

A PROTOTYPE OF A MACHINE LEARNING WORKFLOW TO CLASSIFY LAND USE FROM  
THE HOUSING MARKET DYNAMICS. PART OF A LONGITUDINAL ANALYSIS OF  
HOUSING SALES IN THE GREATER TORONTO-HAMILTON AREA.

by

Stepan Oskin

A thesis submitted in conformity with the requirements  
for the degree of Master of Applied Science  
Graduate Department of Civil Engineering  
University of Toronto

© Copyright 2019 by Stepan Oskin

# Abstract

A prototype of a machine learning workflow to classify land use from the housing market dynamics.  
Part of a Longitudinal Analysis of housing sales in the Greater Toronto-Hamilton Area.

Stepan Oskin  
Master of Applied Science  
Graduate Department of Civil Engineering  
University of Toronto  
2019

There is ample evidence of the role of land use and transportation interactions in determining urban spatial structure. The new data sources introduced by increased digitization of human activity, such as Teranet's dataset of real estate sales, offer opportunities for development of integrated urban models or studies conducting longitudinal analysis of changes in land value distributions. To facilitate this, data from various sources (Census, TTS, etc.) needs to be merged at the land parcel level to enhance datasets with additional attributes, while maintaining the ease of data storage and retrieval. In addition to that, accurate land use information needs to be added to Teranet records to allow separating sales data by major property types. This master's thesis proposes a prototype of a workflow to augment Teranet's dataset with data from multiple sources and use machine learning to classify land use at the transaction level based on the housing market dynamics.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Background information</b>	<b>3</b>
2.1	Complexity of urban systems and "wicked" problems . . . . .	3
2.2	Transportation-land use cycle . . . . .	3
2.3	Evolution of LUT models . . . . .	5
2.4	Longitudinal Housing Market Research conducted by UTTRI . . . . .	6
2.5	POLARIS: electronic system of land registration in Ontario . . . . .	7
2.6	Teranet's dataset, its challenges and the proposed solution . . . . .	7
2.7	Chapter summary . . . . .	9
<b>3</b>	<b>Spatial and temporal relationships between data sources</b>	<b>10</b>
3.1	Description of the data sources used . . . . .	10
3.2	Spatial relationships between data sources . . . . .	11
3.3	Temporal relationships between data sources . . . . .	12
3.4	Chapter summary . . . . .	14
<b>4</b>	<b>Data preparation</b>	<b>15</b>
4.1	Tidy data and database normalization . . . . .	16
4.2	Introduction of new keys and attributes via spatial and temporal relationships . . . . .	17
4.3	Outliers . . . . .	19
4.4	Engineering new features for the classification algorithm . . . . .	21
4.5	Chapter summary . . . . .	23
<b>5</b>	<b>A prototype of a machine learning workflow to classify land use</b>	<b>24</b>
5.1	Selecting and encoding the target variable . . . . .	25
5.2	Missing values . . . . .	27
5.3	Dimensionality reduction . . . . .	27
5.4	Tuning model hyperparameters . . . . .	31
5.5	Feature scaling . . . . .	33
5.6	Model selection . . . . .	34
5.7	Chapter summary . . . . .	37

<b>6</b>	<b>Model evaluation</b>	<b>38</b>
6.1	Metrics for evaluating model performance . . . . .	38
6.2	Model performance . . . . .	38
6.3	Best performing model: Random Forest . . . . .	40
<b>7</b>	<b>Conculsion</b>	<b>44</b>
	<b>Bibliography</b>	<b>45</b>

# Chapter 1

## Introduction

The fundamental link between transportation and urban form creates a feedback relationship between land development, travel needs, viability of alternative modes, accessibility, and other important characteristics of the urban transportation system. Numerous "top-down" and "bottom-up" models have been designed to analyze and forecast the behaviour of urban regions and interaction of their transportation and land use systems. Since urban systems are complex in nature and require "re-solving" over and over, data science process models present a good fit for this task with their iterative structure.

Increased digitization of human activity, such as introduction of POLARIS land registration system by the Government of Ontario in 1985, produces a wealth of new information that can be used to study interaction between land use and transportation at a fine spatial and temporal scale. Teranet's dataset of real estate transactions presents a wealth of information on the housing market of Ontario and can be used for empirical studies of transportation-land use interaction. However, along with the opportunities, the new data sources also present new challenges. Teranet's dataset has some data quality issues that need to be addressed and might require special skills to work with due to its size. But most importantly, it is very limited in the number of features available for each transaction.

One of the major challenges of working with Teranet's data is the lack of features describing each transaction, namely the type of property being sold. As Teranet's records have timestamps (dates) and coordinates of parcel centroids for each record, they can be joined with other data sources via temporal and spatial relationships, integrity of which needs to be maintained when combining data sources to ensure semantic interoperability. These relationships are implemented through a standardized data preparation workflow using Python via a series of jupyter notebook which was designed as a part of this master's thesis. In addition, since all the available sources of land use information have their limitations, a prototype of a machine learning workflow to classify land use from the housing market dynamics is proposed and tested by this master's thesis.

This master's thesis generally follows the structure proposed by CRISP-DM, a comprehensive data mining methodology and a process model[43] (shown on figure 1.1). CRISP-DM breaks down the life cycle of a data mining project into six phases: business understanding, data understanding, data preparation, modeling, evaluation, and deployment.

The first phase of CRISP-DM is business understanding, which involves such key elements as the definition of goals of the project and their acceptance criteria and the definition of the target variable of analysis[37]. Chapter 2 focuses on the definition of the goals of this master's thesis: it provides

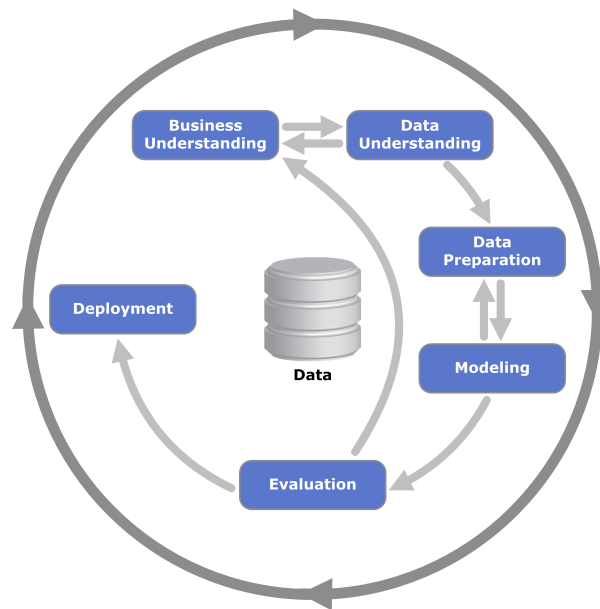


Figure 1.1: Phases of the CRISP-DM Reference Model, adapted from Shearer[43].

the background information on land use and transportation models and some of the ongoing research efforts at UTTRI involving the housing market, outlines the role of Teranet’s dataset in these efforts, the challenges of its usage and the proposed solution. The definition of the target variable and the acceptance criteria for model evaluation are discussed in more detail in chapter 5. Chapter 3 covers the second phase of the process model, data understanding, and discusses the nature of the spatial and temporal relationships of different data sources used in this master’s thesis. Chapter 4 covers the third phase, data preparation, and presents the data preparation workflow designed to implement the relationships introduced in chapter 3. Chapter 5 covers the fourth phase, modeling, and describes a prototype of a machine learning workflow to classify land use from the housing market dynamics. Chapter 6 covers the fifth phase, evaluation, and presents the discussion of the results and chapter 7 presents the conclusion and outlines opportunities for future work.

## Chapter 2

# Background information

This chapter discusses the complex interaction of land use and transportation, provides a brief overview of the history of development of land use and transportation (LUT) models, presents some legal and historical background for Teranet’s dataset of land registration records, and finishes with a discussion of challenges of working with Teranet’s data and the proposed solution.

### 2.1 Complexity of urban systems and ”wicked” problems

In her famous 1961 book, Jane Jacobs[19] described a city as ”a problem in organized complexity”; since then, many other researchers have remarked that urban systems exhibit complex behaviour[5, 8]. Complexity of a system can be defined as a state or quality of being intricate or complicated. For a system to be complex is not necessarily the same as to be complicated; complex systems can be simple, i.e. governed by a single equation. Complexity of a system has to do with the intrinsic ability of a system to surprise us with its behaviour; that the system is hard to understand, despite the mechanics of it being relatively simple.

In 1973, a little over a decade after Jacobs, Rittel and Webber[41] presented a path-breaking conceptualization; this conceptualization characterized urban planning problems as ”wicked” problems: problems which cannot be definitively described and for which it makes no sense to talk of ”optimal solutions”. In their paper, Rittel and Webber stated that such ”wicked” problems are never ”solved”, and that the focus instead becomes on iteratively ”re-solving” the problems over and over. More than 40 years after their original publication, Rittel and Webber’s ideas remain relevant to the policy sciences today: there is an intense interest in the nature of ”wicked” problems and the complex tasks of identifying their scope, viable responses, and appropriate mechanisms and pathways to improvement[13]. Interaction between land use and transportation, which is discussed in the following section, presents a prime example of urban complexities and ”wicked” problems.

### 2.2 Transportation-land use cycle

Among the reasons why transportation and land use interaction is ”wicked” are such aspects as pluralism of expectations among stakeholders, institutional complexity in policy making, and scientific uncertainty[38]. More importantly, there is a fundamental link between transportation and urban form:

urban form has an enormous impact on the type and cost of transportation systems needed to serve residents of a metropolitan area[21]. Transportation, in turn, influences land development and location choices of people and firms, and thus completes the formation of a feedback relationship that Stover and Koepke[46] referred to as a cycle. Interconnections between points (activities) in space can be perceived through the medium of the transportation system[35]

Figure 2.1 illustrates the complex interactions between land use and transportation system as summarized by Miller, Kriger and Hunt[35].

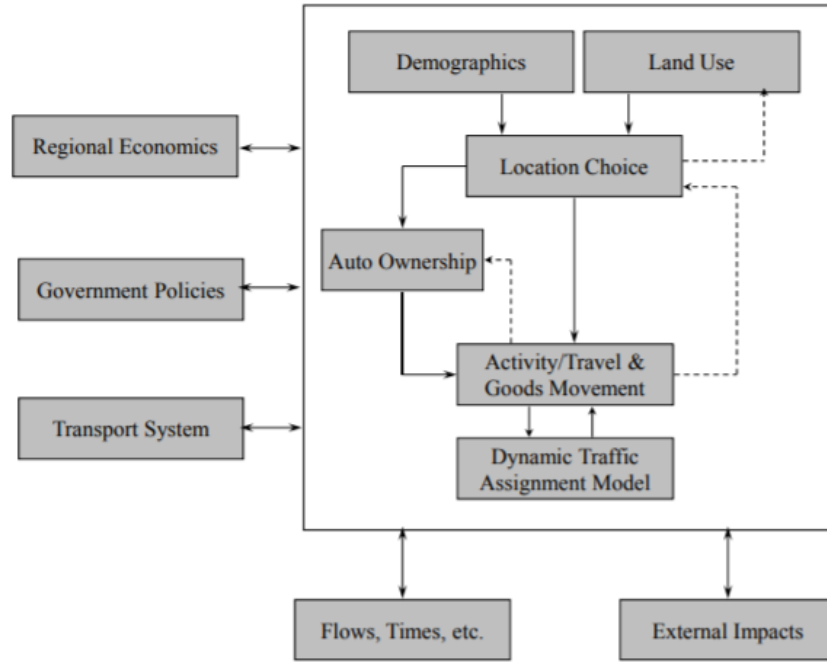


Figure 2.1: An Idealized Integrated Urban Model System, adapted from Miller, Kriger and Hunt[35].

Many different types of models are used in planning, such as demand forecasting models projecting traffic or ridership, or land use models projecting and distributing population and jobs within an area. At an earlier stage of model development, some analysts argued that there is no significant link between transportation and land use, given the near-ubiquity of the transportation (road) network[35]. However, the unprecedented urban growth of the 21st century introduced new challenges for urban systems such as extreme road congestion, equity of access to jobs and services among low-income households, energy scarcity, environmental and GHG impacts from transportation systems and public health impacts of land use patterns[31, 36].

It became apparent that these "transport problems" cannot be solved through transportation policies and investment alone, that the physical design of the city at the "macro" and "micro" scale critically interfaces with the demand for and performance of the transportation system. In addition, to accurately assess the costs and benefits of an expensive long-term transportation infrastructure investment, "feed-back" effects of these investments on urban form, land values, property taxes, quality of life, etc. need to be quantified and included in evaluation and decision making. Thus, today there is a steadily growing recognition within the urban policy field that the interaction between transportation and land use does exist and does matter[31].



In the context of models, integrated urban models (IUMs) aim to capture the complex relationship between urban systems such as transportation and land use more accurately. Integrated land use-transportation models combine travel demand forecasting and land use forecasting functions and recognize that the distribution of population and jobs depends, in part, on transportation accessibility. The reverse is also true, and thus integrated models incorporate feedback relationship between transportation and land use, with economic decisions by households and firms acting as one of the links between the two systems[35].

## 2.3 Evolution of LUT models

The history of treating cities as systems via simulation models of transportation and land use dates back to 1950s when General System Theory and Cybernetics came to be applied in the softer social sciences[5]. The first operational simulation model that truly integrated land use and transportation is considered to be A Model of Metropolis built in 1964 by Ira S. Lowry for the Pittsburgh region based on economic base theory[27]. It was a highly aggregate model based on theories of spatial interaction, such as the gravity model that was popular in quantitative geography and transportation planning at the time[9]. Models based on spatial interaction framework continued to be developed through mid-1980s, until developments in random utility theory allowed researchers to describe choices among discrete alternatives, such as the choice of travel mode, and generate models based on the study of disaggregate behaviour[18].

Figure 2.2 provides the general overview of chronological development of LUT models as summarized by Iacono[18].

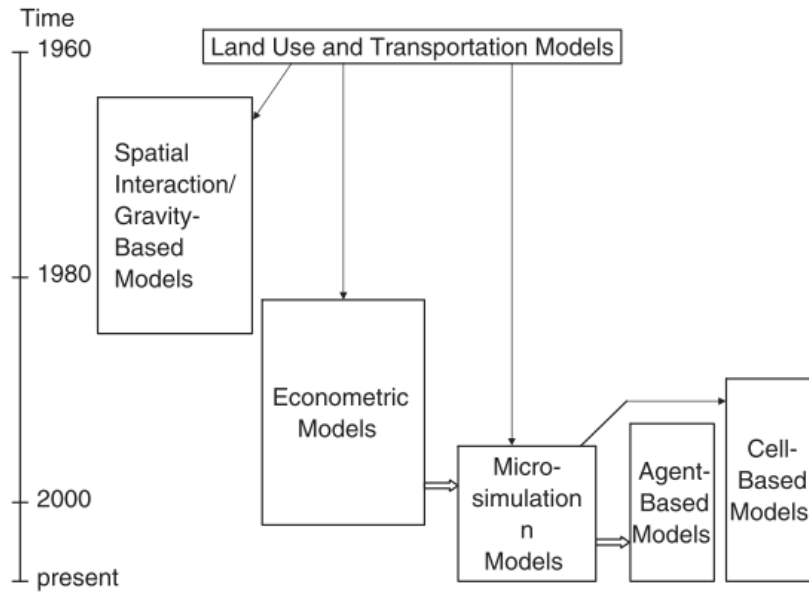


Figure 2.2: Chronological development of LUT models as summarized by Iacono[18].

The modeling paradigm has changed fundamentally in the early 1990s along with the advances in computing power and efficiency of data storage. Urban systems used to be viewed as hierarchical and centrally organized equilibrium structures, or "top-down". Instead, now they were considered to be

structured from the "bottom-up", dynamically retaining their integrity through interactions of numerous microelements[5]. A new broad class of LUT models that could fall under the title of "microsimulation" began to be developed: it included such classes of models as activity-based travel, cell-based models, multi-agent models, and more recently comprehensive urban microsimulation models that reflect the dynamics of changes in the population and the urban environment[18]. "Micro" in the microsimulation implies that the model must be highly disaggregated spatially, socio-economically and in its representation of processes. "Simulation" implies that the model must be numerical, stochastic, have an explicit time dimension, and "evolve" into the end state rather than "solve for it"[33]. An example of such a model has been developed by the University of Toronto ILUTE team; their product is an integrated urban model capable of microsimulating urban demographic evolution, housing markets and travel behaviour over extended periods of time[32].

The Integrated Land Use, Transportation, Environment (ILUTE) model system is a system of disaggregate agent-based microsimulation models for Greater Toronto–Hamilton Area (GTHA). It includes such components as land use, activity/travel, urban economics, auto ownership, demographics and emissions/energy use. The system uses disaggregate models of spatial socioeconomic processes to evolve from a known base case to a predicted end state of the GTHA in 1-year time steps[34]. In ILUTE, the system state of GTHA is defined in terms of the individual persons, households, dwelling units, firms, etc. that collectively define the urban region being modeled. Many markets are of interest within ILUTE, such as housing, labour, commercial, real estate, etc., and are modeled via microsimulation. For example in HoMES, the housing market module of ILUTE, houses are auctioned off one dwelling at a time to interested bidders in a disaggregate implementation of Martinez' Bid Choice theory[42, 29].

Among the major barriers to implementation of integrated urban models since their introduction were such aspects as data hungriness and computational requirements[35]. However, as an increasing amount of aspects of human life becomes digitalized, a wealth of new data is produced and can be used to model and analyze dynamics of urban systems[2, 11]. Furthermore, continuing methodological advances in computing, such as cost-effective High Performance Computing (HPC), detailed GIS-based datasets and machine learning methods, mean that former barriers now represent opportunities for model system development[32]. An example of such digitalization of human activity is the introduction of the Province of Ontario Land Registration Information System (POLARIS) in 1985 by the Government of Ontario[47]. Introduction of POLARIS lead to the creation of an extensive dataset of real estate transactions (land registration records) by the Teranet Enterprises Inc. Due to its high spatial and temporal resolution, Teranet's dataset presents a valuable resource to inform and validate the design of ILUTE and HoMES. In addition, Teranet's dataset plays an important part in the Longitudinal Housing Market Research conducted by UTTRI which is discussed in the following section.

## 2.4 Longitudinal Housing Market Research conducted by UTTRI

There is ample evidence of the role of land use and transportation interactions in determining urban spatial structure. Accessibility and mobility provided by transportation systems drive economic development and impact travel behaviour and location of households and firms. Similarly, urban sprawl and location of activities drive travel demand and the need for building transport networks. Land values in a metropolitan region are an outcome of such land use-transportation interactions.

Researchers at the University of Toronto Transportation Research Institute are conducting longitudinal analysis to identify trends and changes in land value distributions using the Teranet sales data over the 30-year period (1986-2016) in the Greater Toronto-Hamilton Area (GTHA). The aim is to understand the spatial-temporal effects of changes in socio-economic characteristics, transportation accessibility and built-environment on land values. Data from various sources (e.g., Census, Transportation Tomorrow Survey, etc.) is required to be merged at the land parcel level to enhance datasets with additional attributes, while maintaining the ease of data storage and retrieval for analysis as needed. In addition to that, accurate land use information needs to be added to Teranet records to allow separating sales data by major property types.

Brief legal and historical background of Teranet's dataset, main challenges of working with it and the proposed solution are discussed in the remainder of this chapter.

## 2.5 POLARIS: electronic system of land registration in Ontario

All land owned in Canada is registered in a public land registry in the applicable province. Each province and territory in Canada has its own land registry system, whether it is a land titles system, a registry system or a combination of both, with each system having its own rules. The registry system is a public record of documents evidencing transactions affecting land. In the land titles system, the applicable provincial government determines the quality of the title, and essentially guarantees (within certain statutory limits) the title to, and interests in, the property. As of 2015, most common law provinces and territories in Canada were using the land titles system or were in the process of converting from a registry system to a land titles system[30].

As of 2015, the Province of Ontario has largely converted from registry systems to a land titles system. In 1985, the Government of Ontario initiated the Province of Ontario Land Registration Information System (POLARIS) pilot project for the purposes of the conversion between systems and records automation. The Land Registration Reform Act (Ontario)[49] was introduced in 1990 to facilitate electronic search and registration of properties and the automation of paper-based records. POLARIS was built by the Province to house and process electronic land records, which in turn lead to the creation of an extensive dataset of land registration records managed by Teranet Enterprises Inc. Today, POLARIS is the search/registration and property maintenance system for all automated land records in Ontario.

## 2.6 Teranet's dataset, its challenges and the proposed solution

In 1991, the Government of Ontario established a partnership with Teranet, a Toronto-based organization, founded the same year, which provides e-services to legal, real estate, government, financial, and healthcare markets. The partnership was established to convert Ontario's land registration system to a more modernized electronic title system. The project involved taking a 200-year-old paper-based system and creating a database with electronic records for more than five million parcels of land. Teranet converted all qualified Registry properties in Ontario to the Land Titles system and automated existing paper Land Titles parcels. As a result, 99.9% of property in Ontario was parcelized and administered under the Land Titles system. Teranet fully automated the conversion of millions of paper-based documents and records into the Ontario Electronic Land Registration System (ELRS)[48].

Teranet's dataset presents an extensive historical record of real estate transactions recorded in the

Province of Ontario since the beginning of XIX century. However, its available version also suffers from severe lack of features describing each transaction, which makes meaningful analysis or modeling difficult. At the same time, each Teranet transaction has a timestamp (date) and location information (x and y coordinates) and thus can be joined to variety of other geocoded urban data sources, such as Census demographics, Transportation Tomorrow Survey (TTS) and parcel-level land use information. However, joining these data sources together requires additional considerations, as they use different spatial units and are available at different temporal spans, as will be discussed in chapter 3.

One of the major attributes missing from the available version of Teranet’s dataset is the information about the type of property being transacted, with records of various categories of residential, commercial and industrial properties all being mixed together in the same dataset. This introduces a major limitation on how Teranet’s data can be used, limiting the ability to separate the transactions of different property types, identify submarkets of the housing market and see its fine-scale characteristics and dynamics. Detailed parcel-level land use information could offer some degree of filtering and can be joined from such sources as DMTI Spatial’s land use and detailed land use manually collected by the Department of Geography at the University of Toronto.

However, these sources of land use information also have their limitations:

- DMTI’s land use data does not offer any split between subcategories of residential properties and only covers the period of 2001–2014
- land use from the Department of Geography is a lot more detailed and accurate, but has been collected at a single point in time over the summer of 2012 and 2013
- neither of the available land use sources covers the full span of the Longitudinal Housing Market Research conducted by UTTRI (1986-2016)

At the same time, detailed land use from the Department of Geography offers the opportunity to train a machine learning model capable of recognizing certain property types that have characteristically different behavior on the housing market (e.g., high / low volume of transactions, median price ratio, etc.). Teranet records coming from a particular parcel (x and y coordinate pair) or pin (same parcel can include different pins, i.e. apartment complexes) can be combined to produce new features that characterize the behavior of this parcel or pin on the housing market. Combining these new features with spatially-joined variables from Census and TTS can yield a dataset that can be used to train and test a classification algorithm capable of determining the parcel land use based on the nature of its behavior on the housing market at transaction level (recognizing changes of land use with time).

The primary focus of this master’s thesis consists of:

- augmenting Teranet’s dataset with TTS, Census and land use variables, while maintaining the integrity of the spatial and temporal relationships between different data sources and ensuring the ease of storage and retrieval as needed for the Longitudinal Housing Market Research conducted by UTTRI
- investigating the opportunity for the implementation of a machine learning algorithm to classify land use based on the housing market dynamics

## 2.7 Chapter summary

The complex interaction of land use and transportation has been treated via simulation models since 1950s. Over the decades, the modeling paradigm has evolved from the highly aggregated and "top-down" representations of processes, such as gravity-based models, to highly disaggregated and "bottom-up", defining the macro state of the system through countless interactions of its microcomponents. These changes are represented in the family of microsimulation models which can define the state of an urban region in terms of the individual persons, households, dwelling units, firms, etc. and evolve from a know base case into a future state via simulation. It also became more clear with time that land use and transportation systems are deeply interconnected, and thus must be modeled together to incorporate feedback relationship, as is attempted by integrated urban models. An example of the class of integrated microsimulation urban models, Integrated Land Use, Transportation, Environment, or ILUTE, model system has been developed at the University of Toronto for the Greater Toronto–Hamilton Area.

Since their introduction, among the major barriers to the implementation of integrated models were their data hungriness and computational requirements. However, continuing methodological advances in computing and new datasets created through increased digitization of human activity present opportunities for further improvement and validation of ILUTE and its modules. An example of new emerging data sources is Teranet's dataset of real estate sales that was created after the introduction of POLARIS land registration system by the Province of Ontario in 1985. Teranet's dataset also plays an important part in the Longitudinal Housing Market Research conducted by UTTRI. One of the major challenges of working with Teranet's data is the lack of features describing each transaction, namely the type of property being sold. As Teranet's records have timestamps (dates) and coordinates of parcel centroids for each record, they can be joined with other data sources via temporal and spatial relationships that will be discussed in chapter 3.

## Chapter 3

# Spatial and temporal relationships between data sources

As was discussed in section 2.6, one of the main challenges of working with Teranet’s data is the lack of available features. At the same time, Teranet records have timestamps (dates) and location information (x and y coordinates) and thus can be joined to a variety of other urban data sources, such as Census demographics, Transportation Tomorrow Survey (TTS) and parcel-level land use information. However, as will be discussed in this chapter, these data sources use different spatial units and are available at different temporal spans; therefore, special considerations need to be taken when joining data from these sources with respect to their temporal and spatial relationships to ensure semantic interoperability.

Different data sources joined to Teranet’s dataset are described in this chapter, the implementation of the spatial and temporal relationships via the standardized data preparation workflow in Python and a PostgreSQL relational database are described in chapter 4.

### 3.1 Description of the data sources used

The following data sources are combined into the GTHA housing market database that was created as a part of this master’s thesis:

1. Teranet’s dataset

Due to the introduction of POLARIS by the Province of Ontario in 1985 (discussed in section 2.5), Teranet’s dataset includes a complete population of real estate transactions recorded in Ontario from 1985 up to October of 2017 (records prior to 1985 appear to be incomplete, see chapter ??). Since Teranet’s dataset has a high number of records, it can be used to investigate aspects relating to the housing market at a very fine spatial and temporal scale.

2. Select variables from the Census of Canada

One of the major sources of demographic and statistical data in Canada are the datasets collected under the national Census program. Census data provide valuable insight into the latest economic, social and demographic conditions and trends in Canada and is used to plan important public services. Statistics Canada collects every five years the national Census of Canada and disseminates the information by a range of geographic units, also referred to as ”Census geography”[28].

### 3. Select variables from the Transportation Tomorrow Survey (TTS)

Another major source of information for most transportation planning studies concerned with Southern Ontario is the Transportation Tomorrow Survey (TTS), an origin-destination travel survey[14]. The Transportation Tomorrow Survey (TTS), undertaken every five years since 1986, is a cooperative effort by local and provincial government agencies to collect information about urban travel in southern Ontario. TTS represents a retrospective survey of travel taken by every member (age 11 or over) of the household during the day previous to the telephone or web contact. The information collected and the method of collection has remained relatively consistent over the seven surveys and includes characteristics of the household, characteristics of each person in the household, and details of the trips taken by each member of the household, including details on any trips taken by transit[3].

### 4. Land use from DMTI Spatial Inc. by year (2001-2014)

DMTI Spatial Inc., a Digital Map Products company, is a major provider of location based information in Canada. DMTI has been providing industry leading enterprise Location Intelligence solutions for more than a decade to Global 2000 companies and government agencies[16].

### 5. Detailed land use information from University of Toronto's Department of Geography collected in 2012 and 2013

The detailed land-use data provided by University of Toronto's Department of Geography is a combination of parcel boundaries (from Teranet) and manually coded land-use data produced using Google maps and streetviews; it was collected by Prof. Andre Sorensen and Prof. Paul Hess's research project.

## 3.2 Spatial relationships between data sources

Most urban areas are divided into zones or planning areas on the basis of maintaining similar population sizes and following built or natural boundaries like roads or rivers. Census geography follows a certain hierarchy defined by Statistics Canada, with the largest top-level divisions being provinces and territories, and the lowest-tier divisions to which census data is disseminated being Dissemination Areas (DAs)[45]. Statistics Canada defines a dissemination area as a small area composed of one or more neighbouring dissemination blocks, roughly uniform in population size targeted from 400 to 700 persons to avoid data suppression[44].

To simulate the changes in accessibility, metropolitan regions are usually broken down into a set of small geographic zones, similar (or in many cases identical) to the set of zones used for regional travel forecasting. For TTS variables, the finest level of spatial aggregation is that of the Traffic Zone also referred to as Traffic Analysis Zone (TAZ). The Traffic Zone is a polygon which typically falls along the centre line of roads or the natural geographic boundaries[15]. Not as a rule, but the TAZs roughly follow census tract boundaries, which are slightly bigger than DA boundaries. TTS data has been collected for changing TAZ boundaries or in other words, different zone systems due to growing population and expanding extents of the survey in the GTHA region over the years. To make the TTS data consistent for comparing over all years from 1986 to 2016, the Data Management Group (DMG), the custodian of the dataset derived from TTS, made all surveys available in the 2001 zone system, for convenience of

researchers (any zone system could have been chosen for that matter). UTTRI used the 2001 TAZ system to model travel times for the GTHA on EMME for all TTS years based on the origin-destination trip data collected in the survey. The travel time data was used to create further transportation accessibility variables.

Land use data collected by DMTI and by the Department of Geography uses the spatial unit of a parcel polygon. Teranet records have attributes representing x and y coordinates matching parcel centroids.

Below is the summary of the spatial units being used by the different data sources:

- Point data
  - Teranet
- Parcel-level data (polygons)
  - detailed land use from the Department of Geography
  - land use from DMTI
- DA-level data (polygons)
  - Census variables
- TAZ-level data
  - TTS variables

When joining these data sources, difference in spatial units needs to be respected, which can be more challenging when spatially joining polygons with polygons, since it might require area weighted spatial interpolation of data to a common unit of analysis. In addition, polygon-based data can also vary with time, as is the case with DMTI's land use information, which is available by year. To simplify relating different polygon-based data sources with each other, all of them can be brought together at a single level of time-indexed points, such as Teranet transactions. This allows flexibility in combining data from these sources at transaction level while maintaining the integrity of spatial and temporal relationships through polygon-to-point spatial joins. Implementation of these relationships is described in chapter 4, temporal relationships between different data sources are described in the following section.

### 3.3 Temporal relationships between data sources

In addition to using different spatial units, data sources joined with Teranet's dataset are available at different temporal spans:

- Teranet records have individual timestamps (date) on each record
- Census and TTS variables are available once in 5 years
- DMTI's land use data is available by year and covers a time span from 2001 to 2014
- Detailed land use from the Department of Geography was collected at a single point in time during the summers of 2012 and 2013



Temporal matching between Teranet records and DMTI data can be done directly: DMTI land use for each year from 2001 to 2014 can be spatially joined with a subset of Teranet records from the corresponding year. Since the detailed land use provided by the Department of Geography was collected at a single point in time, it can be joined to all Teranet records, keeping in mind that it will be the most accurate around its time of collection in 2012 and 2013.

DATA SOURCES	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017
Census																			
TTS																			
DMTI land use																			
DMTI postal codes																			
DMTI EPOI																			
EMME network calculations																			
Teranet																			
Dept of Geography land use																			
Fuel price in Toronto																			

Figure 3.1: Temporal spans of the data sources used in the GTHA housing market database.

As for the Teranet and Census / TTS variables, they can be matched in a number of ways:

1. Direct match with appropriate Teranet subsets

- match Census / TTS variables only with Teranet records from the corresponding year (for example, Teranet records from 2016 matched with 2016 Census / TTS variables)
- benefits:
  - Census variables would be composed of the actual values produced by the survey
- disadvantages:
  - limited use of Teranet since only Teranet data from Census years can be used

2. Interpolation of discrete Census / TTS variables

- discrete Census / TTS variables can be turned into continuous via interpolation
- benefits:
  - closest temporal match between Teranet and Census / TTS variables
- disadvantages:
  - additional assumptions need to be made for each Census / TTS variable

3. Assign temporal spans to each Census / TTS survey as new features to Teranet records

- each Census / TTS survey can be assigned a temporal span of 5 years representing a group of Teranet records to which its variables can be matched (for example, the Census of 2016 gets assigned Teranet record from 2014 to 2018)
- benefits:
  - avoid interpolation assumptions
- disadvantages:
  - step-change in Census / TTS variables

To avoid additional assumptions and use the recorded values of Census and TTS variables, the third option for matching Teranet records with Census / TTS variables is chosen. Each Census / TTS survey is assigned a 5-year time span centered at the survey year (i.e., 2014–2018 for 2016 survey year) and new foreign keys are introduced to Teranet records to allow matching with 5-year time spans of Census / TTS variables. Implementation of temporal relationships for Teranet records will be described in chapter 4.

Figure 3.1 presents the temporal spans assigned to the data sources for joining with Teranet records.

### 3.4 Chapter summary

Data sources joined to Teranet’s dataset use different spatial units and are available at varying temporal spans. These relationships need to be respected when combining the data sources together to ensure semantic interoperability. The integrity of spatial relationships can be ensured by spatially joining all polygon-based data sources to Teranet points. Temporal relationships between DMTI land use and Teranet are incorporated by performing separate spatial join operations between DMTI land use data for each year and a corresponding Teranet subset. For Census and TTS variables, additional foreign keys are introduced assigning 5-year spans to each Census / TTS survey. These foreign keys indicate which Teranet records should be joined to a particular Census or TTS survey. Implementation of these relationships via a standardized data preparation workflow in Python and a PostgreSQL relational database are described in chapter 4.

## Chapter 4

# Data preparation

The "wicked" nature of transportation and land use interaction introduced in chapter 2 dictates the need to iteratively "re-solve" transportation and land use planning problems instead of focusing on finding some single "optimal solution". This approach resembles the methodologies typically employed for data science projects, where the sequence of steps is iterated over, producing a more meaningful solution on each new iteration of the cycle, as defined by such process models as CRISP-DM[43]. Similarly, data preparation can be followed in a linear manner, but is very likely to be iterative in nature[10].

Data preparation plays a critical role in research projects:

- it is a prerequisite for any meaningful analysis
- data quality and the amount of useful information that it contains can determine the success of application of machine learning algorithms[40]
- it is often required to allow the introduction of constraints necessary for implementation of RDBMS

To facilitate easy modification and replication of the data preparation process for data sources related to the GTHA housing market, a streamlined data preparation workflow using Python via a series of jupyter notebooks has been established as a part of this master's thesis. It accomplishes three main objectives:

- clean Teranet dataset and correct its records for consistency
- introduce new keys that allow efficient joining of other data sources such as Census or TTS, while maintaining the integrity of spatial and temporal relationships that were discussed in chapter 3
- engineer new features that can be used by the machine learning algorithm along with the features from the joined datasets to classify land use, which will be discussed in chapter 5

This chapter introduces the concepts of "Tidy Data" and database normalization and outlines the standardized data preparation workflow for all data sources related to the GTHA housing market.

## 4.1 Tidy data and database normalization

Hadley Wickham in his paper "Tidy Data"[50] formalized the way how a shape of the data can be described and what goal should be pursued when formatting data. The tidy data standard is closely related to Edgar F. Codd's relational algebra and has been designed to facilitate initial exploration and analysis of the data, and to simplify the development of data analysis tools that work well together. As an integral part of his relational model, Codd[12] proposed a process of database normalization, or restructuring of a relational database in accordance with a series of so-called normal forms in order to reduce data redundancy and improve data integrity. Normalization entails organizing the columns (attributes) and tables (relations) of a database to ensure that their dependencies are properly enforced by database integrity constraints. The principles of "tidy data" essentially reformulate Codd's ideas in statistical language.

According to Wickham[50], "tidy data" is a standard way of mapping the meaning of a dataset to its structure. A dataset is "messy" or "tidy" depending on how rows, columns and tables are matched up with observations, variables and types.

In "tidy data":

1. Each variable forms a column.
2. Each observation forms a row.
3. Each type of observational unit forms a table.

This is Codd's 3rd normal form[12], but with the constraints framed in statistical language, and the focus put on a single dataset rather than the many connected datasets common in relational databases. "Messy data" is any other arrangement of the data.

The structure of Teranet's dataset conforms with the "Tidy data" format. Contrary to Teranet, tables with combined selected Census and TTS variables have variables for different Census and TTS years recorded as columns. This needs to be addressed by "melting" these tables and introducing a new attribute 'year' to be used as a part of a composite foreign key when joining with Teranet records.

Census and TTS tables were "melted" into the "tidy data" format:

- each Census / TTS variable now forms a single column
- each value of a variable is indexed by a composite primary key constituting of spatial identifier and year of the survey
- spatial identifiers are 'DAUID' or 'TAZ\_O' (unique identifier for Dissemination Areas or Traffic Analysis Zones introduced in section 3.2)

Introduction of new foreign keys is described in the following section.

## 4.2 Introduction of new keys and attributes via spatial and temporal relationships

As no combination of columns constitutes a candidate key for Teranet records (unique identifier to be used in RDBMS), a surrogate key (artificial unique identifier for RDBMS) is added to the Teranet dataset via a new attribute 'transaction\_id'. Thus, Teranet's dataset fits into a normalized database, with the new attribute 'transaction\_id' as its primary key. As was discussed in section 4.1, the "melted" Census and TTS tables have composite primary keys consisting of a unique spatial identifier ('DAUID' or 'TAZ.O', respectively) and the year of the survey.

To implement the spatial and temporal relationships between the data sources discussed in chapter 3, a number of new foreign keys needed to be introduced to Teranet records. The new foreign keys either represent spatial identifiers (such as 'david' or 'taz\_o', corresponding to DA or TAZ within which a Teranet record is located), or an attribute identifying the year of the Census or TTS survey to which this Teranet record can be joined. Foreign keys representing spatial identifiers are added through a series of spatial joins while foreign keys identifying temporal spans are produced based on temporal relationships established in section 3.3.

New spatial identifiers were introduced to each Teranet record via a series of spatial joins:

1. 9'039'241 Teranet points were joined with 9'182 polygons of Dissemination Areas (DAs) for GTHA used by Census variables
  - Teranet records with coordinates falling outside of GTHA boundary were filtered out
  - From the original 9'039'241 records, 6'803'691 remained in the dataset
  - New foreign keys 'david', 'csduid' and attribute 'csdname' were added to each Teranet record
2. 6,803,691 Teranet points were joined with 1,716 polygons of Traffic Analysis Zones (TAZ) used by TTS variables
  - New foreign key 'taz\_o' was added to each Teranet record
3. 6,803,691 Teranet points were joined with 525 polygons of Forward Sortation Areas (FSA) and 555,668 polygons of postal geography from DMTI's Platinum Postal Geography Suite
  - New foreign keys 'fsa' and 'pca\_id' and attribute 'postal\_code\_dmti' were added to each Teranet record
  - These keys are not currently used for joining any variables, but were added to expand the potential for relating datasets
4. 6,803,691 Teranet points were joined with 1,664,862 polygons of parcel-level detailed land use provided by the Department of Geography
  - New foreign keys 'pin\_lu', 'landuse' and 'prop\_code' were added to each Teranet record
  - Foreign keys 'landuse' and 'prop\_code' are codes that can be converted to land use categories that were used by the Department of Geography for GTA and Hamilton, respectively

- For records from Hamilton, 'prop\_code' was converted to categories used by GTA land use and reassigned to 'landuse', bringing GTA and Hamilton records to a single system of land use categories
5. Subsets of Teranet points were joined with corresponding yearly polygons of parcel-level land use from DMTI
- New attribute 'dmti\_lu' was added to each Teranet record

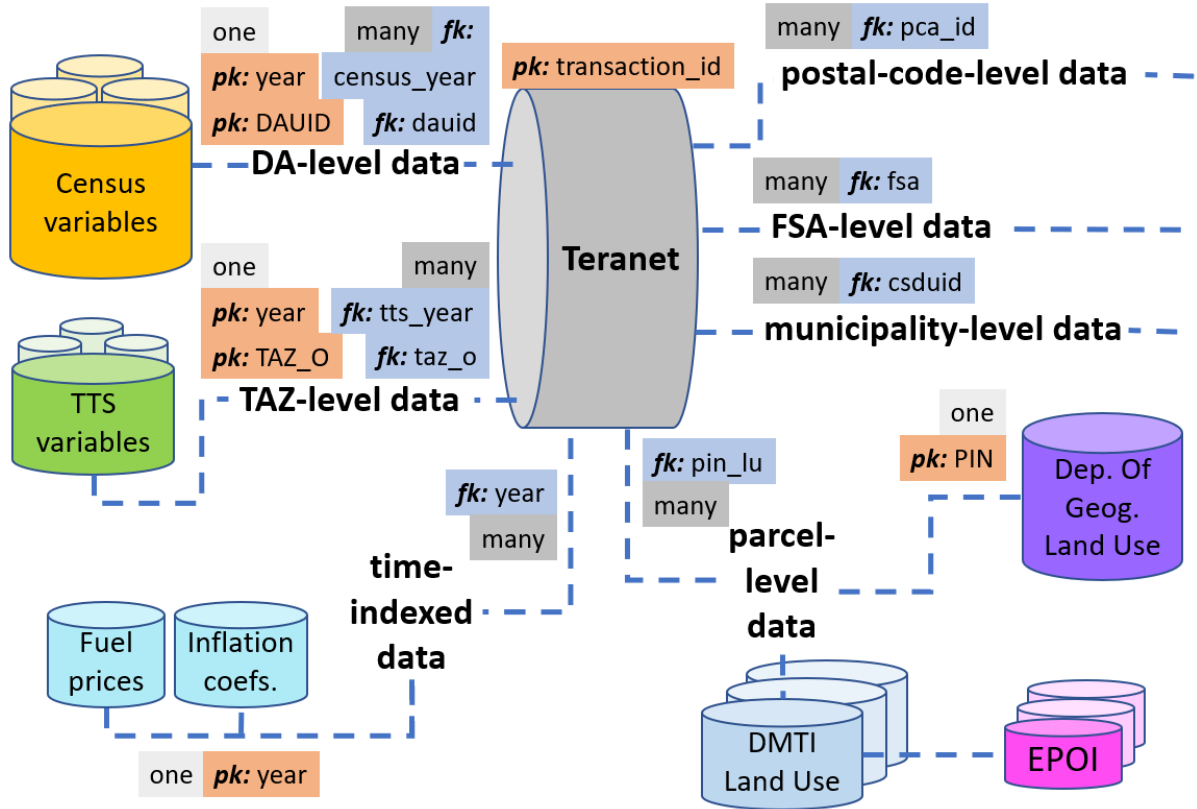


Figure 4.1: Relationships between datasets introduced during data preparation were then used to set up referential integrity constraints of the new PostgreSQL database for GTHA housing market data.

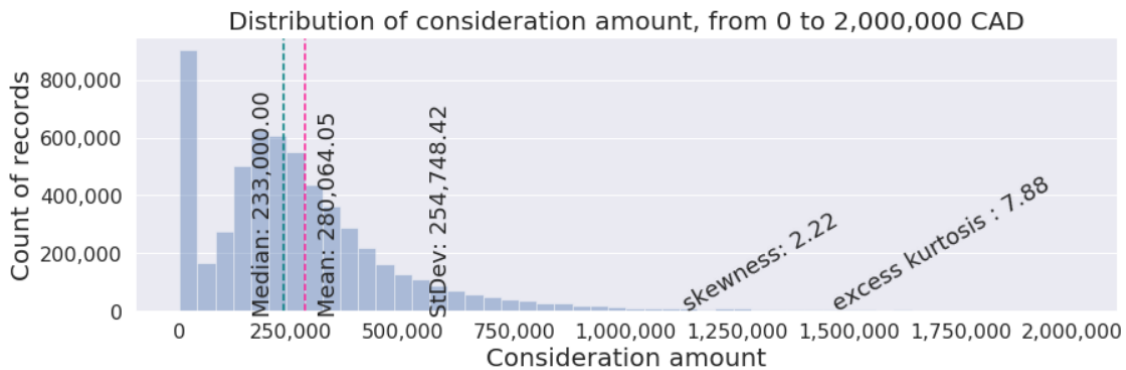
Foreign keys representing temporal identifiers are generated from the registration date of each Teranet record, matching each year of Teranet records to a corresponding 5-year span covered by a Census or TTS survey, as was discussed in section 3.3. Diagram of relationships between datasets and their primary and foreign keys is presented on figure 4.1

Following the steps described above ensures that the integrity of spatial and temporal relationships is maintained when combining attributes from different data sources at Teranet transaction level. For example, Teranet records from 2007 would be spatially joined with DMTI land use data from 2007, and are matched by their attributes 'census\_year' and 'tts\_year' to Census and TTS variables from 2006 Census and TTS surveys. Census and TTS variables can be joined by appropriate 'david' and 'taz\_o' (composite foreign keys are used when joining), and thus all data sources can be spatially and temporally aligned at the level of Teranet transactions.

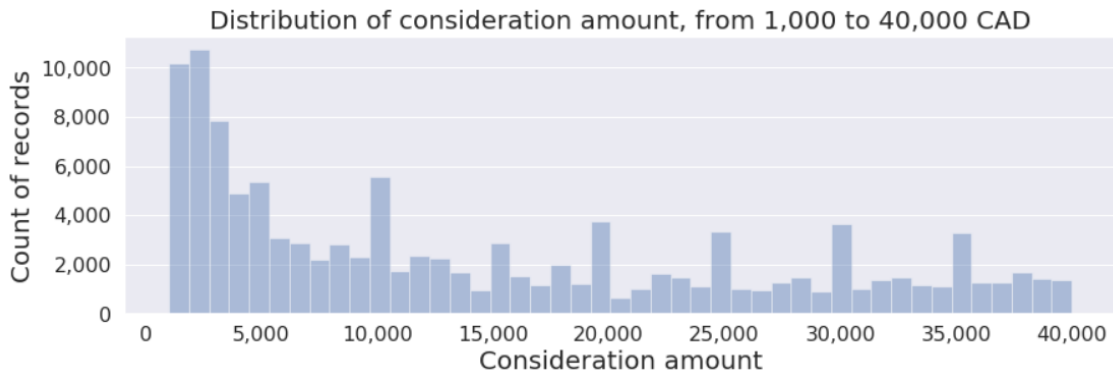
Relationships introduced via the operations described in this section formulate referential integrity constraints that have been used to set up a PostgreSQL database of GTHA housing market data.

### 4.3 Outliers

Outliers present an issue for machine learning as they lead to difficulties in visualizing data and, more importantly, can degrade the predictive performance of an algorithm. In addition, outliers, if left unscaled, can significantly slow down or even prevent the convergence of many gradient-based algorithms, such as Logistic Regression[23]. However, an exact definition of what constitutes an outlier within a given dataset presents a challenge in itself, since it depends on hidden assumptions regarding the data structure and the applied detection method[7]. As defined by Hawkins, an outlier is an observation that deviates so much from other observations as to arouse suspicion that it was generated by a different mechanism[17].



(a) From 0 to 2'000'000 CAD



(b) From 0 to 40'000 CAD

Figure 4.2: Outliers at the bottom end of the price distribution most likely represent gift transactions. Since no clear break can be identified, 10'000 CAD was used as the bottom cut-off threshold for consideration amount to filter Teranet records.

In terms of transaction price distribution, outliers at both the high and the low end are present in Teranet's dataset:

- There is a high number of transactions with very low consideration amounts (starting from a few dollars) that most likely represent transactions recording gifts of property, where some symbolic

consideration amount has been used. A large spike can be seen on the low end of the distribution of consideration amount presented on figure 4.2; this spike corresponds to gift transactions with symbolic prices.

- At the same time, since the dataset includes all types of property transactions, some records have very high values (ranging in the hundreds of millions and billions of dollars) that most likely correspond to transactions recording sales of large commercial and industrial properties or whole residential buildings.

Since low outliers represent transactions that are not useful for analysis, they were removed from the dataset. However, there seems to be no way to establish what constitutes a reasonable bottom cut-off threshold, as there is no criteria available to determine a gift transaction and there is no distinct break in the price distribution. Since there seems to be an exceptionally large spike of transactions with consideration amount under 10'000 CAD, they were considered to be low outliers and were removed from Teranet's dataset, further reducing the number of records from 6,803,691 to 5,188,513.

In case of outliers at the high end of price distribution, they most likely correspond to transactions of expensive commercial and industrial property or whole residential buildings. Since these transactions are useful for research questions concerning commercial and industrial property, they were left in the dataset and instead have been marked with special new attributes "outlier" using different criteria. Since, again, there was no clear criteria available as to what would constitute an outlier, instead of using a single criterion, seven different Boolean variables were added describing whether a record belongs to outliers according to a particular condition or not. For example, feature 'outlier\_y\_10' is a Boolean variable capturing if the price of a record, corrected for inflation, is over 10 times greater than the median price of all records for the corresponding year. Table 4.1 presents the outlier categories that were used and the count of records that each of them marked to belong to outliers.

Category	Count of records
outlier_y_3	251,537
outlier_y_5	137,814
outlier_y_10	84,260
outlier_y_20	53,866
outlier_xy_2	160,662
outlier_xy_4	57,766
outlier_xy_10:	35,778

Table 4.1: Categories of outliers and the count of records that they produced. Categories coded 'y' were compared to median price of all records for that year; categories coded 'xy' were compared to median price of all records for that coordinate pair; a number represents the threshold that the ratio needed to exceed for a record to be considered an outlier.

These new "outlier" attributes were used when filtering records during the assignment of reduced land use classes, which will be discussed in section 5.1; in addition, they were tested as a part of the original full feature set for classification of land use along with the other new features engineered from Teranet's dataset, which will be discussed in section 4.4. However, in the case of classification, all of the new "outlier" features were filtered out in the process of feature selection, which will be discussed in section 5.3. Instead, algorithms performed better with features that numerically represent price, such as 'med\_price\_xy', computed as the median consideration amount (in 2016 dollars) for a particular xy coordinate pair.



Some of the new features that have been engineered from original Teranet attributes, which will be discussed in section 4.4, also contain outliers on the high end of their distributions. For example, in the case of the feature 'xy\_prev\_sales' which captures rolling count of Teranet records from a coordinate pair, count of records from a coordinate pair that corresponds to a house would be in the order of  $10^1$ , while count of records from a coordinate pair that corresponds to an apartment building would be on the order of  $10^2 - 10^4$ . In these cases, outliers contain valuable information that can be used to classify Teranet records by property type, but these outliers still present a challenge for the convergence of gradient-based algorithms. To keep these records in the dataset while facilitating faster convergence, feature scaling was performed and will be discussed in section 5.5.

## 4.4 Engineering new features for the classification algorithm

In addition to producing new keys for joining datasets, a number of the new features was engineered from original Teranet attributes to be tested with the classification algorithm (discussed in chapter 5). These new features were intended to give each Teranet record spatial and temporal "context" of the housing market dynamics by grouping records using different criteria. For example, feature 'xy\_prev\_sales' was computed as the rolling count of Teranet records coming from a particular coordinate pair; feature 'price\_to\_med\_year' captures the ratio of consideration amount of a record to the median consideration amount of all Teranet records for the corresponding year, etc.

The following features have been added to each Teranet record from 1985 to 2017 ('xy' represents 'x' and 'y' coordinates concatenated together as strings, used to group together all records coming from a particular coordinate pair):

- 'price\_2016': consideration amount (in 2016 dollars) using the coefficients from the Inflation Calculator provided by the Bank of Canada[4]
- 'pin\_total\_sales': count of Teranet records grouped by 'pin'
- 'xy\_total\_sales': count of Teranet records grouped by 'xy'
- 'pin\_prev\_sales': rolling count of Teranet records grouped by 'pin'
- 'xy\_prev\_sales': rolling count of Teranet records grouped by 'xy'
- 'xy\_first\_sale': a Boolean variable indicating whether it is the first record coming from this 'xy'
- 'pin\_years\_since\_last\_sale': difference in years since the last record coming from this 'pin'
- 'xy\_years\_since\_last\_sale': difference in years since the last record coming from this 'xy'
- 'xy\_years\_to\_next\_sale': difference in years to the next record coming from this 'xy'
- 'da\_days/years\_since\_last\_sale': difference in days/years since the last sale that occurred on this Dissemination Area
- 'xy\_sale\_next\_6m': a Boolean variable indicating whether there will be another sale on this 'xy' in the upcoming 6 month

- `'pin_price_cum_sum'`: cumulative sum of price (in 2016 dollars) for all Teranet records coming from this `'pin'`
- `'xy_price_cum_sum'`: cumulative sum of price (in 2016 dollars) for all Teranet records coming from this `'xy'`
- `'pin_price_pct_change'`: percentage change of price (in 2016 dollars) from the last Teranet record from this `'pin'`
- `'xy_price_pct_change'`: percentage change of price (in 2016 dollars) from the last Teranet record from this `'xy'`
- `'price_da_pct_change'`: percentage change of price (in 2016 dollars) from the last Teranet record from this Dissemination Area
- `'med_price_xy'`: median price (in 2016 dollars) for all Teranet records from this `'xy'`
- `'med_price_year'`: median price (in 2016 dollars) for all Teranet records for this year
- `'price_to_med_xy'`: ratio of price (in 2016 dollars) to the median price of all records for this `'xy'`
- `'price_to_med_year'`: ratio of price (in 2016 dollars) to the median price of all records for this year
- `'outlier_y_3'`: a Boolean variable marking as outliers all records with the price more than 3 times greater than median for that year
- `'outlier_y_5'`: a Boolean variable marking as outliers all records with the price more than 5 times greater than median for that year
- `'outlier_y_10'`: a Boolean variable marking as outliers all records with the price more than 10 times greater than median for that year
- `'outlier_y_20'`: a Boolean variable marking as outliers all records with the price more than 20 times greater than median for that year
- `'outlier_xy_2'`: a Boolean variable marking as outliers all records with the price more than 2 times greater than median for that `'xy'`
- `'outlier_xy_4'`: a Boolean variable marking as outliers all records with the price more than 4 times greater than median for that `'xy'`
- `'outlier_xy_10'`: a Boolean variable marking as outliers all records with the price more than 10 times greater than median for that `'xy'`

These new features were combined with TTS and Census variables via spatial and temporal relationships that were introduced in chapter 3 and were used to train and test a classification algorithm to classify land use at Teranet transaction level, which is discussed in chapter 5.

Some of the new engineered features, such as `'xy_years_since_last_sale'` and `'xy_years_to_next_sale'` contain missing values in each case of the first or the last Teranet record coming from a coordinate pair. Since classification algorithms cannot handle inputs with missing values, during classification these values have been replaced with the corresponding median, which will be discussed in section 5.2. The recorded version of Teranet's dataset with the new features and keys would still have them as missing.

## 4.5 Chapter summary

The principles of "Tidy data", closely related to Codd's principles of database normalization, formalize the way how a shape of the data can be described and what goal should be pursued when formatting data. The "tidy data standard" facilitates initial exploration and analysis of the data and simplifies the development of data analysis tools that work well together. The structure of Teranet's dataset conforms with the "Tidy data" format. Contrary to Teranet, tables with combined selected Census and TTS variables had variables for different Census and TTS years recorded as columns and thus needed to be "melted". A new attribute 'year' was introduced to these tables to be used as a part of a composite foreign key when joining with Teranet records.

New foreign keys representing spatial identifiers, such as 'dauid' and 'taz\_o', were added to Teranet records via a series of spatial joins. During the first join with Dissemination Area geometry used by Census, Teranet records with coordinates falling outside of GTHA have been filtered out. Furthermore, Teranet records with consideration amount under 10'000 CAD were considered to be outliers at the low end of price distribution and were removed from the dataset. Outliers at the high end of the price distribution were not removed, but instead were marked by seven new Boolean attributes using various criteria to define a top outlier. Finally, new features have been engineered to be tested with classification algorithms in chapter 5.

Characteristics of the raw Teranet dataset:

- 9,039,241 rows
- 15 columns

Characteristics of the Teranet dataset after data preparation described in this chapter:

- 5,188,513 rows
- 75 columns

## Chapter 5

# A prototype of a machine learning workflow to classify land use

As was previously mentioned in section 2.6, one of the major features missing from the available version of Teranet’s dataset is the information about the type of property being transacted, which introduces a major limitation on how Teranet’s data can be used. As was described in chapter 4, parcel-level land use information from DMTI and the Department of Geography has been spatially joined to Teranet records, but these sources of land use information also have their limitations:

- DMTI’s land use data does not offer any split between subcategories of residential properties and only covers the period from 2001 to 2014
- land use from the Department of Geography is a lot more detailed and accurate, but has been collected at a single point in time over the summer of 2012 and 2013
- neither of the available land use sources covers the full span of the Longitudinal Housing Market Research conducted by UTTRI (1986–2016)

To address this issue, detailed land use from the Department of Geography can be used as labelled data to train a machine learning model capable of recognizing certain property types that have characteristically different behavior on the housing market. For example, the proposed model would be able to differentiate a detached house from a condo through such features as high / low volume of transactions, ratio of price to median price for that year, etc. Chapter 4 described the production of a dataset that combines the new features engineered from Teranet data with the variables spatially and temporally joined from Census and TTS surveys. In this chapter, this dataset is used to investigate the opportunity to implement a classification algorithm to determine the parcel land use at Teranet transaction level (land use is determined for each Teranet transaction, recognizing the changes of land use on the same parcel with time) based on the housing market dynamics. This way, a machine learning algorithm could provide a scalable solution to automate a labour-intensive task of collecting parcel-level detailed land use and expand the temporal span for which the land use data collected by Department of Geography can be used with accuracy.

## 5.1 Selecting and encoding the target variable

As was expected, different property types have characteristically different behaviour on the housing market. Figure 5.1 shows the distributions of total count of Teranet records per coordinate pair by the top 10 land use categories used in the Department of Geography’s dataset.

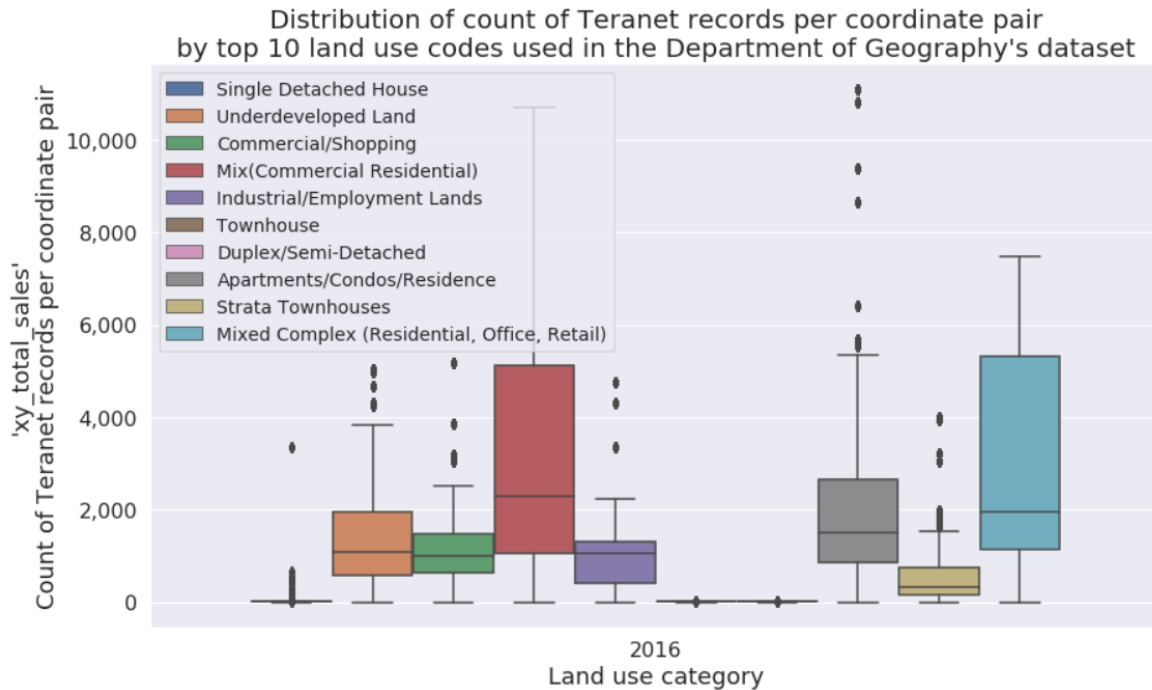
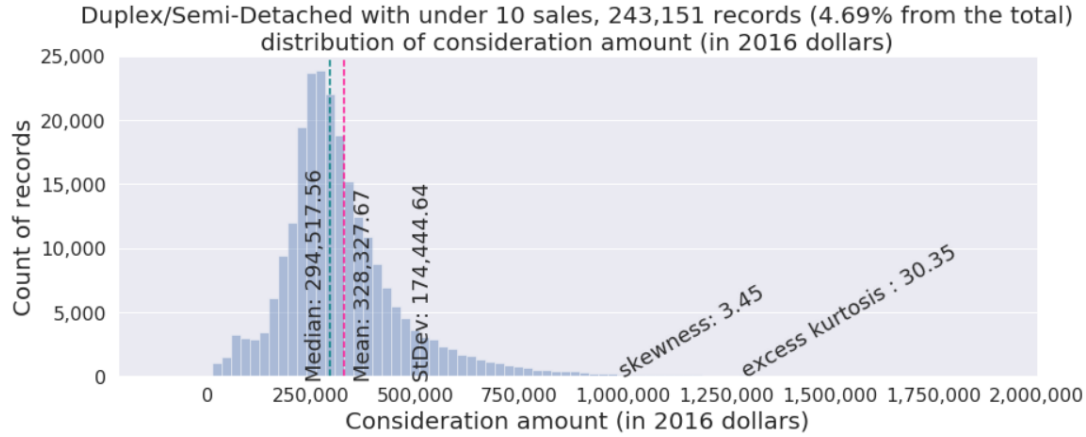


Figure 5.1: Distributions of total count of Teranet records per coordinate pair by the top 10 land use categories. It can be seen that detached houses, duplexes, and townhouses (“collapsed” boxplots) have much lower number of records per coordinate pair. It can also be seen that there are some outliers present on the category of detached houses (first boxplot from the left). These outliers most likely represent mislabelled records, as it is very unlikely that a detached house can have almost 4’000 sales.

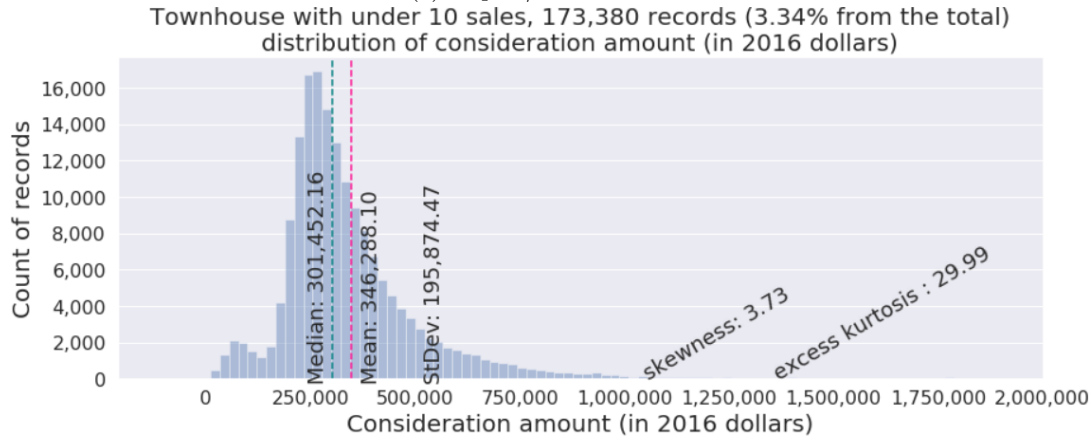
The target variable used for the purposes of investigating the possibility of land use classification using the housing market dynamics was constructed by reducing the land use codes found in the Department of Geography’s land use dataset. Department of Geography’s land use categories that have the highest counts of Teranet records have been used to reduce all different property types to the three major land use classes. As many machine learning algorithms are subject to a frequency bias in which they place more emphasis on learning from data observations which occur more frequently, the three classes were selected to have a comparable number of records between themselves to achieve a more balanced dataset. In addition, chosen groupings combine categories that have a similar distribution of price and count of sales per coordinate pair between land use categories that were grouped together to form a single class. For example, detached and semi-detached houses and townhouses would have a much smaller frequency of transaction and a higher median price per coordinate pair when compared to condos and strata townhouses. Figure 5.2 shows an example of such grouping: in class 2, Duplex/Semi-Detached properties are grouped together with Townhouses.

The three target classes that were introduced are:

- Class 0: "condo", including Apartments/Condos/Residence and Strata Townhouses
- Class 1: "house", including Single Detached Houses, Duplex/Semi-Detached and Townhouses
- Class 2: "other", including Commercial/Shopping, Mix (Commercial Residential), Industrial/Employment Lands, and everything else



(a) Duplex/Semi-Detached



(b) Townhouse

Figure 5.2: Distribution of price (in 2016 dollars) for two out of the three property types that are grouped together under Class 2: House. Both categories have similar distributions of price and records per coordinate pair between themselves and with the Single Detached Houses, and thus all three categories are grouped together to represent a single class "house".

Classification algorithms, a subcategory of algorithms for supervised learning, predict the discrete unordered categorical class labels of new instances based on past observations. A trained algorithm is capable of using a set of rules that were learned from past observations to distinguish the new instances between the possible classes. Once the classes have been determined and assigned, the target variable must be encoded, as it is considered good practice to provide class labels as integer arrays to algorithms to avoid technical glitches and improve computational efficiency. Class labels are not ordinal and classification estimators in the scikit-learn[39] machine learning library in Python treat them as categorical data that does not imply any order.

## 5.2 Missing values

As was mentioned in section 4.4, some of the new features engineered from Teranet attributes, such as 'xy\_years\_since\_last\_sale' and 'xy\_years\_to\_next\_sale' contain missing values in each case of the first or the last Teranet record coming from a coordinate pair. Since classification algorithms cannot handle inputs with missing values, these values were replaced using the following logic:

- For all the records with missing 'xy\_years\_since\_last\_sale' (first record from a coordinate pair), values were replaced with the median 'xy\_years\_to\_next\_sale' for this subset (median time interval between all future sales from this Teranet subset).
- For all the records with missing 'xy\_years\_to\_next\_sale' (last record from a coordinate pair), values were replaced with the median 'xy\_years\_since\_last\_sale' for this subset (median time interval between all past sales from this Teranet subset).

This replacement affects a large number of Teranet records, as there are many coordinate pairs with a low count of records, but this allows the classification algorithm to work on the majority of Teranet data, which improves the generalization of the algorithm and allows classifying land use of most Teranet records. The features that have missing values also have strong predictive power for land use classes and thus are best kept in the input set with the missing values replaced. This replacement is done during the fitting of classification algorithm and is not recorded in the final version of Teranet's dataset with new features and keys, where the missing values are left as missing.

## 5.3 Dimensionality reduction

Quality of the data plays a critical part in the success of application of a machine learning algorithm, and one of important aspects of data quality is the dimensionality of input space. The input space may contain features that are either redundant or irrelevant; highly correlated features also introduce multicollinearity and can make the model unstable, or too sensitive to small changes in the input data. In machine learning, statistics, and information theory, dimensionality reduction is the process of reducing the number of random variables under consideration; there is a number of reasons for reducing dimensionality of a dataset:

1. reducing the number of features improves computational efficiency and reduces training times
2. in case of a low signal-to-noise ratio in the dataset, dimensionality reduction can improve the predictive performance of an algorithm[40]
3. simpler models are easier to interpret[20]
4. excessive complexity of the model could cause overfitting[40]
5. the curse of dimensionality[6] (in the context of machine learning, the curse of dimensionality describes the phenomena where the feature space becomes too sparse for the size of the training dataset)[40]

There is a number of approaches that could be utilized to reduce the dimensionality of the feature space. There are two main categories of dimensionality reduction techniques:

- feature selection, also referred to as Feature Subset Selection, or (FSS), where a subset of the original features is selected
- feature extraction, where a new feature subspace is constructed from information derived from the original feature set

Feature selection can either be performed manually or by utilizing a feature selection algorithm. Exhaustive evaluation of all possible feature subsets is computationally unfeasible even for a moderate number of features. For example, if we are given a feature set with  $m = 64$  features and want to reduce it to  $n = 11$ , an exhaustive evaluation would involve over  $10^{12}$  possible feature sets:

$${}_mC_n = {}_{64}C_{11} = \binom{64}{11} = \frac{64!}{11!(64-11)!} = 743,595,781,824 \quad (5.1)$$

Feature selection algorithms present a practical approach to feature selection at scale; such algorithms combine a search strategy for proposing new feature subsets with an objective function to evaluate these subsets; objective function plays the role of a feedback signal used by the search strategy to choose between candidate subset. Objective functions are divided into three major groups:

- filters
  - evaluate candidate feature subsets by their information content (e.g., inter/intra class distance, mutual information, etc.)
  - advantages: have faster execution since there generally is no iterative computation; have good generality since they evaluate the intrinsic properties of the data.
  - disadvantages: tend to select large feature subsets due to monotonic objective functions
- wrappers
  - use a classifier to evaluate subsets by their predictive accuracy
  - advantages: have better accuracy since they select subsets based on specific interactions between the classifier and the dataset
  - disadvantages: slower to execute since a classifier needs to be re-trained multiple times; selected feature subset will be specific to the classifier that was used to evaluate the candidate subsets.
- embedded methods
  - feature selection is performed as a part of model construction process (i.e., LASSO method for constructing a linear model with L1 regularization)[25]

A combination of different FSS techniques was applied to reduce the dimensionality of the dataset with Teranet, TTS, and Census variables for classification from 64 to 11 input features. The first method utilized was SelectFromModel in scikit-learn, a wrapper FSS algorithm that selects an optimal size and composition of a feature subset by fitting the provided classifier to training data and getting the importance weights from the fit model. SelectFromModel with Random Forest classifier has determined the optimal number of features to be 18; this number of features was provided to the other FSS algorithms to select the 18 best features from the dataset according to each method.



The second FSS method that was applied to the augmented Teranet dataset comes from the filter family of objective functions: a univariate feature selection algorithm `SelectKBest` in `scikit-learn`; this method selects  $k$  highest scoring features based on their scores on the specified univariate statistical test. Univariate FSS techniques can be useful for better understanding of the data, as they examine each feature individually to determine the strength of its relationship with the target variable. However, due to this fact, these methods will not necessarily result in improvement in the generalization of a learning algorithm, and thus their results have been used for feature selection in combination with other methods discussed in this section. The following statistical tests from `scikit-learn` have been selected for scoring features using `SelectKBest` FSS algorithm:

- Chi-squared stats of non-negative (pre-normalized) features for classification tasks
- ANOVA F-value between label/feature for classification tasks
- Mutual information for a discrete target

The third family of FSS algorithms that were applied were the recursive feature elimination algorithms, wrapper algorithms that evaluate candidate feature subsets and recursively eliminate features until the desired number of dimensions is reached. `RFE` class in `scikit-learn` facilitates feature ranking with recursive feature elimination using the provided external classifier that assigns weights to features. Another recursive selection method used was the Sequential Backward Selection algorithm implemented as a Python class by Raschka and Mirjalili[40]. Sequential feature selection algorithms are a family of greedy search algorithms that are used to reduce an initial  $d$ -dimensional feature space to a  $k$ -dimensional feature subspace where  $k < d$ .

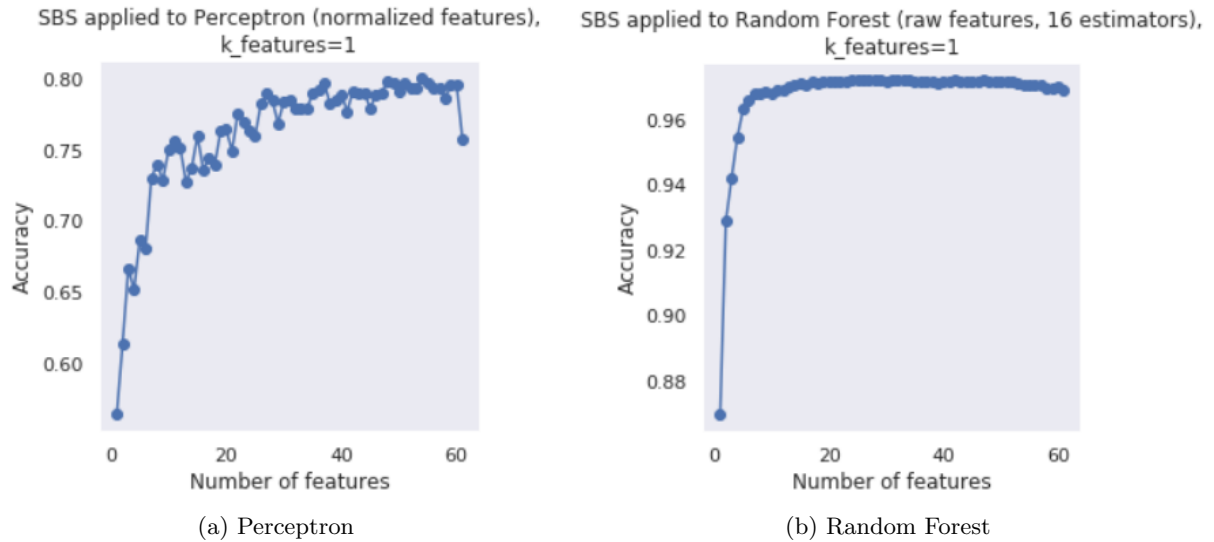


Figure 5.3: Sequential Backward Selection (SBS) is a classic sequential FSS algorithm; on each iteration, it eliminates the feature that results in the minimum drop in performance of the provided classifier. Figure 5.3a shows SBS run on augmented Teranet dataset with Perceptron and figure 5.3b shows SBS with Random Forest. It can be seen that many features can be eliminated without a significant drop in performance of both classifiers, especially in the case of Random Forest.

Sequential Backward Selection (SBS) is a classic sequential FSS algorithm which aims to reduce the dimensionality of the initial feature subspace with a minimum decay in performance of the classifier;

in some cases of model overfitting, SBS can even improve the predictive power of a classifier. On each iteration, a feature is removed using the defined criterion function  $J$  that we want to minimize;  $J$  can simply be defined as the difference in performance of the classifier before and after the elimination of a particular feature. Figure 5.3 presents plots produced by running SBS algorithm on augmented Teranet dataset with Perceptron and Random Forest.

As can be seen from the plots produced by SBS, most features in the augmented Teranet dataset can be eliminated without a significant drop in performance of both linear and tree-based models. As such, final set of features to be used for land use classification was selected by combining the feature subsets selected by each of the FSS methods described above. Features that have been selected by at least four different methods were selected as the final feature subset that was tested with the classification algorithms and used for final classification of land use; figure 5.4 presents the 11 features that were selected using this logic.

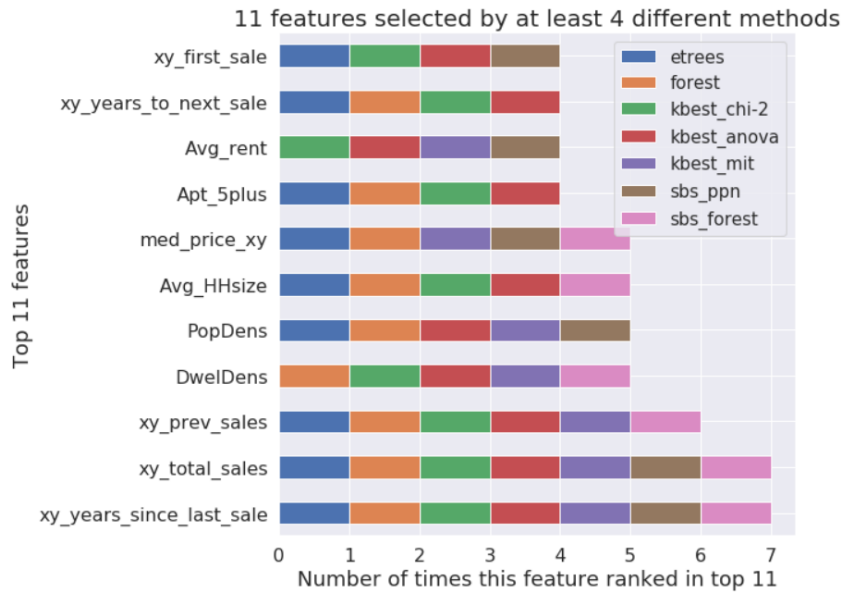


Figure 5.4: FSS results: 11 features that were selected by at least four different feature selection methods. These features were used to select the best-performing model and make the final classification of land use of Teranet records.

From the perspective of embedded feature selection methods, another possible approach to reduce the complexity of a model is to use L1 regularization to penalize large individual weights. L1 regularization introduces a penalty term to the cost function which is taken as the norm of the weight vector defined as:

$$L1 : \|\mathbf{w}\|_1 = \sum_{j=1}^m |w_j| \quad (5.2)$$

L1 regularization is similar to L2 regularization, which is defined as:

$$L2 : \|\mathbf{w}\|_2^2 = \sum_{j=1}^m w_j^2 \quad (5.3)$$

where  $m$  is the number of features.

However, since in the case of L1 regularization, the sum of squares of weights is replaced with the sum of the absolute values of weights, L1 regularization usually yields sparse feature vectors, since most feature weights will be zero[40, 25]. Sparsity of feature vectors can help us get rid of irrelevant features in a high-dimensional dataset; in this context, L1 regularization can be understood as a technique for feature selection[40]. L1 regularization was attempted with Logistic Regression and Linear Support Vector Classifier models in scikit-learn, but did not result in a significant improvement in model performance.

## 5.4 Tuning model hyperparameters

*To determine whether our machine learning algorithm not only performs well on the training set but also generalizes well to new data, we also want to randomly divide the dataset into a separate training and test set. We use the training set to train and optimize our machine learning model, while we keep the test set until the very end to evaluate the final model. Comparing predictions to true labels in the test set can be understood as the unbiased performance evaluation of our model. In scikit-learn a random split into training and test sets can be quickly computed with the `train_test_split` helper function.[40]*

*When evaluating different settings ("hyperparameters") for estimators, such as the  $C$  setting that must be manually set for an SVM, there is still a risk of overfitting on the test set because the parameters can be tweaked until the estimator performs optimally. This way, knowledge about the test set can "leak" into the model and evaluation metrics no longer report on generalization performance. To solve this problem, yet another part of the dataset can be held out as a so-called "validation set": training proceeds on the training set, after which evaluation is done on the validation set, and when the experiment seems to be successful, final evaluation can be done on the test set.[24]*

*However, by partitioning the available data into three sets, we drastically reduce the number of samples which can be used for learning the model, and the results can depend on a particular random choice for the pair of (train, validation) sets. A solution to this problem is a procedure called cross-validation (CV for short). A test set should still be held out for final evaluation, but the validation set is no longer needed when doing CV. In the basic approach, called  $k$ -fold CV, the training set is split into  $k$  smaller sets (other approaches are described below, but generally follow the same principles). The following procedure is followed for each of the  $k$  "folds": A model is trained using of the folds as training data; the resulting model is validated on the remaining part of the data (i.e., it is used as a test set to compute a performance measure such as accuracy). The performance measure reported by  $k$ -fold cross-validation is then the average of the values computed in the loop. This approach can be computationally expensive, but does not waste too much data (as is the case when fixing an arbitrary validation set), which is a major advantage in problems such as inverse inference where the number of samples is very small.[24]*

Figure ?? presents a typical cross-validation workflow to evaluate model performance as summarized in scikit-learn documentation[24].

*A good standard value for  $k$  in  $k$ -fold cross-validation is 10, as empirical evidence shows. For instance, experiments by Ron Kohavi on various real-world datasets suggest that 10-fold cross-validation offers the best trade-off between bias and variance (A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection, Kohavi, Ron, International Joint Conference on Artificial Intelligence (IJCAI), 14 (12): 1137-43, 1995).[22]*

*A slight improvement over the standard  $k$ -fold cross-validation approach is stratified  $k$ -fold cross-validation, which can yield better bias and variance estimates, especially in cases of unequal class pro-*

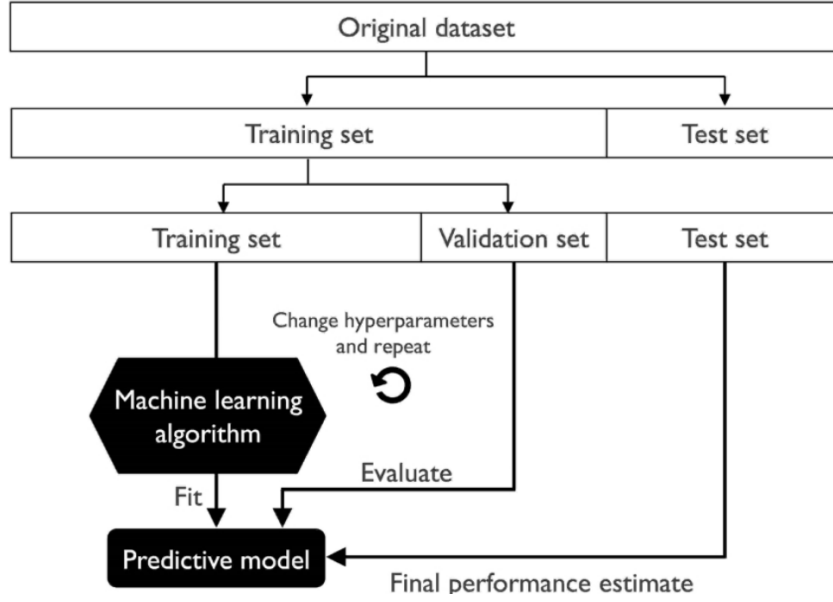


Figure 5.5: Holdout cross-validation as summarized by Raschka and Mirjalili[40].

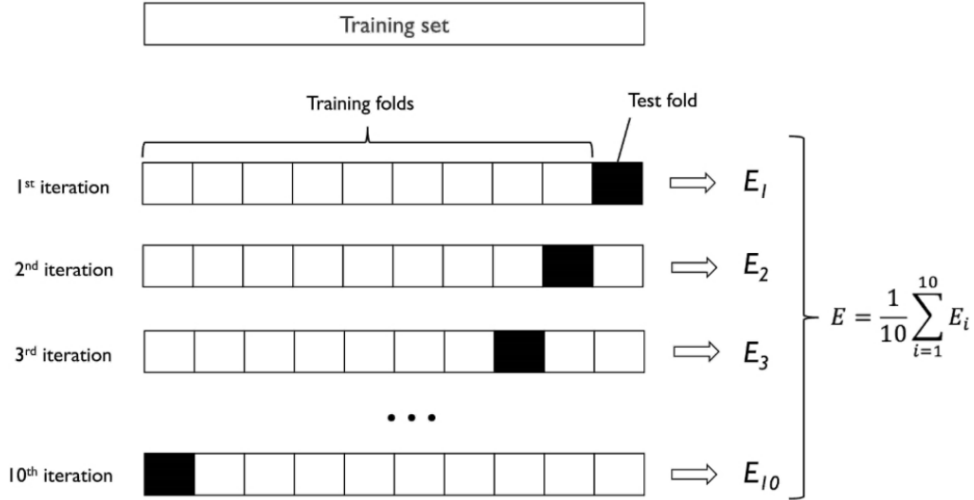


Figure 5.6:  $K$ -fold cross-validation as summarized by Raschka and Mirjalili[40]. In the case of  $K = 10$ , the training dataset is divided into 10 folds, and during the 10 iterations, nine folds are used for training, and one fold is used as the test set for the model evaluation. Estimated performances (for example, classification accuracy or error) for each fold are then used to calculate the estimated average performance  $E$  of the model.

portions, as has been shown in a study by Ron Kohavi (*A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection, International Joint Conference on Artificial Intelligence (IJ-CAI), 14 (12): 1137-43, 1995*). In stratified cross-validation, the class proportions are preserved in each fold to ensure that each fold is representative of the class proportions in the training dataset[22]

For example, each classification algorithm has its inherent biases, and no single classification model enjoys superiority if we don't make any assumptions about the task. In practice, it is therefore essential to compare at least a handful of different algorithms in order to train and select the best performing

model. But before we can compare different models, we first have to decide upon a metric to measure performance. One commonly used metric is classification accuracy, which is defined as the proportion of correctly classified instances.

After we have selected a model that has been fitted on the training dataset, we can use the test dataset to estimate how well it performs on this unseen data to estimate the generalization error. If we are satisfied with its performance, we can now use this model to predict new, future data. It is important to note that the parameters for the previously mentioned procedures, such as feature scaling and dimensionality reduction, are solely obtained from the training dataset, and the same parameters are later reapplied to transform the test dataset, as well as any new data samples—the performance measured on the test data may be overly optimistic otherwise.

randomly partition this dataset into separate test and training datasets is to use the `train_test_split` function from `scikit-learn`'s `model_selection` submodule :

In this context, stratification means that the `train_test_split` method returns training and test subsets that have the same proportions of class labels as the input dataset.

Both the prediction error (ERR) and accuracy (ACC) provide general information about how many samples are misclassified. The error can be understood as the sum of all false predictions divided by the number of total predictions, and the accuracy is calculated as the sum of correct predictions divided by the total number of predictions, respectively:

$$ERR = \frac{\text{Incorrect classifications}}{\text{Total classifications}} \quad (5.4)$$

The prediction accuracy can then be calculated directly from the error:

$$ACC = 1 - ERR = \frac{\text{Correct classifications}}{\text{Total classifications}} \quad (5.5)$$

When dealing with imbalanced data, standard classification metrics do not adequately represent your models performance.

## 5.5 Feature scaling

Raw data rarely comes in the form and shape that is necessary for the optimal performance of a learning algorithm. Thus, the preprocessing of the data is one of the most crucial steps in any machine learning application. Many machine learning algorithms also require that the selected features are on the same scale for optimal performance, which is often achieved by transforming the features in the range  $[0, 1]$  or a standard normal distribution with zero mean and unit variance, as we will see in later chapters.

Indeed many estimators are designed with the assumption that each feature takes values close to zero or more importantly that all features vary on comparable scales. In particular, metric-based and gradient-based estimators often assume approximately standardized data (centered features with unit variances). A notable exception are decision tree-based estimators that are robust to arbitrary scaling of the data.[23]

there are two common approaches to bring different features onto the same scale: normalization and standardization.

In general, learning algorithms benefit from standardization of the data set. If some outliers are present in the set, robust scalers or transformers are more appropriate.

Most often, normalization refers to the rescaling of the features to a range of  $[0, 1]$ , which is a special case of min-max scaling. To normalize our data, we can simply apply the min-max scaling to each feature column, where the new value  $x_{norm}^{(i)}$  of a sample  $x^{(i)}$  can be calculated as follows :

$$x_{norm}^{(i)} = \frac{x^{(i)} - x_{min}}{x_{max} - x_{min}} \quad (5.6)$$

Here,  $x^{(i)}$  is a particular sample,  $x_{min}$  is the smallest value in a feature column, and the largest value. The min-max scaling procedure is implemented in scikit-learn via `MinMaxScaler`.

Although normalization via min-max scaling is a commonly used technique that is useful when we need values in a bounded interval, standardization, or mean removal and variance scaling, can be more practical for many machine learning algorithms, especially for optimization algorithms such as gradient descent. The reason is that many linear models, such as the logistic regression and SVM initialize the weights to 0 or small random values close to 0. Using standardization, we center the feature columns at mean 0 with standard deviation 1 so that the feature columns have the same parameters as a standard normal distribution (zero mean and unit variance), which makes it easier to learn the weights. Furthermore, standardization maintains useful information about outliers and makes the algorithm less sensitive to them in contrast to min-max scaling, which scales the data to a limited range of values.

standardization, which gives our data the property of a standard normal distribution, which helps gradient descent learning to converge more quickly. Standardization shifts the mean of each feature so that it is centered at zero and each feature has a standard deviation of 1. For instance, to standardize the  $j$ th feature, we can simply subtract the sample mean from every training sample and divide it by its standard deviation  $\therefore$ . One of the reasons why standardization helps with gradient descent learning is that the optimizer has to go through fewer steps to find a good or optimal solution (the global cost minimum)

The procedure for standardization can be expressed by the following equation:

$$\mathbf{x}_{std}^{(i)} = \frac{\mathbf{x}^{(i)} - \mu_x}{\sigma_x} \quad (5.7)$$

Here,  $\mathbf{x}$  is a vector consisting of the  $j$ th feature values of all training samples  $n$ , and this standardization technique is applied to each feature  $j$  in our dataset.

we fit the `StandardScaler` and `MinMaxScaler` classes only once—on the training data—and use those parameters to transform the test set or any new data point.

## 5.6 Model selection

A popular class of procedures for solving classification tasks are based on linear models. What this means is that they aim at dividing the feature space into a collection of regions labelled according to the values the target can take, where the decision boundaries between those regions are linear: they are lines in 2D, planes in 3D and hyperplanes with more features.

Trying to understand how the biological brain works, in order to design AI, Warren McCulloch and Walter Pitts published the first concept of a simplified brain cell, the so-called McCulloch-Pitts (MCP) neuron, in 1943 (*A Logical Calculus of the Ideas Immanent in Nervous Activity*, W. S. McCulloch and W. Pitts, *Bulletin of Mathematical Biophysics*, 5(4): 115-133, 1943).

Only a few years later, Frank Rosenblatt published the first concept of the perceptron learning rule based on the MCP neuron model (*The Perceptron: A Perceiving and Recognizing Automaton*, F. Rosen-

blatt, Cornell Aeronautical Laboratory, 1957). With his perceptron rule, Rosenblatt proposed an algorithm that would automatically learn the optimal weight coefficients that are then multiplied with the input features in order to make the decision of whether a neuron fires or not. In the context of supervised learning and classification, such an algorithm could then be used to predict if a sample belongs to one class or the other.

The Perceptron is another simple classification algorithm suitable for large scale learning. By default: It does not require a learning rate. It is not regularized (penalized). It updates its model only on mistakes.

Logistic regression models the probabilities of an observation belonging to each of the  $K$  classes via linear functions, ensuring these probabilities sum up to one and stay in the  $(0, 1)$  range. Logistic regression models are typically estimated by maximum likelihood. Although logistic regression is mostly used as an inference tool in tasks where the goal is to understand the role of input variables in explaining the outcome (it produces easily interpretable coefficients, just like linear regression does), it can also prove to be of significant predictive power.

Linear Discriminant Analysis (LDA) assumes that the joint densities of all features given target's classes are multivariate Gaussians with the same covariance for each class. Quadratic Discriminant Analysis (QDA) relaxes the common covariance assumption of LDA through estimating a separate covariance matrix for each class.

A Naive Bayes Classifier determines the probability that an example belongs to some class, calculating the probability that an event will occur given that some input event has occurred.

let's consider two common ways we'll train a model: tree-based logical rules developed according to some splitting criterion, and parameterized models updated by gradient descent.

Machine learning algorithms can be grouped into parametric and nonparametric models. Using parametric models, we estimate parameters from the training dataset to learn a function that can classify new data points without requiring the original training dataset anymore. Typical examples of parametric models are the perceptron, logistic regression, and the linear SVM. In contrast, nonparametric models can't be characterized by a fixed set of parameters, and the number of parameters grows with the training data. KNN belongs to a subcategory of nonparametric models that is described as instance-based learning. Models based on instance-based learning are characterized by memorizing the training dataset, and lazy learning is a special case of instance-based learning that is associated with no (zero) cost during the learning process.

Machine learning algorithms can be grouped into parametric and nonparametric models. Using parametric models, we estimate parameters from the training dataset to learn a function that can classify new data points without requiring the original training dataset anymore. Typical examples of parametric models are the perceptron, logistic regression, and the linear SVM. In contrast, nonparametric models can't be characterized by a fixed set of parameters, and the number of parameters grows with the training data. Two examples of non-parametric models that we have seen so far are the decision tree classifier/random forest and the kernel SVM.

An important point to be summarized from the famous No Free Lunch Theorems (NFL)[51, 52] by David H. Wolpert is that no single classifier works best across all possible scenarios, as there is a lack of a priori distinctions between learning algorithms.

For example, each classification algorithm has its inherent biases, and no single classification model enjoys superiority if we don't make any assumptions about the task. In practice, it is therefore essential to compare at least a handful of different algorithms in order to train and select the best performing

model. But before we can compare different models, we first have to decide upon a metric to measure performance. One commonly used metric is classification accuracy, which is defined as the proportion of correctly classified instances.

In practice, it is always recommended that you compare the performance of at least a handful of different learning algorithms to select the best model for the particular problem; these may differ in the number of features or samples, the amount of noise in a dataset, and whether the classes are linearly separable or not.

Another simple yet more powerful algorithm for linear and binary classification problems: logistic regression. Note that, in spite of its name, logistic regression is a model for classification, not regression.

Logistic regression is a classification model that is very easy to implement but performs very well on linearly separable classes. It is one of the most widely used algorithms for classification in industry. Similar to the perceptron and Adaline, the logistic regression model in this chapter is also a linear model for binary classification that can be extended to multiclass classification, for example, via the OvR technique.

Decision trees can build complex decision boundaries by dividing the feature space into rectangles. Decision tree classifiers are attractive models if we care about interpretability. As the name decision tree suggests, we can think of this model as breaking down our data by making a decision based on asking a series of questions. Based on the features in our training set, the decision tree model learns a series of questions to infer the class labels of the samples. Using the decision algorithm, we start at the tree root and split the data on the feature that results in the largest Information Gain (IG). In an iterative process, we can then repeat this splitting procedure at each child node until the leaves are pure. This means that the samples at each node all belong to the same class. In practice, this can result in a very deep tree with many nodes, which can easily lead to overfitting. Thus, we typically want to prune the tree by setting a limit for the maximal depth of the tree. As we can see, the information gain is simply the difference between the impurity of the parent node and the sum of the child node impurities — the lower the impurity of the child nodes, the larger the information gain. Now, the three impurity measures or splitting criteria that are commonly used in binary decision trees are Gini impurity ( $I_G$ ), entropy ( $I_H$ ), and the classification error ( $I_E$ ).

Random forests have gained huge popularity in applications of machine learning during the last decade due to their good classification performance, scalability, and ease of use. Intuitively, a random forest can be considered as an ensemble of decision trees. The idea behind a random forest is to average multiple (deep) decision trees that individually suffer from high variance, to build a more robust model that has a better generalization performance and is less susceptible to overfitting.

1. Draw a random bootstrap sample of size  $n$  (randomly choose  $n$  samples from the training set with replacement).
2. Grow a decision tree from the bootstrap sample. At each node:
  - a. Randomly select  $d$  features without replacement.
  - b. Split the node using the feature that provides the best split according to the objective function, for instance, maximizing the information gain.
3. Repeat the steps 1-2  $k$  times.
4. Aggregate the prediction by each tree to assign the class label by majority vote.

In most implementations, including the `RandomForestClassifier` implementation in `scikit-learn`, the size of the bootstrap sample is chosen to be equal to the number of samples in the original training set,



which usually provides a good bias-variance tradeoff. For the number of features  $d$  at each split, we want to choose a value that is smaller than the total number of features in the training set. A reasonable default that is used in *scikit-learn* and other implementations is  $d = \sqrt{m}$ , where  $m$  is the number of features in the training set.

Although we are growing a very small random forest from a very small training dataset, we used the `n_jobs` parameter for demonstration purposes, which allows us to parallelize the model training using multiple cores of our computer (here two cores).

*k*-nearest neighbor (KNN) classifier is particularly interesting because it is fundamentally different from the learning algorithms that we have discussed so far. KNN is a typical example of a lazy learner. It is called lazy not because of its apparent simplicity, but because it doesn't learn a discriminative function from the training data, but memorizes the training dataset instead. KNN belongs to a subcategory of nonparametric models that is described as instance-based learning. Models based on instance-based learning are characterized by memorizing the training dataset, and lazy learning is a special case of instance-based learning that is associated with no (zero) cost during the learning process.

The KNN algorithm itself is fairly straightforward and can be summarized by the following steps:

1. Choose the number of  $k$  and a distance metric.
2. Find the  $k$ -nearest neighbors of the sample that we want to classify.
3. Assign the class label by majority vote.

The principle behind nearest neighbor methods is to find a predefined number of training samples closest in distance to the new point, and predict the label from these. The number of samples can be a user-defined constant (*k*-nearest neighbor learning), or vary based on the local density of points (radius-based neighbor learning). The distance can, in general, be any metric measure: standard Euclidean distance is the most common choice. Neighbors-based methods are known as non-generalizing machine learning methods, since they simply “remember” all of its training data (possibly transformed into a fast indexing structure such as a Ball Tree or KD Tree). Being a non-parametric method, it is often successful in classification situations where the decision boundary is very irregular.[26]

Neighbors-based classification is a type of instance-based learning or non-generalizing learning: it does not attempt to construct a general internal model, but simply stores instances of the training data. Classification is computed from a simple majority vote of the nearest neighbors of each point: a query point is assigned the data class which has the most representatives within the nearest neighbors of the point.

## 5.7 Chapter summary

*asd*

## Chapter 6

# Model evaluation

### 6.1 Metrics for evaluating model performance

### 6.2 Model performance

*Model hyperparameters were tuned for each classifier using k-fold cross validation, as was described in section 5.4. After the tuning, performance of the following models have been compared:*

- *Perceptron learning algorithm, learning rate  $\eta = 0.5$ , maximum iterations=5, features transformed with quantile transformation (uniform PDF)*  
*model code: ppn\_qu\_eta0.5\_maxiter5*
- *Logistic regression, L2 regularization, regularization parameter  $C = 0.1$ , maximum iterations=100, features transformed with quantile transformation (uniform PDF)*  
*model code: lr\_l2\_c0.1\_maxiter\_100*
- *Logistic regression, L1 regularization, regularization parameter  $C = 0.1$ , maximum iterations=100, features transformed with quantile transformation (uniform PDF)*  
*model code: lr\_l1\_c0.1\_maxiter\_100*
- *Linear Discriminant Analysis, features transformed with quantile transformation (uniform PDF)*  
*model code: lda*
- *Quadratic Discriminant Analysis, features transformed with quantile transformation (uniform PDF)*  
*model code: qda*
- *Linear Support Vector Classification, L2 regularization, regularization parameter  $C = 0.1$ , maximum iterations=100, features transformed with quantile transformation (uniform PDF)*  
*model code: lsvc\_l2\_c0.1\_maxiter\_100*
- *Linear Support Vector Classification, L1 regularization, regularization parameter  $C = 0.1$ , maximum iterations=100, features transformed with quantile transformation (uniform PDF)*  
*model code: lsvc\_l1\_c0.1\_maxiter\_100*

- *Gaussian Naive Bayes*

*model code: nb*

- *Decision Tree, Gini impurity criterion, maximum depth of the tree=25, unscaled features*

*model code: tree25*

- *Random Forest, Gini impurity criterion, number of estimators=50, unscaled features*

*model code: forest50*

- *K-Nearest Neighbors, Manhattan distance metric, number of neighbors=4, standardized features*

*model code: knn\_p1\_k4*

Four different subsets of data have been used to test the model: train, test, and two additional validation subsets. The train and test subsets represent Teranet records from 2011 to 2014 randomly sampled into 70% train and 30% test subsets; these are the primary subsets that were used for training and tuning the hyperparameters and then evaluating the performance of classifiers on unseen test data. The two additional validation subsets are composed of Teranet records from 2010 and 2015. Since the land use information from Department of Geography was collected in 2012 and 2013, it can be less accurate for these subsets; these subsets have not been used for model training or selection, but are presented as an additional reference for the generalization of performance of classifiers. Figure 6.1 presents model performance on train, test, and two additional validation subsets.

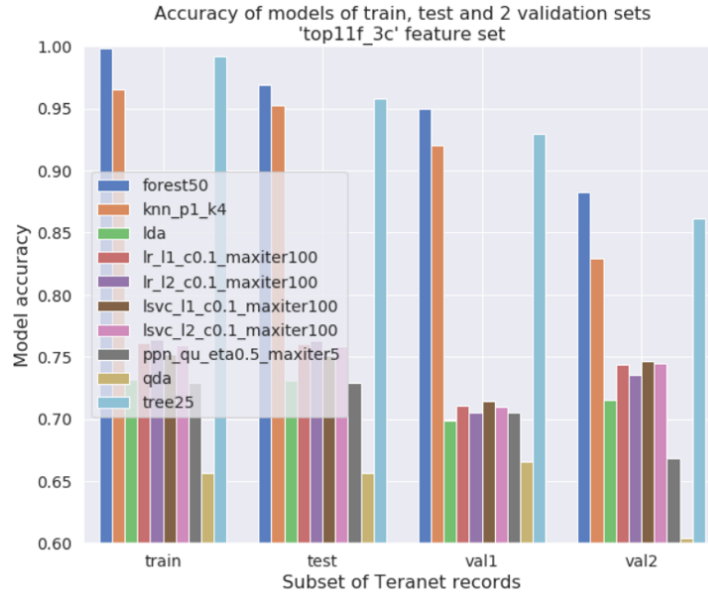


Figure 6.1: Model performance (accuracy) on train, test, and two additional validation subsets. Additional validation subsets are composed of Teranet records from 2010 and 2015; since land use information (target variable) can be less accurate for records in these subsets, they have been used as an additional reference for evaluating the generalization of the classifier performance.

*Different feature scaling techniques have a strong effect on model fit times and predictive performance of linear models, as can be seen on figure 6.2.*

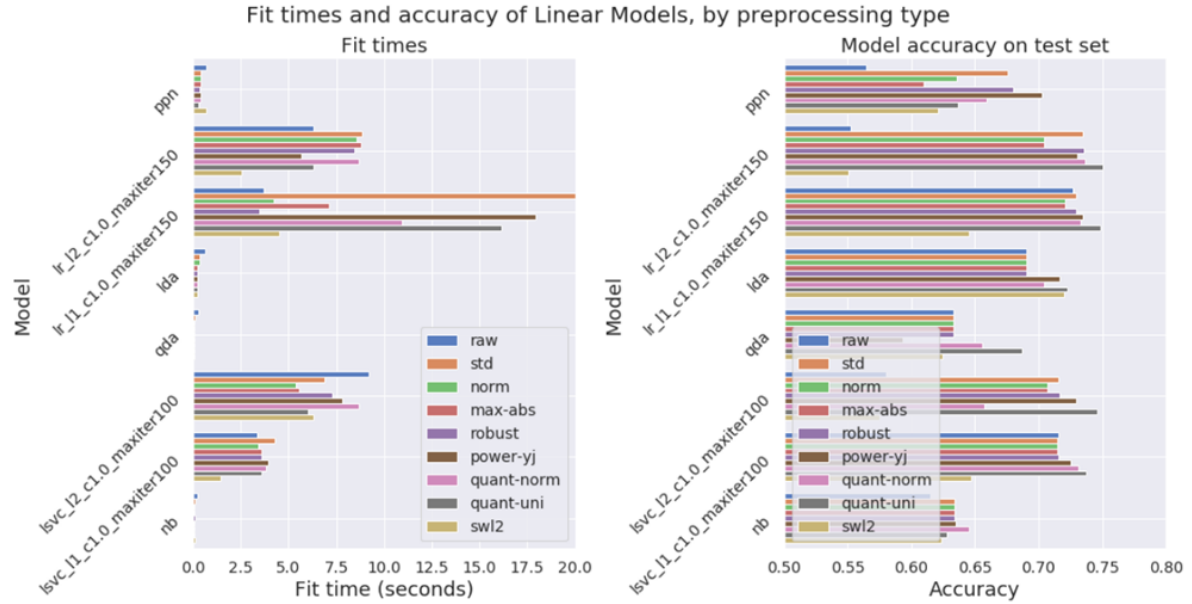


Figure 6.2: Fit times and accuracy of linear models, by feature scaling technique. Different scaling techniques have a strong effect on the performance of linear models.

Similar to linear models, distance-based algorithms, such as *K-Nearest Neighbors*, also are strongly affected by feature scaling. In contrast, tree-based models, such as *Decision Tree* and *Random Forest*, are scale-invariant, and thus have stable performance with most feature scaling techniques. Figure 6.3 presents fit times for *Decision Tree*, *Random Forest*, and *K-Nearest Neighbors* with Manhattan and Euclidean distance metrics.

### 6.3 Best performing model: Random Forest

Random Forest as can be seen from the plots presented in section 6.2, Random Forest with 50 estimators and Gini impurity criterion showed the best results in terms of accuracy and fit times on all subsets. Figure 6.4 presents the classification report showing all main model performance metrics for the best performing model: Random Forest with 50 estimators using 'gini' criterion.

Figure 6.5 presents confusion matrices with and without normalization for Random Forest on the test set.

Figure 6.6 presents feature importance for the best performing model.

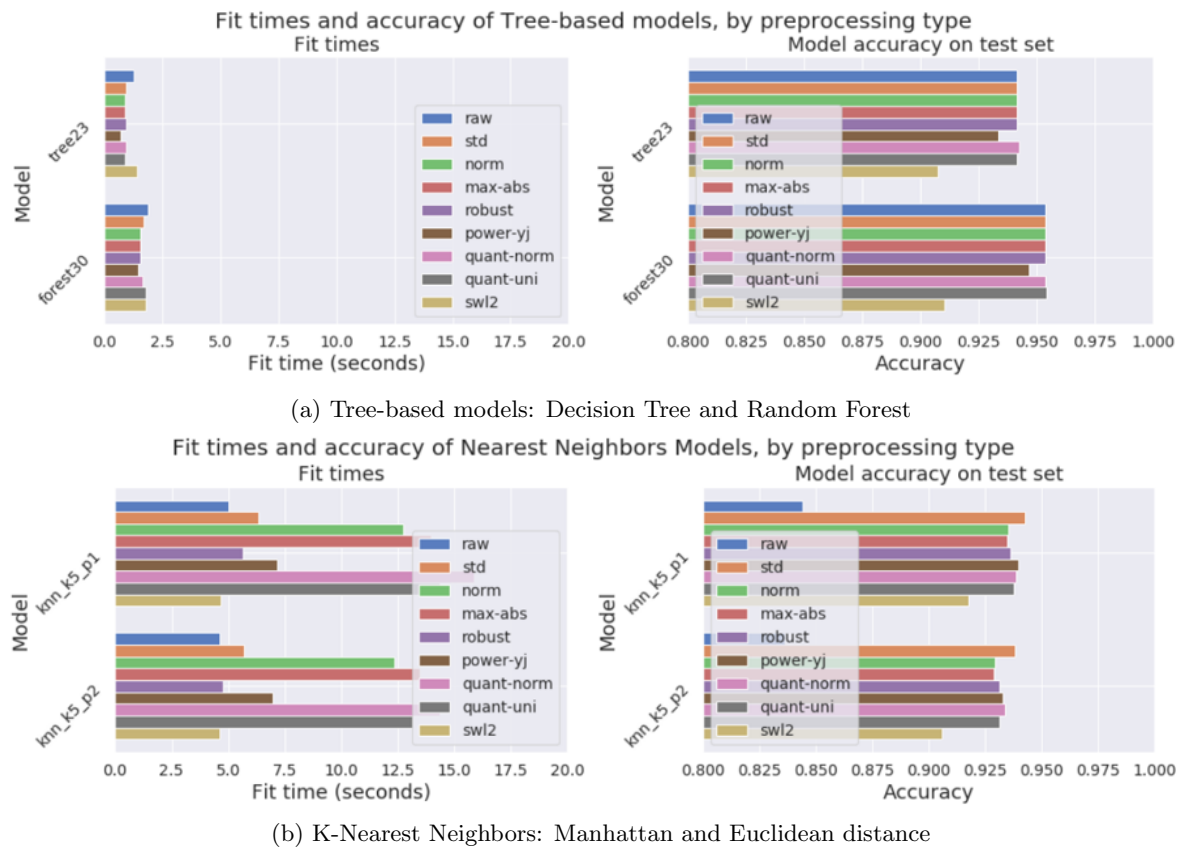


Figure 6.3: Fit times and accuracy for tree-based models and K-Nearest Neighbors, by feature scaling technique. Distance-based algorithms, such as K-Nearest Neighbors are affected by feature scaling, while tree-based models are scale-invariant.

	condo	house	other	accuracy	macro avg	weighted avg
precision	0.996800	0.940099	0.975766	0.968755	0.970888	0.969202
recall	0.992179	0.973011	0.949980	0.968755	0.971723	0.968755
f1-score	0.994484	0.956272	0.962700	0.968755	0.971152	0.968812
support	66871.000000	86664.000000	103079.000000	0.968755	256614.000000	256614.000000

Figure 6.4: Best model performance on test set: classification report for Random Forest with 50 estimators using 'gini' criterion.

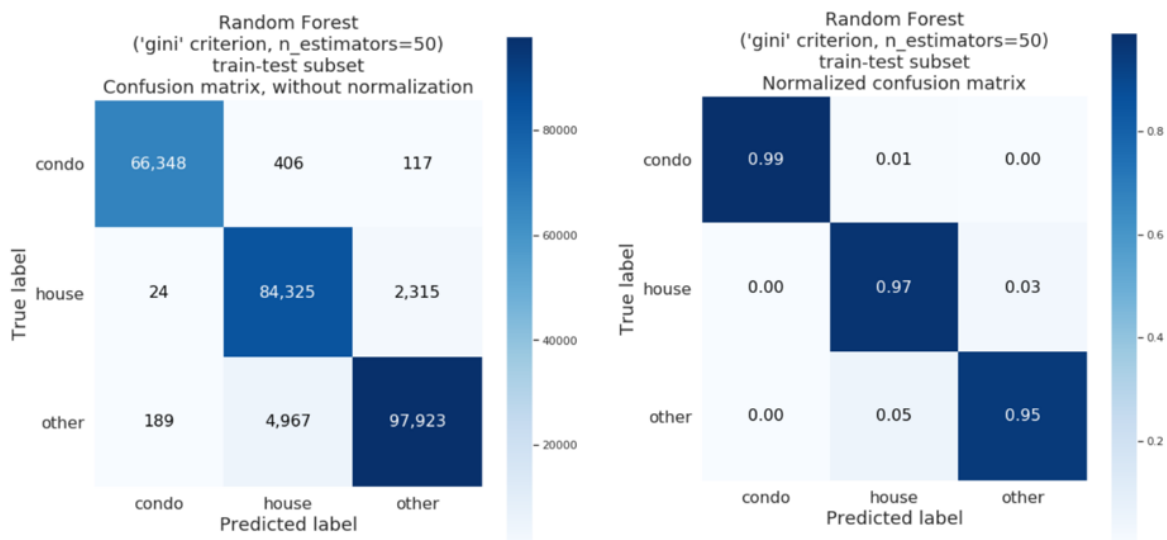


Figure 6.5: Best model performance on test set: confusion matrices with and without normalization for Random Forest with 50 estimators using 'gini' criterion.

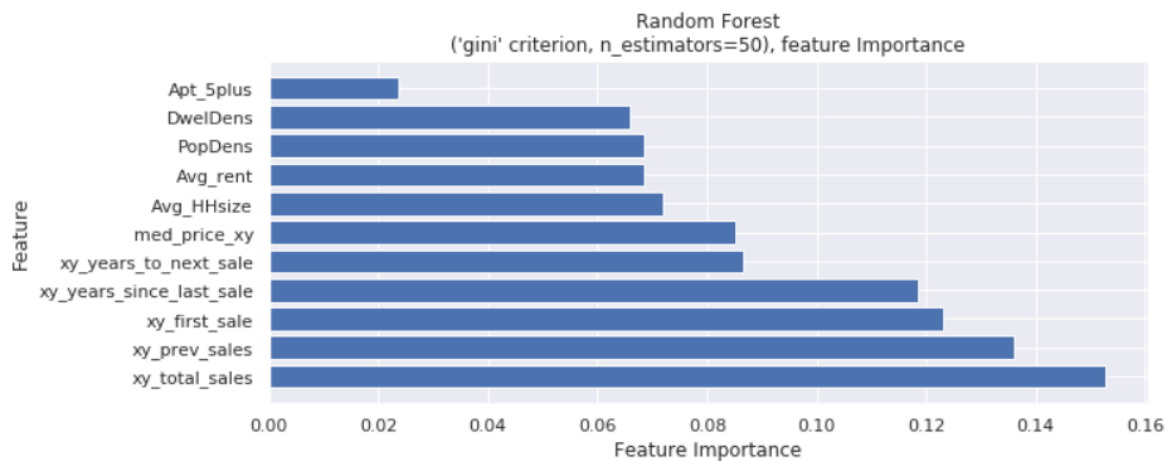


Figure 6.6: Random Forest with 50 estimators using 'gini' criterion: feature importance.

## Chapter 7

# Conculsion

*One of the major attributes missing from Teranet data is the type of property being transacted, or land use information for the parcel where a transaction is recorded. Detailed parcel-level land use from the Department of Geography and DMTI land use data have been spatially joined to each Teranet record. However, since both of these data sources have their limitations, detailed land use data from Department of Geography has been used to train an algorithm capable of classifying land from the housing market dynamics.*

*To augment Teranet's dataset, new variables were engineered from its native attributes to capture the housing market dynamics at the parcel level: for example, 'xy\_total\_sales' was computed as the total count of Teranet records coming from a particular coordinate pair; 'med\_price\_xy' represents the median price of all records coming from a coordinate pair, etc. To augment Teranet data with demographic and transport information, the new Teranet features were spatially and temporally joined with Census and TTS variables recorded at the level of a Dissemination Area and TAZ zone, respectively. Finally, the augmented Teranet dataset has been tested with machine learning algorithms, attempting to classify land use for each Teranet record, thus recognizing land use changes with time.*

*The new Teranet features capturing housing market dynamics have show to have strong predictive power when classifying land use. When joined with Census variables at the level of Dissemination Areas, new features engineered from Teranet's dataset allowed the classification of land use with a high level of accuracy. Random Forest model that was trained using random 70% sample of all Teranet records with new features from 2011 to 2014 stratified by target classes ("condo", "house", or "other") achieved 97% of accuracy on the test subset composed of the remaining 30% of records from 2011 and 2014.*

# Bibliography

- [1] *F. J. Anscombe. Graphs in Statistical Analysis. The American Statistician, 27(1):17–21, 1973.*
- [2] *Daniel Arribas-Bel. Accidental, open and everywhere: Emerging data sources for the understanding of cities. Applied Geography, 49:45–53, 2014.*
- [3] *Bess Ashby. TTS 2016 City of Toronto Summary by Ward. Technical report, malatest, Toronto, 2018.*
- [4] *Bank of Canada. Inflation Calculator, 2019.*
- [5] *Michael Batty. Cities as Complex Systems: Scaling, Interactions, Networks, Dynamics and Urban Morphologies. 2008.*
- [6] *Richard Bellman. The Theory of Dynamic Programming. Technical report, The RAND Corporation, 1954.*
- [7] *Irad Ben-Gal. Outlier Detection. In Data Mining and Knowledge Discovery Handbook: A Complete Guide for Practitioners and Researchers, chapter Chapter 1, pages 1–16. Kluwer Academic Publishers, Tel-Aviv, 2005.*
- [8] *Luís M.A. Bettencourt. The origins of scaling in cities. Science, 340(6139):1438–1441, 2013.*
- [9] *Richard J Bouchard and Clyde E Pyers. Use of gravity model for describing urban travel: An ayalysis and critique. Highway Research Record, (88):1–43, 1965.*
- [10] *Jason Brownlee. How to Prepare Data For Machine Learning, 2013.*
- [11] *Cynthia Chen, Jingtao Ma, Yusak Susilo, Yu Liu, and Menglin Wang. The promises of big data and small data for travel behavior (aka human mobility) analysis. Transportation Research Part C: Emerging Technologies, 68:285–299, 2016.*
- [12] *Edgar F Codd. The relational model for database management : version 2. Addison-Wesley Longman Publishing Co., Boston, MA, 1990.*
- [13] *Kate Crowley and Brian W. Head. The enduring challenge of ‘wicked problems’: revisiting Rittel and Webber. Policy Sciences, 50(4):539–547, 2017.*
- [14] *Data Management Group. Data Management Group at the University of Toronto Transportation Research Institute, 2014.*
- [15] *Data Management Group. Survey Boundary Files, 2019.*



- [16] DMTI Spatial Inc. *CanMap ® RouteLogistics User Manual v2014.2*, 2014.
- [17] Douglas M Hawkins. Identification of outliers. *Chapman and Hall, London, UK*, 1980.
- [18] Michael Iacono, David Levinson, and Ahmed El-Geneidy. *Models of transportation and land use change: A guide to the territory*. Journal of Planning Literature, 2008.
- [19] Jane Jacobs. *The Death and Life of Great American Cities*. In New York, volume 71, pages Alexander, C., Ishikawa, S., & Silverstein, M. (19. 1961.
- [20] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. An Introduction to Statistical Learning. *Springer Texts in Statistics*. Springer New York, New York, NY, 2013.
- [21] Eric Damian Kelly. *The Transportation Land Use Link*. Journal of Planning Literature, 9(2):p.128–145, 1994.
- [22] Ron Kohavi. *A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection*. In International Joint Conference on Artificial Intelligence (IJCAI), 1995.
- [23] Scikit learn developers. *Compare the effect of different scalers on data with outliers*, 2019.
- [24] Scikit learn developers. *Cross-validation: evaluating estimator performance*, 2019.
- [25] Scikit learn developers. *L1 Penalty and Sparsity in Logistic Regression*, 2019.
- [26] Scikit learn developers. *Nearest Neighbors*, 2019.
- [27] Ira S Lowry. *A Model of Metropolis*. Technical report, Rand Corporation, Santa Monica CA, 1964.
- [28] Map and Data Library. *Canadian census geography (unit) definitions*, 2019.
- [29] F J Martinez. *The bid-choice land-use model: an integrated economic framework*. Environment and Planning, 24(January 1991):871–885, 1992.
- [30] Heather McKean and Stella Di Cresce. *The International Comparative Legal Guide to: Real Estate 2015, Chapter 6: Canada*. Technical report, Global Legal Group, London, UK, 2015.
- [31] Eric J Miller. *Integrated urban modeling: Past, present, and future*. Journal of Transport and Land Use, 11(1):387–399, 2018.
- [32] Eric J Miller. *The case for microsimulation frameworks for integrated urban models*. Journal of Transport and Land Use, 11(1):1025–1037, 2018.
- [33] Eric J Miller. *Towards the Next Generation of Integrated Urban Models*, 2018.
- [34] Eric J Miller, Bilal Farooq, Franco Chingcuanco, and David Wang. *Historical validation of integrated transport-land use model system*. Transportation Research Record, (2255):91–99, 2011.
- [35] Eric J Miller, David S Kriger, and John Douglas Hunt. *Integrated Urban Models for Simulation of Transit and Land Use Policies Guidelines for Implementation and Use*. 1998.
- [36] Rolf Moeckel. *Constraints in household relocation: Modeling land-use/transport interactions that respect time and monetary budgets*. Journal of Transport and Land Use, 10(1):211–228, 2017.

- [37] Robert Nisbet, Gary Miner, and Ken Yale. *A Data Preparation Cookbook*. In *Handbook of Statistical Analysis and Data Mining Applications*, chapter *Chapter 18*, pages 727–740. Elsevier, 2018.
- [38] Guido Noto, Federico Cosenz, and Carmine Bianchi. Urban Transportation Governance And Wicked Problems: A Systemic And Performance Oriented Approach. *Phd thesis, University of Palermo*, 2015.
- [39] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. *Scikit-learn: Machine Learning in {P}ython*. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [40] Sebastian Raschka and Vahid Mirjalili. *Python Machine Learning*, 2nd Ed. Packt Publishing, Birmingham, UK, 2 edition, 2017.
- [41] Horst W.J. Rittel and Melvin M. Webber. *Dilemmas in a General Theory of Planning*. *Policy Sciences*, 4(2):155–169, 1973.
- [42] Adam Rosenfield, Franco Chingcuanco, and Eric J Miller. *Agent-based housing market microsimulation for integrated land use, transportation, environment model system*. *Procedia Computer Science*, 19:841–846, 2013.
- [43] Colin Shearer. *The CRISP-DM Model: The New Blueprint for Data Mining*. *Journal of Data Warehousing*, 5(4):13–22, 2000.
- [44] Statistics Canada. *Dissemination area (DA)*, 2015.
- [45] Statistics Canada. *Hierarchy of standard geographic units*, 2018.
- [46] Vergil G Stover and Frank J Koepke. *Transportation and Land Development*. *Pearson College Div*, 1988.
- [47] Teranet Enterprises Inc. *About POLARIS*, Teranet, 2019.
- [48] Teranet Enterprises Inc. <https://www.teranet.ca>, 2019.
- [49] The Government of Ontario. *Land Registration Reform Act, R.S.O. 1990, c. L.4*, 1990.
- [50] Hadley Wickham. *Tidy Data*. *Journal of Statistical Software*, 59(10), 2014.
- [51] David H Wolpert. *The Lack of A Priori Distinctions Between Learning Algorithms*. *Neural Computation*, 8(7):1391–1420, 1996.
- [52] David H Wolpert and William G Macready. *No Free Lunch Theorems for Optimization*. *Technical Report 1*, 1997.