

Загрязнение воздуха в городе Сеул, Южная Корея

Проект подготовили:
Ахтырский Василий,
Бычин Ярослав,
Старшова Дарья



Основные цели проекта

1. Выявление общей тенденции изменения качества воздуха в Сеуле с течением времени. Сравнение качества воздуха с нормативными показателями.
2. Анализ изменений в составе воздуха за определенные промежутки времени:
 - i) анализ суточных колебаний изменений концентраций загрязняющих веществ.
 - ii) анализ сезонных колебаний.
3. Исследование качества воздуха в районах Сеула. В частности, выявление самого "грязного" района и самого "чистого" района. Построение карты качества воздуха в Сеуле.
4. Анализ корреляций между загрязняющими веществами.
5. Прогнозирование уровня загрязнения.



Первичная обработка файлов

- Measurement_summary

Основные измерения (SO₂, NO₂, O₃, CO, PM₁₀, PM_{2.5}).

- Measurement_info

Информация о измерениях (дата, станция, статус прибора).

- Measurement_item_info.csv

Нормативные показатели загрязняющих веществ.

- Measurement_station_info

Информация о станциях мониторинга.

Особенности при обработке:

0 - в норме

1 - требуется калибровка

2 - нештатная ситуация

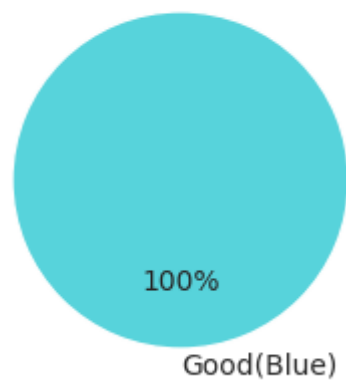
4 - отключение питания

8 - ремонтируется

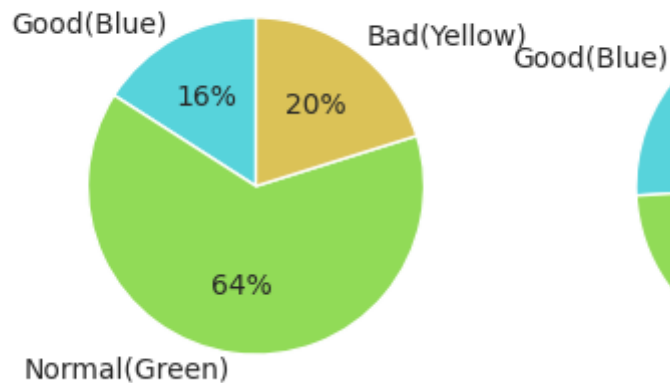
9 - нештатные данные

**Вывод после анализа файлов: 97% - в норме,
3% - некорректные измерения — удаляем их!**

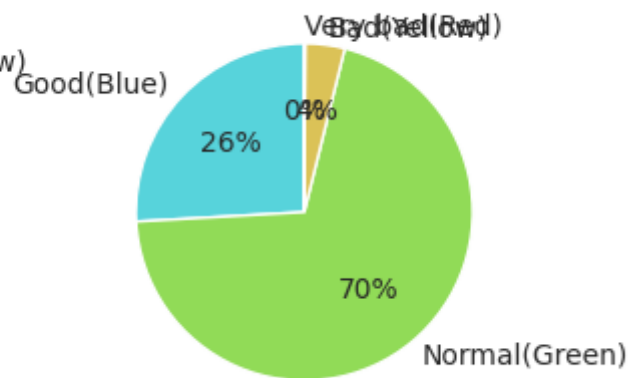
SO2



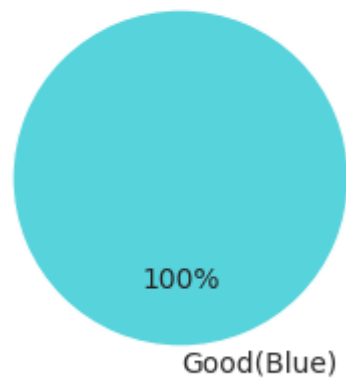
NO2



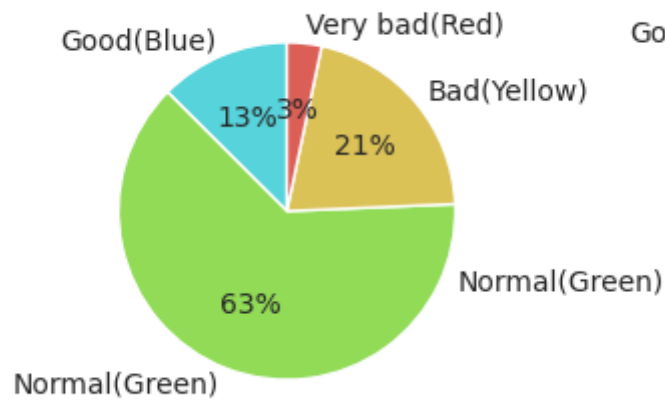
O3



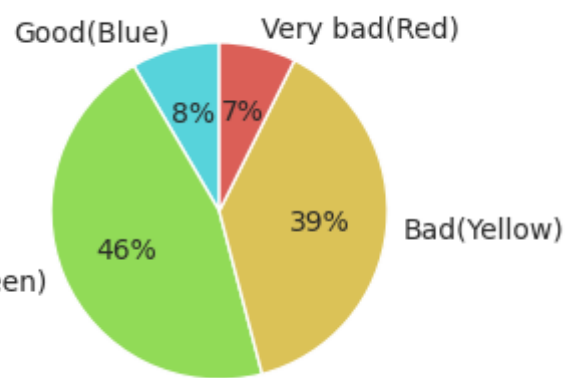
CO

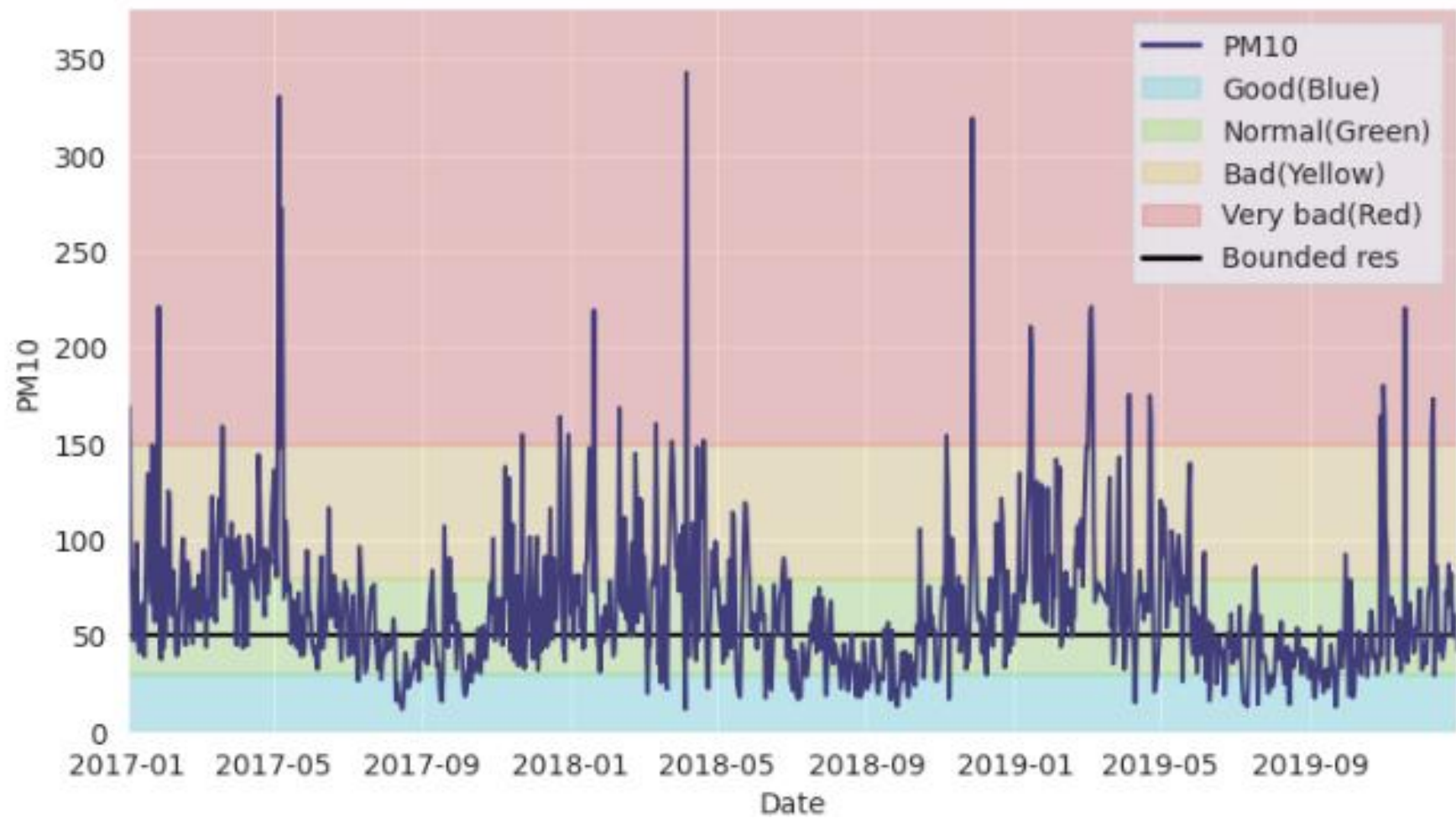


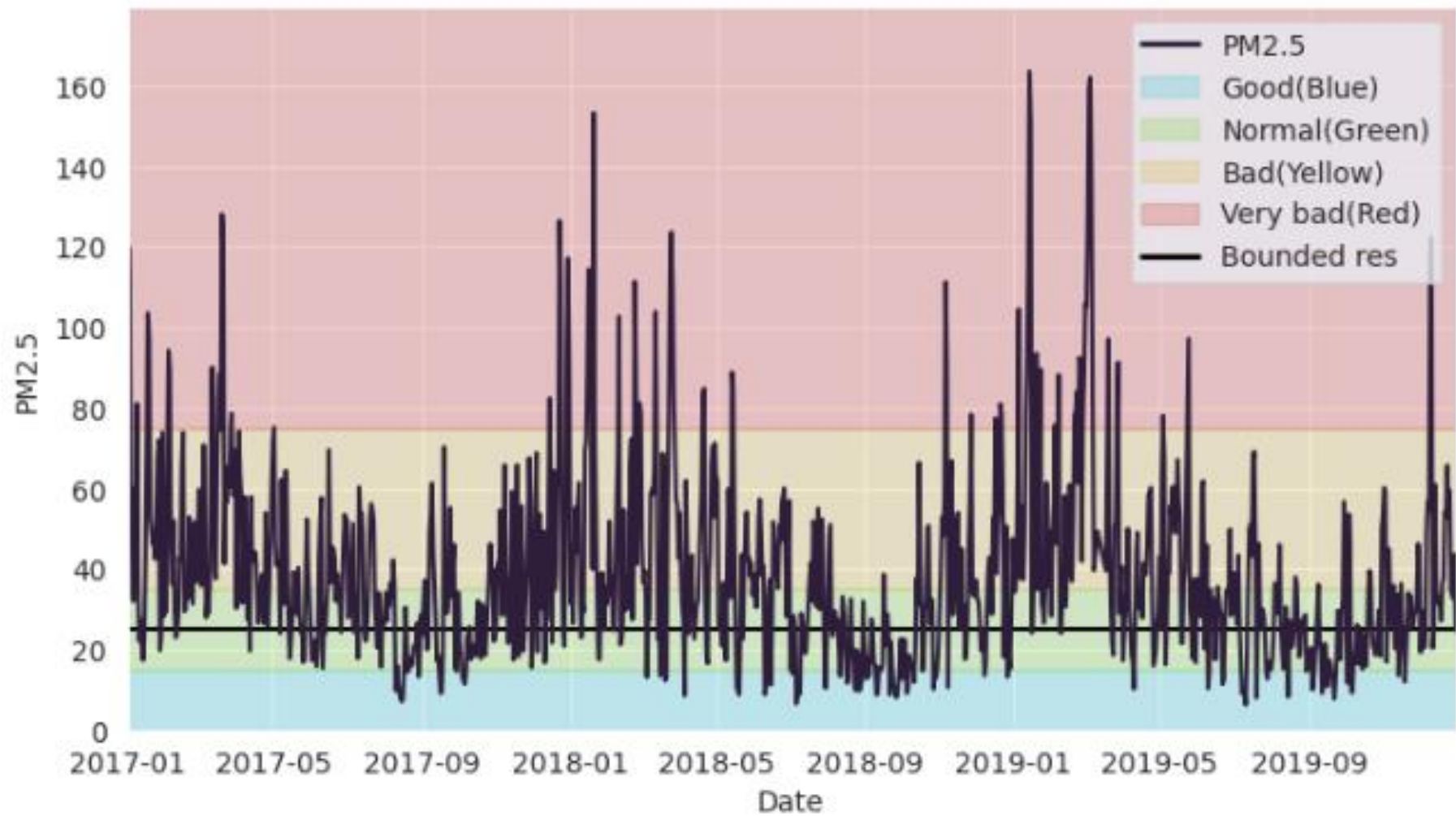
PM10



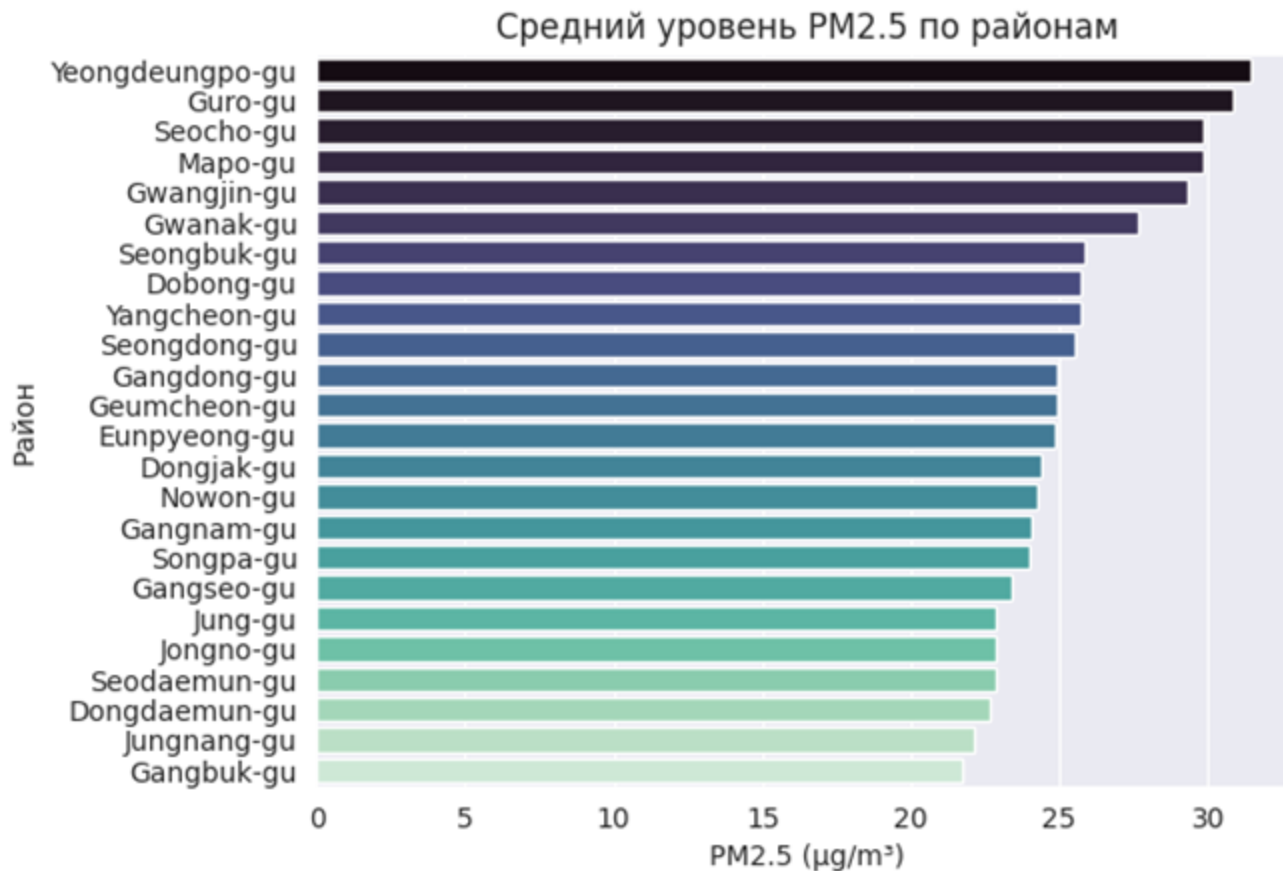
PM2.5



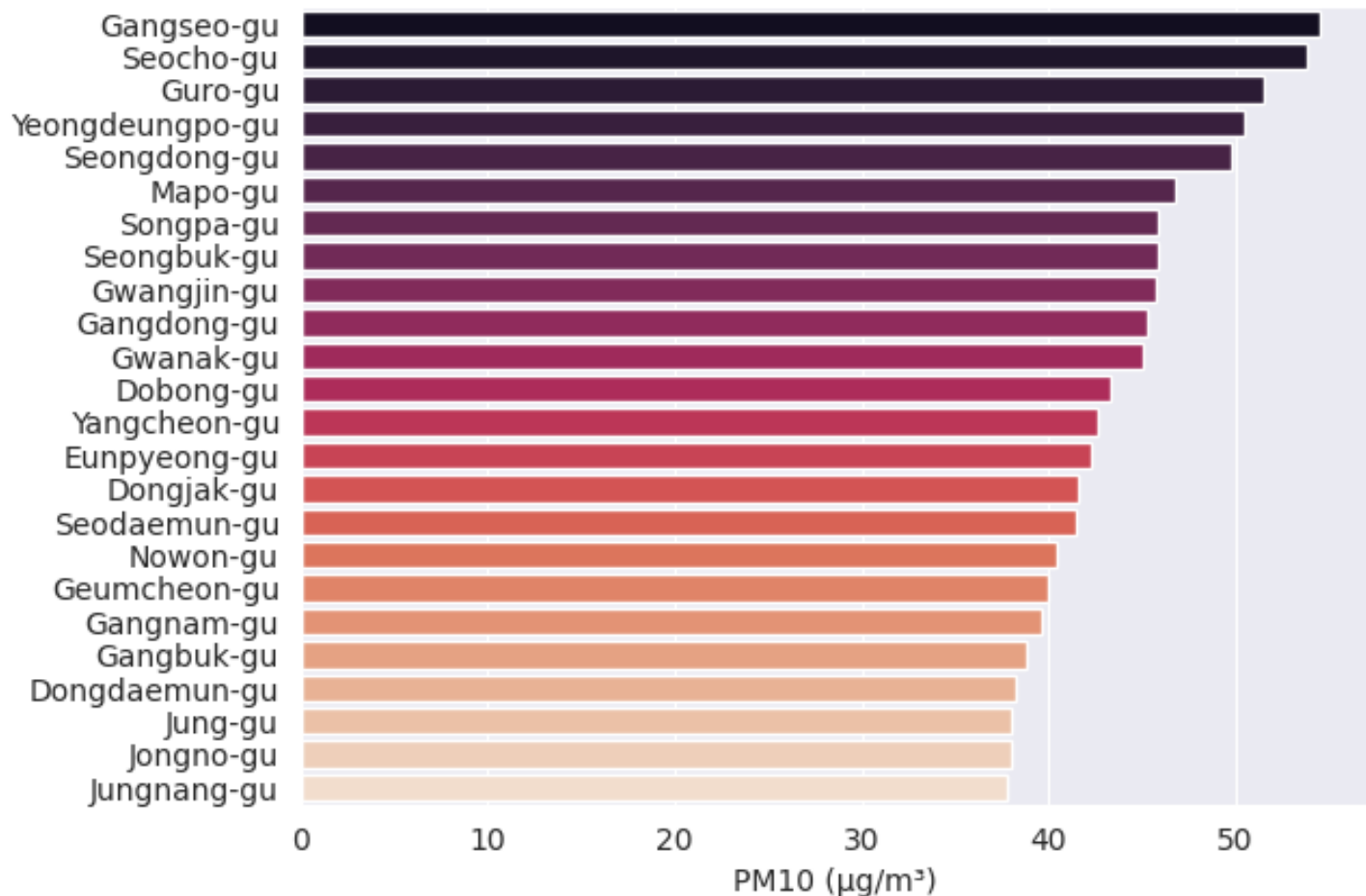




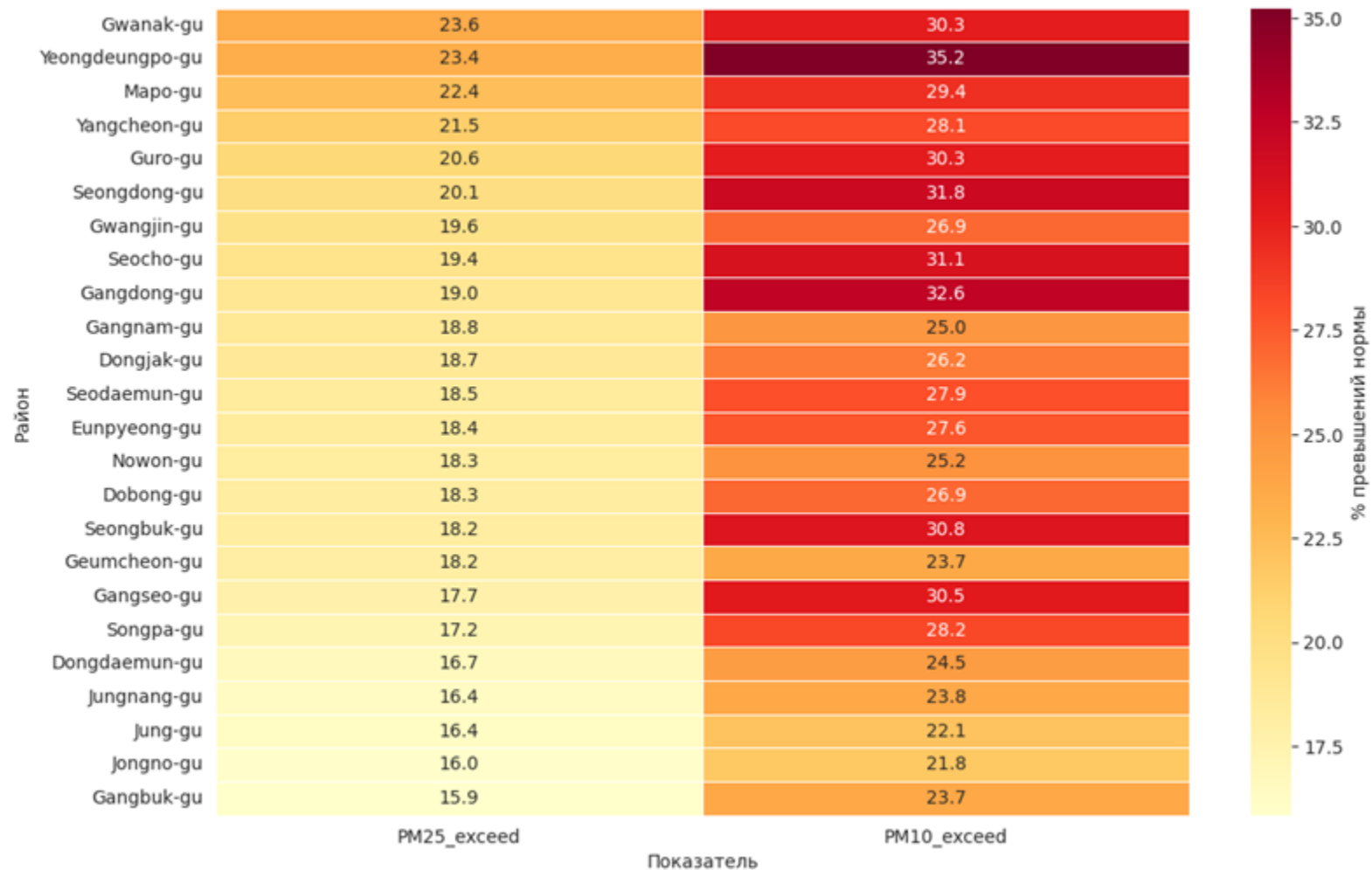
Географический анализ загрязнения воздуха в Сеуле



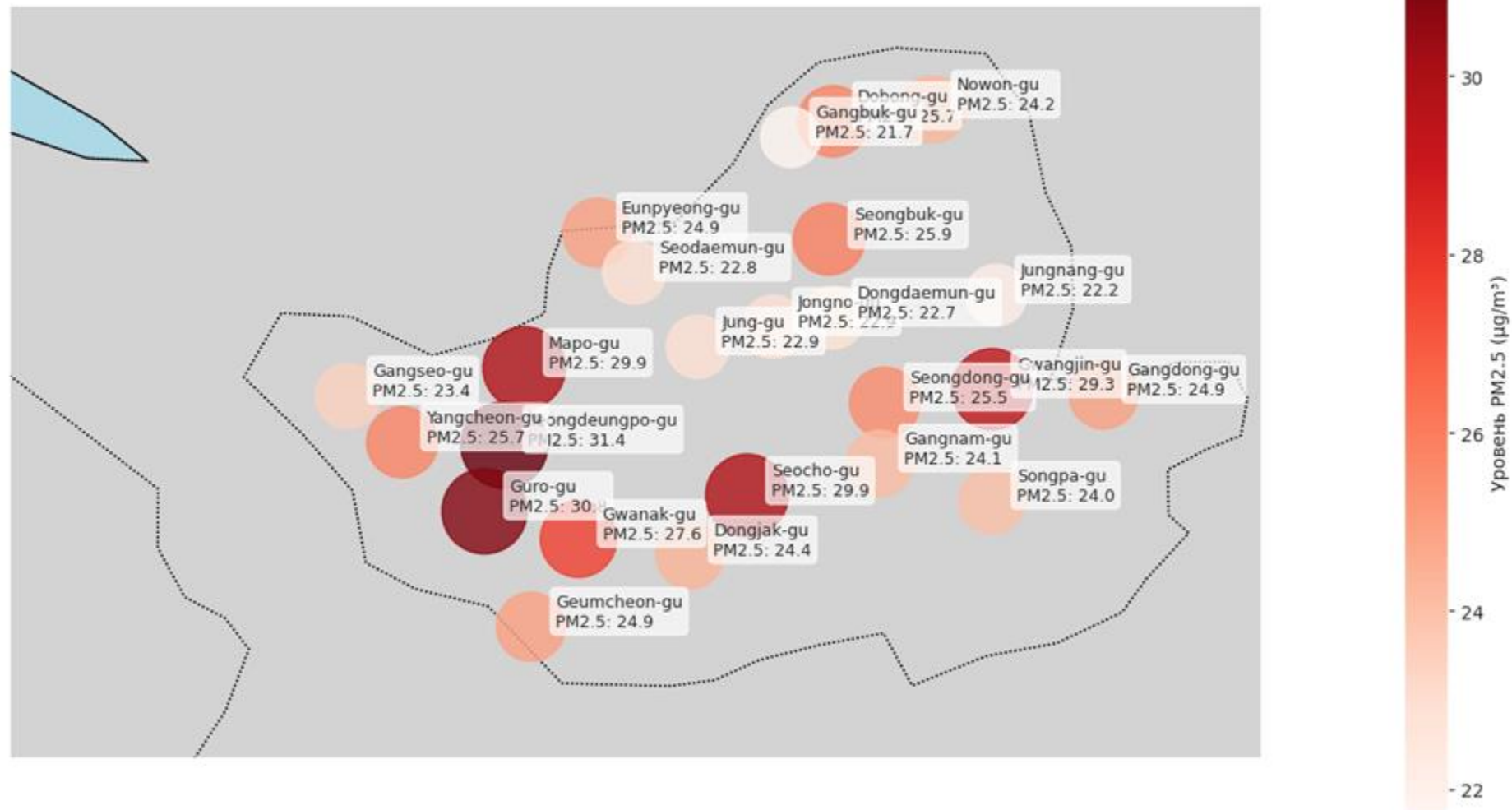
Средний уровень PM10 по районам



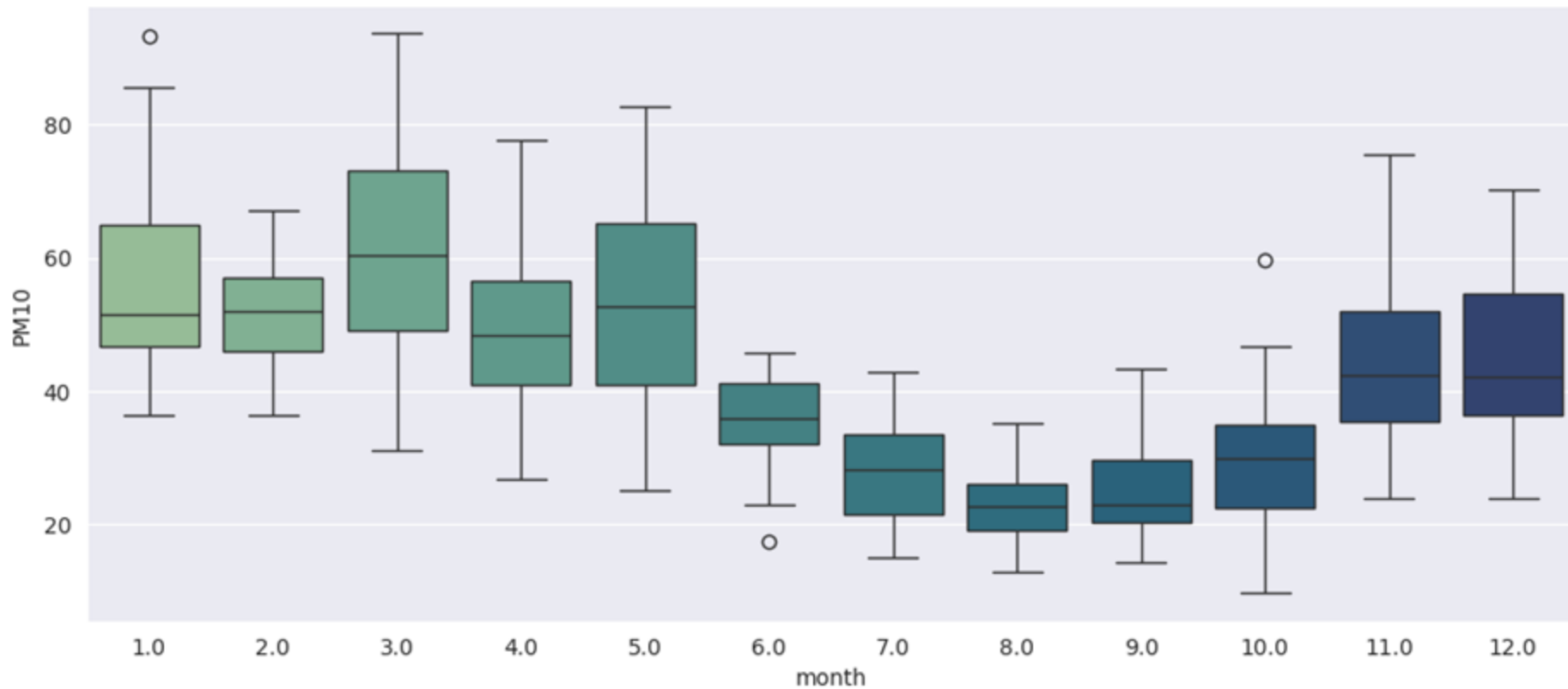
Процент превышений нормативов по районам

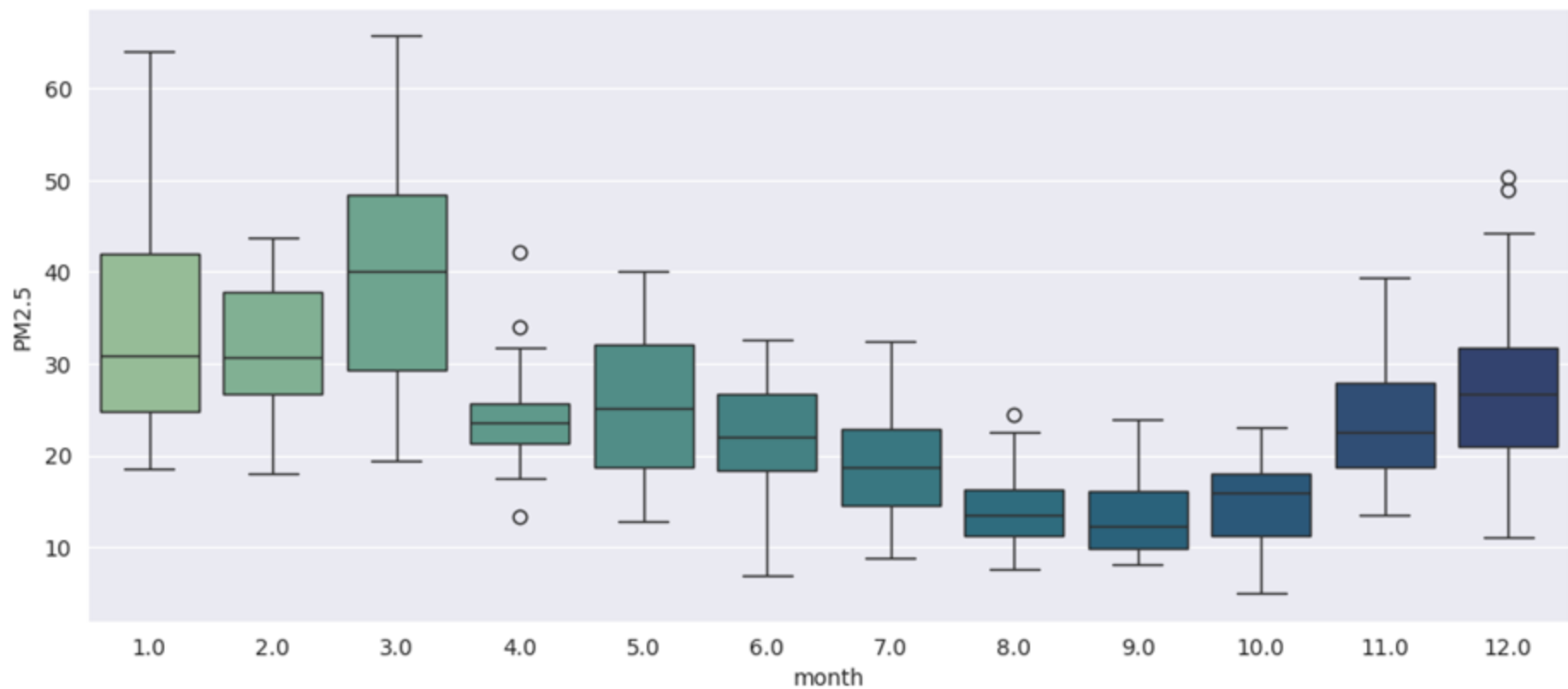


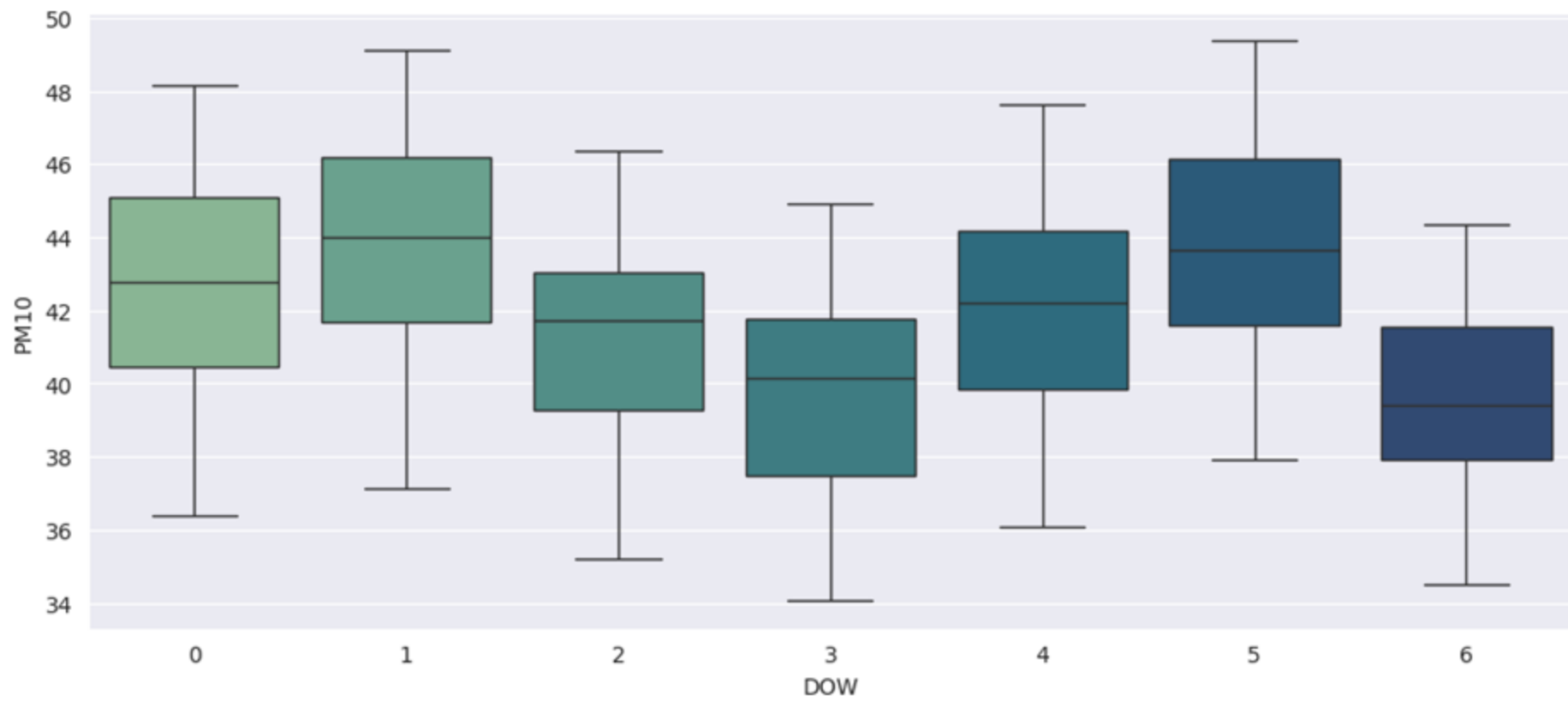
Распределение загрязнения воздуха (PM2.5) по районам Сеула

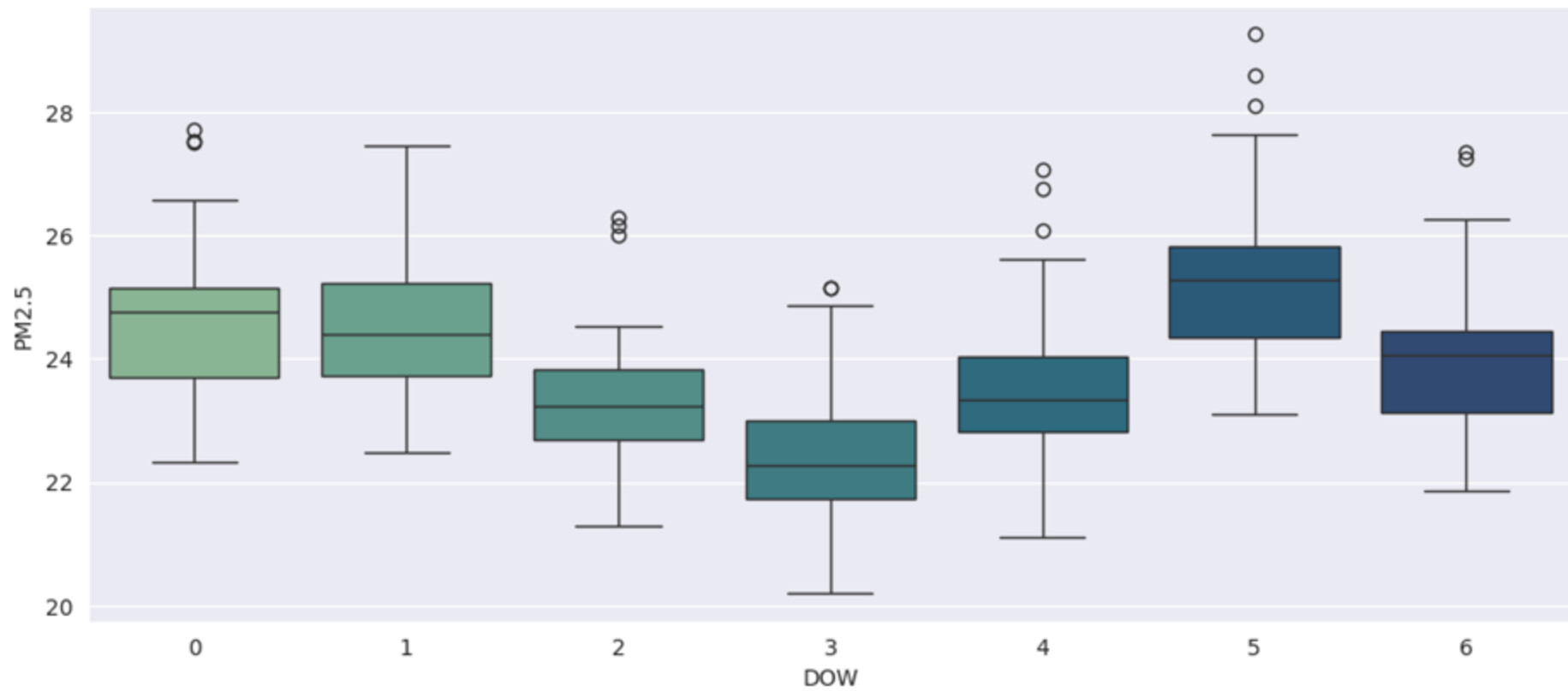


Сезонный анализ загрязнения воздуха в Сеуле

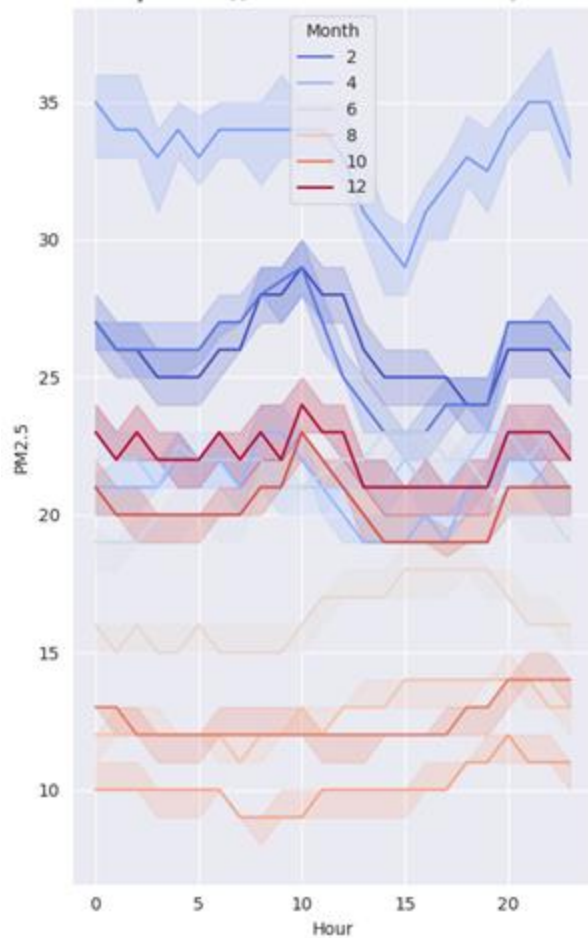




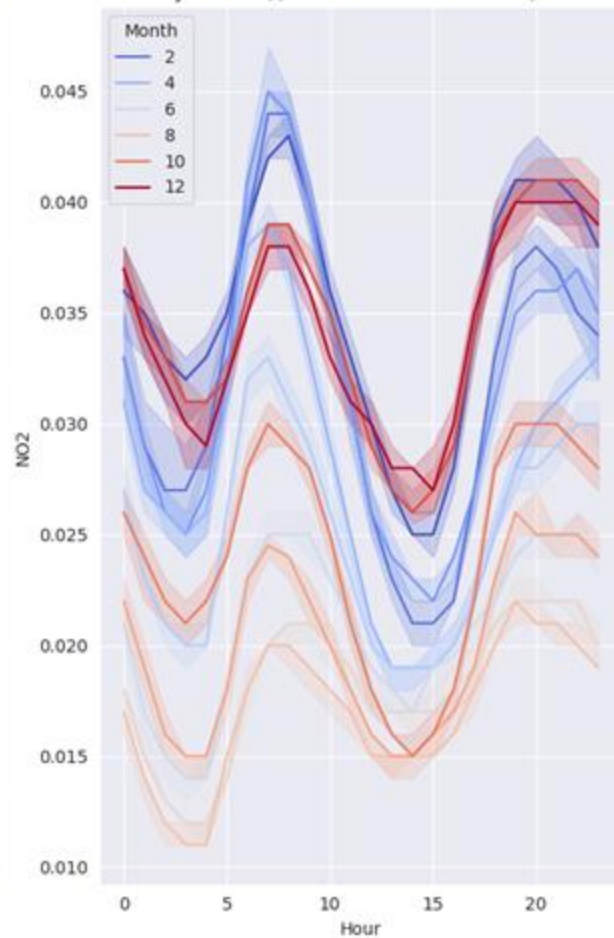




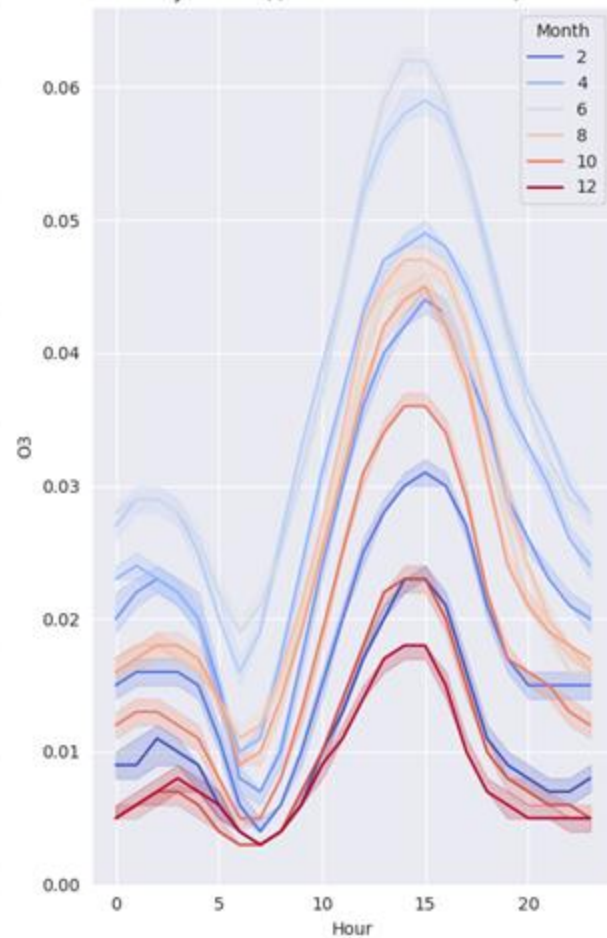
Суточная динамика PM2.5 по месяцам



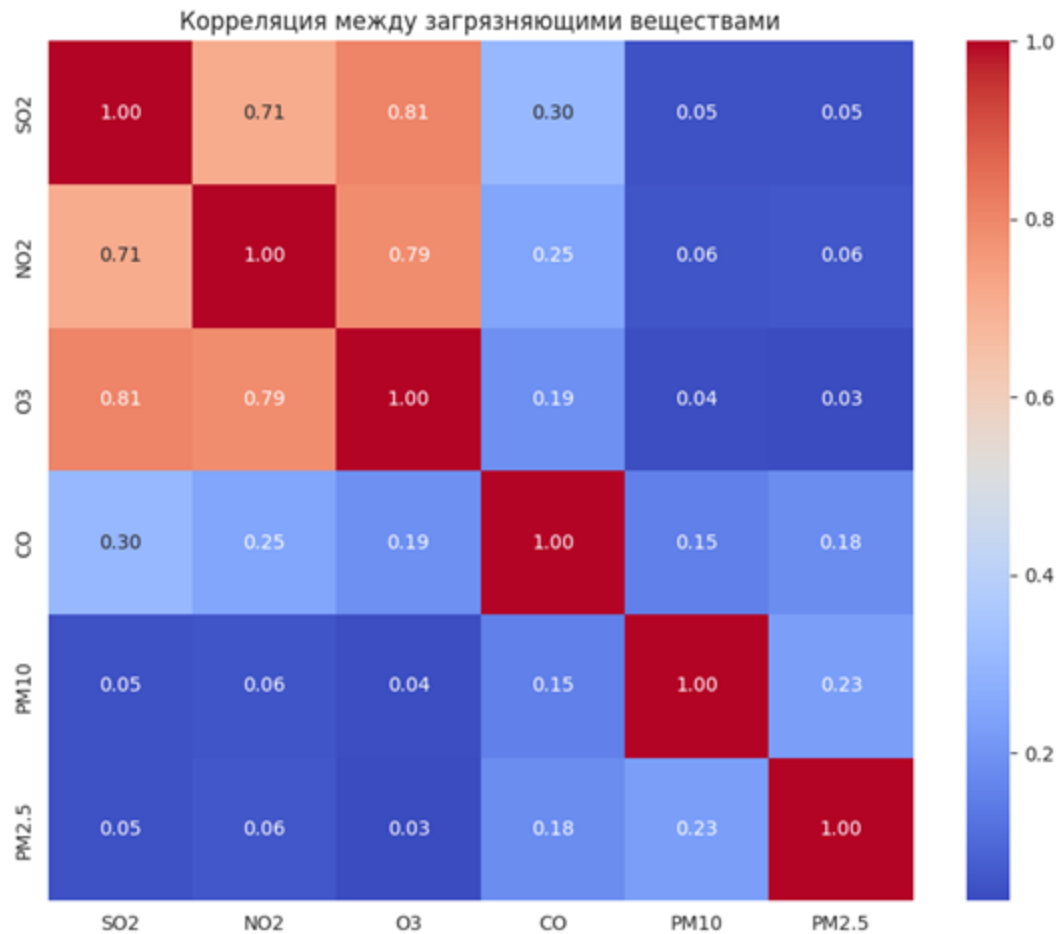
Суточная динамика NO2 по месяцам



Суточная динамика O3 по месяцам



Корреляций между загрязняющими веществами



Задача регрессии

Постановка задачи: предсказание уровня концентрации мелкодисперсных частиц Р.М.2.5

Используемые метрики:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

MSE:

Преимущества:

- Чувствительна к большим ошибкам
- Дифференцируема

Недостатки:

- Чрезмерно штрафует выбросы

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

MAE:

Преимущества:

- Устойчива к выбросам
- Результат в исходных единицах измерения

Недостатки:

- Не дифференцируема в нуле
- Менее чувствительна к большим ошибкам, чем MSE

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

R2:

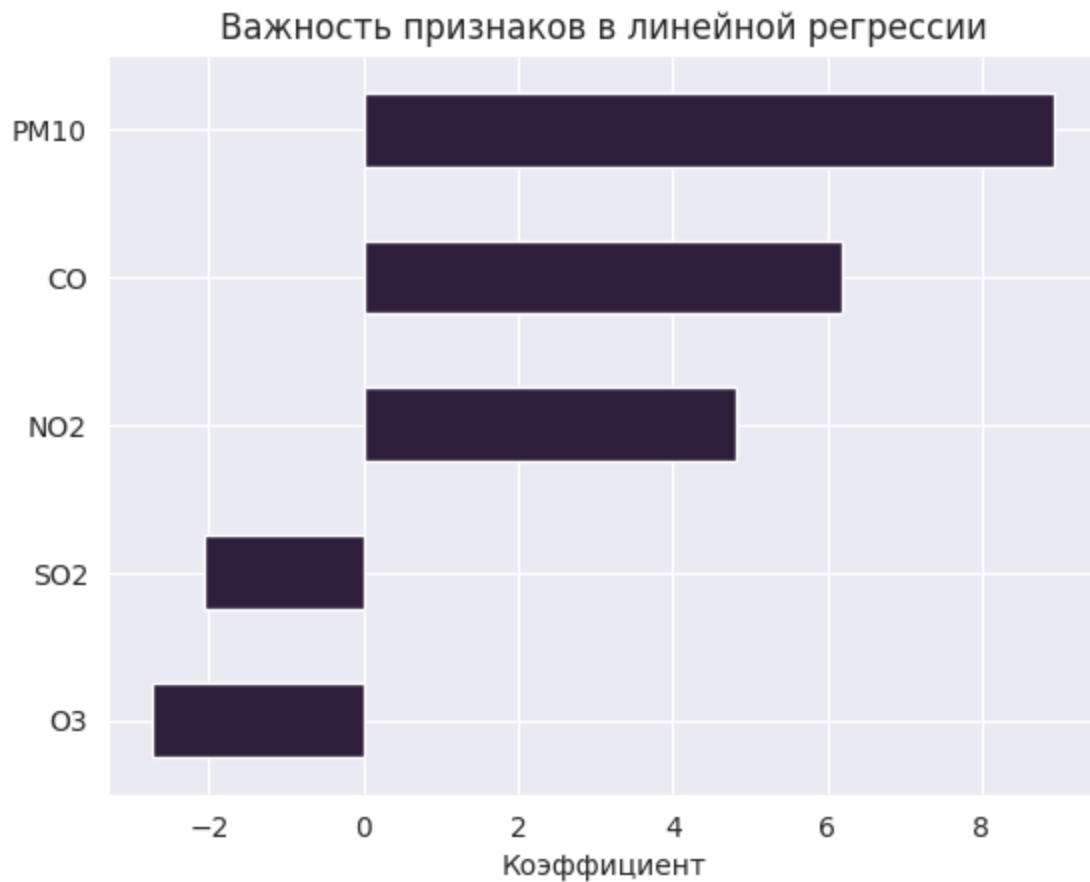
Преимущества:

- Интерпретируема как процент объясненной вариации

Недостатки:

- Может быть отрицательной для плохих моделей
- Чувствительна к выбросам

Линейная регрессия



Линейная регрессия

Работает быстро,
но уступает по
точности для
сложных данных

Хорошо
интерпретируемые
коэффициенты

MSE: 1749.66
MAE: 12.10
 R^2 : 0.07



kNN

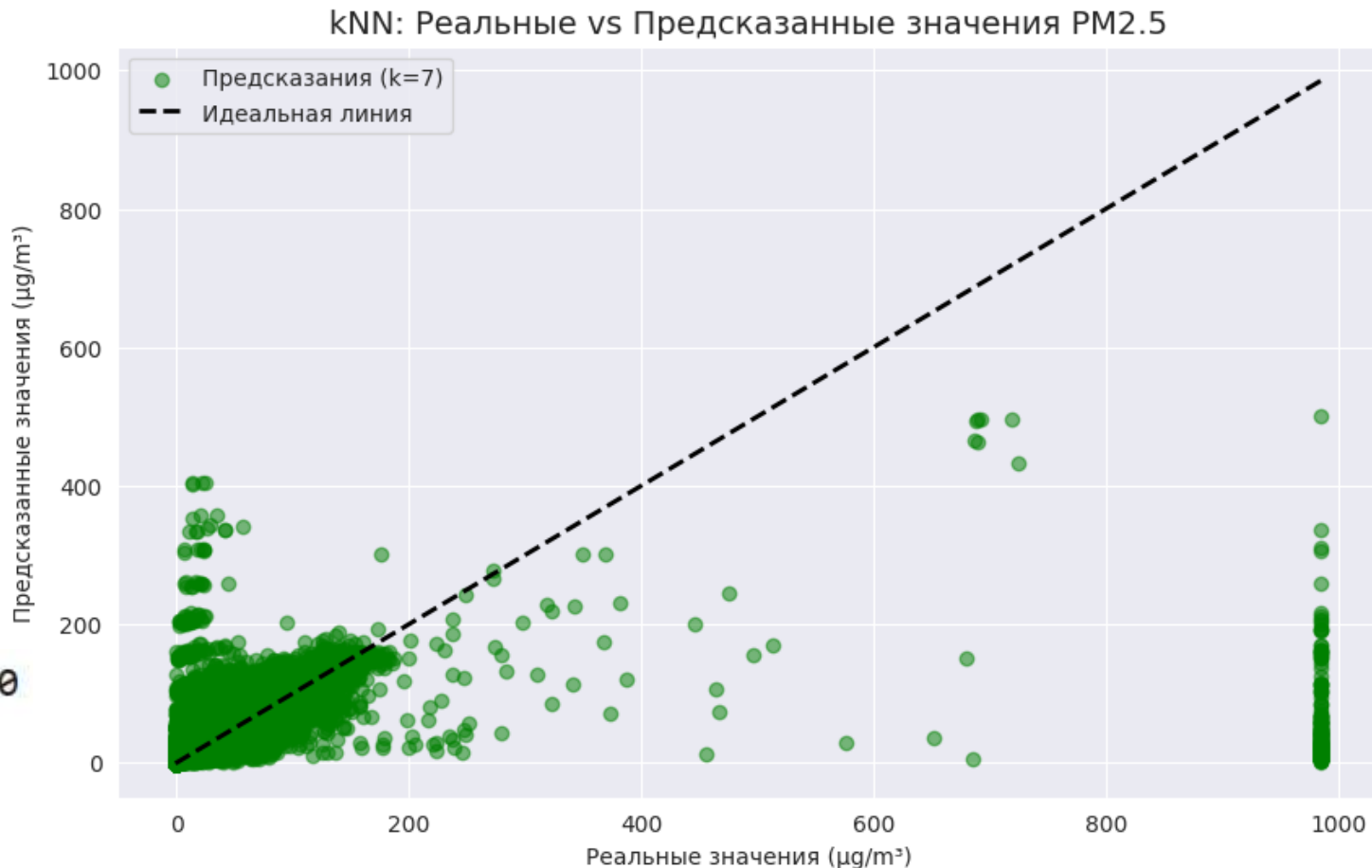
Показывает
улучшение
качества
предсказаний, но
демонстрирует
характерные
"ступеньки" в
прогнозах из-за
природы метода.

Оптимальное k : 20

MSE: 1611.41

MAE: 8.55

R^2 : 0.14



Random Forest

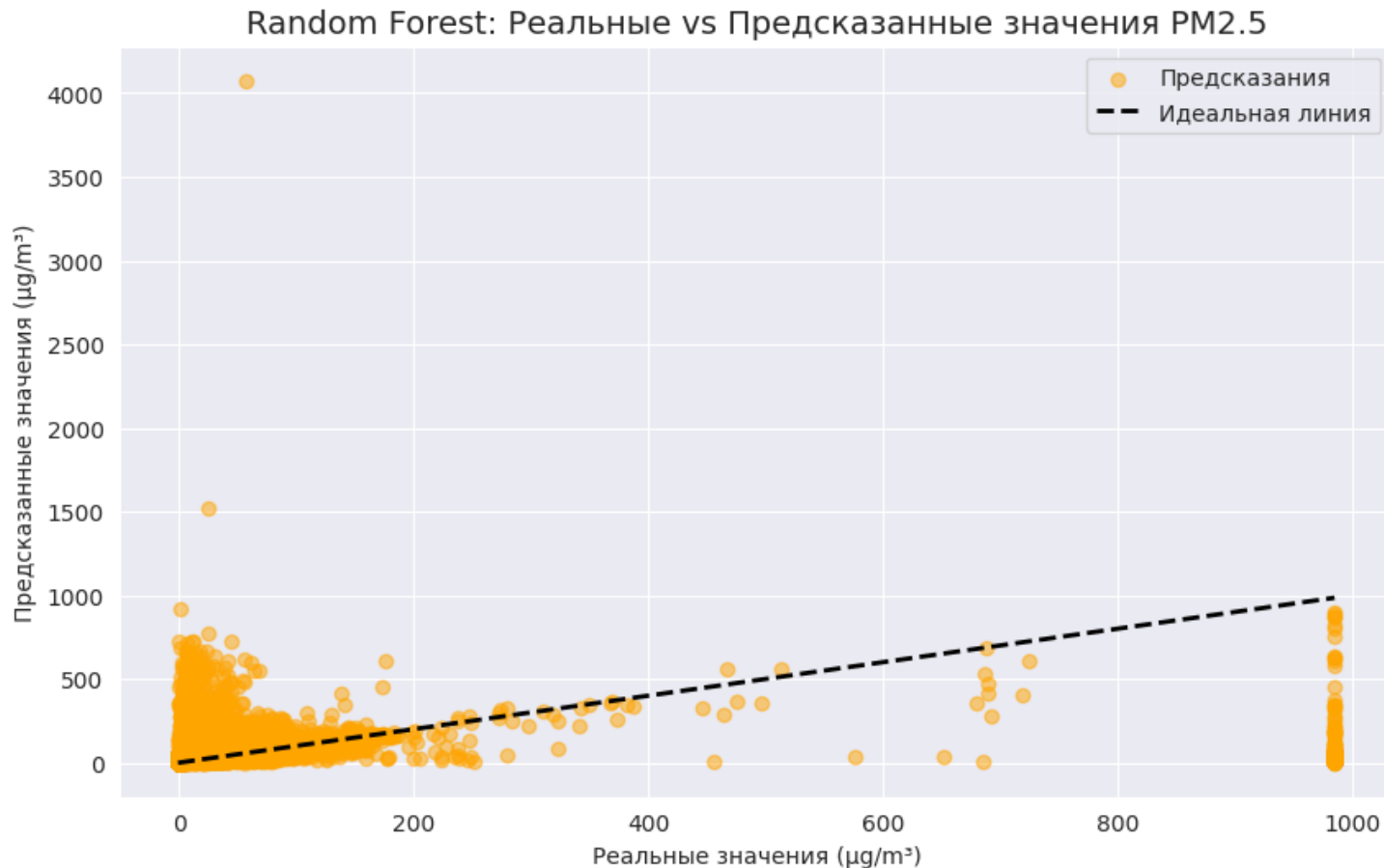
Обеспечивает
наиболее точные
предсказания во
всем диапазоне
значений,
минимальное
систематическое
смещение и
устойчивость к
выбросам

Random Forest:

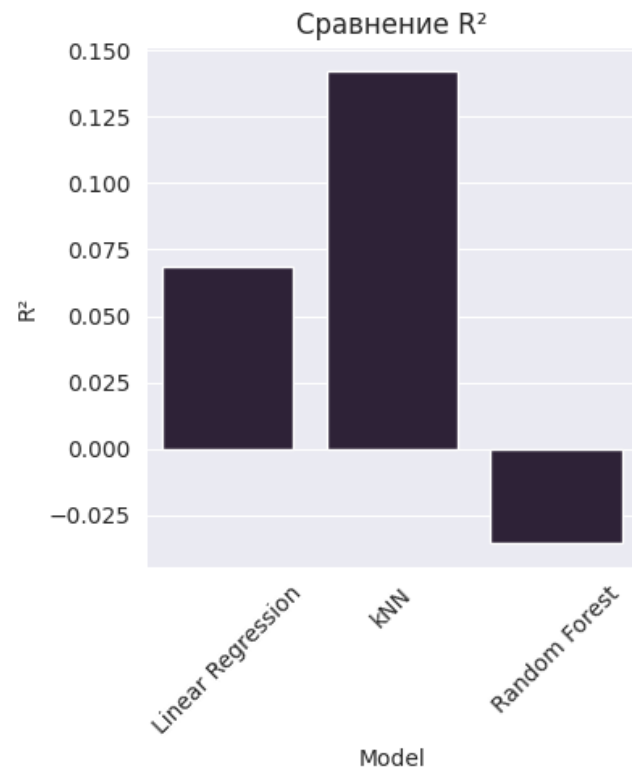
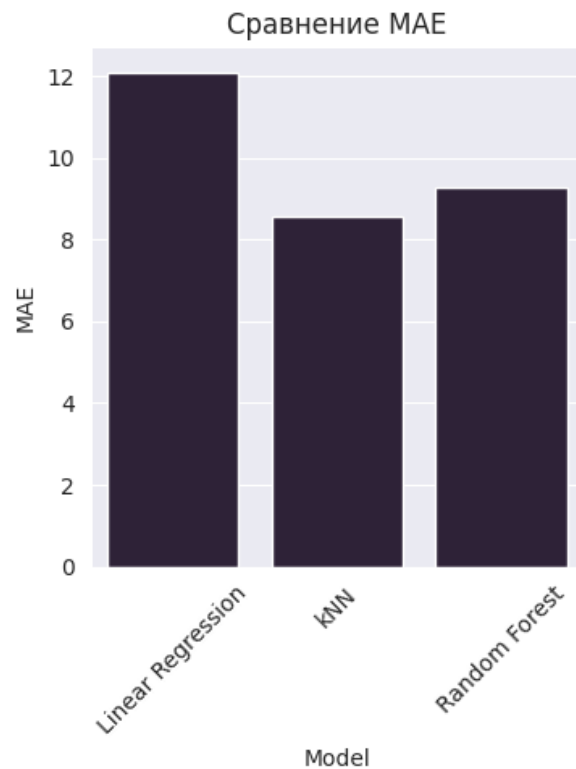
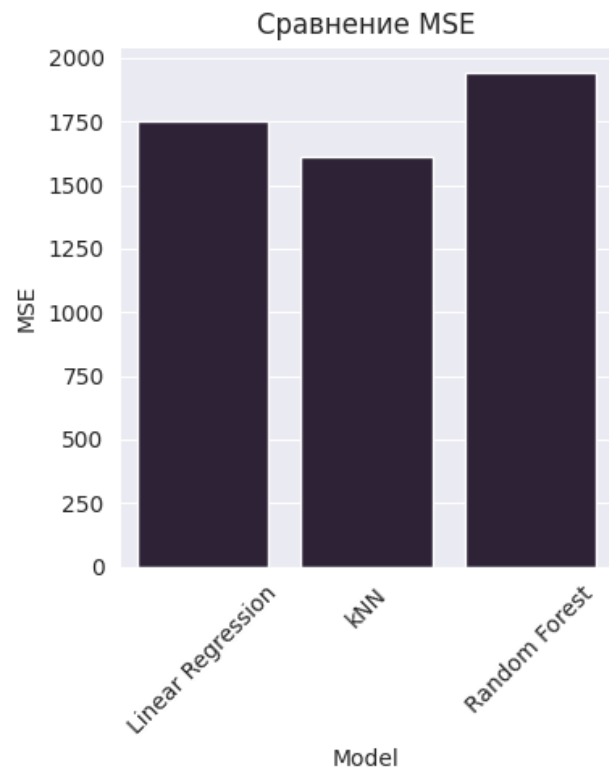
MSE: 1944.89

MAE: 9.28

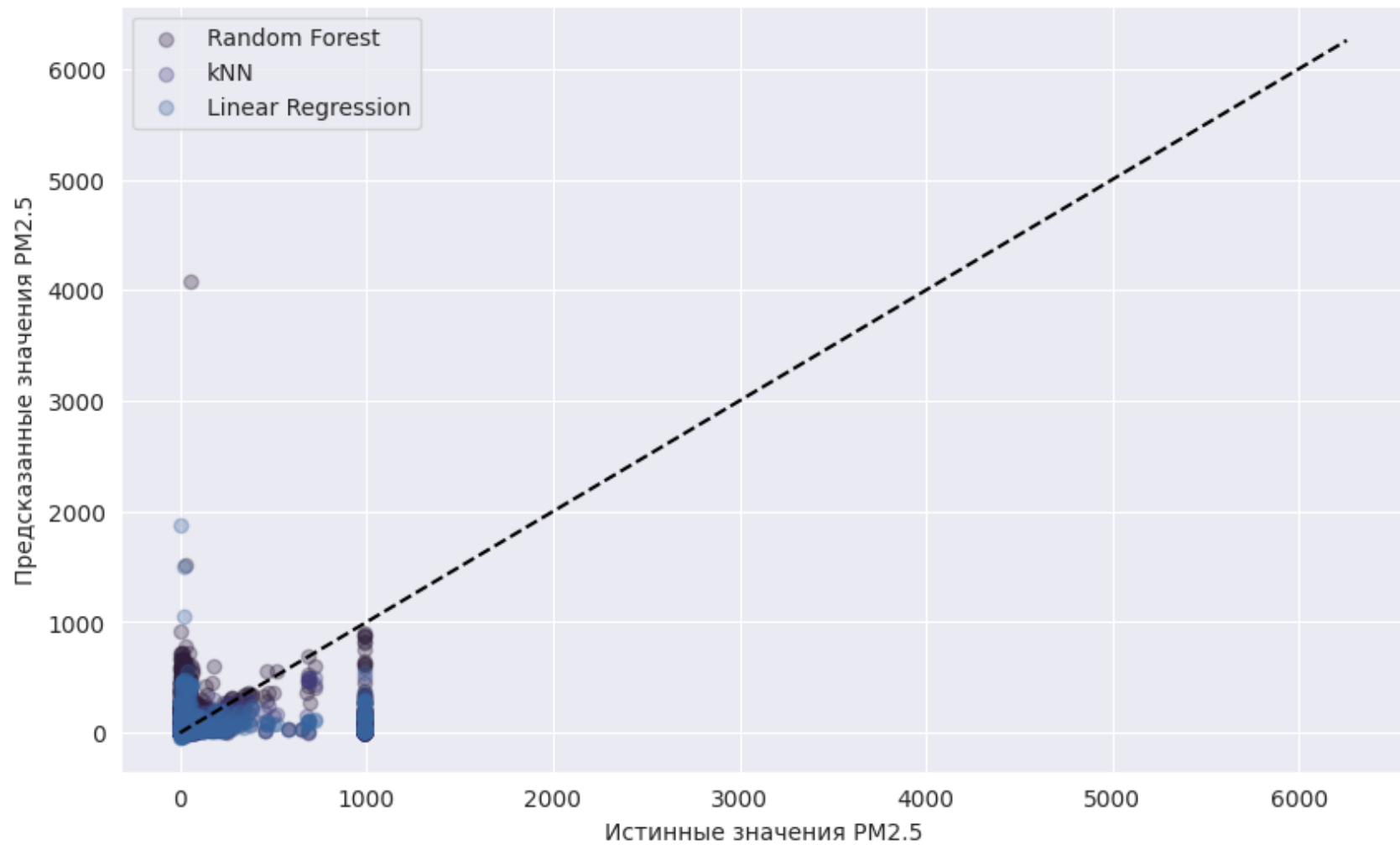
R^2 : -0.04



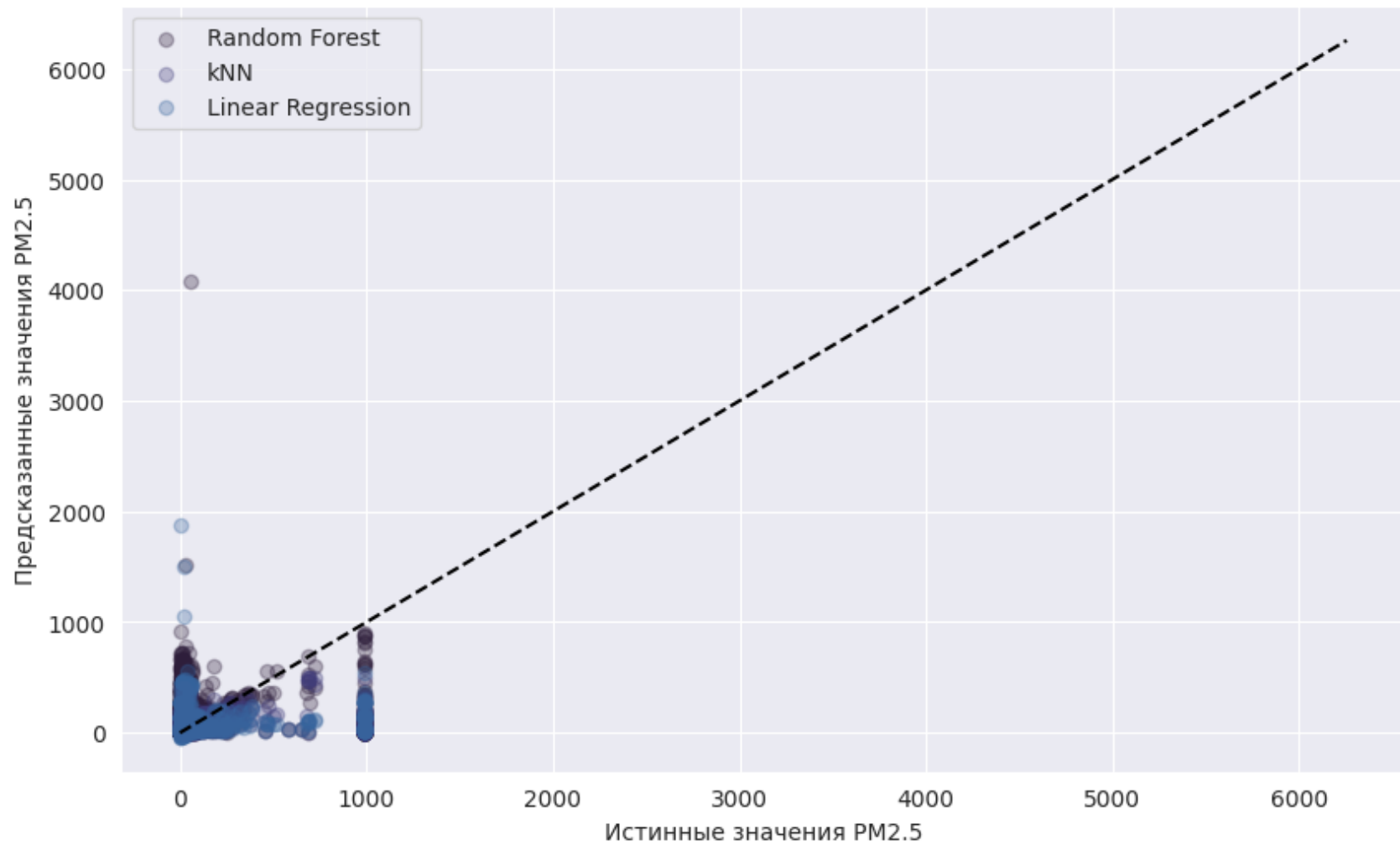
Результаты сравнения



Сравнение предсказаний моделей



Сравнение предсказаний моделей



ЗАДАЧА КЛАССИФИКАЦИИ

Автоматически определять уровень опасности загрязнения воздуха частицами PM_{2.5}, относя каждое наблюдение к одному из заранее определённых классов качества воздуха.

Логистическая регрессия

=== Логистическая регрессия ===

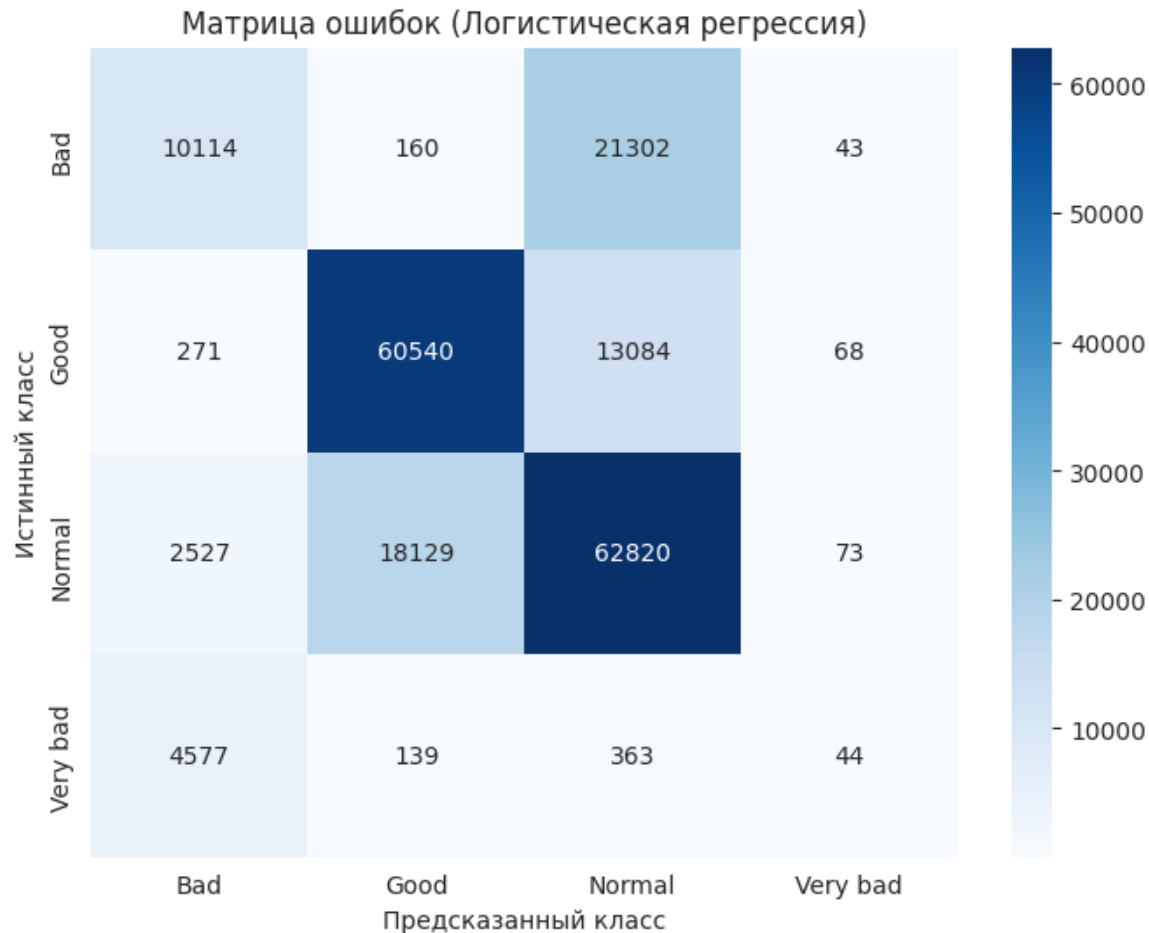
Accuracy: 0.6873371976896229

Precision (macro): 0.545445128703609

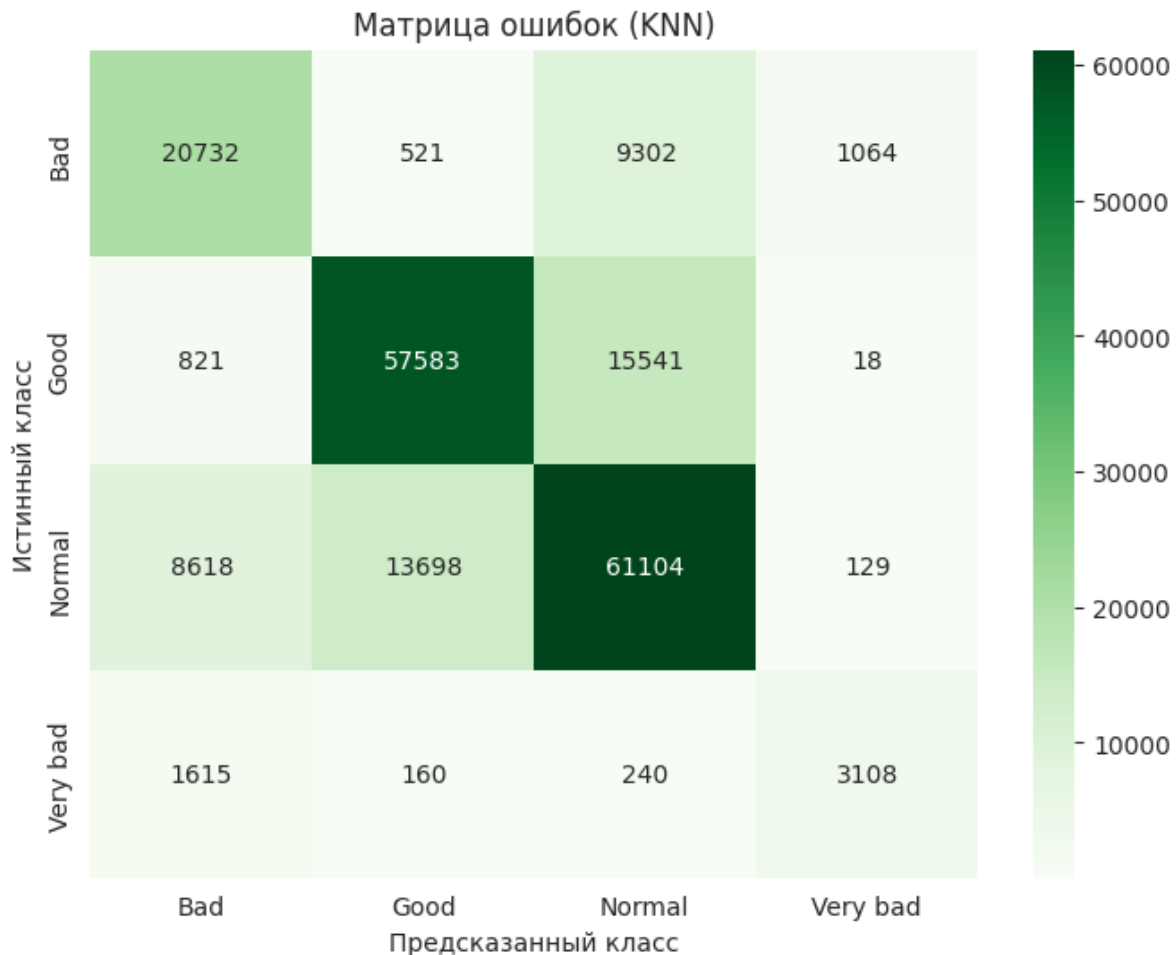
Recall (macro): 0.4747177865187723

F1-Score (macro): 0.47844374500992226

ROC-AUC (ovr): 0.8806502193156305



kNN



=== KNN (k=5) ===

Accuracy: 0.7337146210631441

Precision (macro): 0.7202510236417956

Recall (macro): 0.6930626641225701

F1-Score (macro): 0.7053730398463454

ROC-AUC (ovr): 0.8860601153780329

Random Forest

=== Random Forest ===

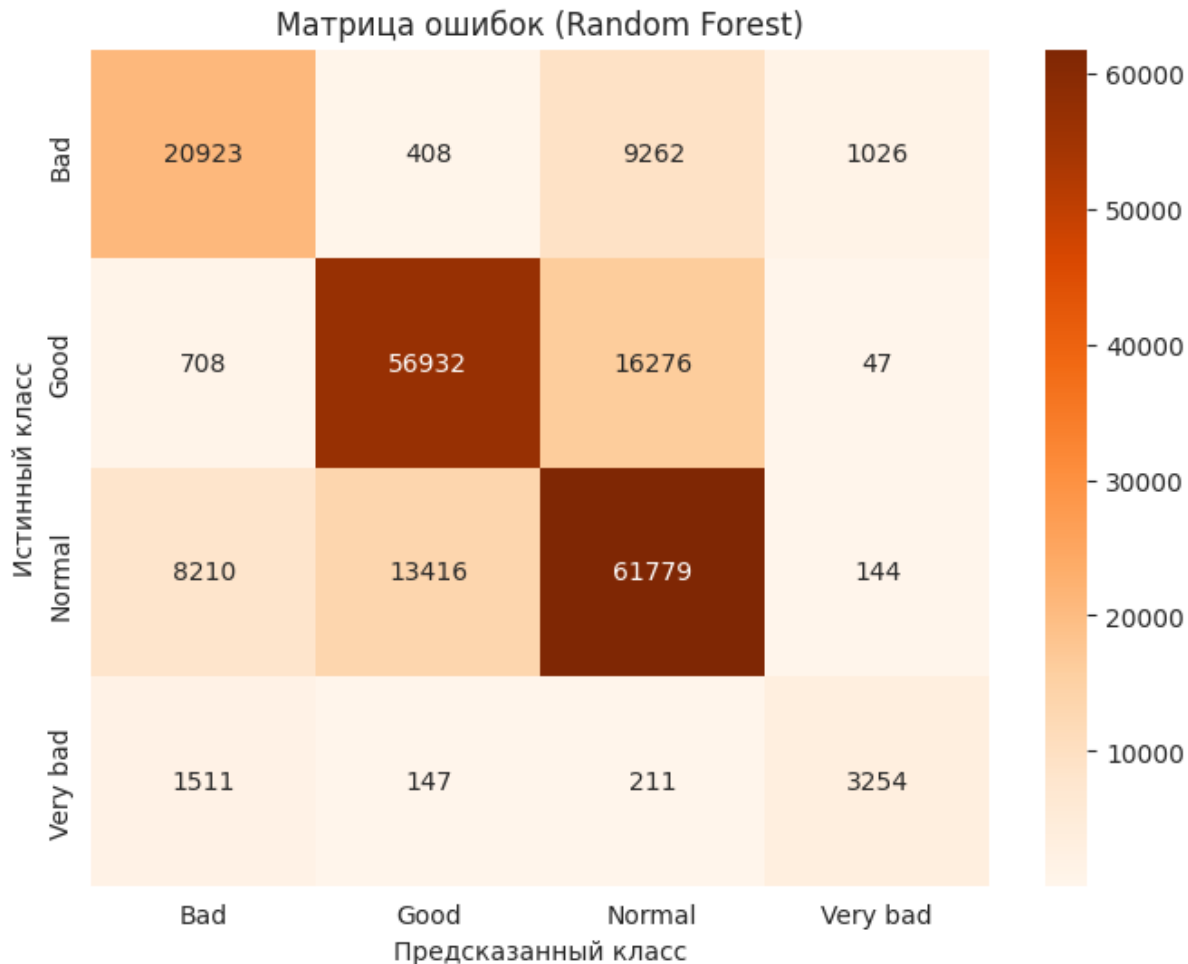
Accuracy: 0.7355730126535361

Precision (macro): 0.7259837835829319

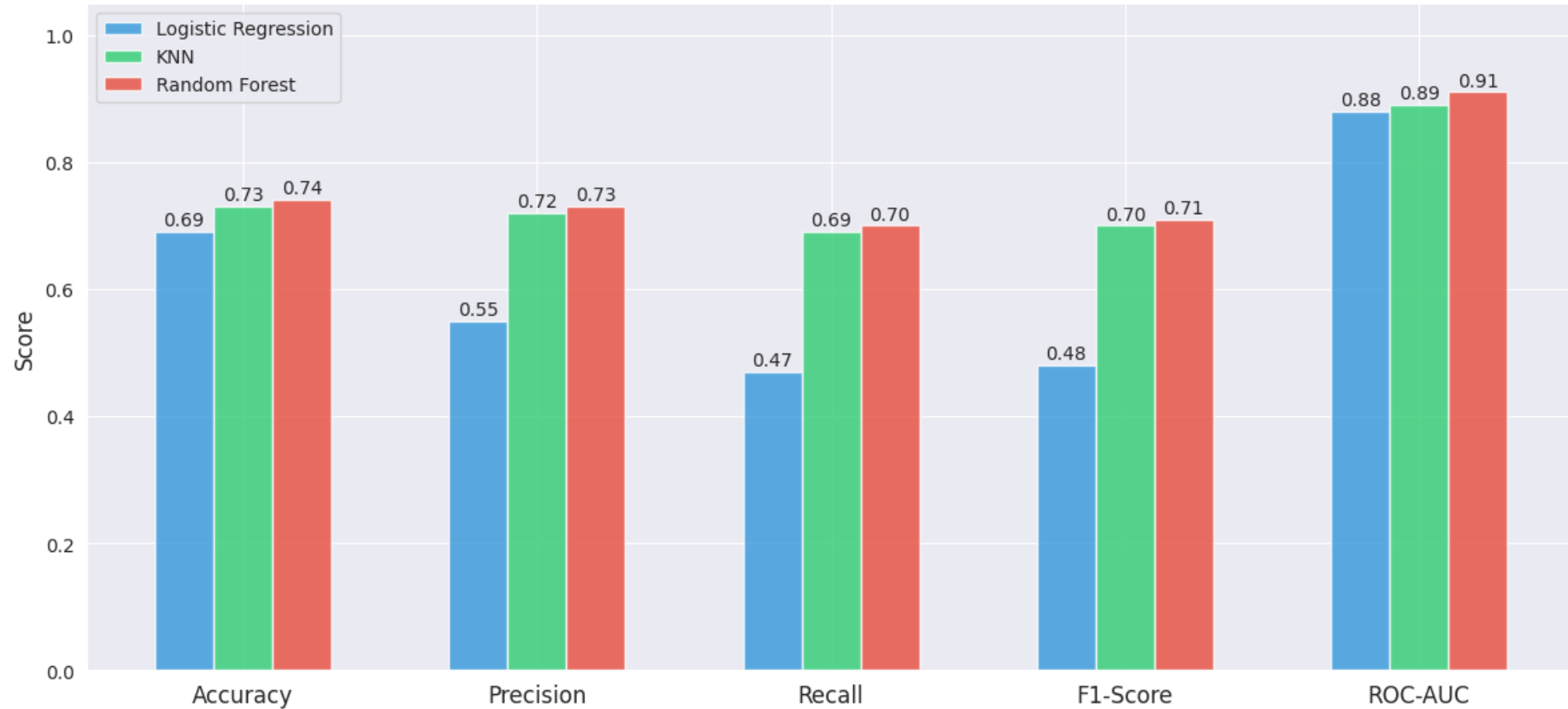
Recall (macro): 0.7015169119534753

F1-Score (macro): 0.7127750663351549

ROC-AUC (ovr): 0.9121539789863153



Сравнение моделей по метрикам классификации



Спасибо за внимание!

