

Solving Cold Start Problem For New Booking Aggregator

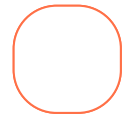
Dmitry Beresnev d.beresnev@innopolis.university

Vsevolod Klyushev v.klyushev@innopolis.university

Nikita Yaneev n.yaneev@innopolis.university

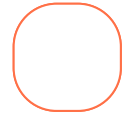
GitHub: <https://github.com/Kadaverciant/bd-s25-project/>

Introduction



Challenge

The online travel accommodation market is highly competitive, and launching a new booking aggregator in a city like Rio de Janeiro presents a critical hurdle: the cold start problem. Without historical user data—bookings, reviews, or behavioral insights—new platforms struggle to rank listings effectively or personalize recommendations, leading to poor user engagement and slow growth.



Our solution

We propose a data-driven approach leveraging publicly available Airbnb data to predict review scores for new, unrated listings. By analyzing features like location, price, amenities, and host attributes, our model generates reliable quality estimates, enabling:

1. Better ranking algorithms from day one
2. Enhanced discoverability of high-quality stays
3. Improved user experience and conversion rates

Introduction

Business Impact

This initiative directly addresses the cold start challenge, driving early-stage bookings, revenue growth, and market penetration. In this presentation, we'll explore our methodology, data insights, and the predictive system designed to ensure a strong launch.



Data collection & Description

Dataset Overview:

- Source: Airbnb (Rio de Janeiro) via Kaggle (public dataset)
- Timeframe: 3 years (2018–2020)
- Size: ~784K records, 108 attributes (2.5 GB)

Key Features:

- Geodata: Housing coordinates
- Rental Attributes: Wi-Fi, TV, beds, etc.
- Host Characteristics: Host-related metrics
- Price: Rental price per listing
- Review Score Rating: Target variable for cold-start prediction

Data Preprocessing Pipeline

Feature Selection:

- Original: 108 columns → Selected: 29 key features
- Criteria: Relevance to rating prediction (e.g., host behavior, property details, pricing)
- Target: review_scores_rating

Data Cleaning:

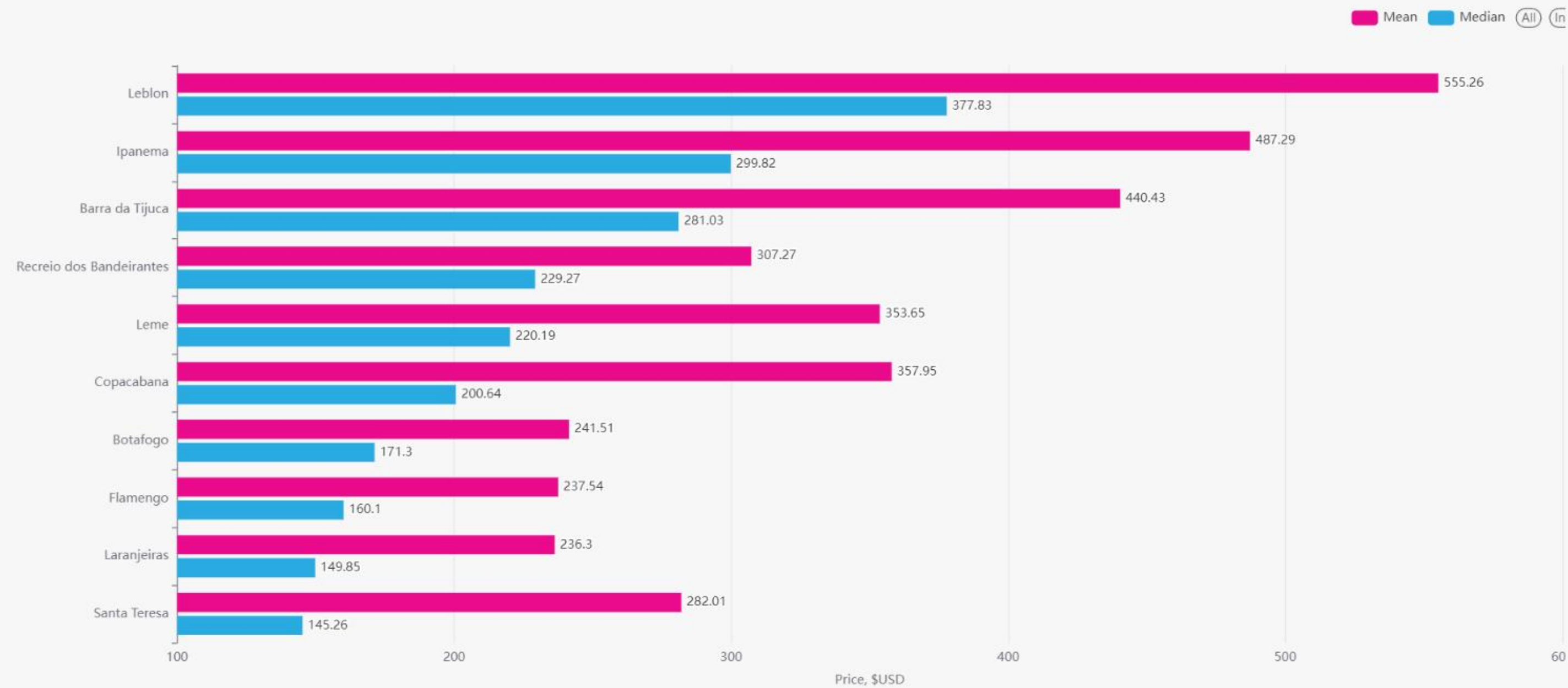
- Handling Missing Values:
- Formatting
- Result: ~379,600 records retained

Feature Engineering:

The amenities field originally contained over 180 unique entries, representing various features of each listing. We selected the 20 most common ones in order to have more actionable information

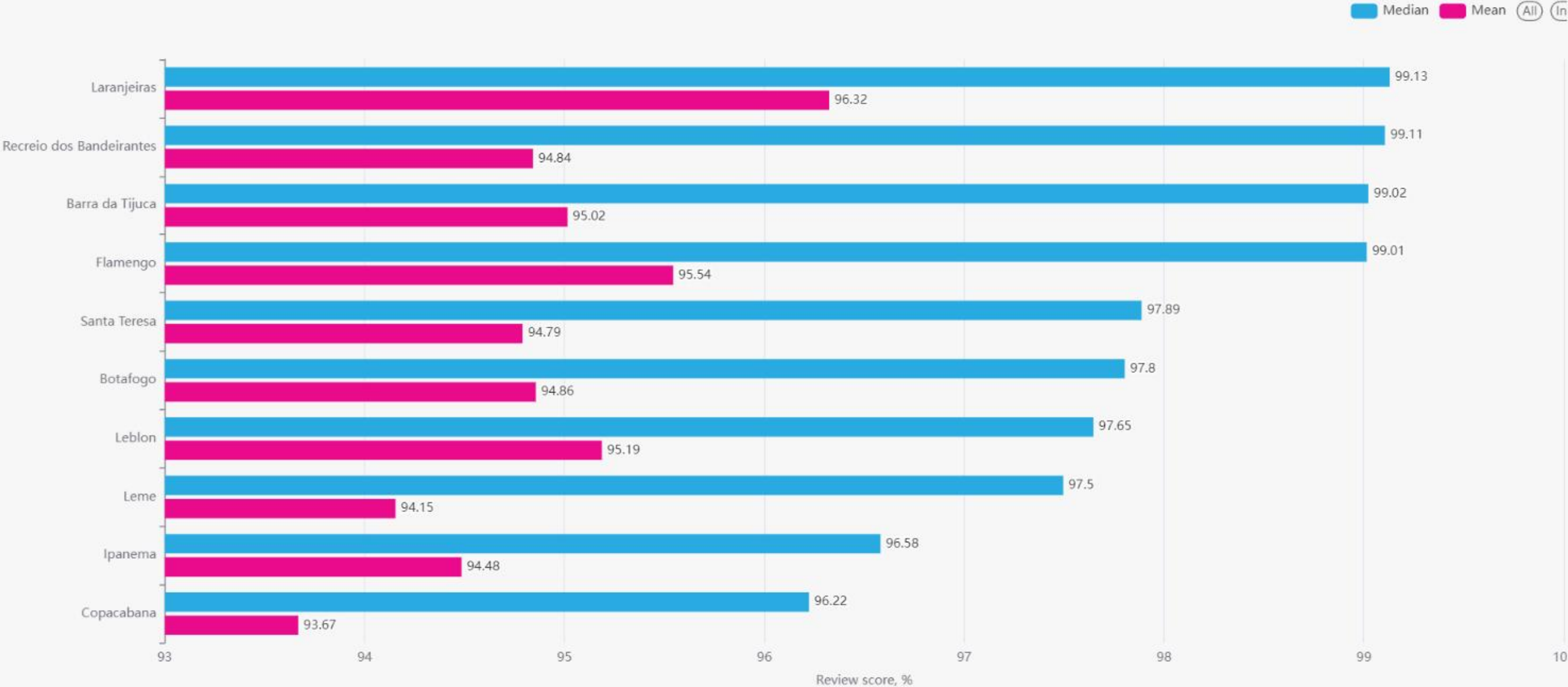
Insight 1.1

Price distribution over 10 most popular neighborhoods

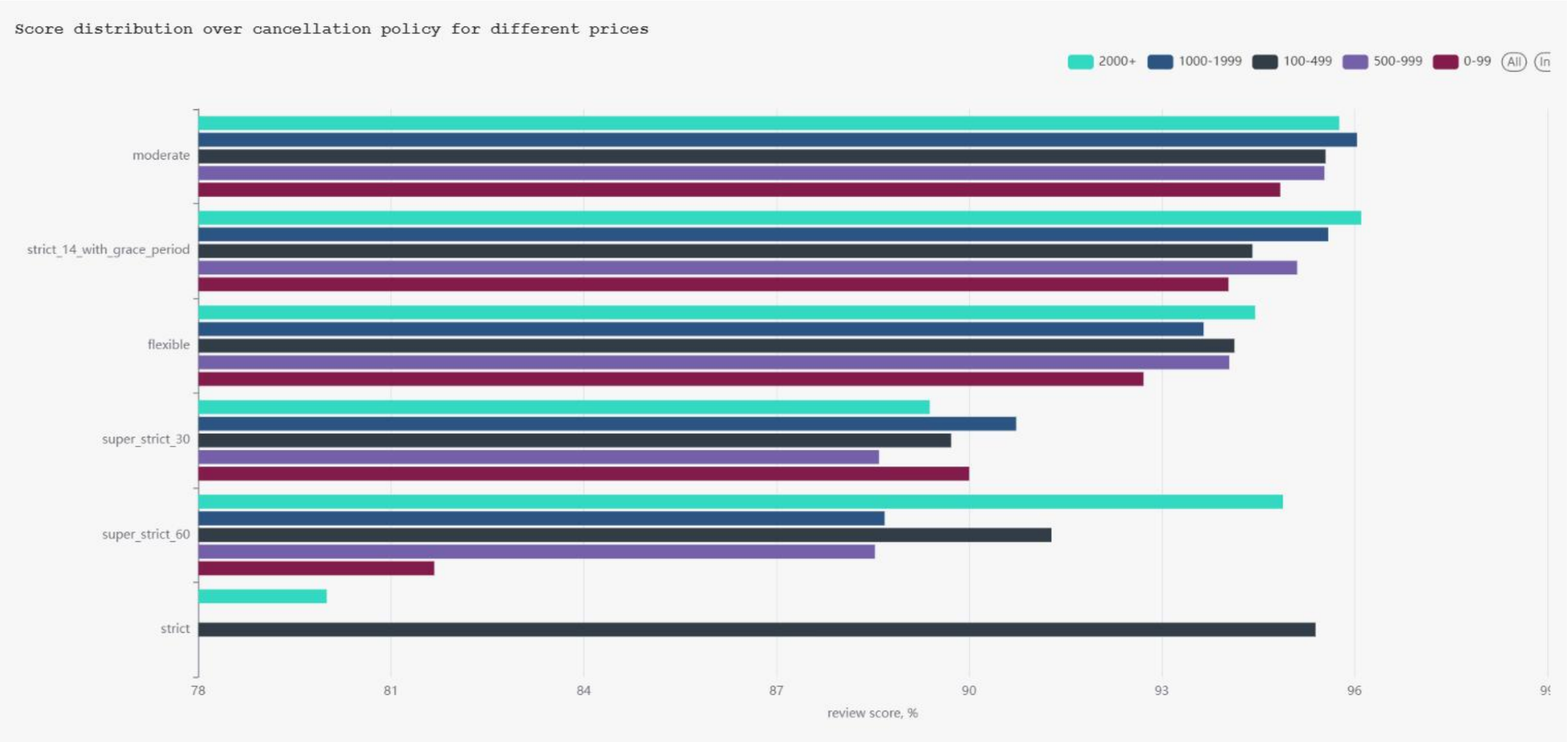


Insight 1.2

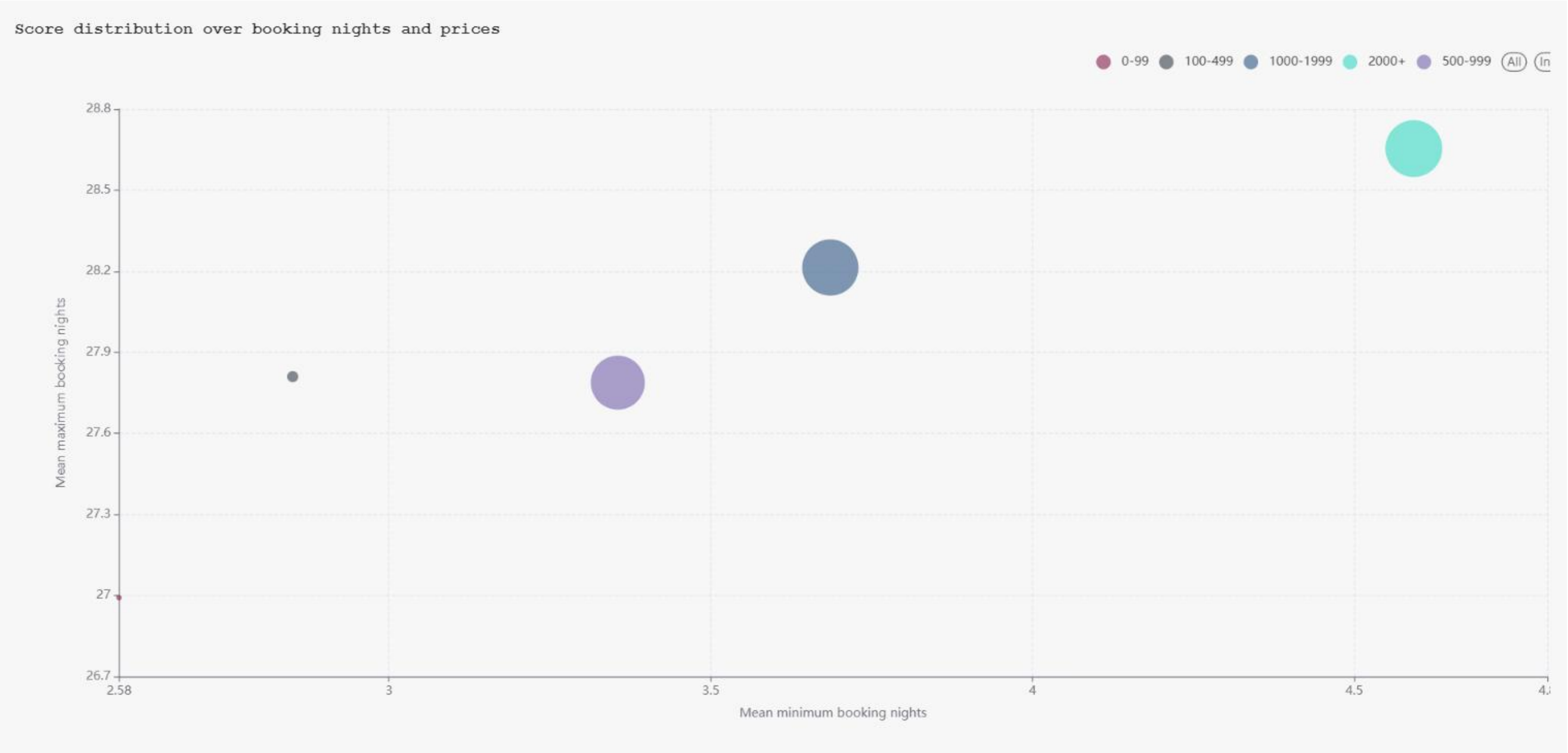
Score distribution over 10 most popular neighborhoods



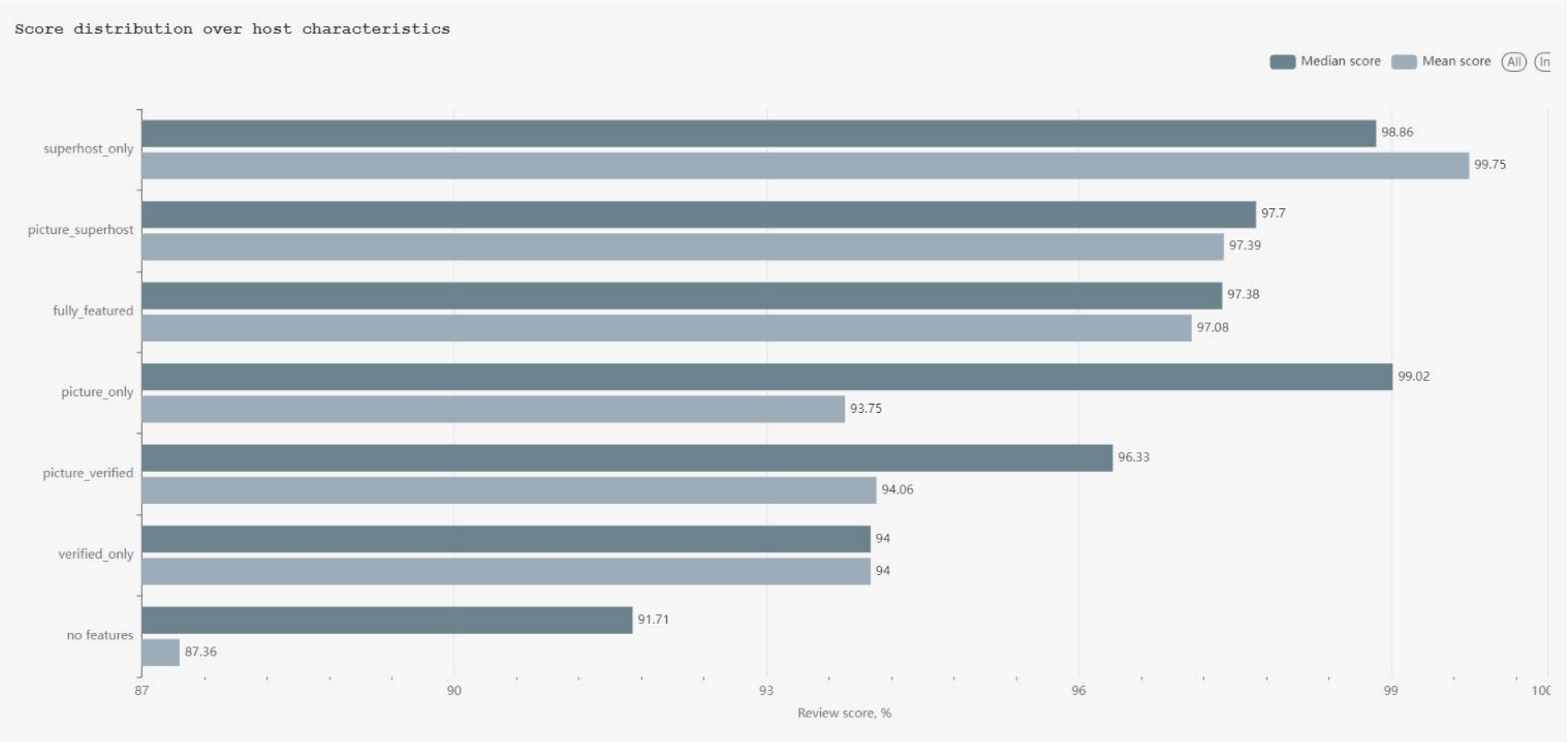
Insight 2.1



Insight 2.2

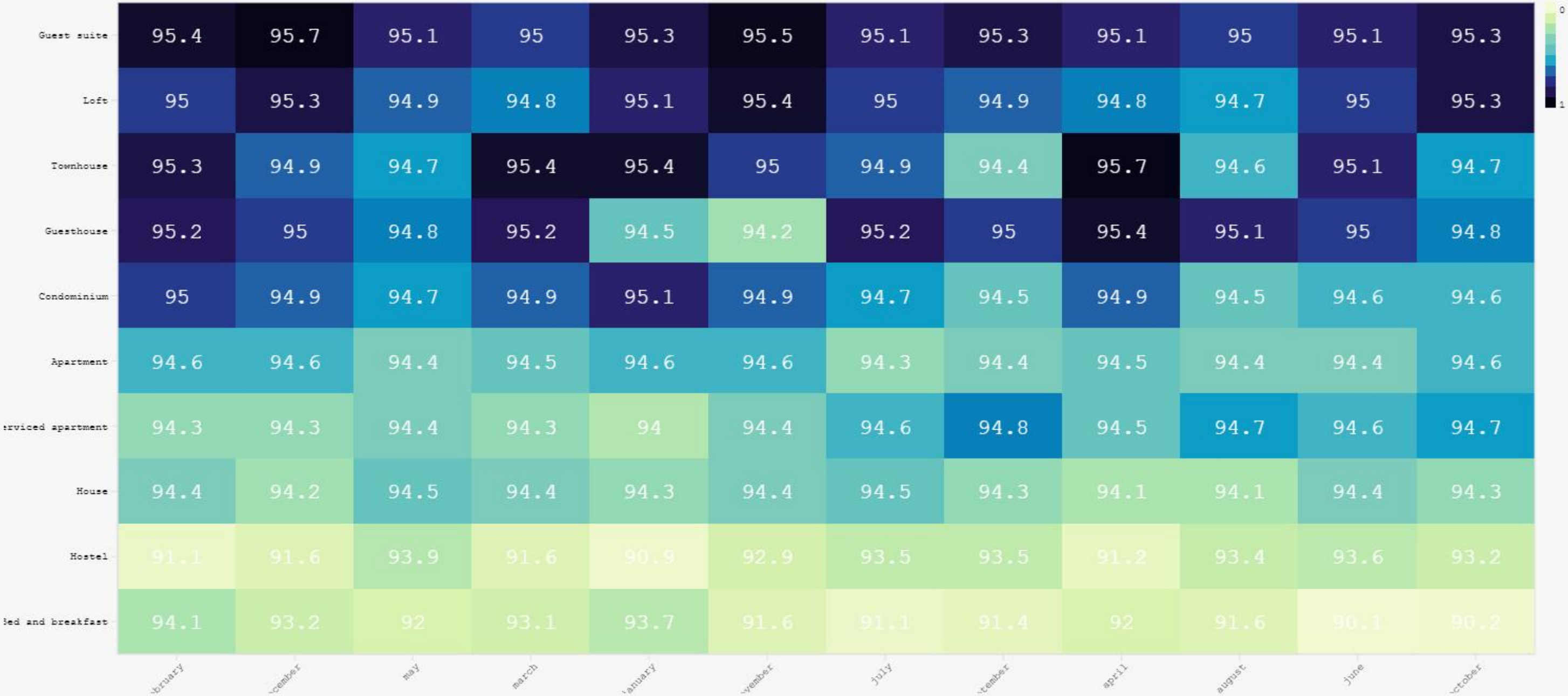


Insight 3



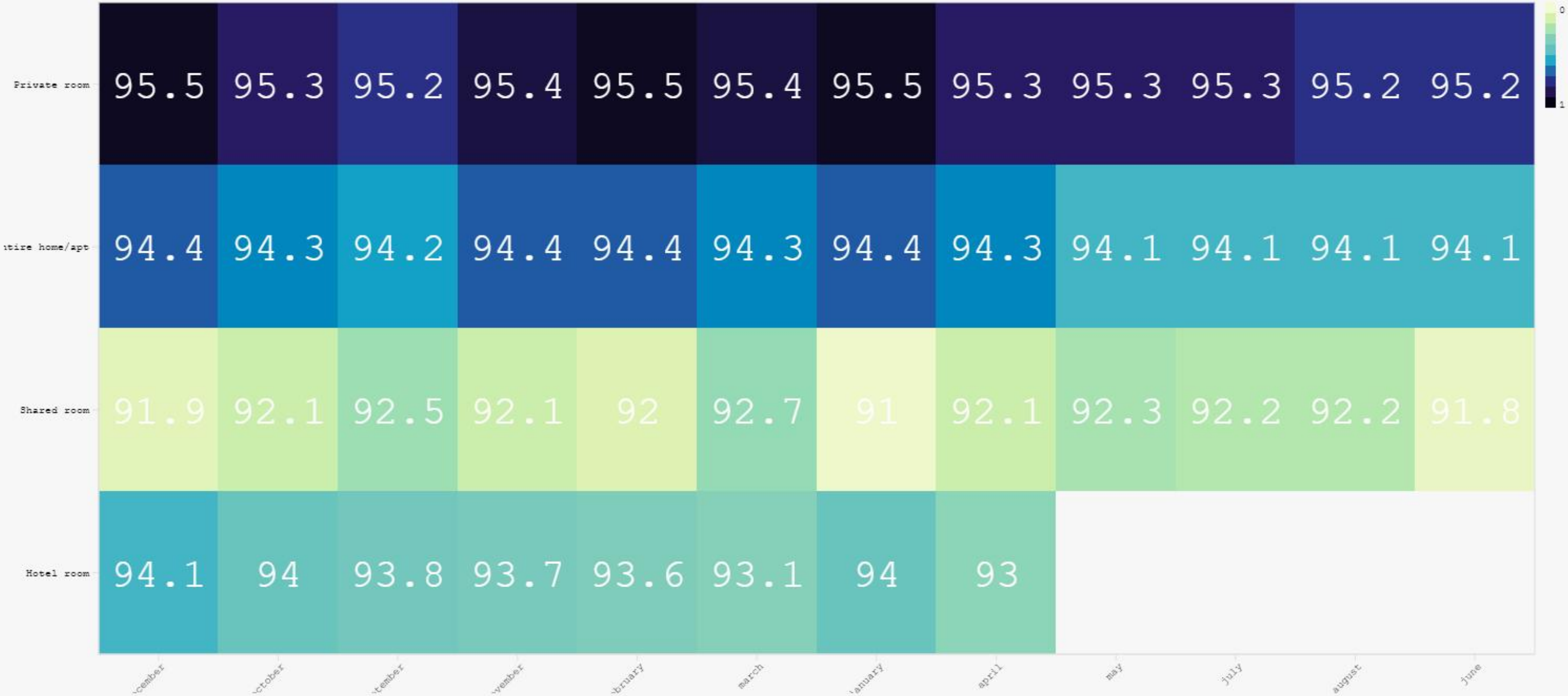
Insight 4.1

Score distribution over months and 10 most popular property types



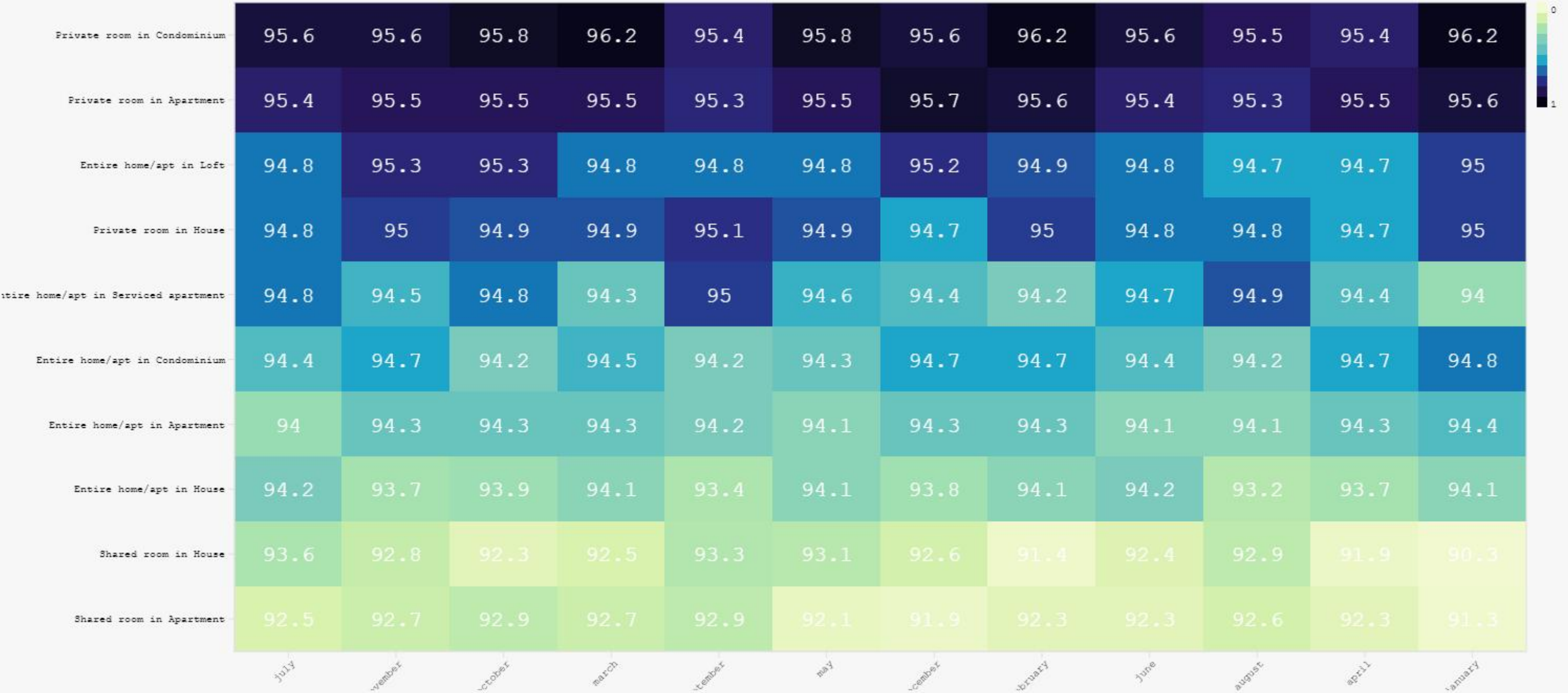
Insight 4.2

Score distribution over months and room types



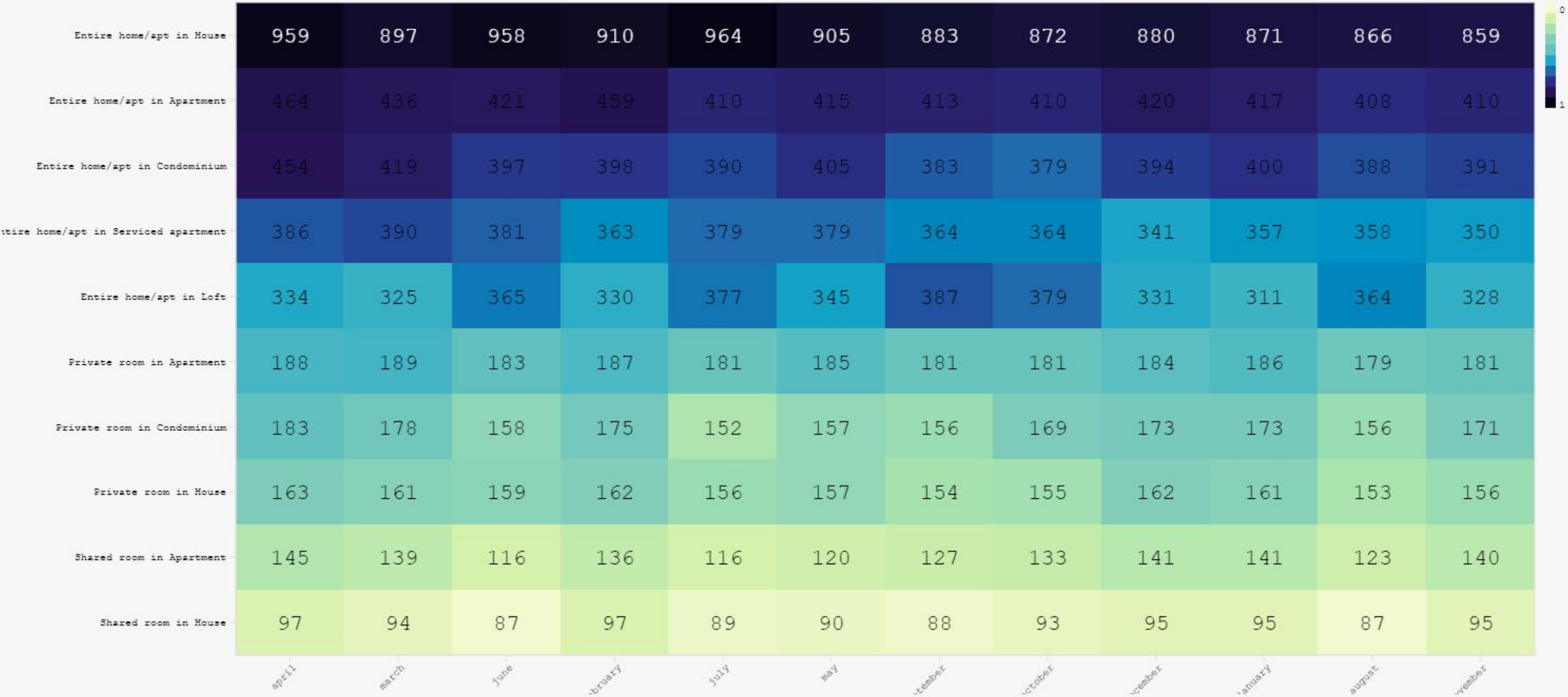
Insight 4.3

Score distribution over months and 10 most popular booking combinations

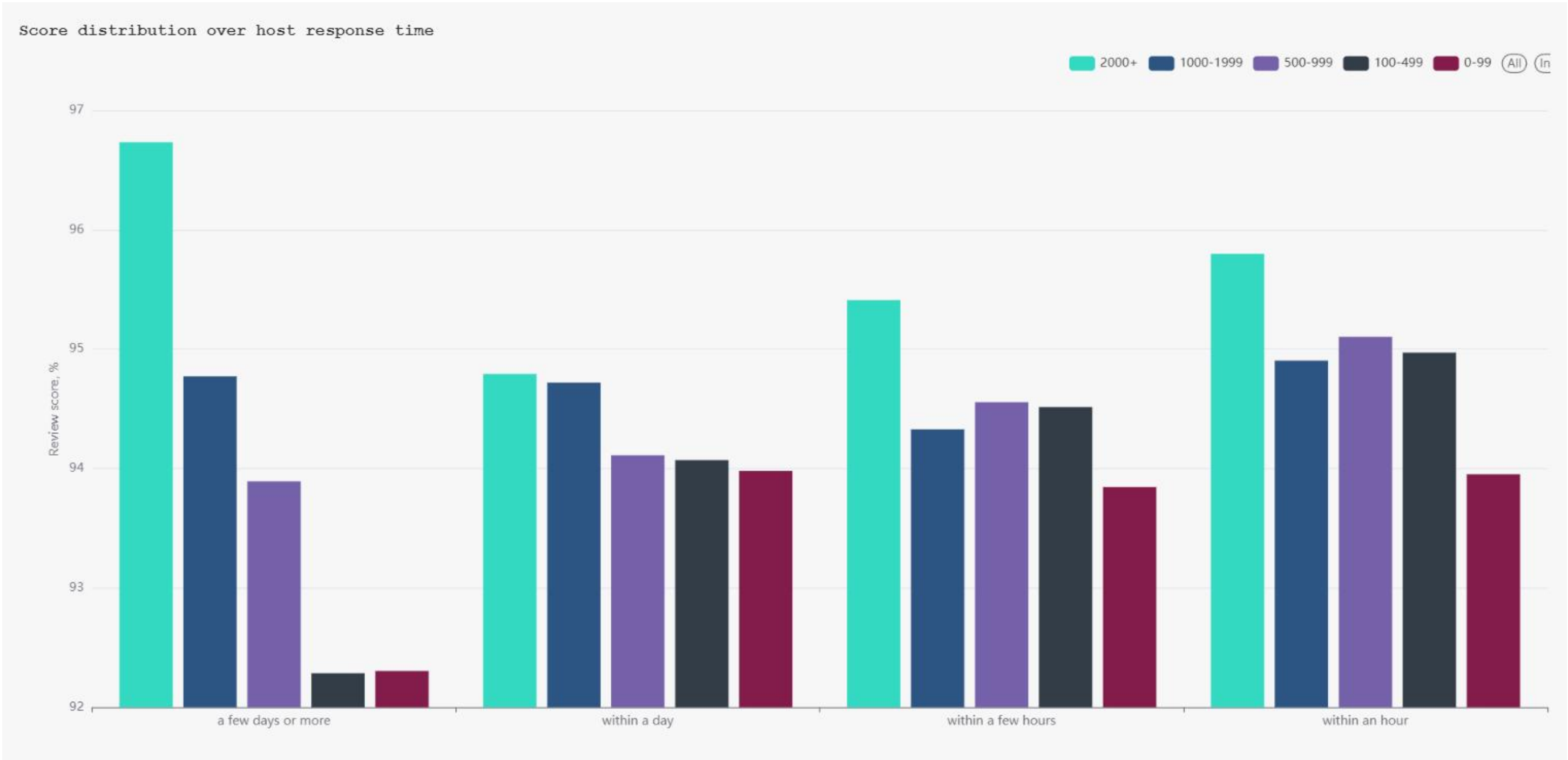


Insight 4.4

Price distribution over months and 10 most popular booking combinations

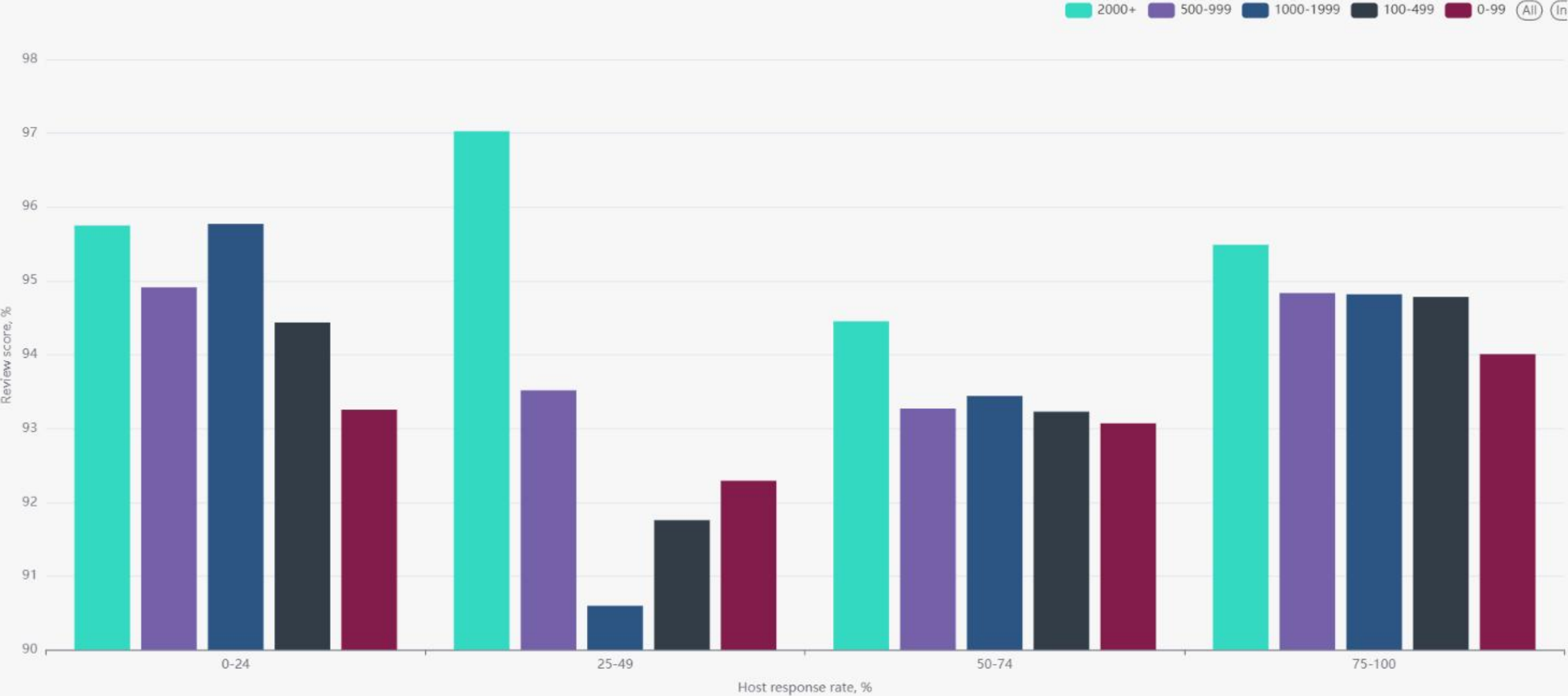


Insight 5.1

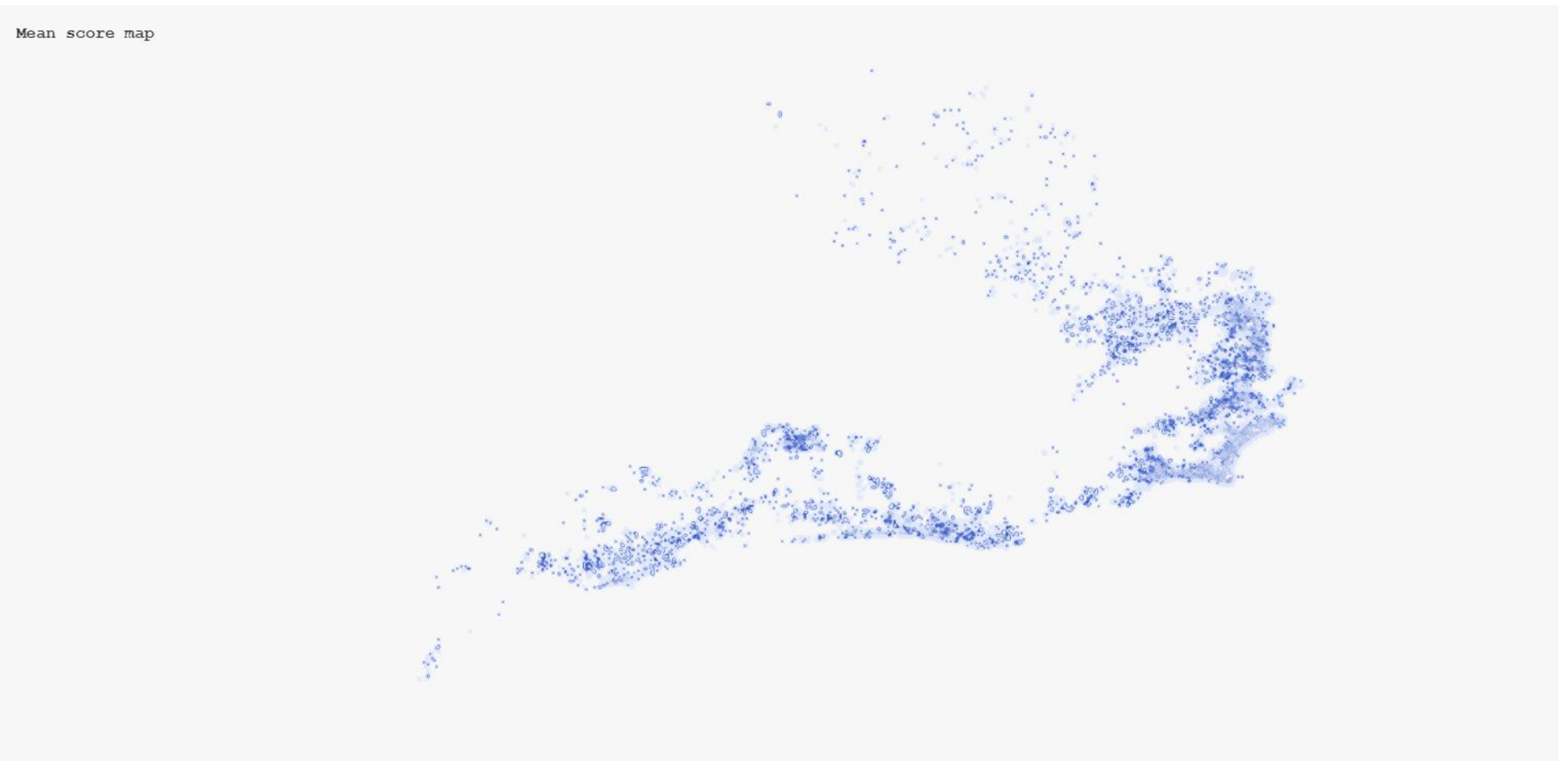


Insight 5.2

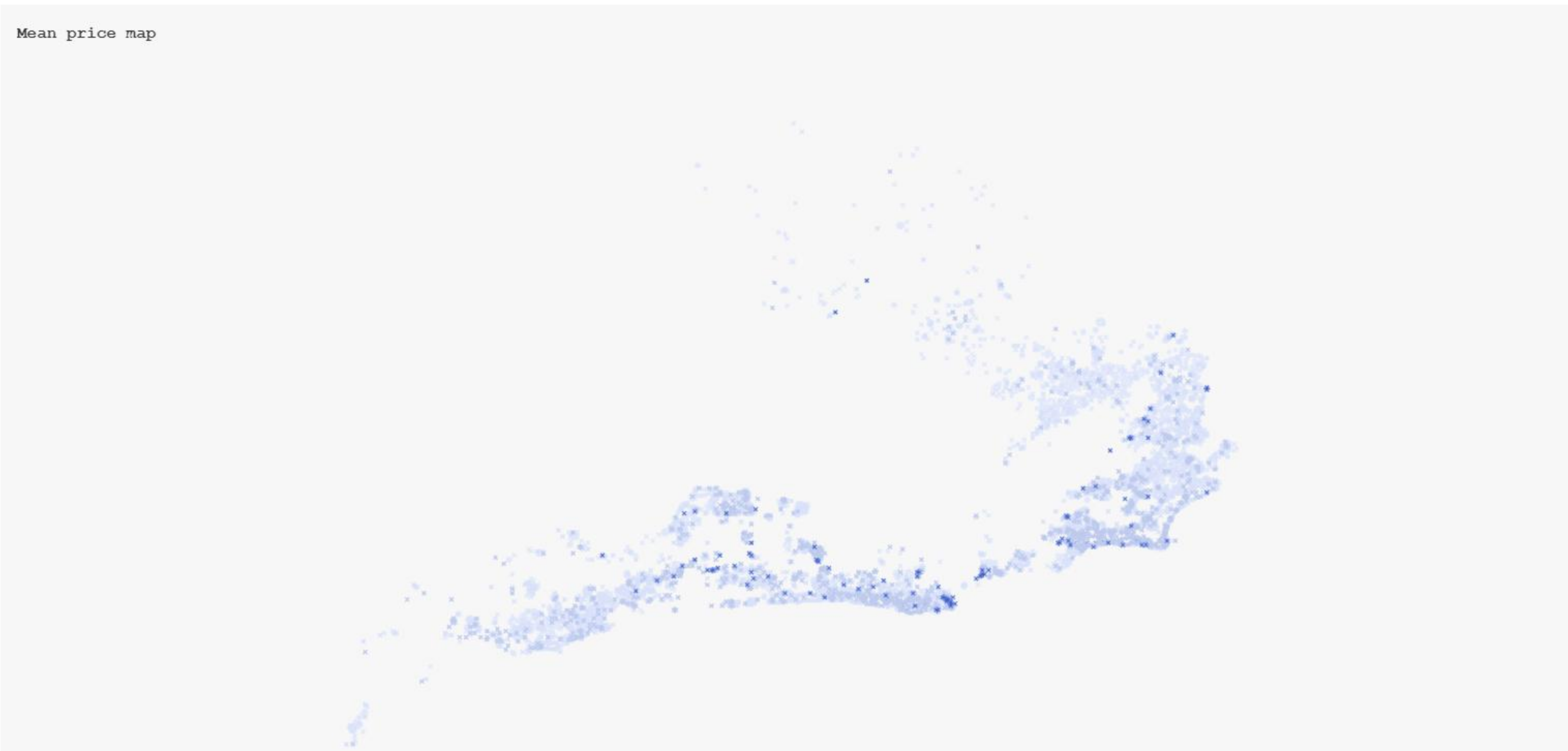
Score distribution over host response rate



Insight 6.1

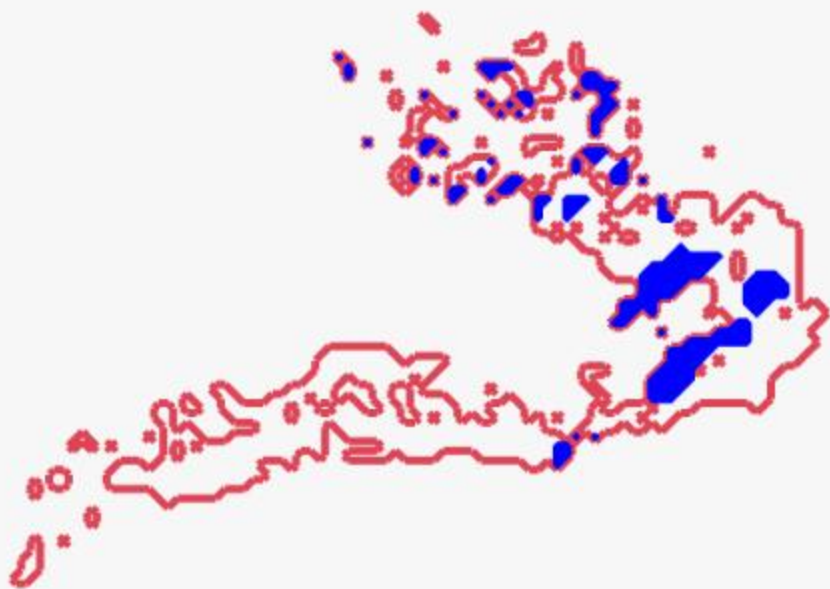


Insight 6.2



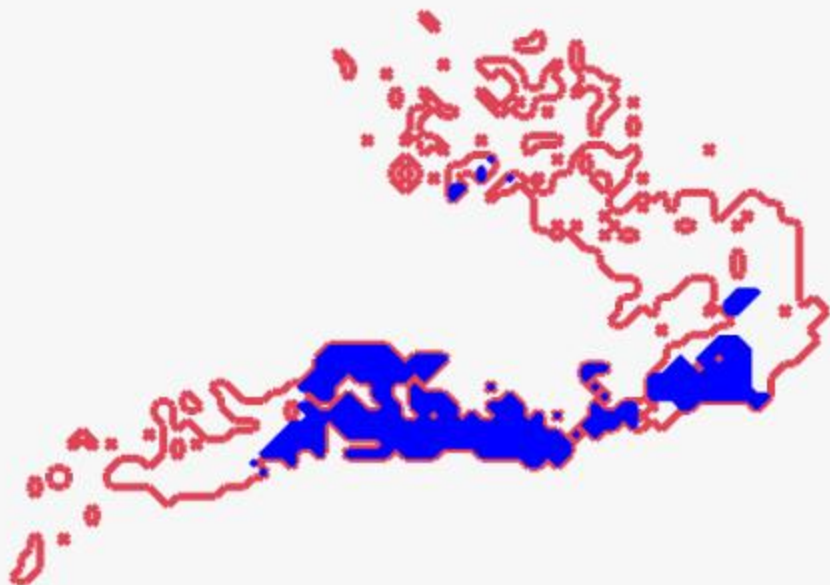
Insight 6.3

Top neighborhoods by mean score



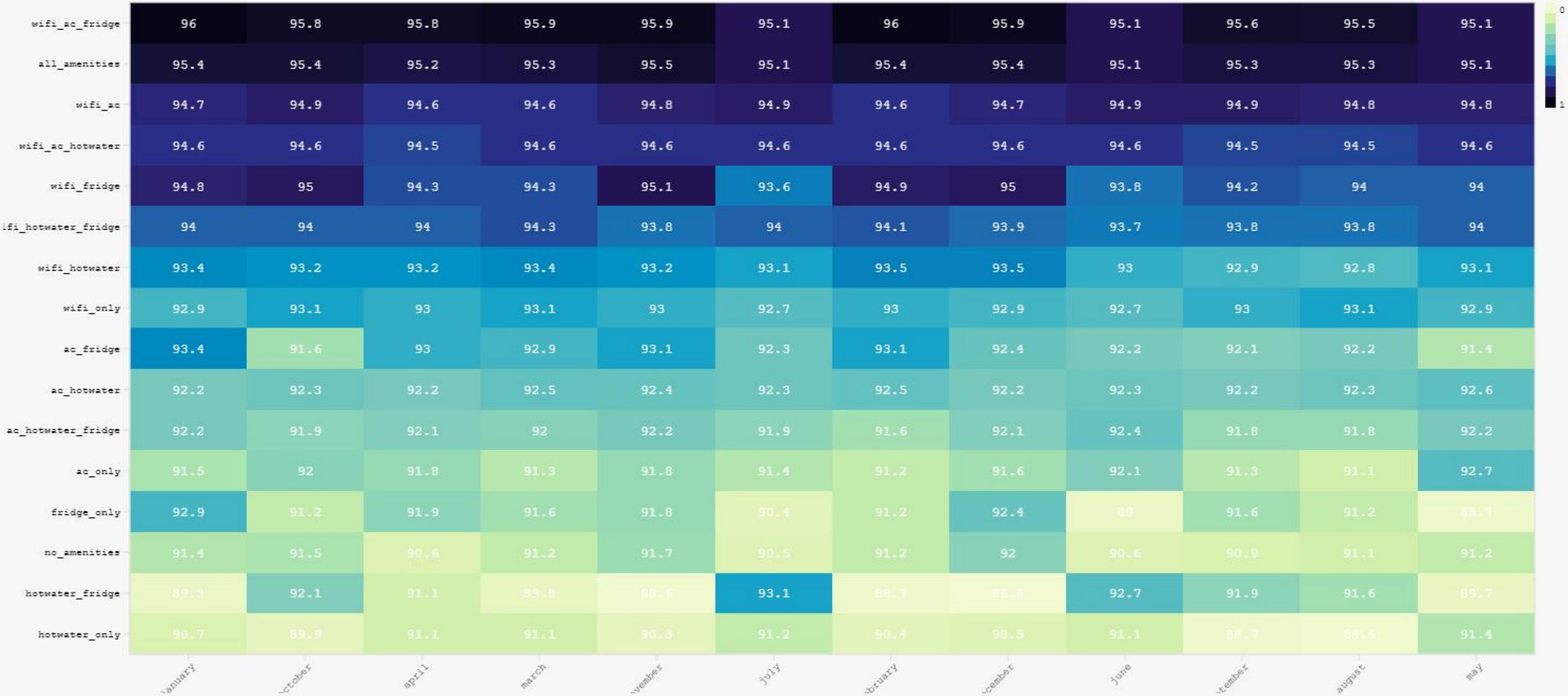
Insight 6.4

Top neighborhoods by mean price



Insight 7

Score distribution over amenities and months



Model Selection & Performance

Why These Models?:

Linear Regression:

- Pros: Fast, interpretable, sets a performance baseline.
- Use Case: Understand key drivers of ratings (e.g., "Price impacts scores linearly").
- Params Grid Search: Reg. Params = [0.01, 0.1, 1] and Elastic Net = [0, 0.5, 1]

Random Forest:

- Pros: Captures complex patterns (e.g., "Superhost status + location boost ratings").
- Outperformed Linear Regression
- Params Grid Search: Num Trees = [5, 10] and Max Depth = [5, 10]

Model	Num Trees	Max Depth	Reg. Params	Elastic Net
Linear Regression	-	-	0.01	0.5
Random Forest	10	10	-	-

Model	RMSE	MAE
Linear Regression	9.06	5.68
Random Forest	8.28	5.34

Discussion

We successfully implemented predictive models, and they demonstrate strong performance

Key Results:

- Random Forest outperformed Linear Regression, achieving a lower MAE and reliably estimating listing ratings—even without prior review data.
- Linear Regression struggled with extreme/high scores, revealing limitations in handling complex patterns.

How Our Model Drives Business Success

- Higher Revenue & Customer Satisfaction
- Accurate Rating Predictions → Ensures new listings (with no reviews) get fair visibility, boosting bookings and trust.
- Reduces "Cold Start" Penalty → Helps high-quality listings rank faster, increasing platform revenue.

Challenges we are faced

We had not reformulated the problem itself during the work. Initially we decided to use only subset of the original dataset.

However, we faced with problem with big number of undefined value: after filtering them the dataset becomes too small to satisfy project requirements. So we have decided to use the whole initial dataset, so after filtering undefined values, it size remains acceptable.

Team Roles

Nikita Yaneev	Vsevolod Klyushev	Dmitry Beresnev
<p>ML - engineer</p> <p>Prepared the data for the model, put together a full-fledged pipeline, trained the models, evaluated them, and identified the best solutions for our task.</p>	<p>Data engineer</p> <p>Implemented data collection, ingestion, preparation and storage pipeline. (Also wrote queries for data bucketing and partitioning.) Tested several compressing techniques and choose the best one.</p>	<p>Data analytic</p> <p>Performed EDA: built seven queries and analyzed the results; Created dashboards for Data Analysis, Insights and Solution models sections; Consult ML and Data engineers on approaches and solution choices based on data nature</p>

**Thanks for your
attention**