

МСПС | Регрессионный анализ

(1)

Постановка задачи

Есть объекты $1 \dots n$, есть признаки Y, X_1, \dots, X_k , меняющиеся на объектах. X_1, \dots, X_k - объясняющие переменные (факторы, признаки)
 Y - отклик. Задача: объяснить Y через X_1, \dots, X_k , т.е. найти такую ф-цию f , что $Y \approx f(X_1, \dots, X_k)$.

Или т.о. наилучшим \hat{f} квадратич. прогнозе, а именно $\min_f E(Y - f(X_1, \dots, X_k))^2 = E(Y | X_1, \dots, X_k)$

Если знаем f в явном виде, то просто модель регрессии (будем считать эту модель в рамках непараметрич. регрессии)

Если знаем f как $\theta_0 + \sum_{i=1}^k \theta_i X_i$, то модель линейной регрессии (для гаусс. случая оптимальна, т.е. $E(Y | X_1, \dots, X_k)$ есть линеар. ф-ция от X_1, \dots, X_k)

Т.е. задача (линейной регрессии) найти по Y и X оценку неизвест. вектора пар-ров θ в такой модели: $Y = X\theta + \varepsilon$, где $X = (X_0, X_1, \dots, X_k)$ - матрица $n \times (k+1)$, $X_0 = (1, \dots, 1)^T$, ε - об. шум.

Знаем из метода наименьших квадратов, что (МНК-метод)
 $\hat{\theta} = \arg \min_{\theta} \|Y - X\theta\|^2 = (X^T X)^{-1} X^T Y$ и сама оц. для Y : $\hat{Y} = X \hat{\theta}$.

Предположение Гаусса-Маркова

$E\varepsilon = 0$, $D\varepsilon = \sigma^2 I_n$, где $I_n = \text{diag}\{1, \dots, 1\}$ - единич. матрица / разм. $n \times n$,

σ^2 - неизвестно. В этом случае $\hat{\theta}$ - оптимальная (с наимен. дисперсией) оценка θ в классе линейных оценок вида BY , где B - матрица (т.е. Гаусса-Маркова)

Знаем, кроме того, что $\hat{\theta}$ - несмещ. и состоят. оц. θ .

Несмещ. и состоят. оц. для σ^2 имеет $\hat{\sigma}^2 = \frac{1}{n-k-1} \|Y - X\hat{\theta}\|^2$

В случае же, если $\varepsilon \sim N(0, \sigma^2 I_n)$, то $\hat{\theta}$ - оптимальная (в классе всех несмещенных оценок, а не только линейных)

Теперь, собственно, к вопросам модели:

- 1) Все ли признаки значимы? (т.е. проверить гипотезу о том, что $\theta_j = 0$)
- 2) Прискачать Y на новом объекте (т.е. дов. инт. построить)
- 3) Проверить, адекватна ли модель.

Итак, к значимости признаков. Проверим гипотезу $H_0: \theta_j = 0$.

Мы знаем, что $\hat{\theta} \sim N(\theta, \sigma^2 (X^T X)^{-1})$, откуда вытекает

$\frac{c^T (\hat{\theta} - \theta)}{\hat{\sigma} \sqrt{c^T (X^T X)^{-1} c}} \sim t_{n-k-1}$, где c - вектор из \mathbb{R}^k . Если хотим проверить гипотезу $H_0: \theta_j = 0$, берём $c_j = (0, \dots, 1, \dots, 0)$ и

получаем $\frac{\hat{\theta}_j - \theta_j}{\hat{\sigma} \sqrt{a_{jj}}} \sim t_{n-k-1}$, где $a_{jj} = ((X^T X)^{-1})_{jj}$, откуда получаем

доверит. интервал для θ_j , откуда и будет следовать критерий.

ДЦ: $(\hat{\theta}_j - u_{1-\frac{\alpha}{2}} \hat{\sigma} \sqrt{a_{jj}}, \hat{\theta}_j + u_{1-\frac{\alpha}{2}} \hat{\sigma} \sqrt{a_{jj}})$. Критерий: если 0 не попал в этот ДЦ, то отвергаем H_0 (в пользу значимости признака j)

(t -критерий Стьюдента)

МОНС | Рефлексивный анализ

Теперь проверим гипотезу о том, что несколько признаков не являются значимыми. Для этого введём коэфф. детерминации R^2 :

$$R^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \hat{y}_i)^2} = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}, \text{ где } RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Логично, что чем меньше RSS , тем лучше мы приближились к истинным y_i , т.е. тем ближе R^2 к 1, тем лучше качество регрессии (позднее об этом ещё поговорим). Кроме того, $R^2 = (\hat{\rho}(y, \hat{y}))^2$, где $\hat{\rho}(y, \hat{y})$ - коэфф. корр. Пирсона между y и \hat{y} .

Критерий Фишера Проверим гипотезу $H_0: \beta_{k_1} = \dots = \beta_{k_q} = 0$.

Пусть RSS_{H_0} построена по модели, где y нас интересует, признаки k_1, \dots, k_q , т.е. $\hat{y}_{H_0} = X_{H_0} (X_{H_0}^T X_{H_0})^{-1} X_{H_0}^T y$, где X_{H_0} - матрица с отсутствующими столбцами k_1, k_2, \dots, k_q , т.е. $(x_1, x_2, \dots, x_{k_1-1}, x_{k_1+1}, \dots)$

Тогда если H_0 верна, то $\frac{(RSS - RSS_{H_0})/q}{RSS/(n-k-1)} \sim F_{q, n-k-1}$

Для проверки гипотезы $H_0: \beta_1 = \dots = \beta_k = 0$ используем такой крит. Фишера:

$$F = \frac{R^2/k}{(1-R^2)/(n-k-1)}, \quad F \sim F(k, n-k-1) \text{ при верной } H_0.$$

2) Для значения отклика на новых наблюд. $x_0: y(x_0) = x_0^T \theta + \epsilon(x_0)$

строится так: $(x_0^T \hat{\theta} - u_{1-\frac{\alpha}{2}} \hat{\sigma} \sqrt{x_0^T (X^T X)^{-1} x_0}, x_0^T \hat{\theta} + u_{1-\frac{\alpha}{2}} \hat{\sigma} \sqrt{x_0^T (X^T X)^{-1} x_0})$ (вместо σ на предыдущей странице стоит $\hat{\sigma}$ и учитываем, что дисп. $\epsilon(x_0)$ равна $\hat{\sigma}^2$)

3) Адекватность модели и model selection.

Видим (по графику, по критериям, по R^2), что выбранная модель плохо данные описывает. Возможные варианты действий:

- 1) Замена переменных. Т.е. применяем к y и x функции. Возможные варианты ф-ций: $\ln x, e^x, e^{-x}, \frac{1}{x}$ и т.д. Здесь важно проверить, что при переходе к новым переменным ϵ_j продолжают оставаться независимыми и, по возможности, были нормально распределены.
- 2) Добавление переменных вида x_i^k , т.е. будем просто $y = \sum \theta_i x_i$, а тут ещё степени x_i добавим. Может привести к переобучению, т.е. неустойчивости решения на новых объектах.
- 3) Убрать перемен. согласно критериям Стьюдента и Фишера.

Как выбрать, какая модель лучше?

а) Информационные критерии Акаике и Шварца (нормальность тут не требуется) Акаике: $AIC = 2k + n [\ln \frac{RSS}{n} + 1]$; Шварц: $BIC = 2k \ln n + n [\ln \frac{RSS}{n} + 1]$. Выбираем ту модель, у которой наименьший AIC или BIC .

Инф. критерии не являются стат. критериями проверки гипотез, к ним не применимо понятие "значимости критерия". Просто числовые характеристики.

б) Критерий RESET Рэйджи

В рамках лм. Gauss. модели проверяется предположение $E\epsilon_i = 0$ (мы же хотим, чтобы наша оценка ~~не~~ отклика была несмещенной!)
Можно выявить с помощью этого крит.

- 1) наличие производных, значимых для регрессии факторов
- 2) неправильную функц. форму регрессии
- 3) наличие корреляции между факторами (мультиколлинеарность)

Правда, заранее неизвестно, что именно выявили.

Как устроен критерий: с помощью МНК подгоняем \hat{y}_i - оценки отклика, потом раскл. вспомогат. модель $y_i = \sum_{j=1}^k \theta_j x_{ij} + f_1(\hat{y}_i)^2 + \eta_i, i=1..n$

и проверяется гипотеза $H_0: \theta_j = 0$. Для этого вычисляем МНК

новой модели, \hat{y}_i подгоняем, и пишем статистику:

при верной H_0
$$F = \frac{RSS_{(1)} - RSS_{(2)}}{RSS_{(2)} / (n - k - 1)} \xrightarrow{d} F_{1, n-k-1}$$

Замечание 1) Почему важно убрать незначимые признаки? Потому что при их убирании оценка $\hat{\theta}$ остается несмещенной, а дисперсия снижается.

в) Стандартный коэфф. детерминации R^2 всегда увелич. при добавлении признаков в модель, поэтому для отбора признаков не нужно использовать. Можно брать приведенный коэфф. детерминации:

$$R_a^2 = 1 - (1 - R^2) \frac{n-1}{n-k-1}$$

2) Пошаговая регрессия. Вообще по регрессионному анализу хороша книга Сэйдера "Линейный регрессионный анализ"

Шаг 0 ~~Регрессия~~ для опис. Уберем модель либо только с константой, либо с 1 переменной. Выбираем 2 α , порога для ~~удаления~~ R_{IN} и R_{OUT} (например, $R_{IN} = 0.05$ и $R_{OUT} = 0.1$). Переб. включ. в модель, если у нее наимен. ~~довер-~~ ~~значение~~ по крит. Фишера (см. стр. 2, $\frac{RSS - RSS_{H_0}}{q} - \text{это}$ $\frac{RSS - RSS_{H_0}}{RSS / (n - k - 1)}$)

Шаг 2N+1 Если есть переменные, при удалении которых ~~статистика~~ крит. Фишера дает $p\text{-знач.} > R_{OUT}$, то удаляем их.

Шаг 2N Если есть переим., при добавл. которых крит. Фишера дает $p\text{-знач.} < R_{IN}$, то добавляем их.

Множеств. проверка гипотез тут не устроены, потому что процедура пошаговая.

Что делать с категориал. переменными? Есть

- 1) Кодировать, ~~присвоив~~ присвоив каждому фактору номер. Давать номер номеру. Но мне сказал программист, что много места ест.
- 2) Заполнить частотами встречаемости категор. признака. Надо это делать аккуратно, чтобы не переобучиться (касается задач прогнозирования)
- 3) Хэширование. См. курс машин. обучения