

<https://doi.org/10.3176/phys.math.tech.1965.4.05>

И. ПЕТЕРСЕН

ПРИМЕНЕНИЕ МЕТОДА ГЛАВНЫХ КОМПОНЕНТОВ ДЛЯ ОПИСАНИЯ ТЕХНОЛОГИЧЕСКИХ ПРОЦЕССОВ С КОРРЕЛИРОВАННЫМИ ВХОДНЫМИ ПАРАМЕТРАМИ

Для построения математического описания технологического процесса на основании статистических данных нормальных режимов предлагается методика, сочетающая статистическую конденсацию коррелированных входных параметров путем выделения главных компонент с регрессионным анализом. Эта методика позволяет получить некоторую модель процесса даже в тех случаях, когда входные параметры сильно коррелированы, или когда число наблюдений меньше числа входных параметров.

Введение

При построении математических моделей статистики реальных технологических процессов приходится часто прибегать к статистическим методам. Если все интересующие нас входные параметры процесса поддаются управлению и имеется возможность провести эксперименты над процессом, то регрессионный анализ результатов оптимально планированных экспериментов позволяет построить модель, обладающую высокой достоверностью. Однако на практике большая часть параметров процесса обычно плохо поддается управлению и проведение экспериментов над производственными процессами не всегда удается. В такой ситуации желательно возможно полнее использовать ту легко доступную информацию, которая накапливается при нормальных режимах.

При «хороших» процессах, которые протекают все время в одинаковом режиме, данные о нормальных режимах не содержат никакой информации, кроме средних значений параметров. Но при более реальных процессах изменяются в некоторых пределах как свойства сырья и внешние условия, так и параметры технологии.

Для нормальных режимов характерна сильная коррелированность, а иногда даже почти чистая линейная зависимость входных параметров процесса. Параметры сырья, внешние условия и многие технологические параметры естественным образом связаны друг с другом. Оператор и система регулирования коррелируют управляемые параметры между собой и с другими параметрами процесса.

Коррелированность входных параметров ограничивает информацию, которую можно получить в результате наблюдения за процессом. Эффекты изменения входных параметров на выходы процесса в таких направлениях (соотношениях), в которых наблюдались заметные изменения, могут определяться с достаточной достоверностью, в то время, как в других направлениях, где параметры мало изменились, эти эффекты трудно отделить от шумового фона. Поэтому регрессионный анализ требует при заметной коррелированности входов тщательного статистического исследования резуль-

татов. Для такого исследования, в частности, придется определить собственные значения и собственные векторы корреляционной матрицы вектора коэффициентов регрессионного уравнения.

С другой стороны, коррелированность входов вызывает ряд вычислительных трудностей. Если наблюдения входов линейно зависимы (например, когда число наблюдений меньше числа параметров), то регрессионный анализ непосредственно совсем не применим. Если входы сильно коррелированы, то матрица нормальных уравнений плохо обусловлена и поэтому коэффициенты регрессионного уравнения и их корреляционная матрица получаются с большими вычислительными ошибками. Эти трудности быстро растут при увеличении числа входов.

В свете указанных обстоятельств для построения моделей процессов по данным нормальных режимов можно рекомендовать следующий путь. Сперва проводится репараметризация процесса с введением новых параметров в виде линейных комбинаций исходных так, чтобы новые параметры были некоррелированы и чтобы они при возможно меньшем их числе учитывали существенную часть изменчивости исходных параметров. Этим требованиям удовлетворяют главные компоненты входных параметров. После репараметризации определяются регрессии выходов над новыми параметрами. Благодаря некоррелированности новых параметров при таком подходе, упрощается статистический анализ регрессионных уравнений, не возникает вычислительных трудностей при решении нормальных уравнений и уменьшение числа параметров упрощает дальнейшее использование регрессионных уравнений. В случае необходимости можно возвращаться к исходным параметрам процесса. Тогда соответствующие уравнения относятся к некоторому подпространству пространства исходных параметров.

Такой подход можно особенно рекомендовать для описания процессов, у которых коррелированность входных параметров носит устойчивый характер, так как в этом случае и новые параметры — главные компоненты — также устойчивы и имеют определенное самостоятельное значение.

1. Главные компоненты входных параметров

Пусть рассматриваемый технологический процесс имеет n входных параметров $\mathbf{x} = (x_1, \dots, x_n)$ и один выходной параметр y . Значения этих параметров в N наблюдениях составляют соответственно $N \times n$ -матрицу (x_{ij}) и $N \times 1$ -матрицу (y_i) . Средние значения параметров в данной серии наблюдений обозначим соответственно $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_n)$ и \bar{y} . Еще введем в рассмотрение диагональную матрицу L масштабов l_j ($j = 1, \dots, n$) входных параметров. На практике масштабные коэффициенты l_j следует выбирать с учетом области изменения входных параметров и их влияния на выходной параметр так, чтобы в этих масштабах входные параметры играли примерно одинаковую роль при описании процесса. Если нет других соображений, то можно в качестве масштабов взять выборочные стандартные отклонения $l_j^2 = \frac{1}{N-1} \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2$.

Обозначим масштабированные и центрированные входные параметры через $\mathbf{z} = (z_1, \dots, z_n)$, так что

$$\mathbf{z} = (\mathbf{x} - \bar{\mathbf{x}}) L^{-1}. \tag{1.1}$$

Матрицу сумм произведений и квадратов переменных \mathbf{z} обозначим через C , так что

$$C = L^{-1} (x_{ij} - \bar{x}_j)' (x_{ij} - \bar{x}_j) L^{-1}. \tag{1.2}$$

Пусть $\lambda_1 > \dots > \lambda_m > \lambda_{m+1} = \dots = \lambda_n = 0$ — собственные значения матрицы C и $f_1, \dots, f_m, f_{m+1}, \dots, f_n$ — соответствующие ортонормированные собственные векторы. Обозначим $(f_{k1}, \dots, f_{kn}) = f_k$ и $F = (f_{kj})$.

$$\text{Тогда} \quad F'F = FF' = E_n. \quad (1.3)$$

Введем теперь новые переменные $u = (u_1, \dots, u_n)$, полагая

$$u = zF', \quad z = uF. \quad (1.4)$$

Переменные u называются главными компонентами [1, 2] переменных z . Они обладают следующим экстремальным свойством: для $k = 1, \dots, n$ u_k имеет максимальную вариацию на данной совокупности значений z (z_{ij}) среди всех линейных функций zc' , векторы коэффициентов которых c удовлетворяют условиям

$$cc' = 1, \quad f_l c' = 0 \quad (l = 1, \dots, k-1). \quad (1.5)$$

Благодаря этому свойству главные компоненты являются в некотором смысле наилучшими линейными функциями для описания изменений в режимах в той серии наблюдений, на основе которой они построены. Главные компоненты некоррелированы и некоторое количество q первых из них позволяет точнее всех других q линейных функций описать изменения в режимах (x_{ij}).

Вариация переменной u_k в данной серии наблюдений равна собственному числу λ_k . Первые q главных компонентов учитывают из полной вариации наблюдений долю

$$Q_q^2 = \frac{\lambda_1 + \dots + \lambda_q}{\lambda_1 + \dots + \lambda_n}. \quad (1.6)$$

С возрастанием q Q_q растет и, в случае сильной коррелированности входных параметров x , быстро приближается к единице. Поэтому можно часто с достаточной для практики точностью считать $Q_q = 1$ для некоторого $q < m$, что позволяет свести исследование влияния n коррелированных (в случае $m < n$ даже линейно зависимых) входных параметров к исследованию влияния меньшего числа q некоррелированных параметров. Основанием для выбора необходимого q на практике является формула (1.6). Имеются и статистические критерии для выбора q в случае нормальности распределения значений параметров x [3].

2. Регрессионный анализ над главными компонентами

Предположим, что входной параметр y является линейной функцией q первых главных компонент. Это предположение, в частности, законно, когда y — линейная функция исходных параметров x и $Q_q = 1$.

Обозначим вектор q первых главных компонент через $\tilde{u} = (u_1, \dots, u_q)$ и $N \times q$ -матрицу значений этих главных компонент в рассматриваемой серии наблюдений через (\tilde{u}_{ik}) ; $q \times n$ -матрицу первых q строк матрицы F обозначим \tilde{F} . Тогда

$$\tilde{F}\tilde{F}' = E_q \quad (2.1)$$

и

$$(\tilde{u}_{ik})'(\tilde{u}_{ik}) = \tilde{D}, \quad (2.2)$$

где \tilde{D} — q -мерная диагональная матрица, диагональными элементами которой являются вариации переменных u_1, \dots, u_q , т. е. $\lambda_1, \dots, \lambda_q$.

Линейная регрессия y над u_1, \dots, u_q приведет к уравнению

$$y = b_0 + b_1 u_1 + \dots + b_q u_q, \quad (2.3)$$

где
$$b_0 = \bar{y} \quad (2.4)$$

и вектор $\mathbf{b} = (b_1, \dots, b_q)$ определяется системой нормальных уравнений

$$(\tilde{u}_{ik})' (\tilde{u}_{ik}) \mathbf{b}' = (\tilde{u}_{ik})' (y_i - \bar{y}). \quad (2.5)$$

Учитывая (2.2), (1.4) и (1.1), получаем

$$\mathbf{b}' = \tilde{D}^{-1} \tilde{F} L^{-1} (x_{ij} - \bar{x}_{.j})' (y_i - \bar{y}) \quad (2.6)$$

или

$$b_k = \frac{1}{\lambda_k} \sum_{j=1}^n \frac{f_{kj}}{l_j} \sum_{i=1}^N (x_{ij} - \bar{x}_{.j}) (y_i - \bar{y}) \quad (k=1, \dots, q). \quad (2.7)$$

Предположим, что ошибки измерения выхода y нормально распределены, независимы и имеют одинаковую дисперсию σ^2 . Корреляционной матрицей вектора $(y_i - \bar{y})$ тогда будет матрица $C_y = \sigma^2 E_N$ и, следовательно, корреляционная матрица C_b вектора \mathbf{b} выражается в виде

$$\begin{aligned} C_b &= \tilde{D}^{-1} \tilde{F} L^{-1} (x_{ij} - \bar{x}_{.j})' \sigma^2 E_N [\tilde{D}^{-1} \tilde{F} L^{-1} (x_{ij} - \bar{x}_{.j})]' = \\ &= \sigma^2 \tilde{D}^{-1} \tilde{F} L^{-1} (x_{ij} - \bar{x}_{.j})' (x_{ij} - \bar{x}_{.j}) L^{-1} \tilde{F}' \tilde{D}^{-1}. \end{aligned}$$

Поэтому на основании (1.2) и вытекающего из определения \tilde{F} равенства

$$\tilde{F} C \tilde{F}' = \tilde{D} \quad (2.8)$$

имеем
$$C_b = \sigma^2 \tilde{D}^{-1}, \quad (2.9)$$

так что коэффициенты b_k некоррелированы и их дисперсии σ_k^2 выражаются в виде

$$\sigma_k^2 = \frac{\sigma^2}{\lambda_k} \quad (k=1, \dots, q). \quad (2.10)$$

При этом для σ^2 можно через остаточную дисперсию получить [4] оценку с $N - q - 1$ степенями свободы.

$$s^2 = \frac{1}{N - q - 1} [(y_i - \bar{y})' (y_i - \bar{y}) - \mathbf{b} \tilde{D} \mathbf{b}'], \quad (2.11)$$

так что с $N - q - 1$ степенями свободы стандартное отклонение коэффициента b_k ($k=1, \dots, q$) оценивается величиной

$$s_k = \frac{1}{\sqrt{\lambda_k}} \cdot \frac{1}{\sqrt{N - q - 1}} \sqrt{\sum_{i=1}^N (y_i - \bar{y})^2 - \sum_{k=1}^q \lambda_k b_k^2}. \quad (2.12)$$

З а м е ч а н и е. Если $m = n$, т. е. матрица $(x_{ij} - \bar{x}_{.j})'(x_{ij} - \bar{x}_{.j})$ невырожденная, и вектор $\mathbf{a} = (a_1, \dots, a_n)$ является вектором коэффициентов регрессионного уравнения

$$y = a_0 + a_1(x_1 - \bar{x}_{.1}) + \dots + a_n(x_n - \bar{x}_{.n}), \quad (2.13)$$

построенного по наблюдениям (x_{ij}) , то регрессионный коэффициент b_k в уравнении (2.3) является проекцией вектора \mathbf{a} на вектор \mathbf{f}_k в метрике, определенной масштабами L :

$$b_k = \mathbf{f}_k L \mathbf{a}'. \quad (2.14)$$

Действительно, с одной стороны, по определению \mathbf{f}_k имеем

$$\mathbf{f}_k C = \lambda_k \mathbf{f}_k$$

и, следовательно,

$$\mathbf{f}_k L [(x_{ij} - \bar{x}_{.j})'(x_{ij} - \bar{x}_{.j})]^{-1} = \frac{1}{\lambda_k} \mathbf{f}_k L^{-1}. \quad (2.15)$$

С другой стороны, вектор \mathbf{a} как решение системы нормальных уравнений определяется в виде

$$\mathbf{a}' = [(x_{ij} - \bar{x}_{.j})'(x_{ij} - \bar{x}_{.j})]^{-1} (x_{ij} - \bar{x}_{.j})'(y_i - \bar{y}). \quad (2.16)$$

Сопоставление (2.16), (2.15) и (2.6) даст (2.14).

3. Описание процесса

Регрессионное уравнение в главных компонентах (2.3) может быть использовано при описании процесса двумя способами. Первый способ заключается в репараметризации процесса переходом к u_1, \dots, u_q как к новым переменным. Переменные \mathbf{x} и $\tilde{\mathbf{u}}$ при этом связаны формулами преобразований, полученными на основании (1.1) и (1.4):

$$\tilde{\mathbf{u}} = (\mathbf{x} - \bar{\mathbf{x}}) L^{-1} \tilde{F}', \quad \mathbf{x} = \bar{\mathbf{x}} + \tilde{\mathbf{u}} \tilde{F} L \quad (3.1)$$

или в более подробной записи

$$u_k = \sum_{j=1}^n \frac{f_{kj}}{l_j} (x_j - \bar{x}_{.j}) \quad (k = 1, \dots, q), \quad (3.2)$$

$$x_j = \bar{x}_{.j} + \sum_{k=1}^q l_j f_{kj} u_k \quad (j = 1, \dots, n).$$

Регрессионное уравнение в этом случае применяется в виде (2.3). Область изменения переменных \mathbf{u} , где оправдано использование уравнения (2.3), обычно определяется по допустимому фидуциальному пределу. В нашем случае такой областью является эллипсоид

$$\frac{u_1^2}{\lambda_1} + \dots + \frac{u_q^2}{\lambda_q} \leq \frac{\Delta^2}{l^2 \sigma^2} - \frac{1}{N} \quad (3.3)$$

(при условии $\frac{\Delta^2}{t^2\sigma^2} - \frac{1}{N} > 0$), где Δ — предельная допустимая ошибка выхода с данной фидуциальной вероятностью и t — соответствующее значение распределения Стюдента.

Второй способ заключается в следующем. Если в уравнении (2.3) подстановкой (3.1) перейти к переменным \mathbf{x} , то получим уравнение вида

$$y = c_0 + c_1(x_1 - \bar{x}_{.1}) + \dots + c_n(x_n - \bar{x}_{.n}), \tag{3.4}$$

где $c_0 = \bar{y}$, и для вектора $\mathbf{c} = (c_1, \dots, c_n)$ имеем

$$\mathbf{c} = \mathbf{b}\tilde{F}\mathbf{L}. \tag{3.5}$$

При этом использование уравнения (3.4) оправдано только при таких \mathbf{x} , для которых соответствующее $\mathbf{z} = (\mathbf{x} - \bar{\mathbf{x}})\mathbf{L}^{-1}$ ортогонально $n - q$ векторам $\mathbf{f}_{q+1}, \dots, \mathbf{f}_n$. Последнее условие определяет q -мерную плоскость, уравнениями которой являются

$$\sum_{j=1}^n \frac{f_{kj}}{l_j} (x_j - \bar{x}_{.j}) = 0 \quad (k = q + 1, \dots, n). \tag{3.6}$$

Описание процесса уравнениями (3.4) и (3.6) явно подчеркивает тот факт, что информация о процессе, имеющаяся в данных наблюдениях, ограничивает описание процесса режимами, которые соответствуют точкам определенной q -мерной плоскости. Для ограничения ошибки уравнения (3.4) выбранной фидуциальной вероятностью можно к уравнениям (3.4) и (3.6) прибавить неравенство, соответствующее (3.3).

4. Пример

Для иллюстрации приведем результаты одного расчета по описанной методике.

Рассматривался процесс с четырьмя входными параметрами x_1, x_2, x_3, x_4 и с двумя выходами Φ и Ψ . Число наблюдений $N = 98$ с оценками средних и стандартных отклонений:

$\bar{x}_1 = 0,947$	$s_{x_1} = 0,0326$	$\bar{\Phi} = 0,682$
$\bar{x}_2 = 0,536$	$s_{x_2} = 0,0187$	$\bar{\Psi} = 0,208$
$\bar{x}_3 = 50,94$	$s_{x_3} = 0,971$	$s_{\Phi} = 0,0359$
$\bar{x}_4 = 91,67$	$s_{x_4} = 53,27$	$s_{\Psi} = 0,0303$

Выборочные корреляционные коэффициенты параметров приведены в следующей таблице:

	x_1	x_2	x_3	x_4	Φ	Ψ
x_1	1	0,068	-0,735	0,544	0,375	-0,455
x_2		1	0,265	-0,130	0,180	-0,090
x_3			1	-0,220	-0,190	0,124
x_4				1	-0,055	-0,567

Обыкновенный регрессионный анализ дал следующие уравнения (в скобках после коэффициента дана оценка стандартного отклонения этого коэффициента на основании остаточной дисперсии):

$$\begin{aligned}\Phi &= -1,017 + 1,053(0,207)x_1 - 0,118(0,208)x_2 + \\ &\quad + 0,0156(0,00610)x_3 - 0,000331(0,0000810)x_4 \\ \Psi &= 1,098 - 0,414(0,162)x_1 - 0,0579(0,163)x_2 - \\ &\quad - 0,00875(0,00477)x_3 - 0,000221(0,0000634)x_4.\end{aligned}$$

Для сравнения с регрессионным анализом над главными компонентами были вычислены также регрессионные уравнения с тремя аргументами без учета того аргумента, который меньше других уменьшает остаточную дисперсию. Для обоих выходов таким аргументом оказался x_2 . Соответствующие уравнения получились следующие:

$$\begin{aligned}\Phi &= -0,923 + 0,989(0,173)x_1 + 0,0137(0,00499)x_3 - \\ &\quad - 0,000312(0,0000735)x_4 \\ \Psi &= 1,146 - 0,446(0,135)x_1 - 0,00973(0,00390)x_3 - \\ &\quad - 0,000212(0,0000576)x_4.\end{aligned}$$

При построении главных компонент в качестве масштабных коэффициентов были взяты соответствующие стандартные отклонения $l_j = s_{x_j}$. Значения q_j получились: $q_1 = 0,717$, $q_2 = 0,881$ и $q_3 = 0,986$ ($q_4 = 1$). Как видно, данные наблюдения не особенно нуждаются в применении рассматриваемой методики — только три главных компонента позволяют охватить существенную часть вариаций входных параметров. Это обстоятельство явствует и из корреляционной матрицы этих параметров: заметно коррелированы лишь x_1 , x_3 и x_1 , x_4 .

Первые три главных компонента получились следующие:

$$\begin{aligned}u_1 &= -19,6(x_1 - 0,947) + 8,75(x_2 - 0,536) + 0,602(x_3 - 50,9) - 0,00885(x_4 - 91,7) \\ u_2 &= 9,03(x_1 - 0,947) + 50,2(x_2 - 0,536) + 0,154(x_3 - 50,9) + 0,00207(x_4 - 91,7) \\ u_3 &= -4,09(x_1 - 0,947) - 7,57(x_2 - 0,536) + 0,565(x_3 - 50,9) + 0,0153(x_4 - 91,7).\end{aligned}$$

Приравнивание к нулю четвертого главного компонента определяет трехмерную гиперплоскость уравнением

$$21,3(x_1 - 0,947) - 14,7(x_2 - 0,536) + 0,595(x_3 - 50,9) - 0,00607(x_4 - 91,7) = 0.$$

Обратное преобразование, ограниченное этой гиперплоскостью, имеет вид

$$\begin{aligned}x_1 &= 0,947 - 0,0209 u_1 + 0,00960 u_2 - 0,00435 u_3 \\ x_2 &= 0,536 + 0,00305 u_1 + 0,0175 u_2 - 0,00264 u_3 \\ x_3 &= 50,9 + 0,568 u_1 + 0,145 u_2 + 0,533 u_3 \\ x_4 &= 91,7 - 25,1 u_1 + 5,87 u_2 + 43,3 u_3.\end{aligned}$$

Регрессионные уравнения над тремя главными компонентами для выходов Φ и Ψ получились следующего вида (в скобках даны оценки стандартных отклонений коэффициентов):

$$\begin{aligned}\Phi &= 0,682 - 0,00516(0,00236)u_1 + 0,00840(0,00331)u_2 - 0,0103(0,00382)u_3 \\ \Psi &= 0,208 + 0,00907(0,00172)u_1 - 0,00758(0,00246)u_2 - 0,0123(0,00278)u_3.\end{aligned}$$

Области, где с фидуциальной вероятностью 95% ожидаемая ошибка меньше, чем $\Delta = 0,005$, определены для Φ и Ψ соответственно неравенствами

$$\frac{u_1^2}{13,0} + \frac{u_2^2}{6,65} + \frac{u_3^2}{4,98} \leq 1, \quad \frac{u_1^2}{18,7} + \frac{u_2^2}{9,50} + \frac{u_3^2}{7,13} \leq 1.$$

В исходных параметрах x_1 , x_2 , x_3 и x_4 описание процесса по использованным наблюдениям получает вид

$$\Phi = 0,628 + 0,219x_1 + 0,454x_2 - 0,00762x_3 - 0,0000935x_4$$

$$\Psi = 0,666 - 0,196x_1 - 0,208x_2 - 0,00265x_3 - 0,000283x_4$$

$$21,3x_1 - 14,7x_2 + 0,595x_3 - 0,00607x_4 - 42,1 = 0.$$

ЛИТЕРАТУРА

1. Hotelling H., J. Educ. Psychol., **24**, 417—441, 498—520 (1933).
2. Андерсон Т., Введение в многомерный статистический анализ, М., 1963.
3. Lawley D. N., Biometrika, **43**, 128—136 (1956).
4. Линник Ю. В., Метод наименьших квадратов и основы теории обработки наблюдений, М., 1958.

Институт кибернетики
Академии наук Эстонской ССР

Поступила в редакцию
5/V 1965

I. PETERSEN

PEAKOMPONENTIDE MEETODI KASUTAMINE KORRELEERITUD SISENDPARAMEETRITEGA TEHNOLOOGILISTE PROTSESSIDE KIRJELDAMISEKS

Artiklis esitatakse meetod staatiliste tehnoloogiliste protsesside matemaatiliste mudelite leidmiseks nende protsesside normaalses režiimides kogutud statistilise vaatlusmaterjali alusel. Meetod põhineb protsessi reparametriseerimisel peakomponentide meetodil ning peakomponentide regressioonanalüüsil.

I. PETERSEN

USING PRINCIPAL COMPONENTS IN DESCRIBING TECHNOLOGICAL PROCESSES WITH CORRELATED INPUT PARAMETERS

A method is proposed of constructing mathematical models of static technological processes on the basis of data obtained from normal operation. The method is based on the reparametrization of input parameters by the method of principal components and on regression analysis of the principal components.