基於實價登錄資料之土地價值預測

何彥南 國立政治大學、莊崴宇 國立政治大學、周捷因 國立台灣科技大學

摘要

本研究將建立一套可用於多樣土地資訊預測的深度學習模型,使用 GMAN (Graph Multi-Attention Network) 作為主要預測模型,透過實價登錄中土地的交易資訊進行土地價值的預測與分析,解決土地價值過往因人為判斷而有主觀差異的問題,使土地價值的資料能有更明顯的評斷標準,進而作為土地利用流程中多種指標的考量之一。

在模型的設計中,我們展現深度學習的靈活性,使模型能投入更多資訊與數據,可與基於參考點的「時間」與「空間」資訊作為輸入預測目標土地之價值,並進行長時段的預測(long term forecast)以提升整體預測資料的準確度與可信度。此外,由於土地價值指標大多皆為時間序列性資料,且包含了地區差異的空間資訊,本研究所使用的方法亦可靈活套用於下同的土地價值指標,如:各地區不同時間下的平均每坪土地的交易價格、土地價格的中位數、變異數等。

我們期望能透過本研究所開發的預測模型,將多種各地區於不同時段的預測土地價值指標,結合土地利用的決定流程,產出對於不同區域與時段下,最適合該土地的利用選擇與開發順序,輔助資產活化決策,並期望能將利益達到最大化。

一、前言

本研究採用從大量土地交易資料,並以時 間序列角度對未來的土地價值的預測。本方法 使用 GMAN (Graph Multi-Attention Network) 為 主要模型,這個模型的特色是可以進行長時段 的預測 (long term forecast),也就是說它可以一 次預測多筆未來的資訊,例如以一個月為單 位,此模型可以預測一個月後、兩個月後的資 訊。此外,我們以 GMAN 模型做基底,對其結 構做調整,讓模型可以使用自行設定的參考點 作為額外的特徵輸入,參考點會與目標土地會 以一定範圍與時間內實價登錄中的土地計算平 均價值作為該點的指標,並將參考點與目標點 的歷史資訊作為輸入去預測未來目標土地價值 指標。在深度學習模型上。本研究先以平均土 地價值為出發點建立一套預測流程作為範例, 可以將其擴充到其他指標上,此指標可以代表 目標土地在一段時間內的區域狀態,像是價格 得平均、中位數、變異數或是交易數量。透過 此方法可以讓決策者依據需求去設計對應指 標,透過我的方法對指標進行預測,提升後續 資產活化決策之品質。

二、研究目的

三、文獻探討

1. 土地/房屋價值指標與分析

政府等公部門對於土地、房屋的開發與運用往往會受到諸多不同的因素所牽絆,從而影響各方決策,S. Geng 等人針對中國廣州市社會住宅原型,揭發社會住宅與城中村基本要,提會住宅原型,揭發社會住宅與城中村基本要素性實明,以促進中國的可持續住宅發展 [1],其中分成識別城中村的基本要素並闡明顯著特徵如空分成,以及綠門對於土地與住房開發,第四考量重點。

此外,A. M. Simarmata 亦透過模糊層次分析法 (FAHP),確定住房開發的土地優先級 [2],可以得知有關住房土地的開發利用、不同 指標對於土地開發的影響程度以及決策者該如 何規劃其優先次序,除了是政府密切關注的議 題,在學術領域中也是一常見且熱門的研究主 題。

而隨著房地產的價值逐漸升溫,越來越多投資客利用買賣房地產的方式進行投資,房價的波動對於土地資產與土地開發運用的影響房們波動對於土地資產與土地開發運用的影響房間波動因素的議題進行討論。M. Shao 以台灣新北市的414間房屋作為樣本,透過線性回歸模型揭示了房價、屋齡與車站距離之間的線性關係[3];Stijn Van Nieuwerburgh 和 Pierre-Olivier Weill透過房地產市場的動態均衡模型,揭露出不同地區的房價跟當地人口的人均收入有關[4];Owen Lamont 與 Jeremy C. Stein 則使用相關數據來分析屋主借貸模式與房價之間的關係[5],不難發現,房價的波動受到生活中各種不同的因素所產絆,連帶著對於土地開發利用的情況也因著各類指標而有所影響。

2. 深度學習

自大數據、人工智慧的崛起,近年來有越 來越多深度學習的應用案例, F. Wang [6] 等人 利用深度學習的方式對房屋價格進行預測,過 去並不常見將房價作為土地開發的考量因素, 原因在於房價的評斷往往會因為估價師的主觀 因素、個人評斷與認知落差而產生偏頗,然而 將房價使用深度學習的模式進行預估,不僅可 以讓預估出來的房價維持在一個較客觀的平均 水平上,更可以透過所建立的模型使用房價作 為對土地的開發運用進行預測的一個項目;此 外,由於房價屬於時間序列性資料,在加上近 幾年於學術領域中有許多與時間序列性相關的 模型出現,如 H. Xu 和 A. Gade [7] 就曾利用結 構化深度神經網絡對房地產進行相關評估、Y. Chen, R. Xue 及 Y. Zhang 亦透過多種機器學習 與深度學習的方法對房價進行預測,實驗結果 亦顯示出使用如貝葉式、SVM 等機器學習方法 對於房價的預測有顯著的成效 [8]。

本專案將使用 GMAN 模型對進行建模預測,過去許多 GMAN 相關的應用多是與交通狀況與風險等預測相關,如 J. Hu 和 L. Chen 提出了一種基於 Multi Attention 的時空圖卷積網絡(MASTGCN)來預測道路網絡上不同位置的長期交通狀況 [9]; Z. Wenjing 和 Z. Gang [10]則使用 GMAN 模型對大豆價格波動進行風險預測,從時間和空間維度對商品的價格波動進行建模,最終結果顯示出使用模型能有效預測大豆價格波動、影響價格波動之風險評估和預警。

由次可見,使用深度學習方式進行土地開發的預測將有助於使效益極大化,而 GMAN 模型結合了時間與空間兩個序列性資料以及比實用的實價登錄屬時間序列性資料以及比賽個所提供係屬於空間序列性之地區資料以及出有單位所提供係屬於空間序列性之地區資料對人類不可,因此本專案將嘗試以 GMAN 模型針對房價進行,,並期望能夠將此方法應用至更多所之一,對望能有預測結果更為精準客觀的同時,成功將利益達到最大化。

四、研究方法

1. 模型介紹

本專案採用 GMAN 模型進行預測,其中 GMAN 模型主要有兩大特點,分別為 Encoder -Decoder 技術和結合 Spatio-Temporal Embedding 的 Multi-Attention 值運算,以下分別針對兩個部 分進行講解。

1.1 Encoder - Decoder

當輸入的時間與預測的時間長度不同時, Encoder - Decoder 的架構可以讓模型透過中間的 TransAtt 更好的轉換(如圖 1),不會因為 input 和 output 的序列長度不同導致模型預測時的偏 差。而同樣的技術也用在 Seq2seq,可以用在翻 譯上,因為在翻譯任務上常常會遇到輸入的字 數和翻譯目標的字數不一樣。

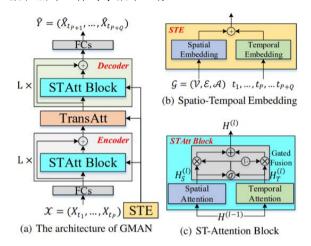


圖 1. GMAN 模型架構圖

1.2 Multi-Attention and Spatio-Temporal Embedding

深度學習中 Attention 機制是一個重要的技術,最早是由 Ashish、Vaswani 等人發表的 "Attention is all you need "中所提出,在隱藏層上再加一層可以讓模型決定要給予隱藏層中的哪個輸出更高的權重,透過疊加好幾層模型可以在不同層中學習到要專注在哪個特徵上藉此優化預測 [11]。而 Multi-Attention 即是將多個Attention 並行在一起處理,最後 concat 在一起做預測,雖然會增加訓練時的參數和時間,但是對模型會更好去學習更深一層的資訊。

STE 的部分由 SE 和 TE 組成,SE 部分是使用站點間直線距離作為 edge 的值,並使用史丹佛提出的 node2vec 將圖的資訊轉換成每個點固定維度的特徵序列,方便後續進入到深度學習做訓練;而 TE 的部分則是將一天分成 T 個時間區間,並將每筆時間序列資料的時間部分轉換成在一天中的哪個時間區間,還有在一個禮拜中的星期幾,將此資訊用 one-hot 的方式編碼,並 concat 作為 TE。

GMAN 模型則是利用此特性同時對輸入的存量資料與 STE 去訓練,並利用 Attention 機制讓模型訓練時能對不同的時間和不同的地理關係中找到重要的部分去訓練(如圖 2)。

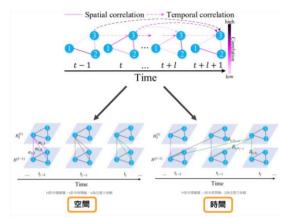


圖 2. GMAN 模型中地理與時間上的注意力機 制示意圖

2. 資料集

2.1 主辦單位提供之土地資料

本研究資料為此次比賽之主辦單位所提供的土地資料,以下簡稱「目標土地」。目標土地包含了:桃園市 56 筆、新北市 37 筆。其中資料內容則含有: 重劃區名稱、鄉鎮市區、地段、地號、土地面積 (m2)、使用分區等欄位。然因本研究會需要爬取大量實價登錄資料,因此,我們會先以桃園市資料作為實作範例,並期望未來在時間與設備允許的情況下,有更多的機會能夠將此研究開發出來的流程套用至更多不同的資料。

2.2 歷年實價登錄資料

本專案使用「內政部不動產成交案件實際資訊資料供應系統」」(如圖 3),此網站包含了全國的不動產買賣、預售屋買賣、不動產租賃等歷年實價登錄的資料。其中,我們將 101 年~111 年第 2 季的資料,批次一季一季的壓縮資料下載下來,經過檢視,我們發現每一季的壓縮檔不以縣市英文代號開頭,而各縣市內的資料皆包含各縣市房地產交易 (a)、預售屋交易(b)、房地產租賃 (c) 這三類資料;再細檢視後發現,這三類資料都有主要檔案和建築 (build)、土地 (land) 和停車位 (park) 等其他附加資訊。

而本次研究,我們主要使用桃園市的房地 產交易資料,並擷取土地交易相關的部分,附 加資訊則未採納使用。



圖 3. 內政部不動產成交案件實際資訊資料供應 系統示意圖

表 1. 實價登錄資料使用欄位與範例

欄位	範例
鄉鎮市區	蘆竹區
交易標的	土地
土地位置建物門牌	中興段 1064 地號
土地移轉總面積平方 公尺	1179
都市土地使用分區	曹辰
非都市土地使用分區	特定農業區
非都市土地使用編定	農牧用地
交易年月日	1010709
交易筆棟數	土地1建物0車位0
總價元	14800000
單價元平方公尺	12553

3

.

¹ https://plvr.land.moi.gov.tw/DownloadOpenData

3. 資料爬取與處理

3.1 資料前處理

我們擷取了近 10 年內於桃園市之所有實價登錄資料 (約 45 萬筆)。其中,由於比賽的目標為「土地」相關,因此此部分只取用「交易標的」為土地的資料,亦即屬於純土地的交易資料 (約 10 萬筆),而非房屋買賣中所附帶之土地交易資料。

於土地識別上則是透過「鄉鎮市區」與 「土地位置建物門牌」兩欄位進行區分是否屬 同一筆土地。在純土地的交易中「土地位置建 物門牌」即為地號,而地號是由地段與土地 號組成,其中段可能會再細分成小段。此部分 透過 regex(正規表達式,Regular Expression) 的方式切分出地段與地號,以供後續查詢土地 位置與新舊地號使用。

此外,由於我們有發現地段上有亂碼問題,經查看過後發現僅限於某些地段上,因此這部分我們選擇直接以手動方式觀察並抓出有亂碼的地段,再將其替換成正常的顯示方式(參考「地籍圖資網路便民服務系統」上的最新資料)。

至於在重複值的處理上,以「鄉鎮市區」、「土地位置建物門牌」、「交易年月日」為基準去掉於同一天重複的交易項目(約1萬筆),而根據觀察大部分資料皆無法與內政部的資料年度匹配,判定係屬於系統上的 bug 導致重複值,所以有關重複值的部分,我們一律不採用於後續的模型開發流程。

3.2 土地爬取經緯度

我們的方法需要知道實價登錄與目標土地 的經緯度位置,才確定哪些實價登錄資料在預 測目標土地有效範圍內。本身實價登錄資料只有包含土地地號,像是「桃園市 楊梅區草湳坡段埔心小段 64-32 地號」,我們必須透過一些查詢服務去查詢對應的經緯度。然而,十年內桃園的實價登錄資料中有交易的土地總數高達6 萬多筆土地,所以有關此部分我們採取python 爬蟲的方式去批量處理。



圖 4. 土地查詢服務

首先我們找到了3個服務,「國土測繪圖 資圖資服務雲2」、「地籍圖資網路便民服務 ³」、「地號 GeoJSON API⁴」。前兩者為政府單 位提供的資料,「國土測繪圖資圖資服務雲」 目前還未尋找到有效的爬取手段,因此先做保 留,而「地籍圖資網路便民服務」的查詢結果 中也沒提供經緯度資訊,經研究可以查詢自動 定位到目標土地後再拖拉地圖,間接獲取地圖 中心的經緯度定位,但是整體流程爬取一筆需 要 10~15 秒,初估 6 萬筆需要 10 天,不符合時 間成本,且在伺服器端短時間內太頻繁查詢會 擋 IP, 雖然可以透過 tor 換 IP 的方式處理,但 是開發時間不足,因此該方向也暫時保留,我 們期望於時間充足的情況下再使用此方法。至 於「地號 GeoJSON API」為民間友人自行架設 的 API 服務,每筆時間只需要1秒左右,速度 快且穩定,但其缺點在於資料只更新至2015 年, 導致有些新增的土地可能會查詢不到, 且 有些土地經緯度似乎不太準確。

為了驗證經緯的準確性,這邊將目標的 56個土地分別到「地籍圖資網路便民服務」、「地號 GeoJSON API」爬取經緯度,並且再用經緯度到 google map 做交叉比對,結果是「地籍圖資網路便民服務」獲取的經緯度較為精準,「地號 GeoJSON API」在有些路段會有明顯偏差。為了更直觀看到兩者的經緯度偏差,這邊使用 geopy 套件底下的 distance.geodesic() 函式計算兩經緯度在實際地面上的直線距離。如圖 5 所示,可以看到大部份的偏差小可將其

² https://maps.nlsc.gov.tw/

³ https://easymap.land.moi.gov.tw/Home

⁴ https://twland.ronny.tw/

直接忽略,然而於「桃園區中埔段」和「蘆竹區大興段」兩者的偏差都達到 800 公尺以上,所以屬於兩路段的目標土地將不會作為本次實作的範例。

在服務來源的選擇上實價登錄資料是選擇用「地號 GeoJSON API」可以在短時間內蒐集到需要的資料,而目標土地則是以「地籍圖資網路便民服務」為主,因為筆數少,所以精準度較為準確且可信度相對較高。

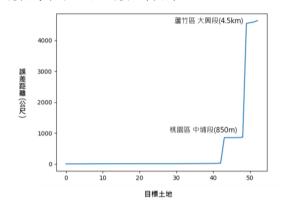


圖 5.「地籍圖資網路便民服務」與「地號 GeoJSON API」經緯度距離偏差

3.3 舊地號轉新地號

在舊地號查詢新地號的部分是使用「桃園地政資訊服務網」中的「新舊地號查詢⁵」服務。主要爬取不到的地段地號,有一部分屬於舊地號,所以需要獲取新地號後,才能查詢的到經緯度。這部分也是使用 python 的 selenium等工具進行爬蟲,可以自動處理批量資料。在這部分總共轉換 1,000 筆左右的土地資訊。



圖 6. 「地籍圖資網路便民服務、地號 GeoJSON API」經緯度距離偏

3.4 土地使用的分類與篩選

在實價登錄資料中的土地使用分區欄位包含「都市土地使用分區」、「非都市土地使用分區」、「非都市土地使用 分區」、「非都市土地使用編定」,這三個欄位共組成277種土地利用組合,而在主要的56 筆土地資料中僅有5種類,大多數為住宅與商業區(圖7)。我們必須找到與目標土地對應的土地使用類型作為區域指標的參考。

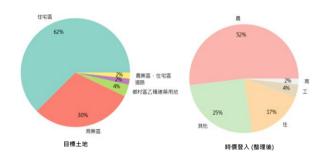


圖 7. 目標土地與實價登錄的使用分區比例

可以看到圖,經過關鍵字的篩選後,我們將實價登錄的交易資料分成農、工、商、住和其他這五種類型,發現農業最多比例佔了一半以上,然而,由於在土地價值上農業用地與商業和住宅的差異性較大,故本次專案不採用農業相關的土地資料,這同時也導致許多實價登錄資料無法使用。這邊主要針對商業區和住宅區,其中商業相關的有1,375 筆,住宅相關的有13,654 筆。以數量來看本專案會優先以資料較為充足住宅區為主要範例。

3.5 資料量與時間區間選擇

在後續計算目標指標上,需要資足夠的資料量,才具有參考性。這部分的目標是希望透過視覺化分析找出適合的資料時間範圍。由圖可以看到2012年6月之前與2022年5月之後的資料起伏大,且一開始的資料量有點不足(<200),所以本研究會以2012年6月~2022年5月的資料為主,這部分共移除494筆交易資料。

⁵

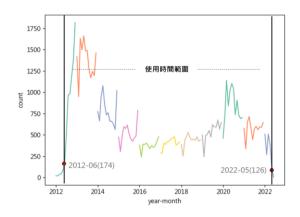


圖 8. 各年月交易量與使用時間範圍

4. 方法

4.1 目標土地的分群

目標的 56 筆土地中,有許多土地屬為鄰近土地,為了減少重複實驗的時間成本,會先將目標土地相近的分成一群,並且以該群的中心經緯度為代表的目標點。分群的方式是使用DBSCAN 並用經位度間的實際距離作為分群的依據。實際距離也是使用 geopy 套件底下distance.geodesic() 計算,eps 以公尺為單位。下圖為 eps (100、300、500) 的分結果,紅色圈為相鄰區域有多個不同群,黑色虛線為有分好的區域。可以看到新屋區、桃園區、流營區的目標土地分群結果較穩定,中壢區、蘆竹區的目標土地則需要更大的容許範圍(esp) 才能分成一群。

最後使用 esp=500 為分群依據,將目標土 地分成 11 群目標,各群的經緯度平均後,可以 獲取代表整個群體的中心點,以下稱為「目標 點」。

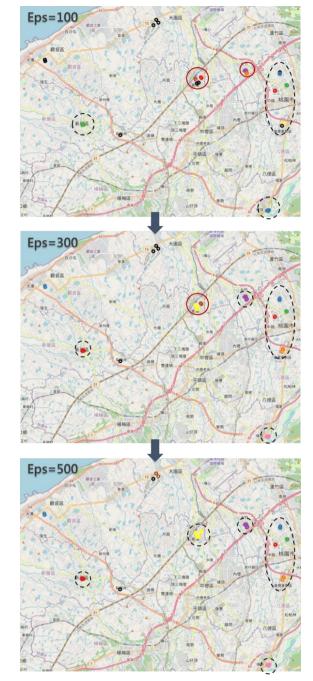


圖 9. DBSCAN 分群結果

4.2 設置各群的參考點

這部分是為了用設置參考點的方式,增加輸入的資訊量,幫助預測。且在可以在模型練完後透過抽取模型最後一層的 Attention 權重,對參考點與目標點進一步分析。

參考點為目標點周圍自定義點,可輔助目標點的預測或是分析與目標點之間的關係。本研究是採取上下左右四個方向向外推 300 公尺的距離作為參考點,下圖為目標點與參考點示意圖。



圖 10. 參考點示意圖

4.3 計算目標點、參考點與交易土地距離

這部分獲取目標點、參考點到所有交易土地經緯距離,這部分實際距離也是使用 geopy 套件底下 distance.geodesic() 計算。如下表會將資料整理成距離矩陣,方便後續篩選土地。(非真實資料)

表 2. 目標點、參考點到各土地距離矩陣

	目標點	參考點 1	~	參考點 n
交易 土地 1	20	50		100
交易 土地2	280	350		500
•••				

4.4 篩選目標經緯度一定範圍內的交易點

以圖 10 和表 2 為例,指標範圍為 300 公尺。如表 3,這步驟會篩選出 300 公尺內所有 交易土地點。

表 3. 篩選 300 公尺內的距離

	目標點	參考點 1	~	參考點 n
交易 土地 1	<	>		~
交易 土地 2	<			
•••				

4.5 計算各月的目標點和參考點的區域指標

從交易資料篩選出在範圍內的土地交易紀錄,並計算每個月的區域指標,如下表。(非真實資料)

表 4. 計算區域指標

	目標點	参考點 1	~ ~	參考點 n	
2012.7	50	100		80	
2012.8	30	60		70	
2022.4	100	50		120	

4.6 使用 GMAN 模型預測未來區域指標

會先將表 4 的時間-地點資料做轉換,這邊使用 window rolling 方式將時間序列資料轉換成輸入特徵和標籤,如圖 11 {t1,t2,...,t8} 代表一連串隨著時間變化的數值,以 window rolling = 3 為例,可以看到對應的輸入資料 (input) 和預測目標 (label)。共可以分成 5 筆資料,其中一部分作訓練一部分作測試。第六筆則是預測未來未知的資料。

	I	Label		
	History1	History2	History3	Target
Train -	t1	t2	t3	t4
	t2	t3	t4	t5
	t3	t4	t5	t6
Test -	t4	t5	t6	t7
	t5	t6	t7	t8
Predict -	t6	t7	t8	unknow

圖 11. 訓練資料格式 (window_size=3)

接著就會將資料輸入到 GMAN 模型中進行,圖 12 為可以看到 window_size=3 資料輸入範例,原始 GMAN 模型輸出的會是目標點與參考點的預測結果,也就是途中 Output 的部分。在最後面進行以下客製化調整:

- 1. 讓模型多一層注意力機制 (Attention), 可以將資訊收斂到目標點的預測上。
- 2. 透過注意力機制 (Attention),其訓練後 的權重可以觀察參考點與目標點之間的 關係。
- 3. GMAN 模型中 Encoder-Decoder 的架構可以讓模型透過中間的 TransAtt 做時間長度上的非線性轉換,讓我們可以一次預測一個以上的時間點資訊 (long term forecast)。

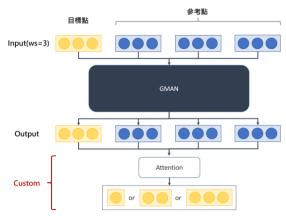


圖 12. 訓練過程與客製化設計

五、預期 AI 模組效益與成果

本研究的目標為結合最新的深度學習技術與歷年實價登錄資料,建立一套基於時間序列與地理資訊的預測分析流程,可以根據需求自由選擇區域指標與參考點,並可以用來預測土地價值、一段時間內的變化性。此外,還可以也透過 Attention機制,去探討目標點與參考點之間的關係程度,為土地的使用方向提供更多有價值的資訊,以增進資產活化的效益。

六、程式碼與相關資料

- 程式碼的部分可以到 github 取得。
 連結如下:
 https://github.com/aaron1aaron2/Land-price-predict-based-on-transaction-records
- 資料可以到 google drive 下載。連結如下:

https://drive.google.com/file/d/1Q-s6My4sy-_LWMJ7JV6ESgQeHUPBGyRr/view?usp=sh aring

註:程式碼的部分目前僅完成資料爬蟲與 前處理的部分。

参考文獻

- [1] S. Geng, H. -W. Chau and S. Yan, (2020), "Identify and Elucidating Urban Village Essentials Through Remodeling and Visualising a Social Housing Prototype in Guangzhou for Sustainable Residential Development in China," 2020 24th International Conference Information Visualisation (IV), pp. 609-613
- [2] A. M. Simarmata, Yennimar, A. M. Husien, M. Harahap and S. Aisyah, (2019), "Determining land priority for Housing Development: using Fuzzy AHP application," 2019 International Conference of Computer Science and Information Technology (ICoSNIKOM), pp. 1-5
- [3] M. Shao, (2022), "Factors Affect the House Price," 2022 14th International Conference

- on Computer Research and Development (ICCRD), pp. 136-139
- [4] Stijn Van Nieuwerburgh, Pierre-Olivier Weill, Why Has House Price Dispersion Gone Up?, The Review of Economic Studies, Volume 77, Issue 4, October 2010, Pages 1567–1606
- [5] Lamont, O., & Stein, J. C. (1997). Leverage and house-price dynamics in US cities.
- [6] F. Wang, Y. Zou, H. Zhang and H. Shi, (2019), "House Price Prediction Approach based on Deep Learning and ARIMA Model," 2019 IEEE 7th International Conference on Computer Science and Network Technology (ICCSNT), pp. 303-307
- [7] H. Xu and A. Gade, (2017), "Smart real estate assessments using structured deep neural networks," 2017 IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computed, Scalable Computing & Communications, Cloud & Big Data Computing, Internet of People and Smart City Innovation (SmartWorld/SCALCOM/UIC/ATC/CBDCom/IOP/SCI), pp. 1-7
- [8] Y. Chen, R. Xue and Y. Zhang, (2021), "House price prediction based on machine learning and deep learning methods," 2021 International Conference on Electronic Information Engineering and Computer Science (EIECS), pp. 699-702
- [9] Zheng, Chuanpan, et al. (2020). "Gman: A graph multi-attention network for traffic prediction." Proceedings of the AAAI conference on artificial intelligence. Vol. 34. No. 01. 2020.
- [10] Z. Wenjing and Z. Gang, (2021), "Temporal and Spatial Attention Network Model Based Evolution Model for Bulk Commodity Price Fluctuation Risk," 2021 IEEE International Conference on Big Data (Big Data), pp. 3284-3289
- [11] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in neural information processing systems, 30.