

## Практика 13. Регрессионные модели характеристик насосных агрегатов

(начало совпадает с практикой 2, новое начиная с части 3)

### Часть 1. Подготовка данных

1. Скачать с сайта курса данные по характеристикам насосных агрегатов (НА).

Данные представляют собой реальные обезличенные данные по напорным характеристикам (QH) и характеристикам мощности ( $\eta$ ) насосных агрегатов нефтеперекачивающих станций. В одном файле содержатся данные по одному НА. В имени файла указывается номер нефтеперекачивающей станции (НПС) и номер агрегата. Каждый файл содержит 3 колонки Q, H,  $\eta$ .

**Если данные по всем вариантам не выложены, нужно взять пример данных (он точно выложен) и начать отлаживать код на нем.**

2. Для зачитки данных по агрегатам воспользуйтесь функцией (есть в архиве с данными):

```
function [Q, H, N] = read_pump_data(station_number, pump_number)
```

функция на вход принимает номер НПС и номер НА (возьмите их своего варианта), на выходе выдает Q, H,  $\eta$ .

**Если данные не выложены, временно берите `station_number = 1` и `pump_number = 1`. По готовности данных, переделать под свой вариант.**

### Часть 2. QH-характеристика и характеристика КПД

1. Построить регрессионную модель с полиномом 0-го, 1-го, 2-го, 3-го порядков.
2. Построить модели на основе аппроксимации из литературы по гидравлике:

- a. Для QH-характеристики

$$H(Q) = a - bQ^2$$

- b. Для характеристики мощности

$$\eta(Q) = k_1Q - k_2Q^2$$

**Для оценок коэффициентов моделей использовать как формулу  $\hat{\beta} = (X^T X)^{-1} X^T Y$ , так и функцию `regress`**

3. Построить график  $H(Q)$ , на который вывести
  - a. исходные данные (вывести отдельными точками, например, кружками)
  - b. предсказания всех моделей из части 2 в диапазоне Q от нуля до  $1.5 * \max(Q)$ , (график вывести линиями)
  - c. легенду для графиков всех моделей (чтобы не путать модели друг с другом)

Отмасштабировать график по H в диапазоне от нуля до  $\max(H)$ , используя функцию `ylim`.

4. Построить график  $\eta(Q)$ , на который вывести
  - a. исходные данные (вывести отдельными точками, например, кружками)
  - b. предсказания всех моделей из части 2 в диапазоне Q от нуля до  $1.5 * \max(Q)$ , (график вывести линиями)
  - c. задать диапазон по кпд. от нуля
  - d. вывести легенду для графиков всех моделей (чтобы не путать модели друг с другом)

Отмасштабировать график по  $\eta$  в диапазоне от нуля до  $\max(\eta)$ , используя функцию `ylim`.

5. Для каждой модели рассчитать вектор регрессионных остатков  $e = Y - X\hat{\beta}$ 
  - a. построить график регрессионных остатков как функции от расхода Q (вывести отдельными точками, например, кружками)
  - b. построить график регрессионных остатков по номерам экспериментов
  - c. вывести гистограмму регрессионных остатков

Провести визуальный анализ графика регрессионных остатков и их гистограмм. Найти подозрительные на ваш взгляд особенности. По возможности обосновать свои подозрения.

### Часть 3. Метрики качества модели

1. Реализовать функцию расчета коэффициента детерминации и сопутствующие функции:

```
function S2 = calc_residual_dispersion(y, y_prediction, k)
```

```
function TSS = calc_TSS(y)
```

```
function ESS = calc_ESS(y, y_prediction)
```

```
function RSS = calc_RSS(y, y_prediction)
```

```
function R2 = calc_R2(y, y_prediction)
```

y – вектор фактических значений

y\_prediction – вектор откликов тестируемой модели

k – порядок тестируемой модели

2. Рассчитать и свести в таблицу следующие показатели по всем моделям
  - a. Несмещенную оценку дисперсии шума
  - b. Коэффициент детерминации  $R^2$  с помощью calc\_R2 и с помощью функции regress. Убедиться, что значения совпадают (при несовпадении найти ошибку в Матлабе или у себя в коде).
3. Провести анализ таблицы, выбрать лучшую модель.

### Часть 4. Анализ регрессионных остатков

1. Реализовать функцию расчета статистики и функцию проверки гипотезы критерия Вальда-Вольфовица

```
function Ns = get_series_count(e)
```

```
function [Nplus, Nminus] = get_signed_values_count(e)
```

```
function [random_residuals, z, Ns] = test_residuals_randomness(e)
```

e – вектор остатков, отсортированных по величине фактора

Ns – количество серий

Nplus, Nminus – соответственно, количество положительных и отрицательных элементов в выборке

random\_residuals – равно единице, если остатки случайны, нулю если нет

z – стандартизированное значение статистики критерия

**Примечание.** При реализации корректно считайте пороги с учетом того, должен ли быть критерий односторонний или двусторонний.

2. Убедиться, что z-статистика совпадает со значением, возвращаемой функцией runstest.
3. Построить графики регрессионных остатков с сортировкой по величине возрастания фактора (расхода), по оси абсцисс – номер эксперимента в отсортированной последовательности. (Можно исправить вывод графика в части 2).
4. Рассчитать и свести в таблицу результат анализа регрессионных остатков по всем моделям (QH и КПД).
5. Выявить наиболее простую модель (QH и КПД), которая прошла тест на случайность остатков.

### Часть 5. Проверка значимости модели

1. Реализовать расчет статистики Фишера, используя реализованный ранее расчет ESS и RSS в виде отдельной функции. Реализовать функцию проверки гипотезы о значимости модели.

```
function [F, n1, n2] = model_Fstat(y, y_prediction, k)
```

```
function significant_model = test_model_significance(F, n1, n2, alpha)
```

y, y\_prediction, k – аналогичны с расчетом  $R^2$

F – значение формируемой F-статистики при проверке значимости модели, n1, n2 – степени свободы распределения данной статистики

significant\_model – равна нулю, если модель незначима, единице если значима

alpha – уровень значимости

2. Убедиться, что расчетное значение F-статистики совпадает с расчетом из функции regress.
3. Рассчитать и свести в таблицу (см. ниже) результат проверки значимости по всем моделям (QH и КПД).

### Часть 6. Проверка значимости оценок коэффициентов модели

1. Реализовать функцию расчета студентизированных оценок коэффициентов регрессии. Реализовать функцию проверки значимости коэффициентов регрессии.

```
function [beta_st, beta_bounds] = ...
    studentize_beta_est(beta_est, K_est, n, alpha)
function significant_coeff = ...
    test_coefficient_significance(beta_st, n, alpha)
```

beta\_est – оценки коэффициентов

beta\_st – студентизированные оценки коэффициентов

beta\_bounds – доверительный интервал оценок коэффициентов регрессии

$$\hat{\beta}_i - t_{\pi}(\alpha) < \beta_i < \hat{\beta}_i + t_{\pi}(\alpha)$$

K\_est – оценка ковариационной матрицы оценок коэффициентов

n – объем выборки, по которой строились оценки коэффициентов

significance\_coeff – вектор из k булевых значений, которые равны единице, если коэффициент значим, иначе равны нулю

alpha – уровень значимости при расчете доверительного интервала и порогов критерия

2. Убедиться, что расчет доверительных интервалов совпадает с расчетом функции regress (выходной аргумент bint [b, bint] = regress(y, X, alpha)).

3. Реализовать альтернативный способ проверки значимости коэффициентов: если ноль лежит внутри доверительного интервал оценки коэффициента, то коэффициент незначим.

```
function significant_coeff = ...
    test_coefficient_significance_alt(beta_bounds)
```

4. Рассчитать и свести в таблицу (см. ниже) результат проверки значимости всех коэффициентов по всем моделям (QH и КПД).

### Результаты практики в табличном виде

Запись результатов расчета выполнить в формате Microsoft Excel с помощью функции xlswrite.

Имя файла указывать в виде models\_ <фамилия>.xls, чтобы не было путаницы.

Коэффициенты модели	$R^2$	Анализ остатков	Оценка дисперсии шума	Значимость модели	Значимость коэффициентов (от $\hat{\beta}_0$ до $\hat{\beta}_{k-1}$ )
---------------------	-------	-----------------	-----------------------	-------------------	---

### Вопросы к защите (будут дополняться)

1. Коэффициент детерминации  $R^2$ 
  - a. Понятие тривиальной (опорной) модели. Эквивалентность МНК тривиальной модели и прогноза по средневыворочному.
  - b. Формулы и содержательный смысл  $RSS, TSS, ESS$
  - c. Доказательство равенства  $TSS = RSS + ESS$
  - d. Формула коэффициента детерминации  $R^2$ . Чему соответствуют значения  $R^2$ , равные 0%, 50%, 100%.
  - e. Как на основе  $R^2$  выбрать лучшую модель и нескольких моделей-кандидатов?
  - f. Какие имеются негативные следствия из того, что  $R^2$  – случайная величина?
2. Статистический критерий серий (Вальда-Вольфовица).
  - a. Случайный ряд, ряд с трендом, ряд с периодичностью. Исходя из определения серии, указать, для какого типа ряда их количество  $N_S$  больше или меньше.
  - b. Формулировки основной и альтернативной гипотезы.

- c. Эскизы распределений статистики  $N_S$  в условиях основной и альтернативной гипотез.
  - d. Формула стандартизированной статистики  $z$
  - e. Распределения статистики  $N_S$  и стандартизированной статистики  $z$  в условиях основной гипотезы (точная форма распределения).
3. Анализ регрессионных остатков
- a. Типичный вид графика регрессионных остатков при заниженном порядке аппроксимации.
  - b. Статистические свойства регрессионных остатков.
  - c. Какие из статистических свойств нарушаются при занижении порядка аппроксимации?
  - d. Применение критерия серия к анализу коррелированности регрессионных остатков. Какой тип ряда с точки зрения критерия серий формируется при занижении порядка аппроксимации?
  - e. Какую модель следует выбрать, исходя из анализа регрессионных остатков.
4. Критерий Фишера для проверки значимости модели.
- a. Распределение Фишера, связь с распределением хи-квадрат.
  - b. Корректно ли применять критерий Фишера для проверки значимости модели, у которой регрессионные остатки не коррелированы?
  - c. Эскиз графиков распределений статистик.
5. Критерий Стьюдента для проверки значимости коэффициентов модели.
- a. Статистика критерия – студентизированные оценки коэффициенты.
6. Что можно сказать о возможности использования построенных моделей в работе?
- a. Будут ли они отражать физику процесса на новых режимах работы?
  - b. Какова физика процесса?
  - c. Чтобы ответить на вопросы, надо указать, какие возможны новые режимы (а какие нет)?
  - d. Чтобы рассуждать о режимах, надо вспомнить, где установлен или может быть установлен насосный агрегат. Это не обязательно магистральный трубопровод.
7. Модель какого порядка лучше выбрать и почему?
8. Есть ли выбросы в ваших данных?
9. Ваши замечания по полученным регрессионным остаткам. Наблюдается ли тренд или колебательность?