

COPPEAD/UFRJ

RELATÓRIO COPPEAD Nº 206

AMOSTRAGEM DESCRITIVA:
UM CONTRA-EXEMPLO ?

EDUARDO SALIBY*

Fevereiro 1988

* Professor adjunto da COPPEAD, Instituto de Pós-Graduação e Pesquisa em Administração/UFRJ. Mestre pela COPPE/UFRJ e PhD pela Universidade de Lancaster, Inglaterra.

1. INTRODUÇÃO

Até o presente, praticamente todos os testes realizados mostraram ser a amostragem descritiva mais eficiente, em termos estatísticos, do que a amostragem aleatória simples em simulação por Monte Carlo; vide por exemplo (Saliby, 1988) e Funchal (1987). A amostragem descritiva, proposta por Saliby (1980), implica numa importante revisão conceitual em simulação pois se baseia numa seleção determinística dos elementos da amostra; assim sendo, ela deve ser minuciosamente testada, identificando-se suas limitações e situações onde possa inclusive levar a piores resultados que a abordagem tradicional. Conforme se poderá ver, ainda que se observe tal fato em casos muito particulares, ele não contradiz a teoria proposta.

O presente estudo refere-se a um problema, especialmente concebido, em que a amostragem descritiva produz resultados menos precisos do que a amostragem aleatória simples. Trata-se de um caso muito particular cujo objetivo inicial era o de servir como contra-exemplo para a amostragem descritiva; no entanto, após uma conceituação mais apurada quanto ao ganho de precisão esperado com a amostragem descritiva, ele teve um efeito oposto. Assim, constatou-se que a comparação entre ambos os métodos amostrais deve ser feita apenas quando o tamanho da amostra de entrada for suficientemente grande a ponto de se evitar eventuais distorções decorrentes de uma corrida subdimensionada. Tomando-se este cuidado, a amostragem descritiva continuará sendo sempre uma melhor opção na simulação por Monte Carlo.

2. O PROBLEMA ESTUDADO

O presente problema foi especialmente sugerido para evidenciar uma situação em que a amostragem descritiva leva a resultados menos precisos que a amostragem aleatória simples. Ele representa, assim, o que se poderia chamar de um caso "patológico". O seu enunciado se segue:

Seja X uma variável aleatória uniformemente distribuída no intervalo unitário. Sejam

$$q_i = (i - 0.5) / 100, \quad i = 1, \dots, 100,$$

as medianas de cada um dos 100 subintervalos consecutivos e de igual amplitude que dividem o intervalo unitário. Define-se uma nova variável aleatória Z , tal que

$$Z = \sum_{j=1}^5 |X_j - q_j|,$$

onde

q_j corresponde ao ponto mediano mais próximo de X_j , ou seja, tal que $|X_j - q_j|$, $i = 1, \dots, 100$, seja mínimo.

O problema consiste em estimar o valor de $E(Z)$ com o Método de Monte Carlo, usando-se ambos os métodos amostrais: a amostragem aleatória simples e a amostragem descritiva.

Antes de aplicar Monte Carlo a este problema, serão apresentados alguns resultados teóricos.

Resultados teóricos

Inicialmente, cabe notar que, independentemente do intervalo considerado,

$X_j - q_j$ é uniformemente distribuída no intervalo $(-1/200, 1/200)$. Conseqüentemente,

$|X_j - q_j|$ é uniformemente distribuída no intervalo $(0, 1/200)$.

Das propriedades da distribuição uniforme, segue-se que

$$E(|X_J - q_J|) = 1/400 \quad e$$

$$\text{Var}(|X_J - q_J|) = 1/(12 \times 200 \times 200) \quad .$$

Então,

$$E(Z) = 5/400 = 0.0125 \quad e$$

$$\text{Var}(Z) = 5/(12 \times 200 \times 200) = 0.00001042 \quad .$$

Com base nestes resultados, será estudada agora a sua solução por Monte Carlo, considerando inicialmente a amostragem aleatória simples.

3. SOLUÇÃO POR MONTE CARLO USANDO AMOSTRAGEM ALEATÓRIA SIMPLES

Neste caso, define-se uma corrida por um conjunto de observações

$$Z_t, \quad t = 1, \dots, N$$

da variável de resposta. Cada corrida origina uma única estimativa

$$\bar{Z} = \sum_{t=1}^N Z_t / N$$

Tem-se que

$$E(\bar{Z}) = E(Z) = 0.0125 \quad e$$

$$\text{Var}(\bar{Z}) = \text{Var}(Z)/N = 5/(12 \times 200 \times 200 \times N)$$

Ilustrando a aplicação de Monte Carlo a este problema, é apresentada no apêndice (A) a listagem do programa DESVIO escrito em TURBO-PASCAL. Juntamente, são também apresentados os resultados de um experimento composto de $M = 50$ corridas independentes, os quais confirmaram os valores teóricos acima.

4. SOLUÇÃO POR MONTE CARLO USANDO AMOSTRAGEM DESCRITIVA

Embora a amostragem descritiva interfira na seleção da amostra de entrada, o resto do processo permanece inalterado. Assim, da mesma forma que no caso anterior, cada corrida é definida por um conjunto de observações

$$ZD_t, \quad t = 1, \dots, N$$

originando uma estimativa

$$\bar{ZD} = \sum_{t=1}^N ZD_t / N$$

Cada observação é calculada a partir de 5 valores de entrada; segue-se que a amostra de entrada associada a uma corrida compõe-se de $n = 5 \times N$ valores. Enquanto estes valores são gerados aleatoriamente com a amostragem aleatória simples, no caso da amostragem descritiva eles são deterministicamente selecionados a partir da relação

$$XD_i = (i - 0.5) / n, \quad i = 1, \dots, n$$

Para este problema em particular, a estimativa em estudo independe da sequência dos valores de entrada. Assim sendo, uma vez fixado o conjunto de entrada, a estimativa \bar{ZD} estará automaticamente determinada, ou seja,

$$\text{Var}(\bar{ZD}) = 0$$

Embora com variância nula, \bar{ZD} irá variar com o número de observações por corrida (N), resultando freqüentemente numa estimativa tendenciosa. Neste caso, para se comparar a amostragem descritiva com a aleatória, deve-se utilizar a relação de precisão

$$\text{EMQ}(\bar{ZD}) / \text{Var}(\bar{Z})$$

onde

$$\text{EMQ}(\bar{ZD}) = [\bar{ZD} - E(Z)]^2$$

é o erro médio quadrático associado a \bar{ZD} . Note-se que, como $\text{Var}(\bar{ZD}) = 0$, apenas o termo do viés foi levado em conta no cálculo de $\text{EMQ}(\bar{ZD})$.

Como um primeiro passo, foi estudado o comportamento da relação de precisão em função do número de observações por corrida (N); isto foi feito com o programa DESVIOAD, também escrito em TURBO-PASCAL. Este programa, apresentado no apêndice (B), é uma versão do programa DESVIO cuja principal diferença reside no emprego da amostragem descritiva em lugar da aleatória simples.

Os dados da tabela 1, resultados do programa DESVIOAD, mostram que o viés associado à amostragem descritiva tem um comportamento oscilante com N ; note-se também que o seu valor é decrescente para $N > 20$. A relação de precisão $EMQ(\bar{ZD})/Var(\bar{Z})$ também oscila, e é também decrescente para $N > 20$; esta relação é máxima (300) para $N = 20$, o que representa a situação mais desfavorável para a amostragem descritiva. Neste caso, como se pode constatar, a amostragem descritiva é bem menos precisa que a amostragem aleatória simples. Ter-se-ia, pois, encontrado um contra-exemplo para a amostragem descritiva?

Embora seja um resultado que surpreenda a primeira vista, nada há de contraditório em relação à teoria da amostragem descritiva. Mesmo reconhecendo que a amostragem descritiva não é sempre mais eficiente do que a amostragem aleatória simples, o que realmente importa para verificar a validade da teoria proposta é o comportamento limite da relação de precisão quando o tamanho da amostra de entrada cresce indefinidamente. A comparação entre ambos os métodos amostrais deve ser feita dentro da premissa de que a corrida de simulação foi corretamente dimensionada, numa situação em que o tamanho da amostra é suficientemente grande a ponto de não mais distorcer os resultados da simulação. Enquanto o tamanho da amostra for insuficiente para eliminar tais distorções, ou seja, quando a corrida for subdimensionada, existe a possibilidade, ainda que remota, da amostragem descritiva ser ineficiente.

Neste exemplo, estudou-se propriedades relacionadas com frações centesimais da distribuição uniforme; portanto, não é de surpreender que para $N = 20$ ($n=100$) os resultados obtidos com a amostragem descritiva sejam piores, pois se dispõe apenas de 1 valor descritivo para cada subintervalo, valor este correspondente ainda a um caso extremo. Trata-se, pois, de uma situação em que a amostra de entrada é muito pequena tendo em vista a variável de resposta em estudo. Assim sendo, voltou-se a atenção para o comportamento da relação de precisão para $N > 20$, quando o viés associado a \bar{ZD} começa a decrescer.

Tabela 1. Viés e relação de precisão em função do número de observações por corrida (N) com amostragem descritiva.

N	\bar{ZD}	Viés(\bar{ZD})	EMQ(\bar{ZD})/Var(\bar{Z})
1	0.02500000	0.01250000	15.00000000
2	0.02500000	0.01250000	30.00000000
3	0.01388889	0.00138889	0.55555556
4	0.00000000	-0.01250000	60.00000000
5	0.02500000	0.01250000	75.00000000
6	0.01388889	0.00138889	1.11111111
7	0.01275510	0.00025510	0.04373178
8	0.01250000	-0.00000000	0.00000000
9	0.01265432	0.00015432	0.02057613
10	0.02500000	0.01250000	150.00000000
11	0.01260331	0.00010331	0.01126972
12	0.01111111	-0.00138889	2.22222222
13	0.01257396	0.00007396	0.00682749
14	0.01275510	0.00025510	0.08746356
15	0.01388889	0.00138889	2.77777778
16	0.01250000	-0.00000000	0.00000000
17	0.01254325	0.00004325	0.00305312
18	0.01265432	0.00015432	0.04115226
19	0.01253463	0.00003463	0.00218691
20	0.00000000	-0.01250000	300.00000000
21	0.01252834	0.00002834	0.00161970
22	0.01260331	0.00010331	0.02253944
23	0.01252363	0.00002363	0.00123284
24	0.01250000	-0.00000000	0.00000000
25	0.01300000	0.00050000	0.60000000
26	0.01257396	0.00007396	0.01365498
27	0.01251715	0.00001715	0.00076208
28	0.01224490	-0.00025510	0.17492711
29	0.01251486	0.00001486	0.00061503
30	0.01388889	0.00138889	5.55555556
31	0.01251301	0.00001301	0.00050351
32	0.01250000	-0.00000000	0.00000000
33	0.01251148	0.00001148	0.00041740
34	0.01254325	0.00004325	0.00610625
35	0.01275510	0.00025510	0.21865889
36	0.01234568	-0.00015432	0.08230453
37	0.01250913	0.00000913	0.00029613
38	0.01253463	0.00003463	0.00437382
39	0.01250822	0.00000822	0.00025287
40	0.01250000	-0.00000000	0.00000000
41	0.01250744	0.00000744	0.00021764
42	0.01252834	0.00002834	0.00323939
43	0.01250676	0.00000676	0.00018866
44	0.01239669	-0.00010331	0.04507889
45	0.01265432	0.00015432	0.10288066
46	0.01252363	0.00002363	0.00246569
47	0.01250566	0.00000566	0.00014448
48	0.01250000	-0.00000000	0.00000000
49	0.01250521	0.00000521	0.00012750
50	0.01300000	0.00050000	1.20000000

5. COMPORTAMENTO ASSINTÓTICO DA RELAÇÃO DE PRECISÃO

Os resultados da tabela 1, assim como dos demais casos estudados, mostraram que a relação de precisão tem seus picos máximos quando o número de observações por corrida é múltiplo de 20, ou seja,

$$N = 20 \times K, \quad K = 1, 2, \dots$$

e, portanto, quando o tamanho da amostra de entrada é múltiplo de 100, ou seja,

$$n = 100 \times K, \quad K = 1, 2, \dots$$

Usando a amostragem descritiva, K ($K \geq 1$) corresponde ao número de pontos equidistribuídos tomados em cada subintervalo centesimal. Representando a situação mais desfavorável para a relação de precisão, serão estudados estes casos ...

Em função das simetrias observadas, os resultados globais podem ser deduzidos do estudo de um único subintervalo. Tomando o primeiro subintervalo como referência, tem-se

$$\overline{ZD} = 5 \sum_{j=1}^K |XD_j - q| / K$$

onde

$$XD_j = (j - 0.5) / n, \quad j = 1, \dots, K \quad (K = n/100 = N/20)$$

e

$q = 1/200$, que corresponde ao ponto médio deste primeiro subintervalo.

Dois casos são considerados para o cálculo de \overline{ZD} , dependendo de K ser ímpar ou par.

Caso em que K é ímpar:

Em função da simetria, tem-se que

$$\sum_{j=1}^K |XD_j - q| = \frac{(K-1)/2}{2} \sum_{j=1}^{(K-1)/2} [(j - 0.5)/(100K) - 1/200]$$

donde,

$$\sum_{j=1}^K |XD_j - q| = \frac{(K-1)/2}{2} \sum_{j=1}^{(K-1)/2} [(K - 2j + 1)/(200K)]$$

Após as devidas simplificações chega-se a

$$\sum_{j=1}^K |XD_j - q| = (K^2 - 1) / (400 K)$$

e, portanto,

$$\bar{ZD} = 5 \sum_{j=1}^K |XD_j - q| / K = (1/80) (1 - 1/K^2) .$$

Segue-se pois que

$$\text{Viés}(\bar{ZD}) = \bar{ZD} - E(Z) = -1 / (80 K^2)$$

e que a relação de precisão

$$\text{EMQ}(\bar{ZD}) / \text{Var}(Z) = 300 / K^2 .$$

Caso em que K é par

Tem-se que

$$\sum_{j=1}^K |XD_j - q| = 2 \sum_{j=1}^{K/2} [(j - 0.5) / (100 K) - 1/200]$$

donde,

$$\sum_{j=1}^K |XD_j - q| = 2 \sum_{j=1}^{K/2} [(K - 2j + 1) / (200 K)] .$$

Após os cálculos, segue-se que

$$\sum_{j=1}^K |XD_j - q| = K / 400$$

e, portanto,

$$\bar{ZD} = 1/80 .$$

Logo, como neste caso

$$\bar{ZD} = E(Z) = 1/80 ,$$

a relação de precisão é nula, ou seja, a amostragem descritiva leva a um valor exato para a estimativa em estudo.

Como verificação dos resultados anteriores, a tabela 2 apresenta os valores de \bar{ZD} e da relação de precisão quando N é múltiplo de 20.

Tabela 2. Valores de \bar{ZD} , do viés e da relação de precisão quando o número de observações por corrida é múltiplo de 20.

N	\bar{ZD}	Viés(\bar{ZD})	EMQ(\bar{ZD})/Var(\bar{Z})
20	0.00000000	-0.01250000	300.00000000
40	0.01250000	-0.00000000	0.00000000
60	0.01111111	-0.00138889	11.11111112
80	0.01250000	-0.00000000	0.00000000
100	0.01200000	-0.00050000	2.40000001
120	0.01250000	-0.00000000	0.00000000
140	0.01224490	-0.00025510	0.87463557
160	0.01250000	-0.00000000	0.00000000
180	0.01234568	-0.00015432	0.41152264
200	0.01250000	-0.00000000	0.00000000
220	0.01239669	-0.00010331	0.22539445
240	0.01250000	-0.00000000	0.00000000
260	0.01242604	-0.00007396	0.13654985
280	0.01250000	-0.00000000	0.00000000
300	0.01244444	-0.00005556	0.08888889
320	0.01250000	-0.00000000	0.00000000
340	0.01245675	-0.00004325	0.06106249
360	0.01250000	-0.00000000	0.00000000
380	0.01246537	-0.00003463	0.04373816
400	0.01250000	-0.00000000	0.00000000

Interpretação

Observando o comportamento da relação de precisão na pior situação, ou seja, quando

$$N = 20 \times K, \quad K = 1, 3, 5, \dots$$

nota-se que ela decresce rapidamente com K (ordem cúbica) e que, já a partir de $K = 7$, ela se torna favorável à amostragem descritiva. Assim, uma vez eliminadas as distorções decorrentes de uma amostra de entrada subdimensionada, a amostragem descritiva apresenta propriedades assintóticas que fazem dela sempre mais precisa que a amostragem aleatória simples.

6. CONSIDERAÇÕES FINAIS

Embora possam ocorrer situações em que a amostragem descritiva seja menos eficiente em termos estatísticos do que a amostragem aleatória simples numa simulação por Monte Carlo, este fato não serve como refutação à teoria da amostragem descritiva. Nossa atenção, para efeito de verificação desta teoria, deve se ater ao comportamento assintótico da relação de precisão, dentro da suposição de que a corrida foi bem dimensionada.

Embora esta consideração quanto ao dimensionamento da corrida pareça ser um artifício encontrado para "salvaguardar" a teoria da amostragem descritiva, deve ser lembrado que todo estudo de simulação subentende uma premissa de convergência estatística. Segundo esta premissa, uma estimativa de simulação deveria sempre ser uma boa aproximação do valor que seria obtido caso a amostra de entrada crescesse indefinidamente. Assim sendo, a comparação de eficiência entre ambos os métodos deve também subentender esta condição de convergência; em outras palavras, não importa o que venha a ocorrer com a relação de precisão entre ambos os métodos para um particular tamanho de amostra, mas sim o seu comportamento limite quando a amostra de entrada cresce indefinidamente.

Antes de se atingir esta condição de convergência existe a possibilidade da amostragem descritiva ser menos eficiente que a amostragem aleatória simples, como de fato ocorreu no exemplo aqui estudado; no entanto, uma análise mais cuidadosa deste "contra-exemplo" mostra que, aumentando-se o tamanho da amostra de entrada, esta relação de eficiência reverte-se totalmente a favor da amostragem descritiva. Este fato também mostra a necessidade de se dimensionar corretamente uma corrida de simulação quando do uso da amostragem descritiva; esta tarefa, no entanto, não traz maiores dificuldades práticas.

De fato, fazendo corridas piloto com amostragem descritiva bastaria que se aumentasse a duração da corrida até que a média das estimativas se estabilizasse ou apresentasse uma variação pequena quando comparada com sua variância. O trabalho adicional requerido por este procedimento seria mínimo, uma vez que um único programa de simulação seria utilizado tanto na fase piloto como no experimento definitivo.

BIBLIOGRAFIA

- FUNCHAL, G. Aplicação da amostragem descritiva na simulação de sistemas não elementares. Rio de Janeiro, IME, 1987. Tese de Mestrado.
- SALIBY, E. Repensando a simulação: a amostragem descritiva. Rio de Janeiro, Atlas/EDUFRJ, 1988.
- . A reappraisal of some simulation fundamentals. Lancaster, Universidade de Lancaster, 1980. Tese de Doutorado.

APÊNDICE A: Programa em TURBO-PASCAL para a simulação do problema teste, utilizando a amostragem aleatória simples, e um conjunto de resultados para $M = 50$ corridas.

```
PROGRAM DESVIO;
```

```
( Estudo do "contra-exemplo" utilizando a amostragem
  aleatoria simples. )
```

```
Type
```

```
  Estatistica = Array[1..3] of real;
```

```
Var
```

```
  Nomearq      : String[12];
  F            : Text;
  I, J, IC, N, M : Integer;
  X, Y, Valor, Z, EZ, VZ : Real;
  S            : Estatistica;
```

```
( Nomearq      = Nome do arquivo de saida
  F            = Variavel auxiliar definindo o tipo de arquivo de saida
  I            = Contador de elementos que compoem uma observacao
  J            = Contador de observacoes numa corrida
  IC           = Contador de corridas num experimento
  N            = Total de observacoes por corrida
  M            = Total de corridas por experimento
  X            = Ponto gerado no intervalo unitario (x100)
  Y            = Distancia do ponto ao centro do sub-intervalo centesimal a
                que pertence
  Valor        = Valor computado para uma observacao
  Z            = Acumulador das observacoes numa corrida e sua media
  EZ           = Valor esperado (teorico) para a estimativa em estudo
  VZ           = Variancia (teorica) para a estimativa em estudo, usando
                amostragem aleatoria simples
  S            = Acumula e resume as estatisticas de cada corrida;
                S[1] : Numero de observacoes
                S[2] : Soma das observacoes e, ao final, sua media
                S[3] : Soma dos quadrados e, ao final, sua variancia )
```

```
Procedure Zera; ( Zera acumuladores para estatisticas )
```

```
Begin
```

```
  S[1] := 0;
  S[2] := 0;
  S[3] := 0;
```

```
End;
```

```
Procedure Calcula; ( Calcula sumario estatistico )
```

```
Begin
```

```
  S[2] := S[2]/S[1];
  S[3] := (S[3] - S[1]*SQR(S[2]))/(S[1] - 1);
```

```
End;
```

```

Procedure Observacao;
( Calcula um valor para a variavel de resposta "Valor" )
Begin
  Valor := 0;
  For I := 1 to 5 do
  Begin
    X := 100*Random;
    Y := 0.01*(X - Trunc(X) - 0.5);
    Valor := Valor + Abs(Y);
  End;
End;

Procedure Corrida;
( Calcula uma estimativa da media (Z) a partir de N observacoes )
Begin
  Z := 0;
  For J := 1 to N do
  Begin
    Observacao;
    Z := Z + Valor;
  End;
  Z := Z/N;
End;

Begin ( Programa principal )
  Assign(F, 'Desvio.rel');
  Rewrite(F);
  N := 20;
  M := 50;
  EZ := 0.0125;
  VZ := 5.0/(200.0*200.0*12*N);
  Zera;
  Writeln(F, 'Resumo do experimento');
  Writeln(F, ' ( M =', M:3, ' Corridas', ' ' :8,
    'N =', N:3, ' Observacoes/corrida )');
  Writeln(F);
  Writeln(F, 'Corrida', ' ' :14, 'Media');
  Writeln(F);
  For IC:= 1 to M do
  Begin
    Corrida;
    Writeln(F, IC:4, ' ' :10, Z:14:8);
    S[1] := S[1] + 1;
    S[2] := S[2] + Z;
    S[3] := S[3] + SQR(Z);
  End;
  Calcula;
  Writeln(F);
  Writeln(F, 'Media', ' ' :9, S[2]:14:8, ' ' :9, 'Media/EZ =', S[2]/EZ:12:8);
  Writeln(F, 'Variância', ' ' :5, S[3]:16:10, ' ' :7, 'Var/VZ =', S[3]/VZ:12:8);
  Writeln(F);
  Close(F);
End.

```

Resumo do experimento
(M = 50 Corridas

N = 20 Observacoes/corrida)

Corrida	Media
1	0.01296782
2	0.01300325
3	0.01292436
4	0.01213212
5	0.01224703
6	0.01173694
7	0.01303181
8	0.01233772
9	0.01165108
10	0.01136738
11	0.01328494
12	0.01163544
13	0.01359086
14	0.01293147
15	0.01285106
16	0.01301777
17	0.01307571
18	0.01307327
19	0.01263832
20	0.01159283
21	0.01304014
22	0.01219887
23	0.01304127
24	0.01225603
25	0.01250921
26	0.01258453
27	0.01168729
28	0.01191494
29	0.01275380
30	0.01317439
31	0.01338206
32	0.01302597
33	0.01243745
34	0.01243175
35	0.01204124
36	0.01284508
37	0.01325612
38	0.01175601
39	0.01241781
40	0.01220866
41	0.01365345
42	0.01331857
43	0.01376526
44	0.01271161
45	0.01130323
46	0.01287229
47	0.01259128
48	0.01248377
49	0.01213313
50	0.01349906

Media 0.01260771
 Variancia 0.0000003884

Media/EZ = 1.00861676
 Var/VZ = 0.74573508

