

20210915

Wednesday, September 15, 2021 4:30 PM

<https://projecteuler.net/archives>

https://projecteuler.net/problem=18

Project Euler .net

About Archives Recent News Register Sign In

### Maximum path sum I

**Problem 18**

By starting at the top of the triangle below and moving to adjacent numbers on the row below, the maximum total from top to bottom is 23.

That is,  $3 + 7 + 4 + 9 = 23$ .

Find the maximum total from top to bottom of the triangle below:

75
95 64
17 47 82
18 35 87 10
20 04 82 47 65
19 01 23 75 03 34
88 02 77 73 07 63 67
99 65 04 28 06 16 70 92
41 41 26 56 83 40 80 70 33
41 48 72 33 47 32 37 16 94 29
53 71 44 65 25 43 91 52 97 51 14
70 11 33 28 77 73 17 78 39 68 17 57
91 71 52 38 17 14 91 43 58 50 27 29 48
63 66 04 68 89 53 67 30 73 16 69 87 40 31
04 62 98 27 23 09 70 98 73 93 38 53 60 04 23

**NOTE:** As there are only 16384 routes, it is possible to solve this problem by trying every route. However, **Problem 67**, is the same challenge with a triangle containing one-hundred rows; it cannot be solved by brute force, and requires a clever method! ;o)

Want to find the maximum path sum

Should be able to feed in the top triangle, and get 23



# Clean Code

Boston University CS 506 - Lance Galletti

## The Problem



## The Problem

"Software Systems get replaced not when they wear out but when they crumble under their own weight because they have become too complex"



## Why Does this Occur?

- english words + grammar != poetry
- Need analysis of the code itself - not just algorithms

## Gems of Clean Code

- Structure
  - Before writing code ask "How will someone use this (or part of this) code?". Minimize side effects
  - Do one thing
- Method
  - Top Down Approach
  - Bottom up Approach
  - Solve the problem first - then improve / refine
- General
  - Boring code is good code: keep it simple
  - Late Binding: start vague and refine (don't commit to specificity)
  - Many functions with small bodies > one function with large body
  - Check soundness by reading your code before testing

## Clean Code

'Clean code' by Robert C. Martin

Github Gist Summary of the Book:

<https://gist.github.com/wojteklu/73c6914cc446146b8b533c0988cf8d29>

# A Look into Functional Programming

"Functional Programs are mathematical expressions that are evaluated and reasoned about much like ordinary mathematical functions. As a result, these expressions are simple to analyze and compose for large-scale programs"



- Combinators on Lists
- An example coding exercise

## Combinators on Lists

$xs = [1, 4, 9, 16, 25]$

- $\text{map}( xs, \lambda(x) \Rightarrow f(x) ) = [f(1), f(4), f(9), f(16), f(25)]$

Ex:  
 $\text{map}( xs, \lambda(x) \Rightarrow x \% 2 ) = [1, 0, 1, 0, 1]$

- Foldleft - idea of consuming the list via an accumulator

Ex:

$\text{foldleft}( xs, \text{init}, \lambda(x, acc) \Rightarrow f(x, acc) )$

- $acc = \text{init}$
- $acc = f(1, acc)$
- $acc = f(4, acc)$

- $acc = f(9, acc)$
- $acc = f(16, acc)$
- $acc = f(25, acc)$

$\text{foldleft}( xs, 0, \lambda(x, acc) \Rightarrow x + acc )$

1. Acc = 0
2. Acc = 0 + 1 = 1
3. Acc = 1 + 4 = 5
4. Acc = 5 + 9 = 14
5. Acc = 14 + 16 = 30
6. Acc = 30 + 25 = 55
7. Return 55

## Combinators on Lists

$xs = [1, 4, 9, 16, 25] \quad ys = [1, 2, 3, 4, 5]$

- $\text{map2}( xs, ys, \lambda(x, y) \Rightarrow f(x, y) ) = [f(1, 1), f(4, 2), f(9, 3), f(16, 4), f(25, 5)]$

Ex:  
 $\text{map2}( xs, ys, \lambda(x, y) \Rightarrow x * y ) = [1, 8, 27, 64, 125]$

- Foldleft2 - consume both lists (like foldleft on zip)

Ex:  $\text{foldleft2}( xs, ys, 1, \lambda(x, y, acc) \Rightarrow acc * (x / y) )$

1. acc = 1
2. acc =  $1 * (1 / 1) = 1$
3. acc =  $1 * (4 / 2) = 2$
4. acc =  $2 * (9 / 3) = 6$
5. acc =  $6 * (16 / 4) = 24$
6. acc =  $24 * (25 / 5) = 120$
7. Return 120

## Example: Max Path Sum

Starting at the top of the triangle and moving down to adjacent numbers below: Find the path from the root to a leaf with the maximum sum.

```
    3  
   7 4  
  2 4 6  
 8 5 9 3
```

## Example: Max Path Sum

Starting at the top of the triangle and moving down to adjacent numbers below: Find the path from the root to a leaf with the maximum sum.

```
  3  
 7 4  
2 4 6  
8 5 9 3
```

## Example: Max Path Sum

Starting at the top of the triangle and moving down to adjacent numbers below: Find the path from the root to a leaf with the maximum sum.

```
  3  
 7 4  
2 4 6  
8 5 9 3
```

## Example: Max Path Sum

Starting at the top of the triangle and moving down to adjacent numbers below: Find the path from the root to a leaf with the maximum sum.

```
    3  
   7 4  
  2 4 6  
8 5 9 3
```

## Example: Max Path Sum

We are going to fold the triangle

Starting at the top of the triangle and moving down to adjacent numbers below: Find the path from the root to a leaf with the maximum sum.

```
    3  
   7 4  
  2 4 6  
8 5 9 3
```

## Example: Max Path Sum

Starting at the top of the triangle and moving down to adjacent numbers below: Find the path from the root to a leaf with the maximum sum.

```
    3  
   7 4  
  2 4 6  
8 5 9 3
```

In the above example the max path sum is  $3 + 7 + 4 + 9 = 23$

## Example: Max Path Sum

Starting at the top of the triangle and moving down to adjacent numbers below: Find the path from the root to a leaf with the maximum sum.

```
    3  
    7 4  
  2 4 6  
8 5 9 3
```

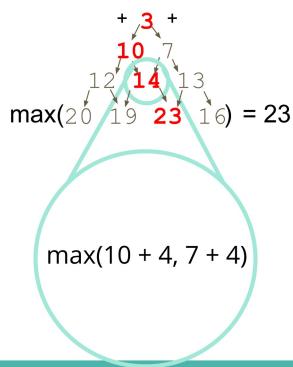
In the above example the max path sum in  $3 + 7 + 4 + 9 = 23$

Brute Force: Find all paths and get the max.

Source: <https://projecteuler.net/problem=18>

## Example: Max Path Sum

Our Algorithm



## Code Using Combinators

Triangle = list of lists. Can we use foldleft?

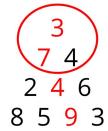
```
max( foldleft( triangle, [], lam(xs, acc) => myfold(xs, acc) ) )
```

```
    3  
    7 4  
  2 4 6  
8 5 9 3
```

## Code Using Combinators

Triangle = list of lists. Can we use foldleft?

```
max( foldleft( triangle, [], lam(xs, acc) => myfold(xs, acc) ) )
```



## Code Using Combinators

Triangle = list of lists. Can we use foldleft?

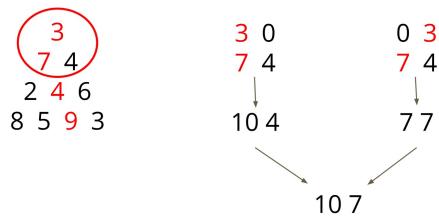
```
max( foldleft( triangle, [], lam(xs, acc) => myfold(xs, acc) ) )
```



## Code Using Combinators

Triangle = list of lists. Can we use foldleft?

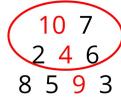
```
max( foldleft( triangle, [], lam(xs, acc) => myfold(xs, acc) ) )
```



## Code Using Combinators

Triangle = list of lists. Can we use foldleft?

```
max( foldleft( triangle, [], lam(xs, acc) => myfold(xs, acc) ) )
```

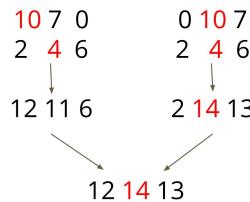
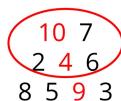


Always take the maximum of the maximum to get to that particular leaf node of intermediary in the triangle

## Code Using Combinators

Triangle = list of lists. Can we use foldleft?

```
max( foldleft( triangle, [], lam(xs, acc) => myfold(xs, acc) ) )
```



Add 0 to right of prev line to get left sum, and 0 to left of prev lin to get right sum

(To handle lists of different lengths...)

## Code Using Combinators

Triangle = list of lists. Can we use foldleft?

```
max( foldleft( triangle, [], lam(xs, acc) => myfold(xs, acc) ) )
```

**myfold(xs, acc) :**

```
option1 = map2( [0, acc...], xs, lam(x, y) => x + y )
option2 = map2( [acc..., 0], xs, lam(x, y) => x + y )
return map2( option1, option2, lam(x, y) => if x > y then x else y )
```

## Python Examples

```
myfolder
└── main.py
    └── data.txt
```

```
main.py
nums = []
with open("../data.txt", "w+") as f:
    lines = f.readlines()
    for line in lines:
        nums += [int(x) for x in line.split(",")]
print(sum(nums))
```

```
data.txt
0,1,2,3,4
```

What is the output of this code?

- a. Prints "10"
- b. Raises an Error
- c. Prints "0"

## Python Examples

```
class Bank:
    def __init__(self, balance):
        self.balance = balance

    def is_overdrawn(self):
        return self.balance < 0

myBank = Bank(100)
if myBank.is_overdrawn():
    print("OVERDRAWN")
else:
    print("ALL GOOD")
```

What is the output of this code?

- a. Prints "OVERDRAWN"
- b. Prints "ALL GOOD"
- c. Raises an Error

## Python Examples

```
for i in range(4):
    print(i)
    i = 10

some_string = "what"
some_dict = {}
for i, some_dict[i] in enumerate(some_string):
    i = 10

print(some_dict)
```

What is the output of this code?

- a. Prints "0"
- b. Prints "0 1 2 3"
- c. Raises an Error

What is the output of this code?

- a. Prints "{}"
- b. Prints "{0: 'w', 1: 'h', 2: 'a', 3: 't'}"
- c. Raises an Error

<https://book.pythontips.com/en/latest/enumerate.html>

## Python Examples

```
row = [""] * 3 # row [", ", ""]  
board = [row] * 3  
print(board) # [[", ", ""], [", ", ""], [", ", ""]]  
print(board[0]) # [", ", ""]  
print(board[0][0]) #  
board[0][0] = "X"  
print(board)
```

What is the output of this code?

- a. Prints "[["X", ""], [", ", ""], [", ", ""]]"
- b. Prints "[", ""], [", ", ""], [", ", ""]"
- c. Prints "[["X", ""], [X, ""], [X, ""]]"

## Python Examples

```
funcs = []  
results = []  
for x in range(3):  
    def some_func():  
        return x  
    funcs.append(some_func)  
    results.append(some_func()) # note the function call here  
  
funcs_results = [func() for func in funcs]  
print(results) # [0,1,2]  
print(funcs_results)
```

What is the output of this code?

- a. Prints "[0,1,2]"
- b. Prints "[2,2,2]"
- c. Prints "[]"
- d. Raises an Error

## Python Examples

You can find more such Python examples here [1]

[1] <https://github.com/satwikkansal/wtfpython>

He'll upload code

You can still upload your notes

Go through the Python examples and see what you expect it to produce, and see what is different, and try to understand why

wtfpython is a great repo for understanding these kind of things

Try to implement this workflow when you can

Tendency to make things perfect from the getgo, but the priority is to get something done

```
git checkout -b triangle
git log
mkdir 03-triangle
cd

# Missed some
# Can use VS code, pycharm, spyder, whatever
# Create a triangle.txt. file that is the triangle from the exercise
# Create a triangle.py file that contains code
```

```
1 # Basically wrote the highest level function first,
2 # and worked backwards
3 # To convince you that it will work, now read from
4 # bottom up
5
6 def read_triangle(path_to_file):
7     # Input in a string representing the file location
8     # returns a list of list of ints
9     f = open(path_to_file, "r")
10    lines = f.readlines()
11    triangle = []
12    for line in lines:
13        # Split and make into int with list
14        # comprehension
15        triangle.append([int(x) for x in line.split(
16            sep=' ')])
17    # Close file since we didn't use with open
18    f.close()
19    return triangle
20
21
22
23
24
25
26
27 def add_levels(x, y):
28     if len(x) != len(y):
29         raise(ValueError("lengths not equal"))
30     res = []
31     for i in range(len(x)):
32         res.append(x[i] + y[i])
33     return res
34
35
36
```

```
File - C:\Users\Wolfs\PycharmProjects\cs506_20210915\triangle.py
37 def fold_level(cur, nxt):
38     #
39     # Need to know what is the left and the right
40     # Add 0 to right of prev level to get the left
41     # Add 0 to left of prev level to get the right
42     # Try to have descriptive names that are not
        misleading
43     # Add 0 to left side:
44     right_options = add_levels([0] + cur, nxt)
45     # Add to right side:
46     left_options = add_levels(cur + [0], nxt)
47     max_options = elementwise_max(right_options,
        left_options)
48     return max_options
49
50
51 # Write this next function with a top down approach
52 def fold(triangle):
53     # Returns list of max path sums
54     # note: does not work when triangle length is 1 (
        can enhance later for edge case)
55     current_level = triangle[0]
56     for level in triangle[1:]:
57         current_level = fold_level(current_level,
        level)
58     return current_level
59
60
61 # or "./triangle.txt"
62 # print(read_triangle("triangle.txt"))
63 print(max(fold(read_triangle("triangle.txt"))))
64 # Result: 23, it works
65 # Big triangle: 1074
66 print(max(fold(read_triangle("triangle2.txt"))))
```

Page 2 of 2

```
23
1074
Process finished with exit code 0
```



## Introduction

Boston University CS 506 - Lance Galletti

Now starting the actual data science

## Data Representation

How we represent data is linked to what information we are able to retrieve from it.

## Data Representation - Records

m-dimensional points / vectors

Example: (name, age, balance) -> ("John", 20, 100)

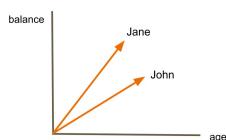
Typically stored as a tuple representing a vector

## Data Representation - Records

m-dimensional points / vectors

Example: (name, age, balance) -> ("John", 20, 100)

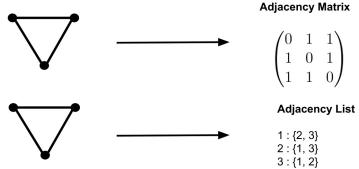
Now you can do anything in the realm of linear algebra on these two vectors



## Data Representation - Graphs

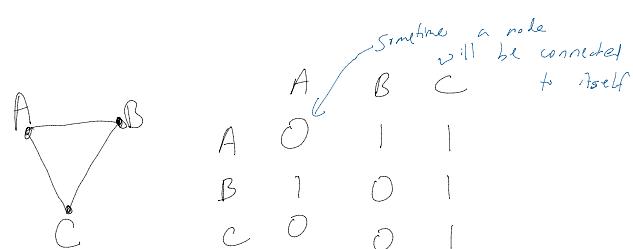
Nodes connected by edges

Example:

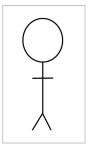


- We'll see more of this at the end of the semester
- Downward facing triangle
  - o One way to represent this is through an adjacency matrix

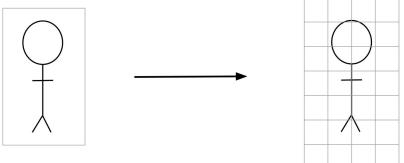
Adjacency list may be more space efficient



## Data Representation - Images

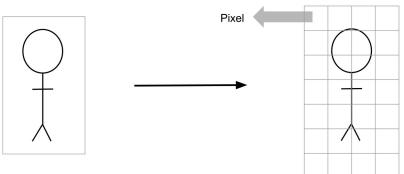


## Data Representation - Images



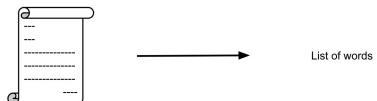
## Data Representation - Images

- Matrix or list of pixels



## Data Representation - Text

- text is usually a list of words



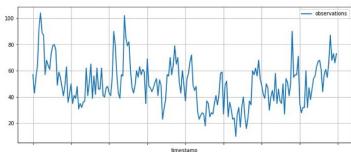
## Data Representation - Strings

- DNA sequence

DNA seq (A T G C C G T A ...) → list of characters

## Data Representation - Time Series

List of data at specific intervals of time



- Can be any of the data we saw before
  - o Video, images
  - o or simple values like temp as a function of time
    - Metal plate example Lec0 was a time series

## Types of Learning

- Supervised Learning
- Unsupervised Learning

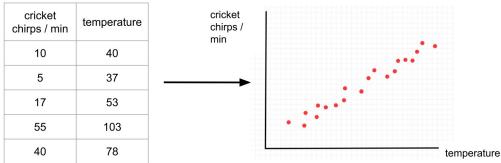
- We will tackle two different types of learning
  - o Unsupervised first and is most of the course
  - o Supervised learning
    - What people think of for machine learning?

## Supervised Learning

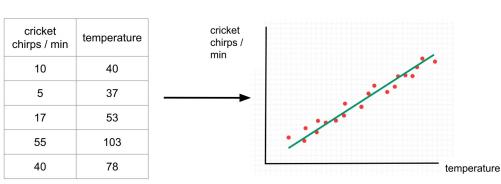
cricket chirps / min	temperature
10	40
5	37
17	53
55	103
40	78

- Total nonsense data

## Supervised Learning

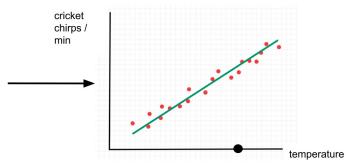


- Want to create some sort of model (green line)
- For a temperature, what is your expected number of cricket chirps



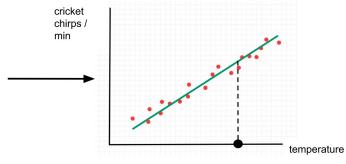
## Supervised Learning

cricket chirps / min	temperature
10	40
5	37
17	53
55	103
40	78



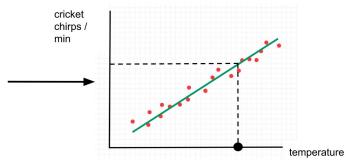
## Supervised Learning

cricket chirps / min	temperature
10	40
5	37
17	53
55	103
40	78



## Supervised Learning

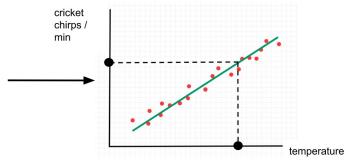
cricket chirps / min	temperature
10	40
5	37
17	53
55	103
40	78



## Supervised Learning

- Will see regression again in this course

cricket chirps / min	temperature
10	40
5	37
17	53
55	103
40	78



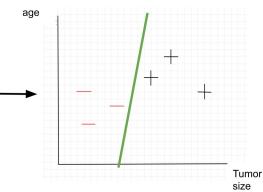
This type of supervised learning is referred to as regression

## Supervised Learning

age	tumor size	malignant
20	12	0
22	15	1
47	20	1
59	2	1

## Supervised Learning

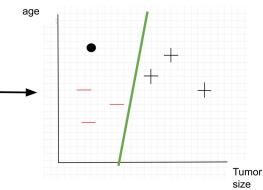
age	tumor size	malignant
20	12	0
22	15	1
47	20	1
59	2	1



- Again simplified unreal dataset

## Supervised Learning

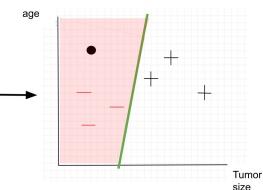
age	tumor size	malignant
20	12	0
22	15	1
47	20	1
59	2	1



- If we have a new point, do we have a malignant or a benign tumor

## Supervised Learning

age	tumor size	malignant
20	12	0
22	15	1
47	20	1
59	2	1

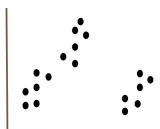


This type of supervised learning is referred to as classification

- Want to be able to classify new datapoints

## Unsupervised Learning

Goal: Find interesting structure in the data



- Unsupervised
  - o Less interested in the specific label, but is interested in the inherent structure of the data
    - Like clusters

## Unsupervised Learning

Goal: Find interesting structure in the data



- Clustering
  - o Assign each point to a given cluster
  - o Will cover this a lot

This type of unsupervised learning is referred to as clustering

## Unsupervised Learning

Dataset: Collection of Articles

Question: Are these articles covering the same topics?

- Will collection of articles or books
  - o High level concepts: Space, ...
  - o Similarity of concepts shows up in the structure of the article
  - o No one told us what topic the article was, we are grouping articles into topics

## Distance & Similarity

Boston University CS 506 - Lance Galletti

- Starting with some unsupervised tools

### Data

$$\begin{matrix} n \text{ data points} \\ \left\{ \begin{pmatrix} x_{11} & \dots & x_{1j} & \dots & x_{1m} \\ \vdots & \ddots & \vdots & & \vdots \\ x_{i1} & \dots & x_{ij} & \dots & x_{im} \\ \vdots & & \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nj} & \dots & x_{nm} \end{pmatrix} \right. \\ m \text{ features} \end{matrix}$$

- $m \times n$  attributes?

### Feature Space

From our data we can generate a **feature space** of all possible values for the set of features in our data.

name	age	balance
Jane	25	150
John	30	100

- Name, age, balance
  - o Name shouldn't matter in a data science model
    - Maybe it would in some models
  - o Focus on the numerical features

### Feature Space

From our data we can generate a **feature space** of all possible values for the set of features in our data.

name	age	balance
Jane	25	150
John	30	100

- What are all the possible values age can have
- What are all the possible values balance can have

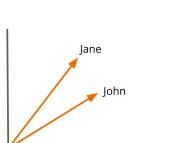
### Feature Space

From our data we can generate a **feature space** of all possible values for the set of features in our data.

name	age	balance
Jane	25	150
John	30	100

- Plot our vectors in this feature space

Our feature space is the Euclidean plane



## Distance

In order to uncover interesting structure from our data, we need a way to **compare** data points.

A **dissimilarity function** is a function that takes two objects (data points) and returns a **large value** if these objects are **dissimilar**.

A special type of dissimilarity function is a **distance** function

- Key point for unsupervised learning:
  - o Really fundamental thing in the dataset
  - o Are these points similar, dissimilar, etc
  - o Want to be able to compare the datapoints in some way
  - o Don't need a distance function, but we could make a dissimilarity function
  - o Distance is a more mathematical function

## Distance

$d$  is a distance function if and only if:

- $d(i, j) = 0$  if and only if  $i = j$
- $d(i, j) = d(j, i)$
- $d(i, j) \leq d(i, k) + d(k, j)$

We don't **need** a distance function to compare data points, but why would we prefer using a distance function?

- What is a distance function?
  - o e.g. Square root of the sum of squares
- For a distance function to be defined, it needs 3 properties
  - o if and only if (iff) statement required in both directions
- Why would we prefer a distance function?
  - o Always positive
  - o Easy to understand and graph
  - o Is intuitive

## Minkowski Distance

For  $x, y$  points in  $d$ -dimensional real space

I.e.  $x = [x_1, \dots, x_d]$  and  $y = [y_1, \dots, y_d]$

$$p \geq 1 \quad L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

When  $p = 2 \rightarrow$  Euclidean Distance

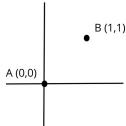
When  $p = 1 \rightarrow$  Manhattan Distance

- WARNING: Math (lol)
- Hopefully you've taken linear algebra...
- $d$ -dimensional points have  $d$  features/components
- Minkowski function
  - o Should have seen this before, review if needed
  - o Take the  $p$ -th root

$$\begin{aligned} & \sum_{i=1}^3 |x_i - y_i| = 1 + 2 + 3 \\ & \sum_{i=1}^d |x_i - y_i|^p = |x_1 - y_1|^p + |x_2 - y_2|^p + \dots + |x_d - y_d|^p \\ & \hookrightarrow \text{will see this a lot} \\ & \sqrt{|x_1 - y_1|^2 + |x_2 - y_2|^2 + |x_3 - y_3|^2} \\ & \hookrightarrow \text{Euclidean Distance} \end{aligned}$$

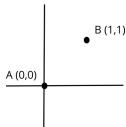
## Example

$d = 2$



## Example

$d = 2$



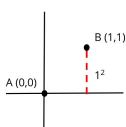
$p = 2$

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

- Standard Euclidean distance you should already know

## Example

$d = 2$

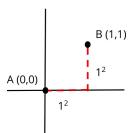


$p = 2$

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

## Example

$d = 2$

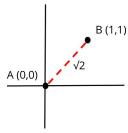


$p = 2$

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

## Example

$d = 2$

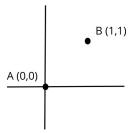


$p = 2$

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

## Example

$d = 2$



$p = 1$

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

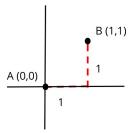
- if  $p=1$

- This is Manhattan distance

- Can only move along x and y and not directly

## Example

$d = 2$

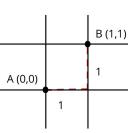


$p = 1$

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

## Example

$d = 2$

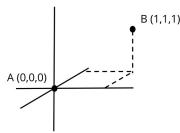


$p = 1$

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

## Example

$d = 3$

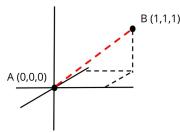


$p = 2$

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

## Example

$d = 3$



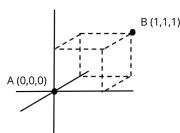
$p = 2$

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

- Euclidean: can go directly

## Example

$d = 3$



$p = 1$

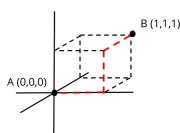
$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

- Manhattan distance

- o Move along the 3D grid edges
- o Multiple ways to get to B, but all length 3

## Example

$d = 3$



$p = 1$

$$L_p(x, y) = \left( \sum_{i=1}^d |x_i - y_i|^p \right)^{\frac{1}{p}}$$

## Minkowski Distance

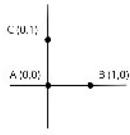
Is  $L_p$  a distance function when  $0 < p < 1$ ?

- At beginning we said  $p \geq 1$

- o Would it still be a distance function?
  - Can you find a counter example?

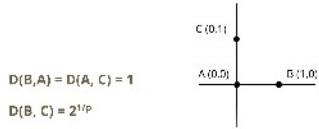
## Minkowski Distance

Is  $L_p$  a distance function when  $0 < p < 1$ ?



## Minkowski Distance

Is  $L_p$  a distance function when  $0 < p < 1$ ?

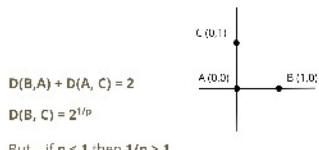


$$D(p,c) = (1+1)^{1/p}$$

if  $p < 1$   
then  $1/p > 1$

## Minkowski Distance

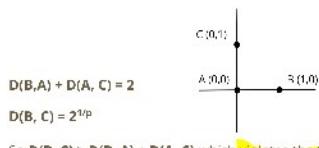
Is  $L_p$  a distance function when  $0 < p < 1$ ?



But... if  $p < 1$  then  $1/p > 1$

## Minkowski Distance

Is  $L_p$  a distance function when  $0 < p < 1$ ?



So  $D(B,C) > D(B,A) + D(A,C)$  which violates the triangle inequality.

## Cosine Similarity

A **similarity** function is a function that takes two objects (data points) and returns a **large value** if these objects are similar.

$$s(x, y) = \cos(\theta)$$

where  $\theta$  is the angle between  $x$  and  $y$

Monday: spark pitches are remote, don't come to class  
Have office hours after this

## Cosine Similarity

To get a corresponding **dissimilarity** function, we can usually try

$$d(x, y) = 1 / s(x, y)$$

or

$$d(x, y) = k - s(x, y) \text{ for some } k$$

Here, we can use

$$d(x, y) = 1 - s(x, y)$$

## Cosine Similarity

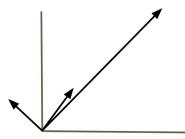
When should you use cosine (dis)similarity over euclidean distance?

When **direction** matters more than **magnitude**

## Cosine Similarity

When should you use cosine (dis)similarity over euclidean distance?

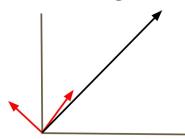
When **direction** matters more than **magnitude**



## Cosine Similarity

When should you use cosine (dis)similarity over euclidean distance?

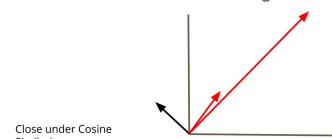
When **direction** matters more than **magnitude**



## Cosine Similarity

When should you use cosine (dis)similarity over euclidean distance?

When **direction** matters more than **magnitude**



## Jaccard Similarity

How similar are the following documents?

	$w_1$	$w_2$	...	$w_d$
x	1	0	...	1
y	1	1	...	0

## Jaccard Similarity

One way is to use the Manhattan distance which will return the size of the set difference

	w <sub>1</sub>	w <sub>2</sub>	...	w <sub>d</sub>
x	1	0	...	1
y	1	1	...	0

$$L_1(x, y) = \sum_{i=1}^d |x_i - y_i|$$

## Jaccard Similarity

One way is to use the Manhattan distance which will return the size of the set difference

	w <sub>1</sub>	w <sub>2</sub>	...	w <sub>d</sub>
x	1	0	...	1
y	1	1	...	0

$$L_1(x, y) = \sum_{i=1}^d |x_i - y_i|$$

Will only be 1 when  $x_i \neq y_i$

## Jaccard Similarity

But how can we distinguish between these two cases?

	w <sub>1</sub>	w <sub>2</sub>	...	w <sub>d-1</sub>	w <sub>d</sub>
x	1	1	1	0	1
y	1	1	1	1	0

Only differ on the last two words

	w <sub>1</sub>	w <sub>2</sub>
x	0	1
y	1	0

Completely different

## Jaccard Similarity

But how can we distinguish between these two cases?

	w <sub>1</sub>	w <sub>2</sub>	...	w <sub>d-1</sub>	w <sub>d</sub>
x	1	1	1	0	1
y	1	1	1	1	0

Only differ on the last two words

	w <sub>1</sub>	w <sub>2</sub>
x	0	1
y	1	0

Completely different

Both have Manhattan distance of 2

## Jaccard Similarity

We need to account for the size of the intersection!

$$JSim(x, y) = \frac{|x \cap y|}{|x \cup y|}$$

$$JDist(x, y) = 1 - \frac{|x \cap y|}{|x \cup y|}$$

**Implement these distance functions in the CS506  
python package**