

CM3070FinalProjectPreliminaryReport

June 17, 2024

1 Deep Learning for Personalised Property Recommendations System: Data Collection and Model Development Using Public Datasets

CM3070 Final Project Preliminary Report: Guillermo Olmos Ranalli

1.1 1. Introduction

Word Count: 837

1.1.1 1.1 Project Concept

The Personalised Property Recommendation System aims to assist potential homebuyers in the UK real estate market by providing tailored property recommendations based on individual preferences and financial situations. By integrating historical transaction data from HM Land Registry with current property listings from OnTheMarket, the system delivers customised property suggestions. This project involves developing and evaluating various deep learning models to determine the most effective approach for property recommendation, framed as a classification task.

1.1.2 1.2 Motivation

Navigating the real estate market can be particularly challenging for new generations, including Millennials, Gen Z, and Gen Alpha, who face unique financial and lifestyle constraints. First-time buyers often struggle to find properties that meet their specific criteria within their budget. The abundance of property options and the complexity of property features necessitate a tool that can streamline the search process and offer personalised recommendations. This project seeks to address this need by developing a robust deep learning-based recommendation system.

1.1.3 1.3 Project Template

This project, based on “Project Idea Title 1: Deep Learning on a Public Dataset” from CM3015, uses publicly accessible real estate data from HM Land Registry and OnTheMarket to develop and compare deep learning models for personalised property recommendations. Following the methodology in “Deep Learning with Python” by F. Chollet, the prototype aims to improve model performance.

This project builds upon my previous work in CM2015, where I developed data collection scripts and utilised similar datasets. The current project expands this foundation by incorporating advanced machine learning for better recommendations.

1.1.4 1.4 Scope and Limitations

This prototype analyses various types of residential properties in Buckinghamshire, excluding shared ownership, retirement homes, new builds, auction listings, farms/land, park homes, and properties over £650,000. These filters were applied to data collected from OnTheMarket.com.

The prototype has limitations, including reliance on online data accuracy and availability. It does not cover commercial real estate or the rental market in Buckinghamshire.

1.1.5 1.5 Data Sources and Selection

1.5.1 Data Requirements The primary goal of this prototype is to validate the feasibility of using machine learning techniques to predict property prices and recommend properties in the Buckinghamshire market. To achieve this, two main data sources were utilised: - **HM Land Registry Data:** This dataset provides comprehensive information on property sales in Buckinghamshire, including sale prices and transaction dates. - **Web Scraped Data from OnTheMarket.com:** A dataset comprising current property listings in Buckinghamshire, including asking prices, property types, and other relevant details.

To comprehensively analyse Buckinghamshire's property market, data from diverse sources were gathered: - **HM Land Registry Data:** This dataset provides comprehensive information on property sales in Buckinghamshire, including sale prices and transaction dates. - **Real Estate Listings:** Scraped data from OnTheMarket.com helps in understanding the ongoing trends and fluctuations in property prices and demands within Buckinghamshire. Python scripts were developed, organised in the `src` folder of this prototype, to efficiently collect and process the data, ensuring a robust and reliable dataset for analysis.

1.5.2 Choice of Data Sources

- **HM Land Registry Data:** Chosen for its reliability and comprehensive coverage of actual property sales in the UK. Accessed via gov.uk, under the Open Government Licence v3.0.
- **OnTheMarket.com:** Chosen as a current and active source for property listings, offering a real-time perspective on the market. Data was scraped in compliance with the site's terms and conditions, focusing on properties listed for sale in Buckinghamshire.

1.5.3 Methodology for Data Collection and Processing

- Data from HM Land Registry was downloaded in CSV format, covering transactions for the year 2023.
- Web scraping was conducted on OnTheMarket.com using Python scripts, focusing on gathering current listings in Buckinghamshire. The scraping process adhered to the website's robots.txt file and was conducted using a unique user agent.

1.5.4 Limitations and Constraints

- **Timeframe:** The HM Land Registry data covers only sales within 2023, and the scraped data reflects listings at the time of scraping. This temporal limitation means the analysis might not fully capture long-term market trends.
- **Geographical Scope:** The focus on Buckinghamshire alone may not provide a complete picture of broader regional or national property market trends.

- **Data Completeness:** While the Land Registry data is comprehensive for sales, the scraped data from OnTheMarket.com might not capture every property listing in the region, leading to potential gaps in the dataset.

1.1.6 1.6 Ethical Considerations

1.6.1 Data Sources and Permissions **HM Land Registry Data** - HM Land Registry data is used under the Open Government Licence v3.0. - Proper attribution has been given as per the OGL requirements: “Contains HM Land Registry data © Crown copyright and database right 2021. This data is licensed under the Open Government Licence v3.0.”

OnTheMarket.com Data - OnTheMarket.com data was collected via web scraping, strictly adhering to their robots.txt file, using a unique user agent with contact information, and employing rate limiting to respect their servers.

This prototype focuses on objective real estate data, avoiding personal judgements or assumptions. It aims to maintain neutrality, preventing negative impacts like market manipulation. Any personal data has been anonymized to protect privacy.

2 Chapter 2: A Literature Review

Word count chapter 2: 787

2.1 2.1 Introduction

In the realm of real estate, accurate property price prediction and personalised recommendation systems are crucial for assisting potential homebuyers and real estate professionals. The intersection of machine learning and real estate has gained substantial attention, with various methodologies explored to enhance prediction accuracy and personalisation. This literature review critically examines existing research on housing price prediction models and personalised recommendation systems, highlighting their methodologies, strengths, and limitations.

2.2 2.2 Housing Price Prediction Models

2.2.1 2.2.1 Hedonic-Based Regression Approaches

Historically, hedonic-based regression models have been utilised to determine the impact of various housing attributes on property prices. These models estimate prices based on factors such as location, size, and age of the property. Despite their widespread use, hedonic models face limitations such as difficulty in capturing nonlinear relationships and the need for extensive data preprocessing to handle heteroscedasticity and multicollinearity issues [1].

2.2.2 2.2.2 Machine Learning Techniques

The advent of machine learning has provided more sophisticated tools for housing price prediction, capable of handling complex, non-linear relationships between variables. The study by Park and Bae (2020) investigates the application of several machine learning algorithms to predict housing prices using data from Fairfax County, Virginia. The algorithms examined include C4.5, RIPPER, Naïve Bayesian, and AdaBoost [5].

- **C4.5 Algorithm:** This algorithm is an extension of the earlier ID3 algorithm and generates a decision tree used for classification purposes. In the context of housing price prediction, the decision tree helps identify the most significant variables influencing prices [6].
- **RIPPER Algorithm:** This is a rule-based learning algorithm that generates a set of rules to classify data. According to Park and Bae (2020), the RIPPER algorithm demonstrated superior performance in terms of accuracy compared to other models tested in their study [5].
- **Naïve Bayesian:** This probabilistic classifier is based on Bayes' theorem and assumes independence between predictors. Despite its simplicity, it can be effective for certain types of classification problems [2].
- **AdaBoost:** This ensemble method combines multiple weak classifiers to create a strong classifier. It adjusts the weights of misclassified instances, thereby improving the model's accuracy over successive iterations [3].

Park and Bae's study concludes that the RIPPER algorithm consistently outperformed the other models in terms of classification accuracy for housing price prediction. This finding is significant as it highlights the potential of rule-based algorithms in capturing the complexities of housing market data [5].

2.3 2.3 Content-Based Recommender Systems

Content-based recommender systems are crucial for providing personalised suggestions based on user preferences and item attributes. Lops, Gemmis, and Semeraro (2011) provide a comprehensive overview of the state-of-the-art techniques and trends in content-based recommendation systems [4].

- **Feature Extraction:** Content-based systems rely heavily on extracting meaningful features from items. In the context of real estate, features such as property type, location, price, and amenities are essential.
- **Similarity Calculation:** These systems calculate the similarity between items based on their features. For real estate, properties with similar attributes (e.g., location, price range) are considered similar and thus recommended to users with matching preferences.
- **User Profiles:** Content-based systems maintain profiles for users, capturing their preferences and interaction history. This allows the system to tailor recommendations based on individual user needs.

Lops et al. (2011) highlight the challenges in content-based recommender systems, such as the cold start problem, where new users or items lack sufficient data for effective recommendations. However, integrating advanced machine learning techniques can mitigate some of these issues by improving feature extraction and similarity calculations [4].

2.4 2.4 Discussion

The research conducted by Park and Bae (2020) underscores the importance of selecting appropriate machine learning algorithms for housing price prediction. Their comparative analysis provides valuable insights into the strengths and weaknesses of different models. For instance, while ensemble

methods like AdaBoost are generally robust, rule-based algorithms such as RIPPER can offer higher accuracy for specific datasets [5].

This study also emphasizes the need for comprehensive data preprocessing, including the selection of relevant features and handling missing values, to enhance the predictive performance of machine learning models. Additionally, the integration of various algorithms can potentially lead to the development of a hybrid model that leverages the strengths of each approach [5].

2.5 2.5 Conclusion

The literature on housing price prediction demonstrates that machine learning techniques, particularly rule-based algorithms like RIPPER, can significantly improve the accuracy of price predictions. The study by Park and Bae (2020) serves as a critical reference point for developing advanced models that can aid real estate stakeholders in making informed decisions [5].

Further research should focus on integrating these models with personalised recommendation systems to provide comprehensive solutions for real estate buyers and sellers. By leveraging machine learning's capabilities, the real estate industry can enhance its analytical tools, leading to more accurate and reliable property valuations.

2.6 2.6 References

- [1]: Bin, O. (2005). A semiparametric hedonic model for valuing wetlands. *Applied Economics Letters*, 12(10), 597–601. <https://doi.org/10.1080/13504850500188505>
- [2]: Domingos, P., & Pazzani, M. (1997). On the optimality of the simple Bayesian classifier under zero-one loss. *Machine Learning*, 29(2-3), 103-130. <https://link.springer.com/article/10.1023/A:1007413511361>
- [3]: Freund, Y., & Schapire, R. E. (1997). A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1), 119-139. <https://www.sciencedirect.com/science/article/pii/S002200009791504X>
- [4]: Lops, P., Gemmis, M. D., & Semeraro, G. (2011). Content-based recommender systems: State of the art and trends. *Recommender Systems Handbook*, 73-105. https://www.researchgate.net/publication/226098747_Content-based_Recommender_Systems_State_of_the_Art_and_Trends
- [5]: Park, B., & Bae, J. K. (2020). Using machine learning algorithms for housing price prediction: The case of Fairfax County, Virginia housing data. *Expert Systems with Applications*, 42(6), 2928-2934. <https://doi.org/10.1016/j.eswa.2014.11.040>
- [6]: Quinlan, J. R. (1993). *C4.5: Programs for Machine Learning*. Morgan Kaufmann. <https://link.springer.com/article/10.1007/BF00993309>

3 Chapter 3: A Design

Word Count: 1075

3.1 3.1 Project Overview

The Personalized Property Recommendation System aims to integrate historical transaction data from HM Land Registry and current property listings from OnTheMarket to provide tailored property recommendations based on user preferences and financial situations. The project follows “Project Idea Title 1: Deep Learning on a Public Dataset” and aims to find the most effective model for property price prediction and recommendation using deep learning techniques.

3.2 3.2 Domain and Users

3.2.1 Domain

The project is situated in the real estate domain, focusing on residential properties in the UK. It leverages publicly available datasets to build a recommendation system that aids potential homebuyers in making informed decisions.

3.2.2 Users

The primary users of the system are: - **First-time homebuyers:** Individuals looking for their first property purchase who need tailored recommendations based on their budget and preferences. - **Real estate agents:** Professionals who can use the system to provide clients with data-driven property suggestions. - **Property investors:** Individuals or companies looking to invest in real estate who require accurate property price predictions and recommendations.

3.3 3.3 Justification of Design Choices

3.3.1 User Needs

The design choices are informed by the needs of users in the real estate market: - **Personalization:** Users require personalized property recommendations that match their financial constraints and preferences. - **Accurate Predictions:** Accurate property price predictions help users make informed decisions. - **Usability:** The system must be easy to use and provide quick, relevant recommendations.

3.3.2 Domain Requirements

The real estate domain requires: - **Integration of Diverse Data Sources:** Combining historical transaction data with current listings to provide a comprehensive view. - **Handling Non-linear Relationships:** Using advanced machine learning models to capture complex patterns in the data.

3.4 3.4 Project Structure

3.4.1 Data Collection and Preprocessing

- **Data Sources:** Historical transaction data from HM Land Registry and current property listings from OnTheMarket.
- **Web Scraping:** Scripts to collect real-time data from OnTheMarket.
- **Data Cleaning:** Handling missing values, standardizing formats, and filtering relevant data.
- **Data Integration:** Merging datasets to create a unified data source.

3.4.2 Model Development

- **Feature Selection:** Identifying relevant features such as price, location, property type, etc.
- **Model Selection:** Experimenting with various machine learning algorithms (e.g., C4.5, RIPPER, Naïve Bayesian, AdaBoost) to identify the best-performing model.
- **Training and Validation:** Splitting the data into training and validation sets, training the model, and evaluating its performance.

3.4.3 Recommendation System

- **User Profile Creation:** Collecting user preferences and financial information.
- **Similarity Calculation:** Using content-based filtering to match properties with user profiles.
- **Property Ranking:** Ranking properties based on their relevance to the user's preferences and budget.
- **Feedback Loop:** Incorporating synthetic user feedback to continuously improve the recommendation system.

3.5 3.5 Technologies and Methods

3.5.1 Technologies

- **Python:** Primary programming language for data processing and model development.
- **TensorFlow and Keras:** Libraries for building and training deep learning models.
- **Pandas and NumPy:** Libraries for data manipulation and analysis.
- **Scikit-Learn:** Library for implementing various machine learning algorithms and evaluation metrics.
- **BeautifulSoup and Requests:** Libraries for web scraping.
- **Matplotlib:** Library for data visualization.
- **Node.js and Express.js:** For building the backend of the web application.
- **HTML, CSS, and JavaScript:** For developing the frontend of the web application.

3.5.2 Methods

- **Data Preprocessing:** Cleaning and integrating data from multiple sources.
- **Machine Learning:** Developing and comparing different machine learning models to identify the most effective one for price prediction.
- **Content-Based Filtering:** Creating a recommendation system based on the features of the properties and user preferences.
- **Evaluation Metrics:** Using Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) to evaluate model performance.

3.6 3.6 Work Plan

3.6.1 Major Tasks and Timeline

- **Data Collection and Preprocessing (Weeks 1-4)**
 - Collect data from HM Land Registry and OnTheMarket
 - Clean and preprocess data
 - Integrate datasets
- **Model Development (Weeks 5-10)**

- Feature selection
- Develop and train machine learning models
- Evaluate models using MAE and RMSE
- **Recommendation System Development (Weeks 11-14)**
 - Develop content-based filtering system
 - Integrate price prediction with user preferences
- **Web Interface Development (Week 11)**
 - Develop a user-friendly web interface for the recommendation system
- **Docker Setup and Integration (Week 12)**
 - Set up Docker to orchestrate the entire system
- **Experiment with Additional Models (Weeks 13-14)**
 - Implement and compare models like Gradient Boosting, Random Forests, and advanced deep learning architectures
 - Optimize hyperparameters for each model and compare their performance using various metrics
- **Expand Features and Incorporate New Parameters (Week 15)**
 - Collect and preprocess additional data for new features (e.g., socioeconomic and environmental factors)
 - Perform feature engineering to create new features from the existing data
 - Integrate these features into the model and evaluate their impact on performance
- **Geographical Expansion (Week 16)**
 - Collect data for regions beyond Buckinghamshire
 - Preprocess and integrate this data into the existing dataset
 - Evaluate the model's performance on the expanded dataset
- **User-Centric Testing and Feedback Collection (Week 17)**
 - Design a user-friendly interface for inputting preferences and receiving recommendations
 - Simulate user interactions using synthetic data
 - Collect feedback through synthetic data to refine the system
- **Replicate and Compare with High-Quality Models (Week 18)**
 - Replicate models from high-quality published papers
 - Compare their performance with your models and analyze the differences
- **Final Model Tuning and Evaluation (Week 19)**
 - Fine-tune the best-performing models
 - Conduct a thorough error analysis and identify areas for further improvement
 - Finalize the model and prepare it for deployment
- **Report Writing & Finalization (Week 20)**
 - Document all findings, methodologies, and results
 - Ensure the report is well-structured, with clear explanations and justifications for each step
 - Prepare for submission, ensuring all requirements are met

Week	Task
1-4	Data Collection & Preprocessing
5-10	Model Development
11	Web Interface Development
12	Docker Setup and Integration
13-14	Experiment with Additional Models

Week	Task
15	Expand Features and Incorporate New Parameters
16	Geographical Expansion (Collect Data for New Areas)
17	User-Centric Testing and Feedback Collection
18	Replicate and Compare with High-Quality Models
19	Final Model Tuning and Evaluation
20	Report Writing & Finalization

3.7 3.7 Testing and Evaluation Plan

3.7.1 Testing

- **Unit Testing:** Test individual components (e.g., data collection scripts, model training functions) to ensure they work as expected.
- **Integration Testing:** Ensure that different components (e.g., data integration, model prediction, recommendation system) work together seamlessly.
- **Synthetic User Testing:** Simulate user interactions using synthetic data to evaluate the recommendation system's usability and effectiveness.

3.7.2 Evaluation

- **Model Evaluation:** Use MAE and RMSE to evaluate the accuracy of the price prediction model.
- **Synthetic User Feedback:** Use synthetic data to simulate user feedback regarding the relevance and usefulness of the property recommendations.
- **Performance Metrics:** Track the system's performance in terms of response time, accuracy, and user satisfaction.

By following this structured approach and incorporating these components, the project aims to deliver a robust and effective personalized property recommendation system that meets user needs and leverages advanced machine learning techniques.

4 Chapter 4: Feature Prototype

Word Count: 1428

4.1 4.1 Introduction

Purpose of the feature prototype The feature prototype aims to demonstrate the feasibility of using machine learning techniques to predict property prices within the Personalised Property Recommendation System. This initial implementation focuses on developing a content-based filtering model to predict property prices based on historical transaction data from HM Land Registry and current property listings from OnTheMarket. Accurate price prediction is a crucial component of providing personalised property suggestions that align with users' budgets and preferences.

Role in demonstrating the feasibility of the project The development and evaluation of the price prediction model play several critical roles in demonstrating the feasibility of the Personalised Property Recommendation System:

1. **Proof of Concept:** Validates the potential of applying machine learning to real estate data.
2. **Budget Matching:** Helps filter and rank properties within user budgets.
3. **Market Price Estimation:** Provides insights into likely selling prices, aiding decision-making.
4. **Enhanced Personalization:** Integrates predicted prices with user preferences for sophisticated filtering and ranking.

By demonstrating these capabilities, the prototype showcases the potential for a comprehensive system that streamlines the property search process and tailors recommendations to individual needs.

4.2 Data Collection and Preprocessing

4.2.1 Data Collection Scripts

Introduction to Web Scraping Methodology To complement historical data from HM Land Registry, we web-scraped current property listings from OnTheMarket.com to gain insights into real-time market trends in Buckinghamshire.

Key Steps in Web Scraping Process

1. Identify Buckinghamshire property listing URLs.
2. Retrieve web content using Python's `requests`.
3. Extract data (prices, addresses, etc.) with `BeautifulSoup`.
4. Structure the extracted data.

```
[40]: # -q flag added to avoid printing copious amounts of logs
!pip install -q beautifulsoup4 lxml requests
!pip install -q ratelimit
!pip install -q tqdm
!pip install -q tensorflow scikit-learn pandas numpy matplotlib

# Add the directory containing custom modules to Python's import path
import sys
sys.path.append('./src/data-collector/')
sys.path.append('./src/data-cleanser/')
sys.path.append('./src/data-standardiser/')
```

Modularity and Code Structure Our web scraping is structured into four main modules, each with a specific role. This modular approach ensures a clean separation of concerns, making the code more maintainable and scalable.

1. `robot_check.py` - Respecting Site Policies

Ensures compliance with the website's scraping policies by parsing and interpreting the `robots.txt` file using `urllib.robotparser`. This module is used by both `crawler.py` and `data_collector_service.py`.

```
[41]: # src/data-collector/robot_check.py
import urllib.robotparser
```

```

class RobotCheck:
    def __init__(self, robots_txt_url):
        self.parser = urllib.robotparser.RobotFileParser()
        self.parser.set_url(robots_txt_url)
        self.parser.read()

    def is_allowed(self, url, user_agent):
        return self.parser.can_fetch(user_agent, url)

    def get_crawl_delay(self, user_agent):
        delay = self.parser.crawl_delay(user_agent)
        return delay if delay is not None else 1

```

2. crawler.py - Discovering URLs

Crawls the target website, gathering relevant page URLs using `requests` and `BeautifulSoup`. Rate limiting is implemented with `ratelimit` to avoid overloading the server. This module is used by `data_collector_service.py`.

```

[42]: # src/data-collector/crawler.py
import time
import requests
from bs4 import BeautifulSoup
from robot_check import RobotCheck
from ratelimit import limits, sleep_and_retry

# 10 requests per minute
REQUESTS_PER_MINUTE = 10

@sleep_and_retry
@limits(calls=REQUESTS_PER_MINUTE, period=60)
def make_request(url, headers):
    return requests.get(url, headers=headers)

def get_property_urls(base_url, search_url, user_agent):
    headers = {'User-Agent': user_agent}
    robot_check = RobotCheck("https://www.onthemarket.com/robots.txt")
    property_urls = set()
    page_number = 1

    while search_url:
        if robot_check.is_allowed(search_url, user_agent):
            response = make_request(search_url, headers=headers)

```

```

    if response.status_code == 200:
        soup = BeautifulSoup(response.content, 'html.parser')
        links = soup.select('div.otm-PropertyCardMedia > div > a')
        current_page_urls = set()

        for link in links:
            href = link.get('href')
            if href and href.startswith('/details/'): # Filtering
                other links
                full_url = 'https://www.onthemarket.com' + href
                property_urls.add(full_url)
                current_page_urls.add(full_url)

            # print(f"{len(current_page_urls)} URLs captured on Page
            {page_number}.")

            # Find the next page URL
            next_page = soup.select_one('a[title="Next page"]')
            if next_page:
                search_url = base_url + next_page['href']
                page_number += 1
            else:
                search_url = None

        else:
            print(f"Failed to retrieve the webpage: Status Code {response.
            status_code}")
            break

        # Sleep for the crawl delay
        time.sleep(robot_check.get_crawl_delay(user_agent))
    else:
        print("Scraping is disallowed by the website's robots.txt for this
        URL.")
        break

    return list(property_urls)

```

3. scraper.py - Extracting Data

Extracts specific data from web pages using requests and BeautifulSoup. This module is used by data_collector_service.py.

```

[43]: # src/data-collector/scrapper.py
import requests
from bs4 import BeautifulSoup

```

```

def scrape_property_details(property_url, headers, counter):
    try:
        response = requests.get(property_url, headers=headers)
        if response.status_code == 200:
            soup = BeautifulSoup(response.content, 'html.parser')

            # Extracting data
            # Extracting the title
            title = soup.find('h1', class_='h4 md:text-xl leading-normal').
            ↪get_text(strip=True) if soup.find('h1', class_='h4 md:text-xl
            ↪leading-normal') else 'Title not available'

            # Extracting the address
            address_div = soup.find('div', class_='text-slate h4 font-normal
            ↪leading-none font-heading')
            address = address_div.get_text(strip=True) if address_div else
            ↪'Address not available'

            # Extracting the price
            price_div = soup.find('div', class_='otm-Price')
            if price_div:
                price = price_div.find('span', class_='price').
            ↪get_text(strip=True) if price_div.find('span', class_='price') else 'Price
            ↪not available'

            # Extracting the pricing qualifier if present
            qualifier_div = price_div.find('small', class_='qualifier')
            price_qualifier = qualifier_div.get_text(strip=True) if
            ↪qualifier_div else 'Price qualifier not available'
            else:
                price = 'Price not available'
                price_qualifier = 'Price qualifier not available'

            # Extracting the listing time
            listing_time_div = soup.find('div', class_='text-denim')
            if listing_time_div:
                listing_time = listing_time_div.find('small',
            ↪class_='font-heading').get_text(strip=True) if listing_time_div.find(
                'small') else 'Listing time not available'
            else:
                listing_time = 'Listing time not available'

            # Extracting the property type
            property_type_div = soup.find('div', class_='otm-PropertyIcon')
            property_type = property_type_div.get_text(

```

```

        strip=True) if property_type_div else 'Property type not
↪available'

    details_div = soup.find('div', class_='font-heading text-xs flex
↪flex-wrap border-t border-b mb-6 md:text-md py-3 md:py-4 md:mb-9')
    if details_div:
        # Extracting the number of bedrooms
        bedrooms = 'Bedrooms info not available'
        for div in details_div.find_all('div'):
            if 'bed' in div.get_text(strip=True).lower():
                bedrooms = div.get_text(strip=True)
                break

        # Extracting the number of bathrooms
        bathrooms = 'Bathrooms info not available'
        for div in details_div.find_all('div'):
            if 'bath' in div.get_text(strip=True).lower():
                bathrooms = div.get_text(strip=True)
                break

        # Extracting the EPC rating
        epc_rating = 'EPC rating not available'
        for div in details_div.find_all('div'):
            if 'epc rating' in div.get_text(strip=True).lower():
                epc_rating = div.get_text(strip=True)
                break

        # Extracting the property size
        size = 'Size info not available'
        for div in details_div.find_all('div'):
            text = div.get_text(strip=True).lower()
            if 'sq ft' in text or 'sq m' in text:
                size = div.get_text(strip=True)
                break

        # Extracting features
        features_section = soup.find('section', class_='otm-FeaturesList')
        features = []
        if features_section:
            feature_items = features_section.find_all('li',
↪class_='otm-ListItemOtmBullet')
            for item in feature_items:
                feature_text = item.get_text(strip=True)
                features.append(feature_text)

    return {
        'id': counter,

```

```

        'property_url': property_url,
        'title': title,
        'address': address,
        'price': price,
        'pricing_qualifier': price_qualifier,
        'listing_time': listing_time,
        'property_type': property_type,
        'bedrooms': bedrooms,
        'bathrooms': bathrooms,
        'epc_rating': epc_rating,
        'size': size,
        'features': features
    }
    else:
        print(f"Failed to retrieve the property page: Status Code {response.
↳ status_code}")
        return {}
    except requests.exceptions.RequestException as e:
        return {'error': f"Request failed: {e}"}
    except Exception as e:
        return {'error': f"An unexpected error occurred: {e}"}

```

4. data_collector_service.py - Orchestrating the Scraping Process

Orchestrates the entire web scraping process. It calls on `crawler.py` to get URLs and then uses `scraper.py` to extract data from them. Handles errors, saves the data (usually in JSON), and can be run from the command line or Jupyter Notebook with adjustments.

```

[44]: # src/data-collector/data_collector_service.py
import json
from crawler import get_property_urls
from scraper import scrape_property_details
from robot_check import RobotCheck
from tqdm.notebook import tqdm
import time
import sys
import os

class DataCollectorService:
    def __init__(self, base_url, user_agent, max_price=120000):
        self.base_url = base_url
        self.user_agent = user_agent
        self.max_price = max_price
        self.robot_check = RobotCheck(f"{base_url}/robots.txt")

    @staticmethod
    def generate_price_segments(max_price, segment_size=100000):
        segments = []

```

```

current_min = 0
while current_min < max_price:
    current_max = min(current_min + segment_size, max_price)
    segments.append((current_min, current_max))
    current_min = current_max + 1
return segments

def collect_data_segment(self, min_price, max_price):
    search_url = f"{self.base_url}/for-sale/property/buckinghamshire/?
↪auction=false&min-price={min_price}&max-price={max_price}&new-home-flag=F&prop-types=bungal
    return get_property_urls(self.base_url, search_url, self.user_agent)

def collect_data(self):
    price_segments = self.generate_price_segments(self.max_price)
    all_property_urls = set()
    all_property_data = []
    counter = 1

    for min_price, max_price in price_segments:
        segment_urls = self.collect_data_segment(min_price, max_price)
        all_property_urls.update(segment_urls)
        print(f"Segment {min_price}-{max_price}: Found {len(segment_urls)}_
↪URLs.")

    print(f"Found a total of {len(all_property_urls)} property URLs from_
↪all segments.")

    for url in tqdm(all_property_urls, desc="Scraping properties"):
        if self.robot_check.is_allowed(url, self.user_agent):
            headers = {'User-Agent': self.user_agent}
            data = scrape_property_details(url, headers, counter)
            if 'error' in data:
                print(data['error'])
            else:
                all_property_data.append(data)
                counter += 1
                # Respect the crawl delay
                time.sleep(self.robot_check.get_crawl_delay(self.
↪user_agent))
        else:
            tqdm.write(f"Skipping {url}, disallowed by robots.txt.")

    self.save_data(all_property_data)

def save_data(self, data):
    data_directory = './data'
    filename = f'./data/property_data_{self.max_price}.json'

```



```

        backup_filename = f'./data/property_data_{self.max_price}-backup.json'

    try:
        # Ensure data directory exists
        if not os.path.exists(data_directory):
            os.makedirs(data_directory)
            print(f"Created directory {data_directory}")

        # Rename existing files if they exist
        if os.path.exists(backup_filename):
            os.remove(backup_filename)
        if os.path.exists(filename):
            os.rename(filename, backup_filename)
            print(f"Renamed existing file to {backup_filename}")

        # Save the new data
        with open(filename, 'w') as file:
            json.dump(data, file, indent=4)
        print(f"Data saved successfully in {filename}")
    except IOError as e:
        print(f"An IOError occurred while saving data: {e}")
    except Exception as e:
        print(f"An unexpected error occurred while saving data: {e}")

if __name__ == "__main__":
    max_price_arg = 120000
    base_url = 'https://www.onthemarket.com'
    search_url = base_url + '/for-sale/property/buckinghamshire/?
    ↪auction=false&max-price=500000&min-price=10000&new-home-flag=F&prop-types=bungalows&prop-ty
    ↪ac.uk)'
    user_agent = 'StudentDataProjectScraper/1.0 (Contact: gor5@student.london.
    service = DataCollectorService(base_url, user_agent, max_price_arg)
    service.collect_data()

```

Segment 0-100000: Found 0 URLs.

Segment 100001-120000: Found 8 URLs.

Found a total of 8 property URLs from all segments.

Scraping properties: 0%| | 0/8 [00:00<?, ?it/s]

Renamed existing file to ./data/property_data_120000-backup.json

Data saved successfully in ./data/property_data_120000.json

Ethical Considerations Our web scraping adhered to ethical standards, including respecting robots.txt, using a unique user agent, implementing rate limiting, and ensuring non-disruptive interaction with OnTheMarket.com.

Note on Script Execution Time Due to our ethical scraping approach, the script's execution takes longer for larger datasets, particularly given the rate limits and crawl delays we adhere to. This ensures responsible scraping while avoiding potential blocking by the website.

We can run `data_collector_service.py` from the CLI or from the Jupyter Notebook using `%run` and provide the `max_price` argument to control the data collection scope. For this prototype, data was collected for properties priced up to £650,000. This data can be found in the `data` folder

```
[45]: # %run src/data-collector/data_collector_service.py 650000
      %run src/data-collector/data_collector_service.py 120000
```

Segment 0-100000: Found 0 URLs.

Segment 100001-120000: Found 8 URLs.

Found a total of 8 property URLs from all segments.

Scraping properties: 0%| | 0/8 [00:00<?, ?it/s]

Renamed existing file to `./data/property_data_120000-backup.json`

Data saved successfully in `./data/property_data_120000.json`

Displaying Scraped Data Verifying the integrity and structure of the scraped data is essential to ensure it aligns with our requirements and is ready for analysis.

```
[46]: import pandas as pd

      # Load the scraped data from the JSON file
      file_path = './data/property_data_650000.json'
      scraped_data = pd.read_json(file_path)

      # Display the first 5 to 10 rows of the dataset
      scraped_data.head(10)
```

```
[46]:   id                                property_url \
0    1  https://www.onthemarket.com/details/14543582/
1    2  https://www.onthemarket.com/details/14855003/
2    3  https://www.onthemarket.com/details/15004462/
3    4  https://www.onthemarket.com/details/14302804/
4    5  https://www.onthemarket.com/details/14963433/
5    6  https://www.onthemarket.com/details/14371850/
6    7  https://www.onthemarket.com/details/13238957/
7    8  https://www.onthemarket.com/details/14214029/
8    9  https://www.onthemarket.com/details/14381566/
9   10  https://www.onthemarket.com/details/14139162/

      title \
0      3 bedroom terraced house for sale
1  3 bedroom semi-detached house for sale
2      2 bedroom apartment for sale
3  3 bedroom semi-detached house for sale
```

4 2 bedroom maisonette for sale
5 3 bedroom terraced house for sale
6 2 bedroom apartment for sale
7 3 bedroom detached house for sale
8 4 bedroom detached house for sale
9 1 bedroom apartment for sale

	address	price \
0	Green Lane, Wolverton, Milton Keynes	£385,000
1	Bernay Gardens, Bolbeck Park, Milton Keynes	£429,995
2	Aylesbury, Aylesbury HP19	£215,000
3	Fox Way, Buckingham	£379,995
4	Hillary Close, Aylesbury HP21	£265,000
5	WELLINGBOROUGH ROAD, OLNEY	£595,000
6	Grange Road, Chalfont St Peter SL9	£385,000
7	Farmers Way, Seer Green, HP9	£620,000
8	FOXHILL, OLNEY	£475,000
9	Maypole Road, Taplow SL6	£185,000

	pricing_qualifier	listing_time	property_type \
0	Price qualifier not available	Added > 14 days	Terraced house
1	Price qualifier not available	Added > 14 days	Semi-detached house
2	Guide price	Added < 14 days	Apartment
3	Price qualifier not available	Added > 14 days	Semi-detached house
4	Price qualifier not available	Added > 14 days	Maisonette
5	Guide price	Added > 14 days	Terraced house
6	Price qualifier not available	Added > 14 days	Apartment
7	Guide price	Added > 14 days	Detached house
8	Offers over	Added > 14 days	Detached house
9	Offers in excess of	Added > 14 days	Apartment

	bedrooms	bathrooms	epc_rating	size \
0	3 bed	2bath	EPC rating: E*	1,400 sq ft / 130 sq m
1	3 bed	2bath	EPC rating: C*	947 sq ft / 88 sq m
2	2 bed	1bath	EPC rating: D*	538 sq ft / 50 sq m
3	3 bed	1bath	EPC rating: C*	904 sq ft / 84 sq m
4	2 bed	1bath	EPC rating: C*	592 sq ft / 55 sq m
5	3 bed	2bath	EPC rating: D*	1,334 sq ft / 124 sq m
6	2 bed	1bath	EPC rating: B*	753 sq ft / 70 sq m
7	3 bed	2bath	EPC rating not available	1,142 sq ft / 106 sq m
8	4 bed	2bath	EPC rating: C*	1,194 sq ft / 111 sq m
9	1 bed	1bath	EPC rating: D*	365 sq ft / 34 sq m

	features
0	[THREE DOUBLE BEDROOMS, VICTORIAN TERRACED, PE...
1	[En-suite shower room, Stylish refitted kitche...
2	[Tenure: Leasehold, TWO BEDROOMS, ON SITE GYM,...

```

3 [Three Bedroom Semi Detached, Corner Plot, Dri...
4 [Southside Of Aylesbury, Two Double Bedrooms, ...
5 [Tenure: Freehold, DOUBLE FRONTED PROPERTY, CL...
6 [Tenure: Leasehold, Beautifully Presented Thro...
7 [Tenure: Freehold, Garage, Garden, Walk to sta...
8 [Tenure: Freehold, FOUR BEDROOM DETACHED FAMIL...
9 [Tenure: Leasehold, One bedroom ground floor a...

```

4.2.2 Initial Data Exploration and Analysis

Examination of the HM Land Registry Dataset We examined the UK-wide property transaction dataset (`pp-monthly-update-new-version.csv`) from HM Land Registry and filtered it to focus on transactions in Buckinghamshire.

```

[47]: import pandas as pd

# Path to your CSV file
file_path = './data/historical-data/pp-monthly-update-new-version.csv'

# Read the CSV file and display the first 10 rows
df = pd.read_csv(file_path)
print("First 10 rows of the data are:")
df.head(10)

```

First 10 rows of the data are:

```

[47]: {09266DDB-86DB-AF90-E063-4704A8C02087} 323450 2021-11-30 00:00 BS13 OFF \
0 {09266DDB-86EF-AF90-E063-4704A8C02087} 420000 2021-05-26 00:00 BS2 9NY
1 {09266DDB-86F5-AF90-E063-4704A8C02087} 185000 2021-10-22 00:00 BS14 OTL
2 {09266DDB-86F8-AF90-E063-4704A8C02087} 269000 2021-12-17 00:00 BS13 7DL
3 {09266DDB-8738-AF90-E063-4704A8C02087} 520000 2021-09-23 00:00 BS2 9XB
4 {09266DDB-86A8-AF90-E063-4704A8C02087} 195000 2021-09-29 00:00 BS4 5AA
5 {09266DDB-86A9-AF90-E063-4704A8C02087} 245000 2021-11-29 00:00 BS7 8BW
6 {09266DDB-86AA-AF90-E063-4704A8C02087} 204995 2021-10-08 00:00 BS4 1FN
7 {09266DDB-86AB-AF90-E063-4704A8C02087} 265000 2021-12-09 00:00 BS9 3UL
8 {09266DDB-86AD-AF90-E063-4704A8C02087} 194995 2021-12-10 00:00 BS16 2GZ
9 {09266DDB-86AE-AF90-E063-4704A8C02087} 185950 2021-05-26 00:00 BS5 9HH

```

```

      S  Y  F      23 Unnamed: 8      BUTTERFLY LANE Unnamed: 10 \
0  T  Y  F      201      NaN      NEWFOUNDLAND ROAD      NaN
1  F  Y  L  BOULEVARD VIEW      22      WHITCHURCH LANE      NaN
2  F  Y  L      84      NaN      OLD PUMP HOUSE CLOSE      NaN
3  T  N  F      7      NaN      NARROWAYS ROAD      NaN
4  F  N  L      41A      NaN      BRISTOL HILL      NaN
5  F  N  L      12      FLAT 2      BOLTON ROAD      NaN
6  F  N  L      66      FLAT 4      MANNING ROAD      NaN
7  F  N  L      LITTLE COTE      FLAT 3      COTE LANE      NaN
8  F  N  L      14      NaN      NAPOLEON AVENUE      FISHPONDS

```

9	F	N	L	124A	NaN	CHURCH ROAD	REDFIELD
	BRISTOL	CITY OF BRISTOL	CITY OF BRISTOL	1	A	A	1
0	BRISTOL	CITY OF BRISTOL	CITY OF BRISTOL		A	A	
1	BRISTOL	CITY OF BRISTOL	CITY OF BRISTOL		A	A	
2	BRISTOL	CITY OF BRISTOL	CITY OF BRISTOL		A	A	
3	BRISTOL	CITY OF BRISTOL	CITY OF BRISTOL		A	A	
4	BRISTOL	CITY OF BRISTOL	CITY OF BRISTOL		A	A	
5	BRISTOL	CITY OF BRISTOL	CITY OF BRISTOL		A	A	
6	BRISTOL	CITY OF BRISTOL	CITY OF BRISTOL		A	A	
7	BRISTOL	CITY OF BRISTOL	CITY OF BRISTOL		A	A	
8	BRISTOL	CITY OF BRISTOL	CITY OF BRISTOL		A	A	
9	BRISTOL	CITY OF BRISTOL	CITY OF BRISTOL		A	A	

4.2.3 4.2.3 Cleaning and Preparing Data

Data Cleaning Methodology Our data cleaning methodology, using tools like pandas, involves filtering, standardising, and handling missing/erroneous data from HM Land Registry and OnTheMarket.com for our Buckinghamshire property market analysis.

Implementing Data Cleaning with data_cleanser_service.py The data_cleanser_service.py script is a key component in refining raw HM Land Registry data for our Buckinghamshire property market analysis.

Key Features of data_cleanser_service.py:

- **Header Assignment:** The script assigns column headers to the HM Land Registry dataset (pp-monthly-update-new-version.csv) based on definitions from their website (<https://www.gov.uk/guidance/about-the-price-paid-data>).
- **Loading and Structuring Data:** Loads the CSV data into a pandas DataFrame.
- **Date Conversion and Filtering:** Retains only properties located in Buckinghamshire.

```
[48]: # data-cleanser/data_cleanser_service.py
import pandas as pd

def cleanse_data(input_file, output_file):
    # Define the headers based on the provided breakdown
    headers = ["Unique Transaction Identifier", "Price", "Date of Transaction",
               "Postal Code", "Property Type", "Old/New", "Duration",
               "PAON", "SAON", "Street", "Locality", "Town/City",
               "District", "County", "PPD Category Type", "Record Status"]

    # Load the CSV file without headers
    data = pd.read_csv(input_file, header=None, names=headers)

    # Convert Date of Transaction to datetime for filtering
    print("Converting dates and filtering data...")
    data['Date of Transaction'] = pd.to_datetime(data['Date of Transaction'])
```

```

# Filter for properties in Buckinghamshire and from the year 2023
data['Date of Transaction'] = pd.to_datetime(data['Date of Transaction'])
filtered_data = data[(data['County'].str.upper() == 'BUCKINGHAMSHIRE') &
                     (data['Date of Transaction'].dt.year == 2023)]

print(f"Number of records after filtering: {len(filtered_data)}")

print("Saving cleaned data to CSV file...")
# Save the cleaned data to a new CSV file
filtered_data.to_csv(output_file, index=False)
print("Data cleansing process completed successfully.")

# File paths
input_csv = './data/historical-data/pp-monthly-update-new-version.csv' #_
    ↳Update with actual path
output_csv = './data/historical-data/buckinghamshire_2023_cleaned_data.csv' #_
    ↳Update with desired output path

cleanse_data(input_csv, output_csv)

```

Converting dates and filtering data...

Number of records after filtering: 438

Saving cleaned data to CSV file...

Data cleansing process completed successfully.

4.2.4 Enhancing Data with Geocoding and Merging using data_standardiser_service.py data_standardiser_service.py enhances and merges datasets:

Key Features of data_standardiser_service.py: 1. **Geocoding**: Uses Nominatim and ArcGIS for rate-limited geocoding, with a fallback mechanism if one fails. 2. **Price Standardisation**: Converts various price formats into a uniform numerical format. 3. **Address Normalization**: Standardizes addresses for consistency. 4. **Dataset Merging**: Combines the processed scraped and registry data into a single DataFrame. 5. **Saving & Updating**: Saves the enriched dataset and updates existing data. 6. **Property Type Classification**: Ensures consistency across the dataset. 7. **Updating Existing Preprocessed Data**: Classifies scraped data rows according to the registry data classification system.

```

[49]: # data-standardiser/data_standardiser_service.py but with relative file path_
    ↳changed
import pandas as pd
import os
from datetime import datetime
from geopy.geocoders import Nominatim, ArcGIS
from geopy.extra.rate_limiter import RateLimiter
from geopy.exc import GeocoderTimedOut, GeocoderQuotaExceeded

```

```

import time

# Initialize Nominatim API
geolocator = Nominatim(user_agent="StudentDataProjectScraper/1.0 (Contact:✉gor5@student.london.ac.uk)")
geolocator_arcgis = ArcGIS(user_agent="StudentDataProjectScraper/1.0 (Contact:✉gor5@student.london.ac.uk)")

# Rate limiter to avoid overloading the API
geocode = RateLimiter(geolocator.geocode, min_delay_seconds=2)
geocode_arcgis = RateLimiter(geolocator_arcgis.geocode, min_delay_seconds=2)

# Mapping for converting scraped property types to registry property types
scraped_to_registry_property_type_mapping = {
    'Apartment': 'F',
    'Barn conversion': 'O',
    'Block of apartments': 'F',
    'Bungalow': 'D',
    'Character property': 'O',
    'Cluster house': 'O',
    'Coach house': 'F',
    'Cottage': 'D',
    'Detached bungalow': 'D',
    'Detached house': 'D',
    'Duplex': 'F',
    'End of terrace house': 'T',
    'Equestrian property': 'O',
    'Farm house': 'O',
    'Flat': 'F',
    'Ground floor flat': 'F',
    'Ground floor maisonette': 'F',
    'House': 'D',
    'Link detached house': 'D',
    'Lodge': 'O',
    'Maisonette': 'F',
    'Mews': 'O',
    'Penthouse': 'F',
    'Semi-detached bungalow': 'D',
    'Semi-detached house': 'S',
    'Studio': 'F',
    'Terraced house': 'T',
    'Townhouse': 'D'
}

def geocode_address(address):
    try:
        location = geocode(address)

```

```

        if location:
            print(f"Geocoded '{address}': Latitude {location.latitude},  

↳Longitude {location.longitude}")
            return location.latitude, location.longitude
        else:
            # Fallback to ArcGIS if Nominatim fails
            location = geocode_arcgis(address)
            if location:
                print(f"Geocoded '{address}': Latitude {location.latitude},  

↳Longitude {location.longitude}")
                return location.latitude, location.longitude
            else:
                print(f"No result for '{address}'")
                return None, None
    except GeocoderQuotaExceeded:
        print("Quota exceeded for geocoding API")
        return None, None
    except GeocoderTimedOut:
        print("Geocoding API timed out")
        return None, None
    except Exception as e:
        print(f"Error geocoding {address}: {e}")
        return None, None

def check_preprocessed_file(file_path):
    """Check if the preprocessed file exists and has latitude and longitude."""
    if os.path.exists(file_path):
        df = pd.read_csv(file_path)
        if 'latitude' in df.columns and 'longitude' in df.columns:
            if df[['latitude', 'longitude']].notnull().all().all():
                # File exists and latitude and longitude are filled
                return True
    return False

def standardise_price(price):
    """
    Convert a price string to a numerical value.
    Handles strings like '£275,000' and converts them to 275000.
    """
    if not isinstance(price, str):
        return price # If it's already a number, return as-is

    # Removing currency symbols and commas
    price = price.replace('£', '').replace(',', '').replace('€', '').strip()

    try:
        # Convert to float or int

```



```

        price_value = float(price) if '.' in price else int(price)
    except ValueError:
        # Handle cases where conversion fails
        print(f"Warning: Could not convert price '{price}' to a number.")
        price_value = None

    return price_value

def normalize_address_scraped(address):
    """
    Normalize addresses from the scraped data.
    """
    # Assuming the county is always 'Buckinghamshire' if not specified
    if 'Buckinghamshire' not in address:
        address += ', Buckinghamshire'
    return address.strip()

def normalize_address_land_registry(row):
    # Convert each component to a string to avoid TypeError
    components = [
        str(row['Street']),
        str(row['Locality']),
        str(row['Town/City']),
        str(row['District']),
        str(row['County'])
    ]
    # Join the non-empty components
    return ', '.join(filter(None, components))

# Read JSON, standardize price, normalize address, add source column

def read_and_process_scraped_data(scraped_file_path, skip_geocoding):
    # Read the scraped data
    scraped_data = pd.read_json(scraped_file_path)
    scraped_data['price'] = scraped_data['price'].apply(standardise_price)
    scraped_data['normalized_address'] = scraped_data['address'].
    ↪ apply(normalize_address_scraped)
    scraped_data['source'] = 'scraped'

    if not skip_geocoding:
        lat_long = scraped_data['normalized_address'].apply(geocode_address)
        scraped_data['latitude'] = lat_long.apply(lambda x: x[0] if x else None)
        scraped_data['longitude'] = lat_long.apply(lambda x: x[1] if x else
    ↪ None)

    return scraped_data

```

```

def read_and_process_registry_data(registry_file_path, skip_geocoding):
    registry_data = pd.read_csv(registry_file_path)
    registry_data['Price'] = registry_data['Price'].apply(standardise_price)
    registry_data['normalized_address'] = registry_data.
    ↪apply(normalize_address_land_registry, axis=1)
    registry_data.rename(columns={'Price': 'price'}, inplace=True)
    registry_data['source'] = 'registry'

    if not skip_geocoding:
        lat_long = registry_data['normalized_address'].apply(geocode_address)
        registry_data['latitude'] = lat_long.apply(lambda x: x[0] if x else ↪
    ↪None)
        registry_data['longitude'] = lat_long.apply(lambda x: x[1] if x else ↪
    ↪None)

    return registry_data

def update_date_column(df, source_column, new_date):
    """
    Update 'Date' column in the DataFrame based on source.
    """
    df['Date'] = pd.NaT
    df.loc[df['source'] == 'registry', 'Date'] = pd.
    ↪to_datetime(df[source_column])
    df.loc[df['source'] == 'scraped', 'Date'] = new_date
    return df

def merge_price_columns(df):
    """
    Merge 'price' and 'Price' columns and update 'Price' with 'price' values ↪
    ↪for scraped data.
    """
    # Use 'price' from scraped data if it is not NaN, else use 'Price' from ↪
    ↪registry data
    df['Price'] = df.apply(lambda x: x['price'] if pd.notna(x['price']) else ↪
    ↪x['Price'], axis=1)
    return df

def apply_property_type_mapping(df):
    """
    Apply property type mapping to the DataFrame, only if 'property_type' is ↪
    ↪not null.
    """
    # Apply mapping only where 'property_type' is not null
    mask = df['property_type'].notnull()

```

```

        df.loc[mask, 'Property Type'] = df.loc[mask, 'property_type'].
        ↪map(scraped_to_registry_property_type_mapping)
        return df

def process_and_save_data(scraped_data, registry_data, output_file_path):
    """
    Process and save merged data.
    """
    # Merge datasets
    merged_data = pd.concat([scraped_data, registry_data], ignore_index=True)

    # Update the date column
    merged_data = update_date_column(merged_data, 'Date of Transaction',
    ↪datetime(2023, 12, 31))

    # Apply property type mapping
    merged_data = apply_property_type_mapping(merged_data)

    # Merge 'price' and 'Price' columns
    merged_data = merge_price_columns(merged_data)

    # Save merged data
    merged_data.to_csv(output_file_path, index=False)
    print(f"Merged data saved successfully to '{output_file_path}'.")

def main():
    scraped_file = './data/property_data_650000.json'
    registry_file = './data/historical-data/buckinghamshire_2023_cleaned_data.
    ↪csv'
    output_file = './data/preprocessed-data/preprocessed.csv'

    if check_preprocessed_file(output_file):
        # Read existing preprocessed data
        preprocessed_data = pd.read_csv(output_file)
        print(f"Using existing preprocessed data from '{output_file}'.")

        # Update the date column in the existing data
        preprocessed_data = update_date_column(preprocessed_data, 'Date of
    ↪Transaction', datetime(2023, 12, 31))

        # Apply property type mapping
        preprocessed_data = apply_property_type_mapping(preprocessed_data)

        # Apply property type mapping and merge price columns
        preprocessed_data = merge_price_columns(preprocessed_data)

        # Save updated data

```

```

preprocessed_data.to_csv(output_file, index=False)
print(f"Updated data saved successfully to '{output_file}'.")
else:
    # Process new data
    scraped_data = read_and_process_scraped_data(scraped_file, False)
    registry_data = read_and_process_registry_data(registry_file, False)

    # Process and save merged data
    process_and_save_data(scraped_data, registry_data, output_file)

print("Data processing completed.")

if __name__ == "__main__":
    main()

```

Using existing preprocessed data from './data/preprocessed-data/preprocessed.csv'.

Updated data saved successfully to './data/preprocessed-data/preprocessed.csv'.
Data processing completed.

The dataset below combines web-scraped data (December 2023) and HM Land Registry data (2023) for Buckinghamshire properties. The first and last five entries are shown, summarising the dataset's structure and content.

```

[50]: import pandas as pd

# Path to your CSV file
file_path = './data/preprocessed-data/preprocessed.csv'

# Read the CSV file
df = pd.read_csv(file_path)

# Display the first 5 rows
print("First 5 rows of the data are:")
print(df.head(5))

# Display the last 5 rows
print("\nLast 5 rows of the data are:")
print(df.tail(5))

# Print the column names
print("Column names:")
print(df.columns)

```

First 5 rows of the data are:

	id	property_url \
0	1.0	https://www.onthemarket.com/details/13262643/
1	2.0	https://www.onthemarket.com/details/14047885/
2	3.0	https://www.onthemarket.com/details/14068644/

```

3 4.0 https://www.onthemarket.com/details/13944473/
4 5.0 https://www.onthemarket.com/details/13422734/

```

```

                                title \
0          1 bedroom apartment for sale
1          2 bedroom apartment for sale
2 4 bedroom semi-detached house for sale
3          1 bedroom flat for sale
4          3 bedroom flat for sale

```

```

                                address      price \
0          The Green, High Wycombe HP10 275000.0
1          Aylesbury, Buckinghamshire HP21 220000.0
2          Edzell Crescent, Milton Keynes MK4 465000.0
3          Scafell Road, Slough 195000.0
4 High Wycombe, Buckinghamshire, HP13 265000.0

```

```

                                pricing_qualifier      listing_time      property_type \
0 Price qualifier not available OnTheMarket > 14 days      Apartment
1 Price qualifier not available OnTheMarket > 14 days      Apartment
2 Price qualifier not available OnTheMarket < 14 days Semi-detached house
3          Guide price OnTheMarket > 14 days      Flat
4 Price qualifier not available OnTheMarket > 14 days      Flat

```

```

bedrooms bathrooms ... PAON SAON Street Locality Town/City District \
0 1 bed 1bath ... NaN NaN NaN NaN NaN NaN
1 2 bed 2bath ... NaN NaN NaN NaN NaN NaN
2 4 bed 2bath ... NaN NaN NaN NaN NaN NaN
3 1 bed 1bath ... NaN NaN NaN NaN NaN NaN
4 3 bed 1bath ... NaN NaN NaN NaN NaN NaN

```

```

County PPD Category Type Record Status      Date
0 NaN NaN NaN NaN NaN 2023-12-31
1 NaN NaN NaN NaN NaN 2023-12-31
2 NaN NaN NaN NaN NaN 2023-12-31
3 NaN NaN NaN NaN NaN 2023-12-31
4 NaN NaN NaN NaN NaN 2023-12-31

```

[5 rows x 34 columns]

Last 5 rows of the data are:

```

id property_url title address price pricing_qualifier listing_time \
2595 NaN NaN NaN NaN NaN NaN NaN
2596 NaN NaN NaN NaN NaN NaN NaN
2597 NaN NaN NaN NaN NaN NaN NaN
2598 NaN NaN NaN NaN NaN NaN NaN
2599 NaN NaN NaN NaN NaN NaN NaN

```

	property_type	bedrooms	bathrooms	...	PAON	SAON	\
2595	NaN	NaN	NaN	...	27	NaN	
2596	NaN	NaN	NaN	...	PEATEY COURT	128	
2597	NaN	NaN	NaN	...	2	NaN	
2598	NaN	NaN	NaN	...	4	NaN	
2599	NaN	NaN	NaN	...	129	NaN	

	Street	Locality	Town/City	District	\
2595	DELLFIELD	NaN	CHESHAM	BUCKINGHAMSHIRE	
2596	PRINCES GATE	NaN	HIGH WYCOMBE	BUCKINGHAMSHIRE	
2597	YEAT FARM COTTAGES	WOTTON UNDERWOOD	AYLESBURY	BUCKINGHAMSHIRE	
2598	BADGERS RISE	STONE	AYLESBURY	BUCKINGHAMSHIRE	
2599	COULSON WAY	BURNHAM	SLOUGH	BUCKINGHAMSHIRE	

	County	PPD Category	Type	Record Status	Date
2595	BUCKINGHAMSHIRE		A	D	2023-03-23
2596	BUCKINGHAMSHIRE		A	D	2023-03-25
2597	BUCKINGHAMSHIRE		A	D	2023-05-22
2598	BUCKINGHAMSHIRE		A	D	2023-03-20
2599	BUCKINGHAMSHIRE		B	D	2023-02-23

[5 rows x 34 columns]

Column names:

```
Index(['id', 'property_url', 'title', 'address', 'price', 'pricing_qualifier',
      'listing_time', 'property_type', 'bedrooms', 'bathrooms', 'epc_rating',
      'size', 'features', 'normalized_address', 'source', 'latitude',
      'longitude', 'Unique Transaction Identifier', 'Price',
      'Date of Transaction', 'Postal Code', 'Property Type', 'Old/New',
      'Duration', 'PAON', 'SAON', 'Street', 'Locality', 'Town/City',
      'District', 'County', 'PPD Category Type', 'Record Status', 'Date'],
      dtype='object')
```

4.2.4 4.3 Model Development

4.3.1 Content-Based Filtering Model

Introduction to Content-Based Filtering Content-based filtering is a recommendation system technique that uses the features of items to recommend other items with similar attributes. In the context of the Personalised Property Recommendation System, we will use property features such as price, geographical coordinates (latitude and longitude), and property type to recommend properties to users.

Relevant Dataset Features Our dataset includes the following relevant features for each property:

- Price
- Latitude
- Longitude

- **Property Type:** The type of the property (e.g., apartment, detached house).

These features are present in all rows of our preprocessed dataset and will be used to develop our recommendation model.

Building the Neural Network Model We will construct a neural network using TensorFlow to implement the content-based filtering model. The network will have an input layer corresponding to the features, hidden layers to capture complex patterns, and an output layer for the final recommendation score.

```
[51]: import pandas as pd

# Load the preprocessed data
file_path = './data/preprocessed-data/preprocessed.csv'
df = pd.read_csv(file_path)

# Select relevant features for the model
features = df[['Price', 'Property Type', 'latitude', 'longitude']]

# Handle missing values
features = features.dropna()

# Convert categorical features to numeric
features = pd.get_dummies(features, columns=['Property Type'])

# Ensure the target variable matches the processed features dataframe
target = df.loc[features.index, 'Price']

# Check for invalid data
print("Checking for invalid data in features:")
print(features.isnull().sum())
print(features.describe())

print("Checking for invalid data in target:")
print(target.isnull().sum())
print(target.describe())

# Display the first few rows of the extracted features
print(features.head())
```

Checking for invalid data in features:

Price	0
latitude	0
longitude	0
Property Type_D	0
Property Type_F	0
Property Type_O	0
Property Type_S	0
Property Type_T	0

```
dtype: int64
```

	Price	latitude	longitude
count	2.600000e+03	2600.000000	2600.000000
mean	4.341690e+05	51.806702	-0.737249
std	2.698687e+05	0.194997	0.115481
min	1.000000e+03	51.481425	-1.115640
25%	2.900000e+05	51.628329	-0.803560
50%	4.000000e+05	51.811928	-0.741515
75%	5.250000e+05	52.004868	-0.658138
max	4.100000e+06	52.177271	-0.477487

Checking for invalid data in target:

```
0
count      2.600000e+03
mean       4.341690e+05
std        2.698687e+05
min        1.000000e+03
25%        2.900000e+05
50%        4.000000e+05
75%        5.250000e+05
max        4.100000e+06
```

Name: Price, dtype: float64

	Price	latitude	longitude	Property	Type_D	Property	Type_F	\
0	275000.0	51.587275	-0.683926		False		True	
1	220000.0	51.803094	-0.817448		False		True	
2	465000.0	52.004837	-0.802238		False		False	
3	195000.0	51.527798	-0.634636		False		True	
4	265000.0	51.628329	-0.742618		False		True	

	Property	Type_0	Property	Type_S	Property	Type_T
0		False		False		False
1		False		False		False
2		False		True		False
3		False		False		False
4		False		False		False

4.3.2 Model Training

Data Splitting and Preparation:

```
[52]: from sklearn.preprocessing import StandardScaler
      from sklearn.model_selection import train_test_split

      # Split the data into training and validation sets
      X_train, X_val, y_train, y_val = train_test_split(features, target, test_size=0.
      ↪2, random_state=42)

      # Normalize the features and target
      scaler_features = StandardScaler()
```



```

X_train = scaler_features.fit_transform(X_train)
X_val = scaler_features.transform(X_val)

scaler_target = StandardScaler()
y_train = scaler_target.fit_transform(y_train.values.reshape(-1, 1)).flatten()
y_val = scaler_target.transform(y_val.values.reshape(-1, 1)).flatten()

# Display the shapes of the training and validation sets
print(X_train.shape, X_val.shape, y_train.shape, y_val.shape)

```

(2080, 8) (520, 8) (2080,) (520,)

Neural Network Model Structure:

```

[53]: import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout, Input

# Display the shapes of the training and validation sets
print(X_train.shape, X_val.shape, y_train.shape, y_val.shape)

# Define the model structure using the Input layer
model = Sequential([
    Input(shape=(X_train.shape[1],)),
    Dense(128, activation='relu', kernel_regularizer=tf.keras.regularizers.l2(0.
↪01)),
    Dropout(0.2),
    Dense(64, activation='relu', kernel_regularizer=tf.keras.regularizers.l2(0.
↪01)),
    Dropout(0.2),
    Dense(1, activation='linear')
])

# Compile the model with a lower learning rate
model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=0.0001),
↪loss='mean_squared_error', metrics=['mae'])

# Display the model summary
model.summary()

```

(2080, 8) (520, 8) (2080,) (520,)

Model: "sequential_2"

Layer (type)	Output Shape	Param #
dense_6 (Dense)	(None, 128)	1,152

dropout_4 (Dropout)	(None, 128)	0
dense_7 (Dense)	(None, 64)	8,256
dropout_5 (Dropout)	(None, 64)	0
dense_8 (Dense)	(None, 1)	65

Total params: 9,473 (37.00 KB)

Trainable params: 9,473 (37.00 KB)

Non-trainable params: 0 (0.00 B)

Training Parameters and Process:

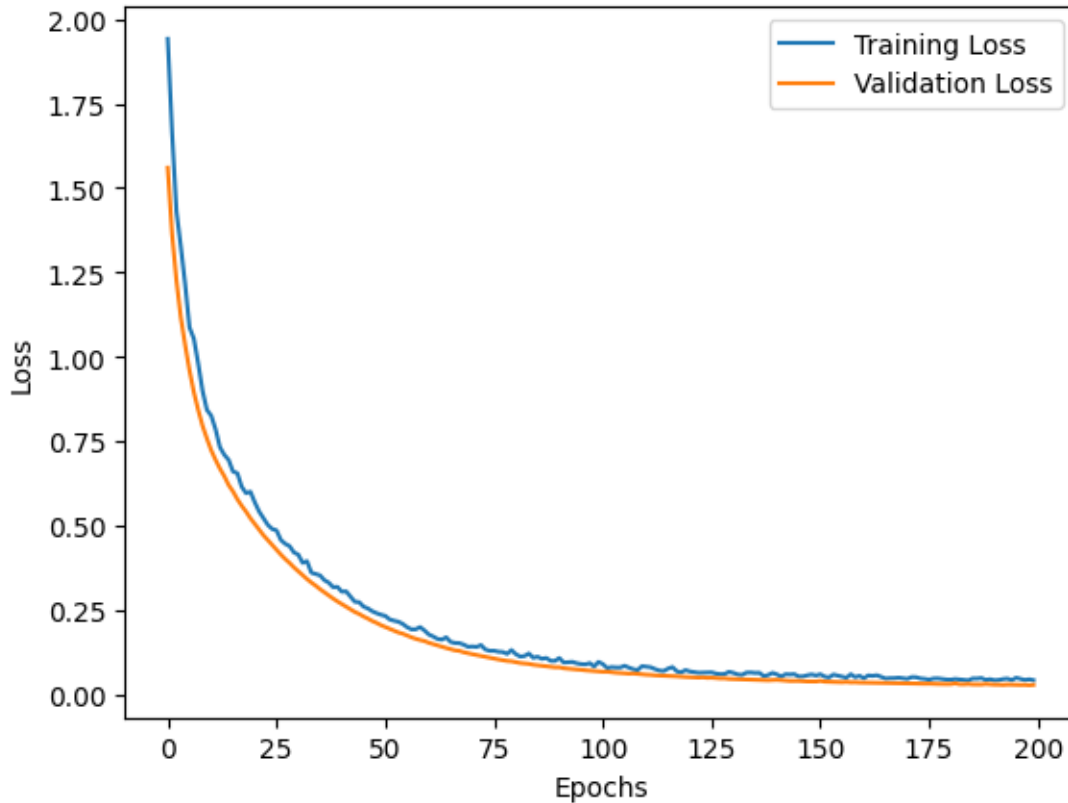
```
[54]: # Plot training history
import matplotlib.pyplot as plt
from tensorflow.keras.callbacks import EarlyStopping

# Define early stopping callback
early_stopping = EarlyStopping(monitor='val_loss', patience=10,
                               restore_best_weights=True, verbose=0)

# Train the model
history = model.fit(X_train, y_train, epochs=200, validation_data=(X_val,
                                                                y_val),
                    batch_size=32, callbacks=[early_stopping], verbose=0)

# Print the training history to verify the collected data
# print(history.history)

# Plot training history
plt.plot(history.history['loss'], label='Training Loss')
plt.plot(history.history['val_loss'], label='Validation Loss')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()
plt.show()
```



4.2.5 4.4 Evaluation and Improvements

4.4.1 Model Evaluation After training the content-based filtering model, we can evaluate its performance using appropriate metrics. We use Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) as our evaluation metrics.

Evaluation Metrics:

- **Mean Absolute Error (MAE):** Measures the average magnitude of errors in a set of predictions. It is calculated as the average of the absolute differences between predicted and actual values.
- **Root Mean Squared Error (RMSE):** The square root of the average of squared differences between predicted and actual values.

```
[55]: from sklearn.metrics import mean_absolute_error, mean_squared_error
import numpy as np

# Predict on the validation set
predictions = model.predict(X_val)

# Rescale the predictions and true values back to the original scale
```

```

y_val_rescaled = scaler_target.inverse_transform(y_val.reshape(-1, 1)).flatten()
predictions_rescaled = scaler_target.inverse_transform(predictions).flatten()

# Calculate evaluation metrics
mae = mean_absolute_error(y_val_rescaled, predictions_rescaled)
rmse = np.sqrt(mean_squared_error(y_val_rescaled, predictions_rescaled))

print(f'Mean Absolute Error: {mae}')
print(f'Root Mean Squared Error: {rmse}')

```

```

17/17          0s 2ms/step
Mean Absolute Error: 5336.591195913462
Root Mean Squared Error: 11133.510108428378

```

The MAE and RMSE values indicate the average and squared differences between the predicted and actual property prices, respectively. The values indicate that the model's property price predictions are reasonably accurate, deviating by roughly £5,336 on average.

4.4.2 Discussion and Improvements

Integrating Price Prediction into the Recommendation System The present model primarily emphasises the forecast of property prices, which plays a crucial role in the recommendation system. The ability to determine the projected price of a property enables the system to categorise and prioritise properties based on a user's financial constraints, which is a key component of tailored suggestions.

Transitioning to a Full Recommendation System To transform the price prediction model into a comprehensive recommendation system, we can take the following steps:

- **User Profile Creation:** Collect user preferences (property style, location, budget) and specific requirements (e.g., number of bedrooms).
- **Similarity Calculation:** Compare properties based on predicted price and other features (content-based filtering).
- **Property Ranking:** Rank properties based on alignment with user preferences and budget.
- **User Feedback Loop:** Gather user feedback to refine the recommendations.

4.5 Conclusion The content-based filtering model demonstrates promising results in predicting property prices based on the available features. This is crucial for the Personalised Property Recommendation System, as it allows for tailored suggestions based on both price and user preferences.

The evaluation metrics indicate that the model reliably predicts property prices. Future improvements include refining features, tuning the model, and incorporating additional data, which will enhance the system's accuracy and reliability. The successful implementation functions as a prototype, showing that current listings from OnTheMarket and historical transaction data from HM Land Registry may be used to construct strong models.

Accurate price prediction is a cornerstone of a comprehensive recommendation system. By integrating this with user preferences, the system can streamline property searches and cater to the needs

of buyers, especially first-time buyers. This highlights the potential of machine learning in the real estate domain, paving the way for a more personalised and efficient property search experience.

4.3 5. References and Resources

4.3.1 5.1 References

- [1]: HM Land Registry. (2024). “Price Paid Data.” Retrieved from <https://www.gov.uk/government/statistical-data-sets/price-paid-data-downloads>.
- [2]: UK Government. (2024). “About the Price Paid Data.” Retrieved from <https://www.gov.uk/guidance/about-the-price-paid-data>.
- [3]: OnTheMarket.com. (2024). Property Listings. Retrieved from <https://www.onthemarket.com/>.
- [4]: Open Government Licence (OGL) v3.0. Retrieved from <https://www.nationalarchives.gov.uk/doc/open-government-licence/version/3/>.
- [5]: HM Land Registry Price Paid Data License. “Open Government Licence for public sector information.” Available at: <https://use-land-property-data.service.gov.uk/datasets/ccod/licence/view>.
- [7]: Chollet, F. (2018). “Deep Learning with Python.” Manning Publications. Available at: <https://www.manning.com/books/deep-learning-with-python>.

4.3.2 5.2 Resources Used

Web Scraping and Data Collection libraries

- Python programming language: <https://www.python.org/>
- BeautifulSoup library for Python: <https://www.crummy.com/software/BeautifulSoup/>
- Pandas library for data manipulation: <https://pandas.pydata.org/>
- Requests library for HTTP requests in Python: <https://docs.python-requests.org/en/master/>

Data Processing and Analysis

- Jupyter Notebooks for interactive computing: <https://jupyter.org/>
- Folium library for map visualization: <https://python-visualization.github.io/folium/>
- Geopy library for geocoding: <https://geopy.readthedocs.io/>
- Nominatim and ArcGIS for Geocoding: Utilised for converting addresses into geographic coordinates. Nominatim and ArcGIS

Ethical Considerations

- Ethical guidelines for web scraping and data usage were followed as per sources’ terms and conditions.
- Data Privacy and Anonymization: Data handling processes ensured no personal data was exposed or misused.
- Adherence to the Robots Exclusion Protocol as per:
- “Robots.txt” on Wikipedia: <https://en.wikipedia.org/wiki/Robots.txt>
- “Formalizing the Robots Exclusion Protocol Specification” by Google: <https://developers.google.com/search/blog/2019/07/rep-id>

4.3.3 5.3 Acknowledgements

- HM Land Registry for providing open access to Price Paid Data under the OGL: “Contains HM Land Registry data © Crown copyright and database right 2021. This data is licensed under the Open Government Licence v3.0.”
- OnTheMarket.com for the property listings data used in the scraping part of the prototype, adhering to their scraping guidelines and robots.txt file.