

Informatik 2 Cheatsheet

Julian Lotzer – jlotzer@student.ethz.ch

Daniel Steinhauser –

dsteinhauser@student.ethz.ch

Modified by: Christian Leser - cleser@ethz.ch

Version: November 2, 2023

Template by Micha Bosshart

1 General Python

1.1 Reference Semantics and Aliasing

Everything is a pointer: l1 and l2 point to the same adress

```
1 l1 = [1, 3, 'hi', -4] #l1 --> [1, 3, 'hi', -4]
2 l2 = l1 #l2 --> [1, 3, 'hi', -4]
```

If a copy is needed, use:

```
1 import copy
2 l2 = copy.copy(l2) #shallow copy (1D-Array)
3 l2 = copy.deepcopy(l2) #deep copy (Multi
  Dimensional Arrays)
```

1.2 Data Types

Python dynamically types variables, which means that the variable type can change during the program's execution

```
1 s = 5
2 print(type(s)) #output: <class 'int'>
3 s = False # <class 'bool'>
4 s = "Hello World" # <class 'str'>
```

To convert the data type:

```
1 s = 5.9 #type(s) = <class 'float'>
2 x = int(s) #type(x) = <class 'int'>
3 y = str(s) #type(y) = <class 'str'>
```

1.2.1 Type Hints

help make code more legible

```
1 def add_integers(l1: list[int]) -> int:
2     ...
3 # after colon: expected input type
4 # after arrow: expected return type
```

1.3 Input and Output

Output

```
1 print("Hello World") #output: Hello World
2 print("Hello", "World") #output: Hello World
3 print("Hello", "World", sep = "--", end = "!")
4 #output: Hello--World!
```

Input

```
1 name = input("Enter name: ") #input returns a
  string
2 print("Hello", name)
3 number = int(input("Enter your number: "))
4 print("Number:", number)
```

1.4 Control Flows (if/else, while, for)

if, else

```
1 x = int(input("Enter a number: "))
2 if x < 5 and x >= 0:
3     print("too small") #if x between 0 and 5
4 elif x < 69 or x == 420:
5     print("nice") #if x is equal to 69 or 420
6 else:
7     print("big number") #x is positive and the ifs/
  elifs conditions don't hold
```

while

```
1 x = 0
2 while x <= 3:
3     x += 1
```

for with range

```
1 l=[3,5,25]
2 for i in l:
3     print(i, end = " ") #output: 3 5 25
```

for with lists

```
1 for i in reversed(range(0,4,1)):
2     print(i, end = " ") #output: 3 2 1 0
```

1.5 Functions

Functions do not have to be declared in a specific order

1.5.1 Function Declaration

```
1 def function(arg1, arg2):
2     ...
3     return value
```

1.5.2 Default arguments

When calling a function with default arguments, it is not necessary to call the function with arguments.

```
1 def specialprint(data="hello world"):
2     print(data)
3 specialprint() #output: hello world
```

1.5.3 Global and local variables

Avoid this kind of code!

global variables can be used within a function if declared before the function call. Also, local variables can be made global:

```
1 def double(x):
2     global result #result is defined globally
3     result = x*2
```

2 Algorithms

2.1 Invariants

The invariant is a condition in an algorithm involving a loop.

Phase	Description	Example (list 'a' is sorted until index 'i')
1. Initialization	the condition is met before the loop	$i = 0, a[:i] = a[0]$
2. Continuation	the condition holds at each iteration	$i = j, a[:j]$
3. Termination	the condition holds at the end of the loop	$i = n, a[:n] = a[n]$

2.2 Divide and Conquer Algorithms

1. Divide Problem into easier to solve subproblems
2. Solve subproblems
3. Reassemble solutions from subproblems

Examples: Sorting Algorithms for lists

- Mergesort
- Quick Sort

3 Python Containers

3.1 Operations on Containers

Number of Elements:

```
1 len(c)
```

Contains element x?

```
1 x in c
```

Iterate over all elements:

```
1 for x in c:
2     print(x)
```

3.2 Sequences (ordered containers)

3.2.1 Sequence Operations (Überarbeiten)

Subscript-Operator l[i]:

```
1 l = [1, 3, 'hi', -4]
2 print(l[2]) #output: hi
```

Enumeration:

```
1 enumerate(iterable, start)
2 #iterable = iterable container (sequence)
3 #start (optional) = (optional) enumerate starts
  counting #at this number, starts at 0 when
  omitting start->
4 #enumerate(iterable)
5 #the enumerate(iterable, start) function returns a
  #tuple: (index, object).
6
7 for index, value in enumerate(1):
8     print(index, value)
9 #output:
10 # 0 1
11 # 1 3
12 # 2 Hi
13 # 3 -4
```

Combine sequences s1 and s2 (zip):

```
1 z = zip(s1,s2)
2 #example:
3 #s1 -> "Lea" "Tim" "Mortis"
4 #s2 -> 22 19 69
5 #z -> ("Lea",22) ("Tim",19) ("Mortis",69)
```

Output with a for loop:

```
1 for name, age in z:
2     print(name,"->",age)
3 # output:
4 # Lea -> 22
5 # Tim -> 19
6 # Mortis -> 69
```

Slicing (partial sequence) of a sequence s:

```
1 partseq = s[start:stop:step]
2 partseq = s[start:stop] #step = 1
3 partseq = s[:stop:step] #start = 0
4 partseq = s[start::step] #stop = len(s)
```

3.2.2 List (mutable)

- ordered data
- Fast Access via index
- Slow for updates

Initialise a list with []

```
1 l = [1, 3, "hi", -4]
2 M = [[-1 for i in range(n)] for j in range(m)]
3 # 2D list or m x n-Matrix filled with -1
```

Common List Operations

```
1 l[i] = value #Access and change item
2 l.pop(i) #Delete element at index i
3 l.remove(v) #Delete element with value v
4 l.insert(i, v) #Insert element with value v at i
5 l.append(value) #Add item at the end
6 del l[i] #Remove item at location i
7 l.reverse()
8 l=[v]*k #Create a list of k elements with value v
```

List Comprehension

Apply a function $f(x)$ to all items in list l:

```
1 l2 = [f(x) for x in l] #z.g. 2*x for f(x)
```

Apply a function $f(x)$ to a range:

```
1 r2 = [f(x) for x in range(1,6)]
```

Apply a function $f(x)$ only to items in list l that satisfy $g(x)$ (filter):

```
1 l3 = [f(x) for x in l if g(x)]
```

Example: Read a sequence of numbers:

```
1 l = [int(x) for x in input("Input: ").split()]
```

3.2.3 Sorting Algorithms for Lists

Selection Sort

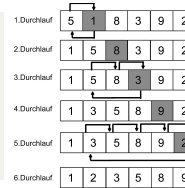
```
5 [5] [6] [2] [8] [4] [1] (i=0) sorted list
6 [1] [6] [2] [8] [4] [5] (i=1) unsorted list
7 [1] [2] [6] [8] [4] [5] (i=2) smallest value
```

```
1 #Input: Array a = (a[0],...,a[n])
2 #Output: Sorted Array a
3 def sort(a):
4     n = len(a)
5     for i in range(n):
```

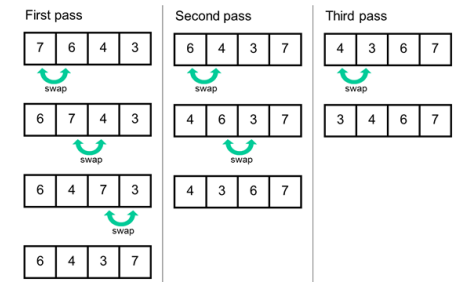
```
6     mini = i
7     for j in range(i+1, n):
8         if a[j] < a[mini]:
9             mini = j
10    a[mini],a[i] = a[i],a[mini]
11    return a
```

Insertion Sort

```
1 def insertionSort(arr):
2     for i in range(1, n):
3         key = arr[i]
4         j = i-1
5         while j >=0 and key <
          arr[j]:
6             arr[j+1] = arr[j]
7             j -= 1
8             arr[j+1] = key
```



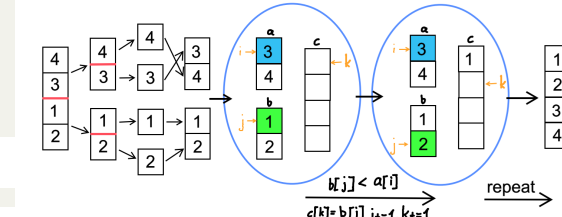
Bubble Sort



```
1 def bubbleSort(arr):
2     n = len(arr)
3     for i in range(n-1):
4         for j in range(0, n-i-1):
5             if arr[j] > arr[j + 1]:
6                 arr[j], arr[j + 1] = arr[j + 1],
          arr[j]
```

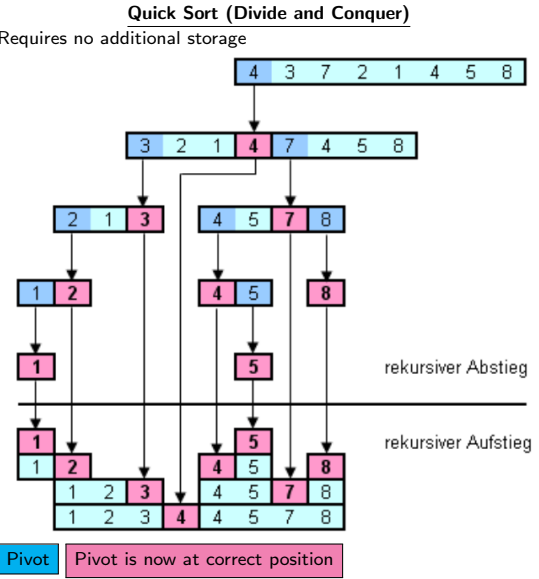
Mergesort (Divide and Conquer)

Requires $\Theta(n)$ additional storage



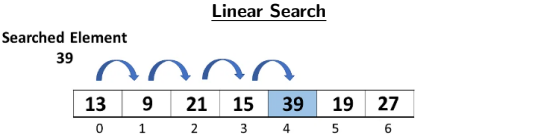
```
1 def merge(a1, a2):
2     b, i, j = [], 0, 0
3     while i < len(a1) and j < len(a2):
4         if a1[i] < a2[j]:
5             b.append(a1[i])
6             i += 1
7         else:
8             b.append(a2[j])
9             j += 1
10    b += a1[i:]
11    b += a2[j:]
12    return b
```

```
14 def merge_sort(a):
15     if len(a) <= 1:
16         return a
17     else:
18         sorted_a1 = merge_sort(a[:len(a) // 2])
19         sorted_a2 = merge_sort(a[len(a) // 2:])
20         return merge(sorted_a1, sorted_a2)
```

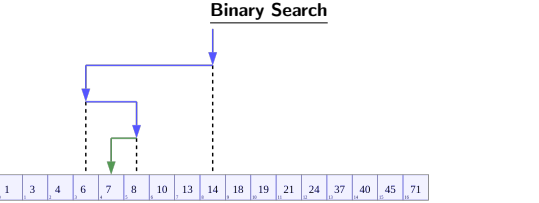


```
1 def partition(a, l, r):
2     p = a[r]
3     j = l
4     for i in range(l, r):
5         if a[i] < p:
6             a[i], a[j] = a[j], a[i]
7             j += 1
8     a[j], a[r] = a[r], a[j]
9     return j
10
11 def quicksort(a, l, r):
12     if l < r:
13         k = partition(a, l, r)
14         quicksort(a, l, k - 1)
15         quicksort(a, k + 1, r)
```

3.2.4 Search Algorithms for Lists



```
1 #Input: array a, key b
2 #Output: index k with a[k] = b or "not found"
3 def LinearSearch(a, b):
4     for i, x in enumerate(a):
5         if x == b:
6             return i
7     return "not found"
```



```
1 #Input: sorted array a, key b
2 #Output : index k with a [ k ] = b or " not found "
3 def bin_search(a, l, r, b):
4     if r < l:
5         return None
6     else:
7         m = (l + r) // 2
8         if a[m] == b:
9             return m
```

```
11 elif b < a[m]:
12     return bin_search(a, l, m-1, b)
13 else: # a[m] > b
14     return bin_search(a, m+1, r, b)
```

3.2.5 Tuple (immutable)

Initialise a tuple with ()

```
1 t = ("a", 0, -6, 3.3)
```

3.2.6 Range (immutable)

Initialise a range

```
1 #range(start, stop, step)
2 #default: start = 0, step = 1, end not included
3 r = range(0, 8, 2) #r -> 0 2 4 6,
```

3.2.7 String (immutable)

Initialise a string with ""

```
1 s = "hello" #s -> 'h' 'e' 'l' 'l' 'o'
2 #you can use both " or ' for strings
```

Common String Operations

```
1 s1 = "hello"
2 s2 = " world"
3 s1[1] #Access element at position 2 -> e
4 s3 = s1 + s2 ##Add strings (concatenating) ->s3 = "hello world"
5 s2=s2.strip() #Remove whitespace at beginning and end -> "world"
6 s = list(s1) #Convert string into list of chars
Example: check if s is a string with content:
1 type(s) == str and len(s.strip()) #False if empty
Example: convert a string to a list of words:
1 s = "Hello World"
2 w = s.split() #-> w = ['Hello', 'World']
3 s.split(seperator, maxsplit)
4 #seperator and maxsplit are optional
5 #s.split() -> split at all whitespaces
6 #seperator = ", " -> split at every ", "
7 #maxsplit = 10 -> split only at first 10 separators
```

3.3 Collections (unordered containers)

3.3.1 Set (non-associative)

Initialise Set

```
1 s = {1, 29, 12}
```

Common Set Operations

```
1 s.add(69) #Add item
2 s.remove(29) #Remove item
3 12 in S #Search for item, returns bool
```

3.3.2 Dictionary (associative)

See 3.2.1 Sequence operations for 'enumerate' and 'zip'

Initialise a Dictionary with {}

A dictionary consists of **tuples (key, value)** as items. For that reason, one can think of it as a list of tuples (Which it is not in reality)

```
1 d = {"Lea":22, "Tim":19, "Mortis":69} #key:value
2 #Merge two lists into one dictionary
3 cities = ["Zurich", "Basel", "Bern"] #list 1
4 code = [8000, 4000, 3000] #list 2
5 d2 = dict(zip(cities,code)) #dictionary D2
```

Common Dictionary Operations

```
1 d["Lea"] = 23 #Change item
2 d["Peter"] = 24 #Add item
3 #no key can exist more than one time
4 del d["Mortis"] #delete item
5 "Tim" in d #Search for key, returns bool
6 d["Tim"] #access value at a key, output: 19
```

Iterate over a Dictionary

```
1 for key in d.keys(): # Iterate over the keys
2     print(key) #Lea Tim Mortis
3 for item in d.items(): # Iterate over the entries
4     print(item) #("Lea",22) ("Tim",19) ("Mortis", 69)
5 for key, value in d.items(): # Iterate over the keys and values
```

```
6 print(" " + value) #Lea 22 Tim 19 Mortis 69
7 for value in d.values(): # Iterate over the values
8     print(value) #22 19 69
```

Dictionary Comprehension

Transform a set into a dictionary by applying $f(x)$ and $g(x)$ on every element in the set:

```
1 d3 = {f(x):g(x) for x in s} #s being a set
2 d4 = {f(x):g(y) for x, y in z.items()} #z being a dict
```

Transform a set into a dictionary, only if the element satisfies $h(x)$:

```
1 d5 = {f(x):g(x) for x in s if h(x)}
```

Example: Multiply the value of every odd key in a dictionary by 2:

```
1 d6 = {k:2*v for k, v in d.items() if k % 2 == 1}
```

4 Python Classes

- entity with a name that contains data and functionality
- bundling of data that belongs together contentwise
- definition of a new type
- every object of this type stores data and offers functionality

```
1 class Class_name_1:
2     def __init__(self, att1, att2): #constructor
3         self.attribute1 = att1
4         self.attribute2 = att2
5
6 #the constructor is called to create object1
7 object1 = Class_name_1(1, 1)
8
9 class Class_name_2: #type definition
10     attribute3 = None
11     attribute4 = Class_name_1(0, 0)
12
13     def method1(self, parameter1):
14         ...
15
16     def __eq__(self, other): #overloads ==
17         return (self.attribute3 == other.attribute3
18                 and self.attribute4 == other.attribute4)
19
20 object2 = Class_name_2() #generate object
21 object3 = Class_name_2()
22 object2.attribute3 = 64 #write attribute
23 print(object_name.attribute1) #read attribute
24 object2 == object3 #True, uses overloaded ==
25 object2.method1(...) #call function
```

4.1 Private attribute

```
1 class ClassName:
2     __private = 0 #private attribute
3
4     def set_privateat(self, a):
5         self.__privateat = a
6
7     def get_private(self):
8         return self.__private
```

4.2 Inheritance

- Child inherits all attributes and methods from Parent
- Child can define additional attributes and methods

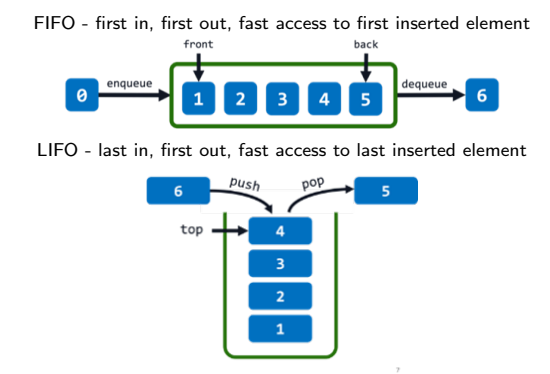
```
1 class Parent:
2     def __init__(self, x1):
3         self.x = x1
4     def double_x(self):
5         self.x *= 2
6
7 class Child(Parent):
8     def __init__(self, x1, y1):
9         Parent.__init__(self, x1)
10        self.y = y1
11    def double_y(self):
12        self.y *= 2
13
14 c = Child(2, 2)
15 c.double_x()
16 print(c.x, c.y) #Output: 4 2
```

4.3 Magical Methods

Operation	Meaning	Magical Method
<	Less than	..lt..
<=	Less than or equal	..le..
>	Greater than	..gt..
>=	Greater than or equal	..ge..
==	Equal to	..eq..
!=	Not equal to	..ne..
+, +=	Addition	..add.., ..iadd..
-	Subtraction	..sub..
*	Multiplication	..mul..
/	Division	..truediv..
//	Integer division	..floordiv..
%	Modulo	..mod..
**	Exponentiation	..pow..
print()	overload print()	..str..
-	Negation	..neg..

5 Data Structures

5.1 FIFO / LIFO



5.2 Search Tree

- sorted data
- all operations fast
- well organized

5.2.1 Tree terminology

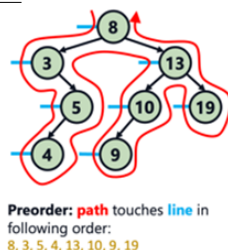
- Root: Highest node has children and no parent
- Order: Maximum number of child nodes (BST order = 2)
- Height: Maximum path length root to leaf
- Degenerated: Every Node has only one child, leading to a linked list structure
- Complete: Every level except for the lowest is filled, lowest level is filled from left to right
- Perfect: Every level is completely filled
- Full: the amount of children of every node equals either the order or 0
- ALV-Tree: balanced binary search tree ($-1 \leq \text{height}(\text{n.left}) - \text{height}(\text{n.right}) \leq 1$)

5.2.2 Binary Tree Traversal

Preorder

```
1 def preorder(node):
2     print(node.key)
3     preorder(node.left)
4     preorder(node.right)
```

If a list 'a' was the inorder traversal of a tree: the tree can always be recreated
⇒ Representation unique



Inorder

```
1 def inorder(node):
2     inorder(node.left)
3     print(node.key)
4     inorder(node.right)
```

Entries are printed in ascending order
If a list 'a' was the inorder traversal of a tree: the tree can be recreated only if the list is in ascending order
⇒ Representation not unique



Postorder

```
1 def postorder(node):
2     postorder(node.left)
3     postorder(node.right)
4     print(node.key)
```

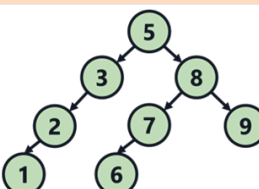
If a list 'a' was the inorder traversal of a tree: the tree can always be recreated
⇒ Representation unique



5.2.3 Binary Search Tree (BST)

A binary search tree is a binary tree which fulfils the following:

1. Every node v stores a key
2. Keys in the left subtree are smaller than v .key
3. Keys in the right subtree are larger than v .key



Implementation

```
1 class SearchNode:
2     def __init__(self, k, l=None, r=None):
3         self.key = k
4         self.left, self.right = l, r
5
6 class Tree:
7     def __init__(self):
8         self.root = None
9     def find(self, key):
10        return findNode(self.root, key)
11    def add(self, key):
12        self.root = addNode(self.root, key)
```

Height

```
1 def height(node):
2     if node == None:
3         return 0
4     else:
5         return 1 + max(height(node.left), height(node.right))
```

Search for Node

```
1 def findNode(root, key):
2     n = root
3     while n != None and n.key != key:
4         if key < n.key:
```

```
5         n = n.left
6     else:
7         n = n.right
8     return n
```

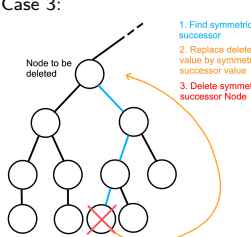
Insert Node

```
1 def addNode(root, key):
2     if root == None:
3         root = Node(key)
4     n = root
5     while n.key != key:
6         if key < n.key:
7             if n.left == None:
8                 n.left = Node(key)
9             n = n.left
10        else:
11            if n.right == None:
12                n.right = Node(key)
13            n = n.right
14    return root
```

Remove Node

Cases:

1. node has no children: set variable to None
2. node has one child: replace node with child
3. node has two children: replace node with symmetric successor – the next biggest element



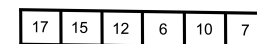
```
1 def findSuccessor(start_node):
2     parent = None
3     node = start_node.right
4     while node.left is not None:
5         parent = node
6         node = node.left
7     return node
8
9 def deleteNode(root, key):
10    if root is None:
11        return None
12    # Find the node to be deleted
13    if key < root.val:
14        root.left = deleteNode(root.left, key)
15    elif key > root.val:
16        root.right = deleteNode(root.right, key)
17    else:
18        # Case 1: Node has no children
19        if root.left is None and root.right is None:
20            root = None
21        # Case 2: Node has one child
22        elif root.left is None:
23            root = root.right
24        elif root.right is None:
25            root = root.left
26        # Case 3: Node has two children
27        else:
28            successor = findSuccessor(root)
29            root.val = successor.val
30            root.right = deleteNode(root.right, \
31                                 successor.val)
32    return root
```

5.2.4 (Max) Heaps

Heaps are only visualised as trees while they are stored as arrays
Useful for quick access to max (/ min) value

- complete binary tree (see 5.2.1 Tree terminology)

- Key of parent is always greater (smaller) than the one of its children



Array corresponding to the max heap

Implementation

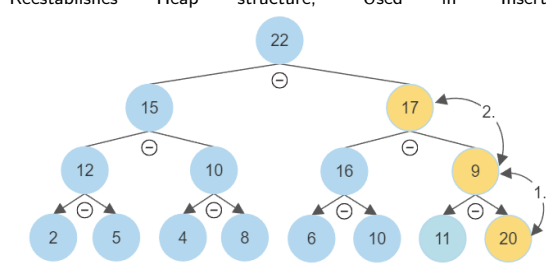
Heap[i = 1] = Root:

- Children of i : $2i+1$, $2i+2$
- Parent of i : $i-1//2$

Height

```
1 H(n) = log2(n + 1)
2 # a is a heap with n elements
3 def height(a):
4     return math.log(len(a) + 1, 2)
```

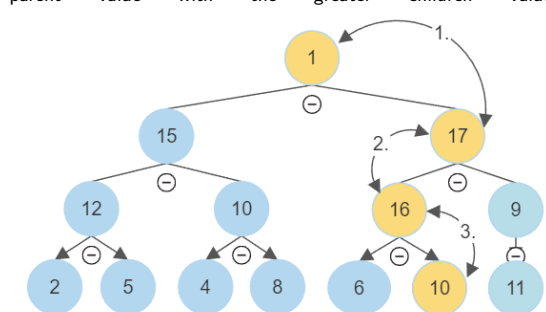
Reestablishes Heap structure, Used in 'Insert'



```
1 def SiftUp(a, m):
2     v = a[m]
3     c = m
4     p = c // 2
5     while c > 0 and v > a[p]:
6         a[c] = a[p]
7         c = p
8         p = c // 2
9     a[c] = v
```

Sift Down

Reestablishes Heap structure, Used in 'Remove' and 'Heapify'
If the parent value is smaller: exchange the parent value with the greater children value



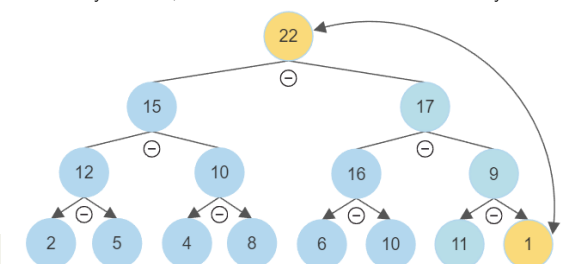
```
1 #i is the index of the element that needs to be
2 sifted down. a is the array. m is the end of
3 the array
4 def SiftDown(a, i, m):
5     while 2*i + 1 < m:
6         j = 2*i + 1
7         if j + 1 < m and a[j] < a[j + 1]:
8             j = j + 1
9         if a[i] >= a[j]:
10            break
11        a[i], a[j] = a[j], a[i]
12        i = j
```

Insert

```
1 def insert(heap, value):
2     heap.append(value)
3     SiftUp(heap, len(heap)-1)
```

Remove Max Value

Change the first and last entry in the array, then delete the last entry. Now, the max value does not exist any more.



Reestablish the Heap condition by applying Sift Down.

```
1 def removeMax(heap):
2     if len(heap) == 0:
3         return None
4     max_val = heap[0]
5     heap[0] = heap[-1]
6     heap.pop()
7     SiftDown(heap, 0, len(heap)-1)
8     return max_val
```

Heap creation / Heapify

Repeatedly apply Sift Down until the array is heapified.
Leaves fulfill the heap condition trivially → only "heapify" the first $n/2$ elements.

```
1 def heapify(a):
2     n = len(a)
3     for i in range(n//2 - 1, -1, -1):
4         SiftDown(a, i, n)
```

Sorting a heap

If "a" is a heap, one can efficiently sort the array:

```
1 #a is a heap
2 def SortHeap(a):
3     n = len(a)-1
4     while n > 0:
5         swap(a[0], a[n])
6         SiftDown(a, 0, n-1)
7         n = n - 1
```

5.3 Linked List

- ordered data
- fast updates in first elements



Implementation

```
1 class Node:
2     def __init__(self, value, Next = None):
3         self.value = value
4         self.Next = Next
5
6 class Linked_List:
7     def __init__(self, head):
8         self.head = head
```

Traversal

```
1 def print_elements(l):
2     current = l.head
3     while current != None:
4         print(current.value)
5         current = current.Next
```

Search

```
1 def search(l, v):
2     current = l.head
3     while current != None:
4         if current.value == v:
5             return current
6         current = current.Next
7     return None
```

Insert

```
1 def insert_after_node(node, value):
2     new_node = Node(value, node.Next)
3     node.Next = new_node
4
5 def insert_after_value(l, value, new_value):
6     node = search(l, value)
7     if node != None:
8         insert_after_node(node, new_value)
```

Remove

```
1 def remove_after_node(node):
2     to_remove = node.Next
3     if to_remove == None:
4         return
5     else:
6         node.Next = to_remove.Next
7
8 def remove_after_value(l, value):
9     node = search(l, value)
10    if node != None:
11        remove_after_node(node)
```

Convert List to Linked List

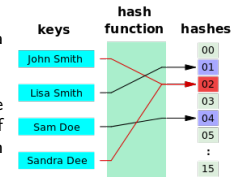
```
1 """Create and return linked list from list a."""
2 def create_from_list(a):
3     l = Linked_List(Node(a[0]))
4     last = l.head
5     for v in a[1:]:
6         last.Next = Node(v)
7         last = last.Next
8     return l
```

5.4 Hash Tables

- fast access to elements
- Unsorted, unordered data
- (O) only in expectation

Given an element 'x' to be stored, a Hash Table 'M' of length 'l' and a hash function 'f(x)':

```
1 M[f(x) % l] = x
```

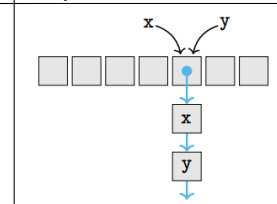


example of collision

5.4.1 Collision handling

Entry at calculated index may already contain an element

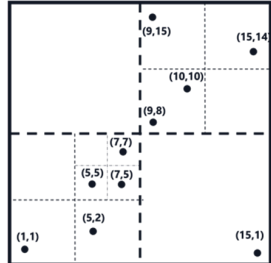
Probing: next available index is chosen



5.5 Quadtree

- Faster searching for points within a rectangle
- Efficient for querying multiple rectangles

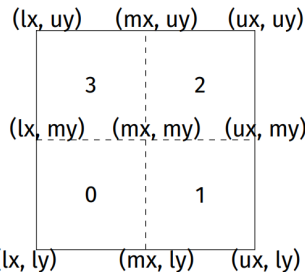
Data containing two values are inserted into squares. The squares have a maximum number of points they may contain. If a new point would exceed that limit, the square is subdivided into four smaller squares.



Implementation

```
1 class QuadTree:
2     def __init__(self, l, u, max_cap):
3         self.l = l #lower left corner coordinate
4         self.u = u #upper right corner coordinate
5         self.m = max_cap
6         self.points = []
7         self.children = None
8         self.count = 0 #points within quadtree
```

Insert



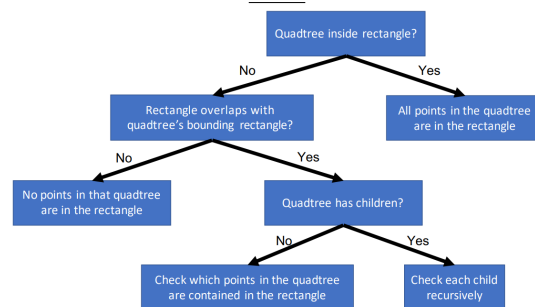
```
1 def subdivide(self):
2     lx, ly = self.l
3     ux, uy = self.u
4     mx, my = (lx + ux) / 2, (ly + uy) / 2
5     self.children = [None, None, None, None]
6     self.children[0] = QuadTree((lx, ly), (mx, my),
7                                 self.m)
8     self.children[1] = QuadTree((mx, ly), (ux, my),
9                                 self.m)
10    self.children[2] = QuadTree((mx, my), (ux, uy),
11                                self.m)
12    self.children[3] = QuadTree((lx, my), (mx, uy),
13                                self.m)
```

```
11 def get_index(self, point):
12     lx, ly = self.l
13     ux, uy = self.u
14     mx, my = (lx + ux) / 2, (ly + uy) / 2
15     px, py = point
16     if lx <= px < mx:
17         if ly <= py < my:
18             return 0
19         elif my <= py <= uy:
20             return 3
21     elif mx <= px <= ux:
22         if ly <= py < my:
23             return 1
24         elif my <= py <= uy:
25             return 2
26     return None
```

```
29 def insert(self, point):
30     if self.children is not None:
31         index = self.get_index(point)
32         if index is not None:
33             self.children[index].insert(point)
34             self.count += 1
35             return
36     self.points.append(point)
37     if len(self.points) > self.m:
38         self.subdivide()
39     for point in self.points:
40         index = self.get_index(point)
```

```
41 if index is not None:
42     self.children[index].insert(point)
43     self.count += 1
44 self.points = []
```

Search



```
1 def rect_intersect(self, lr, ur):
2     def intervals_overlap(a_start, a_end, b_start,
3                           b_end):
4         a_end_in_b = b_start <= a_end <= b_end
5         b_end_in_a = a_start <= b_end <= a_end
6         return a_end_in_b or b_end_in_a
7     self.lx, self.ly = self.l
8     self.ux, self.uy = self.u
9     rect_lx, rect_ly = lr
10    rect_ux, rect_uy = ur
11    return (intervals_overlap(self.lx, self.ux,
12                              rect_lx, rect_ux) and
13            intervals_overlap(self.ly, self.uy,
14                              rect_ly, rect_uy))
```

```
15 def covered_by_rect(self, lr, ur):
16     self.lx, self.ly = self.l
17     self.ux, self.uy = self.u
18     rect_lx, rect_ly = lr
19     rect_ux, rect_uy = ur
20     return (rect_lx <= self.lx and
21             rect_ly <= self.ly and
22             self.ux <= rect_ux and
23             self.uy <= rect_uy)
```

```
26 def count_in_rect(self, lr, ur):
27     if self.rect_intersect(lr, ur):
28         if self.children is None:
29             pts = [p for p in self.points if lr[0]
30                  <= p[0] <= ur[0] and lr[1] <= p[1] <= ur[1]]
31             return len(pts)
32         else:
33             return sum(c.count_in_rect(lr, ur) for
34                          c in self.children)
35     elif self.covered_by_rect(lr, ur):
36         return self.size()
37     else:
38         return 0
```

6 Runtime analysis

6.1 Runtime analysis

6.1.1 upper, lower and tight bound

To measure the performance of an algorithm, we use big O notation. Let g be the relationship time vs input size for an algorithm.

- Upper bound (If g does not grow faster than c*f): $g = O(f)$
- Tight Bound (If g grows about the same as c*f): $g = \Theta(f)$
- Lower Bound (If g does not grow slower than c*f): $g = \Omega(f)$

6.1.2 Asymptotic behaviour of functons

For all functions: $\lim_{n \rightarrow \infty}$

$$\log(n) < \sqrt[n]{n} < n < n \cdot \log(n) < n^2 < a^n < n! < n^n$$

$$n = O(n^2) \rightarrow f(n) = O(g(n)) \rightarrow \lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} \leq C, C \in \mathbb{R}$$

Useful conversions:

$$B^a \cdot B^b = B^{a+b} \quad \log_B(a \cdot b) = \log_B(a) + \log_B(b)$$

$$\frac{B^a}{B^b} = B^{a-b} \quad \log_B\left(\frac{a}{b}\right) = \log_B(a) - \log_B(b)$$

$$(B^a)^b = B^{a \cdot b} \quad \log_B(a^r) = r \cdot \log_B(a)$$

6.1.3 Code Runtime Analysis

Assumptions of 'time cost':

```
1 #Comparisons cost 1
2 if x==1
3 if x>1
4 #Mathematical operations cost 1
5 6 + 4
6 a * b
7 #Assignment cost 1
8 x = 5
```

Example of runtime analysis on selection sort:

$$T(n) = 1 + \sum_{i=0}^{n-1} \left(1 + \left(\sum_{j=i+1}^{n-1} 2 \right) + 1 \right) = \frac{n(n-1)}{2} = \Theta(n^2)$$

Useful Sums

$$\sum_{i=0}^{n-1} 1 = n = \Theta(n)$$

$$\sum_{i=0}^n i = \frac{n \cdot (n+1)}{2} = \Theta(n^2)$$

$$\sum_{i=0}^n i^2 = \frac{1}{6}n(n+1)(2n+1) = \Theta(n^3)$$

$$\sum_{i=0}^n \sum_{j=0}^i j = \frac{1}{12}n(n+1)(2n+4) = \Theta(n^3)$$

Telescoping (recursive code)

Example: Binary search

- Find cases
- Find general expression for $T(n)$
- Find condition for base case $T(1)$

$$1. T(n) = \begin{cases} d & n = 1 \\ T\left(\frac{n}{2}\right) + c & n > 1 \end{cases}$$

$$T\left(\frac{n}{2}\right) = T\left(\frac{n}{4}\right) + c$$

$$2. \Rightarrow T(n) = T\left(\frac{n}{2^i}\right) + i \cdot c$$

$$T\left(\frac{n}{n}\right) = T(1) = d$$

$$3. \Rightarrow 2^i = n \Rightarrow i = \log(n)$$

$$\Rightarrow T(n) = d + \log(n) \cdot c = \Theta(\log(n))$$

6.2 Runtime of specific algorithms

6.2.1 List

Search element: see 6.2.3 Search Algorithms on Lists

Operation	Runtime	
Index Access	$O(1)$	
Insertion	$O(n)$	
Removal	$O(n)$	
Remove last element	$O(1)$	

6.2.2 Sorting Algorithms on Lists

Selection Sort

Case	Description	Runtime
Worst-case	A is reverse sorted	$\Theta(n^2)$
Average-case	-	$\Theta(n^2)$
Best-case	A is already sorted	$\Theta(n)$

Insertion Sort

Case	Description	Runtime
Worst-case	A is reverse sorted	$\Theta(n^2)$
Average-case	-	$\Theta(n^2)$
Best-case	A is already sorted	$\Theta(n)$

Bubble Sort

Case	Description	Runtime
Worst-case		
Average-case		
Best-case		

Mergesort

Case	Description	Runtime
All cases	-	$\Theta(n \cdot \log(n))$

Quick Sort

Case	Description	Runtime
Worst-case	Pivot always min/max value	$\Theta(n^2)$
Average-case	Pivot chosen randomly	$\Theta(n \cdot \log(n))$
Best-case	Pivot always median value	$\Theta(n \cdot \log(n))$

Heap Sort

heapify(a) and HeapSort(a)

Case	Description	Runtime
All cases	-	$\Theta(n \cdot \log(n))$

6.2.3 Search Algorithms on Lists

Linear Search

Case	Description	Runtime
Worst-case	$b = a[n]$ is at end of array	$\Theta(n)$
Average-case	-	$\Theta(n)$
Best-case	$b = a[0]$ at begining of array	$\Theta(1)$

Binary Search

Case	Description	Runtime
Worst-case	b is max / min value	$\Theta(\log(n))$
Average-case	-	$\Theta(\log(n))$
Best-case	b is median value	$\Theta(1)$

6.2.4 Binary Search Tree

Search Node

Case	Description	Runtime
Worst-case	Tree is degenerated	$\Theta(n)$
Average-case	Tree is balanced	$\Theta(\log(n))$
Best-case	key is at root	$\Theta(1)$

Insert Node

Case	Description	Runtime
Worst-case	Tree is degenerated	$\Theta(n)$
Average-case	Tree is balanced	$\Theta(\log(n))$
Best-case	Tree is empty	$\Theta(1)$

Remove Node

Runtime is determined by the Runtime of 'Find symmetric successor'. 'h' is the height of the Tree

Case	Description	Runtime
All cases	-	$\mathcal{O}(h)$

Traversal

Case	Description	Runtime
All cases	-	$\Theta(n)$

6.2.5 Heap

Insert Element

Case	Description	Runtime
Worst case	$\log(n)$ swaps	$\mathcal{O}(\log(n))$

Remove Element

Case	Description	Runtime
Worst case	$\log(n)$ swaps	$\mathcal{O}(\log(n))$

Heapify

Case	Description	Runtime
All cases	-	$\Theta(n)$

Sort a heap

Used in 'Heapsort'. SiftDown traverses $\log(n)$ nodes and is called n times

Case	Description	Runtime
All cases	-	$\Theta(n \cdot \log(n))$

6.2.6 Linked List

Operation	Runtime
Access	$\mathcal{O}(n)$
Search	$\mathcal{O}(n)$
Insertion	$\mathcal{O}(n)$
Insertion at head	$\mathcal{O}(1)$
Removal	$\mathcal{O}(n)$
Removal at head	$\mathcal{O}(1)$

6.2.7 Hash Table

Operation	Best Case	Worst Case
Search	$\mathcal{O}(1)$	$\mathcal{O}(n)$
Insertion	$\mathcal{O}(1)$	$\mathcal{O}(n)$
Removal	$\mathcal{O}(1)$	$\mathcal{O}(n)$

7 Programming Concepts

7.1 Compiled vs Interpreted

Compiled (C++):

- Program code is translated to assembly.
- Assembly is executed.
- Single translation, with optimizations.
- Usually, higher performance

Interpreted (Python):

- Program code executed together with translation.
- Translation is repeated each time.
- Quick and easy to make minor changes.

7.2 Static vs Dynamically Typed

C++ is statically typed:

- Each element has a type defined by the programmer.
- Types used fitting together correctly is checked at compilation, yielding compile time errors (happen during the program itself) if wrong.

Python is dynamically typed:

- Elements have no type in advance.
- At runtime the type is chosen.
- Type changeable at runtime.
- Depending on the type when executing, there may be runtime errors (happen during the program).
- Errors are more difficult to debug, do not happen all the time.

7.3 Generic Programming

The goal of generic programming is to make code as widely usable as possible (no need for new functions for different types). Can be done with templates in C++.

No need to do anything in Python thanks to dynamic typing.

7.4 Functional programming

Pass functions as parameters to functions. Example:

```
1 numbers = [1, 2, 3, 4, 5]
2 def square(x):
3     return x ** 2
4 squared_numbers = list(map(square, numbers))
5 #map function takes square function as a parameter
```

7.4.1 Lambda Expressions

Lambda functions are small functions without a specific name, useful to pass into a function as parameter.

```
1 lambda arguments : expression
2 #Lambda with one argument:
3 n = [1,2,3,4,5]
4 sqrd_numbers = map(lambda x : x**2, n)
5 print(list(sqrd_numbers)) #[1,4,9,16,25]
6 #Lambda with multiple arguments:
7 y = lambda x,y: x*y
8 y(5,3) #15
```

7.4.2 Common Functions

Map

- map(func, it) – applies a function on each element of a container.

```
1 n = [1,2,3,4,5]
2 sqrd_numbers = map(lambda x : x**2, n)
3 print(list(sqrd_numbers)) #[1,4,9,16,25]
```

Filter

- filter(func, it) – removes any elements that don't fulfil a condition.

```
1 n = [1,2,3,4,5]
2 even_numbers = filter(lambda x : x%2==0, n)
3 print(list(even_numbers)) #[2,4]
```

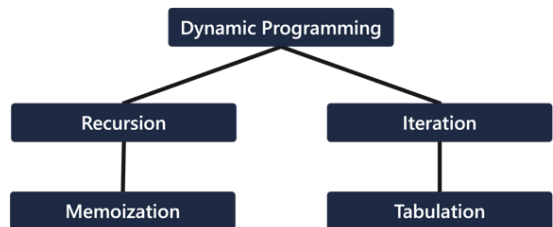
Reduce

- reduce(func,it) – recursively reduce a container to a single value by applying a function to two elements.

```
1 from functools import reduce
2 n = [1,2,3,4,5]
3 sum_numbers = reduce(lambda x,y: x + y, n) #15
```

8 Dynamic Programming (DP)

- “bottom-up” strategy
- iteratively solve smaller problems to solve progressively bigger problems
- Overlapping subproblems
- answers to smaller problems are stored in a table
- Example: Fibonacci



Convert code from recursive to dynamic programming:

- Look for repeating subproblems
- look for an optimal substructure to the solution – how does the answer of the problem depend on the answer of the sub-problems?
- Should the answers of the subproblems be stored in a list? A table?

- Flip the recursive implementation around to implement a bottom-up, iterative solution.

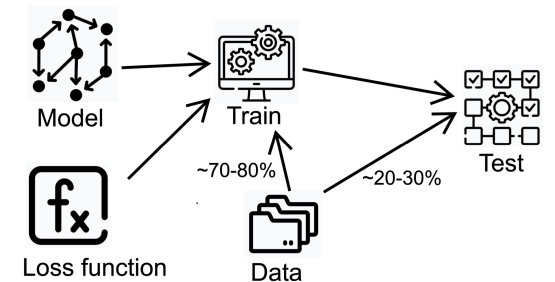
Example: dynamic programming approach to Fibonacci

```
1 F = [None] * (n+1)
2 #border cases
3 F[0] = 1
4 F[1] = 1
5 #Bottom-Up for loop
6 for i in range(2, n+1):
7     F[i] = F[i-1] + F[i-2]
8 fib_n = F[n]
```

9 Machine Learning (ML)

create functions which map inputs to desired outputs by analysing underlying patterns

- Regression: find output value $\in \mathbb{R}$ based on input. Example: given a house has 5 bedrooms, 1000 square meters, what is its price?
- Classification: assign element to a group. Example: Based on color and texture, is a mushroom edible?



```
1 #read data
2 import pandas as pd
3 data = pd.read_csv('data.csv')
4
5 #split data into input and result set
6 y = data["diagnosis"]
7 X = data.drop(columns=["diagnosis"])
8
9 #split data into test and train set
10 from sklearn.model_selection import train_test_split
11 X_train, X_test, y_train, y_test = train_test_split(X, y,
12 test_size=0.3, random_state= 42) #30% of data used for validation
13
14 #choose and train model
15 from sklearn import linear_model
16 model = linear_model.LinearRegression()
17 model.fit(X_train, y_train)
18
19 #validate and test model
20 from sklearn.metrics import accuracy_score
21 y_pred = model.predict(X_test)
22 score = accuracy_score(y_test, y_pred)
23
24 #CodeExpert grading based on y.final
25 X_final = pd.read_csv("X_final.csv")
26 y_final = model.predict(X_final)
27 return y_final
```

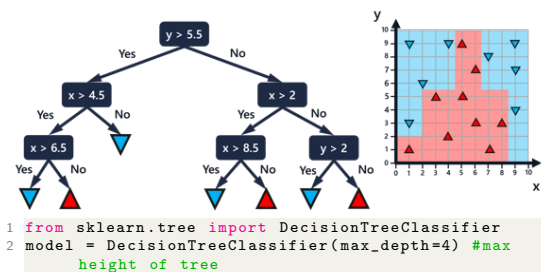
Loss function:

- `# wrongly classified examples`
`# total examples`
- Mean Square Error (MSE)
- L2 (Gauss Loss)

9.1 Models

9.1.1 Decision Tree Classifier

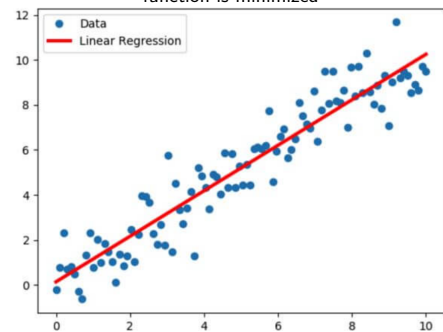
Constructs a decision tree to classify data based on features



9.1.2 Linear

Regression

Constructs a linear model $y = \beta_0 + \sum_{i=0}^n \beta_i x_i$ so that the loss function is minimized

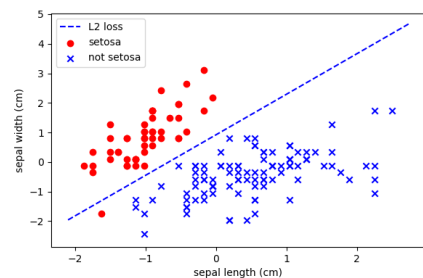


```

1 from sklearn import linear_model
2 model = linear_model.LinearRegression()

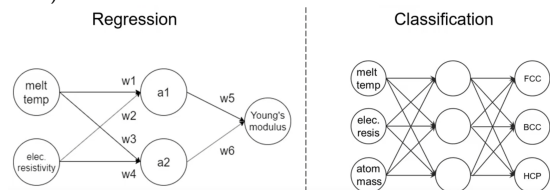
```

Classification



9.1.3 Neural Network

Constructs layers of neurons, where the output of each neuron is made non-linear through an activation function (Sigmoid, Relu, Tanh)



Regression

Usually no activation function in output layer

```

1 from sklearn.neural_network import MLPRegressor
2 #hidden_layer_sizes: array indicating number of neurons per layer
3 #activation: activation function
4 model = MLPRegressor(hidden_layer_sizes = [4,4], activation="relu")

```

Classification

Usually Sigmoid in output layer

```

1 from sklearn.neural_network import MLPClassifier
2 model = MLPClassifier(hidden_layer_sizes = [4,4], activation="logistic")

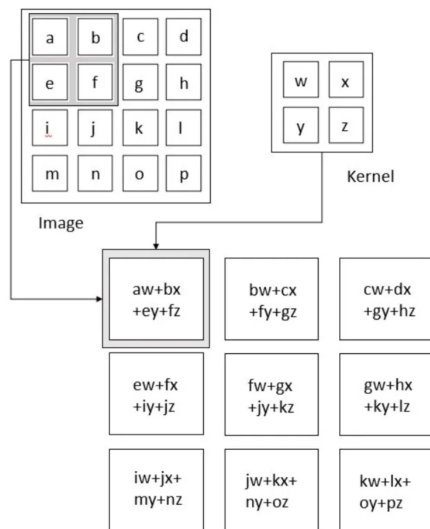
```

Convolutional Neural Network

Softmax can be used in output layer

Image processing, process to reduce data

Example using Joint Matrix (Faltungsmatrix):



One can also choose the max matrix. For the first field in the example this would yield the output: max(a, b, e, f)

9.2 Test and Validate

Accuracy score

#correctly classified / #total. 1 is the best score, 0 is the worst.

```
1 from sklearn.metrics import accuracy_score
```

R2 score

1 is the best score, the more negative the worse.

```
1 from sklearn.metrics import r2_score
```

Mean Squared Error (MSE)

0 is the best score. The bigger the worse.

```
1 from sklearn.metrics import mean_squared_loss
```

9.3 Over-/Underfitting and Cross Validation

- Underfitting: model is too simplistic to capture the underlying patterns in the data
- Overfitting: model memorizes random fluctuations in the data rather than capturing the underlying patterns

Result: No liable predictions from the model

Solution: Cross-Validation

- subdivide data set into smaller parts
- use n-1 parts for training and 1 remaining part for testing, trying out different settings

Fold 1	Testing set		Training set	
Fold 2	Training set	Testing set		Training set
Fold 3		Training set	Testing set	Training set
Fold 4			Training set	Testing set

9.4 Encoding

Ordinal Encoding

Assign number to a string in a dataset:

red → one, green → 2, blue → 3

+ easy, efficient

– In our example, the 'distance' between "red" and "blue" might seem bigger than between "green" and "blue"

Mean Encoding

red	True
red	False
red	True
green	False
green	False

$$\begin{aligned} \text{mean}_{\text{red}} &= \frac{\#(\text{red and True})}{\# \text{red}} \\ &= \frac{2}{3} \\ \text{mean}_{\text{green}} &= 0 \end{aligned}$$

+ meaning behind the encoding (probability of occurrence)

– ML Algorithm can access this probability → Overfitting

One-Hot Encoding

Color	On?	red	green	On?
red	True	1	0	True
red	False	1	0	False
red	True	1	0	True
green	False	0	1	False
green	False	0	1	False

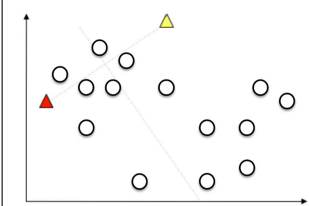
+ Same distance between all properties

– Tables can become very big, leading to slower training

9.5 Clustering

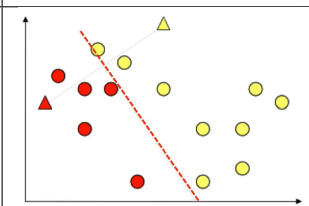
Given: $X = \{x_1, x_2, \dots, x_n\} \in \mathbb{R}^d$

Initialize K centroids (randomly): $z_i \in \mathbb{R}^d \rightarrow Z = \{z_1, z_2, \dots, z_K\}$



Assign each point to centroid it is closest to:

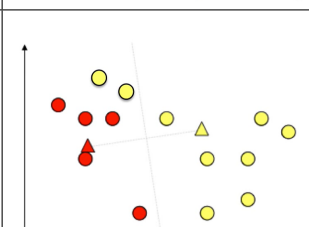
$$c(x) = \text{argmin} ||x - Z_i||$$



Place centroid in the mean of its assigned points:

$$S_j = x_i \in z_j$$

$$z_j = \frac{1}{\text{len}(S_j)} \sum x$$



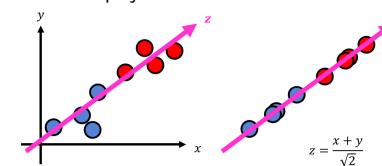
Repeat step 2 and 3 until the centroids have converged, that means no point is assigned to another centroid anymore
Overfitting: Each point has its own centroid
Loss function:

$$L = (\sum ||x_i - z_{x_i}||^2) + \lambda e^K$$

squared distance of point to its centroid + penalty for using too many centroids (prevents overfitting)

Principal Component Analysis (PCA)

Reduce dimension while keeping the maximum information: goal is to have the values as far apart from each other as possible after a linear projection on a line.



10 Numpy

Numpy is a Python package (equivalent to a C++ library) which supports operations with n-dimensional arrays and various computational methods.

Import the Numpy package:

```
1 import numpy as np
```

Now you can refer to functions/classes from Numpy using: "np".

10.1 Python Lists vs. Numpy Arrays

Numpy arrays are like Python lists. Below is a summary of the key differences.

Python Lists	Numpy Arrays
Variable size	Fixed size
Different element types	Single element type
Mathematical operations on single elements only	Mathematical operations on whole arrays
Primarily 1D	Multi-dimensional

10.2 Declaring Numpy Arrays

Using sequences

```
1 a = np.array([1, 2, 3, 4]) #array([1,2,3,4])
2 b = np.array(range(2,10,3)) #array([2,5,8])
3 c = np.array([[1,2],[3,4]]) #comparable to matrix
```

Using random numbers

```
1 R = np.random.random(10) #10 random numbers in [0,1)
2 R = np.random.uniform(-1,1,5) #5 random numbers between -1 and 1
3 R = np.random.randint(1,7,10) #10 random integers between 1 and 6
```

Using arange command

creates an array from the start to the stop value with a given step value, stop is **not** inclusive

```
1 R = np.arange(start, stop, step)
2 R = np.arange(2,10,3) #array([2,5,8])
3 #the same output as np.array(range(2,10,3))
```

Using linspace command

creates a numpy array with num equally spaced elements between start and stop, stop is inclusive:

```
1 R = np.linspace(start,stop,num)
2 #Step size = (stop-start)/(num-1)
3 R = np.linspace(3, 2, 3) #array([3, 2.5, 2])
```

10.3 Numpy Array operations

Common numpy array operations

```
1 a = np.array(10)
2 a.size() #10, return number of elements
3 a[5] #5, access element (1D array)
4 A = np.array([[1,2,3],[4,5,6]])
5 A[1,2] #6, access element (2D array)
6 A[1][2] #6, access element (2D array)
```

To access lines of a 2D array, see slicing.

Slicing

stop is **not** included

- 1D array: $A[\text{start} : \text{stop} : \text{step}]$

```
1 A = np.arange(10) #array([0,1,2,3,4,5,6,7,8,9])
2 A[2:5:2] #array([2,4])
```

- 2D array: $A[\underbrace{\text{start} : \text{stop}}_{\text{row}}, \underbrace{\text{start} : \text{stop}}_{\text{column}}]$

default values: start = 0, stop = len(A) if only one number is given, only the row / column corresponding to that index is taken

```
1 A = np.array([[1,2,3],[4,5,6],[7,8,9]])
2 A[1,:] #array([4,5,6]) #row at index 1
3 A[:,2] #array([3,6,9]) #column at index 2
4 A[0:2,1:3] #array([2,3],[5,6]) # rows 0 and 1, columns 1 and 2
```

Statistics

```
1 a = np.linspace(-4,-3) #array([-4,-3,-2])
2 a.min() #-4, Minimum
3 a.max() #-2, Maximum
4 a.sum() #-9, Sum
5 np.mean(a) #-3, Average
6 np.std(a) #0.81, Standard deviation
```

Mathematical Operations

Most mathematical operations are carried out element-wise:

```
1 A = np.array([[2,3,4],[6,7,6]])
2 B = np.array([[1,9,1],[2,3,9]])
3 C = np.array([[1, 4], [3, 4], [4,6]])
4 A+1
5 #array([[ 3,4,5], [7,8,7]])
6 A * 2
7 #array([[ 4,6,8], [12,14,12]])
8 A ** 4
9 #array([[ 16,81,256], [1296,2401,1296]])
10 np.sin(A)
11 #array([[ 0.909,0.141,-0.756],
12         [-0.279,0.657,-0.279]])
12 A + B
13 #array([[ 3,12,5], [ 8,10,15]])
14 A * B
15 #array([[ 2,27,4], [12,21,54]])
16 np.sum(A, axis = 0) #sum of the columns
17 #= A.sum(axis = 0) -> array([8,10,10])
18 np.sum(A, axis = 1) #sum of the rows
19 #= A.sum(axis = 1) -> array([9,19])
20
21 A @ C #matrix multiplication
22 A.dot(C) #matrix multiplication
23 #array([[ 27, 44], [51, 88]])
24
25 a = np.array([1,2,3])
26 b = np.array([3,4,6])
27 a.dot(b) #29, scalar product
```

Filtering

```
1 a = np.arange(7) #array([0,1,2,3,4,5,6])
2 f = a % 2 == 0 #f = [True, False, True, False, True, False, True]
3 a[f] #array([0,2,4,6])
```

11 Pandas

Pandas is a Python package which supports working with tabulated data.

To import pandas, use:

```
1 import pandas as pd
```

Now you can refer to classes and functions from the package using "pd".

11.1 read CSV file into dataframe

```
1 climate = pd.read_csv("climate.csv", sep=",",
2                       index_col=0, usecols=["time", ...])
3 # "sep" -> what characters values in the csv file are separated by. "index_col" -> what the index column will be. "usecols" -> what columns of the csv data will be #selected.
```

11.2 Pandas Dataframe operations

A dataframe can be thought of as a list within a list supporting access in more meaningful ways compared to using indices.

Example dataframe					Change Index Column			
Unnamed: 0					Unnamed: 0			
time					jan			
feb					feb			
0	0	1864	-7.10	-4.52	time			
1	1	1865	-3.47	-6.25	1864	0	-7.10	-4.52
2	2	1866	-1.31	-0.42	1865	1	-3.47	-6.25
3	3	1867	-3.87	0.56	1866	2	-1.31	-0.42
4	4	1868	-5.46	-1.53	1867	3	-3.87	0.56
...	1868	4	-5.46	-1.53
157	157	2021	-3.56	NaN
					2021	157	-3.56	NaN

table of the climate, entries accessible via index. The left-most column is known as the "index column".

```
1 climate2 = climate.
2   set_index("time")
3   #creates copy
```

table of the climate, entries accessible via time

Rename Columns

```
1 climate = climate.rename(columns={"old_index_name":
2   "new_index_name", ...})
3 #renames the "old_index_name" column to "new_index_name"
```

rename all columns (len(data.columns) = number of columns must be true)

```
1 data.columns = ["Date", "January", "February", ...]
```

Access Dataframe Elements

```
1 climate["jan"].iloc[3] #Access single element
2 climate["feb"] #Access single column (type: Series)
3 climate[["jan", "mar"]] #multiple columns (Dataframe)
4 climate.iloc[3] #Access single row (Series)
5 climate.iloc[1:4] #multiple rows (Dataframe)
6 climate.iloc[4:7,1:2] #Access subtable (Dataframe)
7 #gets rows 4,7 with data only from column 1
8 #Access subtable using index column values and column name (Dataframe)
9 climate2 = climate.set_index("time")
10 climate2.loc[1864:1868, "jan": "mar"]
11 #includes the rows labeled with 1864 until and including #1868, the columns from "jan" until and including "mar"
```

Filter Dataframes

Filter rows:

```
1 climate[climate["jan"]>2]
2 #filters out the rows with values in the "jan" column less than 2
```

- Example: All entries in "jan" with values more than 2:

```
1 climate["jan"][climate["jan"]>2]
```

Dealing with Invalid Data

Convert all the values in a column to numeric:

```
1 data[column] = pd.to_numeric(data[column], errors="coerce")
2 #converts all the values to numeric values. #errors="coerce" -> converts values which cannot be #converted to NaN.
```

Delete all rows containing NaN entries:

```
1 data.dropna(axis = 0, how="any")
2 #how="any" -> delete row if any value is NaN.
3 #how="all" -> delete row if all values are NaN
4 #axis = 1 -> delete column instead of row
```

- Fill all entries containing NaN with a value:

```
1 data.fillna(0) #fill any NaN entries with 0
```

Modify Dataframes

Add a column:

```
1 climate["new_col"] = climate["time"] + climate["jan"]
2 #"new_col" is a new column whose values are #those of the "time" and "jan" column added
```

Delete a column:

```
1 climate = climate.drop(columns=["time"])
2 #delete the "time" column
```

Add a row:

```
1 d = {"mar":34, "jan":23}
2 climate.append(d, ignore_index=True)
3 #adds another row with the values 34 for "mar" and 23 for "jan". Other entries are NaN
```

Delete a row:

```
1 climate = climate.drop(climate.index[0])
2 #deletes row 0
```

Transpose the dataframe:

```
1 climate = climate.T
```

Analyse Data

Sum of all the entries in each column (type: Series):

```
1 climate.sum()
```

Maximum of all the entries in each column (type: Series):

```
1 climate.max()
```

Create a dataframe summarizing the max and sum for each column:

```
1 climate.agg(["max","sum"])
2 #A dataframe containing the same columns as climate
3 #with row 0 containing the max of the column and row 1 #containing the sum of the column. The strings in the
4 #list should be names of valid pandas Series functions.
```

Get statistical information for each column (type: Dataframe):

```
1 climate.describe()
2 #includes a variety of statistical measures
```

Sort a dataframe according to entries in a specific column(s):

```
1 climate = climate.sort_values(["time", "jan"], ascending=False)
2 #sorts the rows by "time" in descending order. If two #entries for "time" are equal, then the rows are sorted #by "jan"
```

Split a dataframe into groups based on a specified column and perform a computation on each group:

```
1 data.groupby("column").sum()
2 #groups data based on the entries for "column" and #calculates the sum for each group.
3 data.groupby("column").max()
4 #groups data based on the entries for "column" and #calculates the max for each group.
```

12 Matplotlib

Matplotlib is a Python package allowing you to visualize a variety of things: from functions to animations. To import matplotlib, use:

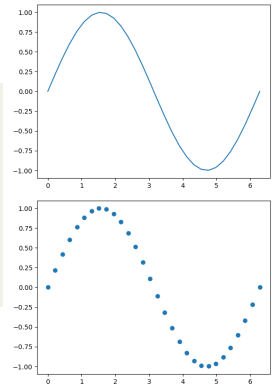
```
1 import matplotlib.pyplot as plt
```

Now you can refer to classes and functions from the package using "plt".

12.1 Line and Scatter Plots

To graph two numpy arrays, one representing the x values and the other representing the corresponding y values:

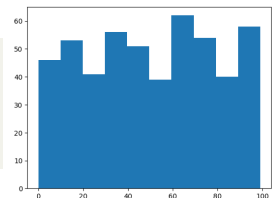
```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 X = np.linspace(0,2*np.pi,30)
5 Y = np.sin(X)
6
7 fig, ax = plt.subplots()
8 ax.plot(X,Y) #line plot
9 ax.scatter(X,Y) #scatter plot
```



12.2 Histogram Plots

To plot a histogram, use the following template:

```
1 fig, ax = plt.subplots()
2 X = np.random.randint(0, 100, 500)
3 # low: 0, high: 100, size: 500
4 ax.hist(X, bins=10)
5 plt.show()
```



12.3 Graph Styling

```
1 fig, ax = plt.subplots()
2 ax.set_title("title") #add title
3 ax.set_xlabel("x label name") #set x-label
4 ax.set_ylabel("y label name") #set y-label
5 ax.legend() #add legend
6 #requires that you labeled your plots, i.e.: when calling
7 #ax.plot(X,Y, label="name of function").
```