

Факторы богатства и успеха

Проект подготовили будущие
топы всех похожих списков:
Стекольников Яна
Губенко Олеся
Журавлев Артемий



Датасет “Billionaires Statistics Dataset” (35 признаков, 2539 строк):

rank: Ранг миллиардера по уровню богатства (место в рейтинге).	date: Дата сбора данных.
finalWorth: Итоговое состояние миллиардера в долларах США (чистый капитал).	state: Штат проживания (для США).
category: Категория или отрасль, в которой работает бизнес миллиардера.	residenceStateRegion: Регион/штат проживания.
personName: Имя миллиардера.	birthYear: Год рождения.
age: Возраст миллиардера.	birthMonth: Месяц рождения.
country: Страна проживания миллиардера.	birthDay: День рождения.
city: Город проживания миллиардера.	cpi_country: Индекс потребительских цен (ИПЦ) в стране миллиардера.
source: Источник богатства миллиардера.	cpi_change_country: Изменение ИПЦ в стране миллиардера.
industries: Отрасли, связанные с бизнес-интересами миллиардера.	gdp_country: ВВП страны миллиардера.
countryOfCitizenship: Страна гражданства миллиардера.	gross_tertiary_education_enrollment: Доля населения с высшим образованием в стране.
organization: Название организации/компании, связанной с миллиардером.	gross_primary_education_enrollment_country: Доля населения с начальным образованием в стране.
selfMade: Указывает, является ли миллиардер self-made.	life_expectancy_country: Ожидаемая продолжительность жизни в стране.
gender: Пол миллиардера.	tax_revenue_country_country: Налоговые поступления в стране.
birthDate: Дата рождения миллиардера.	total_tax_rate_country: Общая налоговая ставка в стране.
lastName: Фамилия миллиардера.	population_country: Население страны.
firstName: Имя миллиардера.	latitude_country: Географическая широта страны.
title: Должность в компании.	longitude_country: Географическая долгота страны.

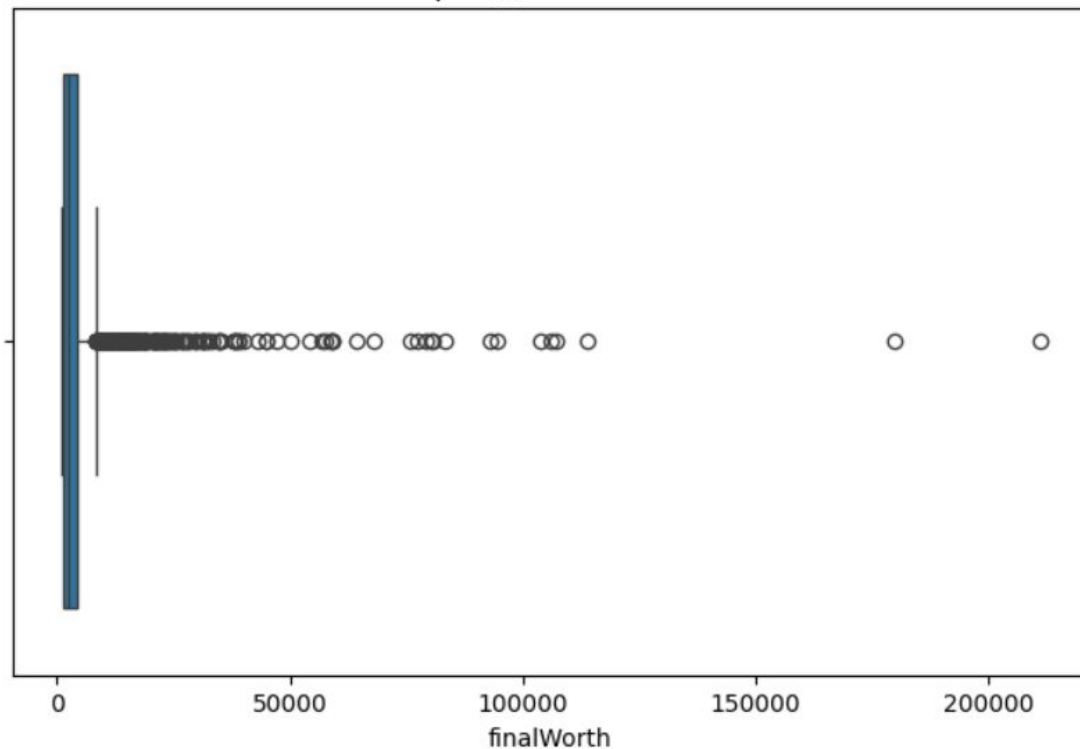
Загрузка и предварительный анализ(EDA)

Пропущенные значения % от общего			title		2301	87.159091
rank	0	0.000000	date	0	0.000000	
finalWorth	0	0.000000	state	1887	71.477273	
category	0	0.000000	residenceStateRegion	1893	71.704545	
personName	0	0.000000	birthYear	76	2.878788	
age	65	2.462121	birthMonth	76	2.878788	
country	38	1.439394	birthDay	76	2.878788	
city	72	2.727273	cpi_country	184	6.969697	
source	0	0.000000	cpi_change_country	184	6.969697	
industries	0	0.000000	gdp_country	164	6.212121	
countryOfCitizenship	0	0.000000	gross_tertiary_education_enrollment	182	6.893939	
organization	2315	87.689394	gross_primary_education_enrollment_country	181	6.856061	
selfMade	0	0.000000	life_expectancy_country	182	6.893939	
status	0	0.000000	tax_revenue_country_country	183	6.931818	
gender	0	0.000000	total_tax_rate_country	182	6.893939	
birthDate	76	2.878788	population_country	164	6.212121	
lastName	0	0.000000	latitude_country	164	6.212121	
firstName	3	0.113636	longitude_country	164	6.212121	

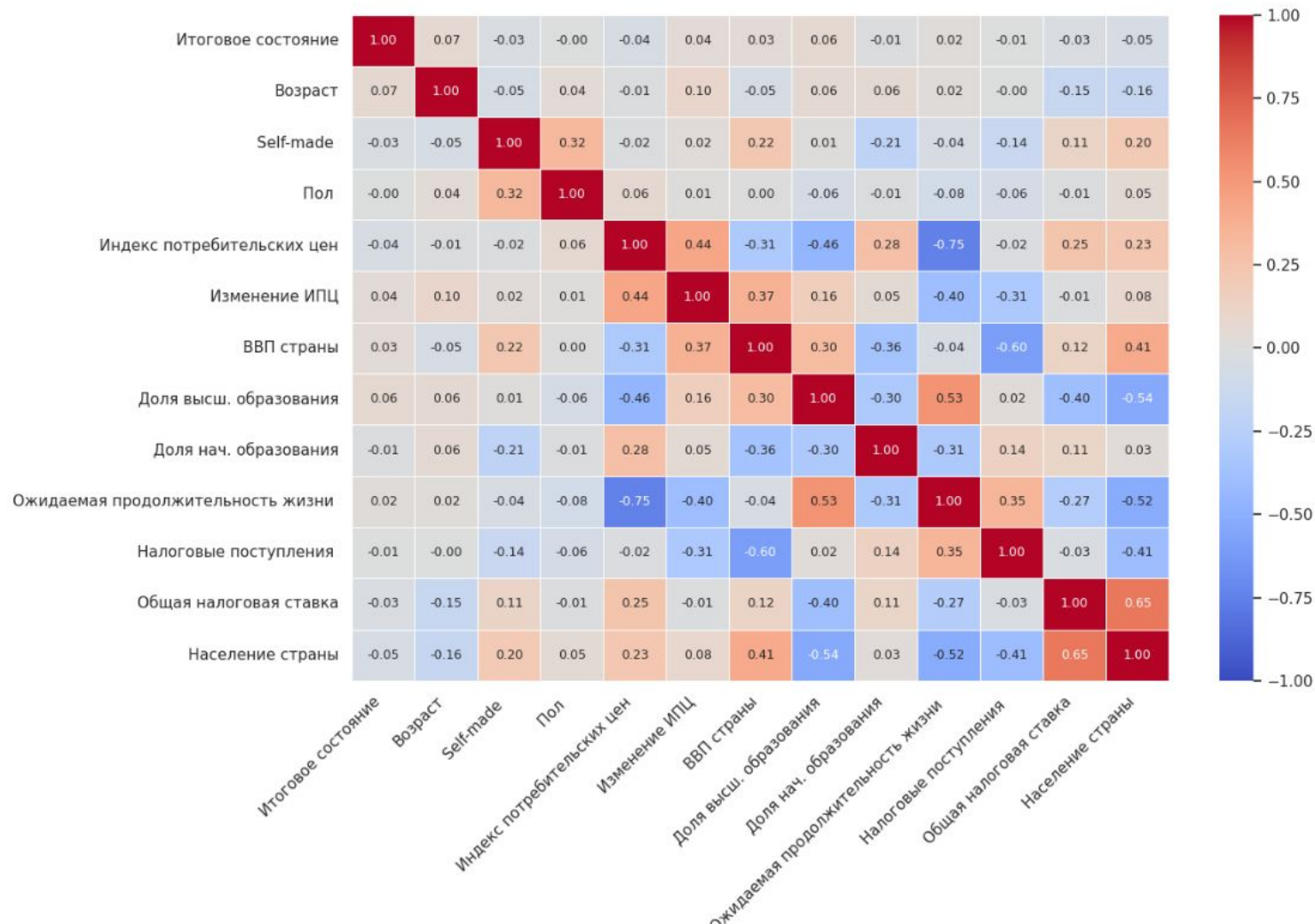
Исследование распределений ключевых признаков

Анализ распределения богатства

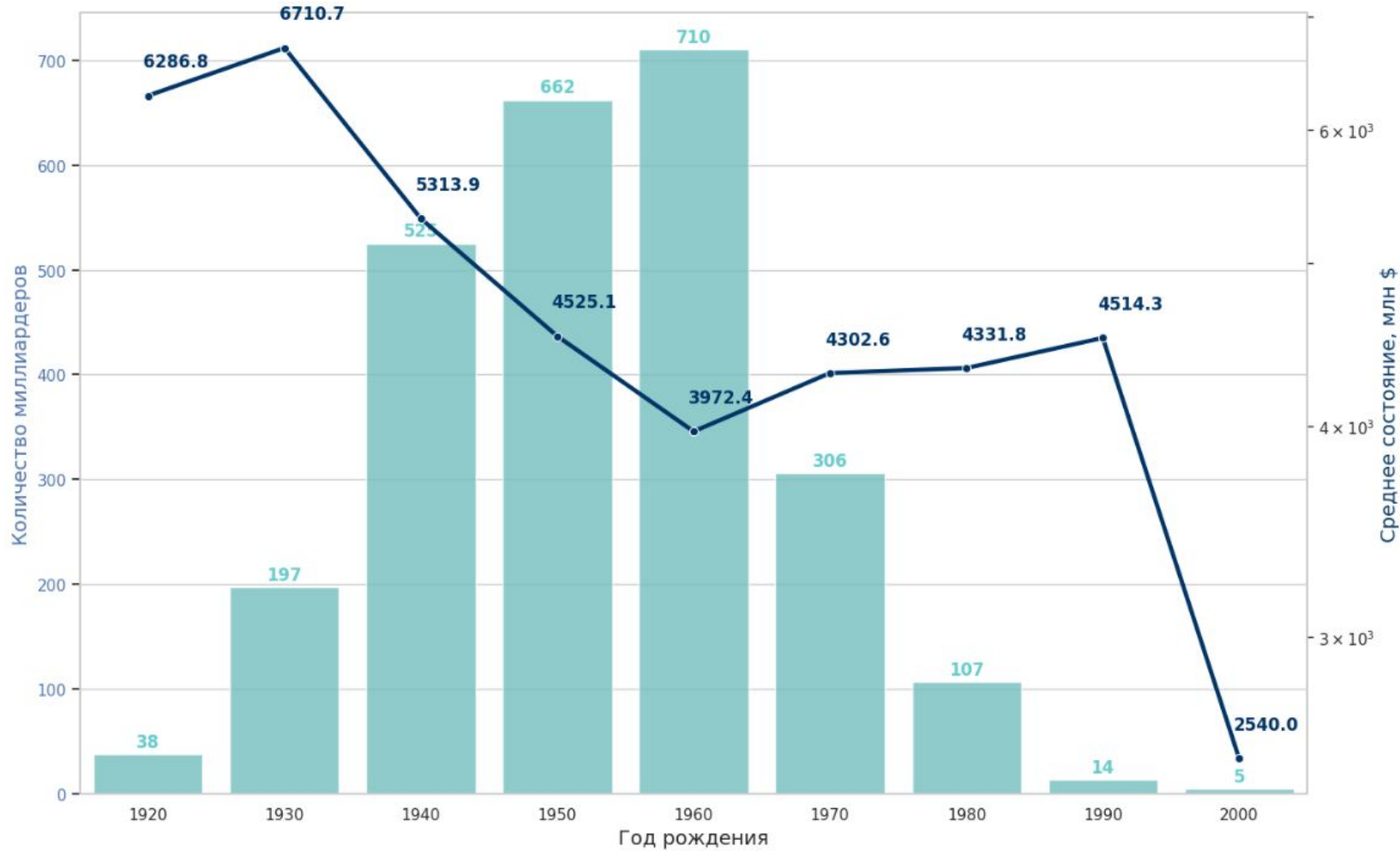
Boxplot для состояния



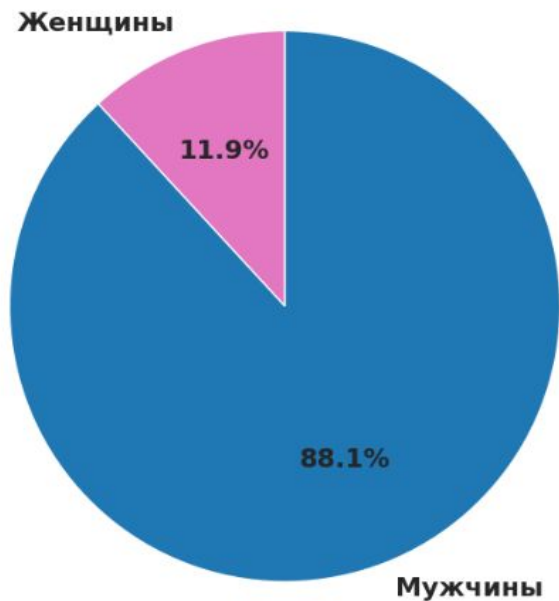
Корреляционная матрица



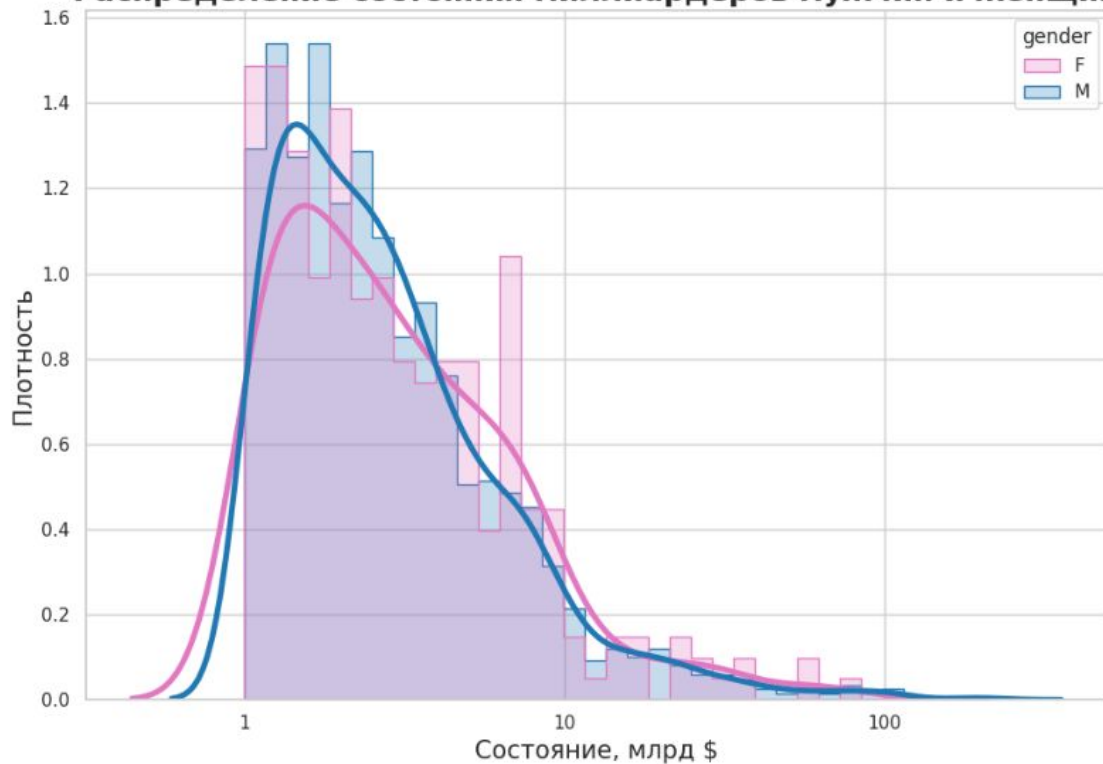
Анализ возраста миллиардеров: количество и среднее состояние



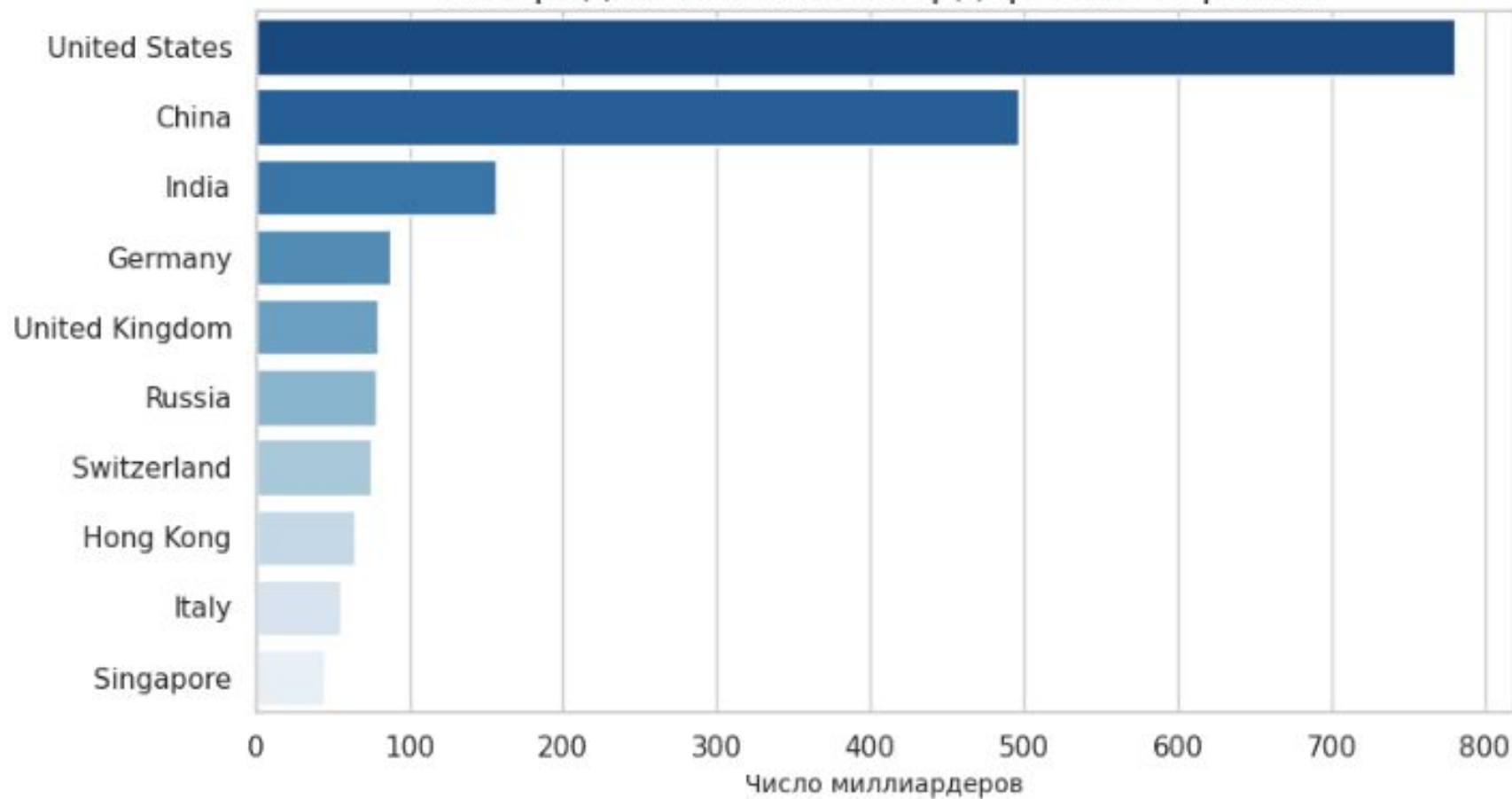
Соотношение по полу

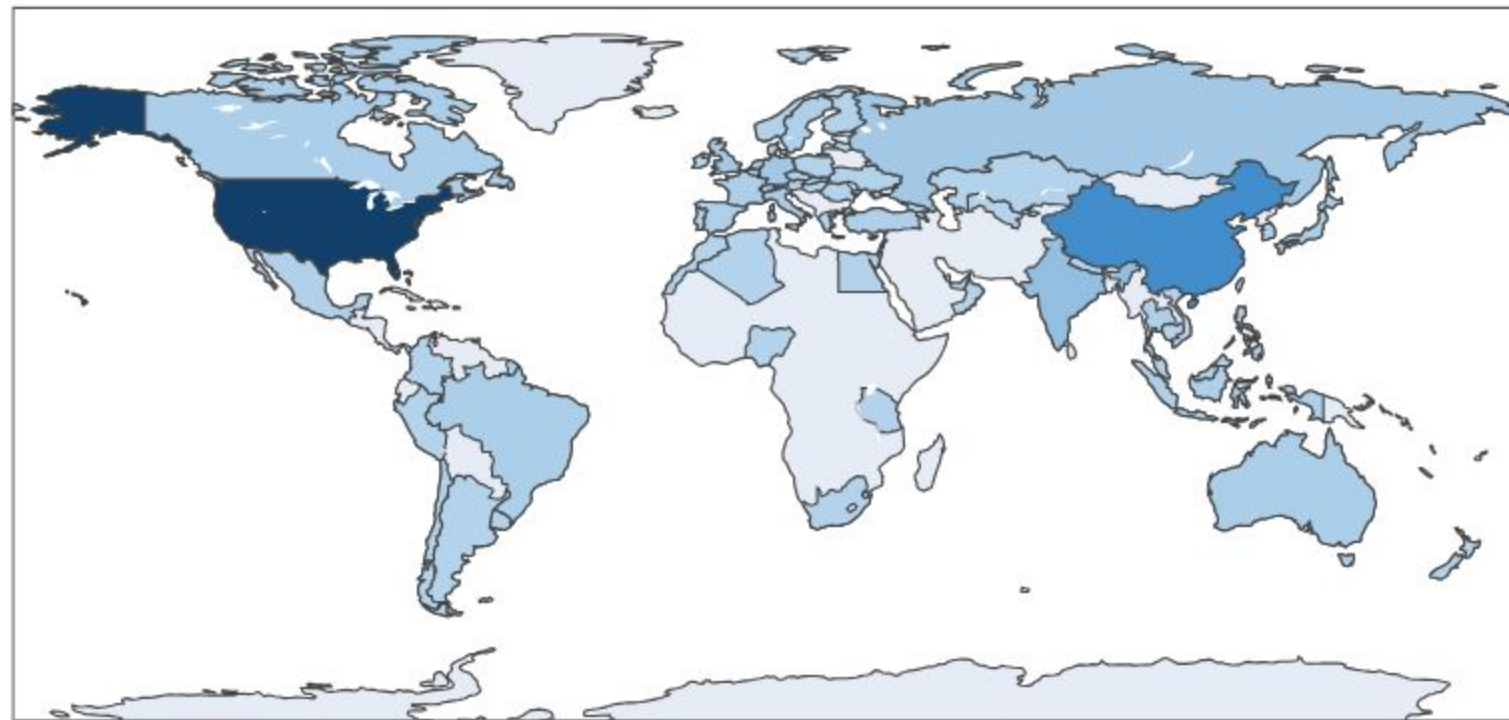


Распределение состояния миллиардеров мужчин и женщин

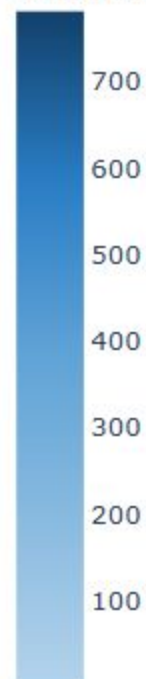


Распределение миллиардеров по странам

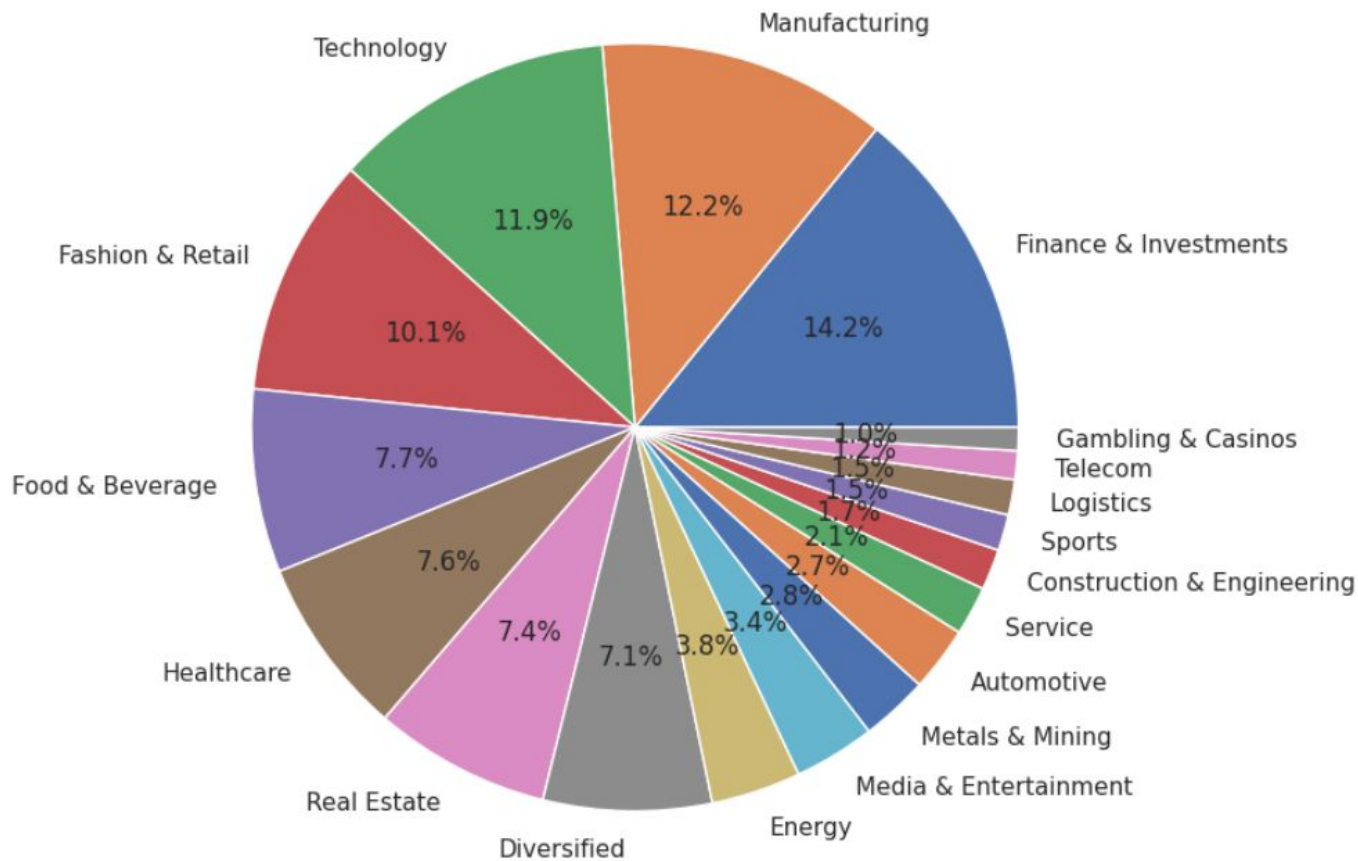




Количество
миллиардов

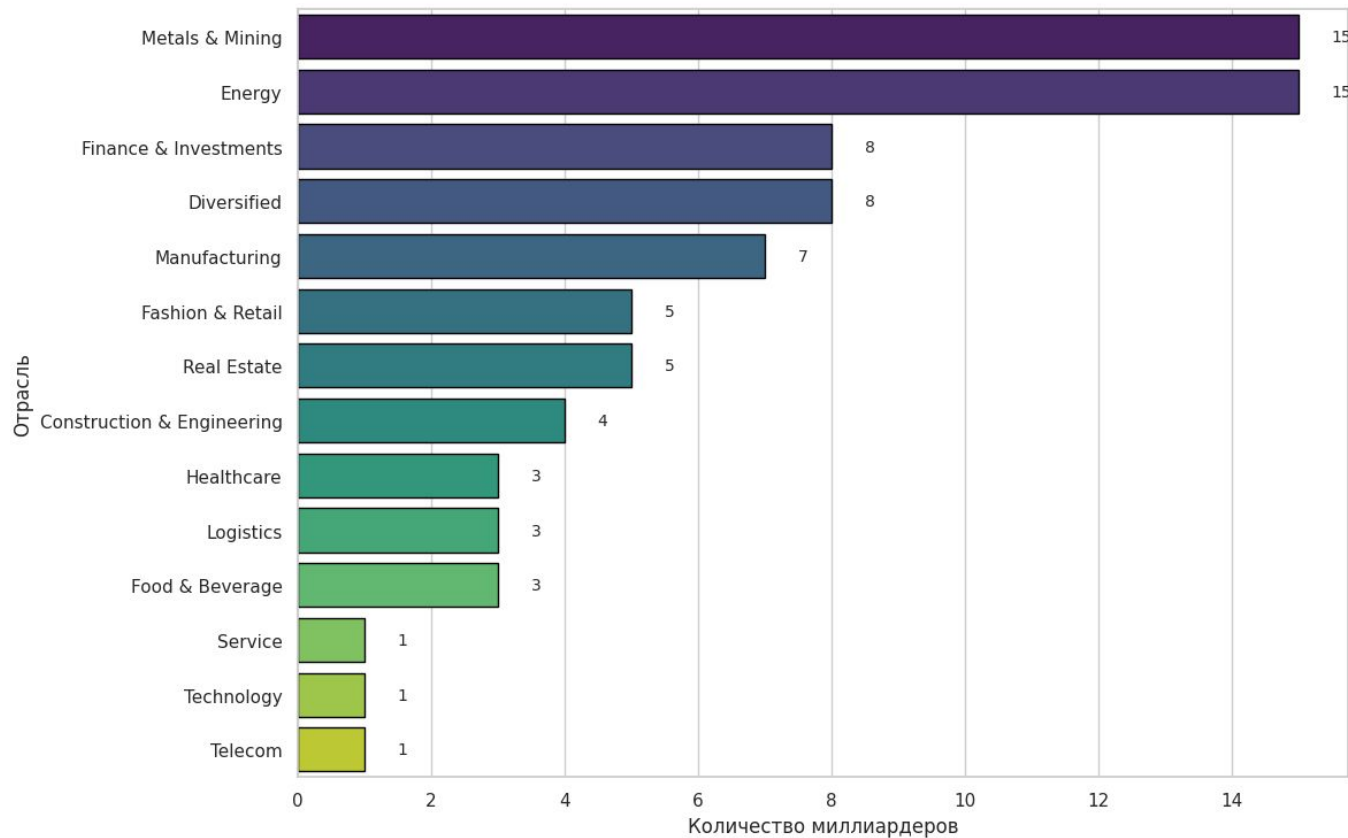


Распределение отраслей среди миллиардеров

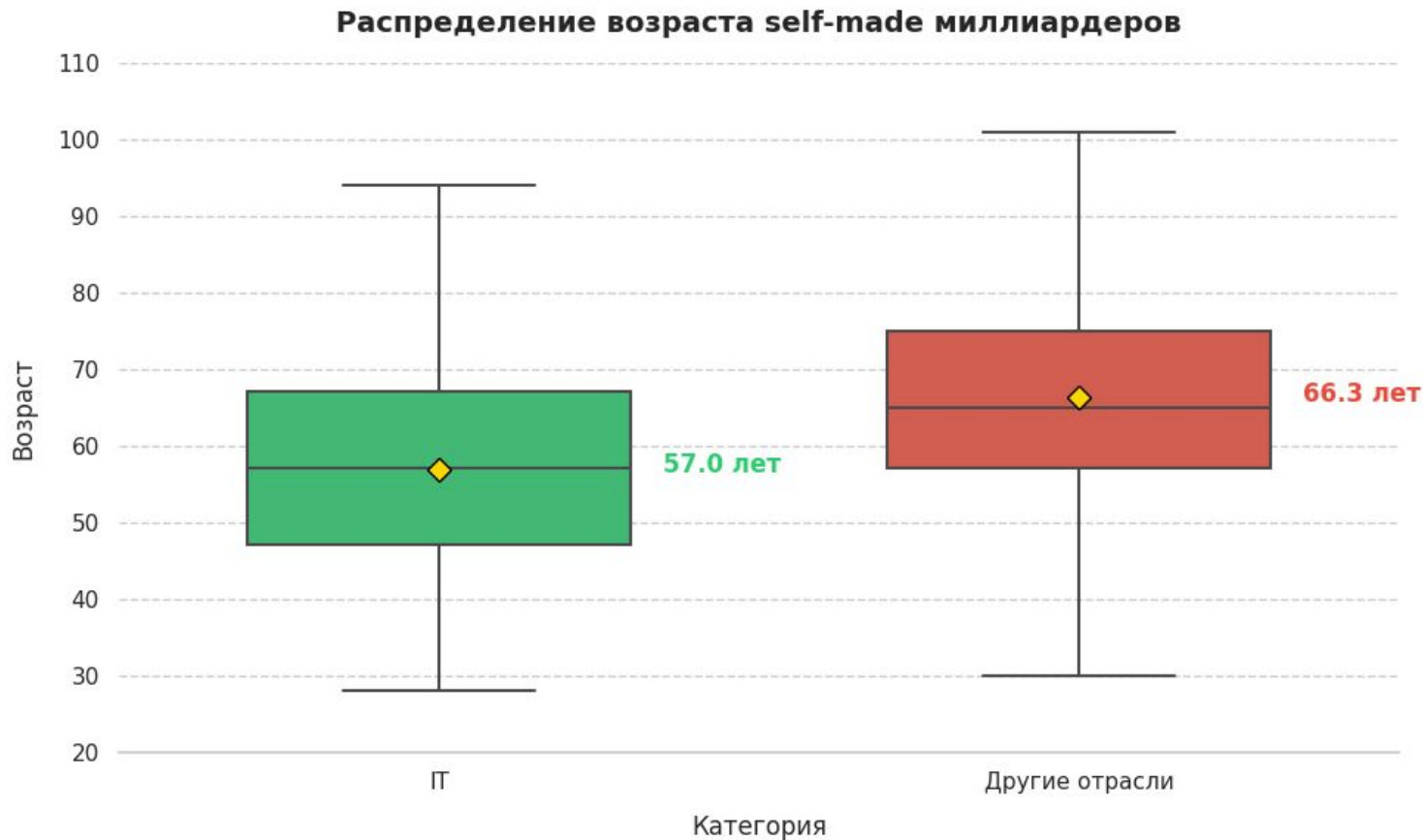


“Большинство русских сверхбогатых людей заработали состояние в металлургической промышленности”

Распределение российских миллиардеров по отраслям



“Средний возраст self-made миллиардеров в IT ниже, чем средний возраст self-made миллиардеров в традиционных отраслях”



“В среднем состояние IT-миллиардеров больше, чем состояние миллиардеров из других областей”

Для проверки данной гипотезы мы используем метод bootstrat.

Суть метода:

Бутстреп — это компьютерный статистический метод, который позволяет оценить точность измерений и проверить гипотезы путем многократного случайного повторного отбора данных с заменой из исходной выборки.

Как работает:

1. Из исходных данных берутся случайные подвыборки того же размера (с возможностью повторения значений)
2. Для каждой подвыборки вычисляется нужная статистика (среднее, разница средних и др.)
3. Процесс повторяется тысячи раз (обычно 10 000+)
4. По полученному распределению статистик определяют:
 - Точечную оценку параметра
 - Доверительные интервалы
 - Статистическую значимость

$M(IT)$ — среднее состояние IT-миллиардеров в генеральной совокупности,
 $M(Other)$ — среднее состояние миллиардеров из других отраслей в генеральной совокупности.

Нулевая гипотеза (H_0):

$H_0: M(IT) \leq M(Other)$

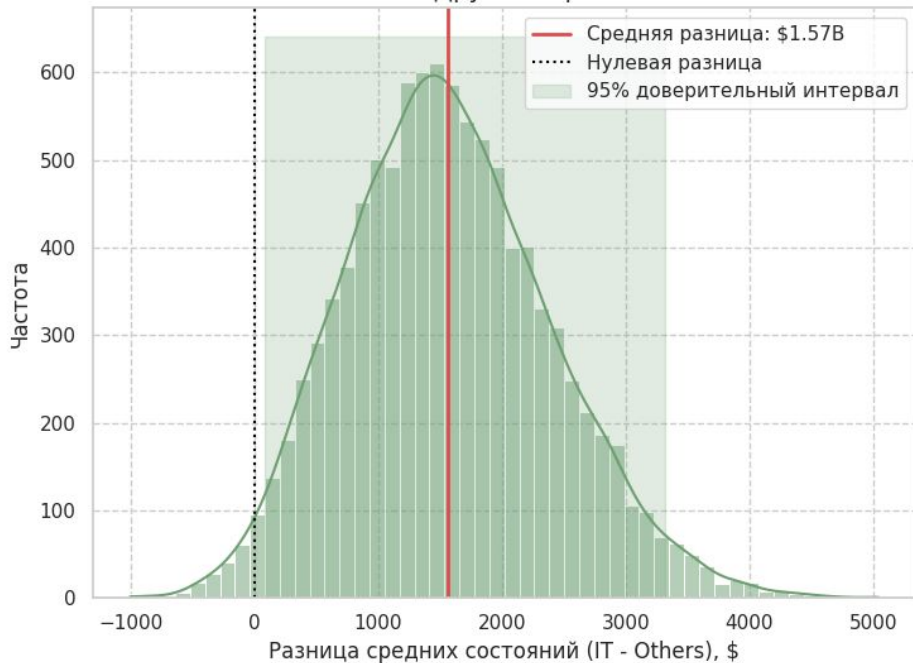
«Среднее состояние IT-миллиардеров не превышает среднее состояние миллиардеров из других отраслей».

Альтернативная гипотеза (H_1):

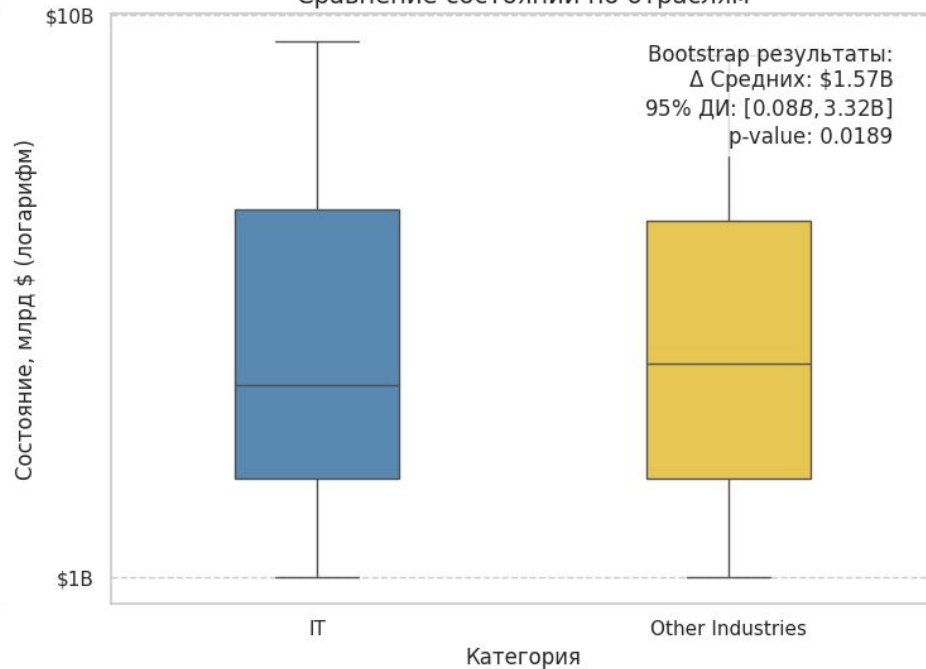
$H_1: M(IT) > M(Other)$

«Среднее состояние IT-миллиардеров больше среднего состояния миллиардеров из других отраслей».

Бутстреп-распределение разницы средних
IT vs Другие отрасли



Сравнение состояний по отраслям

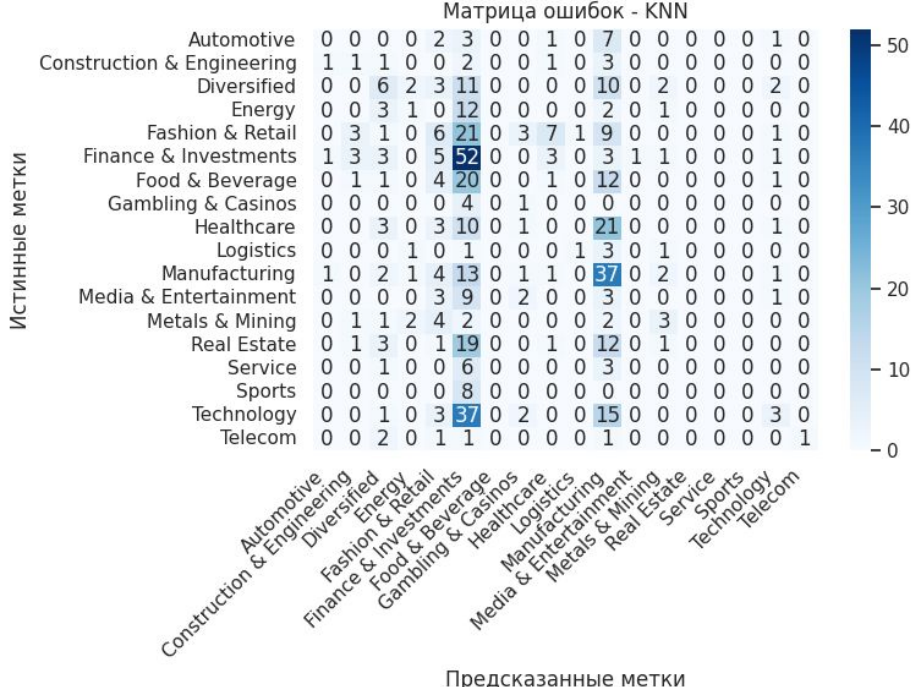
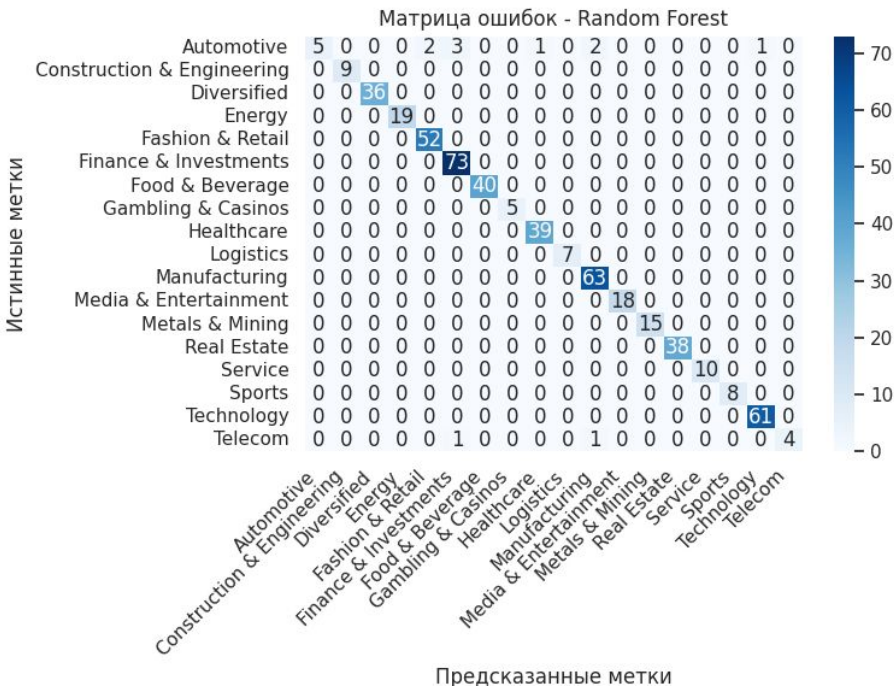


Вывод: Отвергаем Н0. IT-миллиардеры богаче ($p < 0.05$).

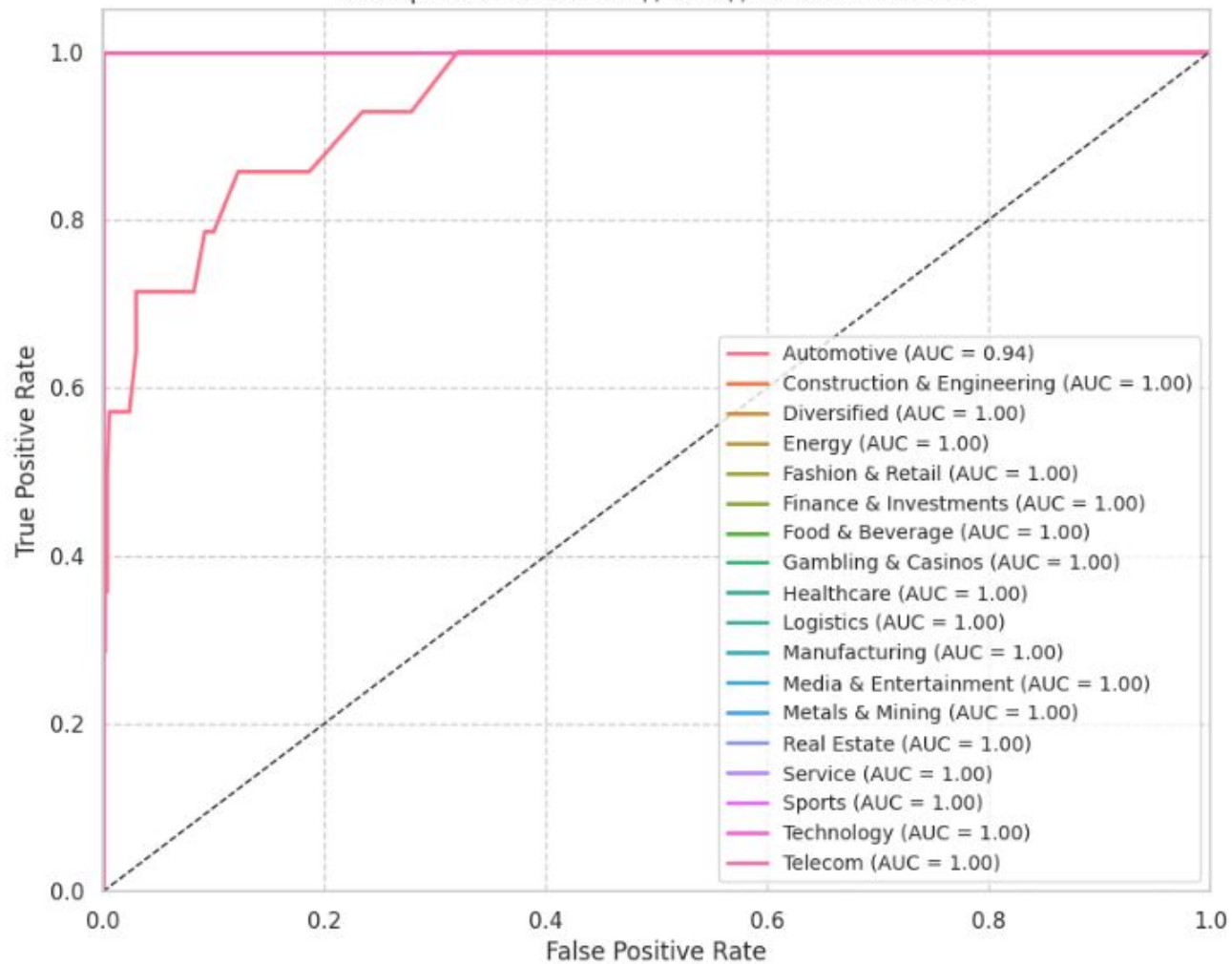
Классификация отрасли - Random Forest + kNN

Лучшие параметры Random Forest: {'max_depth': None, 'min_samples_split': 2, 'n_estimators': 200}

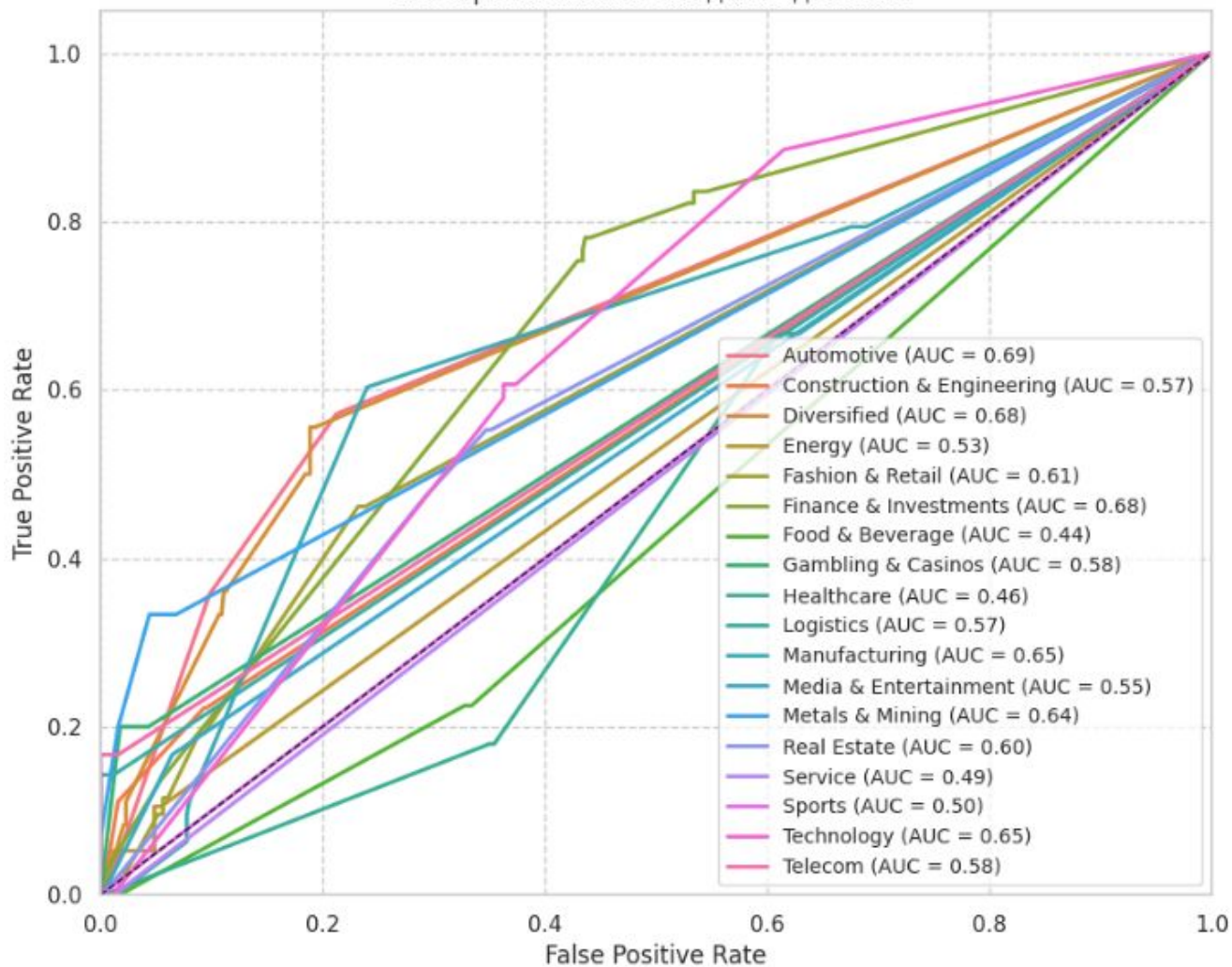
Лучшие параметры KNN: {'knn__n_neighbors': 3, 'knn__weights': 'uniform'}



ROC-кривые по классам для модели Random Forest



ROC-кривые по классам для модели KNN



Прогнозирование состояния

Используем:

Linear Regression:

$R^2 = -0.045$,

RMSE = 5204,

MAE = 3350

Random Forest:

$R^2 = -0.107$,

RMSE = 5355,

MAE = 3341

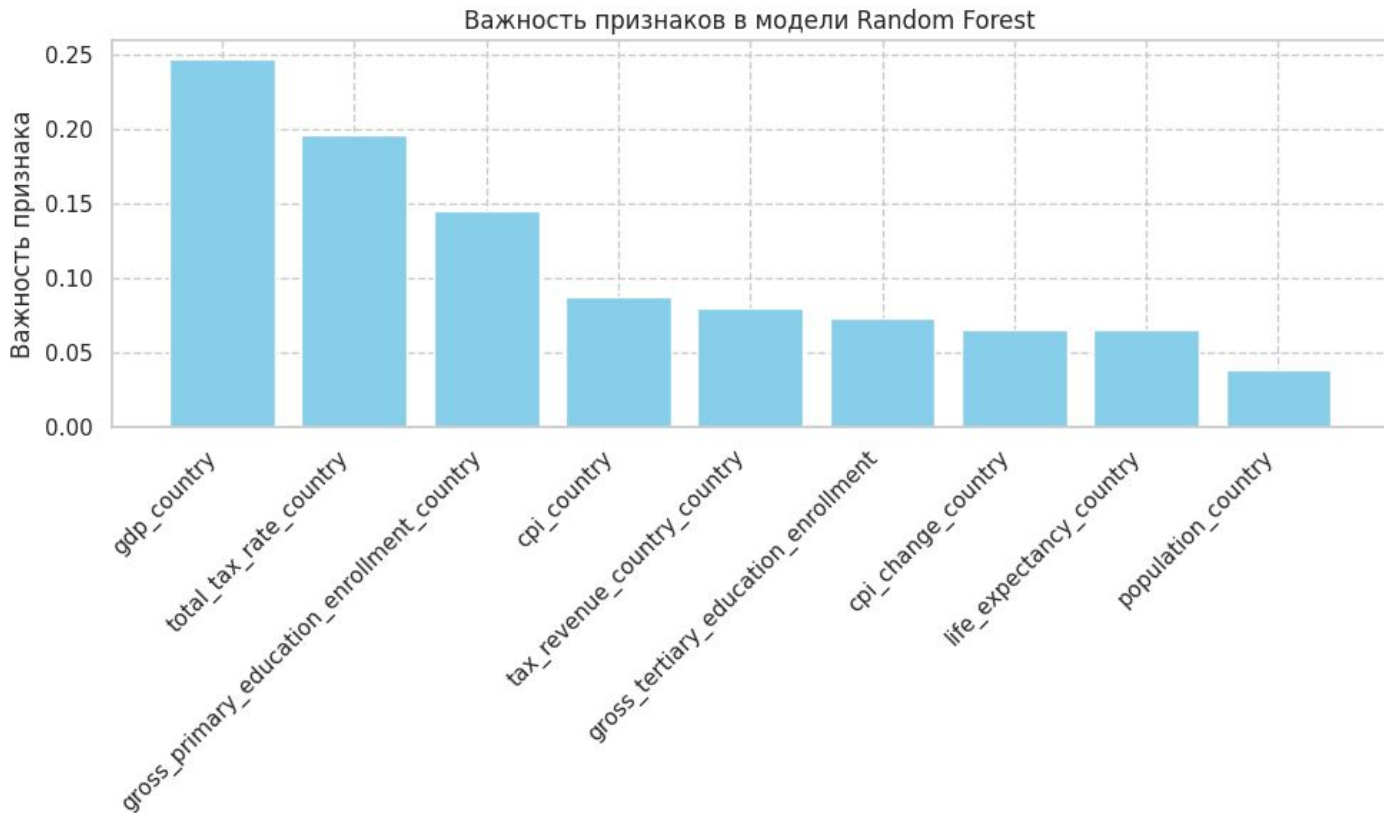
Лучшие параметры

RandomForest:

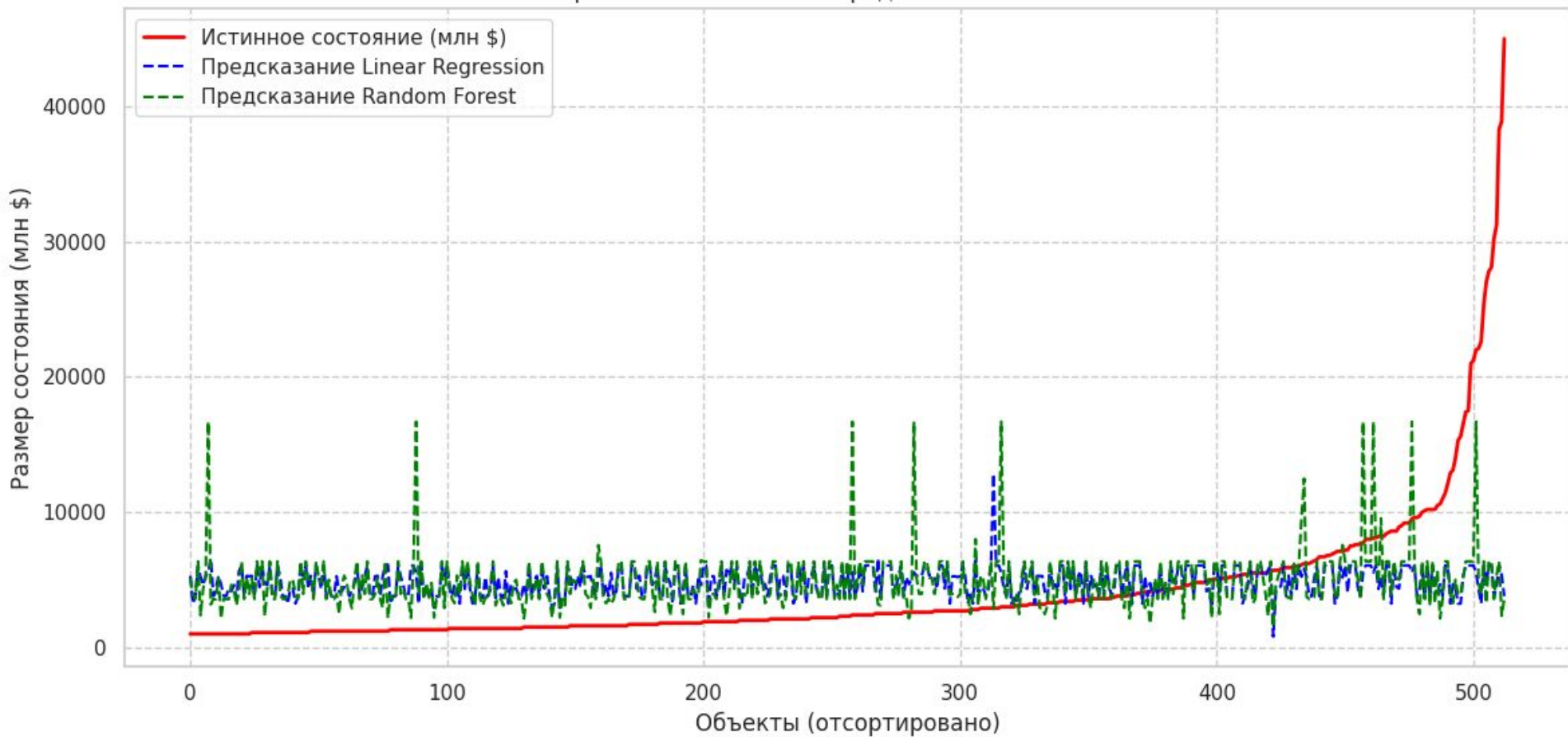
'max_depth': 10,

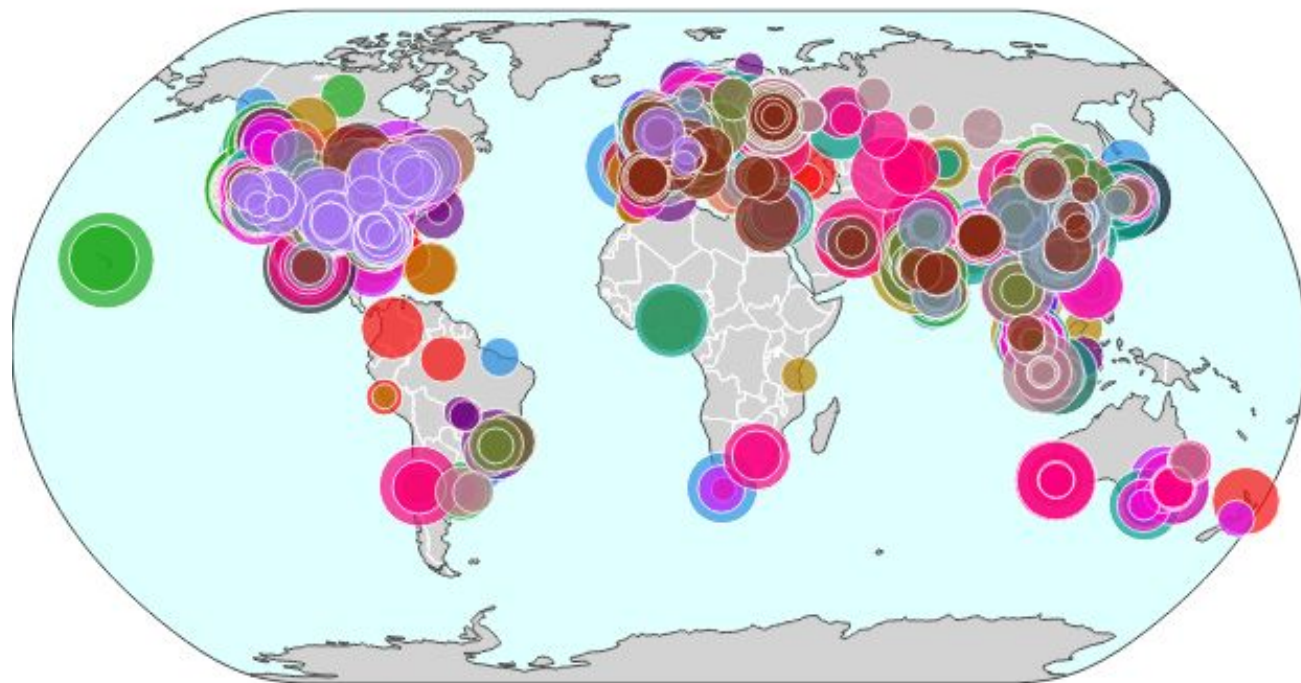
'min_samples_split': 10,

'n_estimators': 200



Сравнение истинных и предсказанных значений





- Fashion & Retail
- Automotive
- Technology
- Finance & Investments
- Media & Entertainment
- Telecom
- Diversified
- Food & Beverage
- Logistics
- Gambling & Casinos
- Manufacturing
- Real Estate
- Metals & Mining
- Energy
- Healthcare
- Service
- Construction & Engineering
- Sports

HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise)

HDBSCAN группирует точки данных на основе их плотности:

- **Кластер** — область с высокой плотностью точек.
- **Шум (Noise)** — точки в разреженных областях, не принадлежащие ни одному кластеру.

В отличие от DBSCAN, который использует фиксированные параметры `eps` (радиус окрестности) и `min_samples` (минимальное число соседей), HDBSCAN:

- Строит **иерархию кластеров** для разных уровней плотности.
- Автоматически выбирает оптимальные кластеры из этой иерархии.

1. Построение графа взаимной достижимости

- Для каждой точки вычисляется **расстояние до k -го ближайшего соседа** (core distance).
- Строится взвешенный граф, где ребра — это **взаимная достижимость** (mutual reachability distance):

$$d_{\text{mreach}}(a, b) = \max(\text{core}_k(a), \text{core}_k(b), d(a, b))$$

где $d(a, b)$ — обычное расстояние между точками.

2. Построение минимального остовного дерева (MST)

- Граф преобразуется в **минимальное остовное дерево**, чтобы выявить иерархическую структуру.

3. Построение иерархии кластеров

- Дерево "разрезается" на разных уровнях плотности, формируя **дендрограмму**.

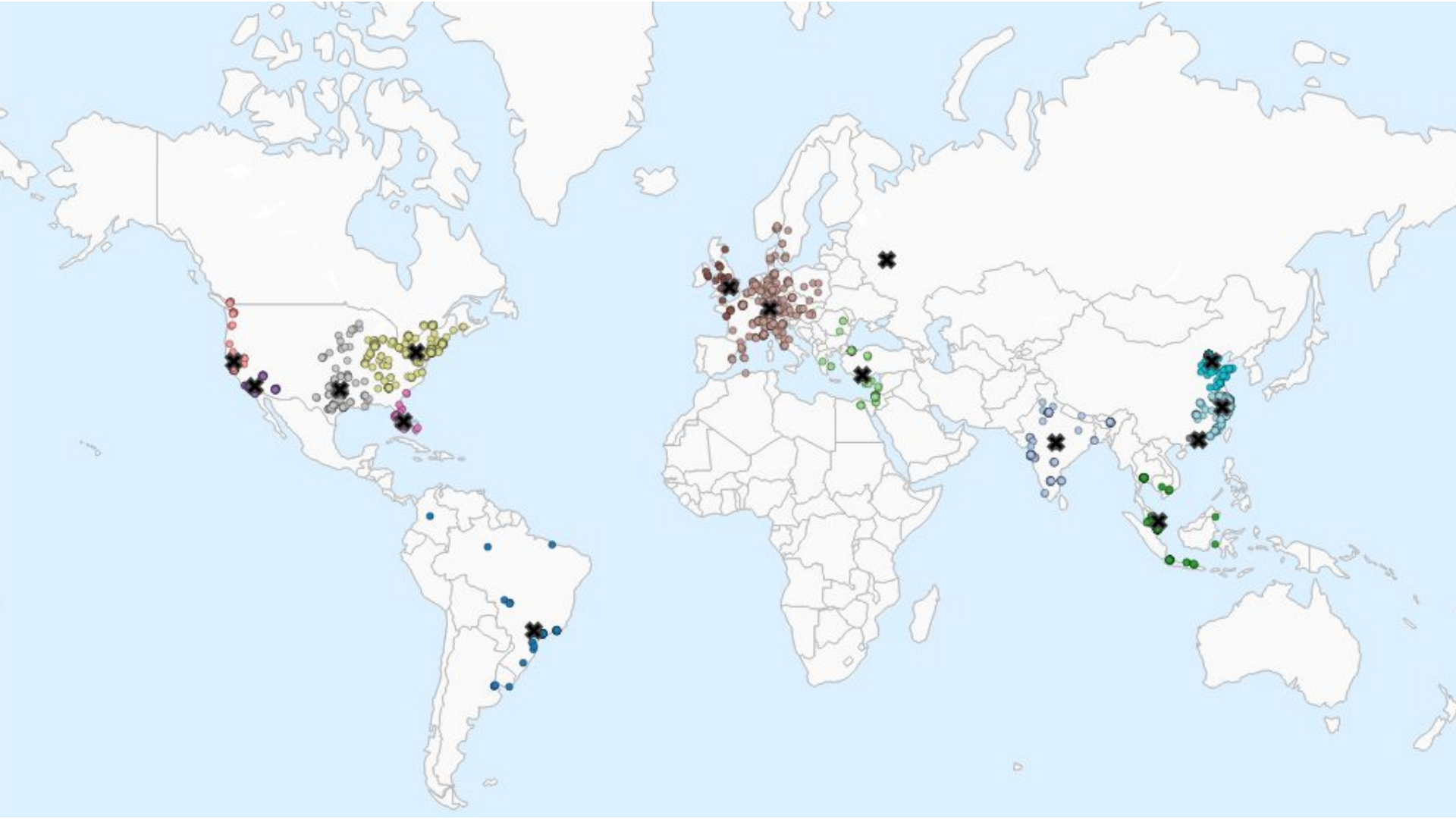
4. Выбор устойчивых кластеров

- Используется метрика **стабильности** (stability), чтобы выбрать наиболее устойчивые кластеры из иерархии.
- Кластеры с максимальной стабильностью сохраняются, остальные отбрасываются.

5. Пометка шума

- Точки, не вошедшие в устойчивые кластеры, помечаются как **шум (-1)**.







ВСЁ!