



浙江大学
ZHEJIANG UNIVERSITY

Forecasting Inter-destination Tourism Flow Within City With Social Media Data via Hybrid Learning model

--- benchmark, interpretation and model

Speaker: Hanxi Fang

Instructor: Prof. Song Gao

CONTENTS



Introduction



Relevant Works



Data & Methods



Results & Conclusion



Discussion



PART ONE

Introduction



Introduction

<https://www.mafengwo.cn/xc/10065/>

来自游记《穷游走北京：一个人的北京，再见了我的大学》

北京10日游

🕒 时间: 2015-10-10 ⌚ 出行天数: 10天 💰 人均费用: 4000RMB

- DAY 1 簋街 > 北京师范大学
- DAY 2 故宫 > 圆明园 > 天安门广场 > 清华大学...
- DAY 3 南锣鼓巷 > 文字奶酪店(南锣鼓巷店) > 创可贴8特色T恤店 > 什刹海酒吧街...
- DAY 4 北京798艺术区 > 三里屯太古里
- DAY 5 鸟巢 > 水立方 > 中国国家图书馆
- DAY 6 地坛公园 > 王府井小吃街 > 王府井步行街 > 北京全聚德(王府井店)
- DAY 7 北京大学
- DAY 8 颐和园 > 东来顺饭庄(西直门店)
- DAY 9 天坛 > 帽儿胡同 > 吉事果(南锣鼓巷店) > 南锣鼓巷...
- DAY 10 北京 > 广州

When planning their trip how people decide which scenic spots to go in the same trip?

Specifically, why people will go to both scenic spot A and B in the same trip?

- Near?
- Of the same type?
- They are both popular?
- Drop by visiting?
-



The nature of ITF:
A quantitative metric to represent
the flow intensity between
tourism destinations.

We extract over 20000 trips from them.

[illegible]



Usage of ITF



For tourism companies:

- Help to design a tourism recommending system and combined tickets for scenic spots.

for city planners:

- Offer important guide to the design of city tour bus.
- Offer guide to city planning, especially tourism planning that regard the whole city as a system.

for scholars:

- Define the interaction between different scenic spot, thus help us to understand the spatial pattern of tourism within city



Characteristics of ITF

ITF is hard to get

The formation of ITF is complex and hard to explain



Contribution

first

We put forward the concept of ITF, which quantitatively describe the interaction between different scenic spots within the same city; and create a benchmark about it from multiple sources of social media data

second

We create a hybrid GNN-based learning model that can predict the ITF . And the result has a mean absolute error less than 53%

third

We give an interpretation of the relationship between features of scenic spots and ITF.

The features include both explicit features of a single scenic spots, and the graph structure features of the interaction graph.



PART TWO

Relevant Works



Tourism Flow Research



01

Crampon, L.J. and Tan, K.T. (1973)

Put forward a regression model to predict international tourism flow between countries

02

Yang, X. A. , et al.(2020)

captures the nationality and movement patterns of foreign tourists to South Korea, and use a community detection algorithm partitions based on tourism flow between cities

03

Seok, H., Barnett, G.A. & Nam, Y(2021)

Use social network data to analyze international tourism flow.

The research scale are always the flow between cities or even countries.



Flow Generation Problem



01

Gravity model, deep gravity model...

02

Random forest, XGBoost...

03

SIGCN, RFGCN...

... ..

The researches are mainly focus on predicting residence-work commute flow, the data are usually less noisy with low variance.



PART THREE

Data & Methods


















Data source

Browser tabs: 北京市文化和旅游局_通知公告 | 目的地旅游攻略 - 马蜂窝 | 北京景点介绍北京旅游景点北京 | +

Address bar: mafengwo.cn/jd/10065/gonglv.html

Navigation bar: 首页 | 行程线路 | 景点 | 酒店 | 机票 | 当地玩乐 | 旅游度假 | 社区 | 餐饮 | 地图

Category bar: 全部景点 | 皇城古迹 | 本周人气玩乐热榜 | 夏日京郊游 | 世界遗产 | 特色公园 | 新地标建筑 | 长城 | 亲子游

				
故宫	八达岭长城	天安门广场	天坛	颐和园
				
圆明园	慕田峪长城	景山公园	南锣鼓巷	恭王府
				
北海公园	北京798艺术区	什刹海	北京杜莎夫人蜡像馆(前门)	雍和宫

共20页 | 1 2 3 4 5 6 后一页 末页



MaFengWo:

<http://www.mafengwo.cn/>

Amap:

<https://www.amap.com/>

Ctrip:

<https://ctrip.com/>

Qunar:

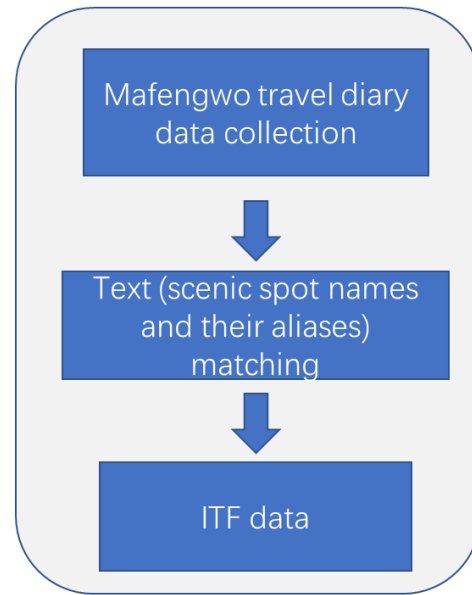
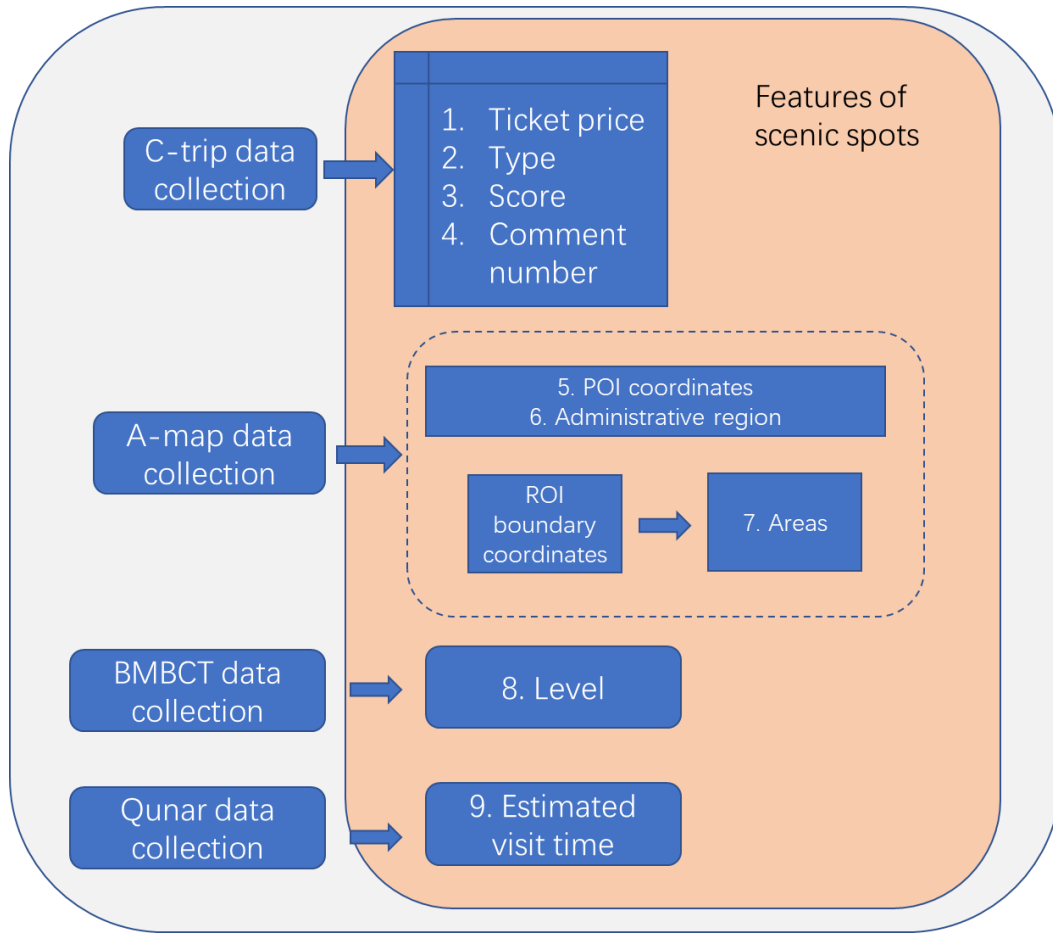
<https://travel.qunar.com/>

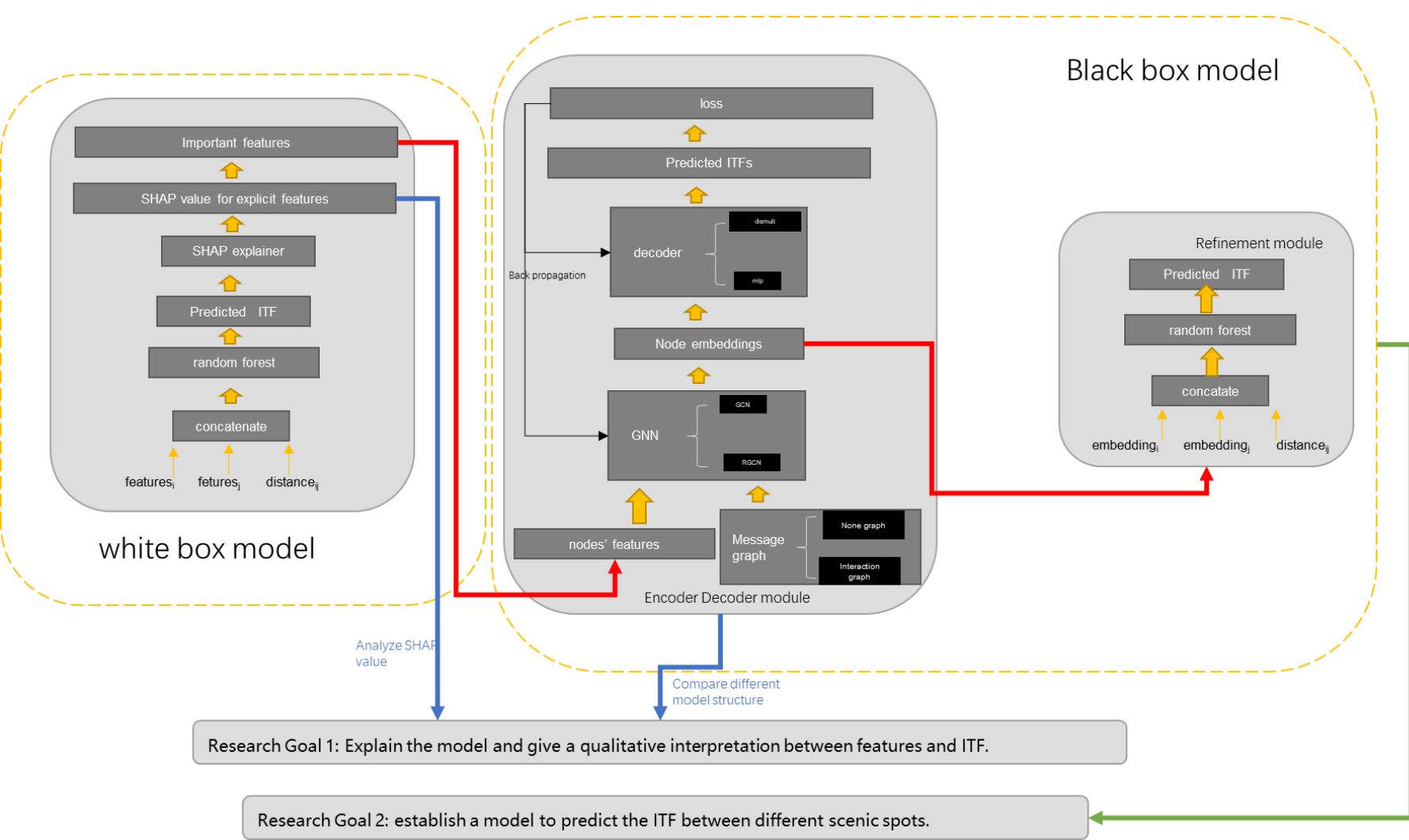
Beijing Municipal Bureau of Culture and Tourism:

<http://whlyj.beijing.gov.cn/>



Data Collection Procedure







PART FOUR

Results & Conclusion



Random forest---Multicollinearity testing

The correlation between continuous variables: Use VIF to test, VIF less than 10 means no obvious Multicollinearity.

Feature	VIF
Score	4.48
Price	1.4
Comment number	1.31
Area	1.12
Level	4.12
Estimated time	4.24

VIF values

The correlation between categorical variables: Use Cramér's Vs to test, Cramér's Vs less than 5 means no obvious Multicollinearity.

	adname	Type
Adname	1.0	0.4
Type	0.4	1.0

Cramér's Vs correlation value

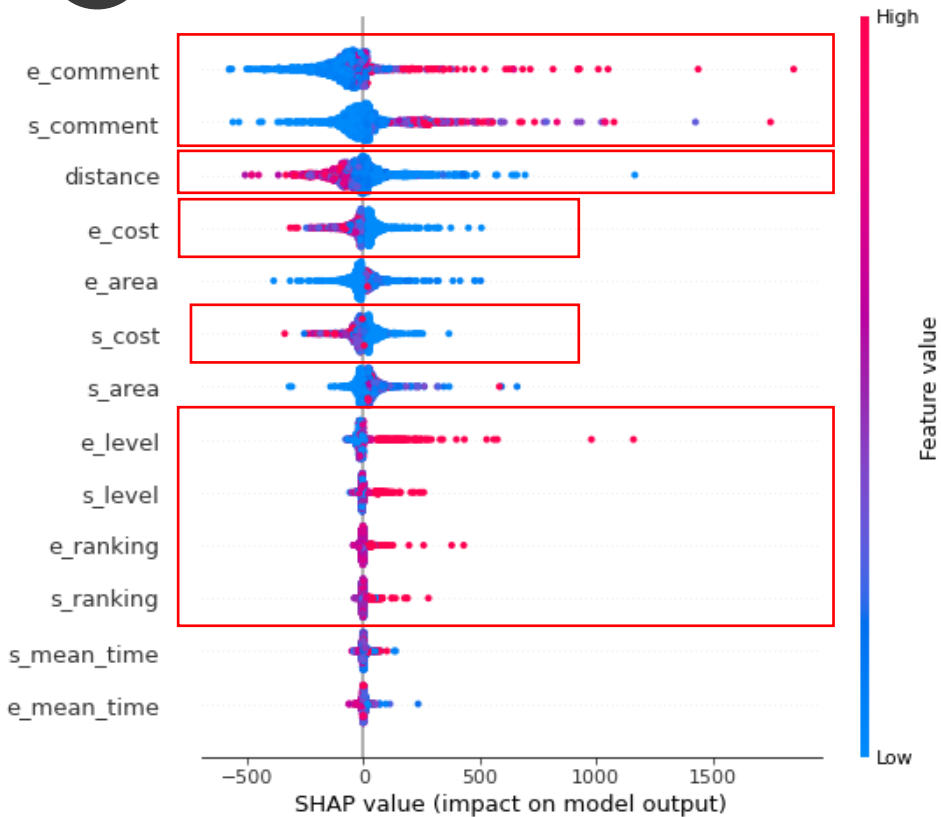
The correlation between categorical variable and continuous variables : Use One-way ANOVA test to get a p-value which means the probability of no correlation

	Score	Price	Comment number	Level	Area	Estimated time
Type	0.0	0.0	0.0	0.0	0.0	0.0
adname	1.62e-148	4.60e-288	5.60e-95	2.07e-194	0.0	0.0

P-value of one-way ANOVA test



Random forest + shap



S_: features of the start point; E_: features of the end point

Problem of random forest:

1. Hard to consider the corresponding features of start spot and end spot jointly.
2. It views the prediction of all the ITF as irrelevant regression problem, thus cannot take interaction graph structure of the spots into consideration.
3. The method doesn't make sure that the ITF between A to B is equal to the ITF between B and A; which is true in reality.

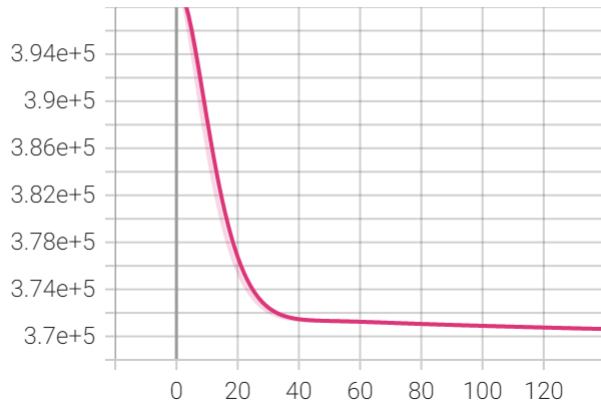
Test mape: 1.915



Deep Gravity

train

tag: Loss/train



Practically, the deep gravity model is hard to train for this problem!

Problem of random forest:

1.

solved

2.

It views all the ITFs as irrelevant regression problem, thus cannot take spatial structure of the spots into consideration.

3.

The method doesn't make sure that the ITF between A to B is equal to the ITF between B and A; which is true in reality.

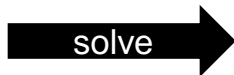
BUT!!!Test mape:

5.7362



A Perfect Solution: GNN based models

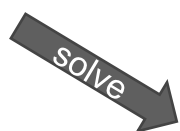
GNN based model use the node features (the spot features) and the graph structure to generate an embedding for each node.



Problem of random forest:

1. It views all the ITFs as irrelevant regression problem, thus cannot take spatial structure of the spots into consideration.
2. Hard to consider the corresponding features of start spot and end spot jointly.
3. The method doesn't make sure that the ITF between A to B is equal to the ITF between B and A; which is true in reality.

After getting the embeddings of both node A and B, we can use DISMILT as the decoder to get the predicted ITF.



Dismilt: a function (2 embeddings \rightarrow 1 scalar value)

Embedding of A: (1,2,3,4,5)

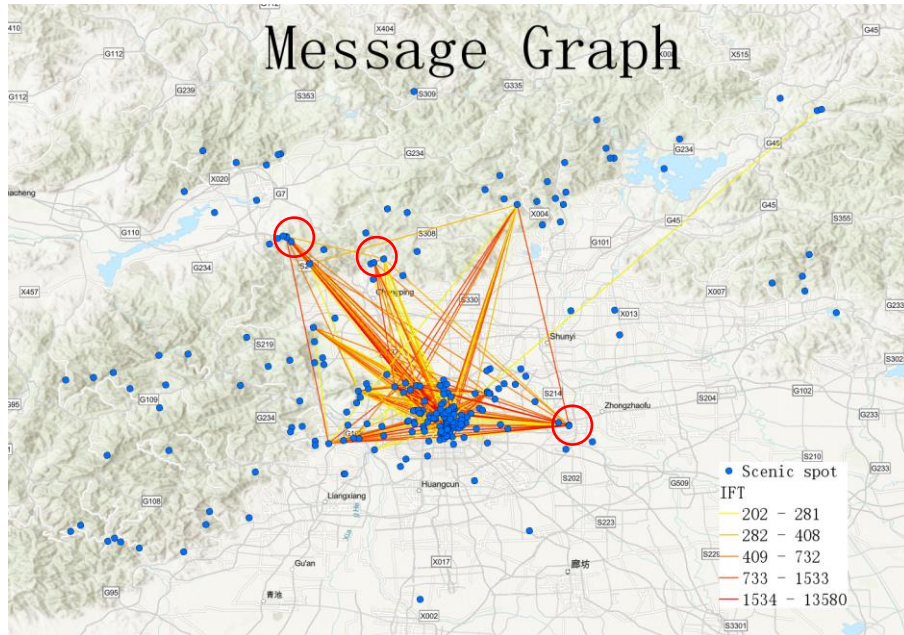
Embedding of B: (4,5,3,4,5)

Result = $(1 * t_1 * 4) + (2 * t_2 * 5) + (3 * t_3 * 6) + (4 * t_4 * 4) + (5 * t_5 * 5)$

t_1, t_2, t_3, t_4, t_5 are trainable parameters.



GNN based models---message graph

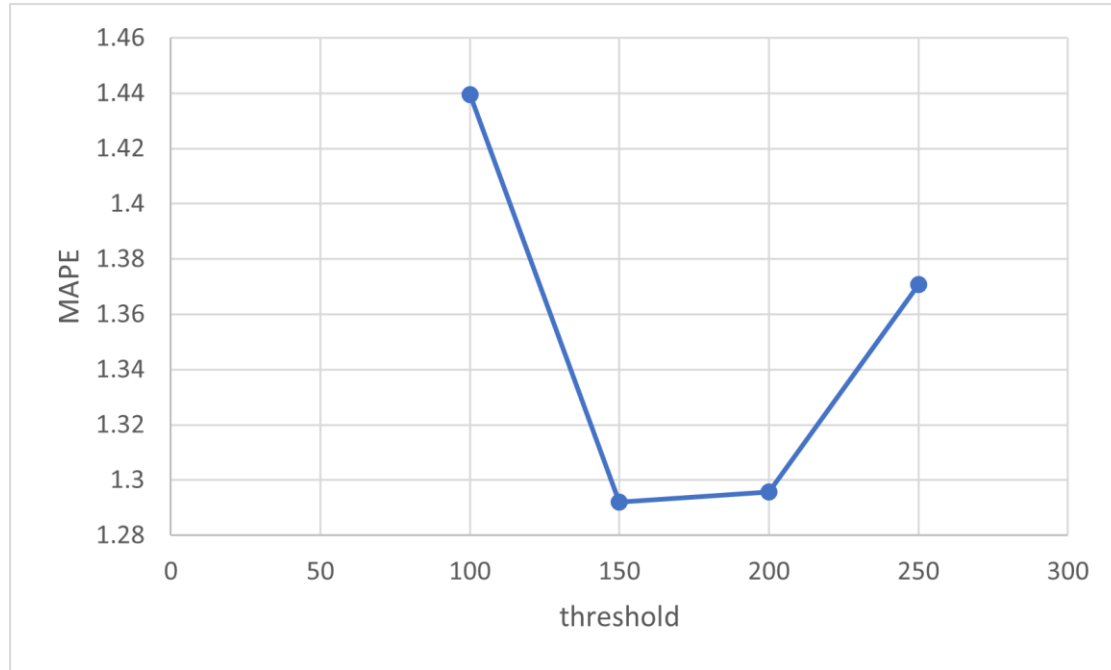


The message graph is a unweighted undirected graph. The existence of an edge between two spots A and B means there exists a strong interaction between A and B.

In our study, if there are more than 150 ITF between A and B in the training data, then there is an edge.

Note that when there is no initial information about any ITF between any two points, you can use other criteria to judge if there exists an edge, like commercial tour plans from tourism company; Since in fact the only thing you need to know is if there exist a large number of visitors who will visit both A and B, instead of the accurate number of the ITF.

GNN based models---threshold sensitive analysis





GNN based models---results

id	MODEL	USE_GRAPH	USE_DISMULT	MAPE	MSE
1	pure_rf	no	yes	1.9147	112085.5
2	deep_network(deep_gravity)	no	yes	5.7362	351221.8
3	rgcn+mlp	yes	no	5.7326	351222.1
4	gcn+mlp	yes	no	5.7321	351222.0
5	rgcn+dismult(SIGCN)	yes	yes	1.3517	13711.59
6	gcn+dismult	yes	yes	1.7882	33560.63
7	rgcn+with_only_self_edges	no	yes	2.8495	84514.1797

Conclusion:

- 1/2 vs 5/6: gcn based model can be much better than RF or Deep_gravity.
- 3/4 vs 5/6: dismult makes sure symmetric result and manual correspondence of features, and is a key to the result.
- 7 vs 5: take the spatial interaction structure into account is also useful and important.

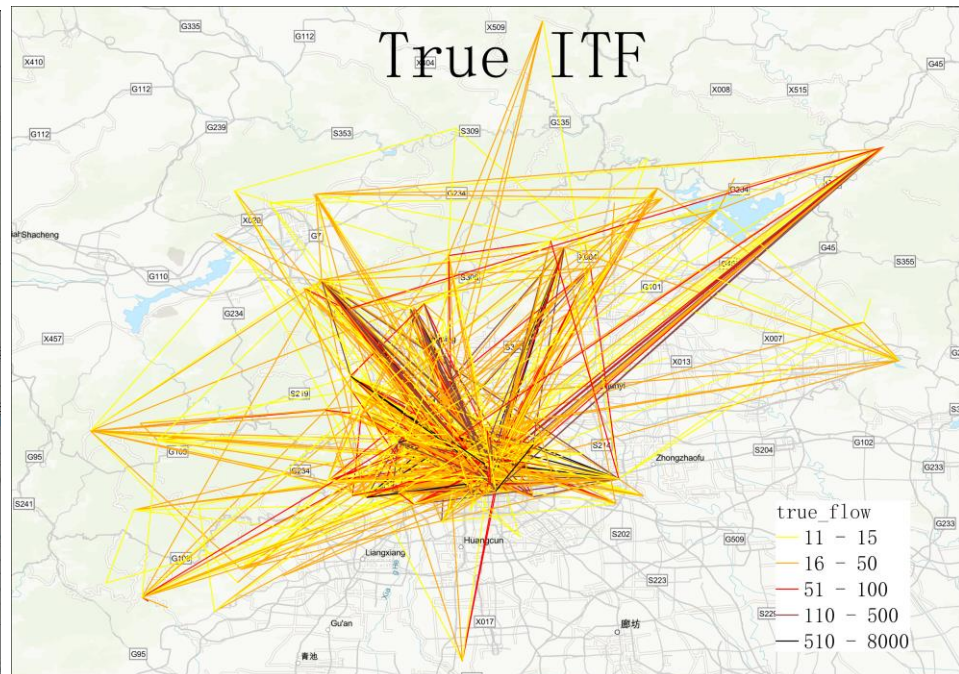
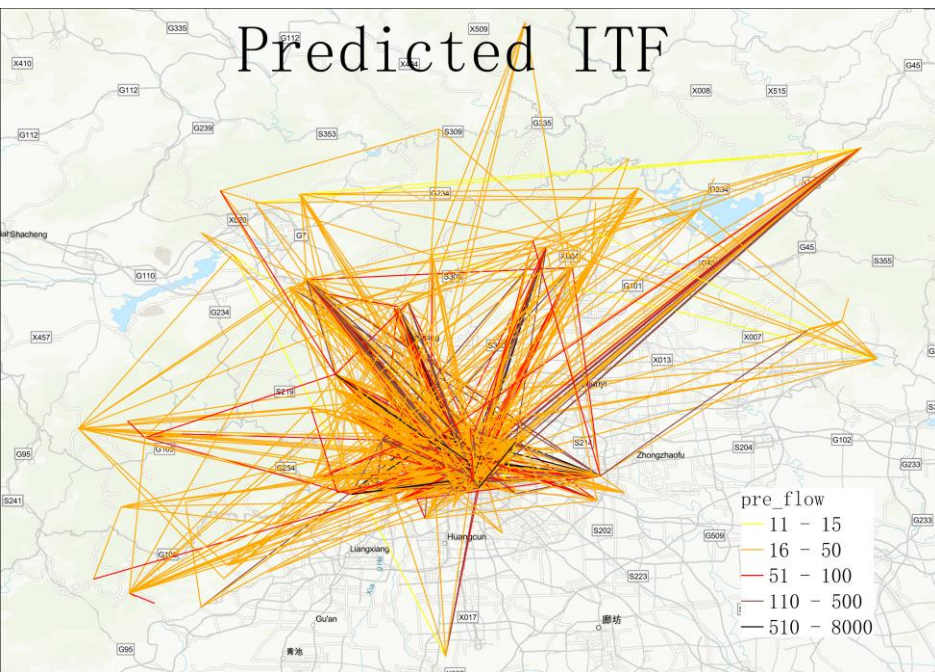


GCN+dismult+RF_refinement

ID	MODEL	USE RF as refinement	MSE	MAPE
1	rgcn+dismult(SIGCN)	no	13711.59	1.3517
2	gcn+dismult	no	33560.63	1.7882
3	rgcn+dismult+RF	yes	21406.1004	0.549
4	gcn+dismult+RF	yes	17436.8982	0.5292



GCN+dismult+RF_refinement





Conclusion

- The factors that influence the value of ITF between two spots **can be divided by two aspects**. The first aspect is **the features of the two spots**. According to the SHAP value of the random forest model, these mainly include the popularity(represented by comment number), ticket price, scores and level of both spots, and the distance between them. However, other features like tour spots' area and estimated visiting time have little impact.

The second aspect is **the interaction graph structure** influence the ITF. Taking account of the features of other scenic spots that have strong connection with the enquired spots by GNN can significantly improve the performance of the prediction model.

- For GNN based model, after getting the embedding for each spot, **the usage of DISMULT decoder and the random forest as a refinement can significantly improve** the performance of the model.

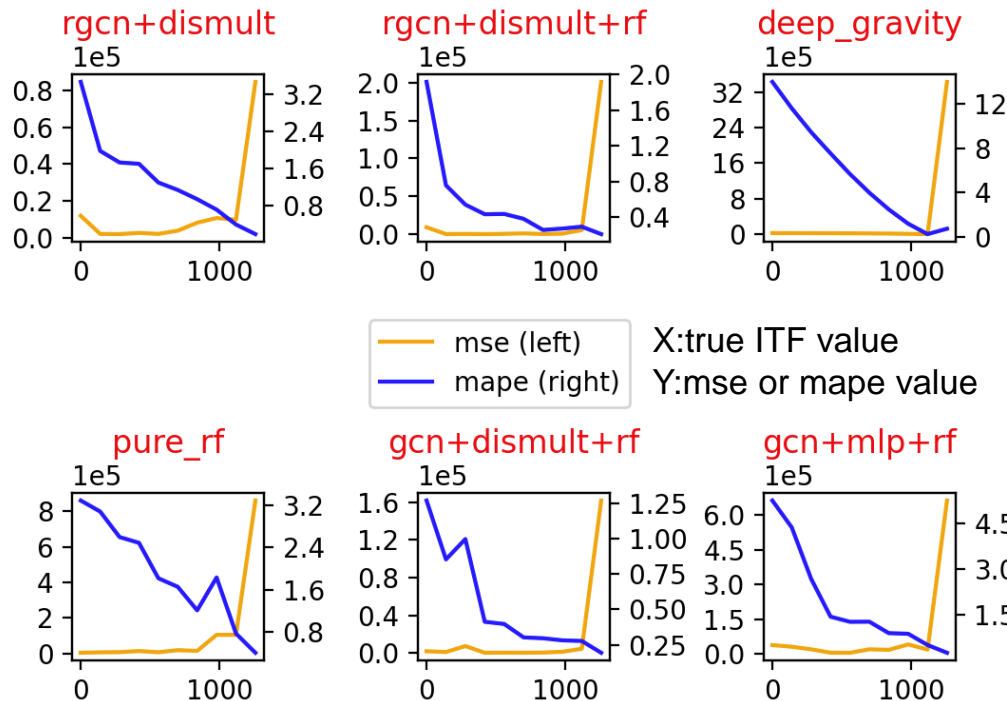


PART FIVE

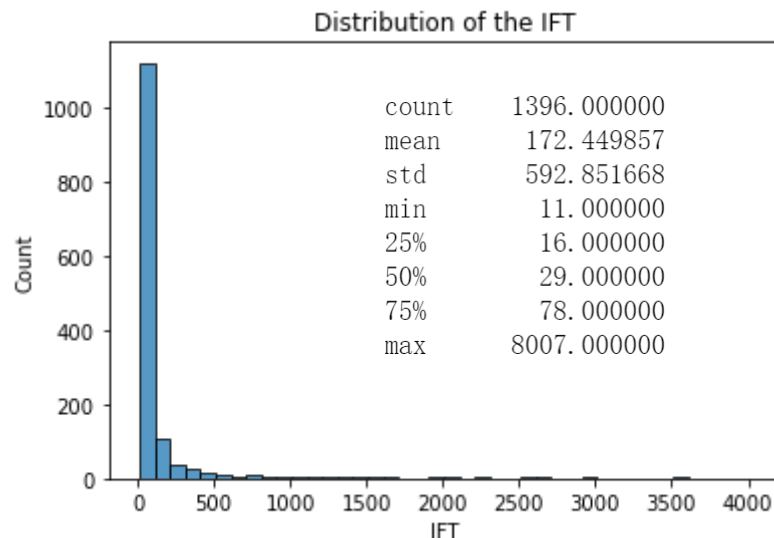
Discussion



The drawbacks of the model



Due to the strong variance of ITF, different ITF value level has different precision.





Other Things to be discovered deeper

1. How to get the information about where the user has gone from travel diary on social media with more precision.
2. How to explain how the structure of interaction graph influence the value of ITF.
3. More features can be used, like the feature embedding extracted from comment content data.

添加标题

THANKS!