

Profiling and analysis of NYC Open Data

...

Sai Likhith
Kartikeya Negi
Reuben Varghese

Overview

Data has gained prominence in the last few years

Big data problems are the norm

Understanding and solving them is the challenge for this generation of engineers

Key Objectives

Simulate a set of real-world big data problems

Solve these problems meaningfully, using techniques learnt in class

NYC Open data

Generic profiling

Semantic profiling

Data analysis

Understanding the problem

Generic profiling

Data comes with very less description (metadata)

Noise

Unclassified data types

Sparse/wrong/missing data

Semantic profiling

What real world entities (semantics) does our data contain?

Data analysis

What can our data tell us about the real world that it models?

Generic Profiling

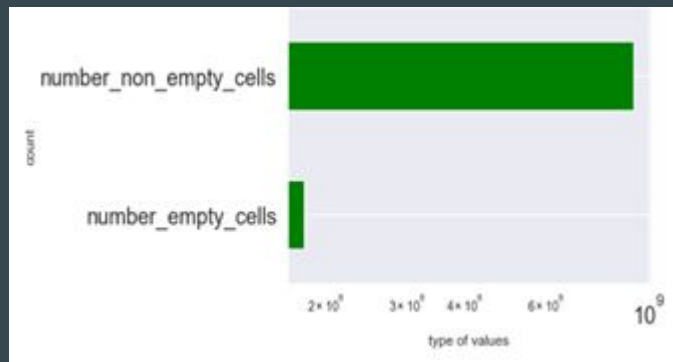
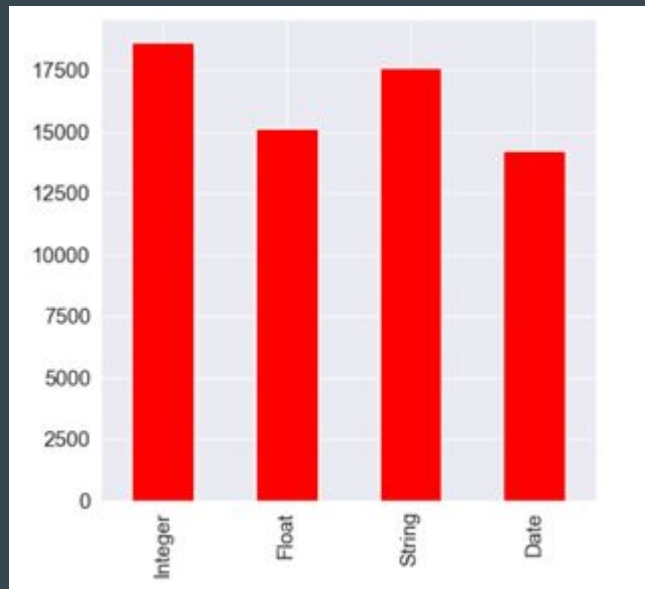
Structure profiling

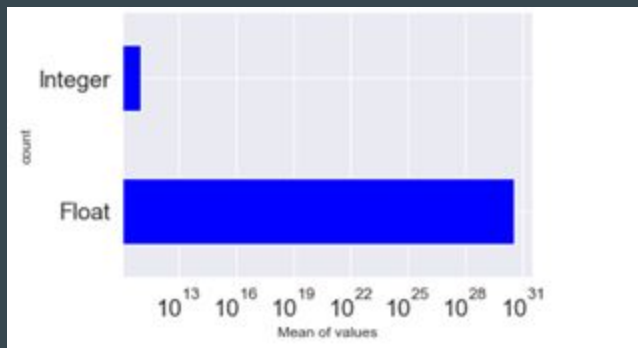
- data types (text, integer(long), time/date, float)

Content profiling

- full/nan values,
- missing values,
- unique values
- frequent values
- max/min/mean/std. dev

Analysis





??

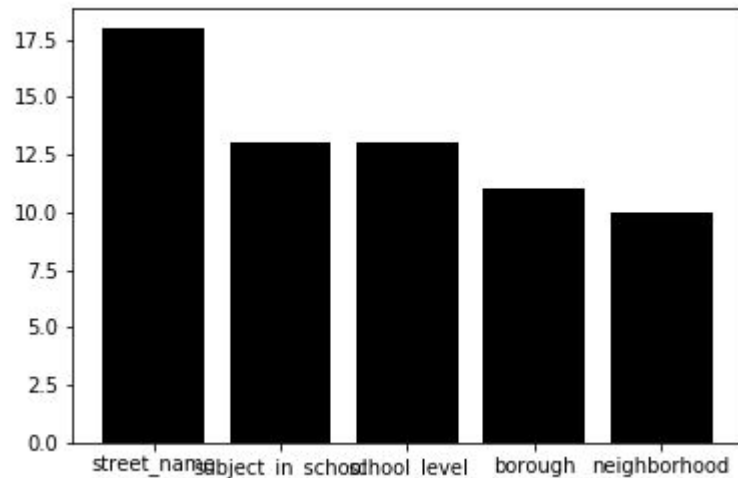
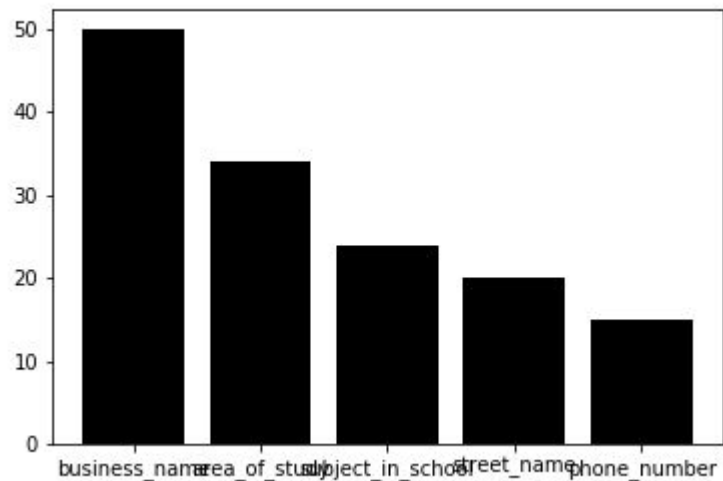
Semantic Profiling

Column semantics using column names: ground truth

Functions defined/ontologies created for each semantic type:

- Regular expressions (phone number, location, etc.)
- NLP model: name entity recognition library (person name, business name, etc.)
- Fuzzy string matching (parks, building classification, vehicle type)
 - Online ontologies

Analysis



Data analysis

Gaining real world information from our model of the world (data)

NYC 311 entries (November 2019)

Most frequent complaints by Borough

Why 1 month?

- Most recent/relevant info
- Analyse the effect of thanksgiving on most frequent complaints

