# Lecture 5: Arrays

---

## Arrays

Last lecture:  read in huge quantities of data.
This lecture:  store and manipulate huge quantities of data.

Arrays.
- Organized way to store huge quantities of data.
    - 52 playing cards in a deck
    - 5 thousand Princeton undergrads
    - 1 million characters in a book
    - 4 billion nucleotides in a strand of DNA
    - 73 billion Google queries per year
    - $6.02 \times 10^{23}$ particles in a mole

Today's applications.
- Data analysis.  (histogram)
- Data processing. (shuffling, sorting)
- Scientific applications. (DLA simulation)

---

## Arrays in Java

Arrays are built into Java.

- Essential property:  can directly access an element given its index.

- Declare and initialize using `[]` and `{}`.

    ```java
    int a0 = 3, a1 = 1, a2 = 4, a3 = 1, a4 = 5;
    int a5 = 9, a6 = 2, a7 = 6, a8 = 5, a9 = 3;
    ```

    vs.

    ```java
    int[] a = { 3, 1, 4, 1, 5, 9, 2, 6, 5, 3 };
    ```

- To access element i of array named `a`, use `a[i]`

---

## Choosing a Random Student

Simple application:  store related data as a group, and select random item.

```java
public class RandomStudent {
   public static void main(String[] args) {
      String[] names = { "Clelia Zacharias", "Hannah Xu",
                          "Virginia Wylly",   "Wendy Wu",
                          "Ashely Wolf",      "Eric Whitman",
                          "Will Weidman",     "Sharon Weeks",
                          "Mary Wathall",     "Sarah Wang",
                          "Michael Wang",     "Madeleine Walsh"
            12          };
             ↓
      int N = names.length;
      int r = (int) (Math.random() * N);   ← integer between 0 and 11
      System.out.println(names[r]);
   }
}
```

## California Runoff Election '04

135 candidates on ballot for governor of California.
- Alphabetical order prejudiced against Jon Zellhoefer.
- One solution: in each district, randomize the order in which the candidates appear.

| name |
| --- |
| Iris Adam |
| Brooke Adams |
| Cruz Bustamante |
| Gary Coleman |
| Larry Flynt |
| Georgy Russell |
| Arnold Schwarzenegger |
| Peter Ueberroth |
| Jon Zellhoefer |

⟹

| name |
| --- |
| Peter Ueberroth |
| Gary Coleman |
| Arnold Schwarzenegger |
| Brooke Adams |
| Jon Zellhoefer |
| Georgy Russell |
| Cruz Bustamante |
| Iris Adam |
| Larry Flynt |

---

## Creating Arrays in Java

How to declare an "empty" array.

- Declare using `[]`.

```
double a0, a1, a2, a3, a4, a5, a6, a7, a8, a9;
```

vs.

```
double[] a = new double[10];
```

- Allocate memory using `new`.

- All array elements are auto-initialized to:
  - zero for numeric types
  - `null` for String

---

## Shuffling

```java
public class Shuffle {
    public static void main(String[] args) {
        int N = Integer.parseInt(args[0]);
        String[] a = new String[N];

        for (int i = 0; i < N; i++)          ← read in and store data
            a[i] = StdIn.readString();



                    SHUFFLE        ▷



        for (int i = 0; i < N; i++)          ← print data
            System.out.println(a[i]);
    }
}
```

---

## Shuffling

Shuffle an N-element array.
- In $i^{th}$ iteration:
  - choose random integer r between 0 and i
  - swap values in positions r and i
- Need random access to arbitrary element ⟹ use arrays.

Property: after $i^{th}$ iteration, array positions 0 through i contain random permutation of elements 0 through i.

```java
for (int i = 0; i < N; i++) {
    int r = (int) (Math.random() * (i + 1));
    String swap = a[r];                    ↑
    a[r] = a[i];                    between 0 and i
    a[i] = swap;
}
                                       shuffle
```

## Gambler's Problem Revisited

Flip a fair coin N times and plot distribution of number of heads.

- Use `freq[i]` to record number of times you get exactly `i` heads.

```java
public class Flip {
    public static void main(String[] args) {
        int N = Integer.parseInt(args[0]);
        int[] freq = new int[N + 1];

        for (int i = 0; i < 10000; i++) {
            int heads = 0;              flip N coins
            for (int j = 0; j < N; j++)
                if (Math.random() < 0.5)
                    heads++;            flip one coin
            freq[heads]++;              increment counter
        }
    }
}
```
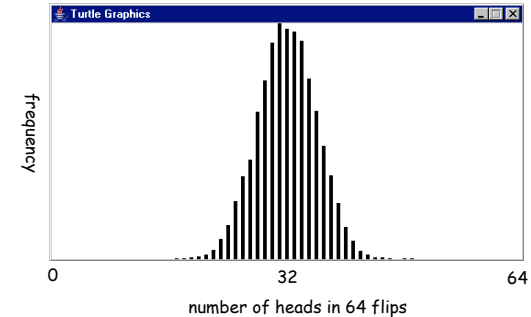
- Add graphic commands to plot.

---

## Gambler's Problem Revisited

What is distribution of number of heads?

- "Bell curve."
- Approximately Gaussian (stay tuned).
- Mean = N / 2, variance = N / 4.    ← 95% confidence interval: N/2 ± √N

```
% java Flip 64
```



number of heads in 64 flips

---

## Benford's Law

Examine listing of statistical data.

- Compute frequency count of leading digit.
  - Ex: leading digit of 456789 is 4.
- Print fraction of occurrences of each digit 1 - 9.
- What is distribution?  11.11% each?  Something else?

Use 10-element array `count`.

- `count[i]` counts number of times `i` is leading digit.   ← count[0] is always 0
- `N` counts total number of items processed.
- Print ratio for each `i`.

---

## Benford's Law

```java
public class Benford {
    public static void main(String[] args) {
        int[] count = new int[10];
        int N = 0;                      auto-initialized to 0

        while (!StdIn.isEmpty()) {
            int x = StdIn.readInt();

            while (x >= 10) x = x / 10;   ← compute leading digit

            count[x]++;  ← increment bin x
            N++;
        }

        for (int i = 1; i < 10; i++)
            System.out.println(i + ": " + 1.0 * count[i] / N);
    }                                          ↑
}                                      avoid integer division
```

## Benford's Law

Newcomb (1881).  Tables of logarithms.

Benford (1938).
- River area.      Population.
  Newspaper.    Specific heat.
  Pressure.       Atomic weight.
  Drainage.       Reader's Digest.
  Baseball.        Black body.
  Death rates.    Addresses.

- Scale invariant!

Hill (1996).
- Distribution of distributions.

```
% more pu-files.txt
96796
4171208          ← 332,952 file sizes
5830
34343656
...

% java Benford < pu-files.txt
1: 0.307882217256544147
2: 0.192502222254258872
3: 0.1302139647757034
4: 0.098659866688771955
5: 0.07445217328623946
6: 0.05945601768423076
7: 0.05162606021288354
8: 0.04417153223287441
9: 0.04103594512121867
```
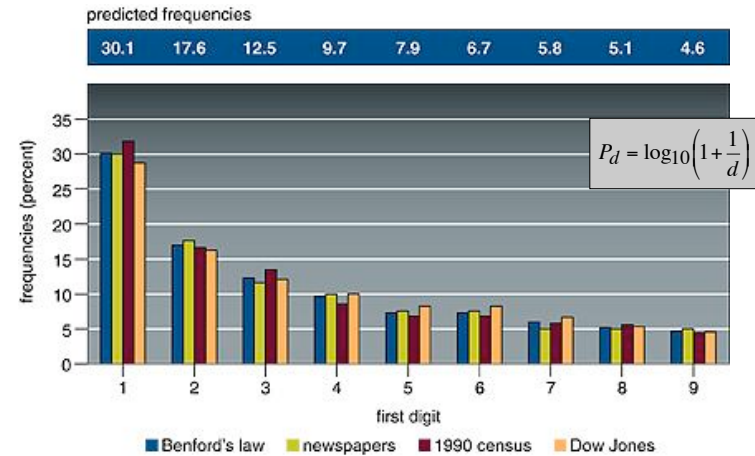
---

## The First-Digit Phenomenon



| predicted frequencies | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 30.1 | 17.6 | 12.5 | 9.7 | 7.9 | 6.7 | 5.8 | 5.1 | 4.6 |

$$P_d = \log_{10}\left(1+\frac{1}{d}\right)$$

■ Benford's law   ■ newspapers   ■ 1990 census   ■ Dow Jones

Reference: *The First-Digit Phenomenon* by T. P. Hill, in *American Scientist*, July-August 1998.

---

## Sorting

Goal:  given N items, rearrange them in increasing order.

Applications.
- Sort a list of names.
- Find duplicates in a mailing list.
- Find the median.
- Identify statistical outliers.
- Data compression.
- Computer graphics.
- Computational biology.

| name |
|---|
| Hauser |
| Hong |
| Hsu |
| Hayes |
| Haskell |
| Hanley |
| Hornet |

⇒

| name |
|---|
| Hanley |
| Haskell |
| Hauser |
| Hayes |
| Hong |
| Hornet |
| Hsu |

---

## Insertion Sort

Insertion sort an N-element array.
- In $i^{th}$ iteration:
  - read $i^{th}$ value
  - repeatedly swap $i^{th}$ value with the one to its left if smaller

Property:  after $i^{th}$ iteration, array positions 0 through i contain original elements 0 through i in increasing order.

```java
for (int i = 0; i < N; i++) {
    for (int j = i; j > 0; j--) {
        if (x[j-1] > x[j]) {
            double swap = x[j];      // swap x[j] and x[j-1]
            x[j] = x[j-1];
            x[j-1] = swap;
        }
    }
}
```

sort array of real numbers

## Linear System of Equations

Linear system of equations.
- N linear equations in N unknowns.
- Matrix notation:  find x such that Ax = b.

$$
\begin{aligned}
0\,x_0 + 1\,x_1 + 1\,x_2 &= 4 \\
2\,x_0 + 4\,x_1 - 2\,x_2 &= 2 \\
0\,x_0 + 3\,x_1 - 15\,x_2 &= 36
\end{aligned}
\qquad
A = \begin{bmatrix} 0 & 1 & 1 \\ 2 & 4 & -2 \\ 0 & 3 & -15 \end{bmatrix}, \quad b = \begin{bmatrix} 4 \\ 2 \\ 36 \end{bmatrix}
$$

Among most fundamental problems in science and engineering.
- Linear regression.
- Kirchoff's current law.
- Polynomial and spline interpolation.
- Linear and nonlinear optimization.
- Numerical solution to differential equations.
- Fluid flow,
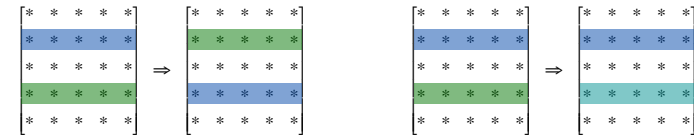- Leontief model of economic equilibrium.

---

## Gaussian Elimination

Gaussian elimination.
- Among oldest and most widely used solutions.
- Repeatedly apply row operations until system is *upper triangular*.
- Solve "trivial" upper triangular system.

Row operations.
- Exchange any two rows.
- Add a multiple of one row to another.

Key invariant.   Row operations preserve solutions.

---

## Gaussian Elimination:  Forward Elimination

$$
\begin{aligned}
0\,x_0 + 1\,x_1 + 1\,x_2 &= 4 \\
2\,x_0 + 4\,x_1 - 2\,x_2 &= 2 \\
0\,x_0 + 3\,x_1 + 15\,x_2 &= 36
\end{aligned}
$$

Interchange row 0 and 1.

$$
\begin{aligned}
2\,x_0 + 4\,x_1 - 2\,x_2 &= 2 \\
0\,x_0 + 1\,x_1 + 1\,x_2 &= 4 \\
0\,x_0 + 3\,x_1 + 15\,x_2 &= 36
\end{aligned}
$$

Subtract 3x row 1 from row 2.

$$
\begin{aligned}
2\,x_0 + 4\,x_1 - 2\,x_2 &= 2 \\
0\,x_0 + 1\,x_1 + 1\,x_2 &= 4 \\
0\,x_0 + 0\,x_1 + 12\,x_2 &= 24
\end{aligned}
$$

---

## Gaussian Elimination:  Back Substitution

Back substitution.  *Upper triangular* systems are easy to solve.

$$
\begin{aligned}
2\,x_0 + 4\,x_1 - 2\,x_2 &= 2 \\
0\,x_0 + 1\,x_1 + 1\,x_2 &= 4 \\
0\,x_0 + 0\,x_1 + 12\,x_2 &= 24
\end{aligned}
$$

- Equation 2:  $x_2 = 24/12 = 2$.
- Equation 1:  $x_1 = 4 - x_2 = 2$.
- Equation 0:  $x_0 = 2 - 4x_1 + 2x_2 = -1$.

```java
for (int i = N-1; i >= 0; i--) {
    double sum = 0.0;
    for (int k = i+1; k < N; k++)
        sum += A[i][k] * x[k];
    x[i] = (b[i] - sum) / A[i][i];
}
```

$$
x_i = \frac{1}{A_{ii}}\left[ b_i - \sum_{i+1}^{N-1} A_{ik}\, x_k \right]
$$

## Gaussian Elimination: Forward Elimination

Forward elimination. Apply row operations to make upper triangular.

Pivot. Zero out entries below pivot $A_{ii}$.

$$A_{jk} = A_{jk} - \frac{A_{ji}}{A_{ii}} A_{ik}$$

$$b_j = b_j - \frac{A_{ji}}{A_{ii}} b_i$$

$$
\begin{bmatrix}
* & * & * & * & * & * \\
0 & * & * & * & * & * \\
0 & 0 & * & * & * & * \\
0 & 0 & * & * & * & * \\
0 & 0 & * & * & * & * \\
0 & 0 & * & * & * & *
\end{bmatrix}
\Rightarrow
\begin{bmatrix}
* & * & * & * & * & * \\
0 & * & * & * & * & * \\
0 & 0 & * & * & * & * \\
0 & 0 & 0 & * & * & * \\
0 & 0 & 0 & * & * & * \\
0 & 0 & 0 & * & * & *
\end{bmatrix}
$$

```
for (int i = 0; i < N; i++) {
   for (int j = i + 1; j < N; j++)
      b[j] -= (A[j][i] / A[i][i]) * b[i];
   for (int j = i + 1; j < N; j++)
      for (int k = N - 1; k >= i; k--)
         A[j][k] -= (A[j][i] / A[i][i]) * A[i][k];
}
```

---

## Gaussian Elimination: Partial Pivoting

Observation. Previous code fails spectacularly if pivot $A_{ii}$ = 0.

Partial pivoting. Swap row i with the row that has biggest entry in column i among unreduced rows j > i.

$$
\begin{bmatrix}
* & * & * & * & * & * \\
0 & * & * & * & * & * \\
0 & 0 & 0 & * & * & * \\
0 & 0 & 3 & * & * & * \\
0 & 0 & 9 & * & * & * \\
0 & 0 & 2 & * & * & *
\end{bmatrix}
$$

```
// find pivot row
int max = i;
for (int j = i + 1; j < N; j++)
   if (Math.abs(A[j][i]) > Math.abs(A[max][i]))
      max = j;

// swap rows i and max of A and b
double[] T = A[i]; A[i] = A[max]; A[max] = T;
double   t = b[i]; b[i] = b[max]; b[max] = t;
```

---

## Gaussian Elimination: Pathologies

Degeneracy. Partial pivot on a value close to zero.
- System is overdetermined: no solutions.
- System is underdetermined: many solutions.

Numerical stability. Floating point roundoff error swamps computation.
- Partial pivoting helps control roundoff error.
- Pathological instances exist that blow up partial pivoting.

Ill-conditioning. Some problems are inherently unsuitable for floating point solution techniques.
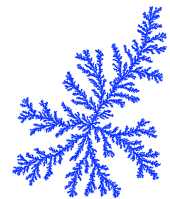
Scientific computing. Much of hard work in designing numerical algorithms is addressing such pathologies.

---

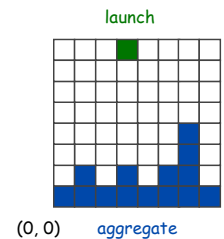## Diffusion Limited Aggregation

Diffusion limited aggregation (DLA).
- Models formation of an aggregate on a surface.
  - growth of lichen on rocks
  - growth of coral reef
  - generation of polymers out of solutions
  - path of electrical discharge
  - urban settlement
  - carbon deposits on walls of a cylinder of Diesel engine

Monte Carlo simulation.
- Launch particle from *launch site*.
- Particle randomly wanders through 2-D grid until
  - it comes in contact with another particle ⇒ sticks to *aggregate*
  - it enters *kill zone*
- Repeat.

launch

(0, 0)    aggregate

## Two Dimensional Arrays in Java

Two dimensional arrays.

```
double a00, a01, a02, a03, a10, a11, a12, a13;
```

vs.

```
double[][] a = new double[2][4];
```

- To access element (i, j) of array named `a`, use `a[i][j]`.

```java
public class DLA {
    public static void main(String[] args) {
        int N = Integer.parseInt(args[0]);          ← N-by-N grid
        StdDraw.create(N, N);
        int launch = N - 10;                          ← launch near top row

        boolean[][] dla = new boolean[N][N];          ← is site i-j occupied?
        for (int x = 0; x < N; x++)
            dla[x][0] = true;                          ← initialize aggregate at row 0
```

25

## Diffusion Limited Aggregation

```java
while (!done) {
    int x = (int) (N * Math.random());          ← launch from random column near top row
    int y = launch;

    while (x < N - 2 && x > 1 && y < N - 2 && y > 1) {
        double r = Math.random();
        if       (r < 0.25) x--;                 particle not in kill zone
        else if  (r < 0.50) x++;                 ← random step
        else if  (r < 0.65) y++;
        else                y--;

        if (dla[x-1][y] || dla[x+1][y] || dla[x][y-1] || dla[x][y+1]) {
            dla[x][y] = true;
            StdDraw.go(x, y);                    check for contact with
            StdDraw.spot();                      a neighboring particle
            StdDraw.pause(10);
            if (y > launch) done = true;         aggregate reaches top
            break;                               
        }                                        breaks out of innermost while loop
    }
}
```
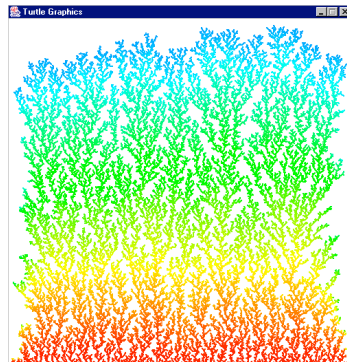
26

## Diffusion Limited Aggregation

Refinements.
- Use diagonals as neighbors, instead of just horizontal and vertical.
- Color particles in launch order, according to rainbow.

```
% java DLA 500
```



27

## Summary

Arrays.
- Organized way to store huge quantities of data.
- Almost as easy to use as primitive types.
- Can directly access an element given its index.

Caveats:
- Need to fix size of array ahead of time.
- Don't forget to allocate memory with `new`.
- Indices start at 0 not 1.
- Out-of-bounds to access `a[-1]` or `a[N]` of N element array.
    - in Java: `ArrayIndexOutOfBoundsException`
    - in C: "ghastly error"

"You're always off by 1 in this business."  - J. Morris

28