

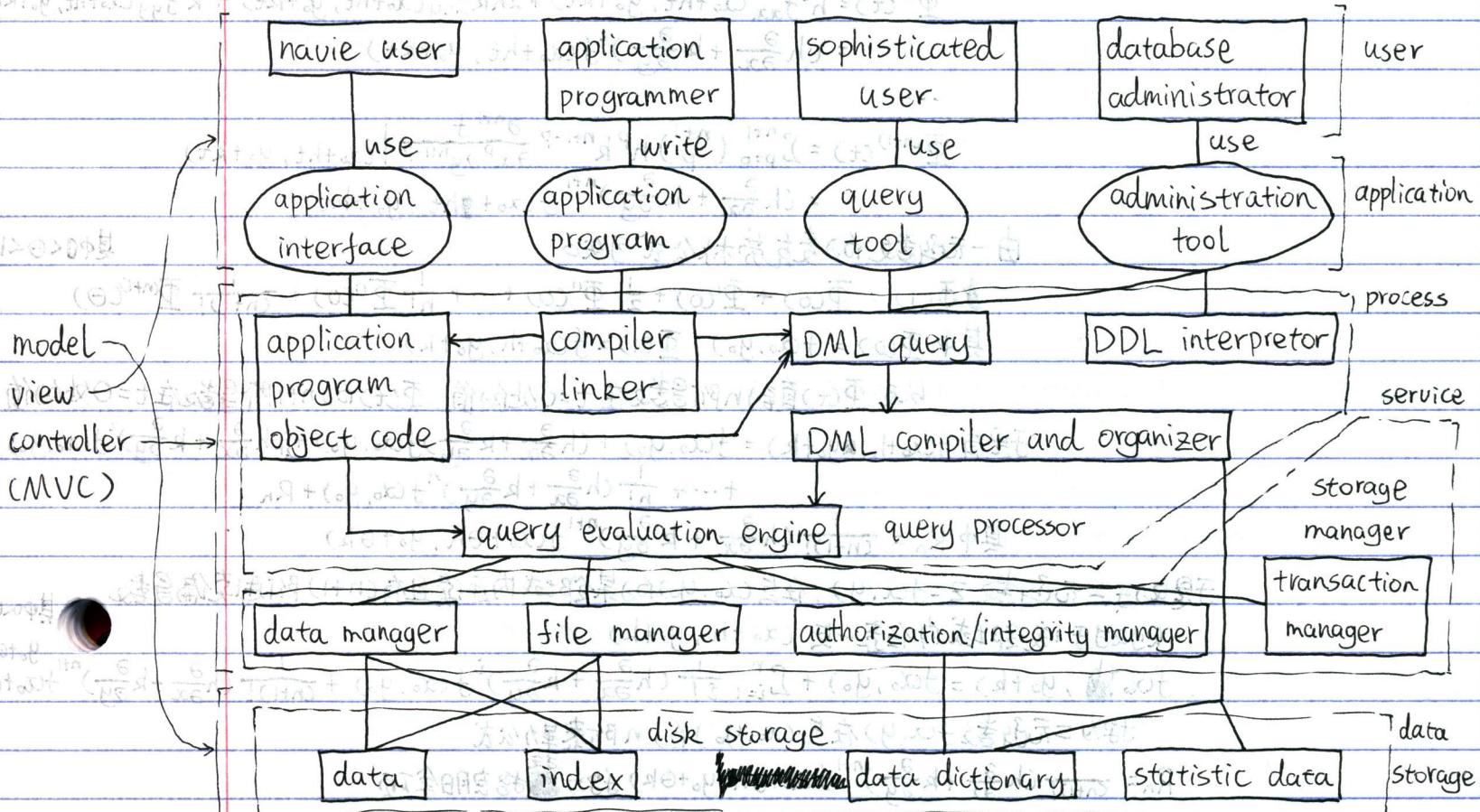
Python - P46

2019 - end of class

- knowledge pattern: characterization (summarizing data) + discrimination (comparing data)
- association: measuring frequency of data occurrence
- classification: modeling data as class
- clustering: high group similarity, low external similarity
- data operation: cleaning: to remove noise and inconsistent data
- integration: where multiple data source may be combined
- selection: where data relevant to task retrieved from database
- transformation: where data transformed into form appropriate for mining
- mining: intelligent method applied in order to extract data pattern
- evaluation: to identify truly interesting pattern
- presentation: where visualization of knowledge is presented

data definition language (DDL): to specify database schema

data manipulation language (DML): to express query and update



Python - P47

attribute type

nominal : relating to name, no meaningful order, not quantitative

binary : nominal attribute with only two categories / states

symmetric : equally valuable, carry same weight

asymmetric : not equally important

ordinal : meaningful order or ranking, no magnitude between successive values

numeric : quantitative, represented in integer / real value

interval scale

ratio scale

ratio : unable to calculate ratio

able to calculate ratio

zero point : arbitrary zero point

Absolute zero / character of origin

mean : arithmetic mean

geometric/harmonic mean

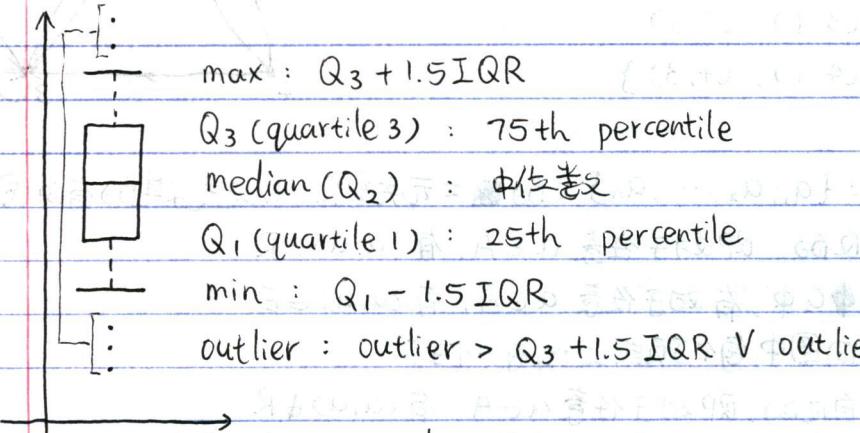
measurement : size and magnitude as +

size and magnitude as one factor

+ multiple factors of defined unit

of one defined unit in terms of another

box plot



四分位距 (Interquartile, IQR) : $Q_3 - Q_1$, 为一种稳健统计 (robust statistic)

四分位差 (quartile deviation, QD) : $\frac{Q_3 - Q_1}{2}$

data issue

incomplete / missing : lacking attribute value / certain attribute of interest

/ containing only aggregate data

inaccurate / noisy : containing error / value deviate from expected

inconsistent : containing discrepancy, inconsistency in department codes

handling missing data : ignore value / tuple

fill missing value manually

using global constant

using measure of central tendency for attribute

using attribute mean / median for all samples

using most probable value

Python - P48

data smoothing technique for noisy data

binning : smooth sorted data value by consulting neighborhood

regression: conform data value to function

outlier analysis : require to find out value differ considerably from other values

data integration : match or relate object from different sources

entity identification : not all entities identified consistently

functional dependency and referential constraint in source system

may not match in target system

correlation analysis : measure how attribute vary depending on another

data reduction , dimensionality reduction : attribute subset selection:

attribute construction : small set of more useful attribute

derived from original set

numerosity reduction : data replaced by alternative, smaller representation

using parametric / nonparametric model

data transformation : aggregation : summary operation applied to data

normalization : data scaled to fall within specified range

generalization : low-level "primitive" data replaced by

high-level concept through concept hierarchy

knowledge discovery : extraction of knowledge from data by detecting interesting pattern

understandable / valid / useful / novel

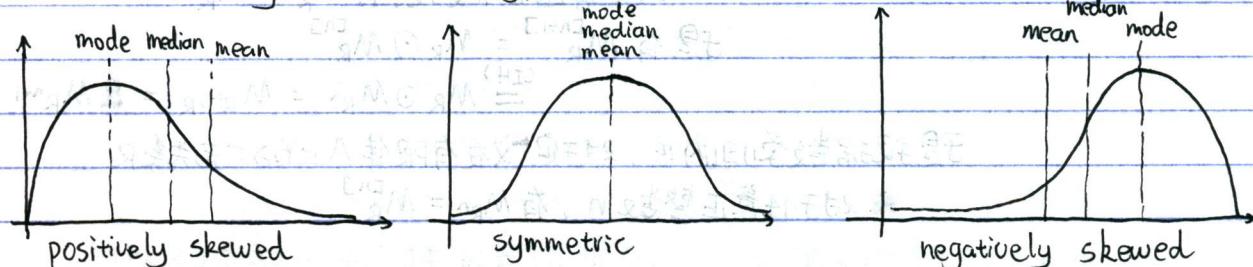
data characteristic , volume : staggering amount of data being generated

jerk or jounce - acceleration of acceleration of amount of data

complexity : structure / spatial / temporal

dimensionality : number / type / scale of attribute

centrality



Python - P49

standard

P49 - comparison

max (iterable, *[, key, default])
min (iterable, *[, key, default])
max (arg1, arg2, *args [,key])
min (arg1, arg2, *args [,key])
return largest/smallest item in iterable or of two or more argument
if one positional argument provided, should be iterable
the largest/smallest in the iterable returned
if two or more ~~one~~ positional arguments provided
the largest / smallest of positional arguments returned
key : specify one-argument ordering function like used for list.sort()
default: specify object to return if provided iterable empty
ValueError raised if iterable empty and default not provided
if multiple items are maximal / minimal
return the first one encountered
consistent with other sort-stability preserving tool

sorted(c iterable, key = keyfunc, reverse = True/False) [0]
heapq.nlargest(1, iterable, key = keyfunc)
heapq.nsmallest(1, iterable, key = keyfunc)

pow (x, y [,z]), return x to the power y
if z present, return x to the power y, modulo z, more efficiently than
two-argument form pow(x,y) equivalent to power operator x ** y
if y negative, z must be omitted
if z present, x and y must be of integer type, and y must be non-negative

repr (object), return string containing printable representation of object
make attempt to return string yield object with same value when passed to eval()
otherwise, string enclosed in angle bracket contain name of type of object
together with additional information including name and address
class can control function return by defining __repr__ method

Python - P50

std.out

std.out = <stream object>

`print(*objects, sep=' ', end='\n', file=sys.stdout, flush=False)`

print objects to text stream file, separated by sep, followed by end
sep/end/file/flush must be given as keyword argument if present
all non-keyword argument converted to string like str(>) →
written to stream, separated by sep, followed by end
sep/end must be string, can be None to use default value
if no objects given, just write end
file must be object with method write(string)
if not present or None, use ~~None~~ sys.stdout
print() can not be used with binary mode file object
use file.write(...) instead
output buffered usually determined by file
if flush is True, stream flushed forcibly

`str(object='')`

`str(object=b'', encoding='utf-8', errors='strict')`

return string version of object, or empty string if object not provided
if neither encoding nor errors given

`str(object) return object.__str__()`

"informal" or nicely printable string representation of object

fall back to returning repr(object) if object not have method __str__()
if at least one of encoding or errors given

object should be bytes-like object, so bytes / bytearray

if object is bytes / bytearray object

`str(object, encoding, errors)` 等价于 bytes. ~~encoding, errors~~ decode

otherwise bytes object underlying buffer object obtained before calling bytes.decode
passing bytes object to str(), without encoding or errors

fall under first case of returning informal string representation

also -b command-line option to Python

`reversed(seq)`, return reverse iterator

seq must be object has __reversed__() method

or support sequence protocol

starting at 0

`len()` method and `__getitem__(index)` method with integer argument

Python - P51

property - ~~containing~~

class property(fget=None, fset=None, fdel=None, doc=None) return property attribute
fget : function for getting attribute value
fset : function for setting attribute value
fdel : function for deleting attribute value
doc : create doc string for attribute
otherwise will copy docstring of fget (if exist)

so class C:

```
class C:  
    def __init__(self):  
        self._x = None  
    def getX(self):  
        return self._x  
    def setX(self, value):  
        self._x = value  
    def delX(self):  
        del self._x  
    x = property(getX, setX, delX, "it's _x")
```

can create read-only property using `property` as decorator

@property decorator turn method into getter for read-only attribute with same name
decorator take docstring of method as docstring of property

decorator create copy of property with corresponding accessor function

sets to decorated function with same name as original property

so in class C:

```
@property  
def x(self):  
    """ it's _x """  
    return self._x  
  
@x.setter  
def x(self, value):  
    self._x = value  
  
@x.deleter  
def x(self):  
    del self._x
```

Python - P52

STUDY

TEST - examination

point

round (number [, ndigits]), return <number> rounded to <ndigits> precision after decimal
if ~~ndigits~~ omitted or None, return nearest integer to input
for built-in type supporting round()

value rounded to ~~ndigits~~ closest multiple of 10 to power minus ndigits
if two multiple equally close, rounding done toward even choice

↳ round(0.5), round(-0.5) → 0

↳ round(1.5), round(-1.5) → 2

any integer value (positive/0/negative) valid for ndigits

if ndigits omitted or None, return value is integer

otherwise/return value same type as <number>

for general Python object number

delegate to number.__round__

sorted (iterable, *, key = None, reverse = False)

return new sorted list from item in iterable

optional argument key, reverse must be specified as keyword argument

key: specify function of one argument

used to extract comparison key from each element in iterable

default value None: compare element directly

reverse: if set to True, list element sorted as if comparison reversed

use functools.cmp_to_key() to convert old-style cmp function to key function

built-in sorted() function guaranteed to be stable

guarantee not to change relative order of element compare equal

class slice(stop)

class slice(start, stop[, step]) range(start, stop, step)

return slice object representing set of indices specified by

Start, step default to None

slice object have read-only data attribute start, stop, step

return argument values, have no other explicit functionality

used by Numerical Python and third party extension

slice object generated when extended indexing syntax used

↳ a[start : stop : step] or a[start : stop, i]

itertools.islice() for alternate version return iterator

Python - P53

graph theory

Floyd-Warshall

graph centrality, eccentricity (偏心率): 对于连通图 $G = (V, E)$ 的点, $v \in V$, 记点 v 的偏心率 $E(v)$ 为 v 到其他顶点的距离最大值

$$\text{即有 } E(v) = \max \{ d_{uv} \mid u \in V \wedge u \neq v \}$$

radius (半径): minimum eccentricity

$$\text{即 } r(G) = \min \{ E(v) \mid v \in V \}$$

diameter (直径): maximum eccentricity

$$\text{即 } d(G) = \max \{ E(v) \mid v \in V \}$$

(注: 直径也是连通图中最远的两点之间的距离)

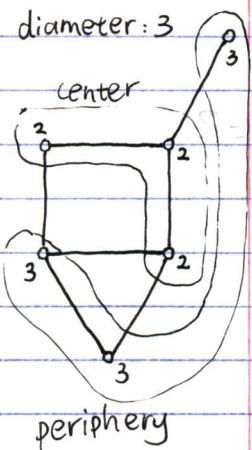
通常需要求任意两点之间的最短路径长度, 则长度的最大值为图的直径

center node (中心点): 偏心率等于 r 的顶点

即对于 $v \in V$, 有 $E(v) = r(G)$, 则 v 为 G 的中心点

periphery node (边缘点): 偏心率等于 d 的顶点

即对于 $v \in V$, 有 $E(v) = d(G)$, 则 v 为 G 的边缘点



eccentricity centrality: minimize maximum distance

closeness centrality: minimize distance to farthest node

betweenness: shortest paths between all pairs

prestige centrality: measure importance or rank of node

recursively measure indegree of node

converge to dominant: eigenvector / eigenvalue

power iteration approach: ratio of max entry in iteration k

to max entry of iteration $k-1$

yield eigenvalue estimate

pagerank centrality: prestige computation for web search

probability of random surfer landing on page

recursively depend on inbound pagerank

fully defined as combination: normalized prestige

random jump

probability of jump to any other page any means

Python - P54

2017-10-18

HITS

(hyperlink-induced topic search), 超文本敏感话题搜索

对于 hub and authority factor score

对于用户输入关键词后，算法对返回的匹配页面计算两种值

枢纽值(hub score)：所有导出链接指向页面的权威值之和

权威值(authority score)：所有导入链接所在页面的枢纽值之和

两种值相互依存，相互影响

to begin ranking, $\text{H}_p \text{ auth}(p) = 1$ and $\text{hub}(p) = 1$

in order to calculate hub/authority score of each node

repeated iteration of Authority Update / Hub Update applied

k-step application of Hub-Authority algorithm

entail applying for k times first Authority Update and then Hub Update

Authority Update Rule: authority score is sum of hub score of page p point to p

$$\text{H}_p \text{ auth}(p) = \sum_{q \in P_{\text{to}}} \text{hub}(q), P_{\text{to}}: \text{all page link to page } p$$

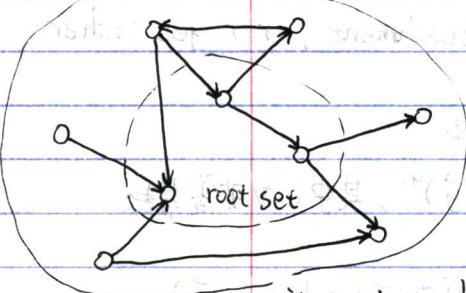
Hub Update Rule: hub score is sum of authority score of page p point to

$$\text{H}_p \text{ hub}(p) = \sum_{q \in P_{\text{from}}} \text{auth}(q), P_{\text{from}}: \text{all page } p \text{ link to}$$

normalization: final hub-authority score of node determined

after infinite repetition of algorithm

base set



as directly and iteratively applying Update Rule lead to diverging value

necessary to normalize matrix after every iteration

value obtained from this process eventually converge

$$\text{H}_p \text{ auth}(p) = \text{auth}(p) / \sqrt{\sum_q [\text{auth}(q)]^2}$$

$$\text{H}_p \text{ hub}(p) = \text{hub}(p) / \sqrt{\sum_q [\text{auth}(q)]^2}$$

iterative algorithm based on linkage of document on the web

like Page, Brin's PageRank

by search item

query dependent: hub/authority score resulting from link analysis influenced

as corollary, executed at query time, not at indexing time

with associated hit on performance accompany query-time processing

not commonly used by search engine

compute two scores per document (hub/authority) as opposed to single score

processed on small subset of 'relevant' document not all document (PageRank)

'Focused subgraph' or base set

Python - P55

real graph model : small world property : short path between node pair
scale logarithmically with number node

ultra-small-world much less than $\lg(n)$

scale-free property : scale-free degree distribution
follow power law

intuitively : vast majority of nodes - very small degree

few "hub nodes" with high degree

clustering effect : two nodes likely to be connected if share neighbour
hierarchical clustering of nodes

strong hub with weak peripheral : sparsely connected node ; smaller degree

strong hub with strong cluster : part of higher clustering area

few hub nodes : higher degree

going to small nodes to make short connect higher clustering area

or short, going to small nodes to make short connect higher clustering area

random graph model : Erdős-Rényi : small-world, no clustering

Watts-Strogatz : high clustering, not small

Barabási-Albert : preferential attachment, rich get richer

polynomial kernel : homogeneous : $K_q(\vec{x}, \vec{y}) = \phi(\vec{x})^T \phi(\vec{y}) = (\vec{x}^T \vec{y})^q$

inhomogeneous : $K_q(\vec{x}, \vec{y}) = \phi(\vec{x})^T \phi(\vec{y}) = c + \vec{x}^T \vec{y})^q$, 其中 c 为非零常数

l_2 norm : for $\phi(\vec{x})$ in feature space, $\|\phi(\vec{x})\|_2^2 = \phi(\vec{x})^T \phi(\vec{x}) = K(\vec{x}, \vec{x})$

$\Rightarrow \|\phi(\vec{x})\| = \sqrt{K(\vec{x}, \vec{x})}$

distance : for $\phi(\vec{x}_i), \phi(\vec{x}_j)$ in feature space,

$$\|\phi(\vec{x}_i) - \phi(\vec{x}_j)\|_2^2 = \|\phi(\vec{x}_i)\|_2^2 + \|\phi(\vec{x}_j)\|_2^2 - 2 \phi(\vec{x}_i)^T \phi(\vec{x}_j)$$

$$\text{根据余弦相似度公式 } S = K(\vec{x}_i, \vec{x}_i) + K(\vec{x}_j, \vec{x}_j) - 2 K(\vec{x}_i, \vec{x}_j)$$

$$\Rightarrow S(\phi(\vec{x}_i), \phi(\vec{x}_j)) = \|\phi(\vec{x}_i) - \phi(\vec{x}_j)\|_2^2 = \sqrt{K(\vec{x}_i, \vec{x}_i) + K(\vec{x}_j, \vec{x}_j) - 2 K(\vec{x}_i, \vec{x}_j)}$$

$$\text{且有 } \frac{1}{2} (\|\phi(\vec{x}_i)\|_2^2 + \|\phi(\vec{x}_j)\|_2^2 - \|\phi(\vec{x}_i) - \phi(\vec{x}_j)\|_2^2) = \phi(\vec{x}_i)^T \phi(\vec{x}_j) = K(\vec{x}_i, \vec{x}_j)$$

$$\text{mean : } \|\vec{\mu}_{\phi}\|_2^2 = \vec{\mu}_{\phi}^T \vec{\mu}_{\phi} = [\frac{1}{n} \sum_{i=1}^n \phi(\vec{x}_i)]^T [\frac{1}{n} \sum_{i=1}^n \phi(\vec{x}_i)]$$

$$= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \phi(\vec{x}_i)^T \phi(\vec{x}_j) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K(\vec{x}_i, \vec{x}_j)$$

$$\text{variance : } \sigma_{\phi}^2 = \frac{1}{n} \sum_{i=1}^n \|\phi(\vec{x}_i) - \vec{\mu}_{\phi}\|_2^2 = \frac{1}{n} \sum_{i=1}^n [K(\vec{x}_i, \vec{x}_i) - \frac{2}{n} \sum_{j=1}^n K(\vec{x}_i, \vec{x}_j) + \frac{1}{n^2} \sum_{a=1}^n \sum_{b=1}^n K(\vec{x}_a, \vec{x}_b)]$$

$$= \frac{1}{n} \sum_{i=1}^n K(\vec{x}_i, \vec{x}_i) - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n K(\vec{x}_i, \vec{x}_j)$$

Python - P56

@staticmethod , transform method into static method ~~define name as string~~ ~~call first self~~

method form is function decorator; using @wrapper syntax

so class C:

 @staticmethod

 def f(c):

 print('call C.f()')

注意 Static method not receive implicit first argument

static method can be called either on class or on instance

so C.f() → 'call C.f()'

C().f() → 'call C.f()'

static method in Python similar to Java / C++

possible to call staticmethod() as regular function

needed where need reference to function from class body

want to avoid automatic transformation to instance method

so class C:

 builtins.open = staticmethod(open)

sum (iterable, /, start=0), sum start and item of iterable from left to right

iterable's item : normally number

start value not allowed to be string

good alternative to sum() for some use cases

preferred fast way to concatenate sequence of string by calling ''.join(seq)

add floating value with extended precision by calling math.fsum()

concatenate series of iterable by using itertools.chains()

vars ([object]), return __dict__ attribute for module, class, instance ~~or any other object~~

or any other object with __dict__ attribute

module / instance have updateable __dict__ attribute

other object may have write restriction on __dict__ attribute

class use types.MappingProxyType to prevent direct dictionary

without argument, vars() act like locals()

locals dictionary only useful for read since update to locals dictionary ignored