

shahinrahbariasl.medium.com

Data Engineering Cheat Sheet - Shahin Rahbariasl - Medium

Shahin Rahbariasl

12–15 minutes



Photo by [Ovan](#) from [StockSnap](#)



Data engineering is all about designing, developing, testing and delivering offerings using leading edge and/or proven technologies. Most companies nowadays work Agile, which makes it possible to have a collaborative environment to understand stakeholder requirements and develop valuable features. Data engineering also includes designing, coding, and testing innovative component-level software solutions. You always need to ensure that the implemented solutions are unit tested and ready to be integrated into their product.

As a data engineer, you're in charge of gathering and collecting the data, storing it, do batch processing or real-time processing on it, and serve it via an API to data scientist who can easily query it. Data engineers may work closely with data architects (to determine what data management systems are appropriate) and data scientists (to determine which data are needed for analysis). They often wrestle with problems associated with database integration and messy, unstructured data sets.

Here's a list of minimum topics that you need to have deep knowledge about as a data engineer. Keep in mind the intention of this article is not to explain the subjects in details, but to list important points about them. As an #IBM engineer, I catch myself using most of these tips pretty often. Hopefully, they'll be useful for you as well.

Data Structures

[Array](#)

A container object that holds a fixed number of values of a single type. Arrays are good for indexing, but not efficient for searching, inserting, and deleting.

Linear arrays: One dimensional arrays declared with a fixed

size.

Dynamic arrays: Hold reserved space for expansion (if it gets full, it copies itself to a larger array).

Multi dimensional arrays: Multi dimensional or nested arrays of arrays.

Linked List

A linear collection of data elements whose order is not given by their physical placement in memory. Instead, each element points to the next. Linked lists are fast for insertion and deletion, but slow at indexing and searching.

Each node contains a value and a pointer to the next node.

Doubly linked list nodes have pointers to the previous node as well. **Circularly linked list** has a tail (the last node), references the head (the first node).

Stack (LIFO) usually implemented with linked lists (the head being the only place for insertion and removal).

Queues (FIFO) can be implemented with linked lists as well (doubly linked list that only removes from head and adds to tail).

Time Complexity:

- Indexing: Linked Lists: $O(n)$
- Search: Linked Lists: $O(n)$
- Optimized Search: Linked Lists: $O(n)$
- Insertion: Linked Lists: $O(1)$

Hash Table or Hash Map

A data structure that implements an associative array abstract

data type, a structure that can map keys to values. A **hash table** uses a **hash** function to compute an index, also called a **hash** code, into an array of buckets or slots, from which the desired value can be found.

Hashing

Uses a **hash function** accept a key and return an output unique only to that specific key. An input data gets passed to the hash function and a much smaller data gets generated which can be a unique address in memory for that data.

- Hashing is not encryption
- Hashing functions are typically one-way
- Information can be lost when hashing
- Two objects that are equal should return the same hash
- The same hash may also result from different objects (**Hash Collision**)
- Hash tables are very fast!

Hash collisions are when a hash function returns the same output for two distinct inputs. Linked lists can be used to implement list of objects for the same hash.

Time Complexity:

- Indexing: Hash Tables: $O(1)$
- Search: Hash Tables: $O(1)$
- Insertion: Hash Tables: $O(1)$

Binary Tree

A **tree** data structure in which each node has at most two

children, which are referred to as the left child and the right child.

Binary Search Tree:

- A binary tree that uses comparable keys to assign which direction a child is.
- Left child has a key smaller than it's parent node.
- Right child has a key greater than it's parent node.
- There can be no duplicate node.

Time Complexity:

- Indexing: Binary Search Tree: $O(\log n)$
- Search: Binary Search Tree: $O(\log n)$
- Insertion: Binary Search Tree: $O(\log n)$

Sorting

Heap Sort

Heap sort is a comparison based sorting technique based on [Binary Heap](#) data structure. It is similar to selection sort where we first find the maximum element and place the maximum element at the end. We repeat the same process for the remaining elements.

Heap sort algorithm:

1. Build a max heap from the input array.
2. At this point, the largest item is stored at the root of the heap. Replace it with the last item of the heap and reduce the size of heap by 1. Finally, heapify the root of the tree.
3. Repeat step 2 while size of heap is greater than 1.

```

2 public class HeapSort {
3     public void sort(int arr[])
4     {
5         int n = arr.length;
6
7         // Build heap (rearrange array)
8         for (int i = n / 2 - 1; i >= 0; i--)
9             heapify(arr, n, i);
10
11        // One by one extract an element from heap
12        for (int i = n - 1; i > 0; i--) {
13            // Move current root to end
14            int temp = arr[0];
15            arr[0] = arr[i];
16            arr[i] = temp;
17
18            // call max heapify on the reduced heap
19            heapify(arr, i, 0);
20        }
21    }
22
23    // To heapify a subtree rooted with node i which is
24    // an index in arr[], n is size of heap
25    void heapify(int arr[], int n, int i)
26    {
27        int largest = i; // Initialize largest as root
28        int l = 2 * i + 1; // left = 2*i + 1
29        int r = 2 * i + 2; // right = 2*i + 2
30
31        // If left child is larger than root
32        if (l < n && arr[l] > arr[largest])
33            largest = l;
34
35        // If right child is larger than largest so far
36        if (r < n && arr[r] > arr[largest])
37            largest = r;
38
39        // If largest is not root
40        if (largest != i) {
41            int swap = arr[i];
42            arr[i] = arr[largest];
43            arr[largest] = swap;
44
45            // Recursively heapify the affected sub-tree
46            heapify(arr, n, largest);
47        }
48    }
49
50    /* A utility function to print array of size n */
51    static void printArray(int arr[])
52    {
53        int n = arr.length;
54        for (int i = 0; i < n; ++i)
55            System.out.print(arr[i] + " ");
56        System.out.println();
57    }
58
59    // Driver code
60    public static void main(String args[])
61    {
62        int arr[] = { 12, 11, 13, 5, 6, 7 };
63        int n = arr.length;
64
65        HeapSort ob = new HeapSort();
66        ob.sort(arr);
67
68        System.out.println("Sorted array is");
69        printArray(arr);
70    }
71 }
72

```

Source: [GeeksforGeeks](https://www.geeksforgeeks.org/heap-sort/)

Time Complexity: Time complexity of heapify is $O(\log n)$. Time complexity of building max heap is $O(n)$ and overall time complexity of Heap Sort is $O(n \log n)$.

Quick Sort

QuickSort is a Divide and Conquer algorithm. It picks an element as pivot and partitions the given array around the picked pivot. The key process in quickSort is partition(). Target of partitions is, given an array and an element x of array as pivot, put x at its correct position in sorted array and put all smaller elements (smaller than x) before x, and put all greater

elements (greater than x) after x.

```
36
37 /* The main function that implements QuickSort()
38 arr[] -> Array to be sorted,
39 low -> Starting index,
40 high -> Ending index */
41 void sort(int arr[], int low, int high)
42 {
43     if (low < high)
44     {
45         /* pi is partitioning index, arr[pi] is
46         now at right place */
47         int pi = partition(arr, low, high);
48
49         // Recursively sort elements before
50         // partition and after partition
51         sort(arr, low, pi-1);
52         sort(arr, pi+1, high);
53     }
54 }
55
```

```
1 // Java program for implementation of QuickSort
2 class QuickSort
3 {
4     /* This function takes last element as pivot,
5     places the pivot element at its correct
6     position in sorted array, and places all
7     smaller (smaller than pivot) to left of
8     pivot and all greater elements to right
9     of pivot */
10    int partition(int arr[], int low, int high)
11    {
12        int pivot = arr[high];
13        int i = (low-1); // index of smaller element
14        for (int j=low; j<high; j++)
15        {
16            // If current element is smaller than the pivot
17            if (arr[j] < pivot)
18            {
19                i++;
20
21                // swap arr[i] and arr[j]
22                int temp = arr[i];
23                arr[i] = arr[j];
24                arr[j] = temp;
25            }
26        }
27
28        // swap arr[i+1] and arr[high] (or pivot)
29        int temp = arr[i+1];
30        arr[i+1] = arr[high];
31        arr[high] = temp;
32
33        return i+1;
34    }
35 }
36
```

Source: [GeeksforGeeks](https://www.geeksforgeeks.org/quick-sort/)

Worst case: $O(n^2)$. when the partition process always picks greatest or smallest element as pivot.

Best case: $O(n \log n)$. when the partition process always picks the middle element as pivot.

Merge Sort

Merge sort divides the input array into two halves, calls itself for the two halves, and then merges the two sorted halves.

```
1 /* Java program for Merge Sort */
2 class MergeSort
3 {
4     // Merges two subarrays of arr[].
5     // First subarray is arr[l..m]
6     // Second subarray is arr[m+1..r]
7     void merge(int arr[], int l, int m, int r)
8     {
9         // Find sizes of two subarrays to be merged
10        int n1 = m - l + 1;
11        int n2 = r - m;
12
13        /* Create temp arrays */
14        int L[] = new int[n1];
15        int R[] = new int[n2];
16
17        /*Copy data to temp arrays*/
18        for (int i = 0; i < n1; ++i)
19            L[i] = arr[l + i];
20        for (int j = 0; j < n2; ++j)
21            R[j] = arr[m + 1 + j];
22
23        /* Merge the temp arrays */
24
25        // Initial indexes of first and second subarrays
```



```

26     int i = 0, j = 0;
27
28     // Initial index of merged subarray array
29     int k = l;
30     while (i < n1 && j < n2) {
31         if (L[i] <= R[j]) {
32             arr[k] = L[i];
33             i++;
34         }
35         else {
36             arr[k] = R[j];
37             j++;
38         }
39         k++;
40     }
41
42     /* Copy remaining elements of L[] if any */
43     while (i < n1) {
44         arr[k] = L[i];
45         i++;
46         k++;
47     }
48
49     /* Copy remaining elements of R[] if any */
50     while (j < n2) {
51         arr[k] = R[j];
52         j++;
53         k++;
54     }
55 }

```

```

57 // Main function that sorts arr[l..r] using
58 // merge()
59 void sort(int arr[], int l, int r)
60 {
61     if (l < r) {
62         // Find the middle point
63         int m = (l + r) / 2;
64
65         // Sort first and second halves
66         sort(arr, l, m);
67         sort(arr, m + 1, r);
68
69         // Merge the sorted halves
70         merge(arr, l, m, r);
71     }
72 }

```

Source: [GeeksforGeeks](https://www.geeksforgeeks.org/merge-sort/)

Time complexity of Merge Sort is $\theta(n\log n)$ in all 3 cases (worst, average and best) as merge sort always divides the array into two halves and takes linear time to merge two halves.

Bubble Sort

Bubble sort repeatedly swaps the adjacent elements if they are in wrong order.

```

2 class BubbleSort
3 {
4     void bubbleSort(int arr[])
5     {
6         int n = arr.length;
7         for (int i = 0; i < n-1; i++)
8             for (int j = 0; j < n-i-1; j++)
9                 if (arr[j] > arr[j+1])
10                {
11                    // swap arr[j+1] and arr[j]
12                    int temp = arr[j];
13                    arr[j] = arr[j+1];
14                    arr[j+1] = temp;
15                }
16     }
17 }

```

Source: [GeeksforGeeks](https://www.geeksforgeeks.org/bubble-sort/)

Worst and Average Case: $O(n^2)$. Worst case occurs when array is reverse sorted.

Best Case: $O(n)$. Best case occurs when array is already sorted.

Searching

Breadth First Search

An [algorithm](#) for traversing or searching [tree](#) or [graph](#) data structures. It starts at the [tree root](#) and explores all of the neighbor nodes at the present depth prior to moving on to the nodes at the next depth level.

It uses the opposite strategy of [depth-first search](#), which instead explores the node branch as far as possible before being forced to backtrack and expand other nodes.

- It finds every node on the same level, most often moving left to right.
- While doing this it tracks the children nodes of the nodes on the current level.
- When finished examining a level it moves to the left most node on the next level.
- The bottom-right most node is evaluated last (the node that is deepest and is farthest right of it's level).
- Optimal for searching a tree that is wider than it is deep.
- Uses a queue to store information about the tree while it traverses a tree.
- Because it uses a queue it is more memory intensive than depth first search.
- The queue uses more memory because it needs to store pointers

Time Complexity:

- Search: Breadth First Search: $O(V + E)$

- E is number of edges
- V is number of vertices

Depth First Search

An [algorithm](#) for traversing or searching [tree](#) or [graph](#) data structures. The algorithm starts at the [root node](#) and explores as far as possible along each branch before backtracking.

- It traverses left down a tree until it cannot go further.
- Once it reaches the end of a branch it traverses back up trying the right child of nodes on that branch, and if possible left from the right children.
- When finished examining a branch it moves to the node right of the root then tries to go left on all it's children until it reaches the bottom.
- The right most node is evaluated last (the node that is right of all it's ancestors).
- Optimal for searching a tree that is deeper than it is wide.
- Uses a stack to push nodes onto.
- Because a stack is LIFO it does not need to keep track of the nodes pointers and is therefore less memory intensive than breadth first search.
- Once it cannot go further left it begins evaluating the stack.

Time Complexity:

- Search: Depth First Search: $O(|E| + |V|)$
- E is number of edges
- V is number of vertices

REST APIs

[Representational state transfer \(REST\)](#) is a [de-facto standard](#) for a [software architecture](#) for [interactive](#) applications that typically use multiple [Web services](#).

- Uniform **stateless** operations
- No session information is retained by the server
- A request should not rely on another request
- **HTTP** methods used in **RESTful** APIs are **GET**, **PUT** (replace resource), **POST**, **PATCH** (update resource), **DELETE**

Advantages:

- Client-server architecture (separation of concern)
- Statelessness
- Cache-ability
- Layered system (eg. proxies should not prevent communication and no update to client/server code)
- Code on demand

Important Design Patterns

1. Singleton

- limit creation of a class to only one object
- Private/Protected constructor
- eg. used in caches, thread pools, registries

2. Factory method

- Produces class objects

- Decouples the code from the implementation class

3. Strategy

- Instead of directly implementing a single algorithm, the code receives runtime instructions to decide which of the group of algorithms to run

4. Observer

- There are observers (interested in updates) and subjects (generate updates)

5. Builder

- Builds more complex objects
- Step by step setting of properties vs factory method that creates objects in one go

6. Adaptor

- Adapter between two interfaces like translators
- eg. converting XML to JSON

7. State

- Allows object to alter its behaviour when its internal state changes
- eg. context for defining states

System Design

1. Outline use cases, constraints, and assumptions
 - Who, how, what, how much data, how many requests
2. High level design
 - Main components and connections
3. Design core components

4. Scale the design

- Load balancer, horizontal/vertical scaling, sharding, caching, etc.

Microservices Architecture

A variant of the [service-oriented architecture](#) (SOA) structural style — arranges an application as a collection of [loosely coupled](#) services. In a micro services architecture, services are fine-grained and the protocols are lightweight.

Strengths

- Independent components
- Easier understanding
- Better scalability
- Flexibility in choosing technology
- Higher level of agility

Weaknesses

- Extra complexity
- System distribution
- Cross cutting concerns
- Testing

When to choose Microservices architecture?

- Expertise availability (DevOps)
- A complex and scalable app
- Enough engineering skills

Monolithic system

In monolithic systems, functionally distinguishable aspects (for example data input and output, data processing, error handling, and the user interface) are all interwoven, rather than containing architecturally separate components.

Strengths

- Less cross cutting concerns
- Easier debugging and testing
- Simpler to deploy
- Simpler to develop

Weaknesses

- Gets complicated over time
- Making changes are hard (tightly coupled)
- Hard/expensive to scale
- New technology barriers

When to choose monolithic?

- Smaller team
- Simpler application
- Quicker launching
- No micro services expertise

Query Optimization Tips

SQL

- Define business requirements and select required fields instead of *select **

- Select more fields to avoid selecting *distinct* which is expensive
- Use inner joins instead of outer joins if possible
- *Where* is more efficient than *Having*
- Use * at the end of the phrase only
- Run analytical queries during off-peak hours
- Order select fields based on **granularity**
- **Index** tables based on int IDs

NoSQL

- Query rewrites based on heuristics (remove unnecessary stuff)
- Index selection (avoid searching whole tables)
- Join reordering
- Join types
- Add index on collections (eg. ascending vs descending)

Good Habits

- Log everything in sys logs
- Always have unit tests
- Validate data before inserting to DB
- Automate everything you do again and again
- Don't delete code, use git
- Instead of deleting in DB, hide from user
- Using too many globals is usually not a good idea
- Skim the whole manual instead of the parts you need

- WTFM: Write The Fucking Manual !
- If you can't verify if something is correctly working, it's probably not working

Interview Coding Canvas

Separate the constraints, limitations, ideas, code, and test cases. For idea generation, try to:

- Simplify the task
- Try examples to find patterns
- Think of suitable data structures