

# Определение скоростей и размеров капель топлива при помощи машинного обучения

Выполнил:

Алешковский Александр Анатольевич

# Техническое задание

- Сбор датасета из спектрограмм 4-х сигналов
- Вырезание окна сигнала, соответствующего капле.
- Определение параметров выреза и способа сглаживания данных
- Дополнение датасета расстоянием между пиками сигналов
- Дополнение датасета столбцами уже имеющихся признаков
- Расчет усредненных характеристик от сигналов
- Выделение геометрических признаков сигналов
- Применение базовой модели catboost для прогнозирования расстояния между пиками

# Входные данные и их обработка



```
1 def finder():
2
3     path = os.getcwd()
4
5     dir_pattern = re.compile(r"^M\d+$")
6     file_pattern = re.compile(r"^M\d+_\d+\.mat$")
7
8     # Список для собранных файлов
9     collected_files = []
10
11     for item in os.listdir(path):
12         item_path = os.path.join(path, item)
13
14         if os.path.isdir(item_path) and dir_pattern.match(item):
15             for file in os.listdir(item_path):
16                 if file_pattern.match(file):
17                     collected_files.append(os.path.join(item_path, file))
18
19     collected_files.sort()
20     print(collected_files)
21
22     return collected_files
```

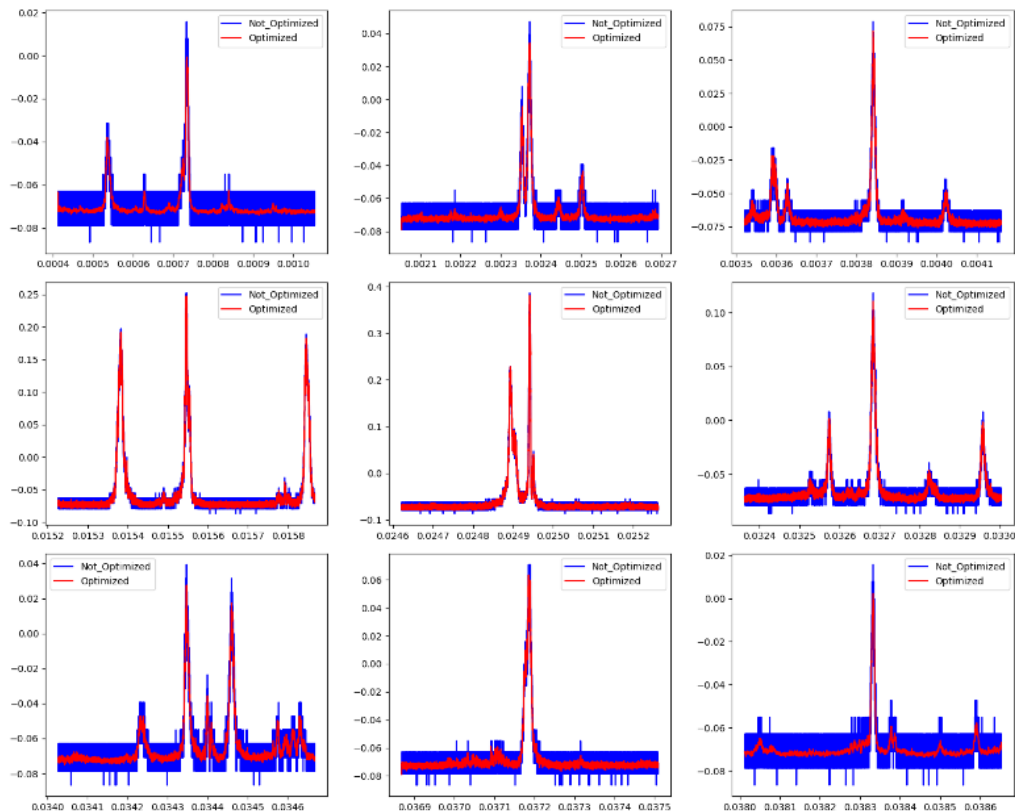
Входные данные представлены набором папок с названием «M{число}», где лежат файлы с названиями «M{Число}\_{Число}.mat»

Файлы поочередно приводятся к формату pandas DataFrame и сглаживаются Low pass фильтром

```
1 def apply_lowpass_filter(data, cutoff_freq=1.5, order=10, sample_rate=5):  
2     '''  
3     Применение фильтра к данным  
4     '''  
5     nyquist_freq = 0.5 * sample_rate  
6     normal_cutoff = cutoff_freq / nyquist_freq  
7     b, a = butter(order, normal_cutoff, btype='low', analog=False)  
8     filtered_data = filtfilt(b, a, data)  
9     return filtered_data
```

Low-pass filter (Фильтр низких частот) - это тип сигнального фильтра, который пропускает сигналы с частотами ниже определенной предела (называемого частотой среза), а блокирует те, которые находятся выше этого предела

# Сбор капель

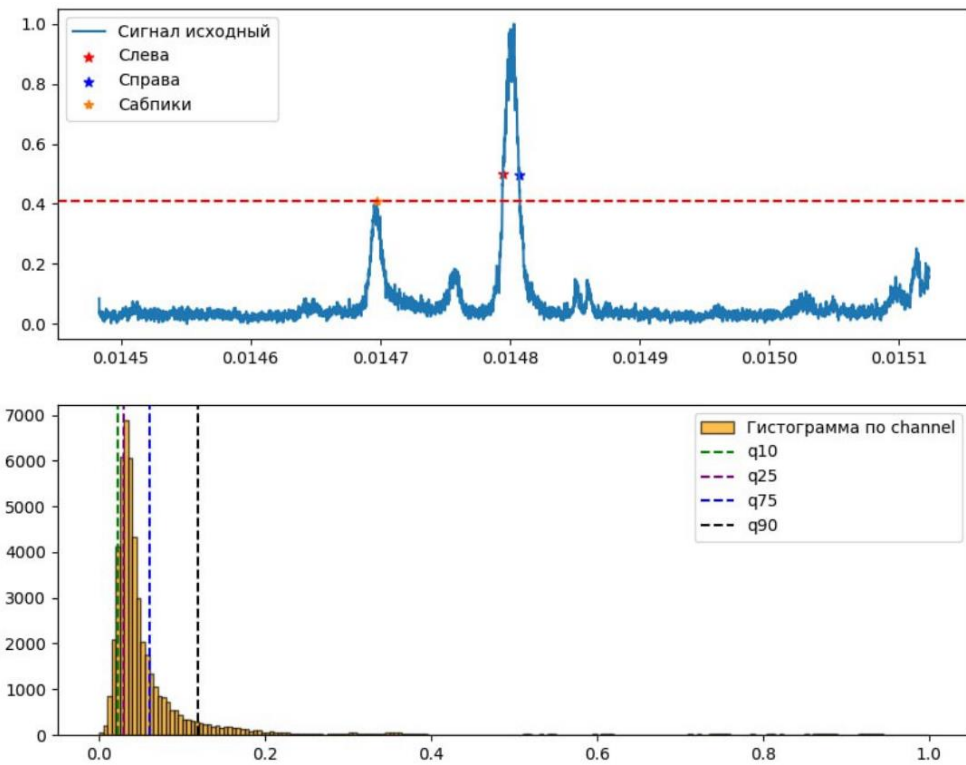


```

1 def raindrops_and_peaks(for_train, height_peak, window_size, butter=False):
2
3     peaks, i = find_peaks(for_train['Channel A'], height=height_peak,
4                             distance=window_size)
5
6     print(str(len(peaks)) + " length of peaks\n")
7
8     # хранение капель
9     setup = np.empty((len(for_train),), dtype=object)
10    print("Создано хранилище капель")
11
12    # поиск пиков и вырез окна
13    for i, peak in enumerate(peaks):
14
15        start = max(peak - window_size // 2, 0)
16        end = min(peak + window_size // 2, len(for_train))
17
18        # Вырезаем окно для каждого канала
19        window = for_train.loc[start:end, ['Time', 'Channel A', 'Channel B', 'Channel
20        C', 'Channel D']].copy()
21
22        # Добавляем окно в массив капелек
23        if butter is True:
24
25            for column in window.columns:
26
27                if column != "Time":
28
29                    height_ = lambda a, b: a if a > b else b
30                    window[column] = apply_lowpass_filter(window[column],
31                                                            height_(abs(window[column].max()), 0.0000000000000001), 1, 2)
32
33            window['ID'] = i
34
35            setup[i] = window
36
37    return setup

```

# Извлечение признаков



```
1 features = {  
2     'ID': ID,  
3     f'Minimum_{col_name[-1]}': channel_values.min(),  
4     f'Maximum_{col_name[-1]}': channel_values.max(),  
5     f'Variance_{col_name[-1]}': feature_calculators.variance(channel_values),  
6     f'Kurtosis_{col_name[-1]}': feature_calculators.kurtosis(channel_values),  
7     f'Skewness_{col_name[-1]}': feature_calculators.skewness(channel_values),  
8     f'Median_{col_name[-1]}': feature_calculators.median(channel_values),  
9     f'Std_{col_name[-1]}': np.std(channel_values),  
10    f'Semimax_{col_name[-1]}': geometricals.semi_max(channel_values),  
11    f'K_{col_name[-1]}': geometricals.k_semi_max(channel_values,  
12    time_values),  
13    f'Q10_{col_name[-1]}': geometricals.q10(channel_values),  
14    f'Subpeak_{col_name[-1]}': geometricals.subpeak(channel_values,  
15    time_values),  
16    f'Semiwidth_{col_name[-1]}': geometricals.semi_width(channel_values,  
17    time_values)[0] if isinstance(  
18    geometricals.semi_width(channel_values, time_values), tuple) else 0  
19 }
```

# Извлечение признаков

- **Minimum**
- **Maximum:** Максимальное значение во временном ряду.
- **Variance:** Дисперсия временного ряда
- **Kurtosis:** Коэффициент эксцесса, измеряющий остроту (высокий или низкий) пика распределения.
- **Skewness:** Коэффициент асимметрии, измеряющий симметрию или асимметрию распределения.
- **Median:** Медиана временного ряда.
- **Std:** Стандартное отклонение временного ряда.
- **Semimax:** Полусумма максимальных значений (геометрический признак).
- **K:** Коэффициент K, связанный с положением и формой временного ряда (геометрический признак).
- **Q10:** Квантиль 10% временного ряда.
- **Subpeak:** Сумма амплитуд подпиков (геометрический признак).
- **Semiwidth:** Полуширина временного ряда (геометрический признак)

# Настройка модели catboost

```
1 def model_rain(features_df):
2
3     columns = ['dtAB', 'dtCD', 'dtAC', 'dtBD']
4
5     #Цели таргета
6     y_train_multi = features_df[columns]
7
8     #Исключение "лишних" обучающих данных
9     X_train = features_df[[col for col in features_df.columns if col not in columns]]
10
11     #Явное указание пула таргета и обучающих данных
12     train_pool = Pool(data=X_train, label=y_train_multi)
13
14     #Мультитаргетная модель
15     model_multi = CatBoostRegressor(iterations=1000, depth=6, learning_rate=0.5,
16     loss_function='MultiRMSE',
17                                     custom_metric=['MultiRMSE'])
18     model_multi.fit(train_pool)
19     predictions_multi = model_multi.predict(X_train)
20
21     return predictions_multi
```



# Руководство пользователя

1. Переместить папки с файлами в директорию с программой, или явно указать путь к файлам в коде программы. (по умолчанию – путь к папке с программой)
  2. Ввести параметры «Размер окна для выреза» (по умолчанию – 50000), «Высота пика» (если не указано – выбирается автоматически с помощью квантиля)
  3. Указать требуется ли создание графиков капель (по умолчанию – создается)
  4. Указать параметры Low pass фильтра (если не указано – будет взято по умолчанию)
  5. \*Указать параметры модели CatBoostRegressor
- \*Напрямую в коде