

Vector Embedding of Code

Code Repository Mining (CRM2020)

Lando Löper
Software Architecture Group
Hasso Plattner Institute Potsdam
02.06.2020

Agenda

1. Motivation
2. Related Work
3. Design Space
4. Data
5. Outlook
6. Discussion

Motivation

What does it do?

```
function(list, target)
```

```
  i := 0
```

```
  for x in list
```

```
    if x = target
```

```
      return i
```

```
    i++
```

```
  return -1
```



- contains
- count
- indexOf
- reverse

Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Motivation

What does it do?

```
function(list, target)
```

```
  i := 0
```

```
  for x in list
```

```
    if x = target
```

```
      return i
```

```
    i++
```

```
  return -1
```



- contains
- count
- indexOf
- reverse

- Reading unlabelled code takes time

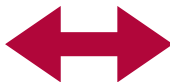
Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Motivation

How similar are they?

```
function(list, target)
  i := 0
  for x in list
    if x = target
      return i
  i++
return -1
```



```
function(list, target)
  i := 0
  for x in list
    if x = target
      i++
  return i
```

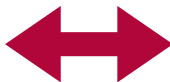
Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Motivation

How similar are they?

```
function(list, target)
  i := 0
  for x in list
    if x = target
      return i
  i++
return -1
```



```
function(list, target)
  i := 0
  for x in list
    if x = target
      i++
  return i
```

- Similar syntax does not mean similar behaviour

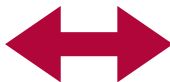
**Vector Embedding
of Code**

Lando Löper, Code
Repository Mining,
02.06.2020

Motivation

How similar are they?

```
function(list, target)
  i := 0
  for x in list
    if x = target
      return i
  i++
return -1
```



```
function(elements, x)
  i := 0
  n := size of elements
  while i < n
    if elements[i] = x
      return i
  i++
return -1
```

Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Motivation

How similar are they?

```
function(list, target)
```

```
  i := 0
```

```
  for x in list
```

```
    if x = target
```

```
      return i
```

```
    i++
```

```
  return -1
```



```
function(elements, x)
```

```
  i := 0
```

```
  n := size of elements
```

```
  while i < n
```

```
    if elements[i] = x
```

```
      return i
```

```
    i++
```

```
  return -1
```

- Different syntax does not mean different behaviour

**Vector Embedding
of Code**

Lando Löper, Code
Repository Mining,
02.06.2020

Motivation

Direct mapping

```
function(list, target)
```

```
  i := 0
```

```
  for x in list
```

```
    if x = target
```

```
      i++
```

```
  return i
```



count

```
function(elements, x)
```

```
  i := 0
```

```
  n := size of elements
```

```
  while i < n
```

```
    if elements[i] = x
```

```
      return i
```

```
    i++
```

```
  return -1
```



indexOf

```
function(list, target)
```

```
  i := 0
```

```
  for x in list
```

```
    if x = target
```

```
      return i
```

```
    i++
```

```
  return -1
```



indexOf

**Vector Embedding
of Code**

Lando Löper, Code
Repository Mining,
02.06.2020

Motivation

What about a different representation?

```
function(list, target)
  i := 0
  for x in list
    if x = target
      i++
  return i
```



```
function(elements, x)
  i := 0
  n := size of elements
  while i < n
    if elements[i] = x
      return i
    i++
  return -1
```



```
function(list, target)
  i := 0
  for x in list
    if x = target
      return i
    i++
  return -1
```



Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Motivation

Indirect mapping

```
function(list, target)
  i := 0
  for x in list
    if x = target
      i++
  return i
```



count

```
function(elements, x)
  i := 0
  n := size of elements
  while i < n
    if elements[i] = x
      return i
  i++
  return -1
```



indexOf

```
function(list, target)
  i := 0
  for x in list
    if x = target
      return i
  i++
  return -1
```



indexOf

**Vector Embedding
of Code**

Lando Löper, Code
Repository Mining,
02.06.2020

Motivation

Embedding

Definition

An embedding is a mapping of a discrete — categorical — variable to a vector of continuous numbers

Benefits

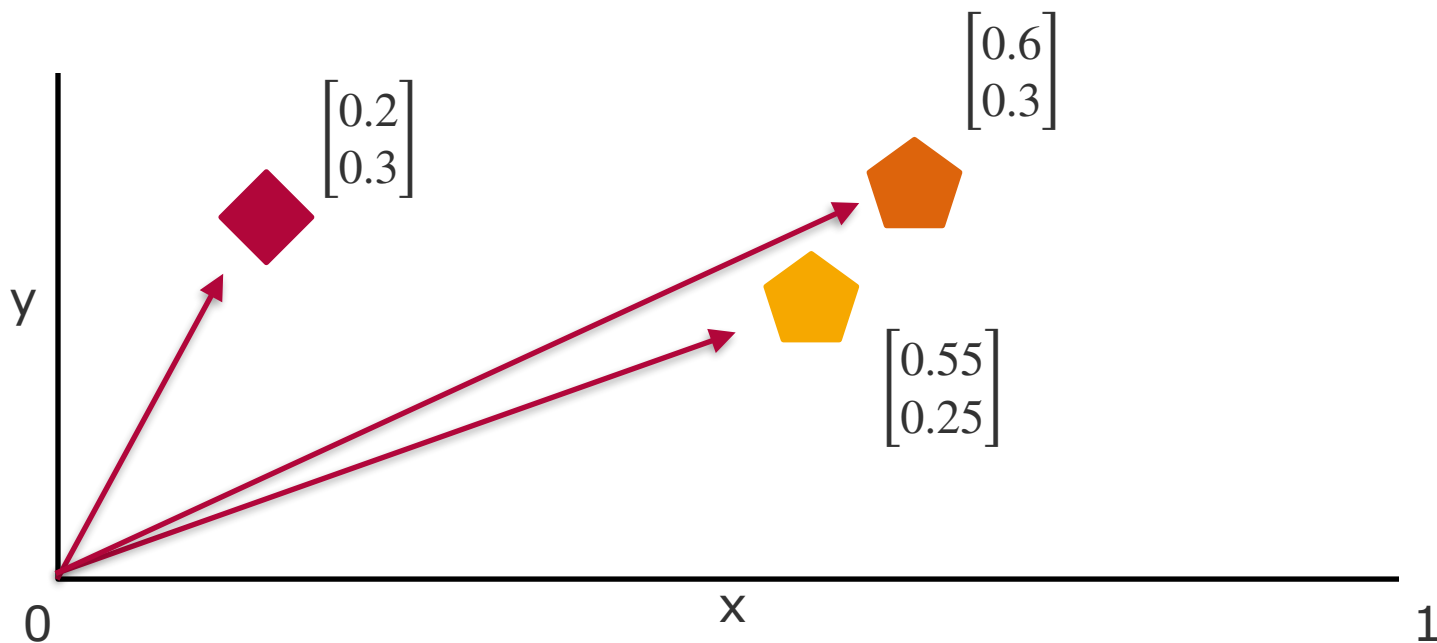
- Finding of nearest neighbours
- Visualisation of relationships

Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Motivation

Code embedding



Vector Embedding of Code

Lando Löper, Code Repository Mining, 02.06.2020

Related Work

Word2Vec

Idea

Words that appear in the same context are similar

Example

- "The grey cat is sitting on the porch."
- "The white dog is sitting in the garage."

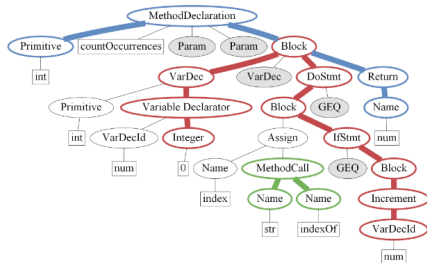
Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Idea

Functions with a similar AST are similar

Example



Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Related Work

DYPRO

Idea

Functions with a similar runtime behaviour are similar

Example

- Sorting Algorithm

Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Related Work

DYPRO

Idea

Functions with a similar runtime behaviour are similar

Example

- Sorting Algorithm

Bubble	Insertion
[5,5,1,4,3]	[5,5,1,4,3]
[5,8,1,4,3]	[5,8,1,4,3]
[5,1,1,4,3]	[5,1,1,4,3]
[5,1,8,4,3]	[5,1,8,4,3]
[1,1,8,4,3]	[5,1,4,4,3]
[1,5,8,4,3]	[5,1,4,8,3]
[1,5,4,4,3]	[5,1,4,3,3]
[1,5,4,8,3]	[5,1,4,3,8]
[1,4,4,8,3]	[1,1,4,3,8]
[1,4,5,8,3]	[1,5,4,3,8]
[1,4,5,3,3]	[1,4,4,3,8]
[1,4,5,3,8]	[1,4,5,3,8]
[1,4,3,3,8]	[1,4,3,3,8]
[1,4,3,5,8]	[1,4,3,5,8]
[1,3,3,5,8]	[1,3,3,5,8]
[1,3,4,5,8]	[1,3,4,5,8]

Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Design Space

	Idea	Assumption	Method
Code as sequence of tokens	<ul style="list-style-type: none"> Learn embeddings from token context 	<ul style="list-style-type: none"> Code can be tokenised 	<ul style="list-style-type: none"> Word2Vec
Code as abstract syntax tree	<ul style="list-style-type: none"> Learn embeddings from syntax 	<ul style="list-style-type: none"> Code can be parsed 	<ul style="list-style-type: none"> Code2Seq
Code as execution traces	<ul style="list-style-type: none"> Learn embeddings from runtime behaviour 	<ul style="list-style-type: none"> Code can be executed Program state can be logged 	<ul style="list-style-type: none"> DYPRO

Vector Embedding of Code

Lando Löper, Code Repository Mining, 02.06.2020

Data

Source Code

- Java (<https://s3.amazonaws.com/code2seq/datasets/java-large.tar.gz>)
- Python (<https://eth-sri.github.io/py150>)
- COSET (<https://arxiv.org/abs/1905.11445>)

Metric

- Use-Case (Classification)
 - Accuracy
 - Precision / Recall / F1

Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Outlook

Problem

Predicting function names from their body.

Current State

- Python dataset
- Training code embeddings with Word2Vec

Next Steps

- Training code embeddings including AST information

Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Discussion

Questions?

- Can you think of other applications for code embeddings?
- Can you think of methods to evaluate code embeddings?

Appendix

- Code Embedding Examples
- Word Embedding Examples
- Visualisations of current state

Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

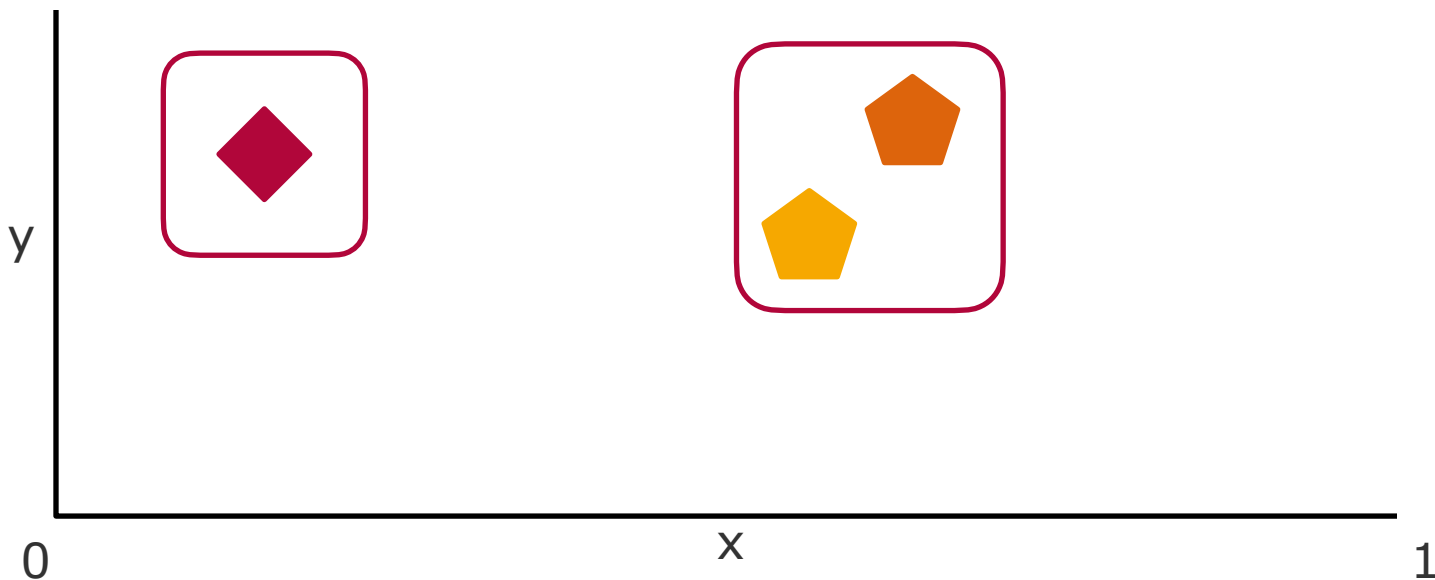


Thank you
for your attention!

Lando Löper
Software Architecture Group
Hasso Plattner Institute Potsdam
02.06.2020

Appendix A

Nearest Neighbour

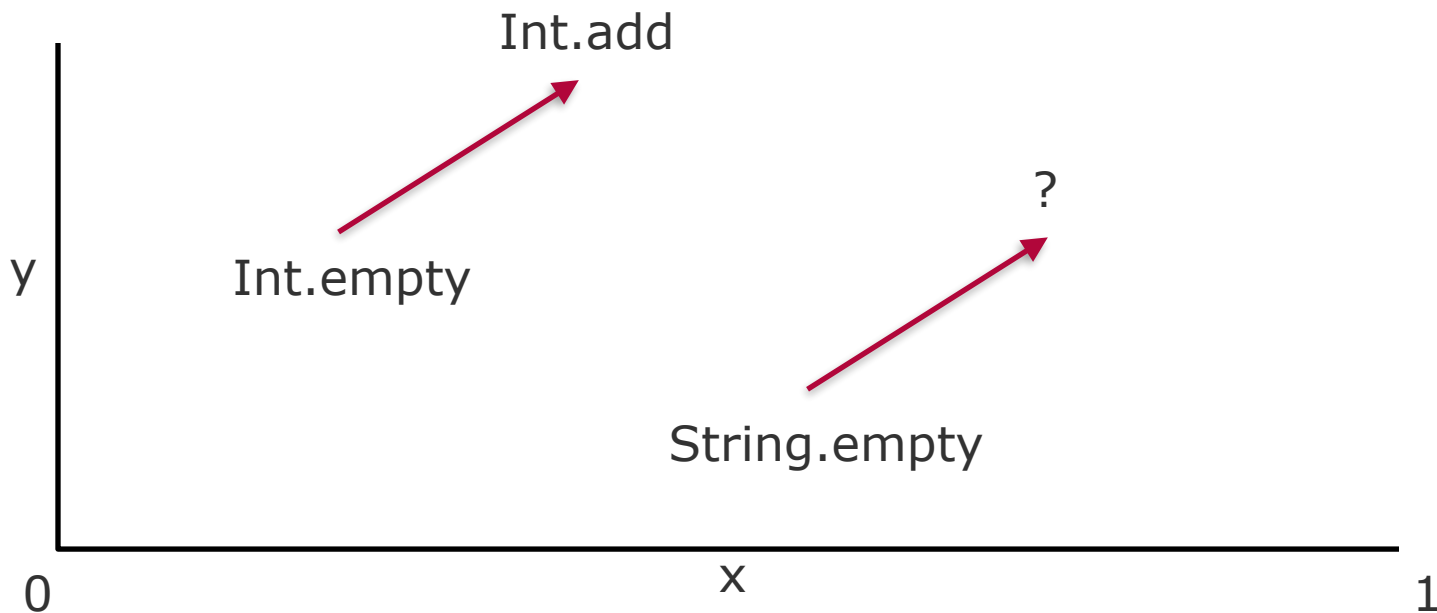


Vector Embedding of Code

Lando Löper, Code Repository Mining, 02.06.2020

Appendix A

Relationships

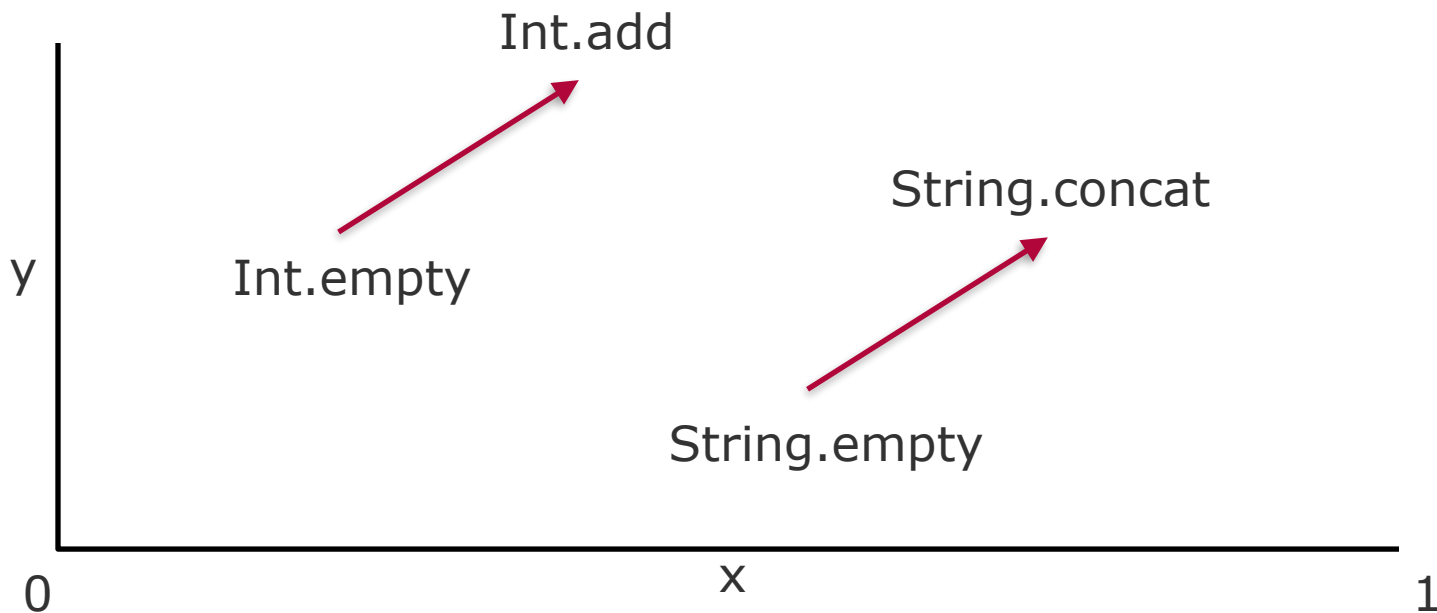


**Vector Embedding
of Code**

Lando Löper, Code
Repository Mining,
02.06.2020

Appendix A

Relationships

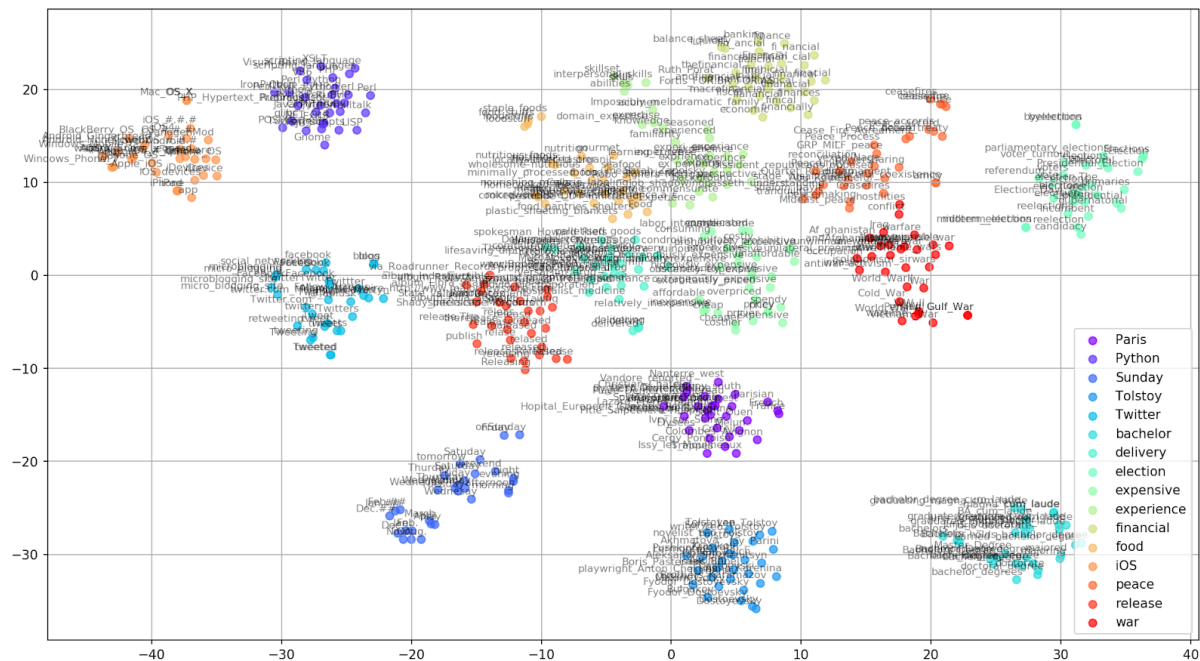


Vector Embedding of Code

Lando Löper, Code
Repository Mining,
02.06.2020

Appendix B

Nearest Neighbour

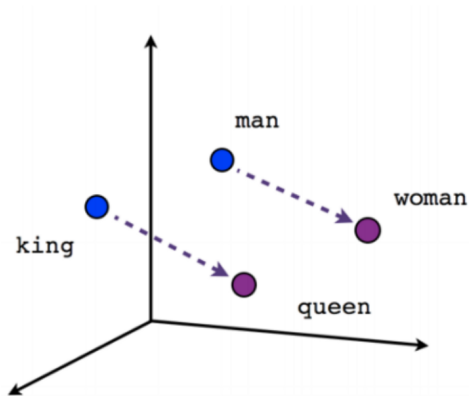


**Vector Embedding
of Code**

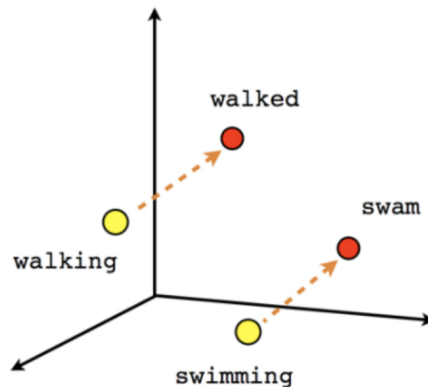
Lando Löper, Code
Repository Mining,
02.06.2020

Appendix B

Relationships



Male-Female



Verb tense

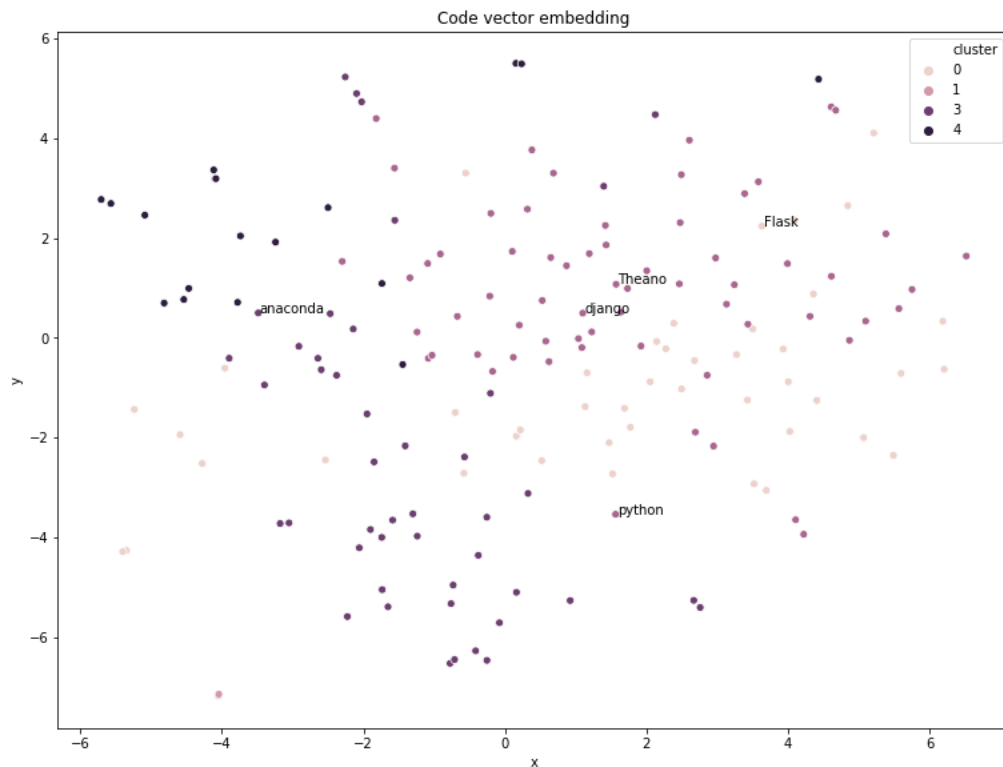
Vector Embedding of Code

Lando Löper, Code Repository Mining, 02.06.2020



28

Appendix C



Vector Embedding of Code

Lando Löper, Code Repository Mining, 02.06.2020