

Python - P57

QUESTION

`super ([type, object-or-type])`

return proxy object that delegate method call to parent / sibling class of type

useful for accessing inherited method that have been overriden in class

`object-or-type`: determine method resolution order to be searched

base class searched for member during lookup in the order

search start from class right after the `<type>`

`mro_()` attribute of `object-or-type` list method resolution search order

used by both `getattr()` and `super()`

dynamic and can change whenever inheritance hierarchy updated

to `object-or-type.mro_(): D -> B -> C -> A -> object`

`super(type, object-or-type)` search `C -> A -> object`

if `object-or-type` omitted, then `super` object returned unbound

if `object-or-type` is `object`

then `isinstance(object-or-type, type)` must be `TRUE`

if `object-or-type` is `type` (useful for class method)

then `issubclass(object-or-type, type)` must be `TRUE`

`super()` used to refer to parent class without naming explicitly, make code maintainable

use closely parallel use of `super()` in other programming language

or to support cooperative multiple inheritance in dynamic execution environment

unique to Python and not found in statically compiled language or only support single inheritance

possible to implement "diamond diagram", multiple base classes implement same method

good design dictate this method have same calling signature in every case

because order of calls determined at runtime

order adapt to change in class hierarchy

can include sibling class unknown prior to runtime

typical superclass call:

`class C(B):`

`def method(self, *args, **kw):`

`super().method(*args, **kw)`

等价于 `super(C, self).method(*args, **kw)`

通过方法重写机制 (MRO) 确保正确调用

Python - P58

Method

func - external

> Super

([type[, object-or-type]])

super() also work for attribute lookup

possible use case called descriptor in parent / sibling class

super() implemented as part of binding process for explicit dotted attribute lookup

↳ super().__getitem__(name)

by implementing own __getattribute__() method

for searching class in predictable order support cooperative multiple inheritance

undefined for implicit lookup using statement/operator, ↳ super().__name__

super() not limited to use inside method aside from zero argument form

two argument form specify argument exactly and make appropriate reference

zero argument form only work inside class definition

as compiler file in necessary detail to correctly retrieve class being defined

as well as accessing current instance for ordinary method

class type(object)

class type(name, bases, dict)

with one argument, return type of <object>

returned value is type object, generally same object as object.__class__

isinstance() built-in function recommended for testing type of object

since take subclass into account

with three argument, return new type object

essentially dynamic form of class statement

name : become __name__ attribute

bases : tuple itemize base classes, become __bases__ attribute

dict : dictionary as namespace containing definition for class body

and copied to standard dictionary to become __dict__ attribute

↳ C = type('C', (object,), dict(a=1, b=2))

class C(object):

a = 1, b = 2

b = 2

有 C = C()

C.a → 1

C.b → 2

Python - P59

itsuseful

1559 - additional

`zip(*iterables)`, return iterator of tuple
where i-th tuple contains i-th element from each of argument sequence/iterable
iterator stops when shortest input iterable exhausted

with single iterable argument, return iterator of 1-tuple
with no argument, return empty iterator

left-to-right evaluation order of iterable guaranteed

possible idiom for clustering data series into n-length group
using `zip(*[iterables]*n)`

repeat same iterator n times, each output tuple has result of n calls to iterator
has effect of dividing input into n-length chunk

`zip()` should only be used with unequal length input

when not care about trailing, unmatched value from longer iterable
or use `itertools.zip_longest()` instead

so def `zip(*iterables)`:

```
def zip(*iterables):
    sentinel = object()
    iterators = [iter(it) for it in iterables]
    while iterators:
        result = []
        for it in iterators:
            elem = next(it, sentinel)
            if elem is sentinel:
                return None
            result.append(elem)
        yield tuple(result)
```

`zip()` in conjunction with * operator can be used to unzip list

so `x, y = [1, 2, 3], [4, 5, 6]`

`list(zip(x, y))` → `[(1, 4), (2, 5), (3, 6)]` to shrink

`x1, y1 = zip(*zip(x, y))` → `[1, 2, 3]`

`list(x1)` → `[1, 2, 3]`

`list(y1)` → `[4, 5, 6]`

`[(1, 4), (2, 5), (3, 6)]`

Python - P60

__import__(name)

`__import__(name, globals=None, locals=None, fromlist=(), level=0)`
invoked by `import` statement.
direct use discouraged in favor of `importlib.import_module()`
function `import module <name>`
potentially using given `<globals>` and `<locals>`
to determine how to interpret name in package context
`<fromlist>` give name of object or submodule
should be imported from module given by `<name>`
standard implementation not use `<locals>` argument at all
use `<globals>` only to determine package context of import statement
`<level>` specify whether to use absolute / relative import
default 0 only perform absolute import
positive value for `<level>` indicate number of parent directory

to search relative to directory of module calling `__import__()`
of form "package.module"

top-level packages returned, not module named by `<name>`
when non-empty `<fromlist>` given

module named by `<name>` returned

to import spam result in bytecode resembling code

`spam = __import__('spam', globals(), locals(), [], 0)`

import `spam.ham` result in

`spam = __import__('spam.ham', globals(), locals(), [], 0)`

return toplevel module because object bound to name by import statement

from `spam` import `eggs, sausage` as `saus`, result in

`_temp = __import__('spam.ham', globals(), locals(), ['eggs', 'sausage'], 0)`

`eggs = _temp.eggs`

`saus = _temp.sausage`

`spam.ham` module returned from `__import__()`

name to import retrieved and assigned to respective name

`importlib.import_module()` simply import module by name

Python - P61

cluster analysis, preprocessing: summarization / preprocess for regression, PCA, classification

compression: image processing / vector quantization

search: finding K-nearest neighbour

localizing search to small number of cluster

detection: outlier

quality: high intra-class similarity: cohesive within cluster

low inter-class similarity: distinct between cluster

clustering method: similarity measure / implementation

ultimately ability to discover hidden pattern

issues: scalability to cluster all data instead of sample

ability to deal with different type of attribute

constraint from user / domain

interpretability and usability

ability to deal with noisy data / high dimensionality

approach: partitioning: construct and evaluate (K-means, K-medoids, CLARANS)

hierarchical: decompose set of data object (Diana, Agnes, BIRCH, CAMELEON)

density-based: use connectivity / density function (DBSCAN, OPTICS, DenClue)

grid-based: build multi-level granularity structure (STING, WaveCluster, CLIQUE)

distinction: exclusive/non-exclusive: belong to multiple clusters or not

fuzzy/non-fuzzy: belong to every cluster with weight $0 \leq w_i \leq 1$, $\sum w_i = 1$

probabilistic clustering: partial data or all data (complete)

heterogeneous/homogeneous: cluster of widely different size / shape / density

type: separated: closer to all point in own cluster

centered: closer to center of own cluster

centroid: average of all points

medoid: most "representative" point in cluster

contiguous: closer to ≥ 1 point in own cluster

dense: closer to higher density than lower

Python - P62

partitioning, partitioning database D of n objects into set of k clusters

$$\text{minimize } \sum_{i=1}^k \sum_{x \in C_i} [d(p, c_i)]^2$$

global optimal: exhaustively enumerate all partitions

heuristic method: K-means / K-medoids

K-means: each cluster represented by center of cluster

K-medoid / PAM (partition around medoids): represented by one object in cluster

K-means

Select K points as initial centroids.

REPEAT

form K clusters by assigning all points to ~~one~~ closest centroid

recompute centroid of each cluster

UNTIL centroids don't change in the iteration

initial centroid often chosen randomly

cluster produced vary from one iteration to another

centroid typically mean of point in the cluster

'closeness' measured by Euclidean distance, cosine similarity, correlation

converge for common similarity measure

most of convergence happen in first few iteration, until relatively few point change

evaluation: usual measure is SSE (sum of squared error)

$$SSE = \sum_{i=1}^k \sum_{x \in C_i} [\text{dist}(m_i, x)]^2$$

easy way to reduce SSE is to increase K

good clustering with smaller K have lower SSE than poor clustering with higher K

centroid: multiple runs: probability not ~~one~~ ~~one~~ useful

sample and use hierarchical clustering to determine initial centroid

select > k initial centroids and then select among >k initial centroids

select most widely separated

emptiness: choose point from cluster with highest SSE

typically split cluster and reduce overall SSE for clustering

repeat several times for several empty clusters

Python - P63

K-means

weakness : applicable only when mean defined

need to specify number of cluster K in advance

unable to handle noisy data and outlier

not suitable to discover cluster with non-convex shape

sensitivity: extreme value may substantially distort distribution of data

K-medoids take most centrally located object in cluster

hierarchical clustering, use distance matrix as clustering criteria

not require number of cluster K as input, but need termination condition

terminates until desired number of clusters found

AGNES agglomerative nesting use single-link method and dissimilarity matrix

single-link method: smallest distance between element in one cluster and one in other

merge node with least dissimilarity

go on in non-descending fashion

eventually all nodes belong to same cluster

DIANA divisive analysis, inverse order of AGNES

eventually each node form cluster on own

density clustering, clustering based on density (local cluster criterion)

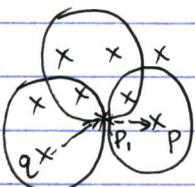
such as density-connected

major feature: discover cluster of arbitrary shape / handle noise / one scan
need density parameter as termination condition

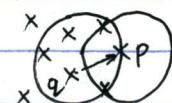
Eps: maximum radius of neighbourhood

MinPts: minimum number of points in Eps-neighbourhood of point

$N_{Eps}(p) = \{q \text{ belongs to } D \mid dist(p, q) \leq Eps\}$



marked draw



directly density-reachable: point p directly density-reachable from point q

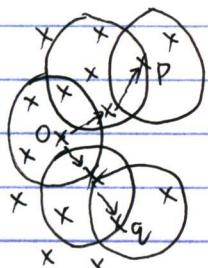
p belongs to $N_{Eps}(q)$ and core point condition: $|N_{Eps}(q)| \geq \text{MinPts}$

density-reachable: point p density-reachable from point q

chain of points $p_1 = q_1, p_2, \dots, p_n = p$, p_i directly density-reachable from p_i

density-connected: point p density-connected to point q

point O: both p and q, density-reachable from O



Python - P64

DBSCAN

density based spatial clustering of application with noise

rely on density-based notion of cluster: maximal set of density-connected point
discover cluster of arbitrary shape in spatial database with noise

OPTICS ordering point to identify clustering structure

produce special order of database w.r.t. density-based clustering structure

cluster-ordering contain info equivalent to density-based clustering

corresponding to broad range of parameter setting

good for both automatic and interactive cluster analysis

including finding intrinsic clustering structure

can be represented graphically or using visualization technique

assessing tendency, assess if non-random structure exist in data

by measuring probability data generated by uniform data distribution

Hopkins Statistic: test spatial randomness by statistic test

extrinsic quality, $Q(c, c_g)$ for clustering C given ground truth C_g ,

cluster homogeneity: the purer, the better

cluster completeness: should assign object belong to same category

rag bag: putting heterogeneous object into pure cluster more penalized than into

small cluster preservation: splitting small category into pieces is more harmful

scatter coefficient = avg intra cluster distance / avg inter cluster distance

intrinsic quality, cluster cohesion: measure how closely object related in cluster

measured by within cluster sum of square

cluster separation: measure how distinct / separated one cluster from other

measured by between cluster sum of square

silhouette coefficient: combine both cohesion and separation, but for individual point

validity aspect: determining clustering tendency of set of data

comparing result of cluster analysis to externally known result

evaluating how well result of cluster analysis fit data without reference to external

comparing result of different set of cluster analysis to determine the best