

# Лабораторная работа 1 “Предобработка данных”

12 сентября 2025 г.

## Цель работы

Исследование влияния различных методов обработки пропущенных значений на качество линейной регрессионной модели и сравнение эффективности методов предобработки данных.

## Задачи

1. Реализовать функцию фильтрации данных методом Савицкого-Голея
2. Реализовать функцию нормализации данных методом IQR
3. Создать функцию для генерации пропущенных значений в датасете
4. Реализовать функцию восстановления пропусков двумя методами:
  - Заполнение средним значением по столбцу
  - Заполнение средним значением по скользящему окну
5. Провести сравнительный анализ качества моделей:
  - Модель, обученная на исходных данных без пропусков
  - Модель, обученная на данных с восстановлением средним
  - Модель, обученная на данных с восстановлением скользящим средним
6. Визуализировать результаты и сделать выводы

## Исходные данные

- Датасет: House Prices или аналогичный регрессионный датасет
- Целевая переменная: `target`
- Факторы: `feature1`, `feature2`, ..., `featureN`
- Процент пропусков: 20% от общего числа наблюдений

## Требования к отчету

- Описание датасетов и выбранных переменных
- Обоснование выбора методов предобработки
- Визуализация результатов на каждом этапе
- Сравнительный анализ метрик качества моделей
- Выводы и рекомендации по использованию методов

## Импорт библиотек

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error, r2_score
5 from sklearn.model_selection import train_test_split
6 from scipy.signal import savgol_filter
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9
10 # For displaying Russian characters in matplotlib
11 plt.rcParams['font.family'] = 'DejaVu Sans'
12 plt.rcParams['axes.unicode_minus'] = False
```

Листинг 1: Import necessary libraries (Python)

```
1 # install.packages(c("dplyr", "ggplot2", "signal", "zoo", "caret", "tidyr"))
2
3 library(dplyr)
4 library(ggplot2)
5 library(signal)
6 library(zoo)
7 library(caret)
8 library(tidyr)
9
10 # For displaying Russian characters
11 Sys.setlocale("LC_CTYPE", "russian")
```

Листинг 2: Import necessary packages (R)

## Функция фильтрации Савицкого-Голея

```
1 def savitzky_golay_filter(data, window_size=11, poly_order=3):
2     """
3     Apply Savitzky-Golay filter to data
4
5     Parameters:
6     data: array-like - input data
7     window_size: int - window size (must be odd)
8     poly_order: int - polynomial order
9
10    Returns:
11    filtered_data: array - filtered data
12    """
13    if len(data) < window_size:
14        window_size = len(data) if len(data) % 2 == 1 else len(data) - 1
15
16    return savgol_filter(data, window_size, poly_order)
```

Листинг 3: Savitzky-Golay filter implementation (Python)

```
1 savitzky_golay_filter <- function(data, window_size = 11, poly_order = 3) {
2     """
3     Apply Savitzky-Golay filter to data
4
5     Parameters:
6     data: vector - input data
7     window_size: integer - window size (must be odd)
8     poly_order: integer - polynomial order
9
10    Returns:
11    filtered_data: vector - filtered data
12    """
13    if (length(data) < window_size) {
14        window_size <- ifelse(length(data) %% 2 == 1, length(data), length(data) - 1)
15    }
16
17    sgolayfilt(data, p = poly_order, n = window_size)
```

18 }

Листинг 4: Savitzky-Golay filter implementation (R)

## Функция нормализации IQR

```
1 def iqr_normalization(data):
2     """
3     Normalize data using IQR method
4
5     Parameters:
6     data: array-like - input data
7
8     Returns:
9     normalized_data: array - normalized data
10    """
11    Q1 = np.percentile(data, 25)
12    Q3 = np.percentile(data, 75)
13    IQR = Q3 - Q1
14    median = np.median(data)
15
16    # Protection against division by zero
17    if IQR == 0:
18        return data - median
19    else:
20        return (data - median) / IQR
```

Листинг 5: IQR normalization implementation (Python)

```
1 iqr_normalization <- function(data) {
2     """
3     Normalize data using IQR method
4
5     Parameters:
6     data: vector - input data
7
8     Returns:
9     normalized_data: vector - normalized data
10    """
11    Q1 <- quantile(data, 0.25, na.rm = TRUE)
12    Q3 <- quantile(data, 0.75, na.rm = TRUE)
13    IQR <- Q3 - Q1
14    med <- median(data, na.rm = TRUE)
15
16    # Protection against division by zero
17    if (IQR == 0) {
18        return(data - med)
19    } else {
20        return((data - med) / IQR)
21    }
22 }
```

Листинг 6: IQR normalization implementation (R)

## Функция создания пропусков

```
1 def introduce_missing_values(data, missing_percentage=0.2, random_state=42):
2     """
3     Create missing values in data
4
5     Parameters:
6     data: DataFrame - input data
7     missing_percentage: float - percentage of rows with missing values (0-1)
8     random_state: int - for reproducibility
9
10    Returns:
11    data_with_missing: DataFrame - data with missing values
12    """
```

```

13 np.random.seed(random_state)
14 data_with_missing = data.copy()
15
16 # Select rows for missing values
17 n_rows = int(len(data) * missing_percentage)
18 rows_with_missing = np.random.choice(data.index, n_rows, replace=False)
19
20 for idx in rows_with_missing:
21     # Select random columns for missing values
22     n_missing = np.random.randint(1, len(factors)//2 + 1)
23     cols_to_miss = np.random.choice(factors, n_missing, replace=False)
24     data_with_missing.loc[idx, cols_to_miss] = np.nan
25
26 return data_with_missing

```

Листинг 7: Missing values generation (Python)

```

1 introduce_missing_values <- function(data, missing_percentage = 0.2, seed = 42) {
2   ""
3   Create missing values in data
4
5   Parameters:
6   data: data.frame - input data
7   missing_percentage: numeric - percentage of rows with missing values (0-1)
8   seed: integer - for reproducibility
9
10  Returns:
11  data_with_missing: data.frame - data with missing values
12  ""
13  set.seed(seed)
14  data_with_missing <- data
15
16  n_rows <- round(nrow(data) * missing_percentage)
17  rows_with_missing <- sample(1:nrow(data), n_rows)
18
19  for (i in rows_with_missing) {
20    n_missing <- sample(1:(length(factors)%/2 + 1), 1)
21    cols_to_miss <- sample(factors, n_missing)
22    data_with_missing[i, cols_to_miss] <- NA
23  }
24
25  return(data_with_missing)
26 }

```

Листинг 8: Missing values generation (R)

## Функция восстановления пропусков

```

1 def impute_missing_values(data, method='mean', window_size=5):
2   ""
3   Impute missing values
4
5   Parameters:
6   data: DataFrame - data with missing values
7   method: str - imputation method ('mean', 'moving_average')
8   window_size: int - window size for moving average
9
10  Returns:
11  imputed_data: DataFrame - data with imputed values
12  ""
13  data_imputed = data.copy()
14
15  for col in data.columns:
16    if data[col].isna().any():
17      if method == 'mean':
18        # Fill with column mean
19        data_imputed[col] = data[col].fillna(data[col].mean())
20      elif method == 'moving_average':
21        # Fill with moving average
22        data_imputed[col] = data[col].fillna(
23          data[col].rolling(window=window_size, min_periods=1).mean()

```

```

24         )
25
26     return data_imputed

```

Листинг 9: Missing values imputation (Python)

```

1 impute_missing_values <- function(data, method = "mean", window_size = 5) {
2     ""
3     Impute missing values
4
5     Parameters:
6     data: data.frame - data with missing values
7     method: character - imputation method ('mean', 'moving_average')
8     window_size: integer - window size for moving average
9
10    Returns:
11    imputed_data: data.frame - data with imputed values
12    ""
13    data_imputed <- data
14
15    for (col in names(data)) {
16        if (any(is.na(data[[col]]))) {
17            if (method == "mean") {
18                data_imputed[[col]] <- ifelse(is.na(data[[col]]),
19                                              mean(data[[col]], na.rm = TRUE),
20                                              data[[col]])
21            } else if (method == "moving_average") {
22                # Use zoo for moving average
23                ma <- rollapply(data[[col]], window_size, mean, na.rm = TRUE,
24                               fill = NA, align = "right", partial = TRUE)
25                data_imputed[[col]] <- ifelse(is.na(data[[col]]), ma, data[[col]])
26            }
27        }
28    }
29
30    return(data_imputed)
31 }

```

Листинг 10: Missing values imputation (R)

## Основной код выполнения

```

1 # Load and prepare data
2 # data = pd.read_csv('your_dataset.csv')
3 # factors = ['feature1', 'feature2', 'feature3', ...] # feature list
4 # target = 'target_variable' # target variable
5
6 # Assume data is already loaded and prepared
7 # data contains factors and target
8
9 print("Initial data size:", data.shape)
10 print("Number of factors:", len(factors))
11
12 # Apply Savitzky-Golay filter to factors
13 print("Applying Savitzky-Golay filter...")
14 for col in factors:
15     data[col + '_filtered'] = savitzky_golay_filter(data[col].values)
16
17 # Normalize factors using IQR method
18 print("Applying IQR normalization...")
19 for col in factors:
20     data[col + '_normalized'] = iqr_normalization(data[col].values)
21
22 # Create missing values in data
23 print("Creating missing values...")
24 data_missing = introduce_missing_values(data[factors + [target]], 0.2)
25 print("Number of missing values after creation:", data_missing.isna().sum().sum())
26
27 # Impute missing values using two methods
28 print("Imputing missing values with mean method...")
29 data_imputed_mean = impute_missing_values(data_missing, 'mean')

```

```

30
31 print("Imputing missing values with moving average method...")
32 data_imputed_ma = impute_missing_values(data_missing, 'moving_average', 5)
33
34 # Function for training and evaluating model
35 def train_and_evaluate(X, y, test_size=0.2, random_state=42):
36     X_train, X_test, y_train, y_test = train_test_split(
37         X, y, test_size=test_size, random_state=random_state
38     )
39
40     model = LinearRegression()
41     model.fit(X_train, y_train)
42
43     y_pred = model.predict(X_test)
44     mse = mean_squared_error(y_test, y_pred)
45     r2 = r2_score(y_test, y_pred)
46     mae = mean_absolute_error(y_test, y_pred)
47
48     return model, mse, r2, mae, model.coef_
49
50 # Train on original data
51 print("Training on original data...")
52 model_orig, mse_orig, r2_orig, mae_orig, coef_orig = train_and_evaluate(
53     data[factors], data[target]
54 )
55
56 # Train on data imputed with mean
57 print("Training on mean-imputed data...")
58 model_mean, mse_mean, r2_mean, mae_mean, coef_mean = train_and_evaluate(
59     data_imputed_mean[factors], data_imputed_mean[target]
60 )
61
62 # Train on data imputed with moving average
63 print("Training on moving average-imputed data...")
64 model_ma, mse_ma, r2_ma, mae_ma, coef_ma = train_and_evaluate(
65     data_imputed_ma[factors], data_imputed_ma[target]
66 )
67
68 # Compare results
69 results = pd.DataFrame({
70     'Model': ['Original', 'Mean', 'Moving Average'],
71     'MSE': [mse_orig, mse_mean, mse_ma],
72     'R^2': [r2_orig, r2_mean, r2_ma],
73     'MAE': [mae_orig, mae_mean, mae_ma]
74 })
75
76 print("Model comparison:")
77 print(results)
78
79 # Analyze coefficient differences
80 coef_comparison = pd.DataFrame({
81     'Factor': factors,
82     'Original': coef_orig,
83     'Difference_Mean': coef_mean - coef_orig,
84     'Difference_Moving_Average': coef_ma - coef_orig
85 })
86
87 print("\nModel coefficient differences:")
88 print(coef_comparison)

```

Листинг 11: Main analysis process (Python)

```

1 # Load and prepare data
2 # data <- read.csv('your_dataset.csv')
3 # factors <- c('feature1', 'feature2', 'feature3', ...) # feature list
4 # target <- 'target_variable' # target variable
5
6 # Assume data is already loaded and prepared
7 # data contains factors and target
8
9 cat("Initial data size:", dim(data), "\n")
10 cat("Number of factors:", length(factors), "\n")
11

```

```

12 # Apply Savitzky-Golay filter to factors
13 cat("Applying Savitzky-Golay filter...\n")
14 for (col in factors) {
15   data[[paste0(col, "_filtered")]] <- savitzky_golay_filter(data[[col]])
16 }
17
18 # Normalize factors using IQR method
19 cat("Applying IQR normalization...\n")
20 for (col in factors) {
21   data[[paste0(col, "_normalized")]] <- iqr_normalization(data[[col]])
22 }
23
24 # Create missing values in data
25 cat("Creating missing values...\n")
26 data_missing <- introduce_missing_values(data[c(factors, target)], 0.2)
27 cat("Number of missing values after creation:", sum(is.na(data_missing)), "\n")
28
29 # Impute missing values using two methods
30 cat("Imputing missing values with mean method...\n")
31 data_imputed_mean <- impute_missing_values(data_missing, "mean")
32
33 cat("Imputing missing values with moving average method...\n")
34 data_imputed_ma <- impute_missing_values(data_missing, "moving_average", 5)
35
36 # Function for training and evaluating model
37 train_and_evaluate <- function(X, y, test_size = 0.2, seed = 42) {
38   set.seed(seed)
39   train_index <- createDataPartition(y, p = 1 - test_size, list = FALSE)
40
41   X_train <- X[train_index, ]
42   X_test <- X[-train_index, ]
43   y_train <- y[train_index]
44   y_test <- y[-train_index]
45
46   model <- lm(y_train ~ ., data = data.frame(X_train))
47   y_pred <- predict(model, newdata = data.frame(X_test))
48
49   mse <- mean((y_test - y_pred)^2)
50   r2 <- cor(y_test, y_pred)^2
51   mae <- mean(abs(y_test - y_pred))
52
53   return(list(model = model, mse = mse, r2 = r2, mae = mae, coef = coef(model)))
54 }
55
56 # Train on original data
57 cat("Training on original data...\n")
58 result_orig <- train_and_evaluate(data[factors], data[[target]])
59
60 # Train on data imputed with mean
61 cat("Training on mean-imputed data...\n")
62 result_mean <- train_and_evaluate(data_imputed_mean[factors],
63                                   data_imputed_mean[[target]])
64
65 # Train on data imputed with moving average
66 cat("Training on moving average-imputed data...\n")
67 result_ma <- train_and_evaluate(data_imputed_ma[factors],
68                                 data_imputed_ma[[target]])
69
70 # Compare results
71 results <- data.frame(
72   Model = c("Original", "Mean", "Moving Average"),
73   MSE = c(result_orig$mse, result_mean$mse, result_ma$mse),
74   R2 = c(result_orig$r2, result_mean$r2, result_ma$r2),
75   MAE = c(result_orig$mae, result_mean$mae, result_ma$mae)
76 )
77
78 cat("Model comparison:\n")
79 print(results)
80
81 # Analyze coefficient differences
82 coef_comparison <- data.frame(
83   Factor = factors,
84   Original = result_orig$coef[-1], # exclude intercept

```

```

85 Difference_Mean = result_mean$coef[-1] - result_orig$coef[-1],
86 Difference_Moving_Average = result_ma$coef[-1] - result_orig$coef[-1]
87 )
88
89 cat("Model coefficient differences:\n")
90 print(coef_comparison)

```

Листинг 12: Main analysis process (R)

## Визуализация результатов

```

1 # Quality metrics visualization
2 fig, axes = plt.subplots(1, 3, figsize=(15, 5))
3
4 metrics = ['MSE', 'R ', 'MAE']
5 values = [results['MSE'].values, results['R '].values, results['MAE'].values]
6 colors = ['lightcoral', 'lightgreen', 'lightblue']
7 titles = ['MSE Model Comparison', 'R Model Comparison', 'MAE Model Comparison']
8 ylabels = ['Mean Squared Error', 'R-squared', 'Mean Absolute Error']
9
10 for i, (ax, metric, value, color, title, ylabel) in enumerate(
11     zip(axes, metrics, values, colors, titles, ylabels)):
12
13     bars = ax.bar(results['Model'], value, color=color, alpha=0.7)
14     ax.set_title(title, fontsize=12)
15     ax.set_ylabel(ylabel, fontsize=10)
16     ax.tick_params(axis='x', rotation=45)
17
18     # Add values on bars
19     for bar in bars:
20         height = bar.get_height()
21         ax.text(bar.get_x() + bar.get_width()/2., height + 0.01 * max(value),
22             f'{height:.3f}', ha='center', va='bottom', fontsize=9)
23
24 plt.tight_layout()
25 plt.show()
26
27 # Coefficient differences visualization
28 plt.figure(figsize=(12, 6))
29 x_pos = np.arange(len(factors))
30 width = 0.35
31
32 plt.bar(x_pos - width/2, coef_comparison['Difference_Mean'],
33     width, label='Difference (Mean)', alpha=0.7, color='orange')
34 plt.bar(x_pos + width/2, coef_comparison['Difference_Moving_Average'],
35     width, label='Difference (Moving Average)', alpha=0.7, color='purple')
36
37 plt.xlabel('Factors', fontsize=12)
38 plt.ylabel('Coefficient Differences', fontsize=12)
39 plt.title('Model Coefficient Differences Compared to Original', fontsize=14)
40 plt.xticks(x_pos, factors, rotation=45, ha='right')
41 plt.legend()
42 plt.grid(axis='y', alpha=0.3)
43 plt.tight_layout()
44 plt.show()
45
46 # Save results
47 results.to_csv('results_comparison.csv', index=False, encoding='utf-8')
48 coef_comparison.to_csv('coefficients_comparison.csv', index=False, encoding='utf-8')
49 print("Results saved to CSV files")

```

Листинг 13: Model comparison visualization (Python)

```

1 # Quality metrics visualization
2 library(ggplot2)
3 library(gridExtra)
4
5 # Prepare data for visualization
6 results_long <- results %>%
7     pivot_longer(cols = c(MSE, R2, MAE),
8         names_to = "Metric", values_to = "Value")

```



```

9
10 # Plots for each metric
11 p1 <- ggplot(results, aes(x = Model, y = MSE, fill = Model)) +
12   geom_bar(stat = "identity", alpha = 0.7) +
13   geom_text(aes(label = round(MSE, 3)), vjust = -0.5) +
14   labs(title = "MSE Model Comparison", y = "Mean Squared Error") +
15   theme_minimal() +
16   theme(axis.text.x = element_text(angle = 45, hjust = 1))
17
18 p2 <- ggplot(results, aes(x = Model, y = R2, fill = Model)) +
19   geom_bar(stat = "identity", alpha = 0.7) +
20   geom_text(aes(label = round(R2, 3)), vjust = -0.5) +
21   labs(title = "R Model Comparison", y = "R-squared") +
22   theme_minimal() +
23   theme(axis.text.x = element_text(angle = 45, hjust = 1))
24
25 p3 <- ggplot(results, aes(x = Model, y = MAE, fill = Model)) +
26   geom_bar(stat = "identity", alpha = 0.7) +
27   geom_text(aes(label = round(MAE, 3)), vjust = -0.5) +
28   labs(title = "MAE Model Comparison", y = "Mean Absolute Error") +
29   theme_minimal() +
30   theme(axis.text.x = element_text(angle = 45, hjust = 1))
31
32 # Combined visualization
33 grid.arrange(p1, p2, p3, ncol = 3)
34
35 # Coefficient differences visualization
36 coef_long <- coef_comparison %>%
37   pivot_longer(cols = c(Difference_Mean, Difference_Moving_Average),
38     names_to = "Method", values_to = "Difference")
39
40 ggplot(coef_long, aes(x = Factor, y = Difference, fill = Method)) +
41   geom_bar(stat = "identity", position = "dodge", alpha = 0.7) +
42   labs(title = "Model Coefficient Differences Compared to Original",
43     x = "Factors", y = "Coefficient Differences") +
44   theme_minimal() +
45   theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
46   scale_fill_manual(values = c("orange", "purple"))
47
48 # Save results
49 write.csv(results, "results_comparison.csv", row.names = FALSE, fileEncoding = "UTF-8")
50 write.csv(coef_comparison, "coefficients_comparison.csv", row.names = FALSE,
51   fileEncoding = "UTF-8")
52 cat("Results saved to CSV files\n")

```

Листинг 14: Model comparison visualization (R)