

УДК 519.711.3

М.В. Потанина, доцент, канд. техн. наук

Севастопольский национальный технический университет

ул. Университетская, 33, г. Севастополь, Украина, 99053

E-mail: mem@sevgtu.sebastopol.ua

СВОЙСТВА СРЕДНЕКВАДРАТИЧЕСКОЙ ОШИБКИ ПРОГНОЗА РИДЖ-РЕГРЕССИИ ДЛЯ ИДЕНТИФИКАЦИИ МОДЕЛИ СЛОЖНОГО ОБЪЕКТА УПРАВЛЕНИЯ ПРИ МУЛЬТИКОЛЛИНЕАРНОСТИ ФАКТОРОВ

Произведено исследование функции среднеквадратической ошибки прогноза (СКОП) для ридж-регрессии на экстремум в зависимости от параметра регуляризации. Доказано наличие локального минимума СКОП, что позволяет решить задачу поиска оптимального параметра регуляризации для идентификации модели сложного объекта управления при мультиколлинеарности факторов.

Ключевые слова: *среднеквадратическая ошибка прогноза, ридж-регрессия, параметр регуляризации, идентификация, мультиколлинеарность.*

Постановка проблемы. Идентификация модели сложного объекта управления (СОУ) при мультиколлинеарности факторов является сложной задачей.

Одним из способов решения проблемы мультиколлинеарности является использование смещенных оценок, позволяющих устранить недостатки традиционного метода наименьших квадратов (МНК). К таким оценкам относятся «гребневые» или ридж-оценки. Но при этом необходимо учитывать, что в настоящее время не существует однозначно рекомендуемого способа выбора параметра регуляризации ридж-регрессии.

Особенностью выбора данного параметра является то, что и смещение и СКОП зависят от неизвестных истинных значений параметров модели.

Предлагается методика исследования СКОП, как критерия оптимальности качества идентификации модели СОУ в заданной области в случае ридж-регрессии, на экстремум в зависимости от параметра регуляризации, что доказало бы возможность автоматизированного вычисления оптимального параметра регуляризации.

Анализ последних исследований и публикаций. В случае обработки данных заранее не спланированного эксперимента при исследовании функционирования СОУ часто встречается явление, называемое мультиколлинеарностью [1,4]. Отрицательные последствия данного явления проявляются в следующем:

- неустойчивость оценок регрессионных коэффициентов. При добавлении или исключении совсем малого количества информации (например, одного наблюдения) может произойти очень сильное изменение оценок коэффициентов. При этом резко уменьшается и точность предсказания по модели;
- численная неустойчивость процедуры оценивания, вызванная ошибками машинного округления и накопления этих ошибок;
- коэффициенты регрессионной модели оказываются сильно коррелированными между собой, что лишает смысла их интерпретацию;
- резко увеличиваются дисперсии оценок коэффициентов.

Все это затрудняет построение прогнозирующих моделей СОУ и делает невозможным использование оценок типа МНК, так как точность их при мультиколлинеарности становится недостаточной [2].

Однако в этом случае можно повысить эффективность оценивания за счет привлечения априорной информации об оцениваемых параметрах уравнений регрессии вероятностного или детерминированного характера. Но это в свою очередь приводит к использованию смещенных оценок, позволяющих уменьшить СКОП по сравнению с МНК-оценками и тем самым получить значения параметров, наиболее близких к истинным.

Смещенные, например, ридж-оценки являются более устойчивыми, чем оценки МНК и имеют меньшее значение СКОП при коррелированности факторов, вызванной мультиколлинеарностью [1].

Главной трудностью при использовании ридж-оценок является выбор параметра регуляризации r .

Процедура регуляризации представляет собой попытку преодолеть последствия плохой обусловленности информационной матрицы, как следствия мультиколлинеарности. Эта идея была обоснована в целом ряде работ как зарубежных, так и отечественных авторов [1, 3].

К сожалению, во всех подходах, рассматриваемых в указанных работах, обычно требуется дополнительная информация, которая редко имеется в распоряжении исследователя и поэтому ряд

вопросов, связанных с идентификацией модели объекта управления в условиях мультиколлинеарности факторов, остаются нерешенными.

Цель статьи. В настоящей работе приводится методика исследования СКОП как критерия оптимальности в случае ридж-регрессии на экстремум в зависимости от параметра регуляризации. Доказано наличие локального минимума СКОП, что позволяет решить задачу автоматического вычисления оптимального параметра регуляризации при смещенном оценивании для идентификации модели СОУ при мультиколлинеарности факторов.

Изложение основного материала исследования с полным обоснованием полученных научных результатов. В работе рассматривается случай, когда на классе возможных моделей задана структура, позволяющая ввести частичный порядок. При таком порядке классы моделей как бы вложены один в другой: $S_1 \subset S_2 \subset \dots \subset S_q$, S_j – модель j класса, q – максимально возможный порядок модели. Так для линейных по параметрам моделей можно задать структуру в зависимости от количества членов модели. В этом случае каждый класс S_j задается следующей моделью

$$\eta_j(\bar{x}, \bar{z}, \bar{\alpha}_j) = f_j^T(\bar{x}, \bar{z}) \bar{\alpha}_j \quad (j=1,2,\dots,q), \quad (1)$$

где $f_j^T(\bar{x}) = [f_{j1}(\bar{x}, \bar{z}), f_{j2}(\bar{x}, \bar{z}), \dots, f_{jn_j}(\bar{x}, \bar{z})]$ – вектор известных функций от векторов \bar{x} и \bar{z} , n_j – число неизвестных параметров модели j , $\bar{\alpha}_j^T$ – вектор неизвестных параметров j -ой модели.

В качестве априорной информации в работе рассматриваются следующие данные: пределы среднего значения функции отклика (η_{\min} и η_{\max}), область возможных значений факторов, максимальная дисперсия случайной составляющей модели σ^2 и закон распределения этой составляющей.

СКОП для любой модели в области прогноза W_1 можно оценить функционалом

$$EI(\bar{\alpha}) = \int_{W_1} (y - \eta(\bar{x}, \bar{\alpha}))^2 P(\bar{x}, \bar{y}) d\bar{x} d\bar{y}, \quad (2)$$

где $P(\bar{x}, \bar{y})$ – совместная плотность распределения переменных \bar{x} и \bar{y} .

Пусть матрица $F_j^T = [f_j(\bar{x}_1), f_j(\bar{x}_2), \dots, f_j(\bar{x}_N)]$ – матрица функций от экспериментальных данных для модели j размером $N \times n_j$ и $Y^T = [y_1, y_2, \dots, y_N]$ – вектор N наблюдений над функцией отклика. Предполагается, что матрица F_j имеет полный ранг.

Если модель линейная по параметрам и используется ридж-регрессия, тогда оценка определится как

$$\bar{\alpha}_{j(RIDG)} = (F_j^T F_j + rI_j)^{-1} F_j^T Y, \quad (3)$$

где r – параметр регуляризации.

Выражение для СКОП модели j в случае линейной параметризации при использовании ридж-регрессии, имеет вид

$$\begin{aligned} EI(j)_{RIDG} = & \bar{\alpha}_j^T (I_j - F_j^T F_j (F_j^T F_j + rI_j)^{-1}) J_j (I_j - F_j^T F_j (F_j^T F_j + rI_j)^{-1}) \bar{\alpha}_j + \\ & + \bar{\alpha}_{q-j}^T J_{q-j} \bar{\alpha}_{q-j} + \bar{\alpha}_{q-j}^T F_{q-j}^T F_j (F_j^T F_j + rI_j)^{-1} J_j (F_j^T F_j + rI_j)^{-1} F_j^T F_{q-j} \bar{\alpha}_{q-j} - \\ & - 2\bar{\alpha}_j^T (I_j - F_j^T F_j (F_j^T F_j + rI_j)^{-1}) J_j (F_j^T F_j + rI_j)^{-1} F_j^T F_{q-j} \bar{\alpha}_{q-j} + \\ & + \sigma^2 sp(F_j B_j J_j B_j F_j^T). \end{aligned} \quad (4)$$

где $B_j = (F_j^T F_j + rI_j)^{-1}$ и $J = \int_{W_1} \bar{f}(x) \bar{f}(x)^T d\bar{x} / \int_{W_1} d\bar{x}$.

Для полной модели q , включающей истинную модель, СКОП в случае ридж-регрессии имеет вид

$$EI(q)_{RIDG} = \bar{\alpha}^T (I_q - F_q^T F_q B_q) J (I_q - B_q F_q^T F_q) \bar{\alpha} + \sigma^2 sp F_q B_q J B_q F_q^T, \quad (5)$$

где $B_q = (F_q^T F_q + rI_q)^{-1}$.

Исследуем СКОП для полной модели q на экстремум в случае ридж-регрессии. Предполагается, что если функция СКОП в случае ридж-регрессии имеет локальный минимум, то он и будет соответствовать оптимальному значению параметра регуляризации.

Продифференцируем выражение (5) по r . Здесь и в дальнейшем будем считать, что $B = B_q$. Имеем

$$\begin{aligned} \frac{dEI}{dr} = & \bar{\alpha}^T F^T FB^{-2} J \bar{\alpha} - \bar{\alpha}^T F^T FB^{-2} JB^{-1} F^T F \bar{\alpha} + \bar{\alpha}^T JB^{-2} F^T F \bar{\alpha} - \\ & - \bar{\alpha}^T F^T FB^{-1} JB^{-2} F^T F \bar{\alpha} - \sigma^2 sp(FB^{-2} JB^{-1} F^T + FB^{-1} JB^{-2} F^T). \end{aligned} \quad (6)$$

Экстремум функции EI в точке r_{OPT} будет при $\frac{dEI}{dr} = 0$.

Очевидно, что это выражение зависит от параметра регуляризации r , от конкретных собранных экспериментальных данных, от истинных коэффициентов и от дисперсии ошибки.

Найдем условие существования минимума для СКОП. Для этого необходимо найти вторую производную выражения СКОП $\frac{d^2 EI}{dr^2}$.

$$\begin{aligned} \frac{d^2 EI}{dr^2} = & -2\bar{\alpha}^T B^{-3} J F^T F \bar{\alpha} + 2\bar{\alpha}^T F^T FB^{-3} JB^{-1} F^T F \bar{\alpha} + \bar{\alpha}^T F^T FB^{-2} JB^{-2} F^T F \bar{\alpha} + \\ & + \sigma^2 sp(2FB^{-3} JB^{-1} F^T + FB^{-2} JB^{-2} F^T) > 0. \end{aligned} \quad (7)$$

Формула (7) будет являться условием существования минимума для СКОП.

Найдем значение оптимального параметра регуляризации r_{OPT} . Для этого раскроем скобки в выражении (7). Имеем $\frac{dEI}{dr} = 0$, тогда, так как в формуле (7) данные слагаемые попарно равны, то

$$\begin{aligned} \bar{\alpha}^T F^T FB^{-1} JB^{-2} F^T F \bar{\alpha} &= \bar{\alpha}^T F^T FB^{-2} JB^{-1} F^T F \bar{\alpha}; \\ \bar{\alpha}^T F^T F JB^{-2} \bar{\alpha} &= \bar{\alpha}^T B^{-2} J F^T F \bar{\alpha}; \\ \sigma^2 sp(FB^{-1} JB^{-2} F^T) &= \sigma^2 sp(FB^{-2} JB^{-1} F^T). \end{aligned}$$

Формулу (7) можно переписать в виде

$$\frac{dEI}{dr} = 2 * (\bar{\alpha}^T B^{-2} J F^T F \bar{\alpha} - \bar{\alpha}^T F^T FB^{-2} JB^{-1} F^T F \bar{\alpha} - \sigma^2 sp(FB^{-2} JB^{-1} F^T)) = 0.$$

Преобразуем это выражение, умножив на B^2 слева и разделив на 2:

$$\frac{dEI}{dr} = \bar{\alpha}^T J F^T F \bar{\alpha} - \bar{\alpha}^T F^T F JB^{-1} F^T F \bar{\alpha} - \sigma^2 sp(F JB^{-1} F^T) = 0.$$

Далее умножаем на B правую часть выражения.

$$\frac{dEI}{dr} = \bar{\alpha}^T F^T B J F \bar{\alpha} - \bar{\alpha}^T F^T F J F^T F \bar{\alpha} - \sigma^2 sp(F J F^T) = 0. \quad (8)$$

Подставим B в формулу (8), имеем

$$\begin{aligned} \bar{\alpha}^T J(F^T F + I_{n,r}) F^T F \bar{\alpha} - \bar{\alpha}^T F^T F J F^T F \bar{\alpha} - \sigma^2 sp(F J F^T) &= 0; \\ \bar{\alpha}^T J F^T F F^T F \bar{\alpha} + r \bar{\alpha}^T J F^T F \bar{\alpha} - \bar{\alpha}^T F^T F J F^T F \bar{\alpha} - \sigma^2 sp(F J F^T) &= 0; \\ r \bar{\alpha}^T J F^T F \bar{\alpha} - \sigma^2 sp(F J F^T) &= 0. \end{aligned}$$

Окончательно имеем

$$r_{OPT} = \sigma^2 sp(F J F^T) / \bar{\alpha}^T J F^T F \bar{\alpha}. \quad (9)$$

Исследование свойств СКОП для ридж-регрессии с помощью имитационных экспериментов.

В качестве тестовой рассматривается модель вида $y = \alpha_0 + a_1 x_1 + \alpha_2 x_2 + \epsilon$, где $\epsilon \approx N(0,1)$ для трех различных вариантов матриц экспериментальных данных X_1, X_2, X_3 , отличающихся степенью корреляции.

Матрицы экспериментальных данных имеют вид

$$X1 = \begin{bmatrix} -1 & -1 \\ -1 & +1 \\ +1 & -1 \\ +1 & +1 \end{bmatrix}; \quad X2 = \begin{bmatrix} +1 & +1 \\ -1 & +1 \\ -1 & +1 \\ -1 & -1 \end{bmatrix}; \quad X3 = \begin{bmatrix} -1 & -1 \\ +1 & -1 \\ +1 & +1 \\ -1 & -1 \end{bmatrix}.$$

Ети матриці задовольняють умову (7). При проведенні імітаційних експериментів: $a1 = 1$, $a2$ змінювався від 1 до 2. Параметр регуляризації r змінювався на відрізку $[0;2]$ з кроком 0,1. В тестових планах присутствує сильна мультиколінеарність. Значення r_{OPT} розраховувалося за формулою (9), а також визначалося з допомогою імітаційного експеримента. Експериментальні дані функції відклику Y визначалися за моделлю, де випадкова складова генерувалася з використанням датчика випадкових чисел. Припускається, що в кожній точці імітаційного експеримента було проведено по 100 спроб.

Для проведення імітаційного експеримента було розроблено програмне забезпечення на мові м-файлів MATLAB.

В таблиці 1 представлені результати визначення оптимального значення параметра регуляризації r_{OPT} за формулою (9) і отриманого експериментальним шляхом, в залежності від варіанта експериментальних даних.

Таблиця 1– Результати визначення для повної моделі

Варіант	Значення коефіцієнтів	Значення r_{OPT} (раховане)	Значення r_{OPT} (експериментальне)
X1	$a2=1$	1	0,99
	$a2=2$	0,4	0,41
X2	$a2=1$	1	0,99
	$a2=2$	0,4	0,39
X3	$a2=1$	0,4	0,40
	$a2=2$	0,2	0,21

Результати теоретичного виводу і імітаційного експеримента практично збігаються. Таким чином, функція СКОП має локальний мінімум, яким і є r_{OPT} .

На рисунках 1–3 зображені графіки залежності СКОП EI від значення коефіцієнта регуляризації r для різних варіантів експериментальних даних.

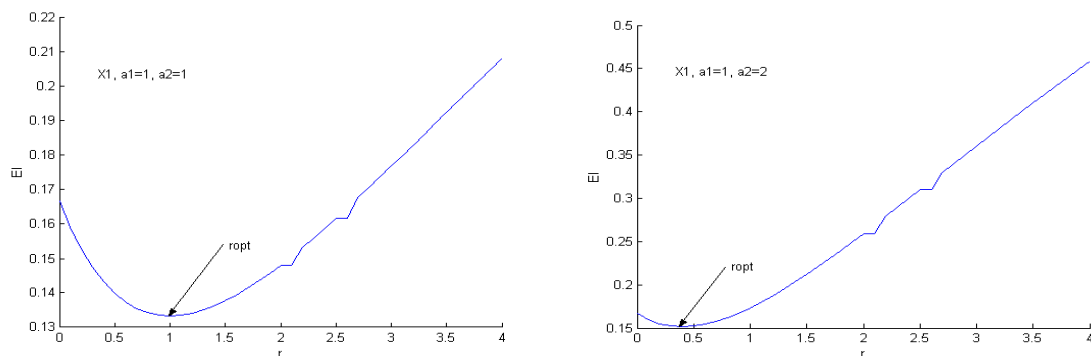


Рисунок 1 – Графік залежності СКОП від значення коефіцієнта регуляризації r для варіанта X1 ($a2=1,2$)

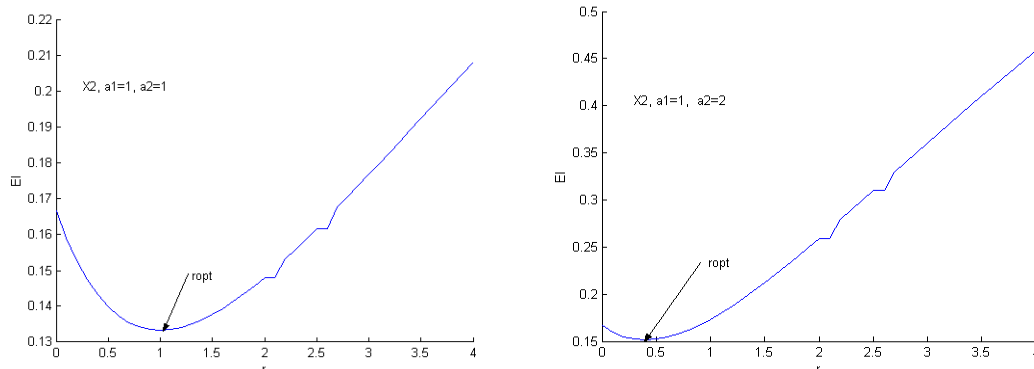


Рисунок 2 – Графік залежності СКОП від значення коефіцієнта регуляризації r для варіанта X2 ($a2=1,2$)

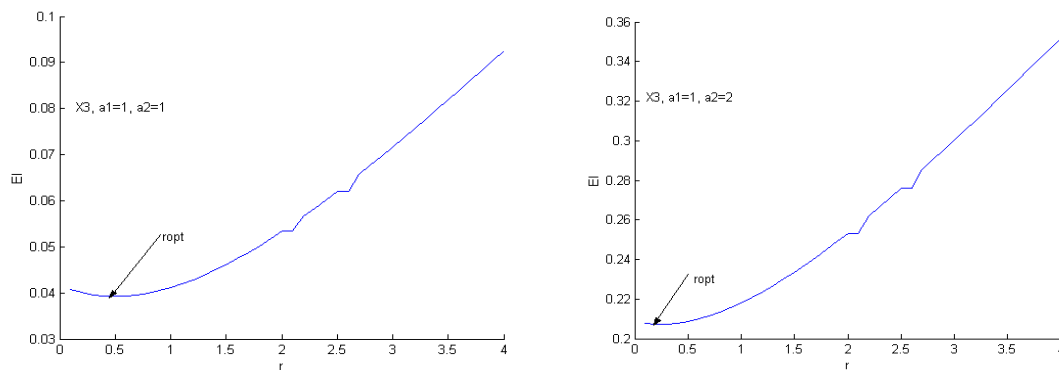


Рисунок 3 – Графік залежності СКОП від значення коефіцієнта регуляризації r для варіанта $X3$ ($a2=1,2$)

Формула (9) може бути використана тільки для обмеженого числа матриць експериментальних даних F , для яких виконується умова (7), наприклад, для матриць з даними в точках границі для змінних. Для інших варіантів матриць F визначення r_{OPT} є оптимізаційною пошуковою задачею. Для таких матриць експериментальних даних пропонується використовувати ітеративний метод пошуку r .

Висновки та перспективи подальших досліджень. В результаті застосування розробленої методики проведено дослідження СКОП СОУ на предмет існування локального мінімуму, що довелося і дозволило визначити значення оптимального параметра регуляризації при зміщеному оцінюванні параметрів регресійних рівнянь. Після проведених імітаційних експериментів видно, що отримане в результаті застосування методики значення оптимального параметра регуляризації r_{OPT} практично не відрізняється від експериментального. Це підтверджує правильність розробленої методики.

Таким чином, можна прийти до висновку, що в розглянутої задачі досягнута поставлена мета – підвищено якість прогнозуючої моделі СОУ в умовах мультиколінеарності, з урахуванням неопределенності структури моделі і методу зміщеного оцінювання параметрів моделі, що свідчить про коректність використаних і розроблених методів управління, створеного програмного забезпечення на мові м-файлів MATLAB, що реалізує запропоновану методику.

Предметом подальших досліджень передбачається вивчення ітеративного методу знаходження оптимального значення параметра регуляризації r_{OPT} для рідж-регресії.

Бібліографічний список використаної літератури

1. Айвазян С.А. Прикладна статистика: Дослідження залежностей: посібник / С.А. Айвазян, І.С. Енюков, Л.Д. Мешалкин; під ред. С.А. Айвазяна. — М.: Финансы и статистика, 1985. — 487 с.
2. Бородин С.А. Эконометрика: учеб. пособие для вузов / С.А. Бородин. — Минск: Новое знание, 2001. — 416 с.
3. Тихонов А.Н. Методы решения некорректных задач / А.Н. Тихонов, В.Я. Арсенин. — М.: Наука, 1974. — 276 с.
4. Кобзарь А.И. Прикладная математическая статистика. Для инженеров и научных работников / А.И. Кобзарь. — М.: Физматлит, 2006. — 816 с.

Поступила в редакцию 17.01.2014 г.

Потаніна М.В. Властивості середньоквадратичної помилки прогнозу рідж-регресії для ідентифікації моделі складного об'єкта управління при мультиколінеарності факторів

Проведено дослідження функції середньоквадратичної помилки прогнозу (СКПП) у разі рідж-регресії на екстремум в залежності від параметра регуляризації. Доведено наявність локального мінімуму СКПП, що дозволяє вирішити задачу пошуку оптимального параметра регуляризації для ідентифікації моделі складного об'єкта управління при мультиколінеарності факторів.

Ключові слова: середньоквадратична помилка прогнозу, рідж-регресія, параметр регуляризації, ідентифікація, мультиколінеарність.

Potanina M.V. The properties of the mean-squared prediction error ridge regression to identify models of complex control object by multicollinearity

The research on the functions of the mean-squared prediction error (MSPE) for ridge-regression to the extremum depending on the regularization parameter. Proved the existence of a local minimum of MSPE, that allows solving a problem of search of the optimal regularization parameter to identify models of complex control object by multicollinearity.

Keywords: mean-squared prediction error, ridge regression, the regularization parameter, identification, multicollinearity.