# Sparse Real Estate Ranking with Online User Reviews and Offline Moving Behaviors

Yanjie Fu[‡], Yong Ge[＋], Yu Zheng[◇], Zijun Yao[‡], Yanchi Liu[†], Hui Xiong[‡*], Nicholas Jing Yuan[◇]

[‡]Rutgers University, Email: {yanjie.fu, zijun,yao, hxiong}@rutgers.edu
[＋]University of North Carolina at Charlotte, Email: yong.ge@uncc.edu
[◇]Microsoft Research, Email: {yuzheng, nicholas.yuan}@microsoft.com
[†]New Jersey Institute of Technology, Email: yl473@njit.edu

*Abstract*—**Ranking residential real estates based on investment values can provide decision making support for home buyers and thus plays an important role in estate marketplace. In this paper, we aim to develop methods for ranking estates based on investment values by mining users opinions about estates from online user reviews and offline moving behaviors (e.g., taxi traces, smart card transactions, check-ins). While a variety of features could be extracted from these data, these features are intercorrelated and redundant. Thus, selecting good features and integrating the feature selection into the fitting of a ranking model are essential. To this end, in this paper, we first strategically mine the fine-grained discriminative features from user reviews and moving behaviors, and then propose a probabilistic sparse pairwise ranking method for estates. Specifically, we first extract the explicit features from online user reviews which express users opinions about point of interests (POIs) near an estate. We also mine the implicit features from offline moving behaviors from multiple perspectives (e.g., direction, volume, velocity, heterogeneity, topic, popularity, etc.). Then we learn an estate ranking predictor by combining a pairwise ranking objective and a sparsity regularization in a unified probabilistic framework. And we develop an effective solution for the optimization problem. Finally, we conduct a comprehensive performance evaluation with real world estate related data, and the experimental results demonstrate the competitive performance of both features and the proposed model.**

*Keywords*—*Real Estate, Sparse Ranking, Online User Reviews, Offline Moving Behaviors*

## I. INTRODUCTION

There are several definitions of estate value according to International Valuation Standards [1]. For instance, market value is defined as the price at which an estate would trade in a competitive Walrasian auction setting. Another example is investment value, which is the value of an estate to one particular investor and may or may not be higher than the market value of the estate. Difference between the investment value and the market value for a particular estate provides the motivation for buyers or sellers to enter the estate marketplace. Thus, providing a ranking of estates based on investment values will greatly help buyers make their purchase decisions.

Which estates have high investment values? While estate industry professionals have used different housing indexes (e.g., price-rent ratio) to approximate the fundamental value of estates, researchers have also used financial time series analysis to investigate the trend, periodicity and volatility of estate prices and assess estate investment potentials [1], [2]. Recent studies have tried to correlate the estate value to the static statistics of urban infrastructure (e.g., the numbers of POIs, the distances to bus stops), because they explicitly reflect the physical facilities of a neighborhood [3], [4]. However, infrastructure statistics is not sufficient for evaluating investment values of estates. Considering the distance to public transit, while an estate near public transit usually leads to high rent and sale price in many cities, there is also possible negative effect when living nearby public transit. For example, the noise and pollution associated with train/bus systems can lower the value of an estate as reported in [5]–[7]. Thus, there is some limitation for using these infrastructure statistics. Moreover, these statistics are often lack of dynamics and hardly reflect the changing pulses of a city.

On the contrary, there are more estate-related dynamic and information-rich data which has been accumulated with the development of mobile, internet and sensor technologies. For example, people may post comments and ratings for POIs (e.g., schools, restaurants and shopping centers, etc.) via mobile apps after their consumptions. Also, the mobility data, such as smart card transactions and taxi GPS traces, comprise both trajectories and consumption records of residents' daily commutes. People's check-ins may reflect the popularity of POIs. If properly analyzed, these data (e.g., user reviews, location traces, smart card transactions, check-ins, etc.) can be a rich source of intelligence for discovering estates of high investment-value.

Indeed, these estate-related dynamic data generated by users could better reflect investment values of estates than urban infrastructure statistics. Generally speaking, if people have better opinions for an estate, the demand for this estate is higher and its investment value will be higher. The challenge is how to uncover people's opinions for an estate. In fact, the opinions of users for an estate can be mined from (1) online user reviews and (2) offline moving behaviors. Specifically, the online reviews (e.g., Zagat/Yelp ratings) contain the explicit opinions for places surrounding an estate. For example, the quality of neighborhood can be partially approximated by the ratings of business venues, such as overall rating, service rating, environment rating, etc. Meanwhile, the offline moving behaviors near an estate not only encode the static statistics of urban infrastructure, but also reflect the implicit "opinions" of residents for a neighborhood. For example, the arriving, transition, and leaving volumes of taxies and buses imply the mobility density of a neighborhood; the average velocity of

---

[*]Contact author.
[1]http://www.ivsc.org/

taxies and buses indicates the degree of traffic congestion or accessibility; the daily frequency of check-ins shows regional popularity and prosperity; the heterogeneity of distributions of check-ins over categories reflects if the facility planning is balanced or not. All these indications by the estate-related dynamic user-generated data comprise the important facets of an estate that home buyers care very much and convey the implicit "opinions" of users for a neighborhood. Therefore, we consider and mine both the explicit opinions from user reviews and the implicit opinions from moving behaviors to enhance the evaluation of estate investment value.

Although we may extract a lot of features from the variety of data sources, these extracted estate-related features usually are correlated and redundant. The feature redundancy results in poor generalization performance. In reality, a small number of good features can determine the ranking of estates based on investment values. Therefore, we explore the sparse learning technique for the ranking of estates. However, classic sparse learning methods use a two-step paradigm, which is basically to first select a feature subset and then learn a ranking model based on the selected features. But the selected feature subset may not be optimal for ranking because the two steps are modelled separately. In contrast, combining sparsity and ranking in a unified model can help to identify the optimal feature subset for better learning an estate ranker, and also have less computational cost in prediction.

Along this line, in this paper, we propose to mine opinions of mobile users and explore the learning-to-rank with sparsity for the investment value based estate ranking. We consider and explore both explicit and implicit opinions that reflect estate investment value by mining online user reviews and offline moving behaviors. Specifically, to capture the opinions of mobile users toward estates, we extract the explicit features from user reviews to reveal user satisfaction of estate neighborhoods. Besides, we measure the traffic volumes with respect to different directions, traffic velocity, functionality heterogeneity, neighborhood popularity, topical profile of estate neighborhoods by mining multi-type mobility data including taxi traces, smart card transactions and check-ins. Moreover, we learn a linear ranking predictor by combining pairwise ranking objective and sparsity regularization in a unified probabilistic framework, which is greatly enhanced by simultaneously conducting feature selection and maximizing estate ranking accuracy. Finally, we conduct comprehensive performance evaluations for the feature sets and models with large-scale real world data and the experimental results demonstrate the competitive performance of our method with respect to different validation metrics.

## II. SPARSE ESTATE RANKING

In this section, we present the proposed system of sparse estate ranking, namely SEK.

### A. The Overview of Sparse Estate Ranking

As shown in Figure 1, our estate ranking system consists of two major components: (1) estate feature extraction and (2) sparse estate ranking.

**Estate Feature Extractions:** As shown in Figure 1, we first collect historical prices of each estate, compute the return
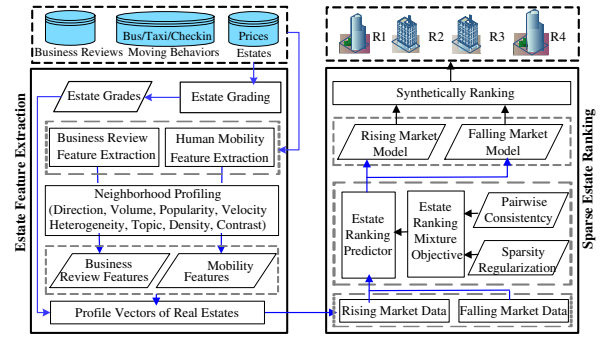


Fig. 1. The framework of the proposed system.

rates [2] of estates and grade estates into five bins/levels in terms of investment returns to prepare labels for training data. The discretization of the estate returns is important because the small difference between estate values in the same value category might be noisy for the ranking model.
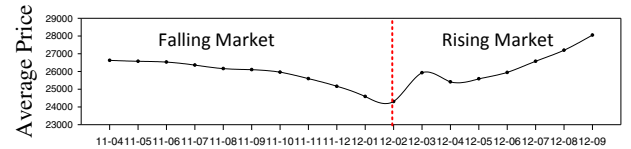


Fig. 2. The rising market period and the falling market period in Beijing.

Specifically, we first calculate the average estate price of a city for each month. For instance, Figure 2 shows the trend of the average estate prices in Beijing. We can see an inflection point in the curve. The point is used to split the time period into two phases, i.e., the rising phrase (from Feb. 2012 to Sept. 2012) and the falling phrase (from Apr. 2011 to Feb. 2012). We then sort estates in rising phase and falling phrase according to their investment returns in the decreasing order as shown in Figures 3 (a) and (d), where the horizontal axis is the order of an estate in the sorted list and the vertical axis represents return rates. As can be seen, the prices of a small number of estates significantly increase or decrease whereas many estates' prices remain stable. In fact, these distributions indicate the power law distribution for estate investment returns. After computing the second order derivatives of these two curves, we find out four inflection points, which show the significant change of return rates as shown in Figures 3 (b) and (e). As a result, we obtain five rating levels for the rising and falling phrases as shown in Figures 3 (c) and (f).

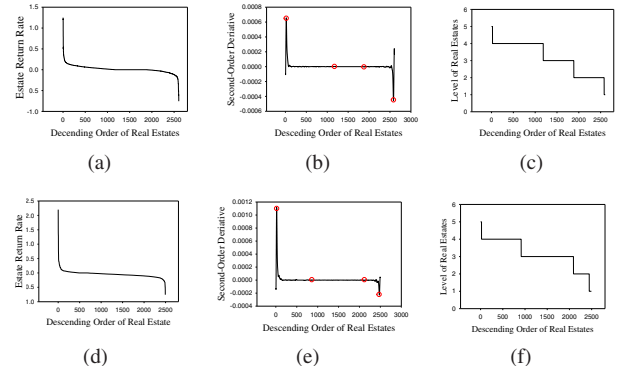

(a)    (b)    (c)

(d)    (e)    (f)

Fig. 3. The grading process of estates.

Next we aim at extracting the features from online user reviews and offline moving behaviors such as taxi traces, smart

---

[2]http://financial-dictionary.thefreedictionary.com/rate+of+return

card transactions, check-ins as shown in Table I. The features from user reviews are summarized by spatial statistics and the features from moving behaviors are derived from multiple angles (e.g., direction, volume, velocity, heterogeneity, topic, contrast, popularity).

TABLE I.    THE EXTRACTED FEATURES.

| Online User Reviews | Offline Moving Behaviors | | |
|---|---|---|---|
| **User Reviews** | **Taxi** | **Bus** | **Check-in** |
| Overall Salification | Arriving Volume | Arriving Volume | Popularity |
| Service Quality | Leaving Volume | Leaving Volume | Topic |
| Environment Class | Transition Volume | Transition Volume | |
| Consumption Cost | Driving Velocity | Bus Stop Density | |
| Functionality Planning | Commute Distance | Smart Card Balance | |

**Sparse Estate Ranking:** We learn a linear ranking predictor by combining a pairwise ranking objective and a sparsity regularization together. By optimizing the overall objective function, we learn the estate ranker by simultaneously conducting feature selection and maximizing ranking accuracy. Two separated models are then built to infer the value-adding and value-protecting ability of an estate in a rising and a falling market respectively. Given a set of estates specified by a user, we extract the features in the same way as we show in Figure 1. Since we do not know whether the market will go up or down, the extracted features are fed into two ranking models respectively to produce the potential ranks of these estates at the current time. Finally, we generate a final score for an estate by aggregating the ranking outputs of these two models.

### B. Estate Feature Extraction

Rather than simply considering the static statistics of urban infrastructure (e.g., the numbers of POIs, the distances to bus stops), we introduce the fine-grained features we have extracted from online users reviews and offline moving behaviors for estate ranking.

*1) Explicit Features from Online User Reviews:* Both prosperity and users' opinion of neighborhood are two important factors determining property investment value. Recent study [8] shows that a strong regional economy usually indicates high housing demand. [9] further points out the word-of-mouth reflects the satisfaction of people toward the quality of a neighborhood. We thus consider to mine the online user reviews of Beijing collected from www.dianping.com. More specifically, for each estate $e_i$, we measure (1) overall satisfaction, (2) service quality, (3) environment class, (4) consumption level, and (5) functionality planning of the neighborhood $r_i$ by mining the reviews of business venues located in $r_i$, $\{p : p \in P \& p \in r_i\}$ in which $P$ is the set of business venues in Beijing.

*Overall Satisfaction:* For each estate $e_i$, we access the overall satisfaction of users over the neighborhood $r_i$. Since the overall rating of a business venue $p$ represents the satisfaction of users, we extract the average of overall ratings of all business venues located in $r_i$ as a numeric score of overall satisfaction. Formally we have:

$$f_i^{OS} = \frac{\sum_{p \in P \& p \in r_i} OverallRating_p}{|\{p : p \in P \& p \in r_i\}|}. \tag{1}$$

*Service Quality:* Similarly, we compute the average of service rating of business venues in $r_i$ and represent the service quality of the neighborhood of $e_i$ by

$$f_i^{SQ} = \frac{\sum_{p \in P \& p \in r_i} ServiceRating_p}{|\{p : p \in P \& p \in r_i\}|} \tag{2}$$

*Environment Class:* The environment class of business venues could reflect whether the neighborhood is high-class or not. Therefore, we extract the average environment ratings as

$$f_i^{EC} = \frac{\sum_{p \in P \& p \in r_i} EnvironmentRating_p}{|\{p : p \in P \& p \in r_i\}|} \tag{3}$$

*Consumption Cost:* Average costs of consumption behaviors in business venues can partially reflect the salary income and neighborhood class. We calculate the average consumption cost of business venues of a targeted neighborhood as a feature.

$$f_i^{CC} = \frac{\sum_{p \in P \& p \in r_i} AverageCost_p}{|\{p : p \in P \& p \in r_i\}|} \tag{4}$$

*Functionality Planning:* A competitive neighborhood usually provides convenient access to diverse facilities, such as living demands (e.g., restaurants, supermarkets, and hospitals), education demands (e.g., schools and libraries), safety demands (e.g., police and fire department) and entertainment demands (e.g., theaters and parks), so that it meets various demands of residents. Shortage of diverse facility would reduce estate investment value. High facility diversity of a neighborhood helps to enhance the attractiveness of its estates. This effect is called mixed/diverse land use which plays an important role in metropolitan realty market. We therefore investigate the distribution of POIs over categories in each neighborhood. A high-class neighborhood is expected to provide balanced and heterogeneous categories of facilities. Hence, we apply an entropy to measure the functionality heterogeneity of a neighborhood. Let $\#(i, c)$ denotes the number of business venues of category $c \in C$ located in $r_i$, $\#(i)$ be the total number of business venues of all categories located in $r_i$. The entropy is defined as

$$f_i^{FP} = -\sum_{c \in C} \frac{\#(i, c)}{\#(i)} \times log \frac{\#(i, c)}{\#(i)} \tag{5}$$

*2) Implicit Features from Offline Moving Behaviors:* Recent study [8] reports different types of transit systems (e.g., taxi, bus) have different impacts on estate values due to their different fares, frequencies, speeds, and scopes of service. Figure 4 (a), (b) and (c) show the density distribution of three types of moving behaviors respectively (i.e., taxi, bus and check-in) in Beijing. Taxi transits are fast, expensive and mainly distributed in central business district (CBD) and financial areas. Bus transits are slow, cheap and mainly distributed in information technology (IT) and education areas. Check-ins reflect a broad range of mobility and are mainly distributed in areas full of attractions, entertainments, and POIs. Since different moving behaviors reflect different geographic preferences and social classes of mobile users, we exploit these three types of moving behaviors to uncover the implicit preference of mobile users toward a neighborhood.

**Taxi-Related Features.** Recent study [8] suggests that the ability to travel within a large metropolitan area in a short time, for example, by taxi, is highly valued by residents. To extract the taxi related features, we measure the arriving volume, leaving volume, transition volume, driving velocity and commute distance of a neighborhood using taxi GPS traces. Let $TT$ denote the set of all taxi trajectories of Beijing, each of which represents a taxi trajectory, denoted by a tuple $< p, d >$ where p is a pickup point and d is a drop-off point.

*Taxi Arriving, Leaving and Transition Volume:* According to [8], most affluent homeowners expect time-saving commute

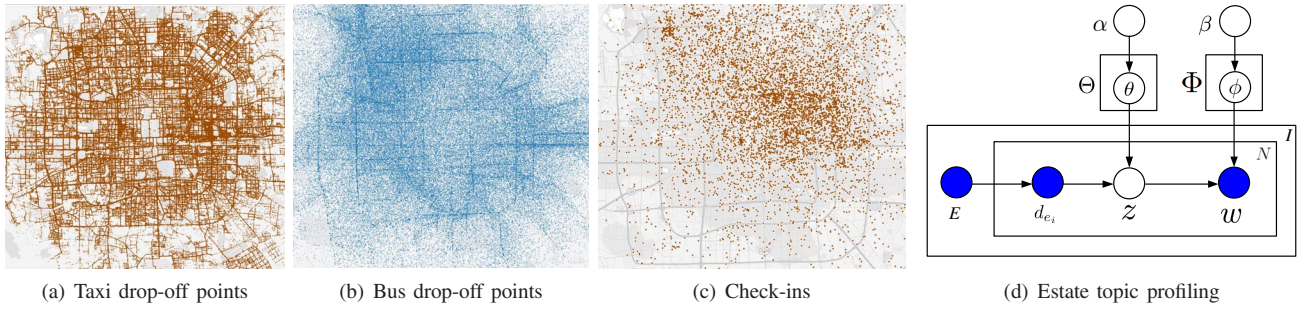(a) Taxi drop-off points     (b) Bus drop-off points     (c) Check-ins     (d) Estate topic profiling

Fig. 4. (a), (b), and (c) respectively show spatial distribution of taxi drop-offs, bus drop-offs and check-ins; (d) illustrates the process of estate topic profiling **using the associated word-of-mouth from check-ins**.

to white-collar jobs downtown and value faster taxies access. Therefore, the arriving, leaving, and transition volumes of taxi mobility reflect the income and social class of residents of the targeted neighborhood. We define a feature as the counted taxi arriving volume of external passengers toward the targeted neighborhood. Formally, the taxi arriving volume is given by

$$f_i^{TAV} = |\{<p,d> \in TT : p \notin r_i \& d \in r_i\}| \tag{6}$$

Similarly, we define a feature as the counted taxi leaving volume from the targeted neighborhood to external venues. Formally, the taxi leaving volume is defined as

$$f_i^{TLV} = |\{<p,d> \in TT : p \in r_i \& d \notin r_i\}| \tag{7}$$

We also define a feature as the taxi transition volume between different venues inside the targeted neighborhood. Formally,

$$f_i^{TTV} = |\{<p,d> \in TT : p \in r_i \& d \in r_i\}| \tag{8}$$

*Taxi Driving Velocity:* According to [8], the value of increased travel velocity and reduced traffic congestion should be reflected in home values. We investigate the average taxi velocity of the neighborhood of each estate, namely $f_i^{TDV}$. Usually, the taxi speed of a neighborhood indicates the accessibility of road network and transportation efficiency. Formally, $f_i^{TDV}$ is given by

$$f_i^{TDV} = \frac{\sum_{p \in r_i \& d \in r_i} dist(p,d)/time(p,d)}{|\{<p,d> \in TT : p \in r_i \& d \in r_i\}|} \tag{9}$$

*Taxi Commute Distance:* Taxi is a kind of expensive but fast transit. Normally, passengers take taxi to the important places (e.g., work place, theater, hotel, etc.) for business or urgent purposes. The shorter distance an estate neighbor is from important places, the more prosperous the neighborhood is, and the higher commute convenience the neighborhood has. A huge part of motivations of trading an estate comes from the incentive of convenient living environment. Formally, the taxi commute distance is defined by

$$f_i^{TCD} = \frac{\sum_{p \in r_i || d \in r_i} dist(p,d)}{|\{<p,d> \in TT : p \in r_i || d \in r_i\}|} \tag{10}$$

**Bus-Related Features.** Most of moderate-income residents choose buses which are cheaper with acceptable speed rather than taxies which are expensive with faster speed [8]. Since most of the residents in a city are middle-class, bus traffic represents the majority of urban mobility. Besides, according to [10], there is a connection between a drop in estate prices and a decreased flow of bus mobility. We thus measure the arriving, leaving and transition volumes of buses in the neighborhood of each estate. Let $BT$ denote the set of all the bus trajectories of Beijing, each of which represents a bus trajectory, denoted

by a tuple $<p,d>$ where p is a pickup bus stop and d is a drop-off bus stop.

*Bus Arriving, Leaving and Transition Volume:* Similar to taxi mobility volume, we also extract the arriving volume, leaving volume and transition volume of buses from smart card transactions. Formally,

$$\begin{aligned} f_i^{BAV} &= |\{<p,d> \in BT : p \notin r_i \& d \in r_i\}| \\ f_i^{BLV} &= |\{<p,d> \in BT : p \in r_i \& d \notin r_i\}| \\ f_i^{BTV} &= |\{<p,d> \in BT : p \in r_i \& d \in r_i\}| \end{aligned} \tag{11}$$

*Bus Stop Density:* Recent work [11] reports that price premiums of up to ten percents are estimated for estates within 300m of more bus stops. In other words, the bus stop density is positively correlated to estate prices. Here, we propose an alternative approach and strategically estimate bus stop density using smart card transactions. In smart card transactions, the ticket fare of a trajectory indeed reflects the number of bus stops in this trajectory. This is because the Beijing Public Transportation Group charges passengers according to the number of stops of each trip. Given the pick-up stop $p$ and the drop-off stop $d$, the trip distance between $p$ and $d$ is fixed in a designed bus route. Then, the ratio of trip distance to bus stop number implicitly suggests in average distance between every two consecutive bus stops. Since the bus stop number of a trip can be approximated by the fare, we compute the ratio of distance to fare for estimating the density of bus stop in a neighborhood. The smaller the distance-fare ratio is, the higher the bus stop density is.

$$f_i^{BSD} = \frac{\sum_{p \in r_i || d \in r_i} dist(p,d)/fare(p,d)}{|\{<p,d> \in BT : p \in r_i || d \in r_i\}|} \tag{12}$$

*Smart Card Balance:* The smart card balances imply the patterns of the consumption and recharge behaviors. If residences always maintain a higher balance in their smart card, this suggests the card holders spend more money on bus travel. The large expense of bus travel implies: (1) residences depend on buses more than other transportation (e.g., subway, taxi), which may indicate that the affiliated neighborhood is lack of subways and taxies; (2) residences travel a longer distance to work, shop and pick up children, and thus need to maintain a high balance. In other words, this place is remote and inconvenient. We thus consider to extract the smart card balance as a feature. Formally,

$$f_i^{SCB} = \frac{\sum_{p \in r_i || d \in r_i} balance(p,d)}{|\{<p,d> \in BT : p \in r_i || d \in r_i\}|} \tag{13}$$

**Check-in Related Features.** Mobile users check in at online location-aware social networks when they walk in an important

place. These check-ins are a significant portion of urban mobility. Estate price is likely high in communities where there are convenient transit stations with good access to retail stores and services [8]. Therefore, check-in behaviors could partially reflect the access convenience to these locations. In our data set, each check-in event can be denoted by a tuple, $< p, t, c > \in CI$, where $p$, $t$, $c$ and $CI$ represent the POI of the check-in, the check-in time stamp, the category of POI, and the set of check-in events, respectively.

*Neighborhood Popularity:* We count the total number of check-ins reported in the neighborhood of each estate as popularity measurement. Formally,

$$f_i^{NP} = |\{< p, t, c > \in CI : p \in r_i\}| \tag{14}$$

*Topic Profile:* The goal of topic distillation is to learn the topic distribution of a neighborhood based on the textual information of check-ins via a two-step approach.

*STEP1: Propagating word-of-mouth from poi to neighborhood.* In check-in data, each POI is associated with textual reviews posted by users. This textual information reflects opinion of users toward this POI. Since each neighborhood is associated with a cluster of POIs, we therefore propose to propagate the word-of-mouth of mobile users from poi to neighborhoods by spatio-textual aggregation using check-in data. We get a cluster of textual posts denoted as $d_{e_i}$ for the neighborhood of each estate $e_i$. We then segment these sentences into words and extract the semantically significant tags for each neighborhood. One reason for propagating word-of-mouth from poi to neighborhood is that the terms associated with a single POI are usually short, incomplete and ambiguous. Moreover, LDA is proven non-effective for short texts. The aggregation process can better learn thousands of mobile users' opinions toward estates in terms of latent topic distributions.

*STEP2: Textual profiling from words to topics.* Next we exploit the LDA model for estate topic profiling by treating each estate neighborhood as a document. In LDA, each document is represented as a probability distribution over topics (document-topic distribution) and each topic is represented as a probability distribution over a number of words (topic-word distribution).In this way, we build an aggregated LDA model as shown in Figure 4(d). Here, the topic distribution of each document $\Pr(z \mid d_{e_i})$ is treated as topical features of estate, where $z$ and $d_{e_i}$ are topic and document respectively. The topic profiling process of the estates is as following:

1. For each topic $z \in \{1, ..., K\}$, draw a multinomial distribution over terms, $\phi_z \sim Dir(\beta)$.
2. For the document $d_{e_i}$ given an estate $e_i$
   (a) Draw a multinomial distribution over topics, $\theta_{d_{e_i}} \sim Dir(\alpha)$
   (b) For each word $w_{d,n}$ in document $d_{e_i}$:
       i. Draw a topic $z_{d,n} \sim Mult(\theta_{d_{e_i}})$
       ii. Draw a word $w_{d,n} \sim Mult(\phi_{z_{d,n}})$

So far, we have extracted two categories of estate features as shown in Table I. We emphasize that the above features are defined in terms of the neighborhood ($r_i$) of each estate, which is parameterized by its radius $d$. Hence, we can extract multiple groups of estate features with respect to different neighborhood radius (e.g., d=0.25,0.5,0.75,1,1.25,...,3km).

## C. Sparse Pairwise Ranking for Estate Appraisal

Here we present the sparse pairwise estate ranker.

**Model Description:** Since many existing learning-to-rank algorithms use linear rankers, we learn a linear ranking predictor. Let $\boldsymbol{x_i}$ denote the M-size vector representation of estate $e_i$ with the above extracted features, $f_i$ denote the predicted estate value, and $y_i$ denote the ground truth estate value, then we have $f_i(\boldsymbol{x_i}; \boldsymbol{w}) = \boldsymbol{w}^\top \boldsymbol{x_i} + \epsilon_i = \sum_{m=1}^M w_m x_{im} + \epsilon_i$, where $\epsilon_i$ is a zero-mean Gaussian bias with variance $\sigma^2$, and $\boldsymbol{w}$ is the weights of features. In other words, $P(y_i|\boldsymbol{x_i}) = \mathcal{N}(y_i|f_i, \sigma^2) = \mathcal{N}(y_i|\boldsymbol{w}^\top \boldsymbol{x_i}, \sigma^2)$ where $\mathcal{N}$ represents normal distribution.

**Objective Function:** While these features indeed capture residents' opinions about estates to be ranked, they usually are inter-correlated and redundant. Thus possible confounders lead to poor generalization performance. To address this issue, we adopt a strategy which simultaneously conducts feature selection while maximizing estate ranking accuracy. Since pairwise ranking strategy is effective with lower complexity comparing with listwise ranking strategy, we combine a pairwise ranking objective and a sparsity regularization term in a unified probabilistic modeling framework.

Next we introduce how to derive the mixture objective of sparse pairwise estate ranking. Let us denote all parameters by $\Psi = \{\boldsymbol{w}, \boldsymbol{\beta}^2\}$ which are the parameters of estate ranker (we will introduce $\boldsymbol{\beta}^2$ in the following), the hyperparamters by $\Omega = \{a, b, \sigma^2\}$ which are the parameters of sparsity regularization, and the observed data by $\mathcal{D} = \{Y, \Pi\}$ where $Y$ and $\Pi$ are the investment values and ranks of $I$ estates respectively. For simplicity, we assume the real estates in $\mathcal{D}$ are sorted and indexed in a descending order in terms of their investment values, which compiles a descending ranks as well. In other words, $i$ is both the index and the ranking order of the given estate $x_i$ By Bayesian inference, we have the posterior probability as

$$Pr(\Psi; \mathcal{D}, \Omega) = P(\mathcal{D}|\Psi, \Omega) P(\Psi|\Omega) \tag{15}$$

First, the term $P(\mathcal{D}|\Psi, \Omega)$ is the likelihood of the observed data collection $\mathcal{D}$, which can be explained as a joint probability of both estate investment values, $P(Y|\Psi, \Omega)$, and estate ranking consistency, $P(\Pi|\Psi, \Omega)$. Here we treat the ranked list of estates as a directed graph, $G = < V, E >$, with nodes as estates and edges as pairwise ranking orders. For instance, edge $i \to h$ represents an estate $i$ is ranked higher than estate $h$. From a generative modeling angle, edge $i \to h$ is generated by our model through a likelihood function $P(i \to h)$. The more valuable estate $i$ is than estate $h$, the larger $P(i \to h)$ should be. On the contrary, the case, in which $i \to h$ but $f_i < f_h$, will punish $P(i \to h)$. Therefore,

$$P(\mathcal{D}|\Psi, \Omega) = P(Y|\Psi, \Omega) P(\Pi|\Psi, \Omega)$$
$$= \prod_{i=1}^I \mathcal{N}(y_i|f_i, \sigma^2) \prod_{i=1}^{I-1} \prod_{h=i+1}^I P(i \to h|\Psi, \Omega) \tag{16}$$

where the generative likelihood of each edge $i \to h$ is defined as Sigmoid$(f_i - f_h)$: $P(i \to h) = \frac{1}{1+exp(-(f_i-f_h))}$.

Second, the term $P(\Psi|\Omega)$ is the prior of the parameters $\Psi$. Here, we introduce a sparse weight prior distribution by modifying the commonly used Gaussian prior, such that a

different and separate variance parameter $\beta_m^2$ is assigned for each weight. Thus, $P(\boldsymbol{w}|\boldsymbol{\alpha}) = \prod_{m=1}^{M} \mathcal{N}(w_m|0, \beta_m^2)$, where $\beta_m^2$ represents the variance of corresponding parameter $w_m$ and $\boldsymbol{\beta}^2 = (\beta_1^2, ..., \beta_M^2)^\top$, each of which is treated as a random variable. Later, an Inverse Gamma prior distribution is further assigned on these hyperparameters, $P(\boldsymbol{\beta}^2|a, b) = \prod_{m=1}^{M} \text{Inverse-Gamma}(\beta_m^2; a, b)$, where a and b are constants and are usually set close to zero. By integrating over the hyperparameters, we can obtain a student-t prior for each weight, which is known to enforce sparse representations during learning by setting some feature weights to zero and avoiding overfitting.

$$P(\Psi|\Omega) = P(\boldsymbol{w}|0, \boldsymbol{\beta}^2)P(\boldsymbol{\beta}^2|a, b)$$
$$= \prod_{m=1}^{M} \mathcal{N}(w_m|0, \beta_m^2) \prod_{m=1}^{M} Inverse - Gamma(\beta_m^2|a, b) \quad (17)$$

**Parameter Estimation:** With the formulated posterior probability, the learning objective is to find the optimal estimation of the parameters $\Psi$ that maximize the posterior. Hence, by inferring Equation 15, we can have the log of the posterior for the proposed model.

$$\mathcal{L}(\boldsymbol{w}, \boldsymbol{\beta}^2|Y, \Pi, a, b, \sigma^2) =$$
$$\sum_{i=1}^{I} \left[ -\frac{1}{2} \ln \sigma^2 - \frac{(y_i - f_i)^2}{2\sigma^2} \right] + \sum_{i=1}^{I-1} \sum_{h=i+1}^{I} ln \frac{1}{1 + exp(-(f_i - f_h))} \quad (18)$$
$$+ \sum_{m=1}^{M} \left[ -\frac{1}{2} \ln \beta_m^2 - \frac{w_m^2}{2\beta_m^2} \right] + \sum_{m=1}^{M} \left[ -(a+1) \ln \beta_m^2 - \frac{b}{\beta_m^2} \right]$$

We apply a gradient descent method to maximize the posterior by updating $w_m, \beta_m^2$ through $w_m^{(t+1)} = w_m^{(t)} - \epsilon \frac{\partial(-\mathcal{L})}{\partial w_m}$ and $\beta_m^{2(t+1)} = \beta_m^{2(t)} - \epsilon \frac{\partial(-\mathcal{L})}{\partial \beta_m^2}$ where

$$\frac{\partial(\mathcal{L})}{\partial w_m} = \sum_{i=1}^{I} \frac{1}{\sigma^2}(y_i - \sum_{m=1}^{M} w_m \cdot x_{im})x_{im} +$$
$$\sum_{i=1}^{I-1} \sum_{h=i+1}^{I} \frac{exp(-(f_i - f_h))}{1 + exp(-(f_i - f_h))}(x_{im} - x_{hm}) + \frac{-w_m}{\beta_m^2} \quad (19)$$

$$\frac{\partial(\mathcal{L})}{\partial \beta_m^2} = \frac{w_m^2 + b}{\beta_m^4} - \frac{3 + 2a}{2\beta_m^2} \quad (20)$$

### D. Ranking Inference

After parameters $\Psi$ are estimated via maximizing the posterior probability, we will obtain the learned model for investment value of estate, i.e., $\mathbb{E}(y_i|\boldsymbol{w}, \boldsymbol{\beta}) = \boldsymbol{x_i w}$ given a rising or falling market period. For a new coming estate $k$, we may predict its investment value accordingly. The larger the $\mathbb{E}(y_k|\boldsymbol{w}, \boldsymbol{\beta})$ is, the higher investment value it has.

For practical usage, we train two ranking models, $g(x)$ and $g'(x)$, for the rising and falling markets respectively. Since we do not predict whether a market will go up or go down, we feed the features of a real estate into two models respectively and generate two value levels, which denote its value-adding and value-protecting abilities in rising and falling markets. To provide a unified ranking to users, the output of these two models can be aggregated as $R = \alpha \cdot g(x) + (1 - \alpha) \cdot g'(x)$.

### III. EXPERIMENTAL RESULTS

We provide an empirical evaluation of the performances of the proposed method on real-world estate related data.

### A. Experimental Data

Table II shows five data sources. The taxi GPS traces are collected from a Beijing taxi company. Each trajectory contains trip id, distance(m), travel time(s), average speed(km/h), pick-up time and drop-off time, pick-up point and drop-off point. Also, we extract features from the Beijing smart card transactions. Each bus trip has card id, time, expense, balance, route name, pick-up and drop-off stops information (names, longitudes and latitudes). Moreover, the check-in data of Beijing is crawled from www.jiepang.com which is a Chinese version of Fourquare. Each check-in event includes poi name, poi category, address, longitude and latitude, comments. Furthermore, we crawl the online business reviews of Beijing from www.dianping.com which is a business review site in China. Each review contains shop ID, name, address, latitude and longitude, consumption cost, star (from 1 to 5), poi category, city, environment, service, and overall ratings. Finally, we crawl the Beijing estate data from www.soufun.com which is the largest real-estate online system in China.

TABLE II. STATISTICS OF THE EXPERIMENTAL DATA.

| Data Sources | Properties | Statistics |
|---|---|---|
| Taxi Traces | Number of taxis | 13,597 |
| | Effective days | 92 |
| | Time period | Apr. - Aug. 2012 |
| | Number of trips | 8,202,012 |
| | Number of GPS points | 111,602 |
| | Total distance(km) | 61,269,029 |
| Smart Card Transactions | Number of bus stops | 9,810 |
| | Time Period | Aug 2012 to May 2013. |
| | Number of car holders | 300,250 |
| | Number of trips | 1,730,000 |
| Check-Ins | Number of check-in POIs | 5,874 |
| | Number of check-in events | 2,762,128 |
| | Number of POI categories | 9 |
| | Time Period | 01/2012-12/2012 |
| Business Review | Number of business POIs | 1472 |
| | Number of reviews | 470846 |
| | Number of users | 159820 |
| Real Estates | Number of real estates | 2,851 |
| | Size of bounding box (km) | 40*40 |
| | Time period of transactions | 04/2011 - 09/2012 |

### B. Baseline Algorithms

To show the effectiveness of our method, we compare our method against the following algorithms. (1) **MART [12]:** it is a boosted tree model, specifically, a linear combination of the outputs of a set of regression trees. (2) **RankBoost [13]:** it is a boosted pairwise ranking method, which trains multiple weak rankers and combines their outputs as final ranking. (3) **Coordinate Ascent [14]:** it uses domination loss and applies coordinate descent for optimization. (4) **LambdaMART [15]:** it is the boosted tree version of LambdaRank, which is based on RankNet. LambdaMART combines MART and LambdaRank. (5) **FenchelRank [16]** beyond traditional ranking methods, we further compare with FenchelRank which is designed for solving the sparse learning-to-rank (LTR) problem with a L1 constraint.

We utilize RTree [3] to index geographic items (i.e., taxi and bus trajectories, checkins, etc.) and extract the defined features. We use Jieba [4] which is a Chinese/English text segmentation module to segment words and extract tags. For traditional LTR algorithms, we use RankLib [5]. We set the number of trees =

---

[3]https://pypi.python.org/pypi/Rtree/
[4]https://github.com/fxsjy/jieba
[5]http://sourceforge.net/p/lemur/wiki/RankLib/

1000, the number of leaves = 10, the number of threshold candidates = 256, and the learning rate = 0.1 for MART. We set the number of iteration = 300, the number of threshold candidates = 10 for RankBoost. We set step base = 0.05, step scale = 2.0, tolerance = 0.001, and slack = 0.001 for Coordinate Ascent. We set number of trees = 100, number of leaves = 10, number of threshold candidates = 256, learning rate = 0.1 for LambdaMART. For FenchelRank, we use the source code[6] provided by the author. We set a=0.01, b=0.01, and $\sigma^2 = 1000$ for our model.

All the codes are implemented in R (modeling), Python (feature extraction) and Matlab (visualization). And all the evaluations are performed on a x64 machine with i7 3.40GHz Intel CPU (with 4 cores) and 24GB RAM. The operation system is Microsoft Windows 7.

*C. Evaluation Metrics*

**Normalized Discounted Cumulative Gain.** The discounted cumulative gain (DCG@N) is given by $DCG[n] = \begin{cases} rel_1 & if\ n = 1 \\ DCG[n-1] + \frac{rel_n}{log_2 n}, & if\ n >= 2 \end{cases}$ Later, given the ideal discounted cumulative gain $DCG'$, NDCG at the n-th position can be computed as NDCG[n]= $\frac{DCG[n]}{DCG'[n]}$. The larger NDCG@N is, the higher top-N ranking accuracy is.

**Precision and Recall.** Since we use a five-level rating system $(4 > 3 > 2 > 1 > 0)$ instead of binary rating, we treat the rating $\geq 3$ as "high-value" and the rating $< 3$ as "low-value". Given a top-N estate list $E_N$ sorted in a descending order of the prediction values, the precision and recall are defined as Precision@N = $\frac{|E_N \bigcap E_{>3}|}{N}$ and Recall@N = $\frac{|E_N \bigcap E_{>3}|}{|E_{\geq 3}|}$, where $E_{>3}$ are the estates whose ratings are greater or equal to three (3).

**Kendall's Tau Coefficient.** Kendall's Tau Coefficient (or Tau for short) measures the overall ranking accuracy. Let us assume that each estate i is associated with a benchmark score $y_i$ and a predicted score $f_i$. Then, for an estate pair $< i, j >$, $< i, j >$ is said to be concordant, if both $y_i > y_j$ and $f_i > f_j$ or if both $y_i < y_j$ and $f_i < f_j$. Also, $< i, j >$ is said to be discordant, if both $y_i < y_j$ and $f_i > f_j$ or if both $y_i < y_j$ and $f_i > f_j$. Tau is given by Tau = $\frac{\#_{conc} - \#_{disc}}{\#_{conc} + \#_{disc}}$.

*D. Correlation Analysis*

We provide a visualization analysis to validate the correlation between the extracted features and estate investment values. We use scatter-plot matrix for correlation analysis. Each non-diagonal chart in a scatter plot matrix shows the correlation between a pair of features whose feature names are listed in the corresponding diagonal charts. Given a set of N features, there are N-choose-2 pairs of features, and thus the same numbers of scatter plots. The dots represent the estates and their colors represent the grades of investment value. For readability, we use $R5 > R4 > R3 > R2 > R1$ (symbol ) to represent $4 > 3 > 2 > 1 > 0$ (number) in Figure 5.

In Figure 5(a), we present the correlation between business review features (overall satisfaction, service quality, environment class, consumption cost) and estate investment value. As can be seen, the R5 estates tend to appear at the top right corner

of all the non-diagonal charts. This implies that if mobile users have higher ratings for estate neighborhoods, estate investment values are the higher. Remind that we mean the heterogenesis of poi planning by the entropy of frequency of categorized POIs. Interestingly, we observe if the heterogenesis of functionality planning is too high or too low, these estates are usually low-value. This can be intuitively explained by the fact that people are willing to live in a community that can meet and balance the needs of their life.

In Figure 5(b), we show the positive correlation between the taxi leaving, arriving and transition volumes of estate neighborhoods and estate investment value. However, the commute distance of taxies has negative correlation with estate investment value. In other words, the shorter the commute distance of taxies is, the higher is the estate investment value. A potential interpretation of this observation is that since taxies are valued by white-collar and business people, the destinations of taxi trajectories usually are important places (e.g., conference centers, business hotels, companies and government organizations, etc). If the commute distance of taxies is short, the targeted neighborhood is close to these important places.

In Figure 5(c), we show the positive correlation between estate investment value and bus related features, such as the leaving, arriving, and transition volumes of buses, bus stop density. Figure 5(d) illustrates that Topic 4 has positive correlation with estate investment value whereas Topic 1,2,3,5 have negative correlation. This validates topic profiling of checkin posts can help discriminate estate values.

The visualization results show the collectiveness of our intuitions for defining and extracting discriminative features
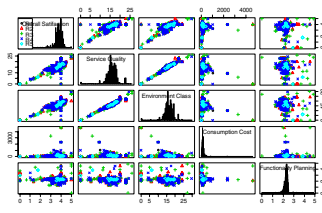
*E. Feature Evaluation*

We evaluate the performances of different features segmented from two perspectives.
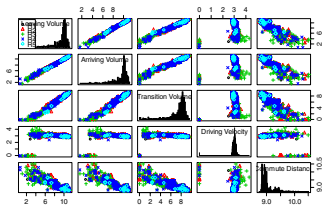
**Evaluation on features of different data sources.** We segment the extracted features in terms of different data sources and investigate which source is more effective for ranking estates. Figure 6 and Figure 7 shows the Tau, NDCG, Precision, and Recall of four feature sets (business reviews, taxi traces, smart card transactions and check-ins) in rising market and falling market respectively. In all cases, we observe the extracted features achieve good performances, yet there are features which are substantially better than others.

Specifically, the check-in features perform best with Tau 0.1046198, NDCGs > 0.75, Precisions > 0.85, and Recalls > 0.24 in rising market, and consistently achieve the best ranking results in falling market. The features of business reviews hold the second place of overall and top-k rankings in rising and falling markets. In sum, business reviews and check-ins performs better than taxi and bus traces. One possible reason is that people's outdoor activities consist (1) moving phrase and (2) attending phrase. Although moving phrase (taxi and bus trajectories) help realize activity attending (check-ins and business reviews), the drop-off points of taxi and bus trajectories are not always the destinations of outdoor activities. Whereas, the locations of check-ins and business reviews usually are the final destinations of people's visits. They reflect direct interaction between users and activities via locations, and thus have semantically richer information
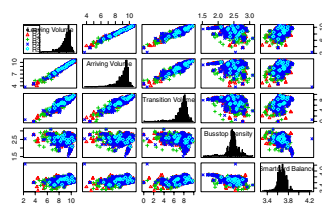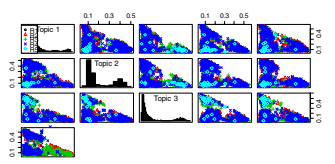
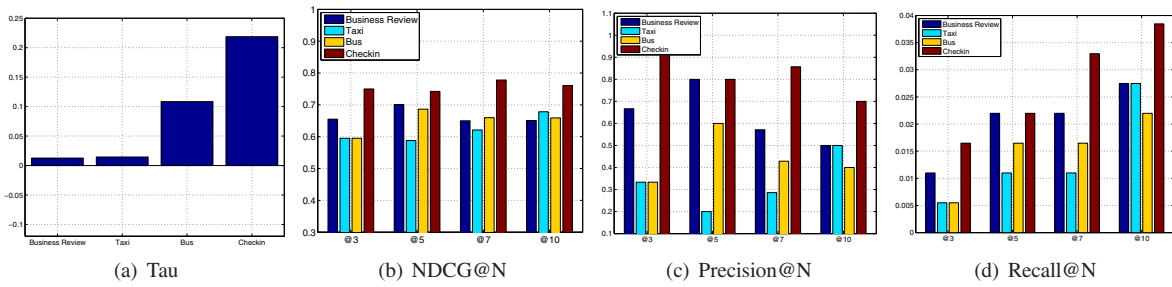(a) Features of business review     (b) Features of taxi traces     (c) Features of bus traces

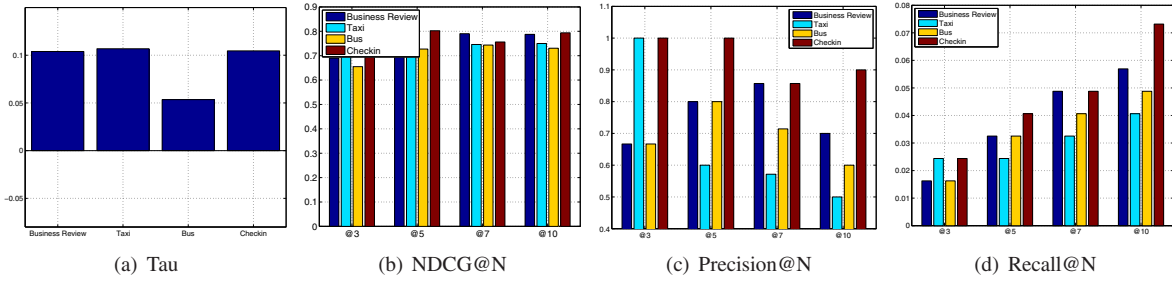Fig. 6. Feature performances of different sources on the rising market dataset.



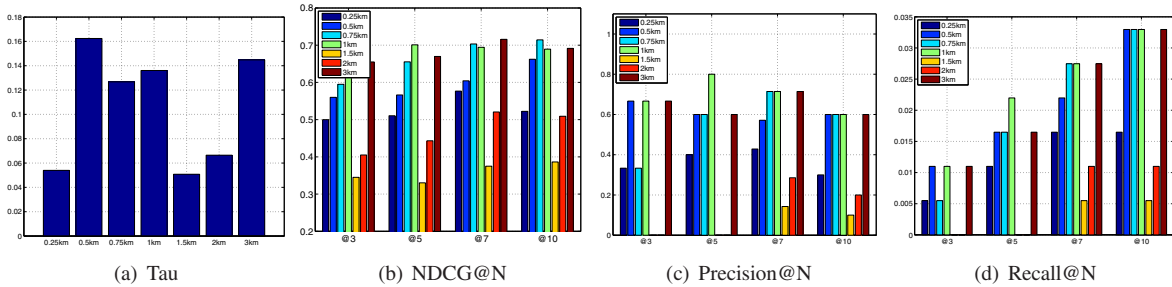Fig. 7. Feature performances of different sources on the falling market dataset.



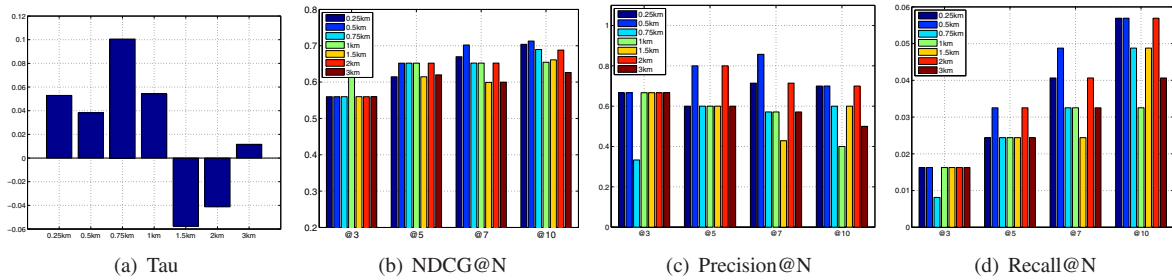Fig. 8. Feature performances of different radius on the rising market dataset.



Fig. 9. Feature performances of different radius on the falling market dataset.

of estate prices [2]. Work [1] checks the volatility of estate price and concludes that low investment-valued estate values relatively volatile. More classic works are based on repeat sales methods and hedonic methods. The repeat sales methods [18] construct a predefined price index based on properties sold more than once during the given period. The hedonic methods [3] assume the price of a property depends on its characteristics and location. Work [1] studies the automated valuation models which aggregate and analyze physical characteristics and sales prices of comparable properties to provide property valuations. More recent works [4], [19] apply general additive mode, support vector machine regression, multilayer perceptron, ranking and clustering ensemble method to computational estate appraisal. In our earlier work [4], we focus on exploiting the mutual enhancement between ranking and clustering to model geographic utility, popularity and influence of latent business area for estimating estate value. Besides, in [4], we identify and jointly capture the geographical individual,

peer, and zone dependencies as an estate-specific ranking objective for enhancing prediction of estate value. However, in this paper, we details comprehensive feature designs that cover most of aspects that have an impact on estate value. Also, we integrate sparsity regularization into pairwise ranking strategy because the extracted features are usually correlated and redundant.

Also, our work can be categorized into Learning-To-Rank (LTR) which includes pointwise, pairwise, and listwise approaches [20]. The point-wise methods [20] reduce the LTR task to a regression problem: given a single query-document pair, predict its score. The pair-wise methods approximate the LTR task to a classification problem. The goal of the pairwise ranking is to learn a binary classifier to identify the better document in a given document pair by minimize average number of inversions in ranking [13], [21]–[23]. The list-wise methods, optimize a ranking loss metric over lists

instead of document pairs [24]. For instance, H. Li et al. propose AdaRank [25] and ListNet [26] and Burges et al. propose LambdaMART [15]. More recent work [16] further learn the ranking model which is constrained to be with only a few nonzero coefficients using L1 constraint and propose a learning algorithm from the primal dual perspective.

Urban computing [27] is a process of acquisition, integration, and analysis of urban data (e.g., sensors, devices, vehicles, buildings, human) to tackle the major issues that cities face. Our work also has a connection with mining mobile, geography and mobility data to tackle issues in urban space. Work [29] identifies emerging patterns with multirelational approach from spatial data. Liu et al. detects spatio-temporal causality of outliers in traffic data [30]. Yuan et al. discovers regional functions of a city using POIs and taxi traces [31] . Heierman et al. mines the device usage patterns of homeowners for smart houses [32] . Paper [33] selects the optimal sites for retail stores by mining Foursquare data. [27] mines the driving route for end users by considering physical feature of a route, traffic flow, and driving behavior.

## V.    Conclusions

In this paper, we aimed to assess estate investment value by mining a variety of user-generated data. We collected a large scale of online user reviews and offline moving behaviors (taxi traces, smart card transactions, and checkins) of mobile users. We index, filter, propagate, distill, aggregate mobile data, and extract the fine-grained features from multiple perspectives (e.g., direction, volume, velocity, heterogeneity, popularity, topic, etc.) for evaluating estate values. However, since the extracted estate features usually are intercorrelated and redundant, we proposed to learn a sparse pairwise ranker, which is mutually enhanced by simultaneously conducting feature selection and maximizing estate ranking accuracy. Finally, the experimental results with real world estate-related data demonstrates the competitive effectiveness of both extracted features and learning models.

## Acknowledgement

## References

[1]  M. L. Downie and G. Robson, "Automated valuation models: an international perspective," 2007.

[2]  L. D. B. Chaitra H. Nagaraja and L. H. Zhao, "An autoregressive approach to house price modeling," 2009.

[3]  L. O. Taylor, "The hedonic method," in *A primer on nonmarket valuation*.   Springer, 2003.

[4]  Y. Fu, H. Xiong, Y. Ge, Z. Yao, Y. Zheng, and Z.-H. Zhou, "Exploiting geographic dependencies for real estate appraisal: A mutual perspective of clustering and ranking," in *KDD'14*, 2014.

[5]  J. Landis, S. Guhathakurta, W. Huang, M. Zhang, and B. Fukuji, "Rail transit investments, real estate values, and land use change: a comparative analysis of five california rail transit systems," 1995.

[6]  D. R. Bowes and K. R. Ihlanfeldt, "Identifying the impacts of rail transit stations on residential property values," *Journal of Urban Economics*, vol. 50, no. 1, pp. 1–25, 2001.

[7]  S. Lewis-Workman and D. Brod, "Measuring the neighborhood benefits of rail transit accessibility," *Transportation Research Record: Journal of the Transportation Research Board*, vol. 1576, no. 1, pp. 147–153, 1997.

[8]  K. Wardrip, "Public transits impact on housing costs: a review of the literature," 2011.

[9]  A. H. b. Hj. Mar Iman al Murshid, "Modelling locational factors using geographic information system generated value response surface techniques to explain and predict residential property values," in *NAPREC Conference*, 2008.

[10]  A. Montanari and B. Staniscia, "From global to local: Human mobility in the rome coastal area in the context of the global economic crisis*," *Belgeo. Revue belge de géographie*, no. 3-4, pp. 187–200, 2012.

[11]  C. D. K. Robert Cervero, "Bus rapid transit impacts on land uses and land values in seoul, korea," *Transport Policy*, vol. 18, pp. 102–116, 2011.

[12]  J. H. Friedman, "Greedy function approximation: a gradient boosting machine," *Annals of Statistics*, 2001.

[13]  Y. Freund, R. Iyer, R. E. Schapire, and Y. Singer, "An efficient boosting algorithm for combining preferences," *The Journal of machine learning research*, 2003.

[14]  D. Metzler and W. B. Croft, "Linear feature-based models for information retrieval," *Information Retrieval*, 2007.

[15]  C. Burges, "From ranknet to lambdarank to lambdamart: An overview," *Learning*, 2010.

[16]  H. Lai, Y. Pan, C. Liu, L. Lin, and J. Wu, "Sparse learning-to-rank via an efficient primal-dual algorithm," *Computers, IEEE Transactions on*, 2013.

[17]  J. Krainer and C. Wei, "House prices and fundamental value," *FRBSF Economic Letter*, 2004.

[18]  R. J. Shiller, "Arithmetic repeat sales price estimators," Cowles Foundation for Research in Economics, Yale University, Tech. Rep., 1991.

[19]  V. Kontrimas and A. Verikas, "The mass appraisal of the real estate by computational intelligence," *Applied Soft Computing*, vol. 11, pp. 443 – 448, 2011.

[20]  L. Hang, "A short introduction to learning to rank," *IEICE TRANSACTIONS on Information and Systems*, 2011.

[21]  C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. Hullender, "Learning to rank using gradient descent," in *ICML'05*, 2005.

[22]  C. Quoc and V. Le, "Learning to rank with nonsmooth cost functions," *NIPS'07*, 2007.

[23]  J. Fürnkranz and E. Hüllermeier, "Pairwise preference learning and ranking," in *Machine Learning: ECML 2003*.   Springer, 2003.

[24]  F. Xia, T.-Y. Liu, J. Wang, W. Zhang, and H. Li, "Listwise approach to learning to rank: theory and algorithm," in *ICML'08*, 2008.

[25]  J. Xu and H. Li, "Adarank: a boosting algorithm for information retrieval," in *SIGIR '07*, 2007.

[26]  Z. Cao, T. Qin, T.-Y. Liu, M.-F. Tsai, and H. Li, "Learning to rank: from pairwise approach to listwise approach," in *ICML'07*, 2007.

[27]  Y. Zheng, L. Capra, O. Wolfson, and H. Yang, "Urban computing: concepts, methodologies, and applications," *ACM TIST*, 2014.

[28]  V. S. Tseng and K. W. Lin, "Efficient mining and prediction of user behavior patterns in mobile web systems," *Information and software technology*, 2006.

[29]  M. Ceci, A. Appice, and D. Malerba, "Discovering emerging patterns in spatial databases: A multi-relational approach," in *PKDD'07*.   Springer Berlin Heidelberg, 2007.

[30]  W. Liu, Y. Zheng, S. Chawla, J. Yuan, and X. Xing, "Discovering spatio-temporal causal interactions in traffic data streams," in *KDD '11*, 2011.

[31]  J. Yuan, Y. Zheng, and X. Xie, "Discovering regions of different functions in a city using human mobility and pois," in *KDD'12*, 2012.

[32]  E. O. Heierman III and D. J. Cook, "Improving home automation by discovering regularly occurring device usage patterns," in *ICDM'03*, 2003.

[33]  D. Karamshuk, A. Noulas, S. Scellato, V. Nicosia, and C. Mascolo, "Geo-spotting: Mining online location-based services for optimal retail store placement," in *KDD '13*, 2013.