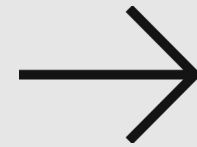


ЦЕНТРАЛЬНЫЙ  
УНИВЕРСИТЕТ

# Прогноз спроса на аренду велосипедов (Bike Sharing)



# Постановка задачи и бизнес-контекст



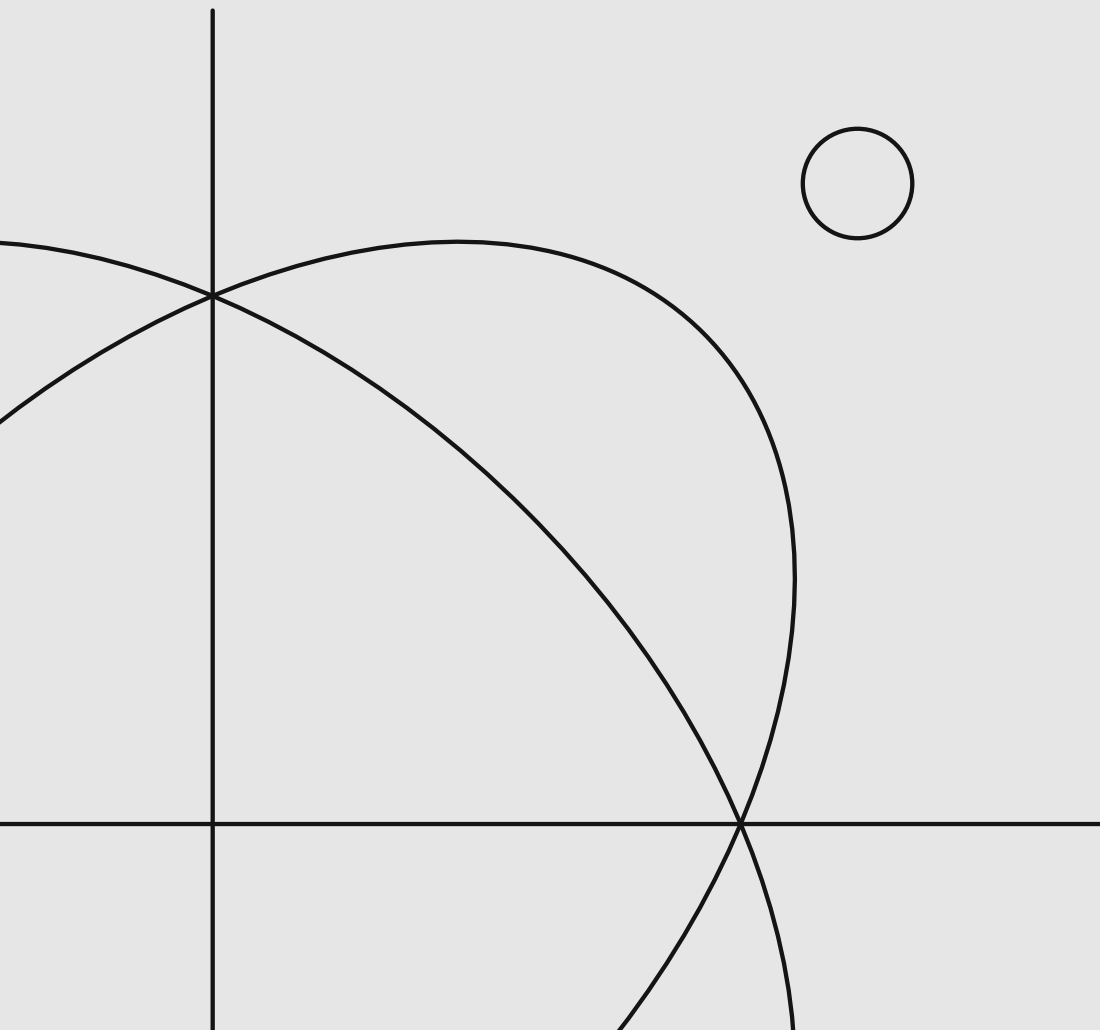
Мы хотим построить модель, которая предсказывает количество аренд велосипедов в конкретный час на основе календарных признаков, погодных условий и характеристик дня.



**Тип задачи:** регрессия

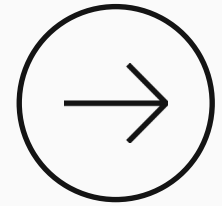
**Зачем это нужно в бизнес-контексте:**

1. обеспечить достаточное число велосипедов в популярных локациях
2. уменьшить ситуации, когда велосипедов не хватает
3. оптимизировать работу машин, которые перевозят велосипеды
4. снизить затраты
5. улучшить сервис и пользовательский опыт

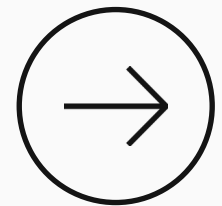




# Описание датасета



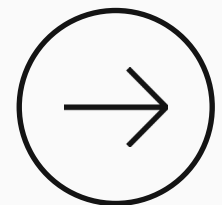
Мы используем Bike Sharing Dataset - это реальные данные из системы Capital Bikeshare (Вашингтон, США) за 2011–2012 годы.



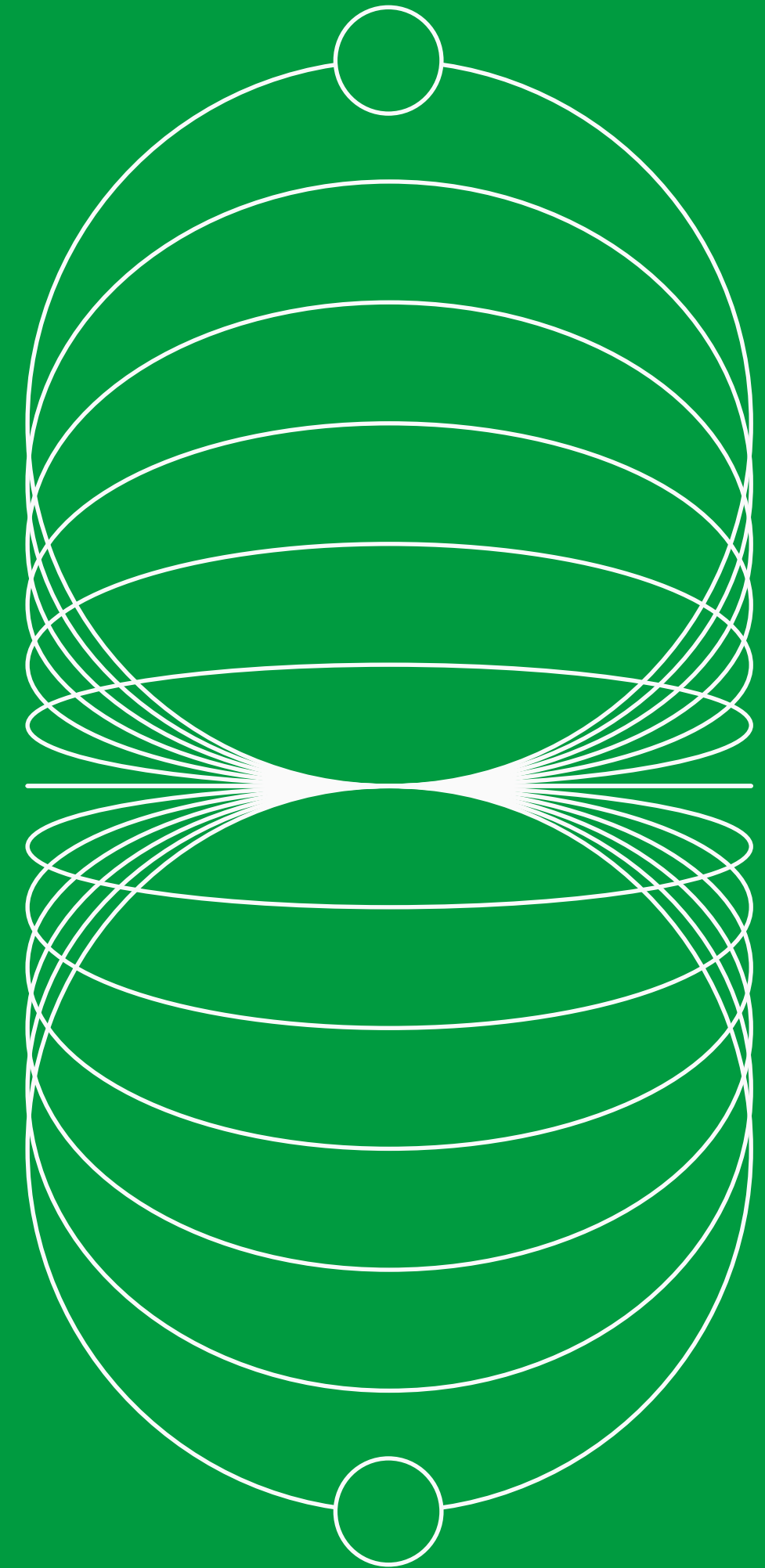
## **Файлы:**

*hour.csv* — данные по часам (17 379 строк)

*day.csv* — данные по дням (731 строк)



Всего 17 признаков, включая таргет



# Основные признаки

## 01 Идентификатор и дата

instant (integer) - порядковый номер записи  
dteday (date) - календарная дата

## 02 Календарные признаки

season (categorical) - сезон года  
yr (categorical) - год наблюдения  
mnth (categorical) - месяц  
hr (categorical) - час суток  
weekday (categorical) - день недели

## 03

holiday (binary) - официальный праздник или нет  
workingday (binary) - рабочий день или нет

## 04 Погодные признаки

temp (continuous) - нормализованная температура  
atemp (continuous) - нормализованная температура «по ощущению»  
hum (continuous) - нормализованная влажность  
windspeed (continuous) - нормализованная скорость ветра

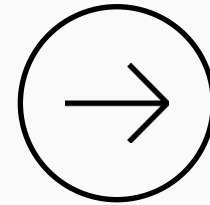
## 05

### Пользовательские признаки

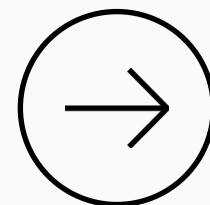
## 06

casual (integer) - количество незарегистрированных пользователей  
registered (integer) - количество зарегистрированных пользователей

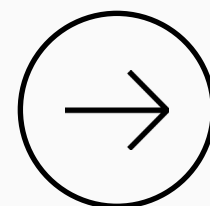
# EDA



Проверили данные на мусор и наличие пропусков и выяснили, что пропуски мусор и пропуски отсутствуют



Распределение целевой переменной `cnt` сильно смещено влево. такая форма распределения это база для данных с неравномерной нагрузкой в течение суток и указывает на присутствие выраженной сезонности и редких пиковых значений.



Признаки `casual` и `registered` практически полностью повторяют целевую переменную `cnt`, поэтому их нельзя использовать в модели, а остальные признаки имеют слабую корреляцию с количеством аренд.

# Построение Baseline

В качестве Baseline мы выбрали: LinearRegression, RandomForest и DummyRegressor и получили такие результаты на следующих метриках:

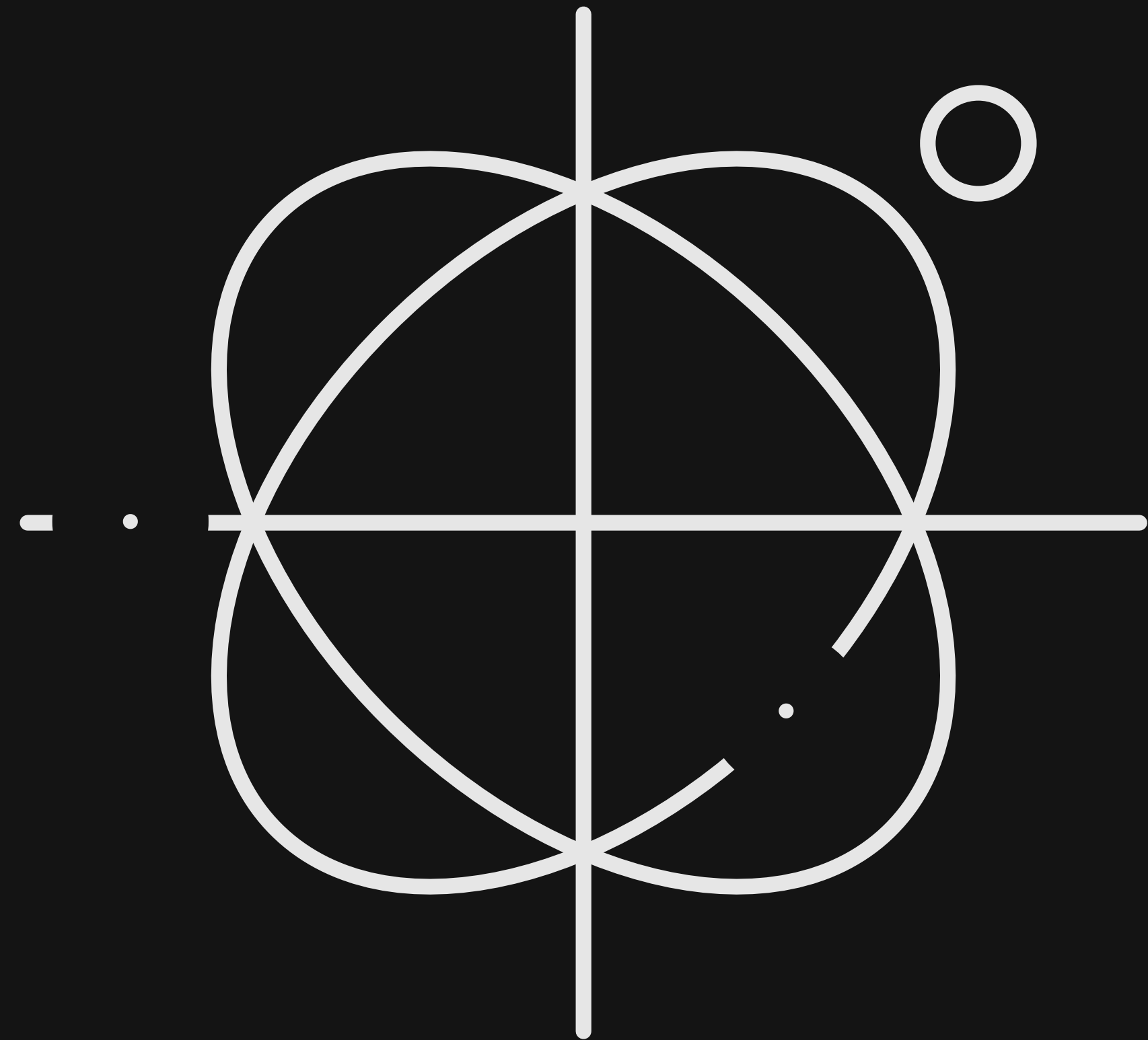
	Model	RMSE	MAE	R2
0	Dummy	178.034917	140.079835	-0.000980
1	Linear Regression	139.201123	104.795881	0.388072
2	Random Forest	40.574918	24.090368	0.948009

# Работа с аномалиями и выбросами



Для поиска сложных выбросов были использованы алгоритмы Isolation Forest и LOF и оба метода выделили около 3% наблюдений как аномальные, что соответствует редким условиям в данных.

Мы выяснили, что найденные аномалии не являются ошибками данных, поэтому удаление их приведёт к потере информации о пиковых и экстремальных условиях, что ухудшит качество модели в ключевых ситуациях



# Генерация признаков



Добавлены флаги выходных и часов пик, отражающие типичное поведение пользователей и различия в спросе в разное время

## погодные условия:

- **is\_good\_weather, is\_bad\_weather** - хорошая погода увеличивает желание кататься, а плохая снижает.
- **is\_cold, is\_comfortable\_temp, is\_hot** - комфортная температура оптимальна для велопогулок, слишком высокие или низкие температуры снижают спрос.
- **is\_high\_humidity, is\_high\_wind** - высокая влажность и сильный ветер делают поездку менее комфортной.
- **is\_ideal\_conditions** - сочетание всех благоприятных факторов должно максимально увеличивать спрос.
- **comfort\_index** - комплексный показатель комфорта, объединяющий все погодные факторы.

## взаимодействия:

- **good\_weather\_weekend** - хорошая погода в выходной создаёт идеальные условия для прогулочных поездок.
- **ideal\_conditions\_weekend** - идеальные условия в выходной должны давать максимальный спрос.
- **good\_weather\_work\_peak** - хорошая погода в рабочий день в час пик влияет на использование велосипедов для поездок на работу.
- **summer\_weekend, warm\_month\_weekend** - тёплое время года + выходной = пик использования как досуг.
- **bad\_weather\_workday** - из-за плохой погоды в будний день люди могут быть вынуждены использовать велосипеды.

## время до пикового часа:

- **hours\_to\_morning\_peak, hours\_to\_evening\_peak, min\_hours\_to\_peak** - спрос может меняться в зависимости от близости к часам пик, когда люди едут на работу или с работы. чем ближе к пику, тем выше спрос.

## время до праздника:

- **days\_to\_holiday, days\_after\_holiday** - в дни перед праздником люди могут больше использовать велосипеды для подготовки или отдыха, а после праздника спрос может падать из-за усталости или возвращения к обычному ритму.
- **near\_holiday** - дни около праздника могут иметь особый паттерн спроса, отличающийся от обычных дней.

## сезонные признаки:

- **is\_summer, is\_winter** - летом спрос максимальный из-за теплой погоды, зимой минимальный.
- **is\_warm\_month, is\_cold\_month** - в тёплые месяцы (май-сентябрь) высокий спрос, в холодные (ноябрь-февраль) низкий.

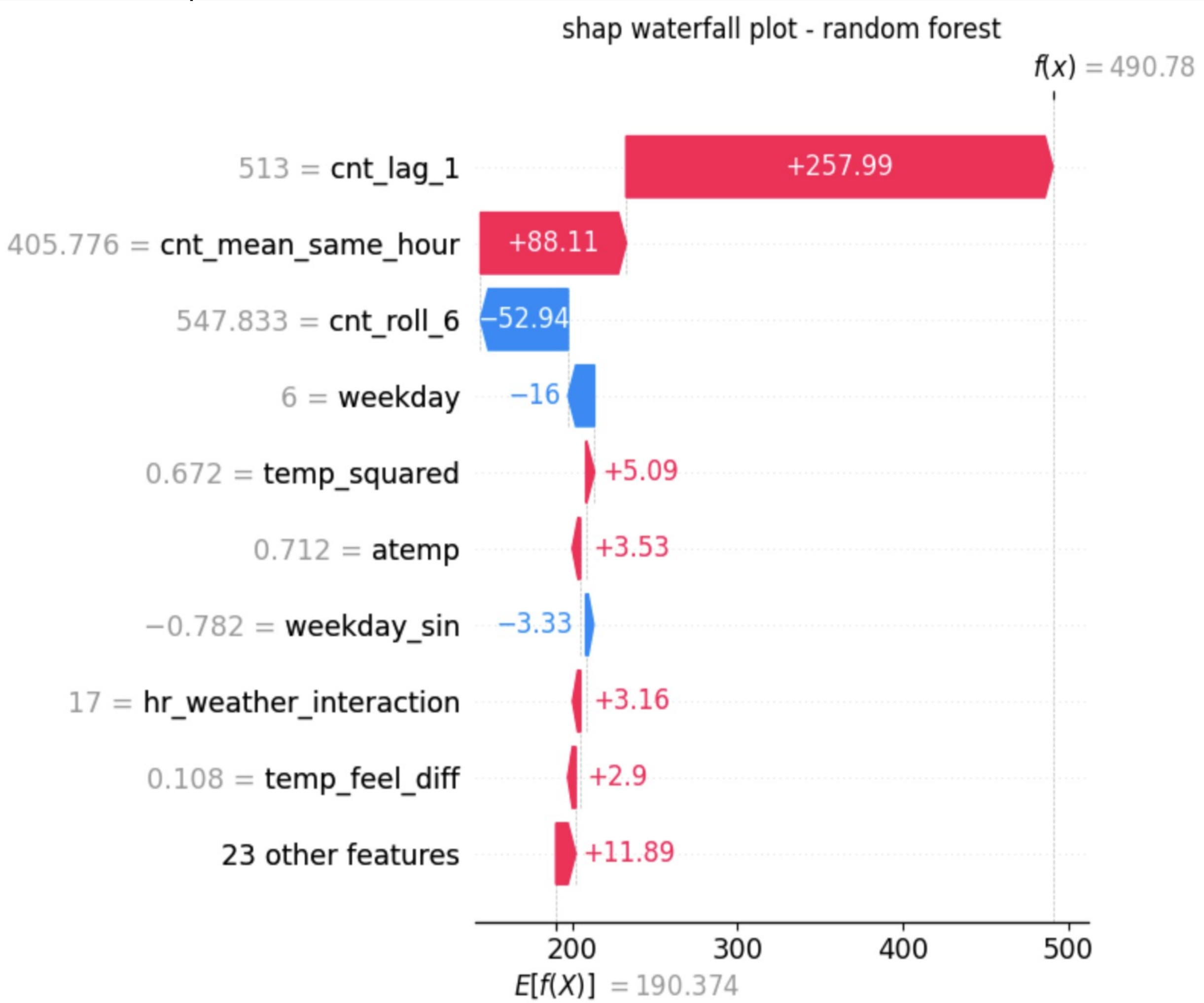
## время суток:

- **is\_morning, is\_afternoon, is\_evening, is\_night** - разные периоды суток имеют разную привлекательность для велопроката в зависимости от целей использования (работа, досуг).



# Важность признаков

Наиболее важными оказались признаки, связанные с историческими значениями спроса (cnt\_lag\_1, cnt\_mean\_same\_hour) и скользящими средними.





# ВЫВОД

По итогу оказалось, что лучшая модель была выбрана в самом начале для нашего датасета и это RandomForest, которая дает такие метрики:  $RMSE = 139.2$ ,  $MAE = 104.8$  и  $R^2 = 0.39$



# Команда



**Егор Шукурлаев**



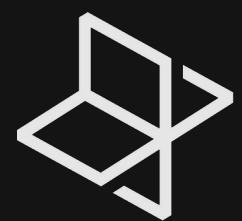
**Мария Чурилова**



**Дмитрий Зейтц**

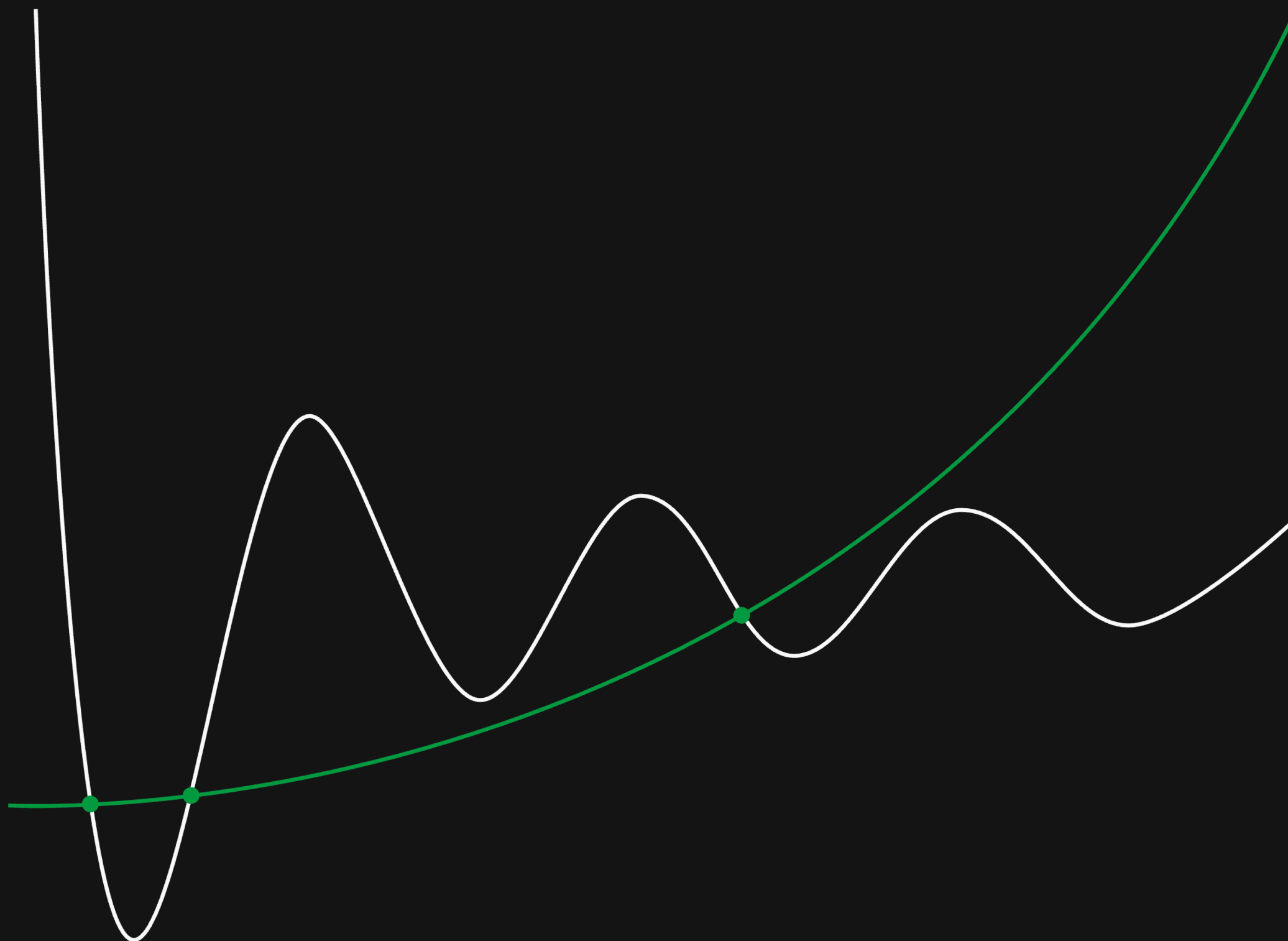


**Матвей Старых**



ЦЕНТРАЛЬНЫЙ  
УНИВЕРСИТЕТ

Спасибо  
за внимание





**Вопросы**