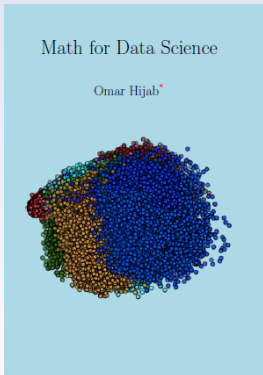# Mathematics for Data Science

Dr. S. M. Moosavi

*smohsenmoosavi2009@gmail.com*

June 23, 2024

The following slides are arranged (with some modifications) based on the book "*Math for Data Science*" by "**Omar Hijab**".

You can follow me on Linkedin. Also, for course materials such as slides and the related python codes, see this Github repository.



Math for Data Science

Omar Hijab*

1 Data Sets

Math for Data

Dr. S. M.
Moosavi

Data Sets

Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

1 Data Sets

# What is a dataset?

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

### Definition 1.1

*Geometrically, a dataset is a sample of $N$ points $x_1, x_2, \cdots, x_N$ in $d$-dimensional space $\mathbb{R}^d$. Algebraically, a dataset is an $N \times d$ matrix.*

Practically speaking, the following are all representations of datasets:

$$\text{matrix} = \text{CSV file} = \text{spreadsheet} = \text{SQL table} = \text{array} = \text{dataframe}$$

### Definition 1.2

*Each point $x = (t_1, t_2, \cdots, t_d)$ in the dataset is a sample or an example, and the components $t_1, t_2, \cdots, t_d$ of a sample point $x$ are its features or attributes. As such, $d$-dimensional space $\mathbb{R}^d$ is feature space.*

### Definition 1.3

*Sometimes one of the features is separated out as the label. In this case, the dataset is a labelled dataset.*

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

The *Iris dataset* contains 150 examples of four features of Iris flowers, and there are three classes of Irises, *Setosa*, *Versicolor* and *Virginica*, with 50 samples from each class.



**Samples**
(instances, observations)

**Petal**

| | Sepal length | Sepal width | Petal length | Petal width | Class label |
|---|---|---|---|---|---|
| 1 | 5.1 | 3.5 | 1.4 | 0.2 | Setosa |
| 2 | 4.9 | 3.0 | 1.4 | 0.2 | Setosa |
| ... | | | | | |
| 50 | 6.4 | 3.5 | 4.5 | 1.2 | Versicolor |
| ... | | | | | |
| 150 | 5.9 | 3.0 | 5.0 | 1.8 | Virginica |

**Sepal**

**Class labels**
(targets)

**Features**
(attributes, measurements, dimensions)

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
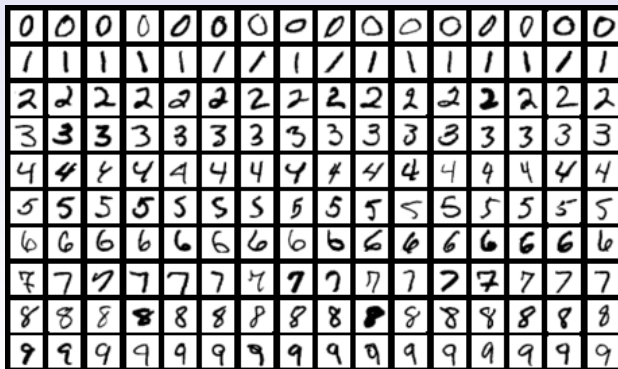Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# MNIST dataset

The *MNIST dataset* consists of $60,000$ images of hand-written digits. There are 10 classes of images, corresponding to each digit $0, 1, \cdots, 9$. We seek to compress the images while preserving as much as possible of the images' characteristics.

Each image is a grayscale $28 \times 28$ pixel image. Since $28^2 = 784$, each image is a point in $d = 784$ dimensions. Here there are $N = 60000$ samples and $d = 784$ features.

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
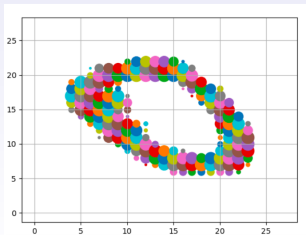Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Exercises

## Exercise 1.1

*Use* sklearn *to download Iris dataset.*

## Exercise 1.2

- *From* keras *read the MNIST dataset.*
- *Let* (train_X, train_y), (test_X, test_y) = mnist.load_data()
- *Let* pixels = train_X[1].
- *Do* for *loops over* i *and* j *in* range(28) *and use* scatter *to plot points at location* (i,j) *with size given by* pixels[i,j]*, then show the following image.*

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Introduction

Suppose we have a population of things (people, tables, numbers, vectors, images, etc.) and we have a sample of size $N$ from this population:

```
l = [x_1,x_2,...,x_N]
```

The total population is the *population* or the *sample space*.

### Example 1.1

The sample space consists of all real numbers and we take $N = 5$ samples from

```
l = [3.95, 3.20, 3.10, 5.55, 6.93]
```

### Example 1.2

The sample space consists of all integers and we take $N = 5$ samples from

```
l = [35, -32, -8, 45, -8]
```

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Introduction

### Example 1.3

The sample space consists of all Python strings and we take $N = 5$ samples from

```
l = ['a2e?','#%T','7y5,','kkk>><</','[[)*+']
```

### Example 1.4

The sample space consists of all HTML colors and we take $N = 5$ samples from

```
1  from random import choice
2  import matplotlib.pyplot as plt
3
4  def hexcolor():
5    return "#" + ''.join([choice('0123456789abcdef') for
                                 _ in range(6)])
6
7  for i in range(5): plt.scatter(i,0, c=hexcolor())
8  plt.show()
```

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Mean

Let $l$ be a list as above. The goal is to compute the sample *average* or *mean* of the list, which is

$$mean = average = \frac{x_1 + x_2 + \cdots + x_N}{N}.$$

In the Example (1.1), the average is

$$\frac{3.95 + 3.20 + 3.10 + 5.55 + 6.93}{5} = 4.546.$$

### Example 1.5

```
1   import numpy as np
2
3   dataset = np.array([3.95, 3.20, 3.10, 5.55, 6.93])
4   print(np.mean(dataset))
5
6   output: 4.546
```

In the Example (1.2), the average is $\frac{32}{5}$. In the Example (1.3), while we can add strings, we can't divide them by 5, so the average is undefined. Similarly for colors: the average is undefined.

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Vector space

A sample space or population $V$ is called a *vector space* if, roughly speaking, one can compute means or averages in $V$. In this case, we call the members of the population "*vectors*".

### Definition 1.4 (Vector space)

*Let $V$ be a set. $V$ is a vector space (over $\mathbb{R}$) if for every $u, v, w \in V$ and $r, s \in \mathbb{R}$:*

1 *vectors can be added (and the sum $v + w$ is back in $V$);*

2 *vector addition is commutative $v + w = w + v$*

3 *vector addition is associative $u + (v + w) = (u + v) + w$;*

4 *there is a zero vector $0$ $(0 + v = v)$;*

5 *vectors $v$ have negatives (or opposites) $-v$ $(v + (-v) = 0)$;*

6 *vectors can be multiplied by real numbers (and the product $rv$ is back in $V$);*

7 *multiplication is distributive over addition $(r + s)v = rv + sv$ and $r(u + v) = ru + rv$;*

8 *$1v = v$ and $0v = 0$;*

9 *$r(sv) = (rs)v$.*

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
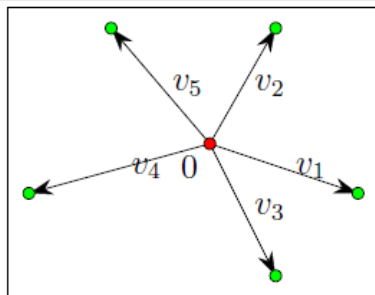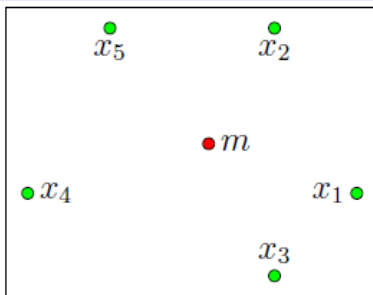Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Centered dataset

### Definition 1.5 (Centered Versus Non-Centered)

If $x_1, x_2, \cdots, x_N$ is a dataset of points with mean $m$ and

$$v_1 = x_1 - m, v_2 = x_2 - m, \cdots, v_N = x_N - m,$$

then $v_1, v_2, \cdots, v_N$ is a centered dataset of vectors where its mean is zero.
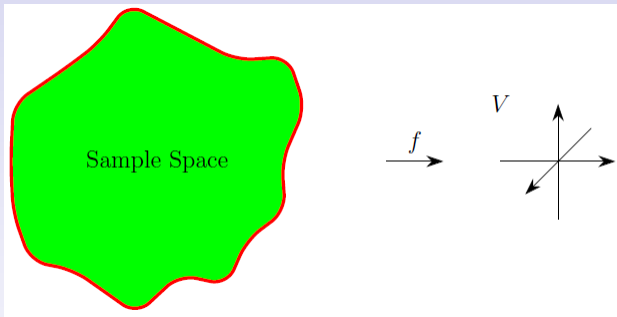
Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Some notes

- When we work with vector spaces, numbers are referred to as *scalars*.
- When we multiply a vector $v$ by a scalar $r$ to get the scaled vector $rv$, we call it *scalar multiplication*.
- The set of all real numbers $\mathbb{R}$ is a vector space.
- The set of all integers $\mathbb{Z}$ is not a vector space.
- The set of all rational numbers $\mathbb{Q}$ is a vector space.
- The set of all Python strings is not a vector space.
- Usually, we can't take sample means from a population, we instead take the sample mean of a *statistic* associated to the population. A statistic is an assignment of a number $f(\text{item})$ to each item in the population. For example, the human population on Earth is not a vector space (they can't be added), but their heights is a vector space (heights can be added). For the Python strings, a statistic might be the length of the strings. For the HTML colors, a statistic is the HTML code of the color.

In general, a statistic need not be a number. A statistic can be anything that "behaves like a number". For example, $f(\text{item})$ can be a vector or a matrix. More generally, a statistic's values may be anything that lives in a vector space $V$.
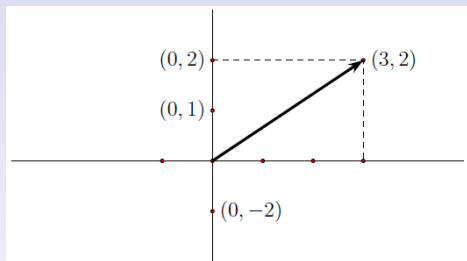
Math for Data

Dr. S. M. Moosavi

Data Sets
  Introduction
  Averages and Vector Spaces
  Two Dimensions
  Complex Numbers
  Mean and Covariance

# Cartesian plane

The *cartesian plane* $\mathbb{R}^2$, also called 2-dimensional real space is a vector space.



For $\mathbf{v}_1 = (x_1, y_1), \mathbf{v}_2 = (x_2, y_2) \in \mathbb{R}^2$ and $t \in \mathbb{R}$ define

- $\mathbf{v}_1 + \mathbf{v}_2 = (x_1 + x_2, y_1 + y_2)$ (Addition).
- $\mathbf{0} = (0, 0)$ (Zero).
- $t\mathbf{v}_1 = (tx_1, ty_1)$ (Scaling).
- $-\mathbf{v}_1 = (-1)\mathbf{v}_1$ (Negative).
- $\mathbf{v}_1 - \mathbf{v}_2 = \mathbf{v}_1 + (-\mathbf{v}_2) = (x_1 - x_2, y_1 - y_2)$ (Subtraction).

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
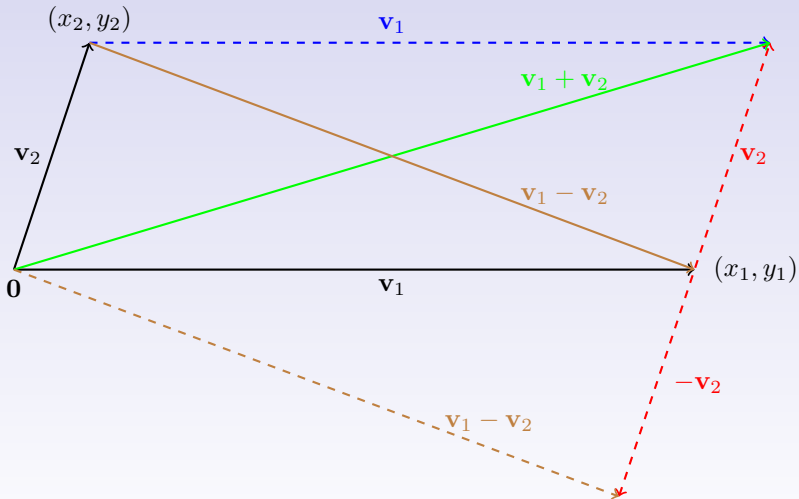Averages and Vector
Spaces
**Two Dimensions**
Complex Numbers
Mean and Covariance

# 2d example

### Example 1.6

```
1   import numpy as np
2
3   v1 = (1,2)
4   v2 = (3,4)
5   print(v1 + v2 == (1+3,2+4)) # returns False
6
7   v1 = [1,2]
8   v2 = [3,4]
9   print(v1 + v2 == [1+3,2+4]) # returns False
10
11  v1 = np.array([1,2])
12  v2 = np.array([3,4])
13  print(v1 + v2 == np.array([1+3,2+4]))
14  # returns [ True   True]
15  print(3*v1 == np.array([3,6]))
16  # returns [ True   True]
17  print(-v1 == np.array([-1,-2]))
18  # returns [ True   True]
19  print(v1 - v2 == np.array([1-3,2-4]))
20  # returns [ True   True]
```

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

## 2d example

For the two-dimensional dataset

$$\mathbf{x}_1 = (1, 2), \mathbf{x}_2 = (3, 4), \mathbf{x}_3 = (-2, 11), \mathbf{x}_4 = (0, 66),$$

or, equivalently,

$$\mathbf{x} = \begin{pmatrix} 1 & 2 \\ 3 & 4 \\ -2 & 11 \\ 0 & 66 \end{pmatrix},$$

the average is

$$\frac{(1, 2) + (3, 4) + (-2, 11) + (0, 66)}{4} = (0.5, 20.75).$$

### Example 1.7

```python
import numpy as np

dataset = np.array([[1,2], [3,4], [-2,11], [0,66]])
print(np.mean(dataset, axis=0))
# returns [ 0.5 , 20.75]
```

# 2d example

## Example 1.8

Generate a 2 dimensional dataset of random points and their mean

```
1   import numpy as np
2   from numpy.random import random as rd
3   import matplotlib.pyplot as plt
4   N = 20
5   dataset = np.array([[rd(), rd()] for _ in range(N)])
6   mean = np.mean(dataset,axis=0)
7   plt.grid()
8   X, Y = dataset[:,0], dataset[:,1]
9   plt.scatter(X,Y)
10  plt.scatter(*mean)
11  plt.annotate('$m$', xy=mean+0.01)
12  plt.show()
```

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Magnitude

## Definition 1.6 (Distance Formula)

If $\mathbf{v}_1 = (x_1, y_1)$ and $\mathbf{v}_2 = (x_2, y_2)$, then the distance between $\mathbf{v}_1$ and $\mathbf{v}_2$ is

$$|\mathbf{v}_1 - \mathbf{v}_2| = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}.$$

The distance of $\mathbf{v} = (x, y)$ to the origin $\mathbf{0} = (0, 0)$ is its magnitude or norm or length

$$r = |\mathbf{v}| = |\mathbf{v} - \mathbf{0}| = \sqrt{x^2 + y^2}.$$

## Example 1.9

For $\mathbf{v}_1 = (1, 2)$ and $\mathbf{v}_2 = (3, 4)$

$$|\mathbf{v}_1| = \sqrt{1^2 + 2^2} = \sqrt{5} \simeq 2.236,$$

$$|\mathbf{v}_1 - \mathbf{v}_2| = \sqrt{(1-3)^2 + (2-4)^2} = \sqrt{4+4} = \sqrt{8} \simeq 2.828.$$

```
1  import numpy as np
2
3  v1 = np.array([1,2])
4  v2 = np.array([3,4])
5  print(np.linalg.norm(v1)) #returns 2.23606797749979
6  print(np.linalg.norm(v1-v2)) #returns 2.
                                  8284271247461903
```

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Polar representation

In terms of $r$ and $\theta$, the *polar representation* of $(x, y)$ is

$$x = r \cos \theta, \quad y = r \sin \theta.$$



The *unit circle* consists of the vectors which are distance 1 from the origin $\mathbf{0}$. When $\mathbf{v}$ is on the unit circle, the magnitude of $\mathbf{v}$ is 1, and we say $\mathbf{v}$ is a *unit vector*. In this case, the line formed by the scalings of $\mathbf{v}$ intersects the unit circle at $\pm \mathbf{v}$.

When $\mathbf{v}$ is a unit vector, then $r = 1$ and $\mathbf{v} = (x, y) = (\cos \theta, \sin \theta)$.

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Polar representation

By the distance formula, a vector $\mathbf{v} = (x, y)$ is a unit vector when

$$x^2 + y^2 = 1.$$

More generally, any circle with *center* $(a, b)$ and radius $r$ consists of vectors $\mathbf{v} = (x, y)$ satisfying

$$(x - a)^2 + (y - b)^2 = r^2.$$

Let $R$ be a point on the unit circle, and let $t > 0$. The scaled point $t\mathsf{R}$ is on the circle with center $(0, 0)$ and radius $t$. Moreover, if $Q$ is any point, $Q + tR$ is on the circle with center $Q$ and radius $t$. It is easy to check that $|t\mathbf{v}| = |t||\mathbf{v}|$ for any real number $t$ and vector $\mathbf{v}$.

From this, if a vector $\mathbf{v}$ is unit and $r > 0$, then $r\mathbf{v}$ has magnitude $r$. If $\mathbf{v}$ is any vector not equal to the zero vector, then $r = |\mathbf{v}|$ is positive, and

$$\left| \frac{1}{r}\mathbf{v} \right| = \frac{1}{r}|\mathbf{v}| = \frac{1}{r}r = 1$$

so $\mathbf{v}/r$ is a unit vector.

# Inner product

## Definition 1.7

Let $\mathbf{v}_1 = (x_1, y_1), \mathbf{v}_2 = (x_2, y_2) \in \mathbb{R}^2$. The inner product or the dot product of $\mathbf{v}_1$ and $\mathbf{v}_2$ is given algebraically as

$$\mathbf{v}_1 \cdot \mathbf{v}_2 = x_1 x_2 + y_1 y_2.$$

From the geometric view, we have:

## Theorem 1.1 (Dot Product Identity)

$$x_1 x_2 + y_1 y_2 = \mathbf{v}_1 \cdot \mathbf{v}_2 = |\mathbf{v}_1||\mathbf{v}_2| \cos \theta,$$

where $\theta$ is the angle between $\mathbf{v}_1$ and $\mathbf{v}_1$.

## Exercise 1.3

Prove the "Dot Product Identity", Theorem (1.1).
Hint: Use Pythagoras' theorem for general triangles.

Math for Data

Dr. S. M. Moosavi

Data Sets
Introduction
Averages and Vector Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# The angle between two vectors

In Python, the dot product is given by numpy.dot and as a consequence of the dot product identity, we have the code for the angle between two vectors:

$$\theta_{\mathbf{v}_1, \mathbf{v}_2} = \arccos\left(\frac{\mathbf{v}_1 \cdot \mathbf{v}_2}{|\mathbf{v}_1||\mathbf{v}_2|}\right).$$

### Example 1.10

Find the angle between the vectors $\mathbf{v}_1 = (1, 2)$ and $\mathbf{v}_2 = (3, 4)$.

```python
1  import numpy as np
2
3  def angle(u,v):
4    a = np.dot(u,v)
5    b = np.dot(u,u)
6    c = np.dot(v,v)
7    theta = np.arccos(a / np.sqrt(b*c))
8    return np.degrees(theta)
9
10 v1 = np.array([1,2])
11 v2 = np.array([3,4])
12 print(angle(v1,v2)) #returns 10.304846468766044 in
                                degree
```

# Cauchy-Schwarz Inequality

Recall that $-1 \leq \cos \theta \leq 1$. Using the dot product identity, we obtain the important inequality:

### Theorem 1.2 (Cauchy-Schwarz Inequality)

*If $u$ and $v$ are any two vectors, then*

$$-|u||v| \leq u \cdot v \leq |u||v|.$$

### Exercise 1.4

*Prove the "Cauchy-Schwarz Inequality".*

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# 2d linear equations system

Consider the homogeneous system

$$\begin{cases} ax + by = 0 \\ cx + dy = 0 \end{cases} \tag{1.1}$$

and let $A$ be the $2 \times 2$ matrix

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}. \tag{1.2}$$

$(x, y) = (-b, a)$ is a solution of the first equation in (1.1). If we want
this to be a solution of the second equation as well, we must have
$cx + dy = ad - bc = 0$.

### Definition 1.8 (Determinant)

*The determinant of A is*

$$\det(A) = \det \begin{pmatrix} a & b \\ c & d \end{pmatrix} = ad - bc.$$

Math for Data

Dr. S. M. Moosavi

Data Sets
Introduction
Averages and Vector Spaces
**Two Dimensions**
Complex Numbers
Mean and Covariance

# 2d linear equations system

### Theorem 1.3 (Homogeneous System)

*When $\det(A) = 0$, the homogeneous system (1.1) has a nonzero solution, and all solutions are scalar multiples of $(x, y) = (-b, a)$. When $\det(A) \neq 0$, the only solution is $(x, y) = (0, 0)$.*

For the inhomogeneous case

$$\begin{cases} ax + by = e \\ cx + dy = f \end{cases} \tag{1.3}$$

we have

### Theorem 1.4 (Inhomogeneous System)

*When $\det(A) \neq 0$, the inhomogeneous system (1.3) has the unique solution*

$$\begin{pmatrix} x \\ y \end{pmatrix} = \frac{1}{\det(A)} \begin{pmatrix} de - bf \\ af - ce \end{pmatrix}.$$

*When $\det(A) = 0$, (1.3) has a solution iff $ce = af$ and $de = bf$.*

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

## 2d linear equations system

When $a^2 + b^2 \neq 0$, a solution is

$$\begin{pmatrix} x \\ y \end{pmatrix} = \frac{1}{a^2 + b^2} \begin{pmatrix} ae \\ be \end{pmatrix}.$$

When $c^2 + d^2 \neq 0$, a solution is

$$\begin{pmatrix} x \\ y \end{pmatrix} = \frac{1}{c^2 + d^2} \begin{pmatrix} cf \\ df \end{pmatrix}.$$

Any other solution differs from these solutions by a scalar multiple of the homogeneous solution $(x, y) = (-b, a)$.

### Exercise 1.5

*Prove the Theorems (1.3) and (1.4).*

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Complex numbers

Roughly speaking, the set of all *complex numbers* is the set of all points in $\mathbb{R}^2$ with different multiplication rule.

## Definition 1.9 (Complex numbers)

*The complex numbers,$\mathbb{C}$, is the set*

$$\mathbb{C} = \{(x, y) \in \mathbb{R}^2\}$$

*with operations*

- *Addition:* $(x_1, y_1) + (x_2, y_2) = (x_1 + x_2, y_1 + y_2)$.
- *Scalar Multiplication:* $t(x, y) = (tx, ty)$
- *Multiplication:* $(x_1, y_1)(x_2, y_2) = (x_1 x_2 - y_1 y_2, x_1 y_2 + x_2 y_1)$.

Then, in $\mathbb{C}$, we have

- zero: $0 = (0, 0)$.
- opposite or additive inverse: $-(x, y) = (-x, -y)$.
- one: $1 = (1, 0)$.

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Example

### Example 1.11

- $(1,2) + (3,4) = (4,6)$.
- $(0,0) + (1,2) = (1,2)$.
- $3(1,2) = (3,6)$.
- $(1,0)(1,2) = (1-0, 2+0) = (1,2)$.
- $(1,2)(3,4) = (3-8, 4+6) = (-5,10)$.
- $(x,0) + (y,0) = (x+y, 0)$.
- $(x,0)(y,0) = (xy, 0)$.

**Note**. By the last two examples, we see that complex numbers with $0$ as their second component act like real numbers in addition and multiplication. So, from now on, we set $x = (x,0)$.

### Example 1.12

- $0 = (0,0)$.
- $1 = (1,0)$.
- $-1 = (-1,0)$.

# Imaginary number

### Definition 1.10 (Imaginary number)

$i = (0, 1)$.

**Note**. Python uses the symbol $j$ for imaginary number.

### Theorem 1.5

*For each $z = (x, y) \in \mathbb{C}$, we can write*

$$z = x + iy.$$

*We call $x$ as the real part of $z$, and $y$ the imaginary part of $z$.*

$$x = Re(z), \quad y = Im(z).$$

**Proof**. $x + iy = (x, 0) + (0, 1)(y, 0) = (x, 0) + (0 - 0, 0 + y) = (x, y)$.

### Theorem 1.6

$i^2 = -1$.

**Proof**. $i^2 = (0, 1)(0, 1) = (0 - 1, 0 + 0) = (-1, 0) = -1$.

Math for Data

Dr. S. M. Moosavi

Data Sets
Introduction
Averages and Vector Spaces
Two Dimensions
**Complex Numbers**
Mean and Covariance

# Example

## Example 1.13

In complex numbers:

- $\sqrt{-1} = i$.
- $\sqrt{-4} = 2i$.
- $(1,2)(3,4) = (1+2i)(3+4i)$

$$= 3 + 4i + 6i + 8i^2$$
$$= 3 + 10i - 8$$
$$= -5 + 10i$$
$$= (-5, 10).$$

- $(1,2)^3 = (1+2i)^3$

$$= (1)^3 + 3(1)^2(2i) + 3(1)(2i)^2 + (2i)^3$$
$$= 1 + 6i + 12i^2 + 8i^3$$
$$= 1 + 6i - 12 - 8i$$
$$= -11 - 2i$$
$$= -(11, 2).$$

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
**Complex Numbers**
Mean and Covariance

# Conjugate

## Definition 1.11 (Conjugate)

*For $z = (x, y) \in \mathbb{C}$, the conjugate is*

$$\bar{z} = (x, -y) = x - iy \in \mathbb{C}.$$

**Some properties**.

- $z + \bar{z} = 2Re(z)$, $z - \bar{z} = 2iIm(z)$.
- $z\bar{z} = Re(z)^2 + Im(z)^2$,

$$\Rightarrow |z| = \sqrt{Re(z)^2 + Im(z)^2} = \sqrt{z\bar{z}}$$

$$\Rightarrow |z|^2 = z\bar{z}.$$

## Example 1.14

For $z = (4, -3) \in \mathbb{C}$:

- $\bar{z} = (4, 3) = 4 + 3i$,
- $z + \bar{z} = 2 \times 4 = 8$, $z - \bar{z} = 2i \times (-3) = -6i$.
- $z\bar{z} = (4)^2 + (-3)^2 = 16 + 9 = 25 \Rightarrow |z| = \sqrt{25} = 5$.
- $z^2 = (4 - 3i)^2 = 7 - 24i$.
- $|z|^2 = 25$.

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Inverse

## Theorem 1.7

*For a non-zero $z \in \mathbb{C}$, the inverse of $z$ is*

$$z^{-1} = \frac{1}{z} = \frac{\bar{z}}{z\bar{z}} = \frac{\bar{z}}{|z|^2}.$$

**Proof**. Firstly, if $z = (x, y)$ then $\frac{1}{z} \in \mathbb{C}$, because,

$$\frac{1}{z} = \frac{x - iy}{x^2 + y^2} = (\frac{x}{x^2 + y^2}, \frac{-y}{x^2 + y^2}) \in \mathbb{C}.$$

Secondly,

$$zz^{-1} = (x + iy)\left(\frac{x - iy}{x^2 + y^2}\right) = \frac{x^2 + y^2}{x^2 + y^2} = 1.$$

## Corollary 1.1 (Division)

*For $z_1 \in \mathbb{C}$ and $0 \neq z_2 \in \mathbb{C}$*

$$\frac{z_1}{z_2} = z_1 z_2^{-1}.$$

Math for Data

Dr. S. M. Moosavi

Data Sets

Introduction

Averages and Vector Spaces

Two Dimensions

Complex Numbers

Mean and Covariance

# Definitions

### Definition 1.12 (Mean-squared distance)

Let $x_1, x_2, \ldots, x_N$ be a dataset, say $D$, in $\mathbb{R}^d$, and let $x \in \mathbb{R}^d$. The mean-square distance of $x$ to $D$ is

$$MSD(x) = \frac{1}{N} \sum_{k=1}^{N} |x_k - x|^2.$$

### Definition 1.13 (Mean)

Let $x_1, x_2, \ldots, x_N$ be a dataset in $\mathbb{R}^d$. The mean or sample mean is

$$m = \bar{x}_N = \frac{1}{N} \sum_{k=1}^{N} x_k = \frac{x_1 + x_2 + \cdots + x_N}{N}.$$

### Theorem 1.8 (Point of Best-fit)

The mean is the point of best-fit: The mean minimizes the mean-squared distance to the dataset.

### Exercise 1.6

Prove the Theorem (1.8).

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Point of Best-fit

## Example 1.15

```python
1  import matplotlib.pyplot as plt
2  import numpy as np
3
4  np.random.seed(1)
5  N = 20
6  dataset = np.array([ [np.random.random(),np.random.random()]
                                      for _ in range(N) ])
7  # Mean
8  m = np.mean(dataset, axis=0)
9  #Random point
10 p = np.array([np.random.random(),np.random.random()])
11
12 plt.grid()
13 X, Y = dataset[:,0], dataset[:,1]
14 plt.scatter(X,Y)
15 for v in dataset:
16    plt.plot([m[0],v[0]],[m[1],v[1]],c='green')
17    plt.plot([p[0],v[0]],[p[1],v[1]],c='red')
18 plt.show()
19
20 # Comparison of MSD of the mean and a random point
21 MSD_m = np.sum(np.abs(dataset-m)**2)/N
22 MSD_p = np.sum(np.abs(dataset-p)**2)/N
23 print(MSD_m, MSD_p) # 0.160478187272121 0.5984208474157081
```

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Point of Best-fit



Figure 1.1: MSD for the mean (green) versus MSD for a random point (red).

# Tensor product

For simplicity, let $u = (a, b)$ and $v = (c, d, e)$ be two vectors.

### Definition 1.14 (Tensor product)

*The tensor product of $u$ and $v$ is the matrix*

$$u \otimes v = \begin{pmatrix} ac & ad & ae \\ bc & bd & be \end{pmatrix} = \begin{pmatrix} cu & du & eu \end{pmatrix} = \begin{pmatrix} av \\ bv \end{pmatrix}$$

### Definition 1.15 (Trace of a matrix)

*The trace of a squared matrix $A$ is the sum of the diagonal entries.*

**Note**. For any vectors $u, v$ and $w$:

- $v \otimes u = (u \otimes v)^t$.

  In square case:

- $\det(u \otimes v) = 0$.
- $trace(u \otimes v) = u \cdot v$.
- $trace(u \otimes u) = |u|^2$.
- $(u \otimes v)w = (v \cdot w)u$.

Math for Data

Dr. S. M.
Moosavi

Data Sets
  Introduction
  Averages and Vector
  Spaces
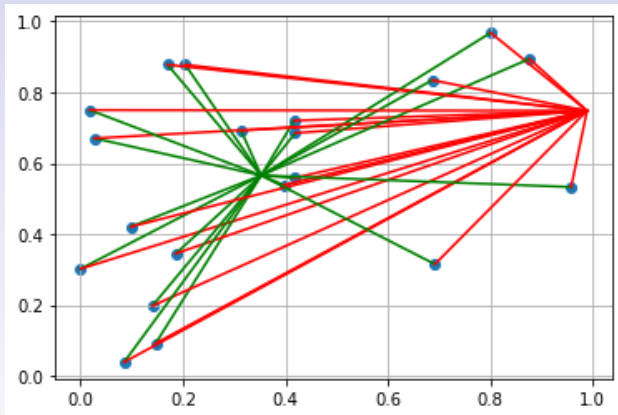  Two Dimensions
  Complex Numbers
  Mean and Covariance

# Covariance

Let $x_1, x_2, \ldots, x_N$ be a dataset in $\mathbb{R}^d$ with $m$ as its mean.

### Definition 1.16 (1d Covariance)

When $d = 1$, the covariance $q$ is a scalar

$$q = \frac{1}{N} \sum_{k=1}^{N} (x_k - m)^2 = MSD(m).$$

In the scalar case, the covariance is called the variance of the scalar dataset.

In general, the covariance is a symmetric $d \times d$ matrix $Q$. We can center the dataset as

$$v_1 = x_1 - m, v_2 = x_2 - m, \ldots, v_N = x_N - m.$$

Then the *covariance matrix* is the $d \times d$ matrix $Q$ as

$$Q = \frac{v_1 \otimes v_1 + v_2 \otimes v_2 + \ldots + v_N \otimes v_N}{N}. \tag{1.4}$$

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Example

## Example 1.16

Suppose $N = 5$ and

$$x_1 = (1, 2), \quad x_2 = (3, 4), \quad x_3 = (5, 6), \quad x_4 = (7, 8), \quad x_5 = (9, 10).$$

Then $m = (5, 6)$ and

$$v_1 = x_1 - m = (-4, -4), \quad v_2 = x_2 - m = (-2, -2),$$
$$v_3 = x_3 - m = (0, 0), \quad v_4 = x_4 - m = (2, 2), \quad v_5 = x_5 - m = (4, 4).$$

Since

$$(\pm 4, \pm 4) \otimes (\pm 4, \pm 4) = \begin{pmatrix} 16 & 16 \\ 16 & 16 \end{pmatrix},$$

$$(\pm 2, \pm 2) \otimes (\pm 2, \pm 2) = \begin{pmatrix} 4 & 4 \\ 4 & 4 \end{pmatrix},$$

$$(0, 0) \otimes (0, 0) = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix},$$

then

$$Q = \begin{pmatrix} 8 & 8 \\ 8 & 8 \end{pmatrix}.$$

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Example

## Example 1.17

```
1  import numpy as np
2
3  def tensor(u,v):
4    return np.array([ [ a*b for b in v] for a in u ])
5
6  np.random.seed(1)
7  N = 20
8  dataset = np.array([ [np.random.random(),np.random.
                            random()] for _ in range(N
                            ) ])
9  # mean
10 m = np.mean(dataset,axis=0)
11 # center dataset
12 vectors = dataset - m
13 # covariance
14 Q = np.mean([ tensor(v,v) for v in vectors ],axis=0)
15 print(Q)
```

Math for Data

Dr. S. M.
Moosavi

Data Sets
Introduction
Averages and Vector
Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Standardized

**Note**. The covariance matrix as written in (1.4) is the *biased covariance matrix*. If the denominator is instead $N - 1$, the matrix is the *unbiased covariance matrix*.

For datasets with large $N$, it doesn't matter, since $N$ and $N - 1$ are almost equal.

In `numpy`, the Python covariance constructor is

### Example 1.18

```
1  import numpy as np
2
3  np.random.seed(1)
4  N = 20
5  dataset = np.array([ [np.random.random(),np.random.
                                    random()] for _ in range(N
                                    ) ])
6  # covariance
7  Q = np.cov(dataset,bias=True,rowvar=False)
8  print(Q)
```

Math for Data

Dr. S. M. Moosavi

Data Sets
Introduction
Averages and Vector Spaces
Two Dimensions
Complex Numbers
Mean and Covariance

# Total variance

## Definition 1.17 (Total variance)

From $trace(u \otimes u) = |u|^2$, if $Q$ is the covariance matrix then

$$trace(Q) = \frac{1}{N} \sum_{k=1}^{N} |x_k - m|^2. \tag{1.5}$$

We call (1.5) the total variance of the dataset. Thus the total variance equals $MSD(m)$.

## Example 1.19

```
1  import numpy as np
2
3  np.random.seed(1)
4  N = 20
5  dataset = np.array([ [np.random.random(),np.random.
                             random()] for _ in range(N
                             ) ])
6  # covariance
7  Q = np.cov(dataset.T,bias=True)
8  print(Q.trace()) # returns 0.16047818727212101
```