

# Predicting Real Estate Sales Using Machine Learning and Spatial Dependence

## Contents

<b>Introduction</b>	<b>1</b>
Motivating Example: Combatting Economic Exclusion . . . . .	1
Geo-spatial Data . . . . .	2
<b>Literature Review</b>	<b>2</b>
Lit Review . . . . .	2
sample citations . . . . .	2
<b>Methodology</b>	<b>2</b>
Methodology Section . . . . .	2
<b>Results</b>	<b>2</b>
<b>Conclusions and Future Research</b>	<b>3</b>
<b>References</b>	<b>4</b>

## Introduction

### Motivating Example: Combatting Economic Exclusion

Predictive modeling using spatial dependence has been employed extensively in recent years, notably in Crime Prediction (Almanie 2015). However, a key deficiency of many spatial models are their use of arbitrarily defined geographic regions, such as zip codes, political districts, police precincts, state lines, neighborhoods, etc. which diminish and obscure potentially valuable insights. Worse yet, many predictive models ignore and exclude spatial dependence, violating one of the basic tenets of geography: the direct relationship between distance and likeness (Miller 2007).

Income inequality may be a central challenge of our time. Researchers at the Urban Institute (Solomon Greene and Lei 2016) recently identified the socio-economic phenomenon of “Economic Exclusion” as one compelling explanation for the recent rise in inequality in the US. Vulnerable populations (disproportionately communities of color, immigrants, refugees, and women), who are displaced by localized economic prosperity enter into a gradual cycle of diminished access to good jobs, good schools, health care facilities, public spaces, etc. This systematic denial causes enduring and self-reinforcing poverty over the course years and even generations, gradually entrenching income inequality and general unrest.

One way to practically combat economic exclusion is to focus on preventing displacement, however, detecting gentrification at an early enough stage can be a daunting task. When an area experiences economic growth, increased housing demands and subsequent affordability pressures can lead to evictions of low-income families and small businesses. Government agencies and nonprofits tend to intervene once displacement is already underway, and after-the-fact interventions can be costly and ineffective. There are a host of pre-emptive actions that can be deployed to stem divestment and ensure that existing residents benefit from new investments. Not unlike medical treatment, early detection is key to success. Consequently, in 2016, the Urban Institute put forth a call for research into the creation of “neighborhood-level early warning and response systems that can help city leaders and community advocates get ahead of neighborhood changes” (2016).

This paper explores novel techniques to predict gentrification in the pursuit of combating displacement and economic exclusion. Modern techniques of data mining, machine learning and predictive modeling are applied to datasets describing property values and sale prices in New York City. We explore the viability of

using spatial lags, i.e., variables created from physically proximate observations, as features in a machine learning predictive model.

## Geo-spatial Data

The world has seen an unprecedented amount of geospatial data produced in recent years (i.e., data that contain information about where an observation exists or happened). Every day in the U.S., federal, state and local government agencies are making their troves of geo-spatially tagged data available for the benefit of the public. Adequate tools to describe, explore and model such data are in short supply for the data-journalists and data-activists who have become modern mechanisms of public service. It is imperative that research be done and tools created to better harness such data for commercial and public good.

## Literature Review

### Lit Review

Much of the research on predicting real estate values has been in service of creating mass appraisal models. Mass appraisal models share many characteristics with predictive machine learning model modeling, primarily in that they are data-driven, standardized methods that employ statistical testing (Eckert 1990).

### sample citations

Sample Citation: (Antipov and Pokryshevskaya 2012) (see: Antipov and Pokryshevskaya 2012, 33–35; also Antipov and Pokryshevskaya 2012, ch. 1 and *passim*)

A minus sign (-) before the @ will suppress mention of the author in the citation. This can be useful when the author is already mentioned in the text:

Antipov says blah (2012).

You can also write an in-text citation, as follows:

Antipov and Pokryshevskaya (2012) says blah.

## Methodology

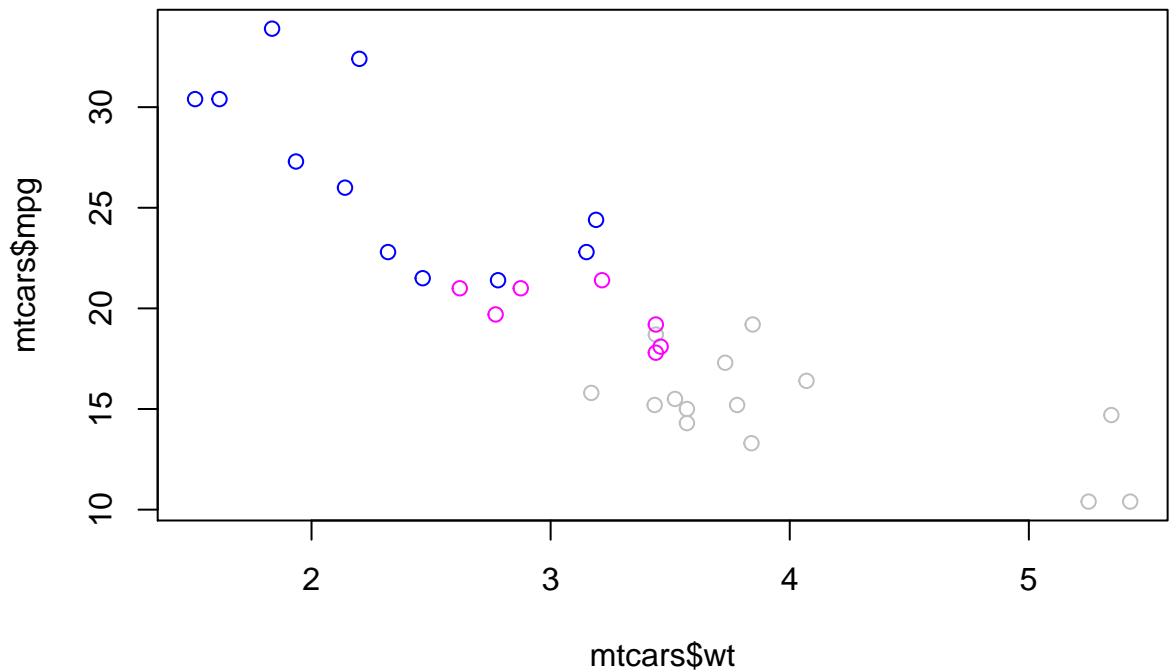
### Methodology Section

Random Forrest has several advantages over traditional geographic weighted regression, among them:

1. Ability to handle large amounts of categorical data without much pre-processing
2. Ability to model in spite of missing values in data
3. Eliminated colinearity as a concern
4. Allows for the introduction of many more variables without requiring penalty for additional predictors
5. Works relatively fast and can be parallelized

## Results

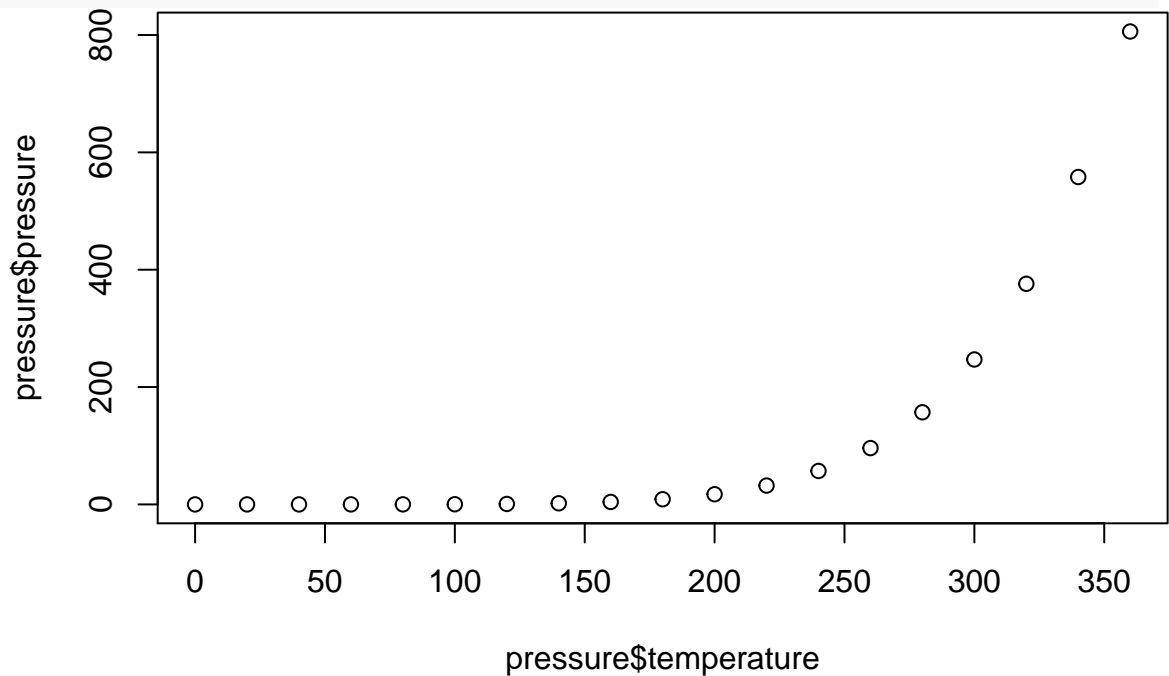
```
plot(mtcars$wt, mtcars$mpg, col = mtcars$cyl)
```



Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.

## Conclusions and Future Research

```
plot(pressure$temperature, pressure$pressure)
```



Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.

## References

- Almanie, R.; Lor, T.; Mirza. 2015. “Crime Prediction Based on Crime Types and Using Spatial and Temporal Criminal Hotspots.” *International Journal of Data Mining & Knowledge Management Process (IJDKP)* 5 (4).
- Antipov, Evgeny A., and Elena B. Pokryshevskaya. 2012. “Mass Appraisal of Residential Apartments: An Application of Random Forest for Valuation and a Cart-Based Approach for Model Diagnostics.” *Expert Systems with Applications*.
- Eckert, J. K. 1990. *Property Appraisal and Assessment Administration*. Chicago, IL.: International Association of Assessing Officers.
- Miller, J.; Aspinall, J.; Franklin. 2007. “Incorporating Spatial Dependence in Predictive Vegetation Models.” *Ecological Modelling* 202 (3): 225–42.
- Solomon Greene, Molly Scott, Rolf Pendall, and Serena Lei. 2016. “Open Cities: From Economic Exclusion to Urban Inclusion.” *Urban Institute Brief*, June. Urban Institute Brief.