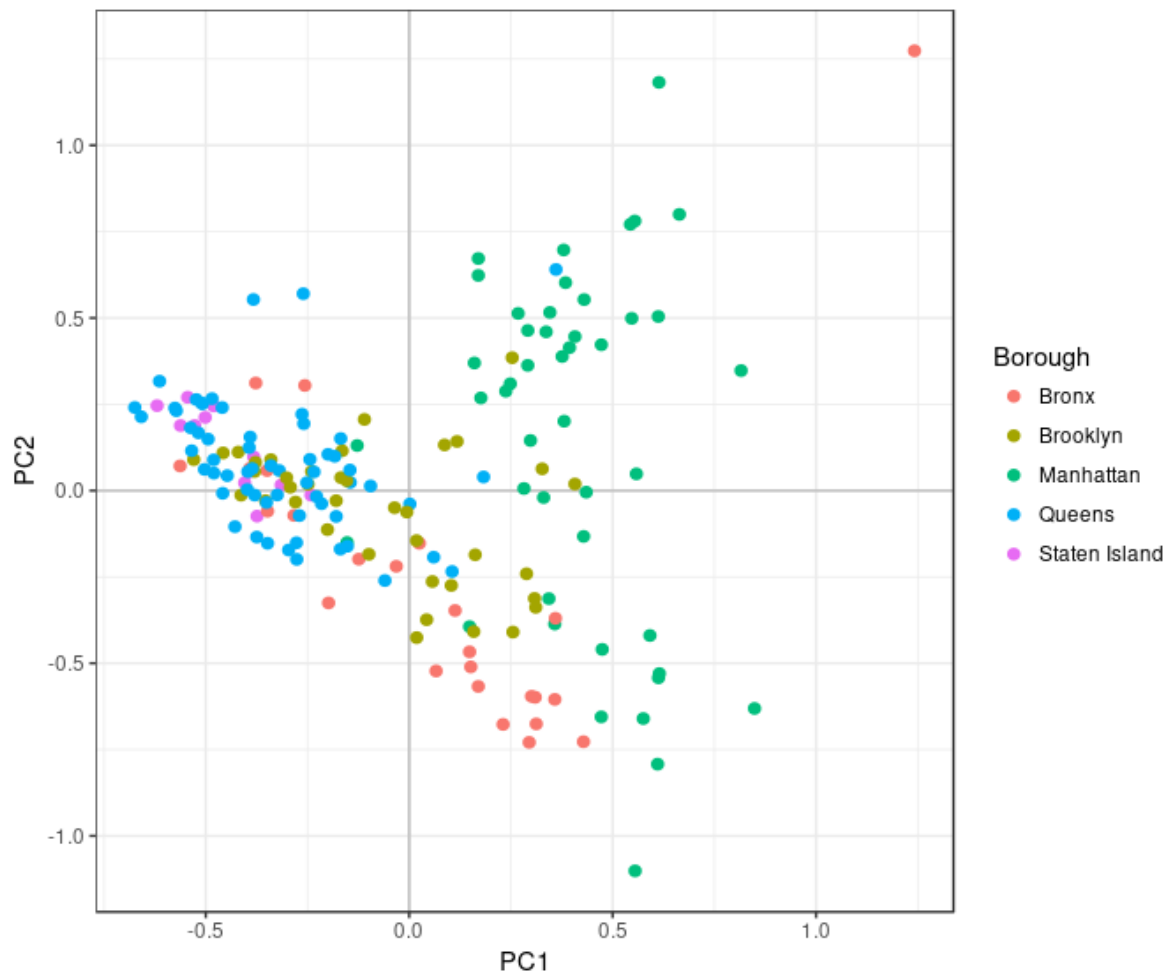# Data ETL – impute / clean

# MVA –  ZIP-SRC contingency table            *(April 2014)*

- Build frequency table
   (row profiles for ZIPs, column profiles for SRCs)


- Observe how 26 row marginals  <  5 / nbr_calls
   → Can we suppress them ?   *(… $\chi^2$-test of independence)*


- Run CA with row marginals as row profile weights
   ( $\chi^2$-*metric)*

# MVA – CA

CA factor map

CA factor map

Row profiles' projection in PC1-2 factorial plane

**CA factor map**

*column profiles (labeled only for cos2>0.4.*

Dim 1 (28.71%)

# MVA – PCA



|  | Dim 1 | | Dim 2 | | Dim 3 | |
|---|---|---|---|---|---|---|
|  | *ctr* | *cos2* | *ctr* | *cos2* | *ctr* | *cos2* |
| **HousCond** | 23.7 | 0.30 | 25.7 | 0.29 | 45.9 | 0.41 |
| **NoiseResid** | 10.4 | 0.20 | 43.2 | 0.69 | - | - |
| **NoiseConst** | - | - | 11.0 | 0.22 | 28.4 | 0.46 |
| **Traffic** | 12.8 | 0.50 | - | - | - | - |
| **ConsumProt** | - | - | - | - | 10.7 | 0.37 |
| **EnvProt** | 11.4 | 0.51 | - | - | - | - |

Eigenvalue mean: 4.2391e-02

# MVA – MCA



**Violations (4699)**
**(April 2014)**

H 42
MH 43
M 48
L 48

**Misdemeanors (21734)**
**(April 2014)**

OC 46
VH 44
H 41
M 50

**Felonies (11156)**
**(April 2014)**

VH 45
H 44
M 45
ML 47

Normalized segmentation (%)

Bronx
Brooklyn
Manhattan
Queens
Staten Island

# MCA factor map

Individuals - MCA

# Hierarchical clustering

- In this section we present an attempt to clusterize our data set.

- This attempt is carried out by applying, in the following order:
  - Probabilistic clustering with k-means replications
  - Hierarchical clustering
  - Clustering consolidation using k-means.

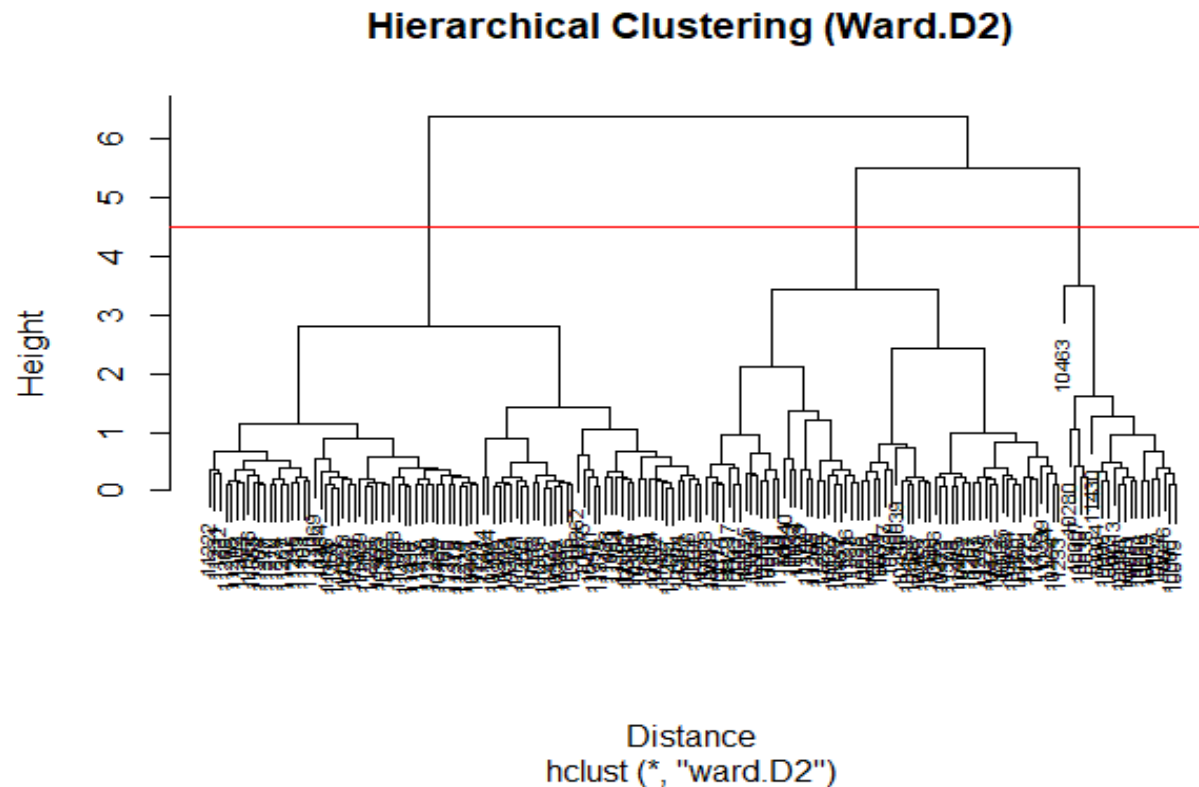# Hierarchical clustering

- Selection of the optimal number of clusters:
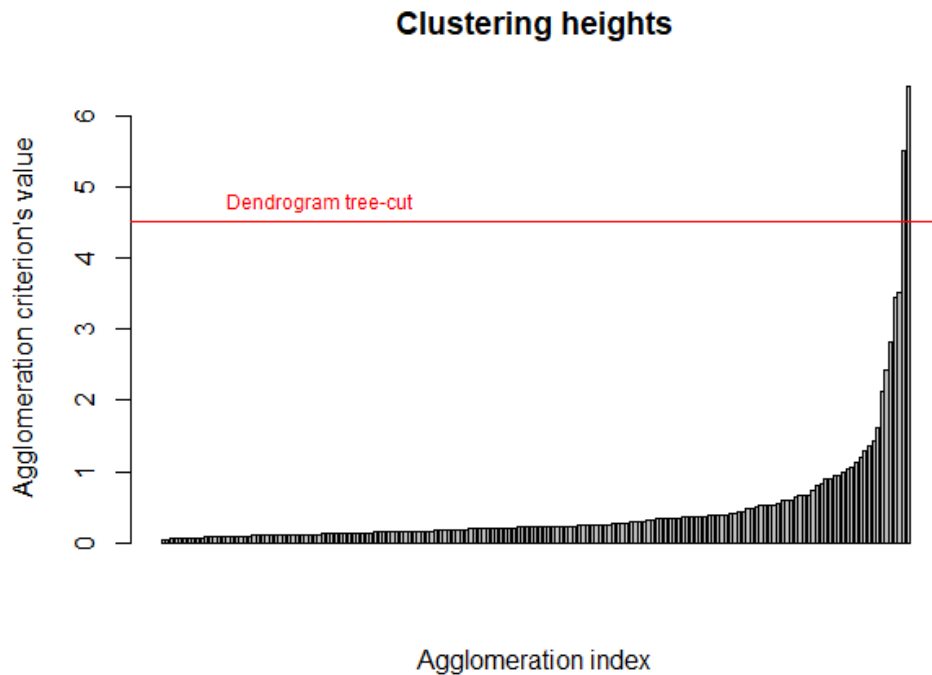


**Figure 13**. Selection of optimal number of clusters

# Hierarchical clustering

- Selection of the optimal number of clusters:

# Hierarchical clustering

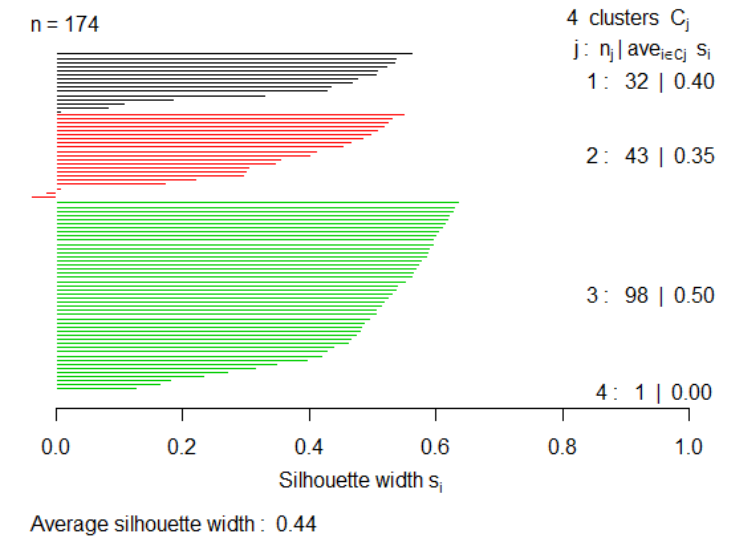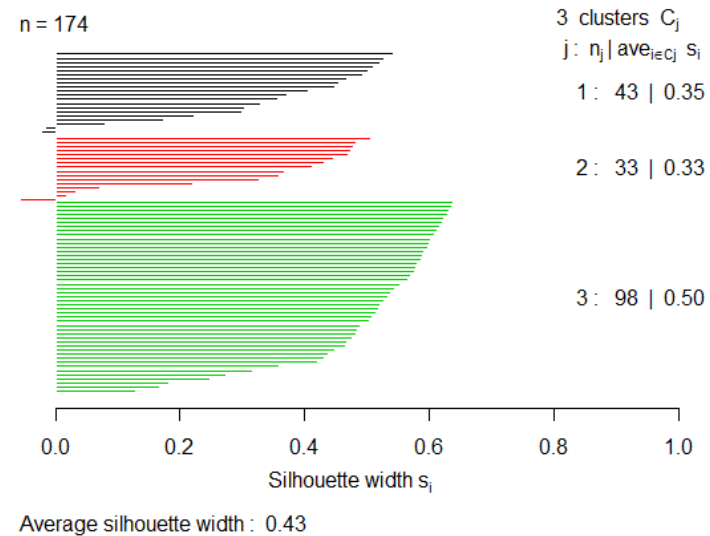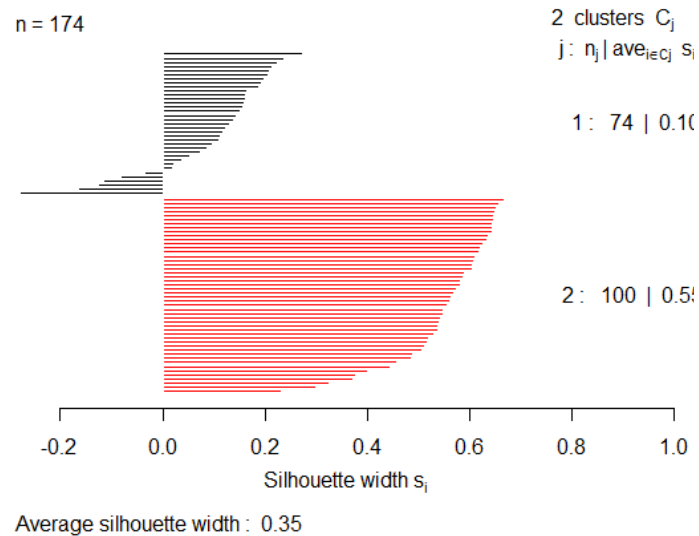- Selection of the optimal number of clusters:



**Clustering heights**

|  | G1 | G2 | G3 |
|------|------|------|-------|
| PC1 | 0,44 | 0,23 | -0,36 |
| PC2 | 0,58 | -0,28 | 0,08 |
| PC3 | 0,35 | 0,00 | -0,06 |
| PC4 | -0,18 | -0,02 | 0,01 |
| PC5 | -0,08 | 0,02 | 0,03 |

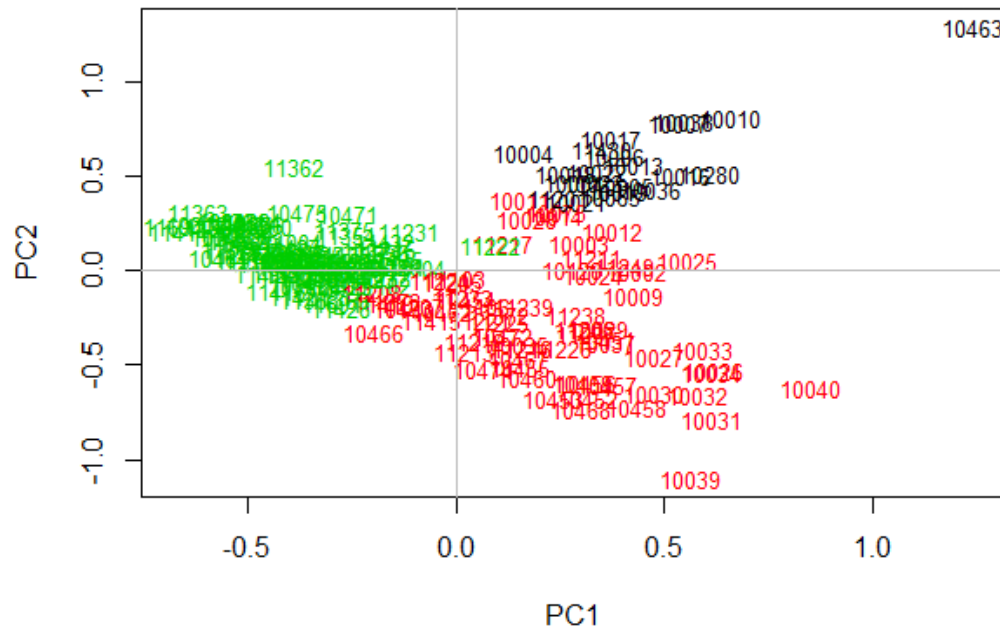**Table 7**. Centroids of the clusters

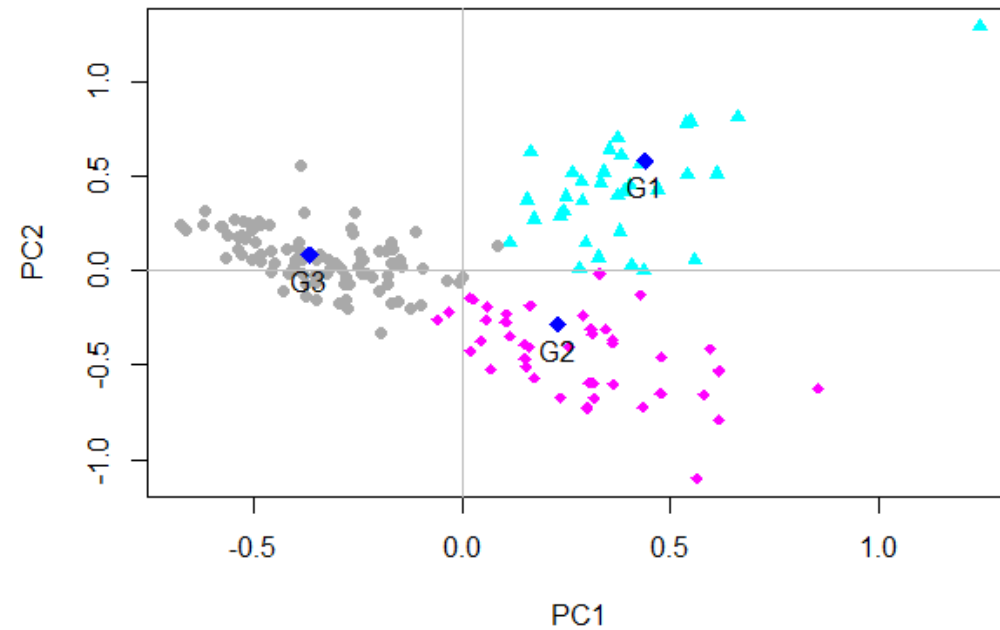# Hierarchical clustering

- Silhouette method:

# Hierarchical clustering
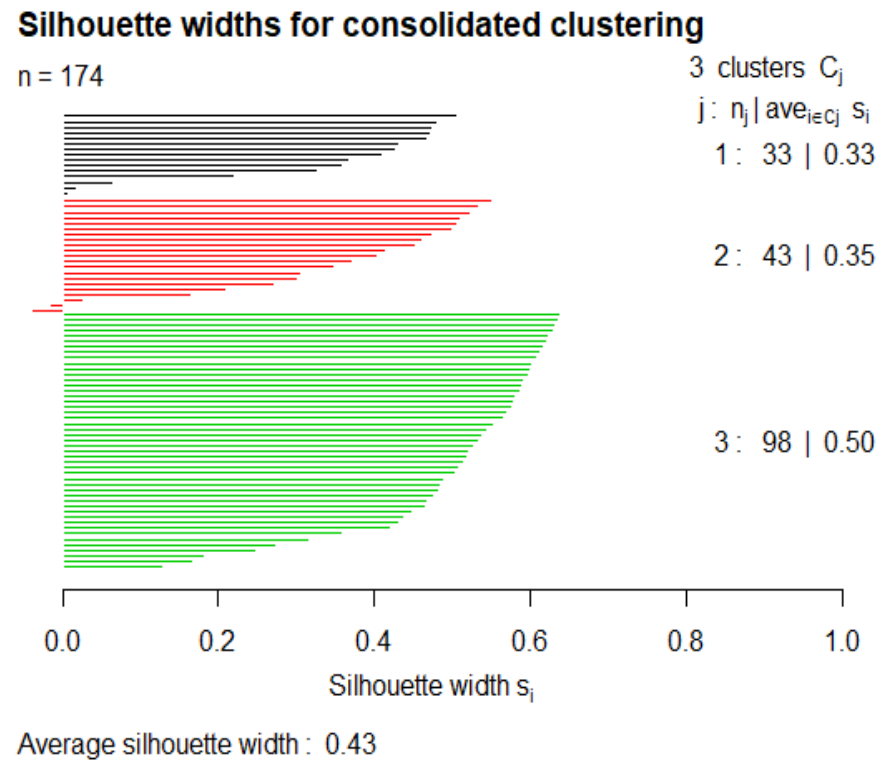
- Visualize partitions:

# Hierarchical clustering

- Silhouette method after consolidation:



**Silhouette widths for consolidated clustering**

n = 174

3 clusters $C_j$

$j : n_j | ave_{i \in C_j} s_i$

1 : 33 | 0.33

2 : 43 | 0.35

3 : 98 | 0.50

Silhouette width $s_i$

Average silhouette width : 0.43

# Hierarchical clustering

- Categorical description to interpret the clusters:
  - We reject the null hypothesis at the risk 0.05 of being wrong when the p-values <0.05.
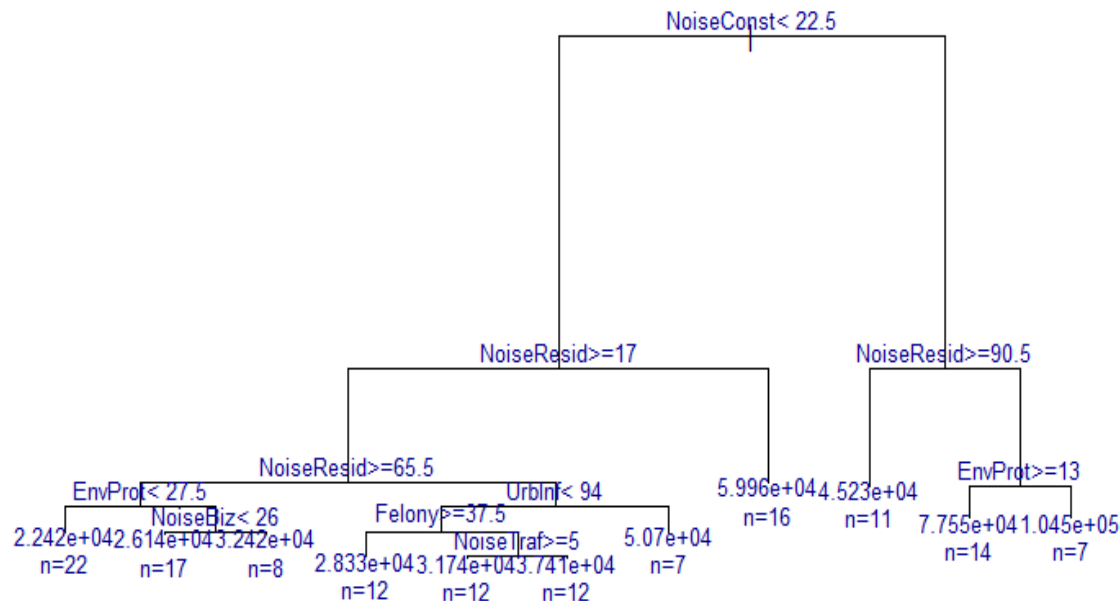  - Variables for which we reject H0 -> meaningful categorization:

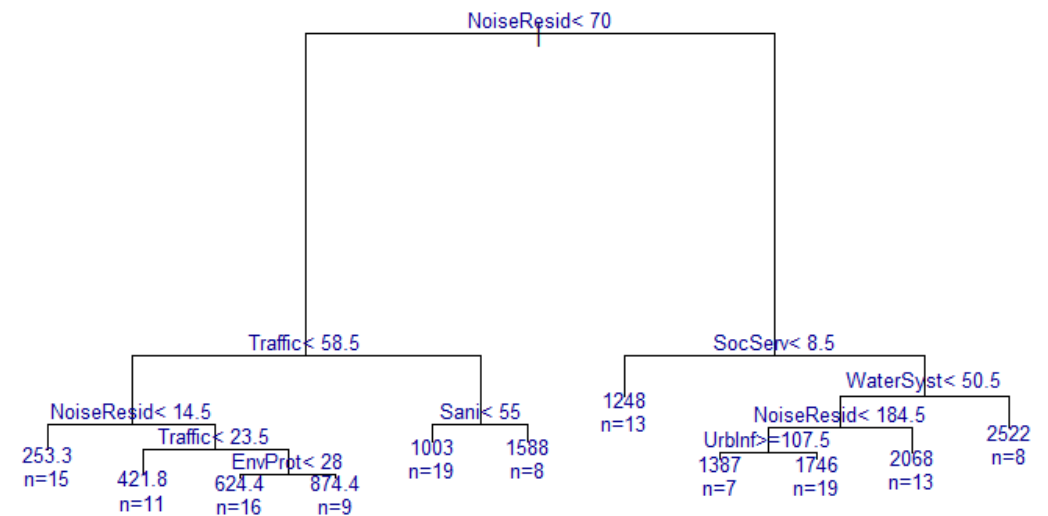| Cat1 | v.test | Mean in category | Overall mean | sd in category | Overall sd | p.value |
|---|---|---|---|---|---|---|
| NoiseConst | 9,145613 | 59,242424 | 17,132184 | 44,958587 | 29,298499 | 5,93E-20 |
| ConsumProt | 8,694289 | 43,303030 | 18,385057 | 25,590582 | 18,236820 | 3,49E-18 |
| UrbInf | 4,484378 | 98,575758 | 67,804598 | 55,909575 | 43,662851 | 7,31268E-06 |
| NoiseBiz | 3,328495 | 33,787879 | 19,885057 | 29,805183 | 26,578197 | 0,000873166 |
| HousCond | 3,277523 | 39,090909 | 13,856322 | 108,106391 | 48,991520 | 0,001047222 |
| NoiseTraf | 2,004039 | 23,393939 | 17,465517 | 17,769465 | 18,823636 | 0,045065931 |
| Sani | -1,960935 | 32,939394 | 41,557471 | 25,178080 | 27,965180 | 0,049886588 |
| NoiseResid | -2,335256 | 58,606061 | 85,896552 | 48,646492 | 74,361347 | 0,019530032 |
| WaterSyst | -2,425274 | 26,696970 | 35,454023 | 16,983923 | 22,975641 | 0,015296833 |
| EnvProt | -3,539835 | 18,424242 | 34,379310 | 13,407439 | 28,680460 | 0,000400377 |
| Traffic | -4,138211 | 33,757576 | 64,637931 | 21,908944 | 47,483221 | 3,50025E-05 |

# Decision trees

- In first place, we build the 2 possible decision trees:
    - One related to each of the 2 decision variables that we have ("medianInc" and "jlBenef").
- Before building the trees we split the dataset in trainng (80% of individuals) and test (20% of individuals).

# Decision trees

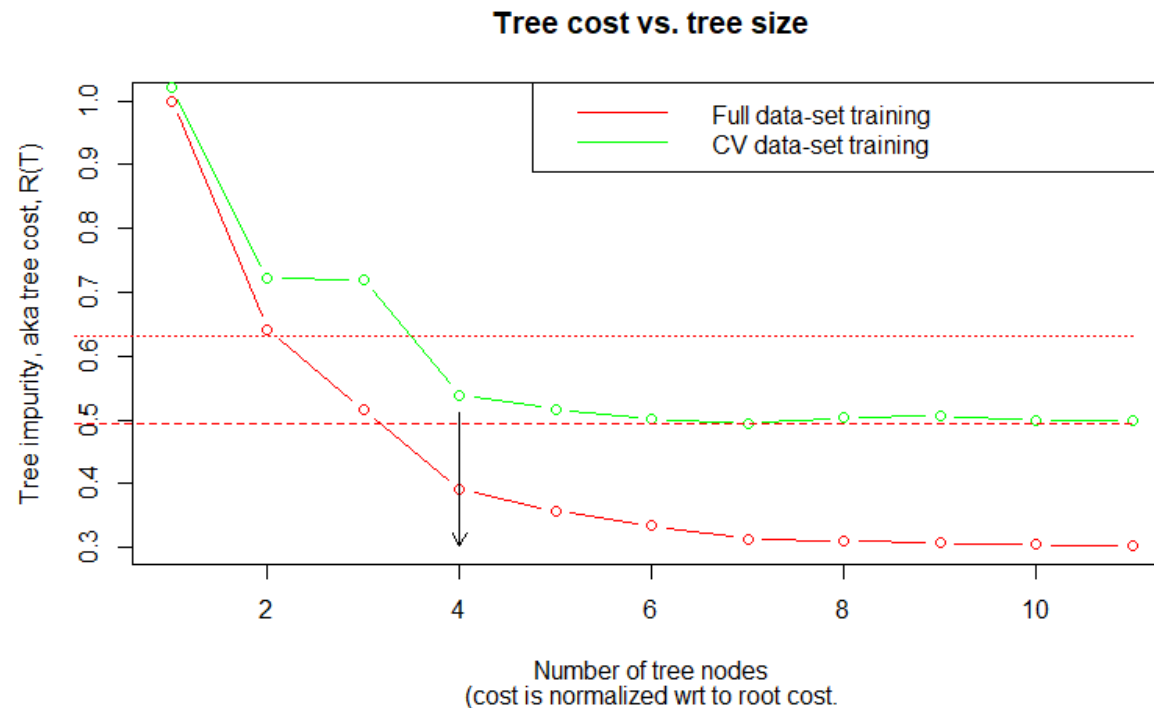Fully grown decision tree for training data-set and "medianInc" as decision variable.

Fully grown decision tree for training data-set and "jlBenef" as decision variable.
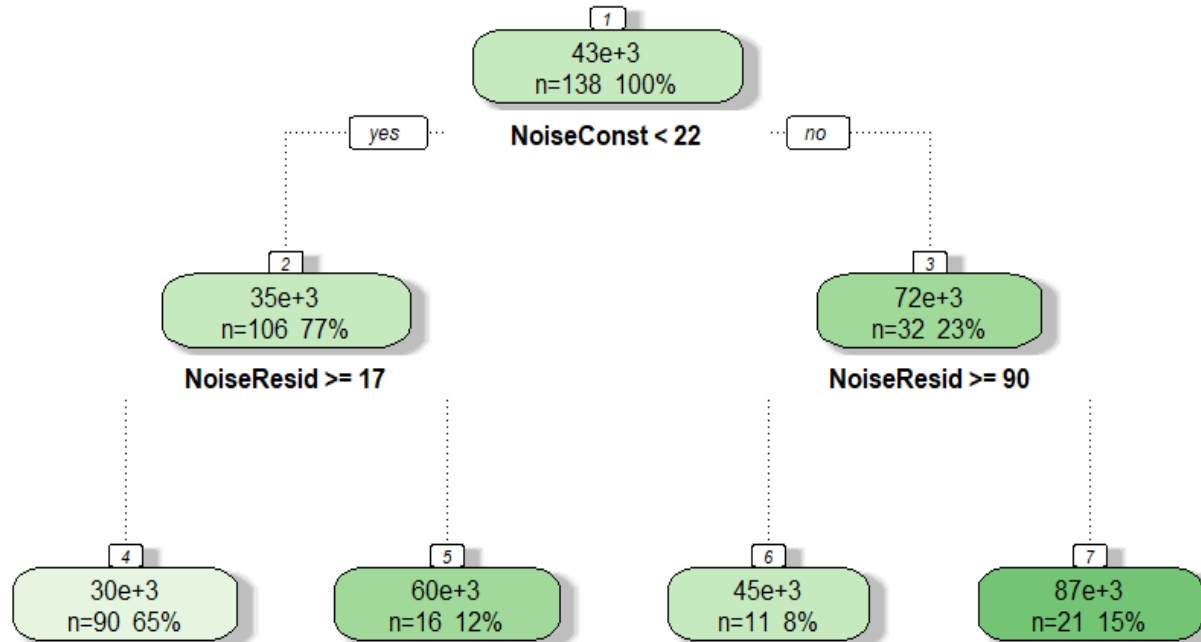
# Decision trees

- CV normalized error mean and the whole data set based training error as a function of tree size:



Tree cost vs. tree size

- Red horizontal dashed line (below) -> minimum tree impurity (MTI) level

- Red dotted line (above) -> MTI + 1.

- Black arrow -> optimum number of nodes for post-pruning.

# Decision trees

- Post-pruning:
  - Optimum complexity parameter -> αopt = 0.03489.
  - Now, we are able to post-prune the decision tree by using -> αopt.

# Decision trees

- Post-pruning:
  - Split rules:

Rule number: 4 [medianInc=30240.9222222222 cover=90 (65%)]

  NoiseConst< 22.5
  NoiseResid>=17


Rule number: 7 [medianInc=86539.2380952381 cover=21 (15%)]
  NoiseConst>=22.5
  NoiseResid< 90.5


Rule number: 5 [medianInc=59963.4375 cover=16 (12%)]
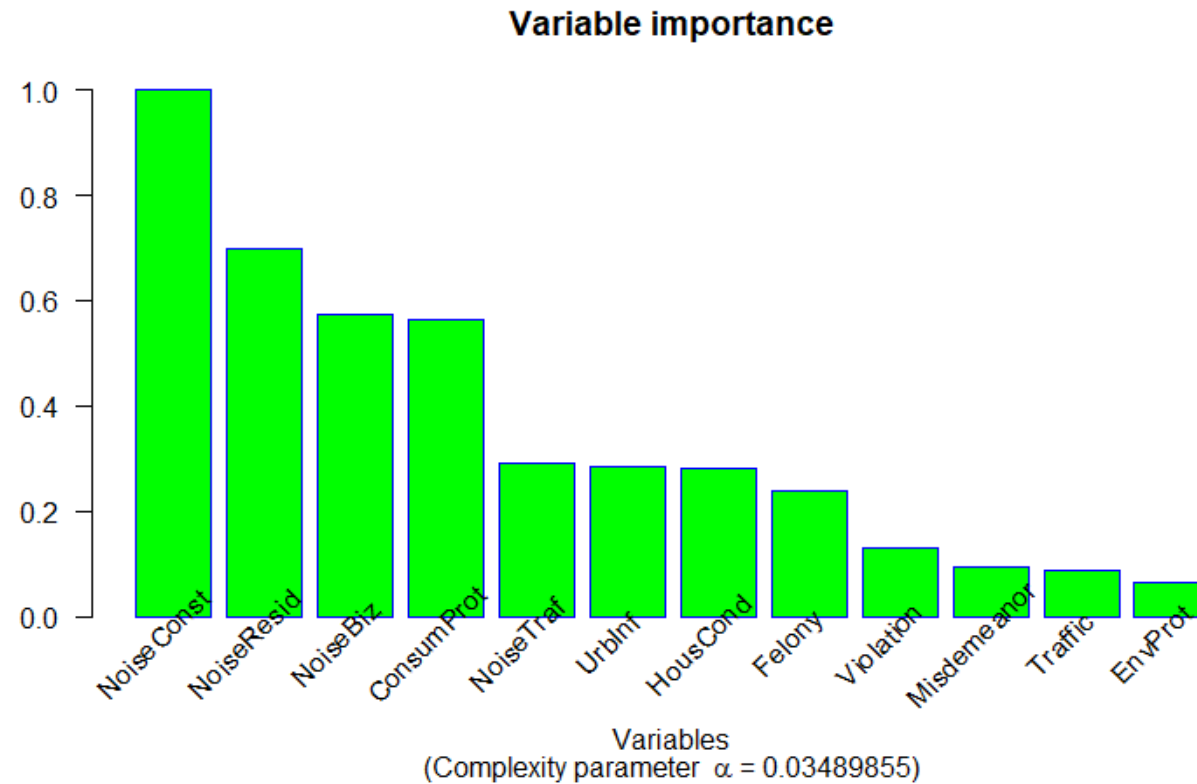  NoiseConst< 22.5
  NoiseResid< 17


Rule number: 6 [medianInc=45225.8181818182 cover=11 (8%)]
  NoiseConst>=22.5
  NoiseResid>=90.5

# Decision trees

- Variable Importance for the optimally pruned decision tree :



Variable importance

# Decision trees

- Predictions:

- Slice of the original results-table which contains the predictions for each value of "medianInc" in the test data set.

| "medianInc" for Test-set | 4 predicted classes | | | |
|---|---|---|---|---|
| | 30240,9222 | 45225,8182 | 59963,4375 | 86539,2381 |
| 17992 | 1 | 0 | 0 | 0 |
| 18164 | 1 | 0 | 0 | 0 |
| 26143 | 1 | 0 | 0 | 0 |
| 26170 | 1 | 0 | 0 | 0 |
| 27102 | 1 | 0 | 0 | 0 |
| 27144 | 1 | 0 | 0 | 0 |
| 27203 | 1 | 0 | 0 | 0 |
| 27303 | 1 | 0 | 0 | 0 |
| 27331 | 1 | 0 | 0 | 0 |
| 27374 | 1 | 0 | 0 | 0 |
| 27898 | 1 | 0 | 0 | 0 |
| 90981 | 1 | 0 | 0 | 0 |
| 92955 | 0 | 0 | 0 | 1 |
| 93056 | 0 | 0 | 1 | 0 |
| 95992 | 0 | 0 | 0 | 1 |
| 97669 | 0 | 0 | 0 | 1 |
| 98024 | 0 | 0 | 0 | 1 |
| 110248 | 0 | 0 | 0 | 1 |
| 128571 | 0 | 0 | 1 | 0 |
| 185593 | 0 | 0 | 0 | 1 |
| ^^^^^^^^^^ | ^^^^^^^^^^ | ^^^^^^^^^^ | ^^^^^^^^^^ | ^^^^^^^^^^ |
| **Total Freq** | **97** | **10** | **15** | **18** |

# Questions?