

# METODY KORELACE VELKÝCH OBJEMŮ DAT

## Modelování a simulace v elektrotechnice

Martin Zlámal

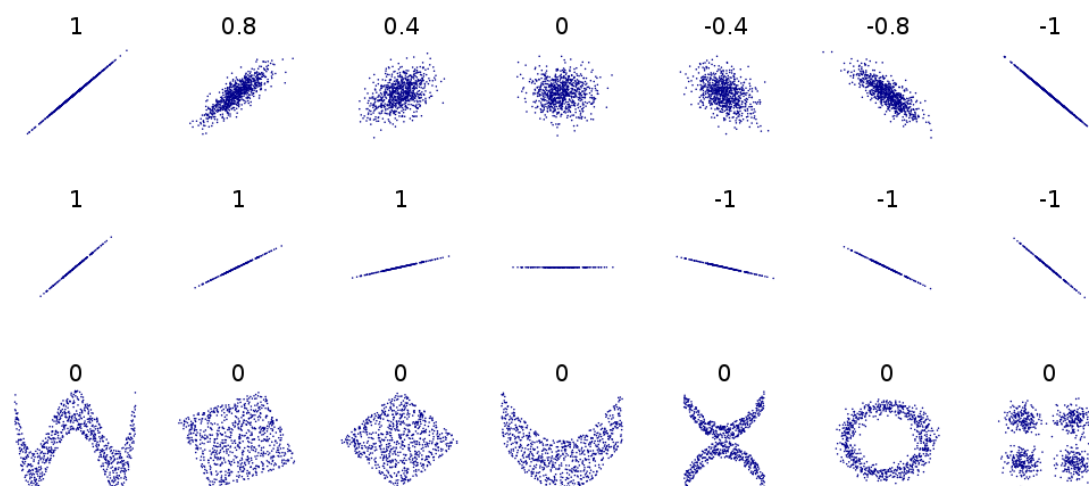
prosinec 2016

## Zadání

Implementujte alespoň dvě rozdílné metody korelace souborů číselných dat [Correlation and dependence, Wikipedia]. Můžete uvažovat například Pearsonova korelační koeficient a Spearmanův koeficient pořadové korelace. Pro vybrané metody implementujte algoritmus v jazyce MATLAB. Následně metody korelace porovnejte na ilustrativním souboru dat (využijte databáze příkladů na Courseware).

## O co jde?

Při práci s větším objemem dat může být složité najít souvislost mezi jedním zdrojem dat a druhým zdrojem dat. V takovém případě se vyplatí využít algoritmů, do kterých vstupují hodnoty  $X$  a s nimi (možná) související hodnoty  $Y$ . Tyto algoritmy následně vrací informaci o míře korelace dat. Pokud např. hodnota  $X$  vždy vzroste o jasně daný přírůstek a hodnota  $Y$  na tento přírůstek zareaguje vždy stejnou odezvou, můžeme bezpečně prohlásit, že mezi těmito soubory dat existuje vzájemný vztah.



Obrázek 1: Korelační koeficienty podle Pearsonova algoritmu

Na výše uvedeném obrázku je vidět přesně toto chování. Prostřední řádka ukazuje na perfektní korelaci - je jednoduché najít vztah mezi osou  $X$  a  $Y$ . U prvního řádku to již není vždy tak jednoznačné, nebo korelace vůbec neexistuje. Stejně tak neexistuje korelace v posledním řádku (nelze najít žádný jasný vztah mezi  $X$  a  $Y$ ).

## Popis korelačních algoritmů

V této práci jsou implementovány dva algoritmy pro zjišťování korelace dat: Pearsonův korelační koeficient a Spearmanův koeficient pořadové korelace. Pearsonův korelační koeficient je roven  $+1$  v případě ideální vzrůstající lineární závislosti dat a  $-1$  v případě klesající lineární závislosti (nepřímá úměra). V ostatních případech

se koeficient nachází v intervalu  $(-1; +1)$  a popisuje míru korelace (0 = data nekoreluje).

A protože každý kód, který nemá napsané testy, je spíše přání než-li skutečný kód, následují krátké unit testy, které testují výše uvedená pravidla a až poté následuje konkrétní implementace:

```
1 function testPearsonCorrelated(testCase)
2     A = [1 2 3 4];
3     B = [1 2 3 4]; % rho === 1.0
4     testCase.verifyEqual(pearson(Correlation, A, B), 1.0);
5 end
```

První parametr testovací funkce je aktuální hodnota, druhý parametr je potom očekávaná hodnota. Celá testovací sada lze spustit pomocí příkazu `run(CorrelationTest)`, jehož výstup by měl vypadat následovně (úspěšný průběh testů):

```
1 >> run(CorrelationTest)
2 Running CorrelationTest
3 .....
4 Done CorrelationTest
5 -----
6
7 Totals:
8     9 Passed, 0 Failed, 0 Incomplete.
9     0.14226 seconds testing time.
```

To je důkaz, že všechny testy fungují. Podobně vypadá test pro klesající závislost (stále Pearson):

```
1 function testPearsonCorrelated(testCase)
2     A = [1 2 3 4];
3     B = [4 3 2 1]; % rho === -1.0
4     testCase.verifyEqual(pearson(Correlation, A, B), -1.0);
5 end
```

Toto chování lze popsat matematickým vzorcem pro výpočet Pearsonova korelačního koeficientu:

$$\rho_{X,Y} = \text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \cdot \sigma_Y} \quad (1)$$

Kde  $\text{cov}(X, Y)$  je kovariance hodnot  $X$  a  $Y$  a  $\sigma$  značí směrodatnou odchylku. Samotný výpočet Pearsonova korelačního koeficientu lze naprogramovat takto:

```
1 function rho = pearson(obj, X, Y)
2     X = X(:); % vectorize
3     Y = Y(:);
4     covariance = cov(X, Y);
5     rho = covariance(1, 2) / (std(X) * std(Y));
```

```

6     if isnan(rho)
7         error('pearson:undefined', 'Correlation coefficient
8             is undefined because the variance of Y is zero.')
9     end
10 end

```

Prvně algoritmus vezme vstupující matice nebo vektory a vyrobí z nich jednu dlouhou posloupnost dat. Následně vypočítá Pearsonův koeficient korelace podle vzorce. Spearmanův koeficient pořadové korelace se chová velmi podobně pouze s tím rozdílem, že se hledá korelace nikoliv hodnot jako takových, ale jejich indexů (pořadí) v souboru dat. Implementace může vypadat třeba takto:

```

1 function rho = spearman(obj, X, Y)
2     X = X(:); % vectorize
3     Y = Y(:);
4     [~, rankX] = sort(X); % calculate Spearman ranking
5     [~, rankY] = sort(Y);
6     rho = obj.pearson(rankX, rankY);
7 end

```

Za povšimnutí stojí, že Spearmanův algoritmus využívá s výhodou Pearsonův (pouze vstupní hodnoty jsou indexy v původním souboru dat). Rozdíl mezi Pearsonovým a Spearmanovým algoritmem je dobře pochopitelný z tohoto testu:

```

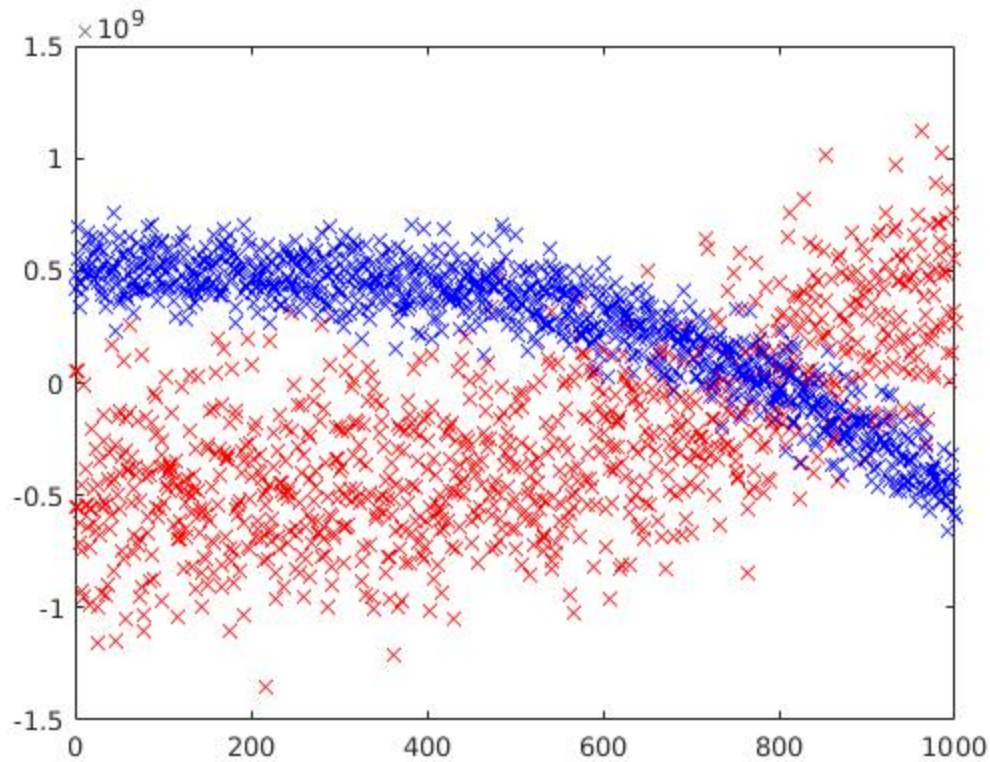
1 function testNonlinearCorrelation(testCase)
2     A = [0 10 101 102];
3     B = [1 100 500 2000];
4     testCase.verifyEqual(pearson(Correlation, A, B), 0.7544);
5     testCase.verifyEqual(spearman(Correlation, A, B), 1.0);
6 end

```

Zatímco Pearsonův koeficient je 0.7544, tak Spearmanův je 1.0. Proč? Pearson hledá korelaci vzhledem k lineárnímu charakteru souboru dat, Spearman však sleduje tendenci dat. Tzn. že pokud vzroste A, vzroste i B. Pokud opět vzroste A, opět vzroste B (nezáleží ani velikosti přírůstku). Z toho plyne, že data korelují a z hlediska Spearmanova algoritmu korelují perfektně.

## Korelace většího objemu dat

Ačkoliv bylo v zadání dáno, že mám použít vzorové matice, nevěděl jsem si s tím rady. Pro ověření algoritmů je totiž zapotřebí dvou odlišných souborů dat o stejné velikosti, mezi kterými se hledá závislost a proto jsou přiložené matice nevhodné. Vytvořil jsem proto třídu `GenerateDemoData`, která k tomuto účelu slouží. Tato třída umí podle předem dané křivky vygenerovat sadu bodů s tím, že body jsou rozptýleny kolem křivky podle normálního rozdělení s určitou variabilitou. Takto může vypadat průběh vygenerovaných hodnot pro  $y = x^3$  s vysokou variabilitou (červené hodnoty) a ten samý průběh pouze vynásobený  $-1$  s menší variabilitou (modré hodnoty):



Obrázek 2: Náhodně vygenerované hodnoty s variabilitou kolem daných křivek

V tomto případě vychází Pearsonův koeficient 0.6378 (pro vzrůstající hodnoty) a  $-0.8591$  pro klesající hodnoty. Spearmanův koeficient se v tomto případě téměř neliší. Z toho je vidět, že korelační koeficient se tím více blíží mezní hodnotě, čím jsou data méně rozptýlena (a tedy korelace dat je jasnější).

U velmi velkých matic je problém s paměťovou náročností. Algoritmus je napsán až příliš naivně a je schopen zabrat celou paměť. Přemýšlel jsem proto, jestli by bylo možné počítat Pearsonův korelační koeficient postupně. To však nejde. Proč ne? Důkaz mého tvrzení lze vyzkoušet takto:

```

1 Ax = [1 2 3 4 5 6 7 8 9 10];
2 Ay = [6 9 8 4 1 2 5 6 3 5];
3
4 Bx = [91 92 93 94 95 96 97 98 99 100];
5 By = [65 69 68 61 63 65 66 67 62 60];
6
7 c = Correlation(false);
8 c.pearson(Ax, Ay) % -0.4160
9 c.pearson(Bx, By) % -0.4973

```

Pokud však vezmeme v úvahu celý soubor dat (tedy A i B dohromady), výsledek je zcela jiný:

```

1 Cx = [1 2 3 4 5 6 7 8 9 10 91 92 93 94 95 96 97 98 99 100];
2 Cy = [6 9 8 4 1 2 5 6 3 5 65 69 68 61 63 65 66 67 62 60];
3 c.pearson(Cx, Cy) % 0.9915

```

Zatímco části souboru dat vycházejí záporně, celý soubor dat poukazuje (podle Pearsona) na velmi dobrou korelaci. Bylo by tedy zapotřebí zvolit méně naivní postup a vybírat podsoubory dat lépe (třeba z obou konců zároveň). Tuto otázku by bylo třeba řešit např. v případě paralelizace, to však není předmětem této práce.

## Závěrem

V této práci se řeší korelace dat, je však potřeba upozornit na to, že z korelace neplyne kauzalita. To že mezi soubory dat existuje určitá souvislost nutně neznamená (ale může znamenat), že změna jedné veličiny přímo způsobuje změnu druhé veličiny. Korelace je nutnou, nikoliv postačující podmínkou kauzality.

Dále je třeba si uvědomit, že výsledný korelační koeficient (zejména Pearsonův) nijak nevypovídá o rozložení testovaných bodů. Mohlo by se totiž zdát, že Pearsonův koeficient ukazuje míru linearitu bodů, to však není pravda. Důkazem je test `testAnscombesQuartet`, který ukazuje čtyři sady bodů, které mají od linearitě daleko a každý vypadá úplně jinak, přesto Pearsonův korelační koeficient vychází vždy stejně. Vizualní pohled na data může být tedy důležitý.

Z předchozího pozorování se tedy Spearmanův koeficient pořadové korelace jeví jako vhodnější protože není citlivý na hodnoty předloženého souboru dat, ale na jejich pořadí a tendenci. Nicméně jeho výpočet je náročnější, protože je potřeba zjistit pořadí prvků v souboru dat. Všechny naprogramované třídy jsou přiloženy na následujících stránkách.

## Implementace korelačních algoritmů

```
1 classdef Correlation
2     properties (SetAccess = private)
3         withGraphics;
4     end
5
6     methods
7         function obj = Correlation(withGraphics) % constructor
8             if nargin == 0
9                 withGraphics = false;
10            end
11            obj.withGraphics = withGraphics;
12        end
13
14        function rho = pearson(obj, X, Y)
15            [X, Y] = obj.vectorize(X, Y);
16            rho = obj.pearsonWithoutVectorize(X, Y);
17        end
18
19        function rho = spearman(obj, X, Y)
20            [X, Y] = obj.vectorize(X, Y);
21            [~, rankX] = sort(X); % calculate Spearman ranking
22            [~, rankY] = sort(Y);
23            rho = obj.pearsonWithoutVectorize(rankX, rankY);
24        end
25    end
26
27    methods (Access = private)
28        function [X, Y] = vectorize(obj, X, Y)
29            X = X(:); Y = Y(:);
30        end
31
32        function rho = pearsonWithoutVectorize(obj, X, Y)
33            covariance = cov(X, Y);
34            rho = covariance(1, 2) / (std(X) * std(Y));
35            if isnan(rho)
36                error('pearson:undefined', 'Correlation coefficient
37                    is undefined because the variance of Y is zero.')
38            end
39            if obj.withGraphics
40                plotmatrix(X, Y)
41            end
42        end
43    end
44 end
```

## Pomocná třída pro generování dat

```
1  classdef GenerateDemoData
2      %GENERATEDEMODATA Class for generation demo data
3      % This class is quite pointless with stupid API, but it generates
4      % random data with value deviation around math curve and that's
5      % pretty useful for correlation testing.
6
7      properties (SetAccess = private)
8          length;
9          deviationScale;
10     end
11
12     methods
13         function obj = GenerateDemoData(length, deviationScale)
14             if nargin == 0
15                 length = 1000;
16                 deviationScale = 1;
17             end
18             obj.length = length;
19             obj.deviationScale = deviationScale;
20         end
21
22         function [X, Y] = exponential(obj)
23             [X, Y] = obj.generateTable(@(x) x^3-0.5e9,
24                 obj.length * 3e5 * obj.deviationScale);
25         end
26         function [X, Y] = quadratic(obj)
27             [X, Y] = obj.generateTable(@(x) -x^3+0.5e9,
28                 obj.length * 1e5 * obj.deviationScale);
29         end
30     end
31
32     methods (Access = private)
33         function value = randn(obj, mu, sigma)
34             value = (randn * sigma) + mu;
35         end
36
37         function [X, Y] = generateTable(obj, func, deviation)
38             for index = 1:obj.length
39                 X(index) = index;
40                 Y(index) = obj.randn(func(index), deviation);
41             end
42         end
43     end
44
45 end
```



## Testy korelačních algoritmů

```
1  classdef CorrelationTest < matlab.unittest.TestCase
2      % run(CorrelationTest)
3
4      methods (Test)
5          function testPearsonCorrelated(testCase)
6              A = [1 2 3 4];
7              B = [1 2 3 4];
8              % actual, expected, ...
9              testCase.verifyEqual(pearson(Correlation, A, B), 1.0,
10                  'AbsTol', 1e-15);
11          end
12
13          function testPearsonAnticorrelated(testCase)
14              A = [1 2 3 4];
15              B = [4 3 2 1];
16              testCase.verifyEqual(pearson(Correlation, A, B), -1.0,
17                  'AbsTol', 1e-15);
18          end
19
20          function testPearsonWithoutCorrelation(testCase)
21              A = [3 4 5 6];
22              B = [1 1 1 1];
23              testCase.verifyError(@() pearson(Correlation, A, B),
24                  'pearson:undefined', 'Correlation coefficient is
25                  undefined because the variance of Y is zero.');
```

```
26          end
27
28          function testPearsonZeroCorrelation(testCase)
29              A = [1 2 1 2];
30              B = [1 1 2 2];
31              testCase.verifyEqual(pearson(Correlation, A, B), 0.0);
32              testCase.verifyEqual(spearman(Correlation, A, B), 0.8);
33          end
34
35          function testSpearmanHalfCorrelation(testCase)
36              % Spearman should be 0.5 but Kendall shoul be 1.0
37              A = [1 2 3];
38              B = [1 3 2];
39              testCase.verifyEqual(pearson(Correlation, A, B), 0.5);
40              testCase.verifyEqual(spearman(Correlation, A, B), 0.5);
41          end
42
43          function testNonlinearCorrelation(testCase)
44              A = [0 10 101 102];
45              B = [1 100 500 2000];
```

```

46         testCase.verifyEqual(pearson(Correlation, A, B), 0.7544,
47             'AbsTol', 1e-4);
48         testCase.verifyEqual(spearman(Correlation, A, B), 1.0,
49             'AbsTol', 1e-15);
50     end
51
52     function testUniformlyRandomCorrelation(testCase)
53         c = Correlation(false); % without graphics (default)
54         A = randi(100, 100); % uniformly distributed pseudorandom integers
55         B = randi(100, 100);
56         testCase.verifyEqual(pearson(c, A, B), 0.05, 'RelTol', 10);
57         testCase.verifyEqual(spearman(c, A, B), 0.05, 'RelTol', 10);
58     end
59
60     function testLargeDemoData(testCase)
61         g = GenerateDemoData(1000, 1);
62         [X, Y] = g.exponential();
63         [A, B] = g.quadratic();
64         c = Correlation(false);
65         testCase.verifyEqual(pearson(c, X, Y), 0.6, 'RelTol', 10);
66         testCase.verifyEqual(pearson(c, A, B), -0.8, 'RelTol', 10);
67     end
68
69     function testAnscombesQuartet(testCase)
70         % https://en.wikipedia.org/wiki/Anscombe%27s\_quartet
71         c = Correlation(false);
72         test = @(x, y) testCase.verifyEqual(pearson(c, x, y), 0.816,
73             'RelTol', 0.001);
74         X1 = [10 8 13 9 11 14 6 4 12 7 5];
75         Y1 = [8.04 6.95 7.58 8.81 8.33 9.96 7.24 4.26 10.84 4.82 5.68];
76         test(X1, Y1);
77         X2 = [10 8 13 9 11 14 6 4 12 7 5];
78         Y2 = [9.14 8.14 8.74 8.77 9.26 8.10 6.13 3.10 9.13 7.26 4.74];
79         test(X2, Y2);
80         X3 = [10 8 13 9 11 14 6 4 12 7 5];
81         Y3 = [7.46 6.77 12.74 7.11 7.81 8.84 6.08 5.39 8.15 6.42 5.73];
82         test(X3, Y3);
83         X4 = [8 8 8 8 8 8 8 19 8 8 8];
84         Y4 = [6.58 5.76 7.71 8.84 8.47 7.04 5.25 12.50 5.56 7.91 6.89];
85         test(X4, Y4);
86     end
87 end
88
89 end

```