

Spl-1 Project Report- 2023

Leverage Score Sampling

Submitted by

Mahir Faisal

BSSE Roll No: 1316

Session- 2020-2021

Supervised by

Dr.Mohammad Shoyaib

Professor

Designation :

Institute of Information Technology



Institute of Information Technology

University of Dhaka

21-05-2023

Table of Contents

1. Introduction.....	3-4
1.2 Background.....	4-5
1.3 Why SVD for Leverage Score Sampling.....	5-6
1.4 Challenges Faced.....	6
2. Description of the Project.....	7-8
2.1 Read Data Matrix From File.....	8
2.2 Transpose of the Matrix.....	9-10
2.3 Multiplication.....	10-11
2.4 Calculate Coefficients.....	11-12
2.5 Bairstow Algorithm.....	12-13
2.5.1 Create & Delete Array.....	13
2.5.2 Calculate Absolute Value.....	13-14
2.5.3 Remove Error.....	14
2.5.4 Calculation r & s.....	15
2.5.5 Root Finding.....	15-16
2.5.5.1 PrintRoot.....	16-18
2.5.5.2 PrintRoot_One.....	18
2.5.6 Power Reduce.....	18-19
2.5.7 Get Coefficients.....	19-20
2.6 Gauss Elimination.....	20
2.6.1 Gauss Elimination Method.....	20-21
2.6.2 Swap Row.....	21
2.6.3 Forward Elimination.....	22
2.6.4 Backward Substitution.....	22-23
2.7 Calculate U & Σ	24
2.8 SVD for Higher Order.....	25-28
2.9 Pseudo Inverse.....	29
2.10 Leverage Score Sampling.....	29-30
2.11 Matrix Projection.....	30-32
3. User Manual.....	34
4. Conclusion.....	34
5. References.....	35

1. Introduction

Leverage Score Sampling is a technique used in data analysis to extract a representative subset of data points from large datasets. The technique is necessary because dealing with large datasets that have a high dimensionality can be computationally expensive and time-consuming. By selecting a representative subset of data points using Leverage Score Sampling, it becomes possible to perform various analyses on the dataset while saving time and computational resources.

Furthermore, **Leverage Score Sampling** allows for the extraction of important features of the dataset. This can be particularly useful when dealing with complex datasets where it is difficult to identify the key features. By identifying the important features of the dataset, it becomes possible to gain insights that may not be possible with other methods.

It has numerous applications in real life. Here are some examples of its importance :

1. Medical Research: In medical research, large datasets are often used to identify risk factors for various diseases. Leverage Score Sampling can be used to extract a representative subset of data points from these datasets, making it easier and faster to identify the risk factors.
2. Finance: In finance, large datasets are used to analyze market trends and predict stock prices. Leverage Score Sampling can be used to extract a representative subset of data points, making it easier to identify patterns and make predictions.
3. Marketing: In marketing, large datasets are used to analyze consumer behavior and develop marketing strategies. Leverage Score

Sampling can be used to extract a representative subset of data points, making it easier to identify trends and develop effective marketing campaigns.

4. Image and Video Processing: In image and video processing, Leverage Score Sampling can be used to compress large images and videos, making them easier to store and transmit.

Overall, **Leverage Score Sampling** is an important technique in data analysis that can help to improve the efficiency and effectiveness of analyses.

1.2 Background

Before we deeply discuss the project, some important terms should be clarified for proper comprehension of the implementation.

➤ **EigenValue** Eigenvalue is a mathematical concept that is commonly used in linear algebra and other fields of mathematics. In simple terms, eigenvalues are a set of scalar values that are associated with a matrix. It is important because it provides information about the properties of a matrix.

➤ **EigenVector** Eigenvector is a nonzero vector that, when multiplied by a matrix, results in a scalar multiple of the original vector. For easy to say, an eigenvector is a vector that remains in the same direction after it is transformed by a matrix. It is used to represent the direction of the greatest variance in a dataset.

➤ **l2 norm** l2 norm also known as the euclidean norm, is a measure of the length or magnitude of a vector in a euclidean space. It is defined as the square root of the sum of square values of the vector component.

➤ **SVD** Singular value Decomposition(SVD) is the process to break down a matrix into its constituent parts. SVD involves decomposing a matrix A into

three matrices : U , Σ & V . U contains the left singular vectors of A where V contains the right singular vectors, & Σ contains the singular values of A .

➤ **Singular Value** It is very important concept in linear algebra. Basically they are positive square roots of the Eigenvalues of a data matrix. They indicate how much each direction in the input space of a linear transformation is a scaled or contracted transformation.

➤ **Leverage Score** Leverage score is a mathematical measure used in linear algebra to determine the importance of each row/column of a matrix. Leverage scores are important because they can be used to identify influential data points or outliers in a dataset.

1.3 Why SVD for Leverage Score Sampling

Singular Value Decomposition (**SVD**) is an important technique used in data analysis, and it plays a crucial role in computing Leverage Score Sampling. It provides a number of benefits, including efficient computation, data compression, dimension reduction & the identification of important features of a dataset.

1. Efficient computation: **SVD** provides an efficient way to decompose a matrix into its constituent parts. This can be particularly useful when dealing with large matrices, as SVD allows for the decomposition to be performed quickly and accurately.
2. Data compression: **SVD** can be used to compress data by reducing the number of dimensions. This can be useful when dealing with large datasets, as it can reduce the amount of memory required to store the data.
3. Dimension reduction: **SVD** can be used to identify the most important features of a dataset by reducing the number of dimensions. This can be particularly useful when dealing with high-dimensional

datasets, as it can make it easier to analyze the data and identify patterns.

4. Leverage Score Sampling: **SVD** is used to compute Leverage Score Sampling by identifying the most important rows of a matrix. This can be useful when dealing with large datasets, as it allows for a representative subset of data points to be extracted while maintaining the important features of the dataset.

1.4 Challenges Faced

While accomplish that project, I had face some challenges that were not too easy to resolve. Some of these are :

- **Finding Appropriate Algorithm** : At the very beginning of this project, it was very difficult to find the algorithms that are needed or related to compute **SVD** (Singular Value Decomposition).
- **Implementation of the Algorithms** : There was a little bit difficulty to implement the Bairstow & Faddeev Leverrier algorithm.
- **Manage Higher order matrix** : For higher order matrix the Faddeev Leverrier shows the coefficients which are very large value. But Bairstow will fail for the large eigenvalues. So, then we had to use another algorithm, Jacobi Rotation.
- **Manage Large Codes** : Working with a large codebase is a tedious task. It becomes very difficult to track changes on different files.
- **Debugging codes** : Frequently, I got segmentation problems or dumped code errors but I was not sure where it happening. Then this issue was instantly solved by debugging the CPP files.

2 Description of the Project

My project is designed to compute Leverage Score Sampling using SVD(Singular Value Decomposition) to efficiently select a subset of rows from a large dataset. Basically, it divides the dataset into important & non-important features. It includes the following attributes:

- Read data matrix from file
- Transpose of the matrix
- Multiply transpose matrix with original matrix

if(dimension< 10) {

- Faddeev Leverrier Algorithm
 - Calculate coefficients of the characteristic equation
- Bairstow Algorithm
 - Create new array
 - Delete array
 - Calculate absolute value
 - Remove error
 - Root finding
 - ➡ Print real & complex root
 - ➡ Print real root only
 - Calculate initial guess r & s
 - Reduce power of the polynomial equation
 - Get coefficient after every iterate
- Gauss Elimination
 - Gauss elimination method
 - forward elimination
 - Back substitution
 - Swap row

}

```

else {
    ➤ SVD For Higher Order
}

➤ Pseudo Inverse
➤ Leverage Score Sampling
➤ Matrix Projection

```

2.1 Read Data Matrix from file

Here, a user can generate random numbers as a matrix. We maintained two for loops that will continuously use the 'rand' method to generate the numbers up to the dimension. And then write the matrix into the file "**spl_task2.txt**". Then read the matrix.

```

1214 for(int i=0;i<Dimen;i++){
1215     for(int j=0;j<Dimen;j++){
1216         ran = rand();
1217         number=ran % (maxn-minn)+minn ;
1218         A[i][j]=number;
1219     }
1220 }
1221 FILE *fp;
1222 fp = fopen("spl_task2.txt","w");
1223 for(int i=0;i<Dimen;i++){
1224     for(int j=0;j<Dimen;j++){
1225         fprintf(fp,"%lf ",A[i][j]);
1226     }fprintf(fp,"\n");
1227 }


```


2.2 Transpose of the Matrix

This method calculates the transpose of the original matrix. If the original matrix is A, then the transpose of the matrix that the method calculates is A^T .

```
750
751 void transpose(double P[Dimen+1][Dimen+1], double Q[Dimen+1][Dimen+1])
752
753 {
754     //double trans_val;
755     for(int i=0; i<Dimen; i++)
756     {
757         for(int j=0; j<Dimen; j++)
758         {
759             Q[i][j] = P[j][i];
760         }
761     }
762     //cout<<endl;
763 }
764
```

The input & output of the function :

 "C:\Users\DELL\Desktop\SPL-1\codes & file\SVD for any order.exe"

```
Given matrix :
41 67 134 100 169 124 78 158 162 64
105 145 81 27 161 91 195 142 27 36
191 4 102 153 92 182 21 116 118 95
47 126 171 138 69 112 67 99 35 94
103 11 122 133 73 64 141 111 53 68
147 44 62 157 37 59 123 141 129 178
116 35 190 42 88 106 40 142 64 48
46 5 90 129 170 150 6 101 193 148
29 23 84 154 156 40 166 176 131 108
144 39 26 123 137 138 118 82 129 141
```

So, the transpose is :

```

41 105 191 47 103 147 116 46 29 144
67 145 4 126 11 44 35 5 23 39
134 81 102 171 122 62 190 90 84 26
100 27 153 138 133 157 42 129 154 123
169 161 92 69 73 37 88 170 156 137
124 91 182 112 64 59 106 150 40 138
78 195 21 67 141 123 40 6 166 118
158 142 116 99 111 141 142 101 176 82
162 27 118 35 53 129 64 193 131 129
64 36 95 94 68 178 48 148 108 141

```

2.3 Multiplication

This function can multiply two different input matrices. In this program, this function multiply A^T with the original matrix A . Here, we use three for loops to compute the multiplication. The main calculative part of this function is : $r[i][j] = p[i][k] * q[k][j]$; where $p[i][k]$ represent the i 'th row & k 'th column of A^T & $q[k][j]$ represent the k 'th row & j 'th column of A .

```

730
731 void multiplication(double p[Dimen+1][Dimen+1], double q[Dimen+1][Dimen+1], double r[Dimen+1][Dimen+1])
732 {
733     double sum = 0;
734     for(int i=0; i<Dimen; i++)
735     {
736         for(int j=0; j<Dimen; j++)
737         {
738             for(int k=0; k<Dimen; k++)
739             {
740                 r[i][j] = p[i][k] * q[k][j];
741                 sum = sum + r[i][j];
742             }
743             r[i][j] = sum;
744             sum = 0;
745         }
746     }
747 }
748
749

```

Let, after multiply the resulting matrix $W = (A^T * A)$. So, the output matrix W , is :

```

So, multiplication between A_t & A is :
120763 42832 95558 112406 99887 110158 90159 118387 96759 97949
42832 46763 56793 47440 59022 50405 58840 64260 37159 40715
95558 56793 134822 116475 116756 113757 89826 136705 100268 91766
112406 47440 116475 152770 126864 122870 106120 144127 130741 126952
99887 59022 116756 126864 154154 126908 112578 148743 130189 108937
110158 50405 113757 122870 126908 131282 83088 128695 117720 104950
90159 58840 89826 106120 112578 83088 127165 127255 86750 89161
118387 64260 136705 144127 148743 128695 127255 168652 132870 120530
96759 37159 100268 130741 130189 117720 86750 132870 136719 116379
97949 40715 91766 126952 108937 104950 89161 120530 116379 115314

```

2.4 Calculate Coefficients

This function calculates the coefficients of the characteristic (polynomial) equation of the matrix which is found after multiplication A^T & A , using the **Faddeev Leverrier Algorithm**. For example if the characteristic equation is $x^3-3x^2+7x-2=0$, then the output of that function is **1,-3,7-2** which are the coefficients of the respective polynomial part.

```

650 void coefficient_calcul()
651 {
652     for(int j=1;j<=Dimen;j++)
653     {
654         for(int i=0;i<Dimen;i++){
655             for(int p=0;p<Dimen;p++){
656                 WM[i][p]=0;
657             }
658         }
659         double trace =0;
660         multiplication(w,M,WM);
661         for(int i=0;i<Dimen;i++){
662             trace += WM[i][i];
663         }
664         trace = ((-1)*trace)/j;
665         coefficients[j]=trace;
666         //trace = 0;
667         for(int i=0;i<Dimen;i++){
668             //
669             for(int p=0;p<Dimen;p++){
670                 if(i == p)
671                     {M[i][p] = WM[i][p] + coefficients[j];}
672                 else
673                     {M[i][p] = WM[i][p];}
674             }
675         }

```

After apply the Faddeev Leverrier Algorithm for the matrix W, the output is :

```

Coefficients of the characteristic equation are :
1 -1.2884e+06 2.45844e+11 -1.95261e+16 7.82218e+20 -1.6467e+25 1.75853e+29 -8.58373e+32 1.59162e+36 -7.98711e+38 5.45314e+40
Solving polynomial equation :
Roots: 80.6040

```

2.5 Bairstow Algorithm

In that part we implement the Bairstow Algorithm to find the roots(**eigenvalues**) of the characteristic equation. Here, we put the coefficients that we found from the above(part-2.4) implementation and then we will find the output of the roots of the polynomial equation.

Input for this part of the project:

```
Coefficients of the characteristic equation are :  
1 -1.2884e+06 2.45844e+11 -1.95261e+16 7.82218e+20 -1.6467e+25 1.75853e+29 -8.58373e+32 1.59162e+36 -7.98711e+38 5.45314e+40  
Solving polynomial equation :  
Roots: 80.6040
```

2.5.1 Create & Delete array

Here, we are allowed to create new arrays. And in the second part we deallocate those arrays.

```
34  
35 void new_arr()  
36  
37 {  
38  
39     a = new double[N];  
40  
41     b = new double[N];  
42  
43     c = new double[N];  
44  
45 }  
46  
47  
48  
49 void del_arr()  
50  
51 {  
52  
53     delete a;  
54  
55     delete b;  
56  
57     delete c;  
58
```

2.5.2 Calculate absolute value

It will return the absolute value of any real number.

```

62
63     double absolute(double x)
64
65     {
66         //return abs value
67
68         if(x<0)
69
70         {
71
72             x *= -1;
73
74         }
75
76         return x;
77     }
78
79

```

2.5.3 Remove Error

When a double value is very close to a round figure, then it will return only the integer part of the value.

```

85         //floating point error
86
87         int integer = val;
88
89         if(absolute(integer - val) <= phi)
90
91         {
92
93             val = (double)integer;
94
95         }
96
97
98
99         return val;
100

```

2.5.4 Calculation r & s

For every iteration we need to find the value of r & s of the quadratic factor x^2+rx+s which initially starts with the value 0.1 & 0.1 respectively. And then every iteration they will change following procedure:

```
281 void cal_r_s()
282
283 {
284
285     //for iteration we need to find r and s
286
287
288
289     dr = (b[0]*c[3] - b[1]*c[2]) / (c[2]*c[2] - c[1]*c[3]);
290
291     ds = (b[1]*c[1] - b[0]*c[2]) / (c[2]*c[2] - c[1]*c[3]);
292
293
294
295     old_r = r ;
296
297     old_s = s;
298
299     r += dr;
300
301     s += ds;
302
```

2.5.5 Root Finding

Basically, in this part we call two different functions **printroot** & **printroot_one** depending on the value of dimension(n) to calculate roots.

```
359 double ratio_s, ratio_r;
360
361
362 if(n == 0)
363 {
364
365
366     cout<<"No such variable.\n Wrong input\n\n";
367
368     exit(0);
369
370 }
371
```

```

371
372     else if(n == 1)
373
374     {
375
376         print_rootOne(a[n], a[n-1]);
377
378     }
379
380     else if(n == 2)
381
382     {
383
384         last = true;
385
386         printRoot(a[n], a[n-1], a[n-2]);
387
388     }
389

```

```

408
409     if(((absolute(b[0]) <= phi) && (absolute(b[1]) <= phi)) || ((absolute(ratio_r) <= phi) || (absolute(ratio_s) <= phi)))
410
411     {
412
413         printRoot(1,r,s);
414
415         if(n == 4)
416
417         {
418
419             last = true;
420
421             printRoot(b[n],b[n-1],b[n-2]);
422
423             break;
424
425         }
426         if(n == 3)
427
428         {
429
430             print_rootOne(b[n], b[n-1]);
431
432             break;
433
434         }
435
436         reduce_eqn();
437
438     }
439

```

2.5.5.1 PrintRoot

First of all we determine whether the root is complex or real. For this we calculate a double value **determine** = $r*r - 4*r*s$. If **determine** < 0, then the roots are complex. Again, whether the roots are pure imaginary or not, it is defined with the value of **r**.


```

158 //pure imaginary number
159
160 if((determine/(2*x)) == 1)
161 {
162
163     //coefficient 1
164
165     cout<<"\tRoot: "<<"i"<<endl;
166
167     cout<<"\tRoot: "<<"-i"<<endl;
168
169 }
170
171 else
172 {
173
174
175
176
177
178     cout<<"\tRoot: "<<(determine/(2*x))<<"i"<<endl;
179
180     cout<<"\tRoot: "<<(determine/(2*x))<<"-i"<<endl;
181
182 }
183
184 }

```

```

189
190 // not a pure imaginary number
191
192 if((determine/(2*x)) == 1)
193 {
194
195     //coefficient 1 , not necessary to print
196
197     cout<<"\tRoot: "<<((-p)/(2*x))<<" + "<<"i"<<endl;
198
199     cout<<"\tRoot: "<<((-p)/(2*x))<<" - "<<"i"<<endl;
200
201 }
202
203 else
204 {
205
206     //there are coefficient
207
208
209     cout<<"\tRoot: "<<((-p)/(2*x))<<" + "<<(determine/(2*x))<<"i"<<endl;
210
211     cout<<"\tRoot: "<<((-p)/(2*x))<<" - "<<(determine/(2*x))<<"i"<<endl;
212
213 }

```

Again, when **determine** >0, then the roots are real.

```
220         else
221
222     {
223
224         //cout<<determine<<endl;
225
226         determine = sqrt(determine);
227
228         //cout<<determine<<endl;
229
230         double first = remove_eror((( -p) - determine)/(2*x));
231
232         double second = remove_eror((( -p) + determine)/(2*x));
233     }
```

2.5.5.2 PrintRoot_one

If existence equation has only one solution, then that function might be helpful.

```
249
250 void print_rootOne(double x, double y)
251
252 {
253
254     // If existence equation has only one solution
255
256     //x *= -1;
257
258     //y *= -1;
259     //double roots[Dimen];
260     int i=0;
261
262
263
264     double root = -(y/x);
265 }
```

2.5.6 Reduce Power

After every iteration if we find two roots dividing with x^2+rx+s , then we have to reduce the equation. For this we may apply following procedure:

```

327 void reduce_eqn()
328
329 {
330
331     //After iteration, found two solution and the equations's power reduce by two.
332
333     //Replace a[] by b[]
334
335     for(int i=0; i<n-1; i++)
336
337     {
338
339         a[i] = b[i+2];
340
341     }
342
343
344
345     n -= 2;
346
347 }
348

```

2.5.7 Get Coefficients

In every iteration we have to define the coefficients of the polynomial equation using the following procedure.

```

446
447 void Get_Coefficient(double pass[Dimen+1], int total)
448
449 {
450
451     new_arr();
452
453     for(int i= total; i>=0; i--)
454
455     {
456
457         a[i] = pass[i];
458
459     }
460
461     n = total;
462
463
464
465
466
467     r = a[n-1]/a[n];
468
469     s = a[n-2]/a[n];
470

```

If we successfully complete the above procedure that means **Bairstow Algorithm**, then it will print all the roots(**eigenvalues**) of the characteristic equation. So, after put the coefficients(part-2.5) of W, then the output is:

```
Solving polynomial equation :  
Root: 80.6949  
Root: 1.07621e+06  
Root: 647.794  
Root: 71216.4  
Root: 2540.28  
Root: 47717.7  
Root: 6736.29  
Root: 41585.5  
Root: 15148.6  
Root: 26518.1
```

2.6 Gauss Elimination

Here, we used Gaussian elimination to calculate the eigenvectors for each eigenvalue. That will perform using **backward substitution**. It will take the eigenvalues found from **part-2.5** as an input and then compute eigenvectors for each eigenvalue.

2.6.1 Gauss Elimination Method

This function will be called for every eigenvalue. It will starts with a matrix W3, which can be found for every eigenvalue **ev[i]** following process: **$W3[i][i] = W[i][i] - ev[i]$**

Then we apply Gaussian Elimination on **W3**.

It will call the **forward_elimination** & **backward_substitution** method.

```

512 // function to get matrix content
513 void gaussianElimination(double w3[Dimen+1][Dimen+1])
514 {
515     /* reduction into r.e.f. */
516     int singular_flag = forwardElim(w3);
517
518     /* if matrix is singular */
519     if (singular_flag != -1)
520     {
521         printf("Singular Matrix.\n");
522
523         /* if the RHS of equation corresponding to
524            zero row is 0, * system has infinitely
525            many solutions, else inconsistent*/
526         if (w3[singular_flag][Dimen])
527             printf("Inconsistent System.");
528         else
529             printf("May have infinitely many "
530                  "solutions.");
531
532         return;
533     }
534
535     backSub(w3);
536 }

```

2.6.2 Swap Row

When the principal diagonal element is zero, it denotes that the matrix is singular. Then we have to swap the greatest value row with the current row.

```

542
543 for (int p=0; p<=Dimen; p++)
544 {
545     double temp = w3[i][p];
546     w3[i][p] = w3[j][p];
547     w3[j][p] = temp;
548 }

```

2.6.3 Forward Elimination

For every iteration we have to define a **pivot** and have to find greater amplitude for pivot if any.

```
557 // Initialize maximum value and index for pivot
558 int i_max = p;
559 int v_max = w3[i_max][p];
560
561 /* find greater amplitude for pivot if any */
562 for (int i = p+1; i < Dimen; i++)
563     if (abs(w3[i][p]) > v_max)
564         v_max = w3[i][p], i_max = i;
565
566 /* if a principal diagonal element is zero,
567  * it denotes that matrix is singular, and
568  */
569 if (!w3[p][i_max])
570     return p; // Matrix is singular
571
572 /* Swap the greatest value row with current row */
573 if (i_max != p)
574     swap_row(w3, p, i_max);
575
576
577 for (int i=p+1; i<Dimen; i++)
578 {
579     /* factor f to set current row kth element to 0,
580     * and subsequently remaining kth column to 0 */
581     double f = w3[i][p]/w3[p][p];
582
583     /* subtract fth multiple of corresponding kth
584     row element*/
585     for (int j=p+1; j<=Dimen; j++)
586         w3[i][j] -= w3[p][j]*f;
587 }
```

2.6.4 Backward Substitution

Basically, here we calculate the eigenvectors for each eigenvalue. We initialize the last element as 1 and then compute others element following backward elimination:

```
603 for (int i = Dimen-2; i >= 0; i--)
604 {
605     x[Dimen-1]=1;
606
607     for (int j=i+1; j<Dimen; j++)
608     {
609
610         x[i] =x[i]- w3[i][j]*x[j];
611     }
612
613     x[i] = x[i]/w3[i][i];
614 }
```

And then we normalize the eigenvectors by following that procedure:

```
624     for(int j=0; j<Dimen; j++)
625     {
626         sum = sum+ pow(x[j],2);
627     }
628     sum = sqrt(sum);
629
630     for(int i=0; i<Dimen; i++)
631     {
632         x[i]=x[i]/sum;
633
634         Vt[i][l] = x[i];
635         //V[l][i] = Vt[i][l];
636
637     }
638     l++;
639
```

After normalize the eigenvectors, we find the transpose matrix of the right singular factor **V**, which is the core element of the **SVD** as SVD stands for:

$$\mathbf{SVD(A)} = \mathbf{U} * \mathbf{\Sigma} * \mathbf{V_T}$$

After successfully completing the above(**part-2.6**) section, we are allowed to find eigenvectors.

Let's see the sample input and output for the above implementation:

If the eigenvalue is **80.6949**, then the eigenvectors(normalized) for this is:

```
Solution for the system:
-0.365186
-0.237957
-0.172109
-0.197151
-0.343623
0.606077
0.188190
0.389213
-0.200025
0.170903
```

2.7 Calculate \mathbf{U} & Σ

Singular matrix Σ is found from the square root of the sorted eigenvalues. And put them as a diagonal element of that matrix.

```
1334     for(int i=0;i<Dimen;i++)
1335     {
1336         S[i][i]=sqrt(eig_val[i]);
1337         SI[i][i]= 1/S[i][i];
1338     }
1339
```

And then the left singular factor \mathbf{U} can be found from the following procedure:

$$\mathbf{U} = \mathbf{A} * \mathbf{V}_T * \mathbf{S}_i ;$$

where \mathbf{S}_i is the inverse of the singular matrix \mathbf{S} . We showed in the above figure the procedure how we compute \mathbf{S}_i .

```
1377
1378     for(int i=0; i<Dimen; i++)
1379     {
1380         for(int j=0; j<Dimen; j++)
1381         {
1382             V[i][j] = Vt[j][i];
1383         }
1384     }
1385
1386
1387
1388     multiplication(Vt,SI,VtSi); // start to calculate U ;
1389     multiplication(A,VtSi,U); // end to calculate U ;
1390
1391     for(int i=0; i<Dimen; i++)
1392     {
```


2.8 SVD For Higher Order

When the dimension of the data matrix is very large then we apply the **Jacobi Rotation** method to calculate the eigenvalues. First of all we find the largest value of the matrix except the diagonal elements. Then store its index to compute the rest of the part.

Let's see the details procedure of that particular part of:

```
782
783   for(i=0; i<Dimen; i++) {
784       for(j=0; j<Dimen; j++) {
785           d[i][j]= w[i][j];
786           if(i==j)
787               s[i][j]= 1;
788           else
789               s[i][j]= 0;
790       }
791   }
792
793   do {
794       flag= 0;
795       p=0; q=1;
796       max= fabs(d[p][q]);
797
798       for(i=0; i<Dimen; i++) {
799           for(j=0; j<Dimen; j++) {
800               if(i!=j) {
801                   if (max < fabs(d[i][j])) {
802                       max= fabs(d[i][j]);
803                       p= i;
804                       q= j;
805                   }
806               }
807           }
808       }
809   }
```

When $d[p][p]$ is equal to $d[q][q]$, then $\theta = \pi/4$ or $-\pi/4$.

if ($d[p][q] > 0$)

$\theta = \pi/4$

else

$\theta = -\pi/4$

```
809
810 if(d[p][p]==d[q][q]) {
811     if (d[p][q] > 0)
812         theta= pi/4;
813     else
814         theta= -pi/4;
815 }
816 else {
817     theta=0.5*atan(2*d[p][q]/(d[p][p]-d[q][q]));
818 }
819
820 for(i=0; i<Dimen; i++) {
821     for(j=0; j<Dimen; j++) {
822         if(i==j) {
823             s1[i][j]= 1;
824             slt[i][j]= 1;
825         }
826         else {
827             s1[i][j]= 0;
828             slt[i][j]= 0;
829         }
830     }
831 }
832
833 s1[p][p]= cos(theta);
834 slt[p][p]= s1[p][p];
835
836 s1[q][q]= cos(theta);
837 slt[q][q]= s1[q][q];
838
839 s1[p][q]= -sin(theta);
840 s1[q][p]= sin(theta);
```

```

844
845     for(i=0; i<Dimen; i++) {
846         for(j=0; j<Dimen; j++) {
847             temp[i][j]= 0;
848             for(p=0; p<Dimen; p++)    temp[i][j]+= s1t[i][p]*d[p][j];
849         }
850     }
851
852     for(i=0; i<Dimen; i++) {
853         for(j=0; j<Dimen; j++) {
854             d[i][j]= 0;
855             for(p=0; p<Dimen; p++)    d[i][j]+= temp[i][p]*s1[p][j];
856         }
857     }
858
859     for(i=0; i<Dimen; i++) {
860         for(j=0; j<Dimen; j++) {
861             temp[i][j]= 0;
862             for(p=0; p<Dimen; p++)    temp[i][j]+= s[i][p]*s1[p][j];
863         }
864     }
865
866     for(i=0; i<Dimen; i++) {
867         for(j=0; j<Dimen; j++)    s[i][j]= temp[i][j];
868     }
869
870     for(i=0; i<Dimen; i++) {
871         for(j=0; j<Dimen; j++) {
872             if(i!=j)
873                 if(fabs(d[i][j]) > zero) flag= 1;
874         }
875     }
876 } while(flag==1);
877

```

After **Jacobi iteration**, the **eigenvalues** will be found from the diagonal element of the **d** matrix.

```

878     printf("\nThe eigenvalues are \n");
879     for(i=0; i<Dimen; i++){
880         eig_val[i]=d[i][i];
881         printf("%8.4lf", eig_val[i]);
882     }
883
884     for(i=0; i<Dimen; i++){

```

After finding the eigenvalues we calculate the **eigenvectors**, **singular matrix**, **left & right singular factor** as a same procedure that we discussed above(part- 2.6.4, 2.7).

So, by calculating the **U**, **Σ** & **V**, the procedure of computing **SVD** is completed.

Let's see the sample input & output of our computed **SVD** :

Input & Output:

Given matrix :

```
41 67 134 100 169 124 78 158 162 64
105 145 81 27 161 91 195 142 27 36
191 4 102 153 92 182 21 116 118 95
47 126 171 138 69 112 67 99 35 94
103 11 122 133 73 64 141 111 53 68
147 44 62 157 37 59 123 141 129 178
116 35 190 42 88 106 40 142 64 48
46 5 90 129 170 150 6 101 193 148
29 23 84 154 156 40 166 176 131 108
144 39 26 123 137 138 118 82 129 141
```

So, the Left Singular Vector U is :

```
0.343597 -0.0376534 -0.216513 -0.465052 -0.0095908 -0.0833276 0.135068 -0.470835 -0.606772 -0.211968
0.294027 -0.682065 0.0308561 -0.00515189 -0.524813 0.0938313 0.0998745 -0.0488022 0.196815 0.0483638
0.341396 0.388314 -0.238589 0.402821 -0.229299 -0.173688 -0.285862 -0.490679 0.362584 0.223153
0.283524 -0.191751 -0.337154 0.127117 0.423984 0.71958 -0.167658 0.0069037 0.0996376 0.0658474
0.275875 -0.162592 0.105822 0.243539 0.275276 -0.302603 -0.523736 0.301371 -0.508009 -0.332654
0.331701 0.126952 0.454483 0.387063 0.259699 0.0828043 0.626002 -0.134869 -0.0789211 -0.00975015
0.271759 -0.11116 -0.538831 0.151689 0.0826119 -0.444339 0.371562 0.380172 0.111524 0.0591917
0.334657 0.48796 -0.0991452 -0.458939 -0.070794 0.166753 0.055263 0.373465 0.178511 0.0794189
0.340691 -0.152885 0.413835 -0.387769 0.332106 -0.266017 -0.18246 -0.0949925 0.496883 0.194943
0.332418 0.171539 0.30396 0.110227 -0.476235 0.20419 -0.154134 0.326885 -0.295839 -0.179598
```

Singular matrix S:

```
1037.41 0 0 0 0 0 0 0 0 0
0 266.864 0 0 0 0 0 0 0 0
0 0 218.444 0 0 0 0 0 0 0
0 0 0 203.925 0 0 0 0 0 0
0 0 0 0 162.844 0 0 0 0 0
0 0 0 0 0 123.08 0 0 0 0
0 0 0 0 0 0 82.0749 0 0 0
0 0 0 0 0 0 0 50.4012 0 0
0 0 0 0 0 0 0 0 25.4518 0
0 0 0 0 0 0 0 0 0 8.98303
```

Right Singular Vector V is :

```
0.294324 0.146865 0.31463 0.35701 0.355005 0.326761 0.287888 0.386264 0.328429 0.304317
0.088919 -0.444081 -0.191207 0.243229 -0.105043 0.219258 -0.579317 -0.210712 0.417875 0.292282
-0.0029268 -0.138651 -0.622818 0.216425 -0.0474735 -0.389686 0.525383 -0.0174958 0.136624 0.311316
0.71791 0.0187673 0.0569307 0.130923 -0.547037 0.0724644 0.0666709 -0.0848892 -0.373187 0.0919311
-0.401974 -0.111631 0.490273 0.445744 -0.408084 -0.396407 -0.0425258 0.187232 -0.0378022 0.157911
-0.276932 0.694158 -0.133823 0.216815 -0.069965 0.229555 -0.105282 -0.381808 -0.133196 0.379925
0.11917 0.318299 -0.00590424 -0.560727 -0.26218 -0.273477 -0.211404 0.359294 0.34883 0.361995
-0.0209576 -0.329193 0.374235 -0.350426 0.234704 -0.0478743 0.199779 -0.423255 -0.209167 0.555119
-0.048898 -0.0974882 -0.271629 0.0688698 0.256111 0.0117714 -0.345679 0.53219 -0.600976 0.282765
-0.365186 -0.237957 -0.172109 -0.197151 -0.343623 0.606077 0.18819 0.389213 -0.200025 0.170903
```

2.9 Pseudo Inverse

The Pseudo Inverse of a matrix is a generalization of the inverse for matrices that may not have an inverse. It allows us to solve systems of equations even when the matrix is not invertible. Following procedure can be used to compute Pseudo Inverse :

```
1458
1459 // calculating Pseudo inverse
1460 multiplication(SI,Ut,SIUt);
1461 multiplication(Vt,SIUt,PI);
1462
1463 cout<<endl<<"Pseudo Inverse of A is :"<<endl;
1464 for(int i=0; i<Dimen; i++)
1465 {
1466     for(int j=0; j<Dimen; j++)
1467     {
1468         cout<<PI[i][j]<<" ";
1469     }
1470     cout<<endl;
1471 }
```

2.10 Leverage Score Sampling

In order to ignore the outlier in the dataset, Leverage Score Sampling is a better way than others. This method computes the leverage score for every row of left singular factor, **U**. And then identifies the highest leverage scores. At last it returns the corresponding row of the original matrix **A**, which is the representative sample of the original data matrix containing the most important features only.

```
1047
1048
1049
1050
1051
1052
    int d = Dimen-1, count = 0;
    for(int i=0;i<Dimen;i++){
        D[i][i] = eig_val[i]/d;
    }
```

```

1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112

    for(int i=0;i<Dimen;i++){
        rowwise_sum = rowwise_sum + lvs[i];
        if(rowwise_sum/total_sum >= 0.95){
            count++;
            break;
        }
        else
            count++;
    }
    cout<<endl<<"Count : "<<count<<endl;
    for(int i=0;i<count;i++){
        for(int j=0;j<Dimen;j++){
            if(lvs[i] == arr[j])
                indx_of_high_lvs[i] = j;
        }
    }
}

```

Leverage scores:

```

Sorted Leverage Scores :
0.0851373 0.0134673 0.00277791 0.00120558 0.000578197 0.000364592 0.0002136 0.000202382 0.000112513 8.71971e-06

Total_sum : 0.104068

Count : 3

Index of High Leverage scores :
9 8 7

```

2.11 Matrix Projection

Projection matrix is the matrix that represents the original matrix which can be created from only sample rows. From the sample row we compute the corresponding **UK**, **ΣK** & **Vk** for projection matrix.

```

1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137

    int ind = count-1, ind2 = count-1, ind3 = 0;
    for(int i=0;i<count;i++){
        for(int j=0;j<Dimen;j++){
            U_k[j][i] = U[indx_of_high_lvs[ind]][j];
            V_k[j][i] = V[indx_of_high_lvs[ind]][j];
            //S_k[i][j] = S[indx_of_high_lvs[ind]][j];
            sampl[i][j] = A[indx_of_high_lvs[ind3]][j];
        }
        ind--;
        ind3++;
    }
    transpose(V_k, Vt_k);
    for(int i=0;i<count;i++){
        S_k[i][i] = S[indx_of_high_lvs[ind2]][indx_of_high_lvs[ind2]];
        ind2--;
    }
}

```

Projection matrix :

```
1192     }
1193     void matrixProjection(double P[Dimen+1][Dimen+1], double Q[Dimen+1][Dimen+1], double R[Dimen+1][Dimen+1], int c) {
1194         for(int i=0; i<Dimen; i++)
1195         {
1196             for(int j=0; j<Dimen; j++)
1197             {
1198                 R[i][j]=0;
1199                 for(int g=0; g<c; g++)
1200                 {
1201                     R[i][j] += P[i][g] * Q[g][j];
1202                 }
1203             }
1204         }
1205     }
1206 }
```

```
1159     int g=0;
1160     for(int i=0; i<count; i++) {
1161         for(int j=0; j<Dimen; j++) {
1162             SkVkt[i][j] = S_k[i][i] * Vt_k[i][j];
1163         }
1164     }
1165
1166     /*cout<<"SKVKT : "<<endl;
1167     for(int i=0; i<count; i++) {
1168         for(int j=0; j<Dimen; j++) {
1169             cout<<SkVkt[i][j]<<" ";
1170         }cout<<endl;
1171     }*/
1172     matrixProjection(U_k, SkVkt, A_k, count);
1173     cout<<endl<<"Reconstructed A_K : "<<endl;
1174     for(int i=0; i<Dimen; i++) {
1175         for(int j=0; j<Dimen; j++) {
1176             cout<<A_k[i][j]<<" ";
1177         }cout<<endl;
1178     }
```

Let's see the sample input & final output of that project :

Input :

```
Given matrix :
41 67 134 100 169 124 78 158 162 64
105 145 81 27 161 91 195 142 27 36
191 4 102 153 92 182 21 116 118 95
47 126 171 138 69 112 67 99 35 94
103 11 122 133 73 64 141 111 53 68
147 44 62 157 37 59 123 141 129 178
116 35 190 42 88 106 40 142 64 48
46 5 90 129 170 150 6 101 193 148
29 23 84 154 156 40 166 176 131 108
144 39 26 123 137 138 118 82 129 141
```

Output(Representative sample) :

```
Sample data Matrix :
144 39 26 123 137 138 118 82 129 141
29 23 84 154 156 40 166 176 131 108
46 5 90 129 170 150 6 101 193 148

Process returned 0 (0x0)   execution time : 0.515 s
Press any key to continue.
```

Here, we see that only three of ten rows contain the most important feature in this dataset. So, we can store this 3x10 matrix instead of the original 10x10 matrix.

3 User Manual

Our project will start from the main function. First of all it generates random numbers & writes them into the file, “ **spl_task2.txt** ”. Then program control goes to the transpose function & calculates the transpose of the matrix. Then it calculates the multiplication of A^T & A . In order to compute **SVD**, we apply two different procedures based on the value of dimension. If $n < 10$, then we use the **Bairstow Algorithm** to calculate the eigenvalues. Before Bairstow, we have to implement the **Faddeev Leverrier** algorithm to find the coefficients of the

polynomial

equation.

```
1274
1275     for(int i=0;i<Dimen;i++)
1276     {
1277         M[i][i]=1;
1278     }
1279     if(Dimen>10){
1280         SVDforHigherOrder();
1281     }
1282     else
1283     {
1284         coefficient_calcul();
1285         cout<<"Coefficients of the characteristic equation are : "<<endl;
1286         coefficients[0]=1;
1287         for(int i=0;i<=Dimen;i++){
1288             cout<<coefficients[i]<<" ";
1289         }
```

In Bairstow, we will use the functions **GetCoefficients**, **reduce_equation**, **root_finding**, **printRoot** etc. After that implementation we will find the eigenvalues.

```
Solving polynomial equation :
Root: 2.24992
Root: 675878
Root: 596.485
Root: 50939.3
Root: 7148.63
Root: 46181.1
Root: 12355
Root: 22007.3
```

If $n > 10$, program control goes to the function **SVDForHigherOrder**. Here, we used the **Jacobi Rotation** method to calculate the eigenvalues. Then program control goes to **Gauss Elimination**. Here, we calculate the eigenvectors for each eigenvalue. Then we calculate the **U**, **Σ** & **V** matrix. After multiplication **U**, **Σ** & **V^T** we find the resulting matrix which is either equal or very close to the original matrix A.

```
So, the complete SVD is :
41.3882 66.2755 133.18 100.355 169.892 124.822 78.0244 157.083
161.896 64.1811 105.113 144.841 80.7271 26.8941 161.129 91.2458
194.9 142.19 27.2285 35.9144 190.767 3.76962 101.978 153.243
91.6139 182.712 21.7596 115.637 117.128 94.2426 47.0167 126.882
171.522 137.016 67.9394 112.518 68.2322 100.05 34.9723 92.7596
102.712 11.5602 122.699 132.758 72.3201 63.2982 140.902 111.715
52.7005 68.5435 147.583 43.7233 61.3307 156.413 37.0185 59.6804
123.103 140.822 128.828 178.106 116.232 35.1753 189.958 41.7698
```

Then program control goes to the function **LeverageScoreSampling**. That function calculates the leverage scores and finally returns the representative sample of the original matrix which contains the most important feature of the dataset.

```
Sample data Matrix :  
144 39 26 123 137 138 118 82 129 141  
29 23 84 154 156 40 166 176 131 108  
46 5 90 129 170 150 6 101 193 148  
  
Process returned 0 (0x0)   execution time : 0.515 s  
Press any key to continue.
```

4 Conclusion

Leverage score sampling using SVD is a game-changer when it comes to analyzing high-dimensional datasets. Its ability to efficiently select representative subsets of data based on leverage scores unlocks new possibilities for deriving actionable insights. By leveraging the mathematical properties of SVD, this technique enhances efficiency, reduces noise, and improves decision-making processes. Whether in finance, marketing, or other fields, leverage score sampling using SVD empowers analysts and decision-makers to make the most of their data.

So, embrace the power of leverage score sampling using SVD and unlock the potential within our high-dimensional datasets

Reference

1. https://math.iitm.ac.in/public_html/sryedida/caimna/transcendental/polynomial%20methods/brs%20method.html, Bairstow Method, 19/05/2023
2. [Faddeev–LeVerrier algorithm - Wikipedia](#), 19/05/2023

3. <https://www.cs.cmu.edu/afs/cs/user/dwoodruf/www/teaching/15859-fall17/scribe8.pdf>, David Woodruff, 19/05/2023
4. [Regular Inverse & Pseudo Inverse Matrix Calculation using Singular Value Decomposition | by Antonius Freenergi | Medium](#), Regular vs Pseudo Inverse, 19/05/2023
5. [Singular Values Decomposition \(SVD\) in C++11 by an Example - CodeProject](#), Arthur V.Ratz, 19/05/2023
6. eigen value & eigen vector.pdf, 19/05/2023
7. Iterative_methods_for_eigenvalue_problem.pdf, 19/05/2023