

Experimentální analýza algoritmů dnes a dnes

1) Typický příklad z roku 1980:

W. Dobosiewicz: An efficient variation of bubble sort

(Information Processing Letters 11 (1980), 5-6)

prezentuje nový algoritmus

teoretickou složitost vyjádřil - ta ostatní není vyčíslena dočtem

jediný výsledek - B. Breyer: Analyzing variants of Shellsort

(Inform. Proc. Letters 49 (2001), 223-224)

dolní odhad oscilované složitosti je $\Omega\left(\frac{n^2}{2^p}\right)$ pro $p \leq \log n - \log \log n$

(pro větší p platí tak, ale to je triviální)

porovnává dobu výpočtu s klasickým Bubblesortem, Shellsortem

a Quicksortem

když se náhodně postupnosti reálných čísel generované a komparace rozdělují

délky postupnosti $n = 10, 50, 100, 500, 1000, 2000, 10000$

čas se měří v milisekundách,

programy byly naprogramovány ve Fortranu a spustily na CDC Cyber 72

výsledky (jako uvedeny v tabulce): nový algoritmus je výrazně rychlejší než klasický Bubblesort i Shellsort a jen o málo pomalejší než Quicksort

I když experimentální analýza může provádět naměřením (když je představení nového algoritmu), přesto tomu něco chybí:

- a) mēre se pausi cās, nīkati kaulichē charakteristig tūdenē
(proci paromānē, proci rjmin apod.)
ēas mēre tjt nīcē ēi minē skresen konkrētīm gēacēm
systēmēm
vaxta na kaulichas analjzu je problematiska
- b) chjpe' komentārē h' ūdajēm o tabulce - jau h' pūmēnē
hodnoty a x kolika mērenē?
- c) chjpe' i h' nījēlementārjgē' dāstlībā' analjza (minimum,
maximum, mediān, smēdātānā' cōstjcha, ...)
- d) chjpe' grafichē' anānomēnē pūblichē' māmērēnēd cāsū
- e) tūdēd se nēlmi bālchē' poslojprōstki (h' mohlō tjt dānō
mōiņostmē' kēkdjstēd' pōcīlācū)
- f) chjpe' interpretāce rjstādū, nāpī. tē māj' algoritmus
je rjchljgē' nē' Quicksort pō bālchē' poslojprōstki dō
1000 pōkū, ocl 2000 je rjchljgē' Quicksort
- g) nējseu formulorāj rādōnē' h'jpotēz, nāpī. tē pōdēc
māmērēnēd cāsū je māj' algoritmus māj pūblichē
2x rjchljgē' nē' Shellsort - mohl tjt mīl kēd' chjpsau
asymptotischou stādōst (tjot se jēm o multiplikātem
konstantē'?)

Pīdānōstki h'ls pūblichas : skūcōnē' pīdclāsenē' mōiņs algoritmus
impīrējgē' māmērē' pō dātē' rjstādū
abklādēnē' pūpīrēd - mōxi rāclānīm dō cāsōpīsu a rjclādām
uplējz 3 mōxē'!

AN EFFICIENT VARIATION OF BUBBLE SORT

Włodzimierz DOBOSIEWICZ

Institute of Informatics, Warsaw University, 00-901 Warszawa, Poland

Received 15 May 1980

Algorithms, sorting

Bubble sort is well known for being one of the worst sorting algorithms. One way to improve its efficiency was described by Batcher [1], however it is so complex and the amount of bookkeeping so huge that this algorithm is not commonly used in single-processor programming.

Bubble sort can be improved in yet another way, which is similar to Shell's version of the insertion sort.

We select a sequence h_1, \dots, h_t , where $h_i > 1$ and $t = O(\log n)$, n being the number of sorted elements.

In pass i ($i \leq t$) the vector to be sorted is traversed from left to right and items distant h_i from each other are compared and exchanged if necessary. After t passes a regular bubble sort should be used to complete the sorting process. This gives the following algorithm:

```

for  $\ell := 1$  to  $t$  do
begin
   $\text{inc} := h_\ell$ ;
  for  $i := 1$  to  $n - \text{inc}$  do
    if  $A[i] > A[i + \text{inc}]$ 
      then  $A[i] \leftrightarrow A[i + \text{inc}]$  fi
end;
 $\ell := n - 1$ ;
while  $\ell > 0$  do
begin
   $k := 0$ ;
  for  $i := 1$  to  $\ell$  do
    if  $A[i] > A[i + 1]$ 
      then  $\{A[i] \leftrightarrow A[i + 1]; k := i\}$ ;
   $\ell := k - 1$ 
end;
```

Obviously, this algorithm works, as the second part of it is just the well-known bubble sort. Its average efficiency is quite good, as illustrated by the Table 1.

The above results were obtained by sorting real numbers generated by a uniform distribution random numbers generator. The tested programs were written in Fortran and run on a CDC Cyber 72. The sequence $h_1 = \lfloor \frac{1}{2}n \rfloor$, $h_{i+1} = \lfloor \frac{3}{4}h_i \rfloor$ was used. The code of Quicksort was copied from [2], while Shellsort was coded after the description given in [3] with increments of form $2^k - 1$. While the version of Quicksort is approximately the best achievable, there is no certainty whether this is true in the case of Shellsort — this algorithm depends on the choice of increments (see [3] for details).

The exact time complexity of the presented algorithm has yet to be determined. The author hopes that this short presentation will encourage some one to work on this problem.

Table 1
Running times of various sorting algorithms (times in ms)

n	Quicksort	Shellsort	Bubblesort	Modified Bubblesort
10	0.86	0.99	0.6	0.58
50	5.84	8.62	14.6	4.62
100	13.38	21.28	57.3	11.87
500	87.81	153.70	1456.6	84.80
1000	199.8	367.0		194.9
2000	441.2	858.2		450.5
10000	2696.8			2743.8

2) Přiblížení o 10 let později:

(4)

M.A. Weiss: Empirical study of the expected running time of Shellsort
(Computer Journal 34(1991), 88-91)

Je to experimentální studie známých algoritmů (Shellsort, Hibbardův, Knuthův a Sedgewickův Shellsort) se známou složitostí a
nejhorším případem, ale neznámou oscilovanou složitostí (to
je dodnes nevyřešený problém).

Dále obsahuje: orientační porovnání podle času s Heapsortem a
standardním a optimalizovaným quicksortem
porovnání s předchozí experimentální studií Shellsortu

Algoritmy jsou detailně popsány, předpokládá se, že jsou známy.

Jsou uvedeny dosavadní teoretické výsledky o složitosti a výsledky
předchozí exp. studie.

teoreticky: Shell $O(n^2)$ nejhorší případ, $O(n^{3/2})$ oscilující
Hibbard, Knuth $O(n^{3/2})$ nejhorší případ
Sedgewick $O(n^{4/3})$ nejhorší případ

experimentálně: Hibbard, Knuth měří $O(n^{1,25})$ a $O(n^{1,28})$
nebo $O(n \log^2 n)$ oscilující případ

↑ experimentálně se měřil čas (pauze pro orientaci) a počet vyjímů
prvků (po účelové analýze).

nečetl se počet porovnání, což je poměrně chyba (viz Načevskij'sova')

Tudíž se náhodné permutace dělají 100 až 12 000 000 prvků.

1997

Prý použil UNIXový náhodný generátor RANDOM (RAND) se normální.

Testy probíhaly na počítačích SUN3 workstation s 8MB pamětí

a VAX 485 s 128 Mb.

Je uveden počet experimentů pro permutace různých délek (řádové desítky tisíc pro permutace délky do 10000, tisíce pro délky do 100000, stovky pro delší).

K naměřeným příkladům jsou uvedeny směřovací odchylky.

Prohlašují se: a) reálné algoritmy na posloupnostech stejných délek

b) hlavní rychlost řešení počtu výměn v závislosti na rozšiřující se velikosti souborů.

Naměřené hodnoty se porovnávají se staršími odhady

$$\Theta(n^{1,26}) \text{ a } \Theta(n \log^2 n)$$

a pokládají se nové markéty křivkami řádu

$$\Theta(n^{5/4}) \text{ pro Jubbardov a Knuthov Shellsort}$$

$$\Theta(n^{4/6}) \text{ pro Sedgewickov}$$

Výsledky porovnání různých aproximací jsou v tabulkách.

Chyby grafické slovním a slovním nýžnou statistickou

hodnotou, např. residuálním součtem čtverců (cga).

Na základě experimentů je vyložena hypotéza: Pro přibližně

posloupnosti délky $O(\log n)$ se složitost v nejhorším případě

řádu $\Theta(n^k)$ transformuje na očištěnou složitost řádu $\Theta(n^{\frac{k+1}{2}})$.

J. Incerpi, R. Sedgewick: Practical variations of Shellsort

(Information Processing Letters 26 (1987/88), 37-43

zobecnění Dobosiewiczova algoritmu (experimentální studie)

Weiss: Prezentace výsledků²

Table 1. Running times (in seconds) of various sorting algorithms

N	Shellsort				Quicksort		
	Shell's	Hibbard's	Knuth's	Sedgewick's	Heapsort	Standard	Optimised*
100	0.0024	0.0022	0.0021	0.0017	0.0042	0.0028	0.0024
1000	0.0354	0.0344	0.0306	0.0293	0.0557	0.0315	0.0259
10000	0.5890	0.5563	0.5000	0.4300	0.7165	0.3677	0.3153
100000	9.230	8.408	8.041	5.730	8.859	4.230	3.588
1000000	141.6	130.3	135.4	71.2	104.7	47.1	41.3

* Recursive with median-of-three partitioning and a cutoff of 10.

Table 2 *Hibbard*

N	Observed exchanges	Our fit	$\Theta(N^{1.26})$ fit	$\Theta(N \log^2 N)$ fit
100	347.6	360.2	400.7	-245.9
500	2950.0	2944.1	3044.4	291.1
1000	7200.8	7139.0	7291.0	2514.6
5000	55831.4	55871.1	55397.2	42232.7
10000	134658.1	134966.1	132673.9	116237.5
50000	1038169.2	1038814.1	1008062.8	1022314.1
100000	2500788.7	2497784.2	2414267.4	2486839.0
500000	19120679.6	19113501.4	18343726.9	18290748.7
1000000	45848881.9	45873300.0	43932444.6	42248509.9
12000000	1058423789.1	1052351112.6	1005897526.8	787733188.2

Table 3

Permutation sizes	Total number of sorts	Machines used
10i, $10 \leq i \leq 99$	16000	16 SUN 3 workstations
100i, $10 \leq i \leq 99$	16000	16 SUN 3 workstations
1000i, $10 \leq i \leq 45$	4000	16 SUN 3 workstations
100000i, $5 \leq i \leq 10$	2000	5 SUN 3 workstations with ≥ 8 Mb
12000000	100	1 VAX 785 with 128 Mb

Table 4

N	Hibbard's		Knuth's		Sedgewick's	
	Average	S.D.	Average	S.D.	Average	S.D.
100	347.6	26.3	432.1	34.3	461.8	37.9
1000	7200.8	361.8	8901.3	446.8	7725.2	192.0
10000	134658.1	6485.6	164960.3	7218.7	108811.7	943.0
100000	2500788.7	116348.3	2965566.6	1331956.0	1409404.5	5281.4
1000000	45848881.9	2.045907.5	53477706.7	2306321.8	17421118.3	36068.8

Table 6

Exchanges						
N	Hibbard's		Knuth's		Sedgewick's	
	Observed	Our fit	Observed	Our fit	Observed	Our fit
100	347.6	360.2	432.1	381.7	461.8	456.3
200	879.5	915.9	1094.9	1088.3	1110.4	1092.4
300	1499.3	1536.8	1863.2	1873.5	1817.1	1810.6
400	2184.5	2215.8	2725.1	2727.8	2598.2	2577.8
500	2950.0	2944.1	3608.3	3640.6	3410.9	3380.3
600	3692.7	3715.1	4583.7	4603.9	4210.7	4210.3
700	4516.7	4523.6	5606.5	5611.6	5054.9	5063.1
800	5370.8	5365.8	6640.8	6659.1	5933.0	5935.4
900	6264.6	6238.4	7729.2	7742.6	6847.4	6824.5
1000	7200.8	7139.0	8901.3	8859.3	7725.2	7772.8
2000	17480.6	17338.6	21395.7	21433.8	17401.2	17360.2
3000	29171.5	29166.6	36035.1	35865.6	27632.1	27700.0
4000	42300.4	42036.6	51670.2	51636.0	38572.8	38503.0
5000	55831.4	55871.1	68128.9	68475.1	49447.3	49652.4
6000	70411.6	70477.2	85966.8	86213.6	61008.9	61079.1
7000	85667.8	85755.5	104575.3	104734.0	72748.0	72737.1
8000	102007.0	101632.5	123897.5	123949.5	84929.4	84593.5
9000	117931.4	118051.2	143972.7	143793.2	96961.9	96624.0
10000	134658.1	134966.1	164960.3	164211.2	108811.8	108809.2
20000	324773.0	325331.3	392459.0	392865.7	236327.7	236779.2
30000	543129.3	543925.2	654918.0	653937.4	327777.2	327092.5
40000	781219.3	783095.2	934984.3	938534.9	512798.8	512196.8
50000	1038169.2	1038814.1	1239432.9	1242001.8	655480.2	655873.8
60000	1309707.4	1308534.4	1562018.2	1561402.6	803536.0	802404.1
70000	1589450.2	1590476.7	1895389.6	1894691.7	951420.6	951311.4
80000	1882107.7	1883317.7	2240797.3	2240350.7	1100403.4	1102256.9
90000	2183821.4	2186028.8	2606280.3	2597201.2	1252909.2	1254986.8
100000	2500788.7	2497784.2	2963566.6	2964298.7	1409404.5	1409304.2
200000	6000092.0	6002478.5	7092600.7	7072978.0	3013151.4	3015182.2
300000	10023436.6	10022280.4	11810720.0	11762684.6	4690850.4	4696427.9
400000	14429076.3	14417206.3	16910384.1	16874618.7	6425108.0	6427022.2
500000	19120679.6	19113501.4	22408676.5	22325341.3	8209540.0	8194562.7
600000	24056494.9	24064514.9	28156209.6	28062164.4	10015607.2	9991722.2
700000	29236932.2	29237678.8	34206881.1	34048349.0	11811384.4	11813644.9
800000	34603915.5	34608781.6	40458926.4	40256606.9	13649817.6	13656870.6
900000	40198106.3	40159017.8	46810624.1	46665758.2	15521352.5	15518805.5
1000000	45848881.9	45873300.0	53477706.7	52258820.7	17421118.3	17397431.6
12000000	1058423789.1	1052351112.6	1209048063.7	1201162111.6	254835978.7	254636241.3

by using a poor random number generator or time lost due to crashes. By running these sorts simultaneously, only about four weeks of actual calendar time was used.

3. Increments which give $O(N^3)$ worst-case bounds

In this section we provide new fits for the running time of Shellsort using the increment sequences suggested by Hibbard and Knuth. The main measure of the running time is the number of exchanges performed by the algorithm. In both cases, the number of exchanges fits $c_1 N^3 + c_2 N^2 + c_3 N + c_4 N^{1/2} + c_5 N^{1/3} + c_6$. In addition to this, there is additional $O(N \log N)$ work corresponding to comparisons which do not require exchanges (one comparison per element per increment). This gives a total running time of $O(N^3)$. Table 6 shows how well our fit compares with the observation.

The number of exchanges observed has a high standard deviation for both of the sequences in this section; in particular, for the sequence recommended by Knuth, the standard deviation is about 5% of the observed value (the claim in Ref. 8 of a standard

deviation of 50 000 for sorting 250 000 numbers is apparently a typo that should read 350 000).

All of the formulas represent simple unweighted least-mean-square fits, and it is certainly possible that some other technique could produce even better fits. The fact that a simple method produces such good fits over such a wide range shows how likely the functional form is. It seems likely that the error in the fit was larger than the probable error in the observation. The actual constants are of little practical interest, probably insignificant and inaccurate, but for the sake of completeness we include them here. For Hibbard's sequence, the fit obtained is

$$1.55376872N^3 - 4.47754N^2 + 47.98950721N - 335.828N^{1/2} + 1140.81404N^{1/3} - 1450.25$$

This is in contrast to the fits $1.21N^{1.26}$ and $0.39N \log N - 2.33N \log N$ previously proposed. For Knuth's sequence, the fit is

$$1.73123269N^3 - 2.28388N^2 + 33.14353571N^{1/2} - 284.216N^{1/3} + 1074.63905N^{1/4} - 1541.59,$$

which contrasts with Knuth's conjecture of

$1.66N^{1.25}$ and a logarithmic form that is not plausible.

4. Sedgewick's increments ($h_i = 9 \cdot 4^i - 9 \cdot 2^i + 1 \cup 4^i - 3 \cdot 2^i + 1$)

These increments exhibit a very low standard deviation which enables us to get fits with much less error than in the previous section. Table 4 compares the average number and standard deviation (s.d.) of exchanges for Hibbard's, Knuth's and Sedgewick's increment sequences on a few file sizes.

Our fit for this sequence is $0.42663452N^3 + 18.49281148N^2 - 61.5728N + 72.69723933N^{1/2} + 105.37474557N^{1/3} - 372.755N^{1/4} + 483.290557LS10$ 115. The $O(N^3)$ term was too insignificant to be reported. The data in tables 5 and 6 shows that this is a very accurate fit. Attempts to fit using other functional forms (including the form in Section 3) give generally poor results.

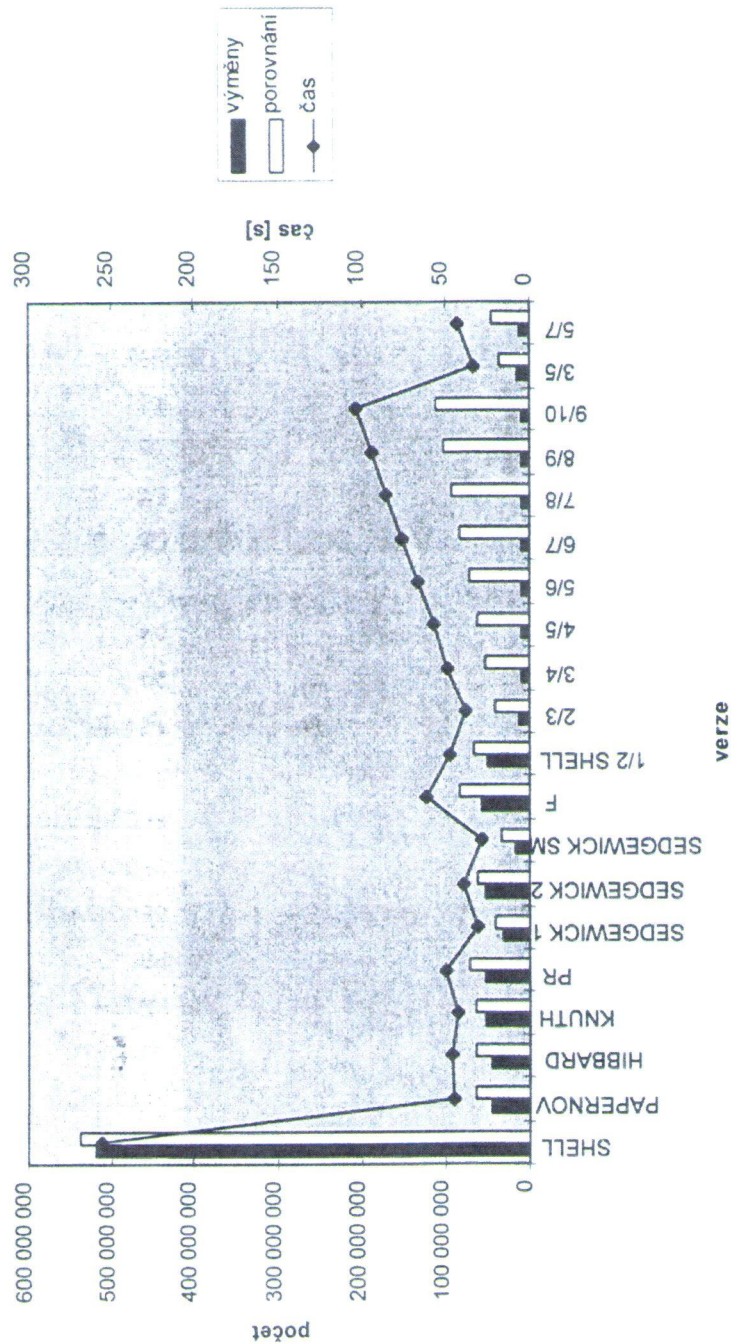
5. Extensions, conjectures and open problems

The natural open problem is to prove any of

Graf 19 - Porovnání měření pro různé verze - permutace délky 1 000 000

Graf zobrazuje průměr naměřených hodnot. V grafu vidíme poměr počtu výměn prvků vůči počtu porovnání prvků a času. U verze Sedgewick 2 bylo vysoké procento prvků po porovnání vyměněno, naopak verze 8/9 a 9/10 prováděly zbytečně moc nevyužitých porovnání (tj. porovnání, po kterých nenásledovala výměna). Naměřený čas koresponduje s počtem porovnání.

Naměřené výsledky pro permutace délky 1 000 000



3) Podstatná skúška 10 letoch:

S. Kanke: The performance of concurrent red-black tree algorithms
(TR 115, Institut für Informatik, Univ. Freiburg, 1998, 45 stran)

Rossahla experimentálne štúdie namierené na najemnejšom pomerní

3 typy tzv. relaxovaných červených-čiernych stromov a klasických
červených-čiernych stromov.

Relaxované vyvážený červených-čiernych strom je binárny vyhľadávací strom

s poněkud slabšími podmienkami na vyváženosť, kvôli čomu s ním
súvisí parametre stromu (a tým i rýchlosť vyhľadávania, vkládania
atď.), ale kompenzuje to menším počtom vyvážacích operácií,
ktoré majú podobu jier a nepriamo súvisia s veľkosťou stromu
a množstvom výpočtu. To je výhodou pri výpočte nákladov a paralelnom
prístupe k dátam - nemusí sa čakať, až sa po každej jednotlivej
operácii strom klasicky vyváži, čo by znamenalo nárast
jeho veľkosti. Štúdiu majú uchovávať ďalšie parametre,
pomocou ktorých je možné v dobe slabšieho prívodu stromu dorovnať
na klasický červených-čiernych strom.

Štúdie obsahujú: úvod do problematiky a reserši desiaradných výskumov

definície daných stromov (standardný červených-čiernych strom

+ algoritmy pre 3 typy relaxovaných vyvážaní)

implementačné detaily (kričov obsluhy vyvážacích príkazov,
paralelný prístup - ramybaš)

popis simulácie paralelného postupu (1, 3, 7, 15, 31 a 63 uzlových
procesov + 1 proces na poradiu pre vyvážacie operácie) (CREW-MIMD)

se sdílenou pamětí a s jednotkovou cenou operací,

optimalizace pomocí semaforů a naměřené

V experimentech bylo provedeno máj 1000 000 slovníkových operací

(search, insert, delete) paralelně na původně vyváženém cívném -
cívém slovní s 1000 000 klíči (vyváženém postupným přidáváním
náhodně vybraných klíčů do původně prázdného slovníku).

předpokládané operace: 0,5 po insert

0,3 po search

0,2 po delete

Experimenty se opakovaly s různými počty procesů.

Byla volena cena naměřených operací 1, 10, 20, 30, 40, 50, 60 jednotek
časů.

Dále následuje obsáhlý popis různých typů experimentů (co se měřilo,
v závislosti na čem, s čím se porovnává) se zprůstřed
dvou- a tříosnými grafy. Některé tabulky naměřených hodnot
se neuvádějí.

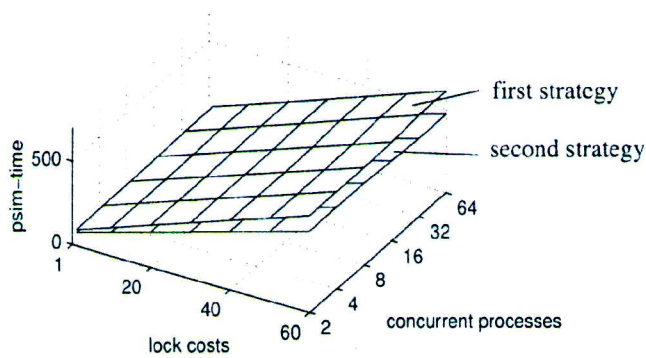
Komentářem jsem slovně popisuje to, co je na obtížné, uvádí typy

Proč to tak dopadlo, dalo se to čekat nebo je to překvapující?

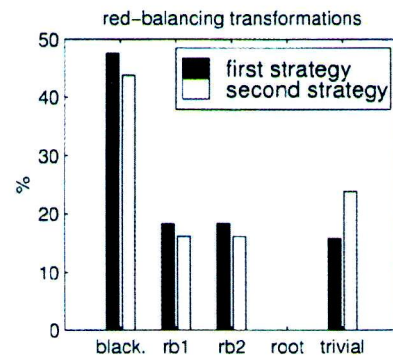
v podstatě chybí, statistické zpracování výsledků tabulí.

Návěr je velmi stručný se smyslem, že relaxované zpracování s
konkurenčním prostředím se osvědčilo.

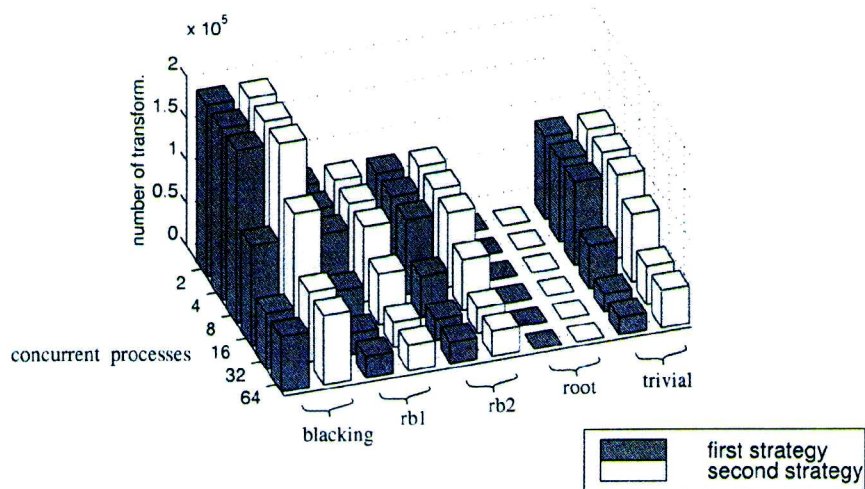
(a)



(b)



(c)



(d)

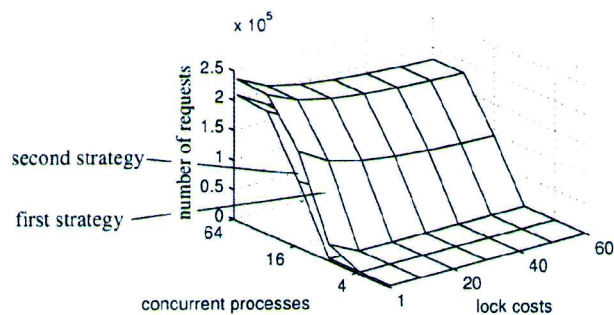


Figure 28: Comparison of the two strategies to handle red-red conflicts, exemplary shown by the chromatic tree: (a) Average time needed to perform a red-balancing transformation. (b) Proportion of the performed red-balancing transformations (using 63 user processes and 1 rebalancing process). (c) Total number of performed red-balancing transformations. (d) Number of unsettled rebalancing requests after the time used to perform 1000000 dictionary operations.

4) O novém lisicileti:

Ch. Baret et al.: Statistical analysis of algorithms: A case study of market-clearing mechanisms in the power industry

(Journal of Graph Algorithms and Applications 7(2003), 3-31)

Jedna ze (kalim) mála studií, kde se výsledky z experimentů statisticky zpracovávají.

Rěší se publicz' problém souvisící s deregulací elektrické pumpy v USA, řešení vztahu mezi producenty a odběrateli v podmínkách, kdy všichni sdílejí jednotnou rozvodnou síť s omezenou kapacitou. Protože je to komplikované, byla se sítěmi státu vytvořena instituce nezávislého systémového gestátora či arbitra, který má rozhodovat, které konkrétní se povolí a které ne - za tím účelem byla stanovena nějaká kritéria. Teoreticky se pak řeší jazyší grafový problém (kdy v síti s omezeními podmínkami?).

Testovaly se 4 algoritmy, a také 3 heuristiky.

Jako síť se zvolila energetická síť státu Poloaoto:

zdroje - elektrárny, jejich umístění a kapacita

potřebiči - odhadly se na základě demografie

kapacitý hran - odhadly se z kapacit zdrojů následovně:

elektrická kapacita potřebičů se rovná celkové kapacitě zdrojů a je rozdělena mezi potřebiče podle počtu obyvatelstva

kapacitý hran se učí simulací - síť se pomocí algoritmy na výpočet max. toku

konflikty se generují podle různých scénářů - jsou popsány 4
 experimentálně se po každém scénáři zkoušejí různé 4 algoritmy.

Byl použit software a formát dat projektu DIMACS.

Měla se dobrá rychlost a kvalita řešení, přičemž kvalita byla
 rozhodující a byla pro ni stanovena jasně vyjádřená míra.
 K testování rozdílů mezi jednotlivými scénáři a algoritmy byla
 použita analýza rozptylu.

Opací je poměrně důležitě popsána metoda analýzy rozptylu - má
 rozdíl od jiných studií, kde se prezentují pouze výsledky rozptylu
 má vyznačené konfliktu a nějaký statistický software jako
 k černé skříňce. Uživatel tak má šanci zjistit, jestli ta použitá
 metoda je vůbec pro jeho data vhodná.

Na ukázkou uvádím pouze část netriviální statistické analýzy

scenarios perform the best.¹

ANOVA has the following three advantages over individual t -tests² when the number of groups being compared is greater than two. See [9] for more details. In our case, we have four algorithms and four scenarios. Standard statistics terminology for a hypothesis that we wish to test, is *null hypothesis*.

- It gives accurate and known type-I error probability.³
- It is more powerful i.e. if null hypothesis is false, it is more likely to be rejected.
- It can assess the effects of two or more independent variables simultaneously.

5.3 Mathematical Model

Quality of Solution: We first describe the experiment for the quality of solution i.e. p_{AS} . We use a two-factor ANOVA model since our experiment involves two factors which are:

1. The algorithms: $\mathcal{A}_i, i = 1, 2, 3$ and 4.
2. The scenario: $\mathcal{S}_j, j = 1, 2, 3$ and 4.

Following classical statistics terminology, we will sometimes refer to algorithms as *treatments* and the scenarios as *blocks*. We will use \mathcal{A} to denote the set of algorithms and \mathcal{S} to denote the set of scenarios. For each algorithm-scenario pair we have *30 observations (or replicates)*. When testing the efficacy of the algorithms, we use 4 algorithms, each having 120 observations (30 for each scenario) from the corresponding population. The design of experiment used here is a *fixed-effect complete randomized block*. *Fixed-effect* because the factors are fixed as opposed to randomly drawn from a class of algorithms or scenarios; the conclusions drawn from this model will hold only for these particular algorithms and scenarios. *Complete* implies that the number of observations are the same for each block. *Randomized* refers to the 30 replicates being drawn randomly. We wish to test the hypothesis:

¹The populations in each of the groups are assumed to be normally distributed and have equal variances. The effect of violation of ANOVA assumptions of normality and homogeneity of variances have been tested in the literature ([10]) and the results show:

- Non-normality has negligible consequences on type-I and II error probabilities unless the populations are highly skewed or the sample is very small.
- When the design is balanced, i.e. the number of observations are the same for each group, violation of homogeneity of variance assumption has negligible consequences on the accuracy of type-I error probabilities.

² t -test checks for the significance of the difference in the means of *two* samples. It can assess whether the difference in sample means is just due to sampling error or they really are from two populations with different means.

³The probability of rejecting a null hypothesis when it is actually true.

Is the mean quality of solution provided by different algorithms the same, against the alternative hypothesis that some or all of these means are unequal?

The model for randomized block design includes constants for measuring the scenario effect (block effect), the algorithm effect (treatment effect) and a possible interaction between the scenarios and the algorithms. The appropriate mathematical model is as follows:

$$X_{ijk} = \mu + \tau_i + \beta_j + (\tau\beta)_{ij} + \varepsilon_{ijk},$$

where X_{ijk} is the measurement (p_{AS}) for the k th sample within the i th algorithm and the j th scenario. τ_i is the algorithm effect. β_j is the scenario effect. $(\tau\beta)_{ij}$ captures the interaction present between the algorithms and the scenarios. ε_{ijk} is the random error. See [8, 9] for further details on ANOVA.

We use S-Plus [15] software to run two-factor ANOVA to test the following three different null hypotheses.

1. Are the means given by the 4 different algorithms equal? The null hypothesis here is, $H_0 : \tau_1 = \tau_2 = \tau_3 = \tau_4$.
2. Are the means given by the 4 different scenarios equal? The null hypothesis here is, $H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4$.
3. Is there any interaction between the two factors? The null hypothesis here is, $H_0 : (\tau\beta)_{ij} = 0$.

The results of two-factor ANOVA are shown in Table 1 and Table 2. In the following discussion, we explain the meaning of each column in Table 1. DF refers to the degrees of freedom, SS refers to the sum of squared deviations from the mean. MS refers to the mean square error, which is the sum of squares divided by the degrees of freedom.⁴

⁴The sum of squares for the algorithm factor can be calculated as:

$$SS_A = nJ\sum_i (\bar{X}_{i..} - \bar{X}...)^2$$

where n is the number of replicates, J is the number of scenarios, $\bar{X}_{i..}$ is the mean of algorithm i across all scenarios and $\bar{X}...$ is the grand mean across all algorithms and scenarios. Recall that in our case $n = 30$ and $J = 4$ yielding a total sample size of 120.

The sum of squares for scenario factor can be calculated as:

$$SS_S = nI\sum_j (\bar{X}_{.j.} - \bar{X}...)^2$$

where as before n is the number of replicates, I is the number of algorithms and $\bar{X}_{.j.}$ is the mean of scenario j across all algorithms. Again, in our case $n = 30$ and $I = 4$.

The sum of squares for algorithms and scenario interaction is:

$$SS_{AS} = n\sum_j \sum_i [\bar{X}_{ij.} - (\bar{X}... + \hat{\tau}_i + \hat{\beta}_j)]^2$$

Here $\bar{X}_{ij.}$ is the mean of observations for the algorithm i scenario j pair. $\hat{\tau}_i$ and $\hat{\beta}_j$ are respectively the estimated least square values of τ_i and β_j . The sum of squares “within” refers to the squared difference between each observation and the mean of the scenario and algorithm of which it is a member. It is also referred as the residual sum of squares. This can be calculated as:

$$SS_W = n\sum_j \sum_i \sum_k (X_{ijk} - \bar{X}_{ij.})^2$$

The p -value gives the smallest level of significance at which the null hypothesis can be rejected.⁵ The lower the p -value, the lesser the agreement between the data and the null hypothesis. Finally the F -test is as follows. To test the null hypothesis, i.e., whether the population means are equal, ANOVA compares two estimates of σ^2 . The first estimate is based on the variability of each population mean around the grand mean. The second is based on the variability of the observations in each population around the mean of that population. If the null hypothesis is true, the two estimates of σ^2 should be essentially the same. Otherwise, if the populations have different means, the variability of the population mean around the grand mean will be much higher than the variability within the population. The null hypothesis in the F -test will be accepted if the two estimates of σ^2 are almost equal.

In a two-factor fixed-effect ANOVA, three separate F -tests are performed: two tests for the factors, and the third for the interaction term. The null hypothesis for the first factor can be written as:

$$H_0^A : \mu_{1..} = \mu_{2..} = \dots = \mu_{j..}$$

which is equivalent to writing: $H_0 : \tau_1 = \tau_2 = \tau_3 = \tau_4$. The F -test is:

$$F_A = \frac{SS_A/(I-1)}{SS_W/IJ(n-1)}$$

and the null hypothesis for the second factor can be written as:

$$H_0^S : \mu_{.1.} = \mu_{.2.} = \dots = \mu_{.j.}$$

which is equivalent to writing: $H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4$. The F -test is:

$$F_S = \frac{SS_S/(J-1)}{SS_W/IJ(n-1)}$$

and the null hypothesis for the interaction term can be written as:

$$H_0^{AS} : (\tau\beta)_{ij} = 0.$$

The F -test is:

$$F_{AS} = \frac{SS_{AS}/(I-1)(J-1)}{SS_W/IJ(n-1)}$$

If this F -ratio is close to 1, the null hypothesis is true. If it is considerably larger – implying that the variance between means is larger than the variance

The total sum of squares is

$$SS_T = SS_A + SS_S + SS_{AS} + SS_W$$

⁵To obtain a p -value for say F_A , the algorithm effect, we would look across the row associated with 3 degree of freedom in the numerator and 464 degrees of freedom in the denominator in the F -distribution table and find the largest value that is still less than the one obtained experimentally. From this value, we obtain a p -value of 0 for F_A .

Source	DF	SS	MS	F-test	p-value
Scenario (Block)	3	0.14	0.05	43.38	0
Algorithm (Treatment)	3	22.78	7.59	6792.60	0
Scenario:Algorithm	9	0.12	0.01	15.90	0
Residuals	464	0.40	.0008		
Total	479	23.45			

Table 1: **Results of Two-Factor ANOVA:** This table shows results of two-factor ANOVA where the factors are algorithms and scenarios. The measurement is the quality of solution, given by p_{AS} . The p -values show that the algorithm effect, scenario effect and the interaction between the algorithms and scenarios are all significant at any level of confidence.

within a population – the null hypothesis is rejected. The F distribution table should be checked to see if the F -ratio is significantly large.

The results in Table 1 show that all the above three null hypothesis are rejected at any significance level. This implies that the performance (measured by p_{AS}) of at least one of the algorithms is significantly different from the other algorithms. Also, different scenarios make a difference in the performance. Finally, the scenarios and the algorithms interact in a significant way. The interaction implies that the performance of the algorithms are different for different scenarios.

5.3.1 Contrasts

The next question of interest is what really caused the rejection of the null hypothesis; just knowing that at least one of the algorithms is different does not help us identify which algorithm is significantly different. To answer this we use a procedure called *contrast*. A contrast C among I population means (μ_i) is a linear combination of the form

$$C = \sum_i \alpha_i \mu_i = \alpha_1 \mu_1 + \alpha_2 \mu_2 + \cdots + \alpha_I \mu_I$$

such that the sum of contrast coefficients $\sum_i \alpha_i$ is zero. In the absence of true population means, we use the unbiased sample means which gives the estimated contrast as:

$$\hat{C} = \sum_i \alpha_i \bar{X}_i = \alpha_1 \bar{X}_1 + \alpha_2 \bar{X}_2 + \cdots + \alpha_I \bar{X}_I.$$

The contrast coefficients $\alpha_1, \alpha_2, \dots, \alpha_I$ are just positive and negative numbers that define the particular hypothesis to be tested. The null hypothesis states that the value of a parameter of interest for every contrast is zero, i.e., $H_0 : C =$