E. Solojentsev, A. Rybakov - RESEARCHES IN IDENTIFICATION OF LOGICAL AND PROBABILISTIC MODELS WITH GROUPS OF INCOMPATIBLE EVENTS

R&RATA # 4
(Vol.1) 2008, December

# RESEARCHES IN IDENTIFICATION OF LOGICAL AND PROBABILISTIC RISK MODELS WITH GROUPS OF INCOMPATIBLE EVENTS

Solojentsev E.D., Rybakov A.V.

●

Institute of Mechanical Engineering Problems of RAS,
sol@sapr.ipme.ru

**Abstract:** In this paper the results of the researches in identification of the logical and probabilistic (LP) risk models with groups of incompatible events are presented. The dependence of the criterion function on several parameters has been investigated. The parameters include: the total number of optimisations, the amplitude of parameters increments, the initial value of the criterion function (CF), the choice of identical or different amplitudes of increments for different parameters, objects risks distribution. An effective technology of defining the global extreme in the identification of LP-risk model for the calculation time, appreciable to practice has been suggested.

**Key words**: risk, logic, probability, model, identification, incompatible events

The logical and probabilistic risk models are almost twice as accurate and have seven times better robustness than other known classification methods [1,2]. However the task of multi-parameter and multi-criteria optimisation for training LP-models is characterised by exclusive difficulty [1-3]. In the process of identification of LP-risk models in business according to statistical data there arise a number of additional features and difficulties [1,2]:

- The criterion function $F_{max}$ (CF) is a number of correctly recognised good and bad objects, i.e. it accepts the integer values and it is stepped;
- CF has some local extrema, and depends on the high number of real positive arguments;
- The derivatives of the criterion function with respect to probabilities $P1_{jr}$ cannot be computed.
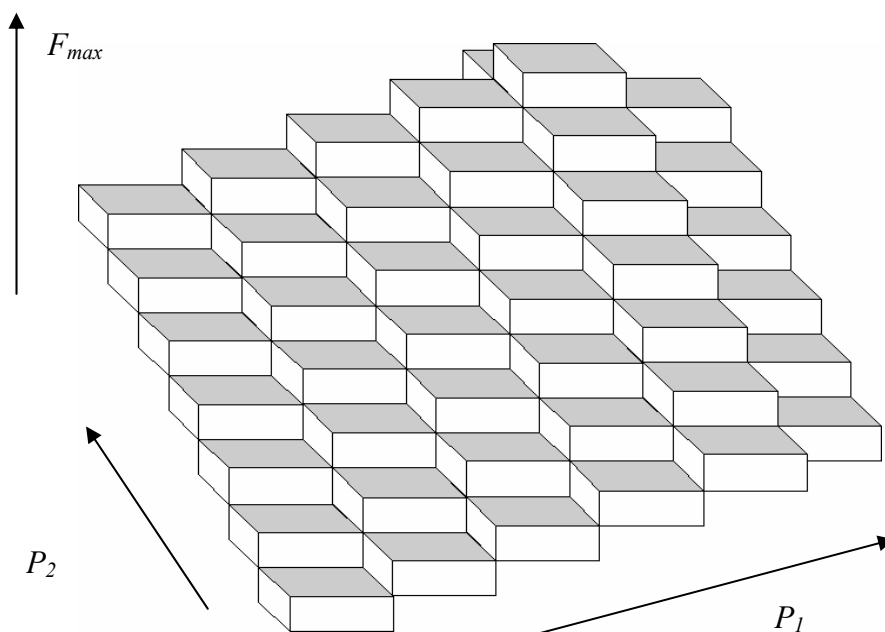


*Fig.1. The stepped changing of the criterion function $F_{max}$ from parameters $P_1$ and $P_2$*

For each event-grade in GIE we consider three probabilities: $W_{jr}$ is the relative frequency of the grade in the objects of the "object-signs" table, $P1_{jr}$ is the probability of the event-grade in GIE, $P_{jr}$ is the probability of the event-grade to be substituted into the probability formula. The sums of the probabilities both $W_{jr}$ and $P1_{jr}$ in GIE equal 1. Connection of these probabilities are considered in [1].

E. Solojentsev, A. Rybakov - RESEARCHES IN IDENTIFICATION OF LOGICAL AND PROBABILISTIC MODELS WITH GROUPS OF INCOMPATIBLE EVENTS

R&RATA # 4
(Vol.1) 2008, December

The criterion function $F_{max}$, presented in Fig.1, depends only on two arguments and changes with steps equal to 2. The platforms have different sizes. The arguments $P1_1$ and $P1_2$ belong to the interval [0,1], but their sizes can differ substantially. While approaching the extreme the platforms decrease in size.

The optimisation can get «stick» at any «platform», not reaching the maximum or crossing the maximum. The character of changing the criterion function in the multivariate space remains the same. Let us remind that the optimisation arguments space dimension for the credit risk LP-model equals 94 [1].

## 1. IDENTIFICATION OF LP-RISK MODELS

The risk object is described by a large number of signs, every sign has several grades. These signs and grades correspond to random events, which lead to a failure [1,2]. The events-signs ($j=1,n$) have logical connections and events-grades for each event-sign ($r=1,Nj$) form groups of incompatible events (GIE).

The identification of the P-risk model consists in the determination of optimal probabilities $P_{jr}$, $r = \overline{1, Nj}; j = \overline{1, n}$, corresponding to events-grades. Let us formulate the identification (training) problem for a B- risk model [1,2 ].

*Available data*: the 'object-signs' table with $N_g$ good and $N_b$ bad objects and the risk B-model;

*Expected results*: to determine the probabilities of $P_{jr}$, $r = \overline{1, Nj}; j = \overline{1, n}$ for events-grades and the acceptable risk $P_{ad}$, dividing the objects into good and bad according the amount of risk.

***We need: to maximise the criterion function, which is the number of correctly classified objects:***

$$(1) \qquad F = N_{bs} + N_{gs} \Rightarrow MAX,$$

where $N_{gs}$ *and* $N_{bs}$ *are the* numbers of objects classified as good and bad using both by statistics and the P- risk model (both estimates should coincide ). From (1) it follows, that the errors or accuracy indicators of the P-risk model in the classification of good $E_g$ and bad $E_b$ objects and in the classification of the whole set $E_m$ are equal:

$$(2) \qquad E_g = (N_g - N_{gs}) / N_g; E_b = (N_b - N_{bs}) / N_b; E_m = (N - F) / N.$$

*Assumed restrictions:*

1) probabilities $P_{jr}$ *and* $P1_{jr}$ *must satisfy the stipulation*:

$$(3) \qquad 0 < P_{jr} < 1, j = \overline{1, n}; r = \overline{1, Nj}.$$

2) the average risks of objects $P_m$ *based* on the P- risk model and on the table $P_{av}$ must be equal; while training the P- risk model we must correct the $P_{jr}$ probabilities on every step of iterative training under the formula

$$(4) \qquad P_{jr} = P_{jr} * (P_{av} / P_m); j = \overline{1, n}; r = \overline{1, Nj}.$$

3) the acceptable risk $P_{ad}$ must be determined with the given ratio of incorrectly classified good and bad objects, because of non-equivalence losses at their wrong classification:

$$(5) \qquad E_{gb} = (N_g - N_{gs}) / (N_b - N_{bs}).$$

## 2. OPTIMISATION IN THE IDENTIFICATION TASK

Identification of the LP- risk model by the random search method is based on the ideas used in the training of neural networks [4]. With reference to the identification task of the LP- risk model, the following formula for the calculation of the changes of events-grades probabilities may be put down:

$$(6) \qquad dP1_{jr} = K_1 * (1 / N_t) * tg(K_3); j = \overline{1, n}; r = \overline{1, Nj},$$

E. Solojentsev, A. Rybakov - RESEARCHES IN IDENTIFICATION OF LOGICAL AND PROBABILISTIC MODELS WITH GROUPS OF INCOMPATIBLE EVENTS

R&RATA # 4
(Vol.1) 2008, December

where: $K_1$ is a coefficient; $N_t$ is the current number of optimisation; $K_3$ is a random number from [- $\pi/2, +\pi/2$], **n** is a number of events-signs, $N_j$ is a number of events-grades in each GIE, i.a. for every event–sign.

In the formula (6) the CF is a current error in training. The number of optimisations $N_t$, before the end of the training process, can be very big. The «tangent» operation is the consequence of the training error distribution recording to Cauchy. Theoretically, this error is distributed according to the normal law, but not spend a lot of time on tabulated values calculation, we use the distribution of the training error under the Cauchy's law. It allows to reduce in 100 times the calculation time, which otherwise, for real problems, would continue for days and weeks.

For failure risk LP-model training the following modification of the formula (6) is suggested [1]:

$$(7) \qquad dP1_{jr} = K_1 * (N_{opt} - N_t) * tg(K_3), \ j = \overline{1, n}; r = \overline{1, Nj},$$

*where*: $N_{opt}$ is the given number of optimisations. The new values of $P1_{jr}$ and $P_{jr}$, obtained at $F > F_{max}$ on every step $N_t$ of optimisation are considered optimal and saved.

*In the LP-risk model identification task, the criterion function cannot exceed the total number of objects in the statistical data. The formula (7) is quite applicable, but the time of calculation is too big (about 10 hours for a session of optimisation).*

To reduce the time of calculation, in the formula (7) the "tangent" operation is eliminated. As a result the following expression is obtained [3]:

$$(8) \qquad dP1_{jr} = K_1 * (N_{opt} - N_t) * K_3, \ j = \overline{1, n}; r = \overline{1, Nj}.$$

Using (7,8) the optimization happens so: if F>Fmax, then we remember the new $P1_{jr}$ and $P_{jr}$. If the criterion function does not strictly increase after the chosen number of trials $N_{mc}$ in Monte-Karlo, then $F_{max}$ is reduced by 2-4 units and optimisation continues.

In spite of the investigation in optimisation, carried out before, where the formulas (7) and (8) were used [1,2], the problem of optimisation in the identification task of LP-risk models is far from the final solution. The following fact proves it. In one of the research with the huge number of optimisations $N_{opt}=245\ 000$ and with the constant, almost optimal, value of the increment $dP1_{jr}$, we obtained $F_{max} = 824$ instead of $F_{max} = 810$ at the usual number of optimisations $N_{opt} \approx 245$. We had to carry out special investigations, the results of which are adduced below.

## 3. INVESTIGATIONS IN IDENTIFICATION / OPTIMISATION

If we generate a random number $K_3$ in the interval [-1, +1], then the absolute values of increments of probabilities $dP1_{jr}$, multiplied by 100, are transformed in percents (%). It is convenient, for practically it solves the problem of the evaluation of probabilities $P1_{jr}$ accuracy. For example, if the increment is $dP1_{jr}=0.0005$, it equals *0.05 %*. We can say that the probability $P1_{jr}$ with the accuracy *0.05 %* is evaluated.

Using the formula (8), in the beginning of optimisation we have the following maximum amplitude of probabilities increments :

$$(9) \qquad AP1_{max} = K_1 * N_{opt}.$$

In the end of optimisation the maximum amplitude of probabilities increments equals 0. Let us designate the current amplitude of probabilities increments as *AP1*. There is an optimal interval *OPT* of the amplitudes increments *AP1*, which position and width are unknown (Fig. 2). For the big values of *AP1* there is a small probability of increasing $F_{max}$, and for small values of *AP1* there is a high probability to stop at the local extreme of the reached value $F_{max}$ (see Fig.1).
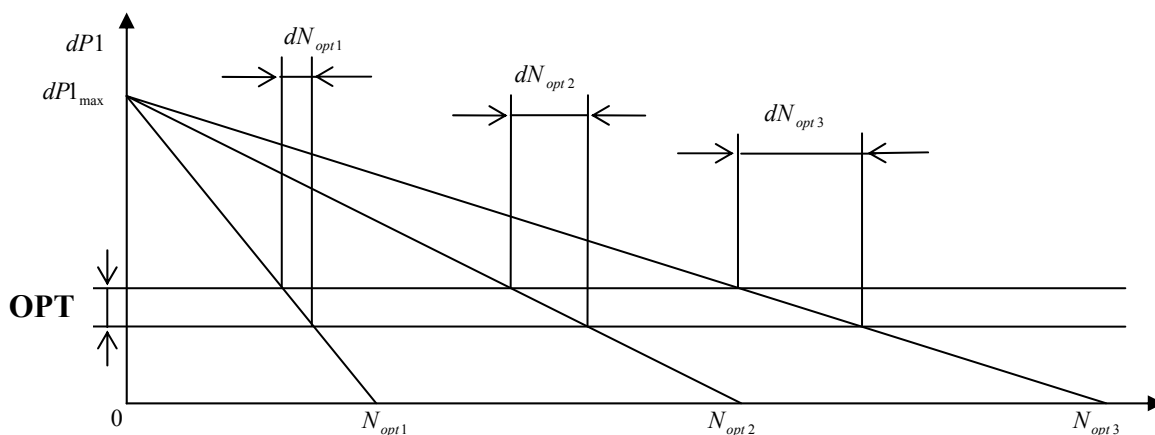
E. Solojentsev, A. Rybakov  -  RESEARCHES IN IDENTIFICATION OF LOGICAL AND PROBABILISTIC MODELS WITH GROUPS OF INCOMPATIBLE EVENTS

R&RATA # 4
(Vol.1) 2008, December

*Fig.2. Graphs of relation between  the  number of optimisations $N_{opt}$*
*and  increments amplitudes AP1*

The optimisation process ( of training the  LP-risk model)  should be long enough in  the optimal OPT interval . The value  of  $dN_{opt}$  duration in  the  optimal *OPT*  interval is  equal

*(10)* $$dN_{opt} \ = \ (OPT * N_{opt} \ ) \ / \ API_{max} \ .$$

It  also depends on  the  number of optimisations  $N_{opt}$  and the  maximum amplitude of the  increment $dP1_{max}$. The  more  $N_{opt}$  is  and  the  less  $API_{max}$  is , the  longer is  the  duration of  $dN_{opt}$. The  purpose of this work  is the  investigation  of    the  dependence of  the  criterion function (accuracy of LP- risk model) on  the  following parameters   in the  training formula (8):
1.  The  number of optimisations $N_{opt}$;
2.  The  increment  minimum amplitude  $API_{min}$, at which the  optimisation is still  possible;
3.  The  initial value of the criterion function $F_{beg}$ ;
4.  The  choice  of  identical or different amplitudes *AP1*   for different grades;
5.  The increment maximum amplitude  $API_{max}$;
6.  Objects risk distribution  in the  statistical data.
    Let us  illustrate it. A question arises,  whether  to choose the    identical or different    values of increments amplitudes *AP1*   for all events-grades ?  In other words, whether the  amplitudes  $API_{jr}$  should depend  on the  values of probabilities  $P1_{jr}$ ?  In the  training formulas of  the  LP-risk model (7) and (8) the increments  amplitudes  $API_{jr}$  are identical  for  all events-grades and  do  not  depend on the values  of their probabilities  $P1_{jr}$. The  increments  $dP1_{jr}$    differ  only  because of the    random simulation of  the   $K_3$ coefficient.
    The model investigations for  the LP-model of the  credit risk were made on  the  PC. The  credit risk structural LP-model has 20 events-signs (correspondingly GIE) and 94 events-grades. The   credit risk L-function  is  [1,2] :
    *(11)* $$Y = X_1 \bigcup X_2 \bigcup ... \bigcup X_{20}$$

Verbally  it can be formulated as follows: a failure occurs, if  any one, or any two, … or all  initiating events happen. After the  orthogonalization  of  the L-function (11) the following P-risk model  for   the evaluation of the  credit  risk  has been  obtained:
    (12) $$P = P_1 + P_2 Q_1 + P_3 Q_1 Q_2 + ....$$

The investigations  were  carried out   in  a set of 1000 credits, 700 of which were good and 300 - bad [5].  For calculation investigations   we  used the   Software , designed in  the   object-oriented languages Visual  C+++  and Java.

E. Solojentsev, A. Rybakov - RESEARCHES IN IDENTIFICATION OF LOGICAL AND PROBABILISTIC MODELS WITH GROUPS OF INCOMPATIBLE EVENTS

R&RATA # 4
(Vol.1) 2008, December

## 3.1 The choice of parameters $N_{opt}$ , $AP1_{min}$ , $F_{beg}$

In comparison with the optimal variant $F_{max}$ = 824, the initial variant had the probabilities $P1_{jr}$ without the last four signs. So the optimisation starts at $F_{beg}$ = 690-760. Such solution allowed to reduce calculation time.

The calculations were made for two values of increments maximum amplitudes: 1) $AP1_{max}$ =0.05 (5 %) , 2) $AP1_{max}$ = 0.1 (10 %) . We used the following numbers of optimisations $N_{opt}$: 150, 300, 500, 750, 1000, 2000, 3000, 4000, 5000, 6000, 7000, 8000.

The results of investigations presented in Table 1 (Var.2-21) and Fig.3 , allow to make the following conclusions:

1) The criterion function $F_{max}$ (column 6 in Table 1 and Fig.3 ) asymptotically increases with the growth of the number of $N_{opt}$ optimisation ;

2) The minimum amplitude $AP1_{min}$ (column 9) equals approximately $0.0025$ ( $0.25$ %); at the smaller values of $AP1_{min}$ the optimisation does not happen and the number of the last optimisation $N_{end}$ (column 10) is less, than the given number of $N_{opt}$ optimisations. It is necessary to modify the law of the change of $AP1$ during the training process , adding the constant line $AP1_{min}$ (Fig.4). It increases the chance to get the greater value of $F_{max}$;

3) The big value of $N_{opt}$ can lead to the disappearance of the B-C line (Fig. 4), which undoubtedly will deteriorate the process of optimisation.

4) The initial value of $F_{beg}$ (column 5) should not be lowered, as it often leads to low final values of $F_{max}$ (Fig. 5) because of the unsuccessful trajectory of optimisation process; in the considered case it is possible to accept $F_{beg}$ =750-760.

Taking into consideration the just made conclusions , instead of the formula (8) the following formula for training the LP-risk model is suggested:

$$(13) \qquad \text{If } AP1 < AP1_{min} \text{ , then } dP1_{jr} = AP1_{min} * K_3 \text{ ,}$$
$$\text{If } AP1 > AP1_{min} \text{ , then } dP1_{jr} = K_1 * (N_{opt} - N_t) * K_3.$$

The optimisation results using the formula (13) under $AP1_{min} = 0.0025$ (0.25 %), different $AP1_{max}$ = 0.098, 0.09, 0.03 (9.8 %, 9 %, 3 %), a rather large number of optimisations $N_{opt}$=5000-12000 and the high $F_{beg}$ =745 in Table 1 ( var. 22-24) are shown. In all variants high values of $F_{max}$ =812-822 have been obtained.

### Table 1. The investigations results in the choice of optimisation parameters

| N | Nopt | $K_1$ | AP1ma | Fbeg | Fmax | dPc | AP1min | Nend | Notes |
|---|------|-------|-------|------|------|-----|--------|------|-------|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| 1 | 2000 | 0.0001 | 0.2 | 776 | 786 | 0.204 | 0.1987 | 20 | |
| 2 | 300 | 0.000165 | 0.05 | 756 | 794 | 0.1969 | 0.00198 | 289 | (3) |
| 3 | 300 | 0.00033 | 0.1 | 712 | 790 | 0.221 | 0.00429 | 288 | (3) |
| 4 | 750 | 0.0000665 | 0.05 | 756 | 802 | 0.1641 | 0.00545 | 669 | (3) |
| 5 | 750 | 0.000133 | 0.1 | 692 | 790 | 0.2052 | 0.01316 | 652 | (3) |
| 6 | 1000 | 0.00005 | 0.05 | 750 | 802 | 0.1867 | 0.00350 | 931 | (3) |
| 7 | 1000 | 0.0001 | 0.1 | 708 | 792 | 0.2174 | 0.01580 | 843 | (3) |
| 8 | 2000 | 0.000025 | 0.05 | 776 | 808 | 0.1595 | 0.00747 | 1702 | (3) |
| 9 | 2000 | 0.00005 | 0.1 | 724 | 798 | 0.1802 | 0.01405 | 1720 | (3) |
| 10 | 3000 | 0.0000166 | 0.05 | 748 | 806 | 0.1867 | 0.00699 | 2581 | (3) |
| 11 | 3000 | 0.000033 | 0.1 | 708 | 806 | 0.1867 | 0.00501 | 2849 | (3) |
| 12 | 4000 | 0.0000125 | 0.05 | 744 | 812 | 0.1945 | 0.00791 | 3368 | (3) |
| 13 | 4000 | 0.000025 | 0.1 | 740 | 802 | 0.2121 | 0.00862 | 3656 | (3) |
| 14 | 5000 | 0.00001 | 0.05 | 754 | 806 | 0.1663 | 0.00556 | 4445 | (3) |
| 15 | 5000 | 0.00002 | 0.1 | 738 | 803 | 0.1586 | 0.00400 | 4801 | (3) |
| 16 | 6000 | 0.000016 | 0.1 | 710 | 810 | 0.1598 | 0.00625 | 5610 | (3) |
| 17 | 6000 | 0.0000183 | 0.109 | 736 | 810 | 0.1618 | 0.00495 | 5730 | (3) |
| 18 | 7000 | 0.0000071 | 0.05 | 764 | 810 | 0.2096 | 0.00407 | 6430 | (3) |
| 19 | 7000 | 0.0000142 | 0.1 | 734 | 810 | 0.1692 | 0.00745 | 6479 | (3) |
| 20 | 8000 | 0.0000062 | 0.05 | 764 | 810 | 0.1755 | 0.00985 | 6425 | (3) |

E. Solojentsev, A. Rybakov - RESEARCHES IN IDENTIFICATION OF LOGICAL AND PROBABILISTIC MODELS WITH GROUPS OF INCOMPATIBLE
EVENTS

R&RATA # 4
(Vol.1) 2008, December

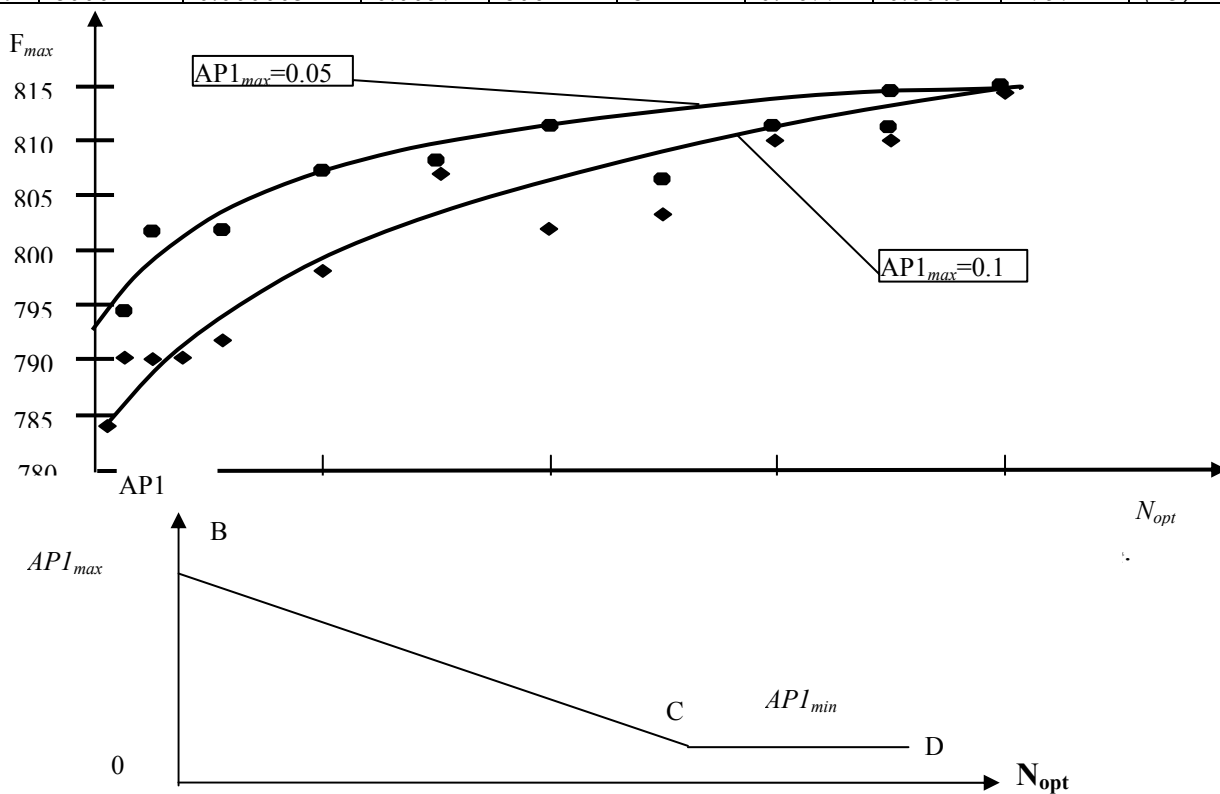| N | Nopt | K₁ | AP1ma | Fbeg | Fmax | dPc | AP1min | Nend | Notes |
|---|------|----|-------|------|------|-----|--------|------|-------|
| 21 | 8000 | 0.0000125 | 0.1 | 718 | 814 | 0.1802 | 0.00286 | 7772 | (3) |
| 22 | 12000 | 0.0000075 | 0.09 | 772 | 812 | 0.1737 | 0.0025 | 11754 | (10) |
| 23 | 8000 | 0.00000375 | 0.03 | 780 | 820 | 0.1526 | 0.0025 | 7662 | (10) |
| 24 | 8000 | 0.00000875 | 0.07 | 744 | 814 | 0.1733 | 0.0025 | 7801 | (10) |
| 25 | 5000 | 0.0000043 | 0.0215 | 812 | 820 | 0.1462 | 0.0025 | 23 | (13) |
| 26 | 5000 | 0.00000043 | 0.0025 | 810 | 824 | 0.1511 | 0.0025 | 34 | (13) |
| 27 | 8000 | 0.00000002 | 0.0025 | 810 | 826 | 0.1538 | 0.0025 | 678 | (13) |
| 28 | 8000 | 0.0000025 | 0.00458 | 806 | 822 | 0.1604 | 0.00609 | 507 | (13) |
| 29 | 8000 | 0.00000312 | 0.00572 | 806 | 822 | 0.1677 | 0.00452 | 1757 | (13) |



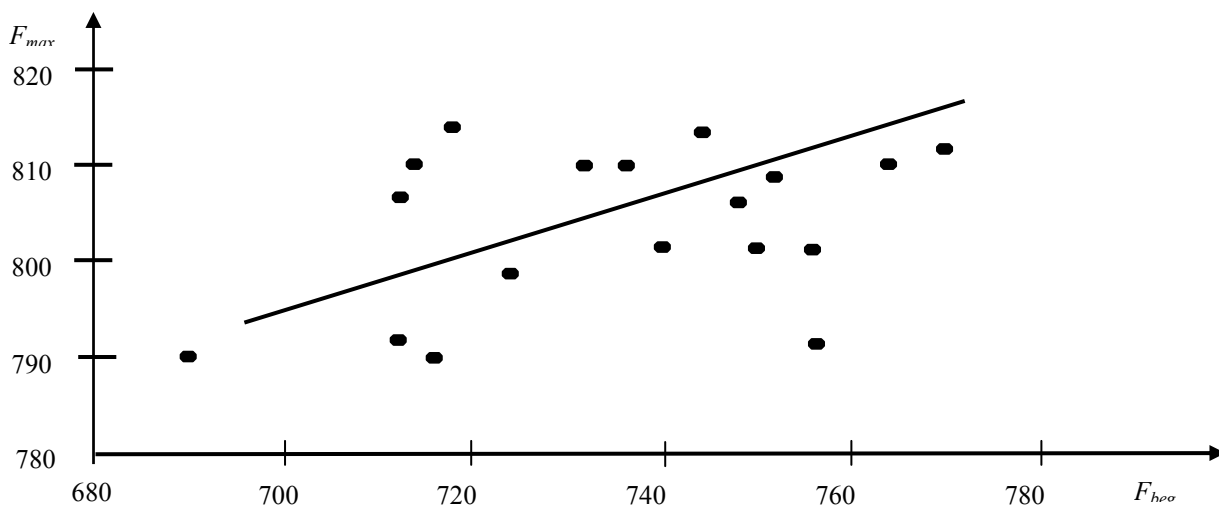Fig.4. The graph of the current amplitude of increment AP1 modification



Fig. 5. Dependence of the criterion function $F_{max}$ on its initial value

E. Solojentsev, A. Rybakov  -  RESEARCHES IN IDENTIFICATION OF LOGICAL AND PROBABILISTIC MODELS WITH GROUPS OF INCOMPATIBLE EVENTS

R&RATA # 4
(Vol.1) 2008, December

### 3.2 Different amplitudes $AP1_{jr}$ of increments for different grades

It should be noted, that the probabilities $P1_{jr}$ depend on: a number of grades in GIE, the frequencies of $W_{jr}$ grades in objects and the grades contributions in the classification errors of objects. In the formula of training the LP-risk model (8) the increments amplitudes $AP1_{jr}$ are identical for all events-grades and do not depend on the magnitude of their probabilities $P1_{jr}$.

Let us change the formula of the training LP-risk model so that it takes into account the value of probability for each grade

(14) $$dP1_{jr} = K_1 * (N_{opt} - N_t) * K_3 * P1_{jr}..$$

Here the amplitudes for every event grade are equal

*(15)* $$AP1_{jr} = K_1 * (N_{opt} - N_t) * P1_{jr}$$

and the formula (14) can be the following:

(16) $$dP1_{jr} = AP1jr * K_3.$$

Let us also put down the formula (14) with the following modification:

(17) $$dP1_{jr} = K_1 * (N_{opt} - N_t) * ((1-a) + a * P1_{jr}) * K_3,$$

where *a* is a coefficient from the interval *[0 < a < 1]*. It determines the formula (8) at *a=0*, the formula (14) at *a=1* and all the modifications at other values of *a*.

In the formula (13) let us take into account the limitations, introduced earlier in the formula (8), and we shall get the following expression for training the LP-risk model:

*(18)* If $AP1_{jr} < AP1_{min}$ , then $dP1_{jr} = AP1_{min}$ ,

If $AP1_{jr} > AP1_{min}$ , then $dP1_{jr} = K_1 * (N_{opt} - N_t) * ((1-a) + a * P1_{jr}) * K_3$ ,

The investigations results in optimisation using the formula (18) at *a=1* ($AP1_{max}$ = 2.15 % , 0.25 % , 0.45 % , 0.57%) are represented in Table 1 (Var.25-29). They show that the high values of the $F_{max}$ =822-826 can be obtained at the limited number of optimisation attempts $N_{end}$ (column 10). Actually the first optimisation already gives the high value of CF ($F_{beg}$=806-810). The optimisation process ends at $N_{end}$ = 23-1750 instead of the given numbers of optimisations $N_{opt}$=5000-8000 (column 6). It seems, that the number of optimisations $N_{opt}$ can be essentially reduced. To verify this hypothesis some extra investigations have been carried out.

The investigations were carried out at small numbers of optimisations $N_{opt}$ = 600, 450, 300, 150, 100, 50 and K1=0.00033, 0.00025, 0.00015, 0.0001. The increments maximum amplitude $AP1_{max}$ varied in an interval 0.5% - 20% from $P1_{jr}$. In Table 2 the CF values and the difference between maximum and minimum risks of objects in the statistics $F_{max}$ / APc are shown. The results of the investigations should be considered as good *($F_{max}$ =810-822)* and completely confirming the effectiveness of the formulas (14), (17) and (18).

Also the investigations of the influence of *a* parameter on the optimisation results have been carried out. It was done at the small numbers of optimisations $N_{opt}$=150 and $K_1$=0.00015. The maximum amplitude of an increment $AP1_{max}$ equals 0.0225* $P1_{jr}$.

Table 2. Values of $F_{max}$ / APc at the small number of optimisations $N_{opt}$ and *a=1*

| Number of optimizations, $N_{opt}$ | $K_1$=0.00033 | $K_1$=0.00025 | $K_1$=0.00015 | $K_1$=0.0001 |
|---|---|---|---|---|
| 600 | 798 / 0.248 | 796 / 0.225 | 810 / 0.180 | 810 / 0.149 |
| 450 | 802 / 0.217 | 804 / 0.187 | 814 / 0.162 | 819 / 0.161 |
| 300 | 810 / 0.146 | 810 / 0.174 | 816 / 0.147 | 820 / 0.162 |
| 225 | 810 / 0.154 | 811 / 0.152 | 818 / 0.148 | 821 / 0.146 |
| 150 | 816 / 0.145 | 820 / 0.156 | 822 / 0.148 | 822 / 0.147 |
| 100 | 818 / 0.146 | 820 / 0.149 | 820 / 0.151 | 820 / 0.153 |
| 50 | 822 / 0.151 | 820 / 0.146 | 820 / 0.152 | 820 / 0.148 |

The investigations results, represented in Table 3, also confirm the effectiveness of the formulas (14),(17) and (18) at *a=1*. Really, at *a=1* $F_{max}$ equals *820*, and at *a=0* $F_{max}$ equals *802*.

E. Solojentsev, A. Rybakov - RESEARCHES IN IDENTIFICATION OF LOGICAL AND PROBABILISTIC MODELS WITH GROUPS OF INCOMPATIBLE EVENTS

R&RATA # 4
(Vol.1) 2008, December

Table 3. Values $F_{max}$ at different values of $a$

| Value $a$ | 0.0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 1.0 |
|---|---|---|---|---|---|---|---|---|---|---|
| Value $F_{max}$ | 802 | 800 | 798 | 804 | 808 | 810 | 808 | 810 | 818 | 820 |

### 3.3 Determination of the amplitude $AP1_{max}$ and the global extreme $F_{max}$

Let us consider again the choice of the increment maximum amplitude of probabilities $AP1_{max}$. The results of the change of $F_{max}$ at the change of $AP1_{max}=K_1*N_{opt}$ in the interval $0.5\text{-}20\%$ of $P1_{jr}$ are represented in Table 2. They demonstrate that the higher is $AP1_{max}$ the less is $F_{max}$. In Fig.6 the dynamics and the results of optimisation for five variants, having $N_{opt}=2000$, are shown:

- Variant 1: $AP1_{max}=0.05(5\%)$, $F_{max}=808$ (Var.8 in Table1), training under the formula (3);
- Variant 2: $AP1_{max}=0.1(10\%)$, $F_{max}=798$ (Var.9 in Table1), training under the formula (3);
- Variant 3: $AP1_{max}=0.05(5\%)$, $F_{max}=820$, training under the formula (14) with a=1;
- Variant 4: $AP1_{max}=0.1(10\%)$, $F_{max}=804$, training under the formula (14) with a=1;
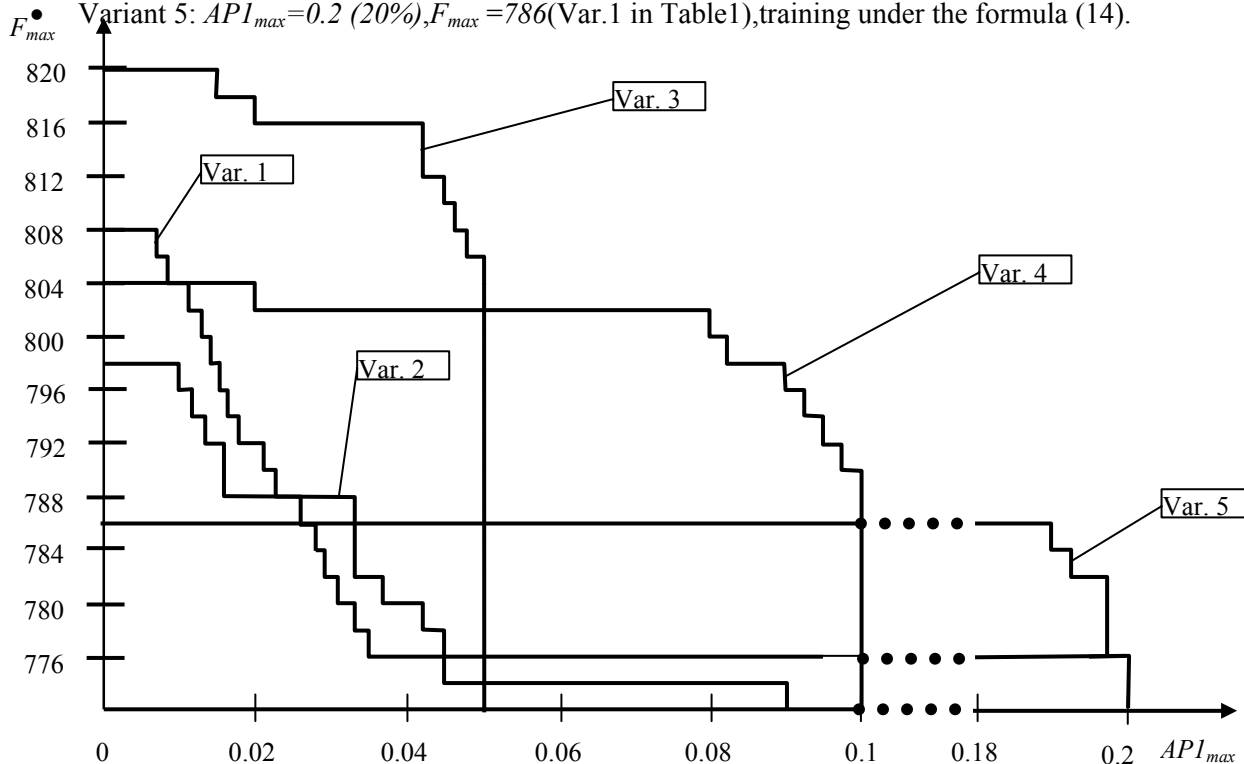- Variant 5: $AP1_{max}=0.2(20\%)$, $F_{max}=786$ (Var.1 in Table1), training under the formula (14).



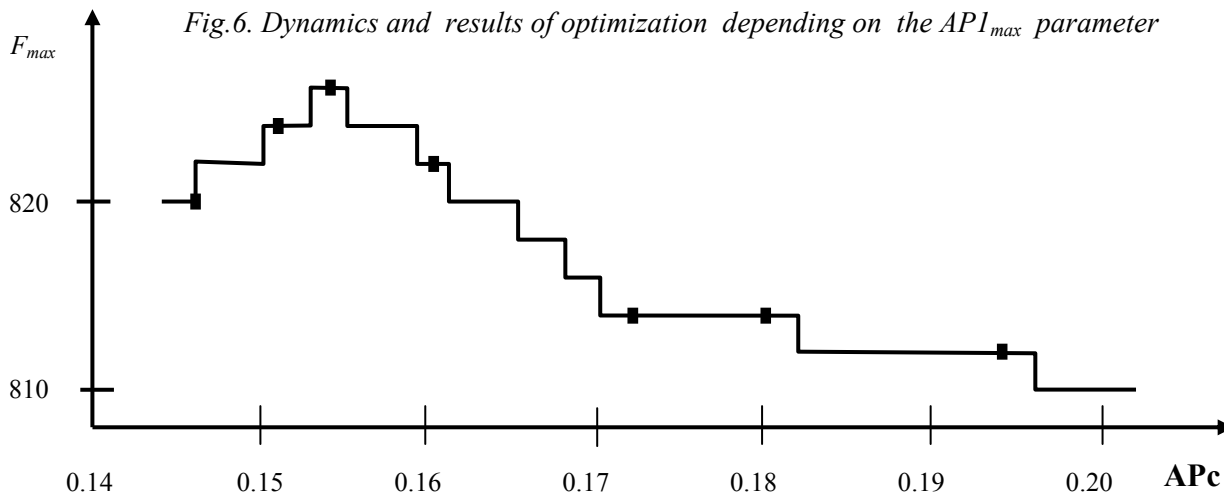Fig.6. Dynamics and results of optimization depending on the $AP1_{max}$ parameter



Fig.7.The connection of parameters $F_{max}$ and APc

E. Solojentsev, A. Rybakov  -  RESEARCHES IN IDENTIFICATION OF LOGICAL AND PROBABILISTIC MODELS WITH GROUPS OF INCOMPATIBLE
EVENTS

R&RATA # 4
(Vol.1) 2008, December

Variants 4 and 5 with  high $AP1_{max}$, despite  using  the effective formula (18)  and $a=1$, have bad training dynamics and  results. In these variants CF  are correspondingly *786* and *804*. The optimisation process  finishes early, *($N_{end}=1608$ and $N_{end}=20$)*. Additional  optimisation attempts  $N_{opt} – N_{end}$  have not increased  CF. This example confirms  that  the  increment amplitude $AP1_{max}$ should not be  more than *0.02 - 0.05 ( 2-5 % )*.

We  check  the  calculation of  the  global extreme of  the CF   by  the graph (Fig.7). The function  $F_{max}$ has an extreme  at   some value  of  the  difference *APc* between the  maximum risk and the minimum risk of objects in statistics [2 ]. This   difference, constructed for variants of computational  investigations , presented  in Table 1 and   2, demonstrates  the robustness (stability) of solutions at  a small dispersion of *APc* in  the  area  of  the  global extreme of  CF.

## 4. CONCLUSION

In the investigations   the following  main results have  been  obtained:
1.   The effective technology of   the  criterion function   global extreme search  in   the tasks   of identification of   LP-risk models under statistical data has been  offered .It permits to solve the task of   multi-parameter  multi-criteria optimisation  with  integer CF for the time, applicable to practice (less than  before).
2.   We suggest to generate in  the  training formula a random number $K_3$  in the interval [-1, +1]. It permits  to  consider the absolute values of  increments $dP1_{jr}$, multiplied  by  100, in percents (%) )  and  to estimate the  accuracy of  probabilities $P1_{jr}$.
3.   In the technology of the CF global extreme search, the following  regularities of  changing   the  CF should  be  used:
   * The  CF   asymptotically increases with the growth of  $N_{opt}$ optimisation number ;
   * The   minimum amplitude  $AP1_{min}$   of probabilities $P1_{jr}$  increments  is established by 2-3 test calculations; at smaller values of $AP1_{min}$  the optimisation does not happen (less than 0.25 %);
   * The  initial CF  $F_{beg}$  should  not  be  lowered , as  low values  more often  result in  low final values of  $F_{max}$ because of  the unsuccessful  trajectory  of  the optimisation  process;
   * Maximum amplitude of increments of  $AP1_{max}$  must  not  exceed  *0.02 - 0.05 (2-5%),* as the  training speed  lows  down  and  the value  of the  CF  $F_{max}$  becomes  less.
4.   For the criterion function global extreme search  new , more effective formulas of  training (14), (17), (18) have been suggested ;  they   use  different amplitudes of  increments for probabilities of different events-grades.
5.   It has  been  confirmed that we can test  the  determination of the global extreme of CF $F_{max}$ by  the graph of change of  $F_{max}$  in  the function of   difference *APc* between  maximum  and  minimum risks of objects  in statistics. The function $F_{max}$  has an extreme at a certain value  of  *APc* .

## REFERENCES

1.   Solojentsev E.D., Karasev V.V.,Solojentsev  V.E.  Logic and probabilistic models  of risk  in banks, business  and  quality / edited  by  E.D.Solojentsev. SPb.: Nauka, 1999.-120 p.
2.  Solojentsev E.D.,  Karassev V.V. (2001) Risk  logic and probabilistic models  in business    and identification of  risk  models. - Informatica 25 (2001) 49-55.
3.  Taha X. Introduce in  research  of  operations . v.1,2. Moscow,  Mir , 1985.
4.  Wasserman  Philips D. Neural   Computing  Theory  and  Practice. ANSA  Research , Inc. VAN NOSTRARD REINOLD, New York, 1990.
5.  Seitz J., Stickel E.  Consumer Loan Analysis Using Neural  Network.   Adaptive Intelligent Systems. Proceed. of the Bankai Workshop, Brussels,    14-19 October, 1992. P.177-189