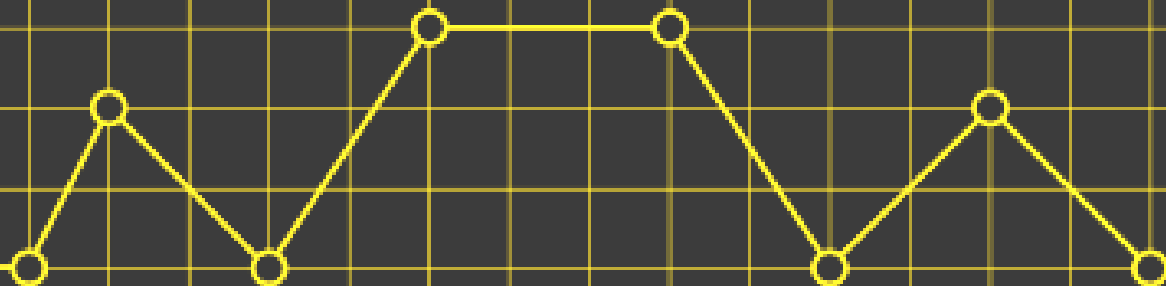


**DANO**

# ИССЛЕДОВАНИЕ

©18 команда

**DANO HAKATON & HSE NN**



ВЫСШАЯ ШКОЛА  
ЭКОНОМИКИ



**ТА** БАНК

**01**

**Предварительный  
анализ данных**

# Data Scientists

**Боронина  
Юля**

**Назаренко  
Даша**

**№18**

**MaLoStь**

**Куликова  
Аня**

**Фахруллина Ками**

Math+Logic+Statistics

Данные

Гипотеза

Обработка

Результаты

Заключение

# Структура базы данных

! часть данных (~0,5% от полученных в каждый день) по полученным клиентами Т-Банка штрафам ГИБДД\* - 97 308 строк, 22 характеристики

Период:

с 28 апреля 2024 г  
по 28 мая 2024 г.

---

\*Одна запись – один уникальный штраф ГИБДД. К характеристикам штрафа добавлена информация о характеристиках водителя и автомобиля на котором совершено правонарушение.

Данные

Гипотеза

Обработка

Результаты

Заключение

# Структура базы данных



Данные

Гипотеза

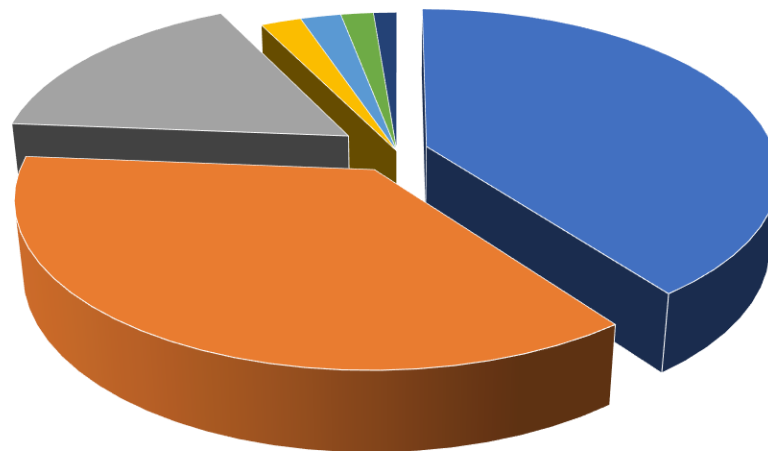
Обработка

Результаты

Заключение

# Структура базы данных

Тип кузова (количество штрафов)



■ Седан

■ Универсал

■ Хэтчбек

■ Купе

■ Фургон

■ Внедорожник

■ Лифтбек

+34 типа малозначимы

Данные

Гипотеза

Обработка

Результаты

Заключение

# Структура базы данных



Данные

Гипотеза

Обработка

Результаты

Заключение

# Структура базы данных

Типы машин

- LADA (BA3)
- KIA
- BMW
- HYUNDAI
- VOLKSWAGEN
- MERCEDES-BENZ
- TOYOTA
- FORD
- SKODA
- NISSAN
- AUDI
- MAZDA
- MITSUBISHI
- RENAULT
- CHEVROLET
- OPEL
- ГАЗ
- LEXUS
- LAND ROVER
- VOLVO





Данные

Гипотеза

Обработка

Результаты

Заклучение

# Структура базы данных

Типы марки+модели

- HYUNDAI SOLARIS
- KIA RIO
- FORD FOCUS
- TOYOTA CAMRY
- VOLKSWAGEN POLO
- LADA (BA3) GRANTA
- LADA (BA3) PRIORA
- SKODA OCTAVIA
- LADA (BA3) 2108
- MERCEDES-BENZ E-КЛАСС
- OPEL ASTRA
- LADA (BA3) VESTA
- KIA CEED
- BMW X5

HYUNDAI SOLARIS	KIA RIO	FORD FOCUS	VOLKSWAGEN POLO	LADA (BA3) PRIORA	LADA (BA3) 2108	OPEL ASTRA	LADA (BA3) VESTA
		TOYOTA CAMRY	LADA (BA3) GRANTA	SKODA OCTAVIA	MERCEDES-BENZ E-КЛАСС	KIA CEED	BMW X5

Данные

Гипотеза

Обработка

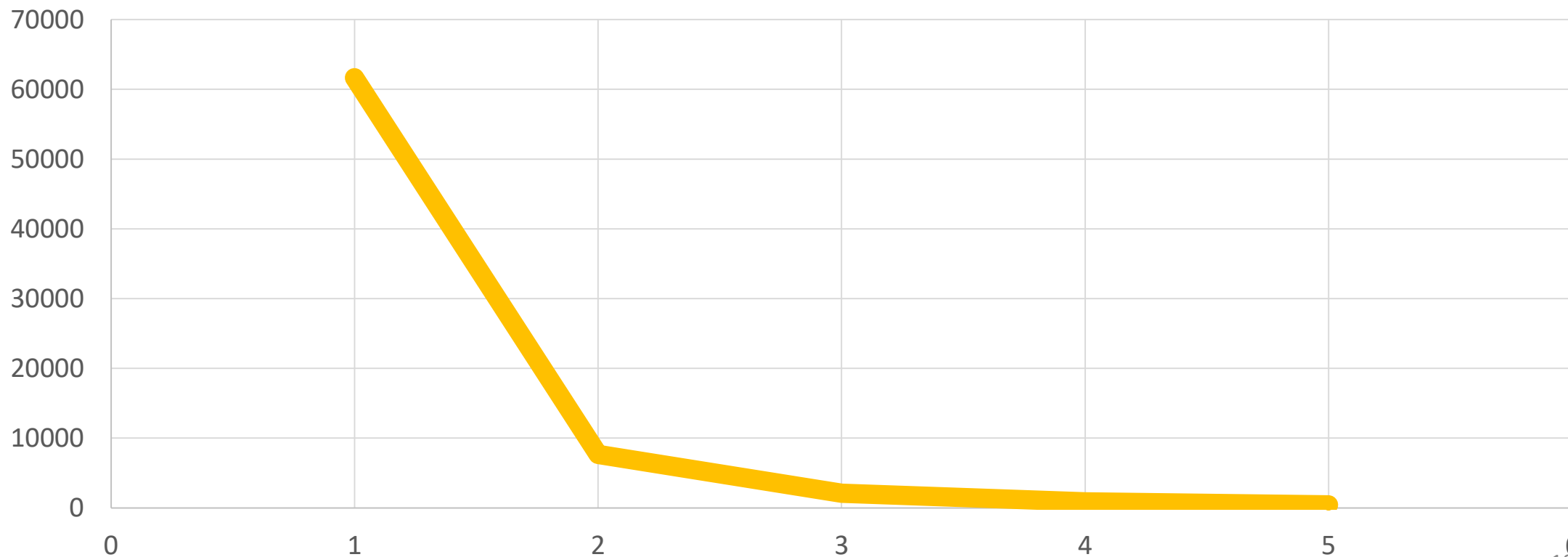
Результаты

Заключение

# Структура базы данных

Статистика количественного показателя  
нарушений

Кол-во человек с таким  
количеством нарушений



**Уникальные водители:** 1297 (с превышением скорости более 5 раз)

Нарушений у водителя<sup>10</sup>

Данные

Гипотеза

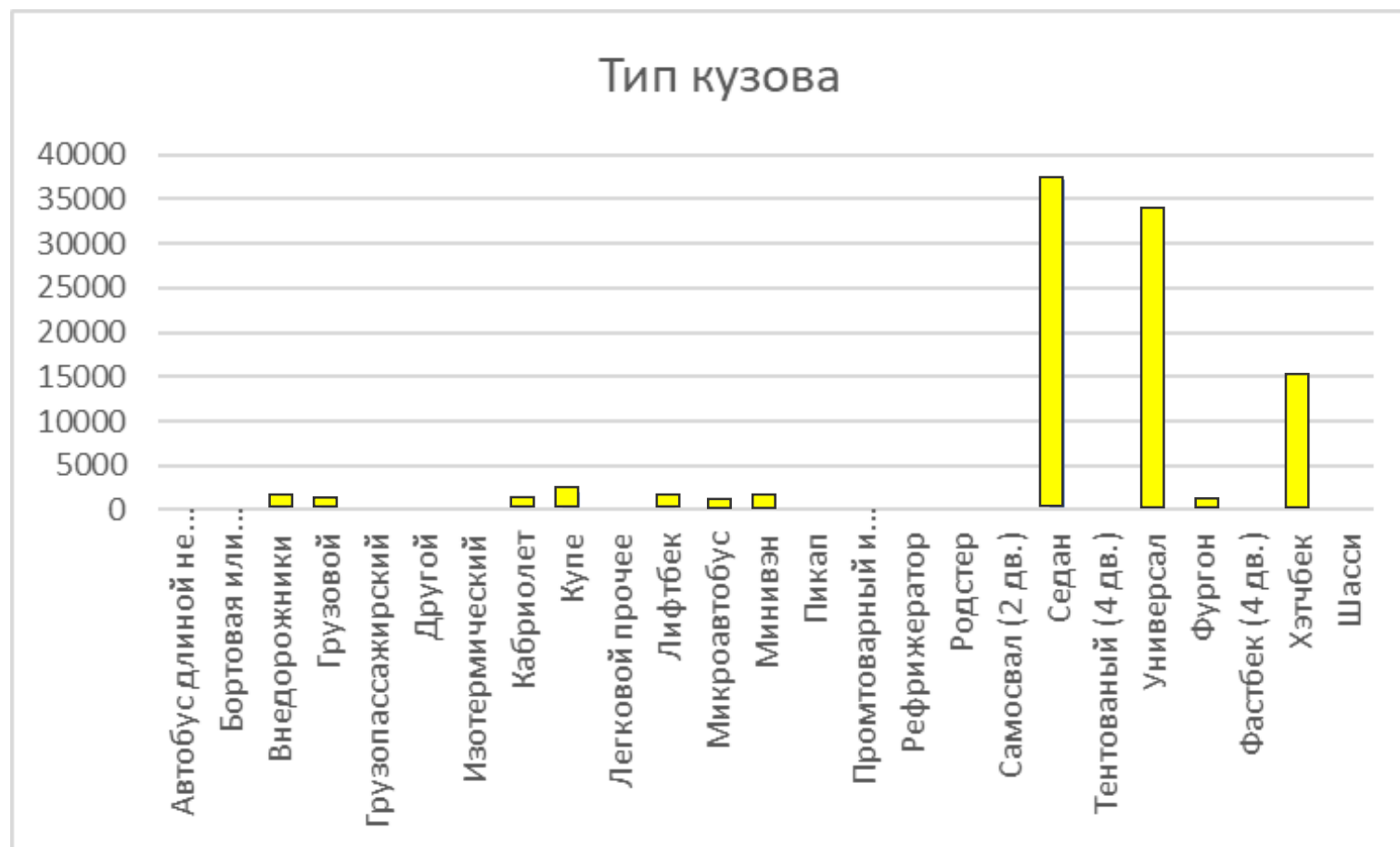
Обработка

Результаты

Заключение

# Структура базы данных

Количество  
штрафов с  
повторениями



Данные

Гипотеза

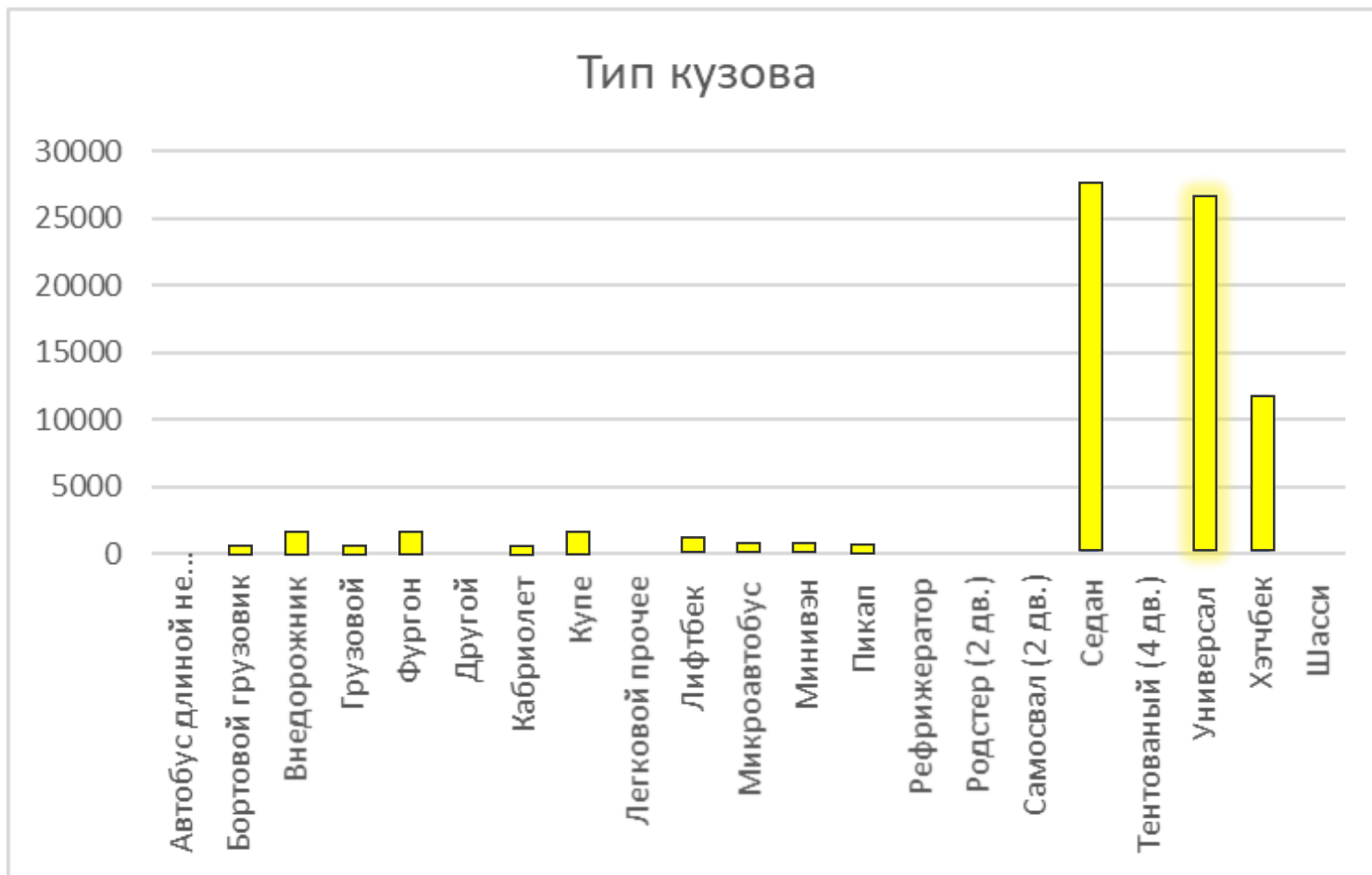
Обработка

Результаты

Заключение

# Структура базы данных

Количество штрафов без повторений (=количество кузовов)



БЫЛО: 9728

СТАЛО: 73555

Данные

Гипотеза

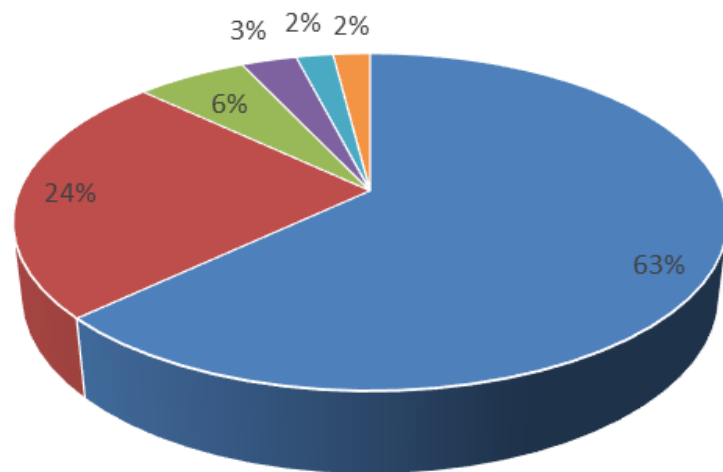
Обработка

Результаты

Заключение

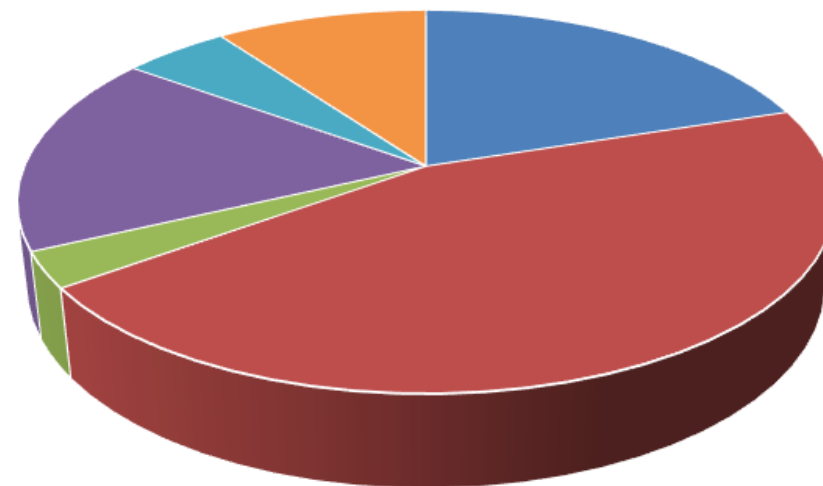
# Структура базы данных

Спрос новых машин на рынке



■ Внедорожники ■ Седаны ■ Лифтбэк ■ Пикапы ■ Универсалы ■ Остальные

Спрос поддержанных машин на рынке



■ Внедорожники ■ Седан ■ Лифтбэк ■ Пикап ■ Универсалы ■ Остальные

Данные

Гипотеза

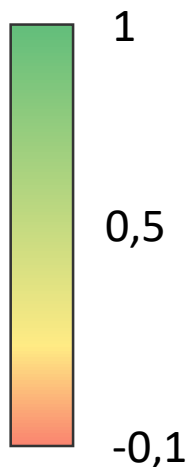
Обработка

Результаты

Заключение

## Предварительный анализ с помощью корреляционной таблицы

	<i>engine_power</i>	<i>car_price</i>	<i>age</i>	<i>children_cnt</i>	<i>income</i>
<i>engine_power</i>	1	0,546196121	0,04852	0,004078043	0,166865959
<i>car_price</i>	0,546196121	1	0,10217	0,005724679	0,181820769
<i>age</i>	0,048522478	0,102166076	1	0,118170086	-0,015317977
<i>children_cnt</i>	0,004078043	0,005724679	0,11817	1	0,024494617
<i>income</i>	0,166865959	0,181820769	-0,01532	0,024494617	1



Зависимости количественных величин:

- engine\_power** – мощность двигателя автомобиля
- car\_price** – стоимость автомобиля (условные единицы)
- age** – возраст владельца
- children\_cnt** – количество детей у владельца
- income** - размер месячного дохода водителя

02

**Постановка**

**гипотезы №1**





	<b>Исследовательский вопрос</b>									<b>Гипотеза</b>																	
	<p>Взаимосвязаны ли характеристики клиентов и их автомобилей с характеристиками правонарушений (время/дата совершения правонарушения, статья правонарушения) и если да, то как?</p>									<p>Соотношение частотности различных категорий штрафов у взрослых с детьми и у взрослых без детей <b>отличаются</b></p>																	

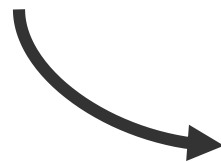




Рождение ребенка значимым шагом для семьи, что заставляет чувствовать больше ответственности



Ответственные водители реже делают то, что может угрожать жизни: нарушают скорость, проезжают на красный и т.д.



Для таких водителей-родителей характерны нарушения, связанные с остановкой, заездом за линию и т.д.

М  
Е  
Х  
А  
Н  
И  
З  
М

**03**

**Проверка**

**ГИПОТЕЗЫ**

Данные

Гипотеза

Обработка

Результаты

Заключение

## Сравнение количества штрафов в зависимости от наличия детей



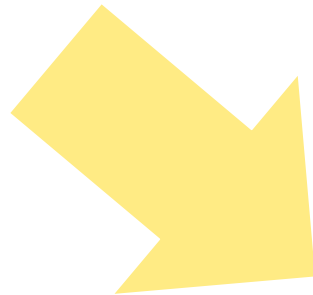


## Сравнение количества штрафов в зависимости от наличия детей



# t-тест

P-value = 0,09



Гипотеза **не**  
подтверждается

**Гипотеза  
опровергнута**

03

# Результаты анализа





Водители совершают схожие нарушения в одинаковом количестве.

Наличие детей (как психологический фактор) не влияет.



Данные

Гипотеза

Обработка

Результаты

Заключение

# Практическая значимость

**Службам  
ГИБДД**

- Проведение психологических тренингов для родителей

# Ограничения и перспективы

## Ограничения

- Номера владельцев, а не непосредственно владельцы

## Перспективы

- Получить данные о конкретных владельцах через систему
- Протестировать социальные критерии на большей выборке
- Провести эмпирическое исследование о вождении автомобилей с типом кузова «Универсал»

# Каких данных не хватило?

- Стаж вождения
- Возраст ребенка (детей)
- Дата покупки машины

Количественных  
данных



**Спасибо!**  
**У нас всё**

© малость педагоги

# Тест Фишера

Статистическая значимость – это количественный показатель, указывающий на то, что полученные результаты не являются случайными и могут быть признаны достоверными.

Она часто используется в маркетинге, когда нужно **проверить чистоту эксперимента и понять, можем ли мы доверять результатам теста.**

02

**Постановка**

**гипотезы №2**



# Рекомендации

## Производителям

- Внедрять передовые технологий безопасности (уведомления о превышении скорости, встроенный навигатор)
- Обучение водителей (привыкание к новой коробке передач)
- Работа над репутацией
- Сотрудничество с органами ГБДД
- Улучшение качества сборки



	<b>Исследовательский вопрос</b>									<b>Гипотеза</b>										
	<p>Взаимосвязаны ли характеристики клиентов и их автомобилей с характеристиками правонарушений (время/дата совершения правонарушения, статья правонарушения) и если да, то как?</p>									<p>Время выписки штрафа зависит от наличия детей в семье: у человека, который имеет детей, большее количество штрафов будет приходиться на отличающееся время</p>										



Данные

Гипотеза

Обработка

Результаты


Заключение

Дети посещают различные развивающие секции и кружки, ходят в школу и детский сад

Многие места находятся вдали от дома, из-за чего родителям нужно возить ребенка, чтобы убедиться в его безопасности на пути

Поэтому возрастает необходимость возить ребенка на автомобиле в утренние и вечерние часы

М  
Е  
Х  
А  
Н  
И  
З  
М



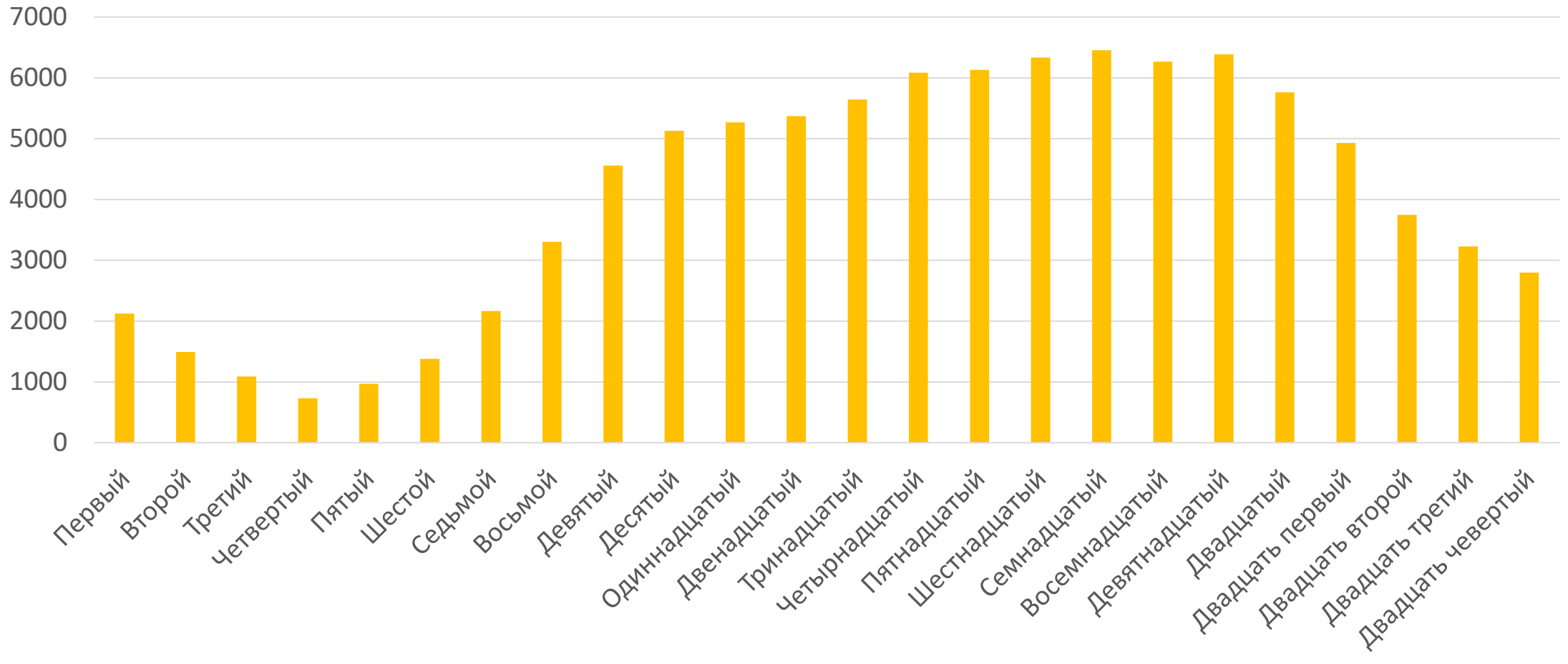
**03**

**Проверка**

**гипотезы №2**



## Общее



час

Данные

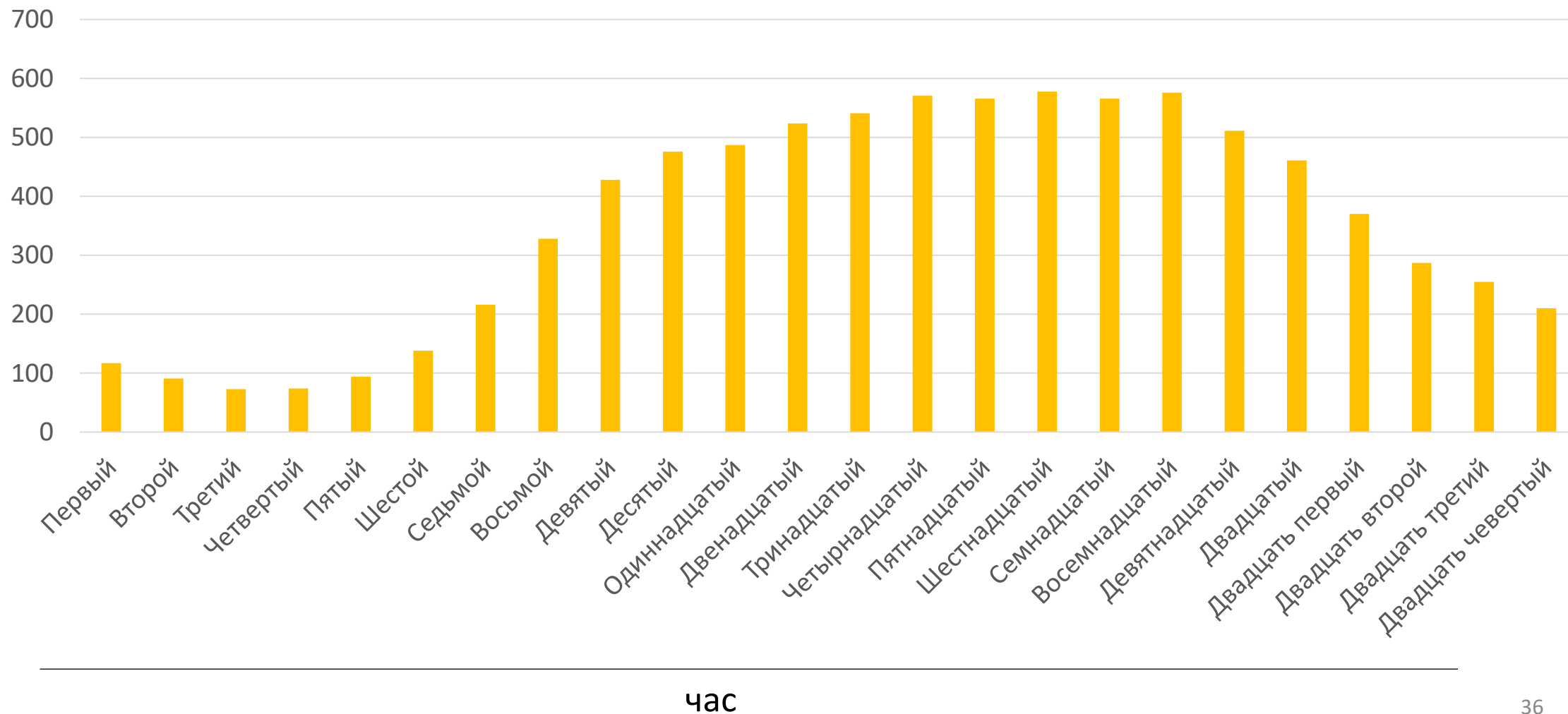
Гипотеза

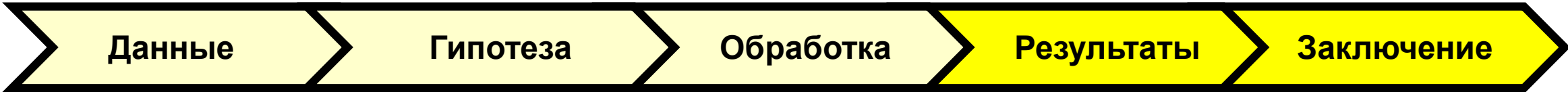
Обработка

Результаты

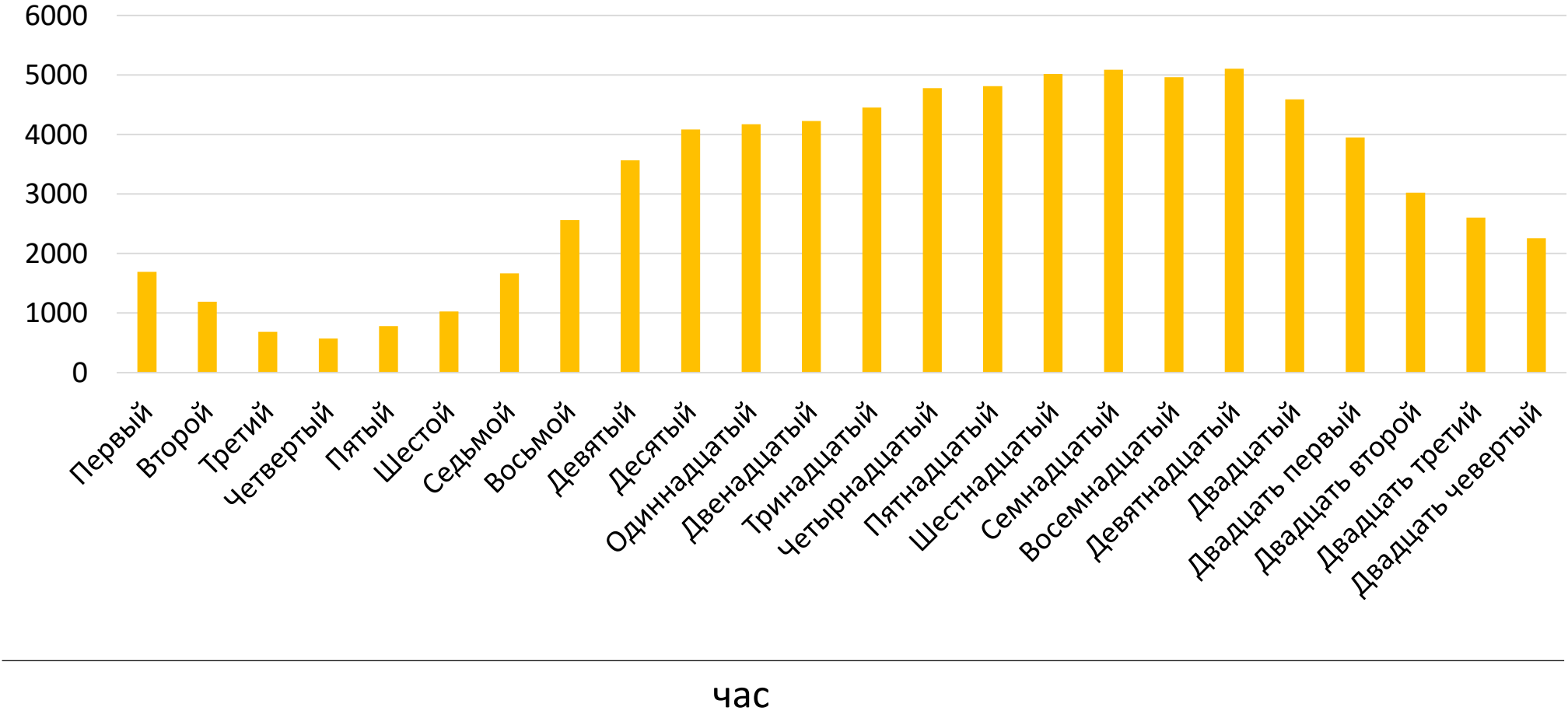
Заключение

## Количество штрафов у водителей с детьми





# Количество штрафов у водителей без детей



# Коэффициент корреляции

$$r_{xy} = \frac{\overline{x \cdot y} - \bar{x} \cdot \bar{y}}{S(x) \cdot S(y)} = \frac{1425847.708 - 355.75 \cdot 3202.5}{185.555 \cdot 1578.947} = 0.978$$

связь между весьма высокая и прямая

**Маленький размер выборки:**

## Точный тест Фишера

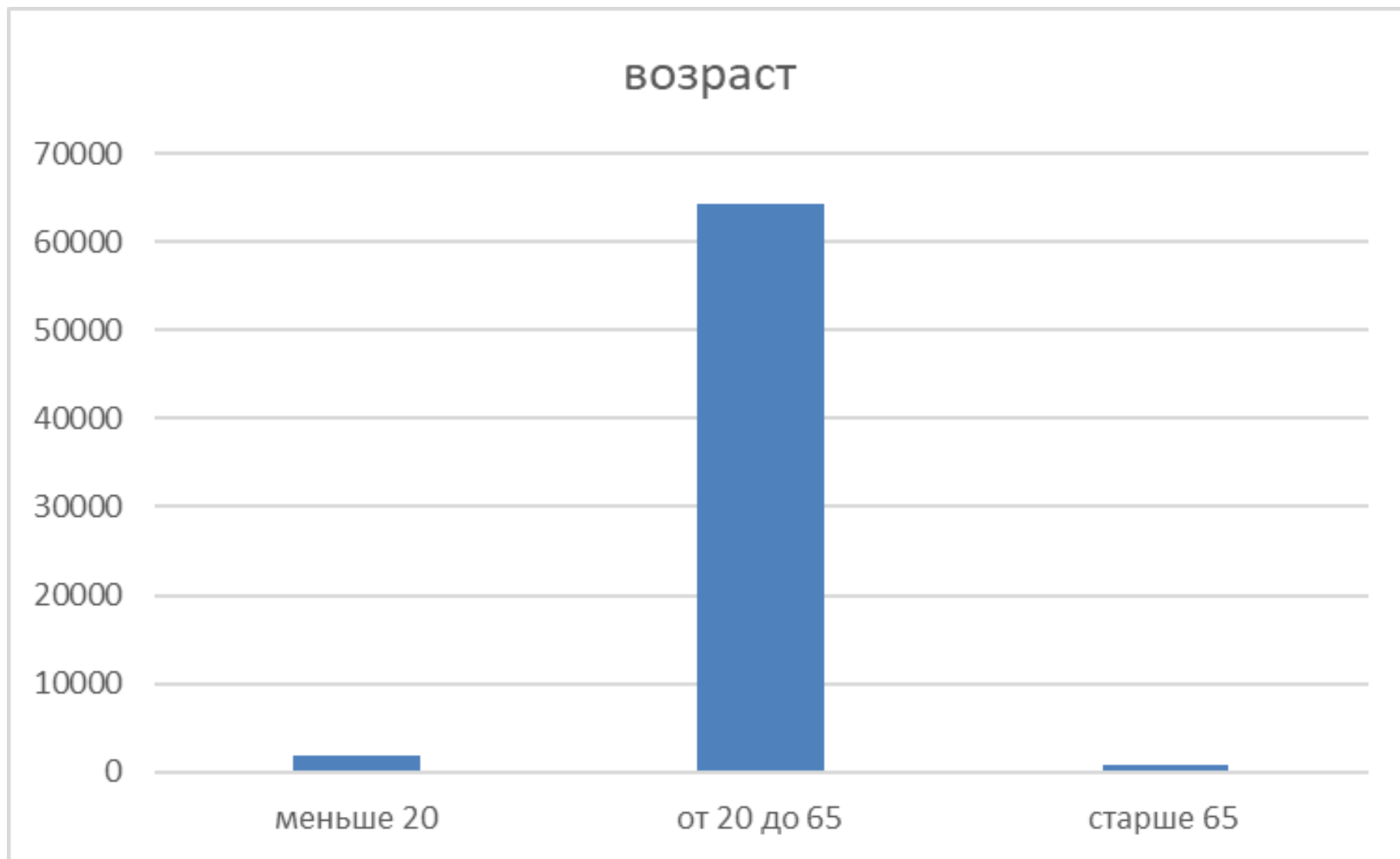
$$F = \frac{\sum (y_x - \bar{y})^2}{\sum (y_i - y_x)^2} \frac{n - m - 1}{m} = \frac{57238946.5551}{2594797.44} \cdot \frac{24 - 1 - 1}{1} = 485.301 \quad F_{табл} = 4.3009$$

$$F > F_{табл}$$

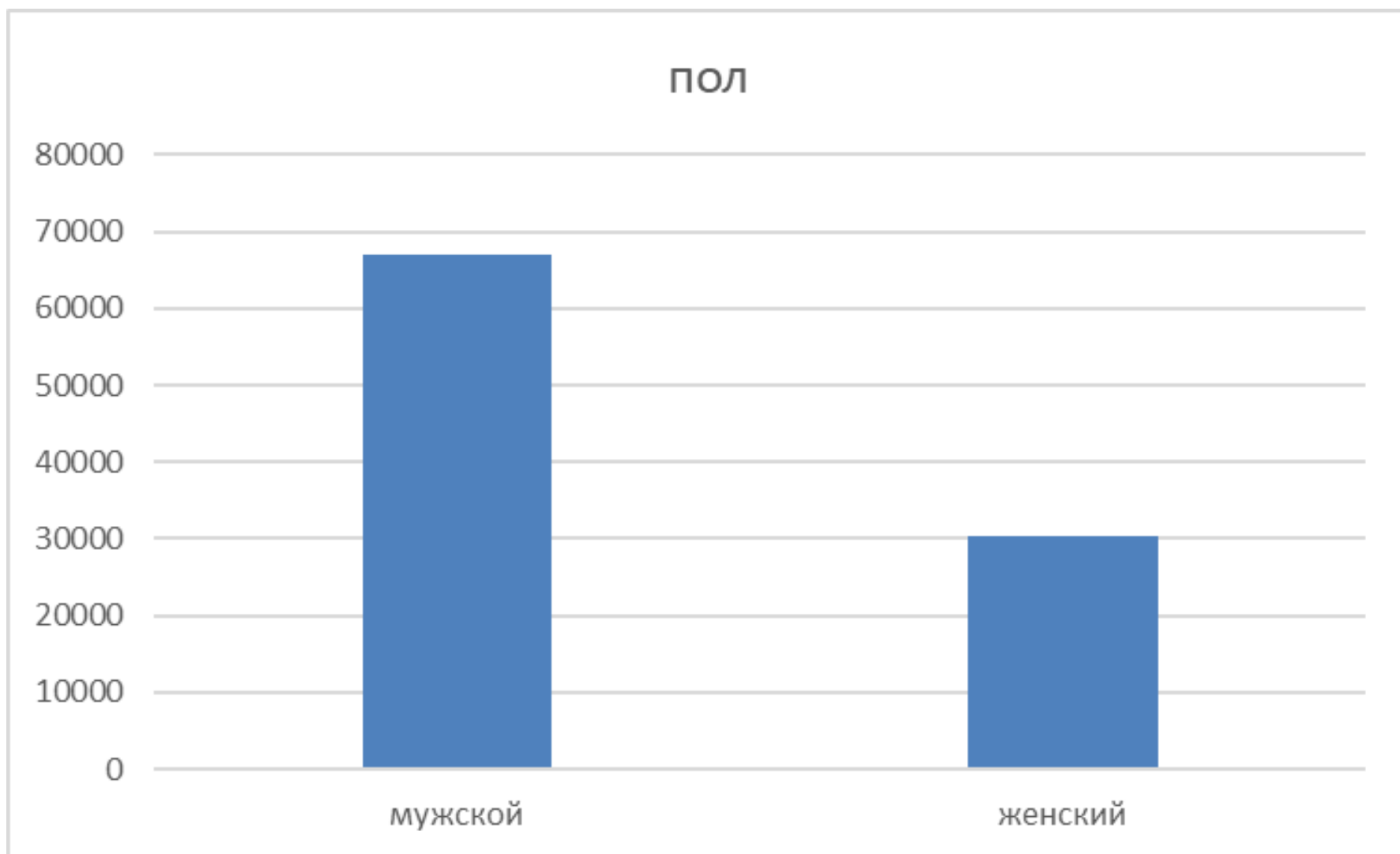


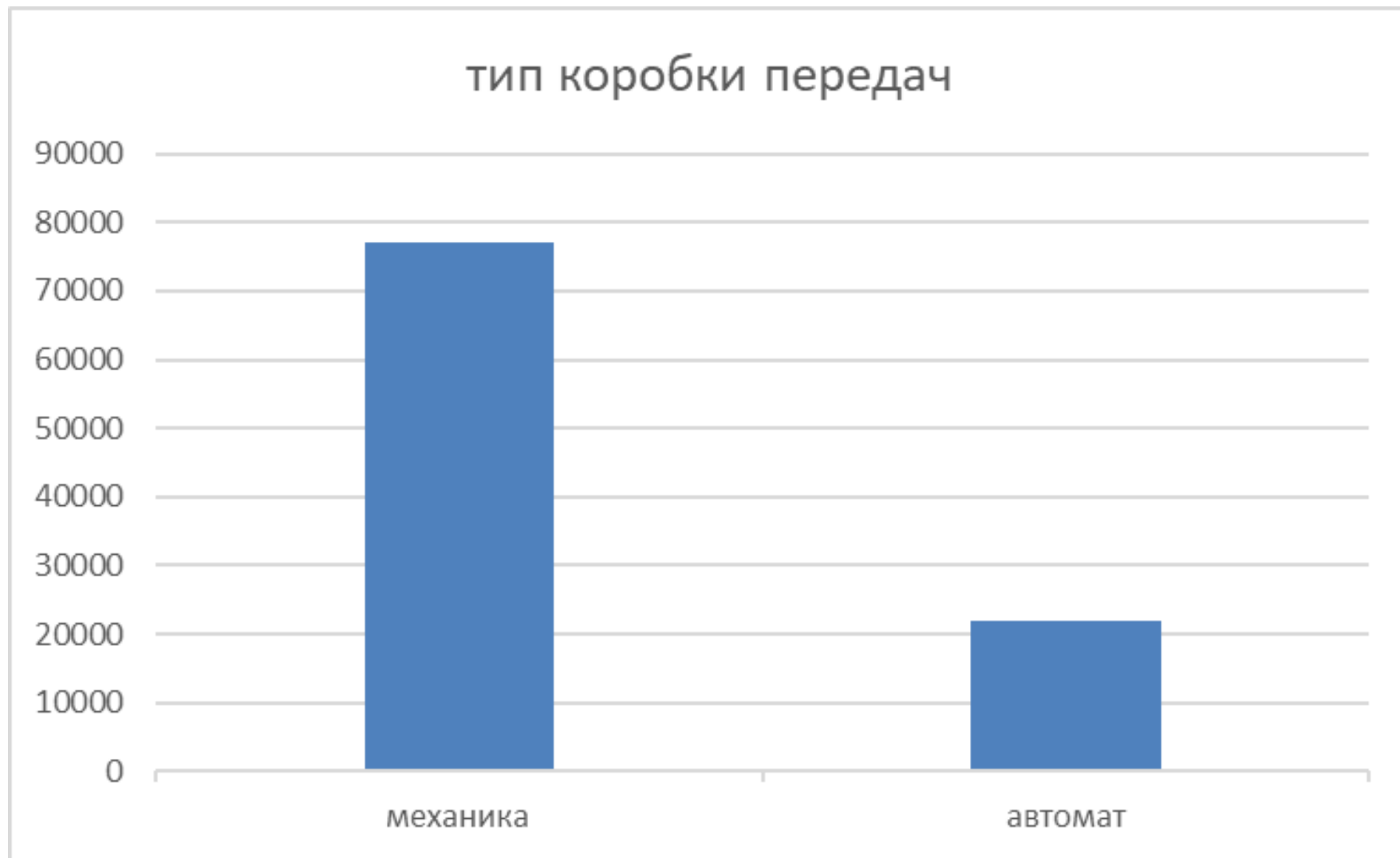
найденная оценка уравнения регрессии  
статистически надежна

**Гипотеза №2  
опровергнута**









# Регрессионный анализ

$$y_x = r_{xy} \cdot \frac{x - \bar{x}}{S(x)} \cdot S(y) + \bar{y} = 0.993 \frac{x - 3.571}{13.186} 12.496 + 3.571 = 0.941x + 0.212$$

Маленький размер выборки:

Связь прямая

## Точный тест Фишера

$$F > F_{\text{табл}}$$

коэффициент детерминации  
статистически значим (найденная оценка  
статистически надежна)

$$F = \frac{0.9852}{1 - 0.9852} \frac{28 - 1 - 1}{1} = 1729.152$$

$$F_{\text{табл}} = 4.2252$$

Данные

Гипотеза

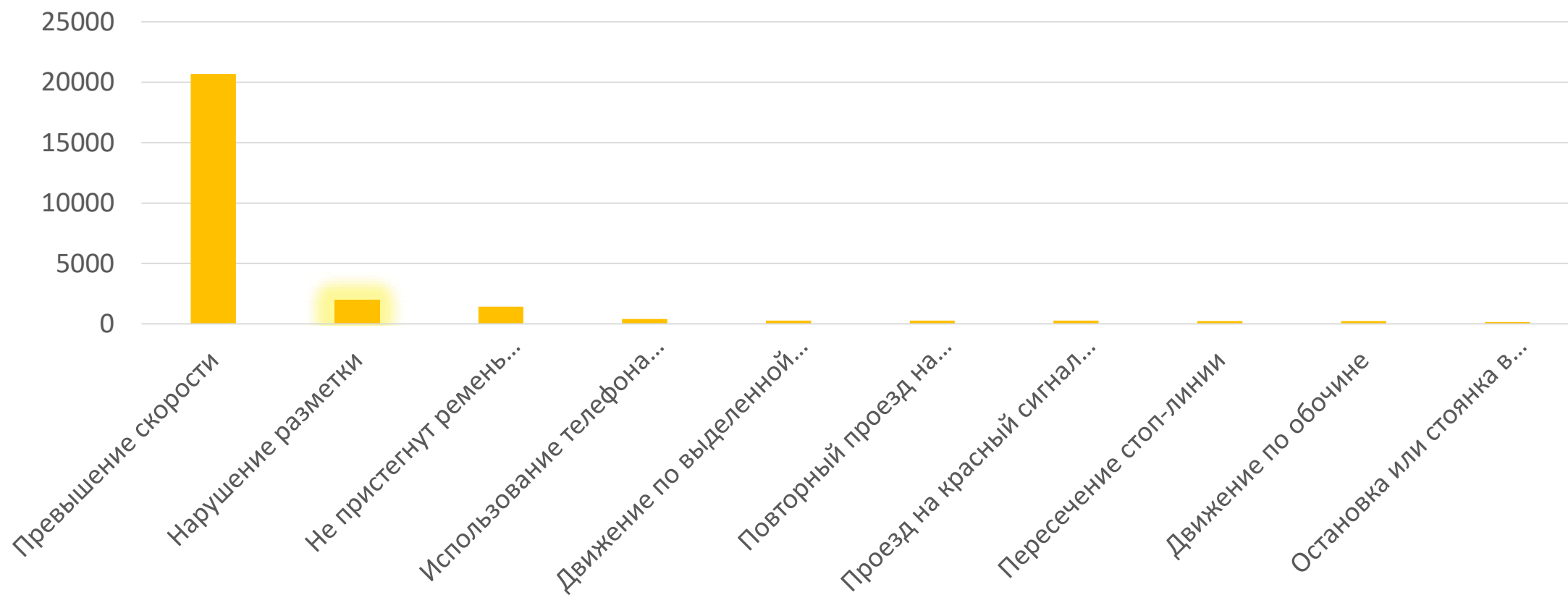
Обработка

Результаты

Заключение

# Структура базы данных

Виды нарушений у универсалов



Данные

Гипотеза

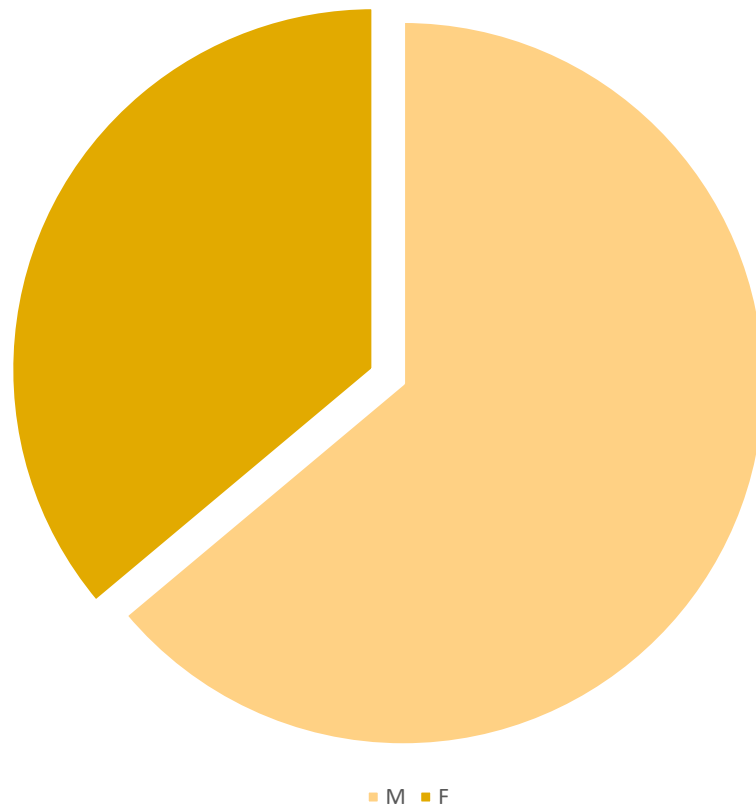
Обработка

Результаты

Заключение

# Структура базы данных

Соотношение полов, пользующихся универсалов



02

**Постановка  
гипотезы №2**





	<b>Исследовательский вопрос</b>									<b>Гипотеза</b>																	
	<p>Взаимосвязаны ли характеристики клиентов и их автомобилей с характеристиками правонарушений (время/дата совершения правонарушения, статья правонарушения) и если да, то как?</p>									<p>Из самых распространенных кузовов автомобилей кузовы «Универсал» имеют больше нарушений</p>																	

Данные

Гипотеза

Обработка

Результаты

Заключение

Большие кузова имеют большую популярность за счет удобства перевозки вещей

При покупке автомобиля с кузовом многие думают о комфорте перевозки, но не о неудобствах управления

Из-за частотными случаями аварий являются нарушения разметки, так как камера считывает малейшие заступы

М  
Е  
Х  
А  
Н  
И  
З  
М



A stylized world map in shades of gray is visible in the background. The map is centered on the Atlantic Ocean, with the Americas on the left and Europe and Africa on the right. The map is composed of solid gray shapes representing continents and oceans.

**02**

**Проверка**

**гипотезы №3**

**Гипотеза №3  
подтверждена**