

Статистические методы в QSAR

Часть 1

Множественная линейная
регрессия

Затронутые темы

- Задача статистического анализа в QSAR
- Множественная линейная регрессия
- Статистические показатели: R , s , F
- Пошаговый вариант множественной линейной регрессии
- Разбиение выборки на обучающую и контрольную
- Процедура скользящего контроля

Задача статистического анализа в QSAR

Целью статистического анализа в QSAR является поиск функции f , связывающей значение свойства y (которое может быть как физико-химическим свойством, так и биологической активностью) со значениями дескрипторов x_1, \dots, x_M , описывающих химическое соединение:

$$y \propto f(x_1, \dots, x_M)$$

Y непрерывное – регрессионный анализ, аппроксимация функции

Y дискретное – дискриминантный анализ, классификация, распознавание образов

Задача статистического анализа в QSAR

$$y^j = f(x_1^j, \dots, x_M^j) + e^j$$
$$y^j = F(c_1, \dots, c_P; x_1^j, \dots, x_M^j) + e^j \quad j=1, \dots, N$$

Принцип максимального правдоподобия

$$L(c_1, \dots, c_P) \rightarrow \max$$

Метод наименьших квадратов

- Выборка является репрезентативной
- Случайная величина ε имеет нормальное распределение
- Наблюдения являются независимыми
- Наблюдения являются равноточными

$$\sum_{j=1}^N (e^j)^2 \rightarrow \min$$

Множественная линейная регрессия – постановка задачи

$$y \propto c_0 + c_1 x_1 + \dots + c_M x_M$$

$$y^j = c_0 + \sum_{i=1}^M c_i x_i^j + e^j \quad j = 1, \dots, N$$

Найти такие значения c_i : $\sum_{j=1}^N (e^j)^2 \rightarrow \min$

Множественная линейная регрессия – решение задачи

$$C = (X^T X)^{-1} X^T Y$$

$$C = \begin{pmatrix} c_0 \\ c_1 \\ \boxtimes \\ c_M \end{pmatrix}$$

$$Y = \begin{pmatrix} y^1 \\ y^2 \\ \boxtimes \\ y^N \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & x_1^1 & \boxtimes & x_M^1 \\ 1 & x_1^2 & \boxtimes & x_M^2 \\ \boxtimes & \boxtimes & \boxtimes & \boxtimes \\ 1 & x_1^N & \boxtimes & x_M^N \end{pmatrix}$$

Регрессионные
коэффициенты

Экспериментальные
значения свойства

Значения дескрипторов

Статистические показатели для МЛР

RSS – сумма квадратов остатков

$$RSS = \sum_{j=1}^N (e^j)^2$$

SS – дисперсия свойства Y

$$SS = \sum_{j=1}^N (y^j - \bar{y})^2 \quad \bar{y} = \frac{1}{N} \sum_{j=1}^N y^j$$

R – коэффициент корреляции

$$R = \sqrt{\frac{SS - RSS}{SS}} \quad 0 < R < 1$$

R_{adj} – скорректированный коэффициент корреляции

$$R_{adj} = \sqrt{R^2 - (1 - R^2) \frac{M}{N - M - 1}}$$

Статистические показатели для МЛР

RMSE – среднеквадратичное значение ошибки $RMSE = \sqrt{\frac{RSS}{N}}$

s – стандартное отклонение $s = \sqrt{\frac{RSS}{N - M - 1}}$

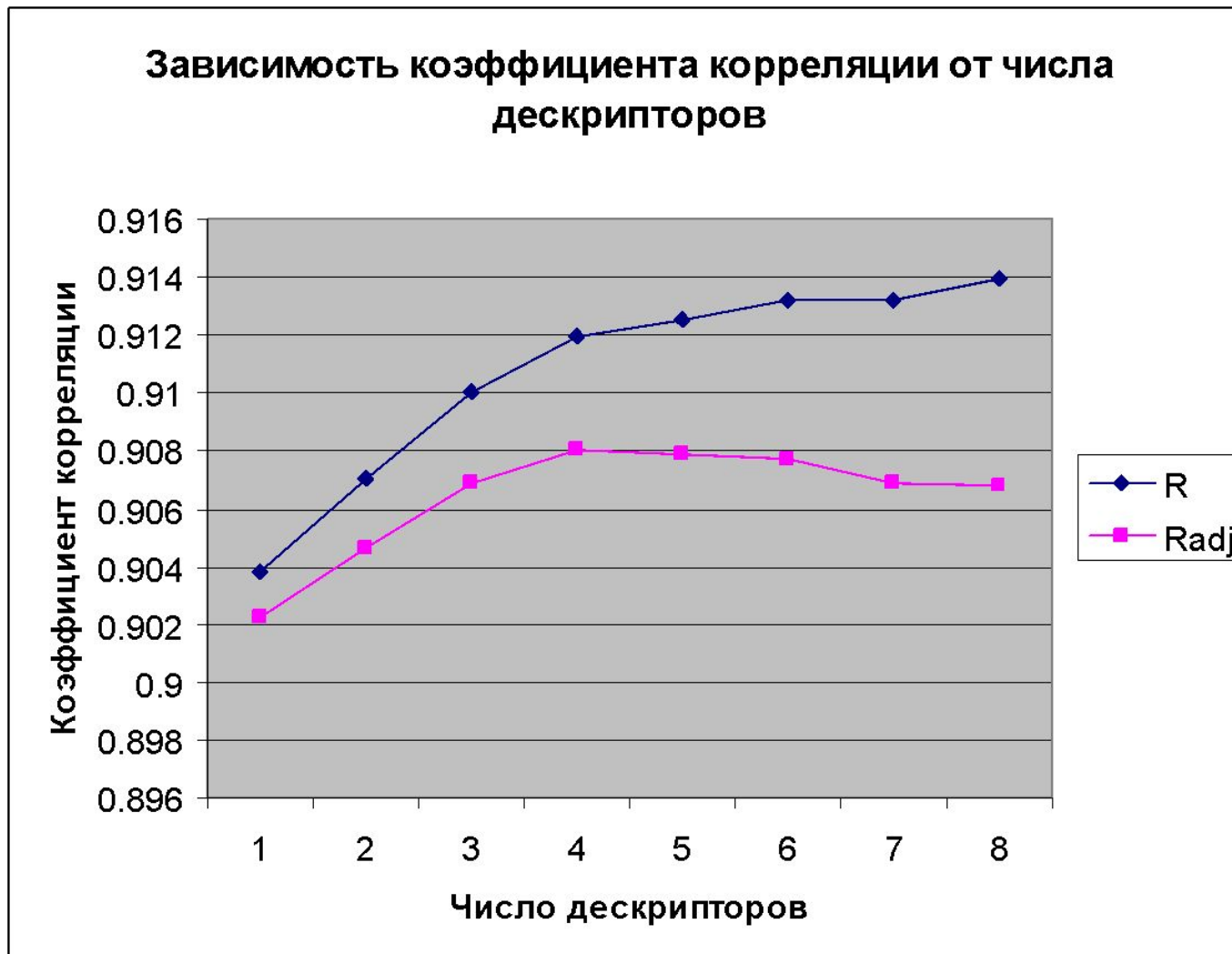
MAE – средняя абсолютная ошибка $MAE = \frac{1}{N} \sum_{j=1}^N |e^j|$

F – критерий Фишера $F = \frac{(SS - RSS)/(M + 1)}{RSS/(N - M - 1)}$

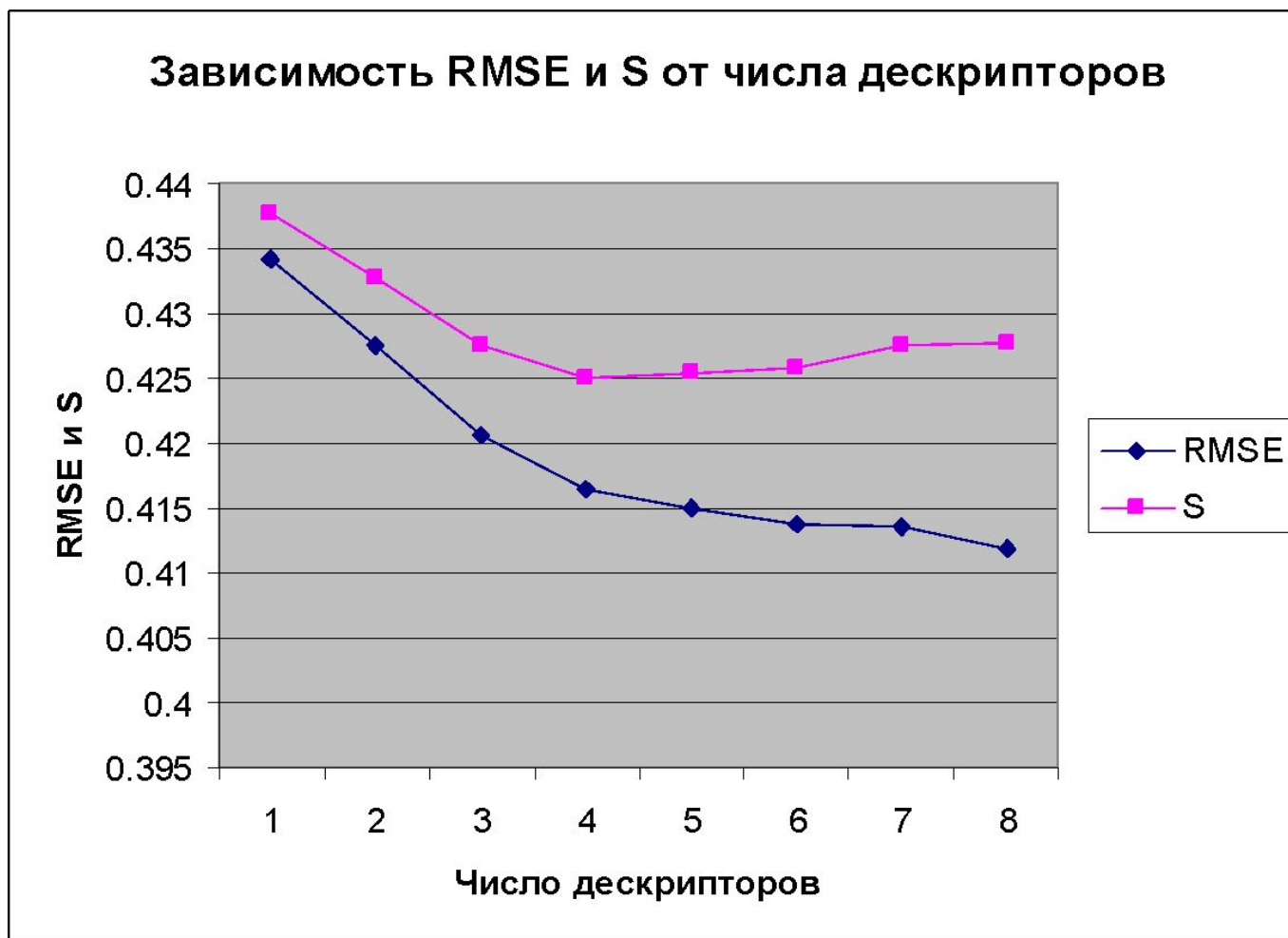
Статистические показатели для МЛР

Показатели описательной способности линейной регрессионной модели	Показатели прогнозирующей способности линейной регрессионной модели
R	R_{adj}
$RMSE$	s
MAE	
	F

Статистические показатели для МЛР



Статистические показатели для МЛР



Статистические показатели для МЛР

Индивидуальный t-критерий (критерий Стьюдента) для дескриптора

$$t_i = \frac{c_i}{s \cdot \sqrt{a_{ii}}}$$

$$\|a_{ij}\| = (X^T X)^{-1}$$

Необходимость отбора дескрипторов

- Проблема мультиколлинеарности дескрипторов и сингулярности матрицы $(X^T X)^{-1}$
- Проблема переопределенности моделей
- Внесение «шума» в модель нерелевантными дескрипторами

Пошаговый вариант множественной линейной регрессии

$$F_{ie} = \frac{(R_2^2 - R_1^2) \cdot (N - M_2 - 1)}{1 - R_2^2}$$

Дескриптор включается в модель, если $F_{ie} = \max \wedge F_{ie} > F_{IN}$

Дескриптор исключается из модели, если $F_{ie} = \min \wedge F_{ie} < F_{OUT}$

Типичные значения порогов: $F_{IN} = 3.84$, $F_{OUT} = 2.7$

Разбиение выборки на обучающую и контрольную

$PRSS_s$ - сумма квадратов остатков при прогнозе

$$PRSS_s = \sum_{j \in S} (e^j)^2$$

PSS_s - дисперсия свойства y на контрольной выборке

$$PSS_s = \sum_{j \in S} (y^j - \bar{y})^2$$

$PRMSE_s$ - среднеквадратичная ошибка на прогнозе

$$PRMSE_s = \sqrt{\frac{PRSS_s}{N_s}}$$

$PMAE_s$ - средняя абсолютная ошибка на прогнозе

$$PMAE_s = \frac{1}{N_s} \sum_{j \in S} |e^j|,$$

Q_s^2 - квадрат коэффициента корреляции на прогнозе

$$Q_s^2 = \frac{PSS_s - PRSS_s}{PSS_s}$$

Процедура скользящего контроля (cross-validation)

1. При μ -кратном скользящем контроле исходная выборка разбивается на μ приблизительно равных частей
2. Каждая из этих частей по очереди объявляется контрольной выборкой
3. Для нее формируется обучающая выборка, состоящая из всех соединений из исходной выборки, в нее не входящих
4. По обучающей выборке строится регрессионная модель
5. По текущей контрольной выборке вычисляется сумма квадратов ошибок PRSSs и сумма абсолютных ошибок PSAEs
6. Пункты 2-5 повторяются для всех μ частей

Процедура скользящего контроля – статистические показатели

$$PRSS_{CV} = \sum_{k=1}^{\mu} PRSS_{S_k}$$

$$PSAE_{CV} = \sum_{k=1}^{\mu} PSAE_{S_k}$$

$RMSE_{CV}$ - среднеквадратичная ошибка прогноза

$$RMSE_{CV} = \sqrt{\frac{PRSS_{CV}}{N}}$$

MAE_{CV} - средняя абсолютная ошибка прогноза

$$MAE_{CV} = \frac{PSAE_{CV}}{N}$$

Q^2_{CV} - коэффициент корреляции для прогноза

$$Q^2_{CV} = \frac{SS - PRSS_{CV}}{SS}$$