
ECE4721J - Homework 5

Methods and Tools for Big Data

Kexuan Huang 518370910126

July 20, 2022



Ex.1 Numerical stability

1. In the big data context, how beneficial would it be to increase precision? For instance what would be the gain of using double instead of float, or move on with multi-precision?¹

The precision can be controlled to avoid accumulated errors. With a double, as the magnitude of the value increases, its precision decreases and this can introduce significant error into the result, which can increase the stability of methods in big data analysis.

2. Generate 100 random 1000×100 matrices X and measure the total time needed for MATLAB to compute:

MATLAB:

```
1 X = randi([0 100000], 1000, 100);  
2  
3 tic  
4 svd(X);  
5 toc  
6  
7 tic  
8 svd(X');  
9 toc  
10  
11 tic  
12 eig(X * X');  
13 toc  
14  
15 tic  
16 eig(X' * X);  
17 toc
```

Results:

```
1 Elapsed time is 0.00286484 seconds.  
2 Elapsed time is 0.00371003 seconds.  
3 Elapsed time is 0.0701489 seconds.  
4 Elapsed time is 0.000877857 seconds.
```

¹stackoverflow

3. Calculations

Let

$$X = \begin{pmatrix} -9 & 11 & -21 & 63 & -252 \\ 70 & -69 & 141 & -421 & 1684 \\ -575 & 575 & -1149 & 3451 & -13801 \\ 3891 & -3891 & 7782 & -23345 & 93365 \\ 1024 & -1024 & 2048 & -6144 & 24572 \end{pmatrix}$$

a) Use MATLAB to determine the eigenvalues of $X + \delta X$, where δX represents a small random perturbation on X . Study the variations over 1000 tests.

MATLAB:

```
1 X = [-9 11 -21 63 -252;
2      70 -69 141 -421 1684;
3      -575 575 -1149 3451 -13801;
4      3891 -3891 7782 -23345 93365;
5      1024 -1024 2048 -6144 24572];
6
7 err_eig = zeros(5, 1);
8 eig_X = eig(X);
9
10 for i = 1:1000
11     dX = eps(X);
12     err_eig = err_eig + eig(X + dX) - eig_X;
13 end
14
15 err_eig
```

Results:

```
1 err_eig =
2
3      7.731538279038739e-01 - 2.561862243142342e+00i
4      7.731538279038739e-01 + 2.561862243142342e+00i
5      2.718376927461463e+00 + 0.000000000000000e+00i
6     -2.132342289641740e+00 - 1.508562891852383e+00i
7     -2.132342289641740e+00 + 1.508562891852383e+00i
```

b) Use MATLAB to determine the singular values of $X + \delta X$, where δX represents a small random perturbation on X . Study the variations over 1000 tests.

MATLAB:

```
1 sum_dX = zeros(5, 5);
2 err_sin = zeros(5, 1);
3 svd_X = svd(X);
4
5 for i = 1:1000
6     dX = eps(X);
7     err_sin = err_sin + svd(X + dX) - svd_X;
8 end
9
10 err_sin
```

Results:

```
1 err_sin =
2
3     1.455191522836685e-08
4     9.305889392408062e-10
5     9.237055564881302e-11
6     3.752553823233029e-11
7     1.691492729682813e-10
```

4. In light of the lectures content, explain your observations.

SVD is much more stable than calculating eigenvalues, but might be a little slower to compute in MATLAB.

Ex.2 Simple SVD calculations

Without using a computer, find the singular values of the matrix

$$X = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 0 & 1 & 2 \end{pmatrix}.$$

$$X^T X = \begin{pmatrix} 1 & 5 & 9 \\ 2 & 6 & 0 \\ 3 & 7 & 1 \\ 4 & 8 & 2 \end{pmatrix} \times \begin{pmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 0 & 1 & 2 \end{pmatrix} = \begin{pmatrix} 107 & 32 & 47 & 62 \\ 32 & 40 & 48 & 56 \\ 47 & 48 & 59 & 70 \\ 62 & 56 & 70 & 84 \end{pmatrix}$$

$$\det(X^T X - \lambda I) = \det \begin{pmatrix} 107 - \lambda & 32 & 47 & 62 \\ 32 & 40 - \lambda & 48 & 56 \\ 47 & 48 & 59 - \lambda & 70 \\ 62 & 56 & 70 & 84 - \lambda \end{pmatrix} = 0$$

$$\lambda^4 - 290\lambda^3 + 12840\lambda^2 - 9600\lambda = 0$$

$$\begin{cases} \lambda_1 \approx 0.76 \\ \lambda_2 \approx 53.54 \\ \lambda_3 \approx 235.70 \end{cases}$$

$$\begin{cases} \sigma_1 \approx 0.87 \\ \sigma_2 \approx 7.32 \\ \sigma_3 \approx 15.35 \end{cases}$$

Ex.3 PCA in Spark

Please refer to `ex3.py`

1. Explain how PCA can be of any help to Krystor?²

PCA:

- Provides the best “perspective” that emphasises similarities and differences in the data
- This new perspective combines the original “characteristics” in order to best summarize the data

As a result, PCA can help Krystor quickly targeting the columns x (sensors) which are most related to the last columns y (hourly electric consumption)

2. How many columns of `sensors1.csv` are necessary to explain 90% of the data? Let n be that number.

```
1 column 0: 0.56
2 column 1: 0.28
3 column 2: 0.15
```

So $n = 3$.

3. Construct the linear model $y = \sum_{i=1}^n \beta_i p_i + x_0 + \varepsilon$, where, $\varepsilon \sim \text{Normal}(0, 1)$, x_0 is the intercept at $x = 0$, and $(p_i)_{1 \leq i \leq n}$ are the first n principal components of `sensors1.csv`.

The linear model is constructed as follow:

```
1 Coefficient: [5.85535522 7.84716081 6.86019518]
2 Intercept: 734.0304003124222
3 R_square: 0.9999719057230781
```

²Lecture Slide

4. Help Krystor determine whether sensors2.csv also contains the output of the sensors in the electric circuit of Reapor Rich's new cinema.

```
1 column 0: 0.52
2 column 1: 0.31
3 column 2: 0.16
4 Coefficient: [5.87137867 7.82779001 6.8058421 ]
5 Intercept: 734.3335911223433
6 R_square: 0.999965257449195
```

Yes, we found column 1~3 as the principal columns with R square value of 0.99, so sensors2.csv also contains the output of the sensors.