

Research Questions:

Can we learn good vector representations (vector *embeddings*) of student code?
 Can we use these representations to reason about student code?

Background:

Many automated tools have been developed to understand and provide feedback to computer programming students [4].
 Neural networks have been used to learn embeddings of words (e.g. word2vec [1] and GloVe [2]) and images (e.g. image auto encoders [3]).

Can we do the same with student code?

Data:

CS1 student code submission to a Python programming problem about modifying lists. A total of 8,476 syntactically correct submissions, split into 5,086 for training, 1,695 validation and 1,695 test.
 Automatically generated 100 unit tests.

Future Work:

Can we use code embeddings to predict student understanding of concepts?

Can we use code embeddings to help give students feedback?

Key References:

- [1] Mikolov et al. (2013) *Distributed Representations of Words and Phrases and their Compositionally*
- [2] Pennington et al. (2014) *GloVe: Global Vectors for Word Representation*
- [3] Cheng et al. (2018) *Deep Convolutional AutoEncoder-based Lossy Image Compression*
- [4] Keuning et al. (2019) *A Systematic Literature Review of Automated Feedback Generation for Programming Exercises*

Acknowledgements: We thank Prof. Andrew Petersen for the data and the helpful feedback. We thank Haotian Yang for some auxiliary model results.

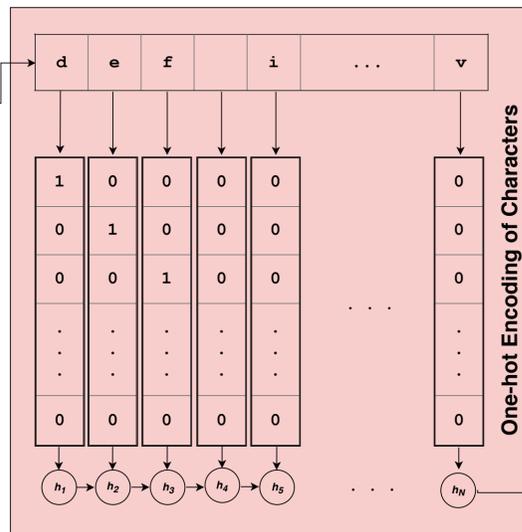
MODEL:

Student Code Submission

```
def insert(lst: List[int], v: int) -> None:
    i = len(lst) - 1
    while i != 0 and lst[i] <= v:
        i = i - 1
    if i == 0:
        lst.append(lst[-1])
        j = len(lst) - 2
        while j != 0:
            lst[j] = lst[j - 1]
            j = j - 1
        lst[0] = v
    elif i > 0:
        lst.append(lst[-1])
        j = len(lst) - 2
        while j > i:
            lst[j + 1] = lst[j]
        lst[i] = v
```

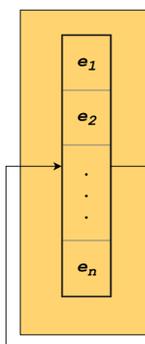
We treat each syntactically-correct student code submission as a sequence of characters.

Recurrent Neural Network (RNN) Encoder

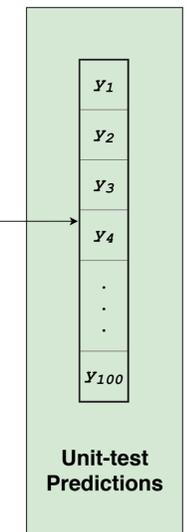


A character-level RNN that uses bi-directional Gated Recurrent Units (GRU). It takes the sequence of characters in a code submission as input, and outputs its vector embedding.

Embedding



Neural Network Decoder



A multi-layer perceptron (MLP) with two fully connected layers. It takes the vector embedding as input, and predicts the pass/fail of each of the 100 unit tests.

Training:

The entire model (encoder + decoder) is trained together to predict the unit test output (pass/fail) of the code submissions in our training set.

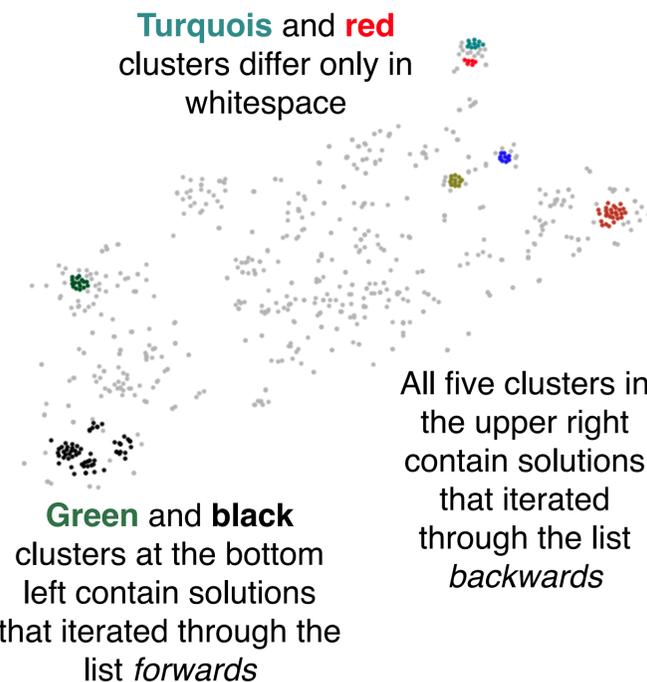
Although the encoder is the only part of the model that we use in our analysis, the decoder helps the encoder learn good embeddings that captures the unit test pass/fail information.

Final Training Accuracy	90%
Final Validation Accuracy	84%
Final Test Accuracy	85%

ANALYSIS:

Code Clustering:

Distances in the embedding space appears to be meaningful. We find seven tight clusters of *correct* solutions (dimensionality reduction via tSNE):



Auxillary (New) Tasks:

Can we use the embeddings to make predictions about other features of the code?

Task 1 Error Identification: Identify whether a code submission throws a Python exception.

Task 2 Error Classification: Identify what type of error a code submission throws. We grouped the errors in to four classes: No error, IndexError, RecursionError or Timeout, and Other.

Input	Test Accuracy	
	Task1	Task2
Unit test	79%	62%
Code	81%	75%
Embedding	91%	87%

Analogies:

Embeddings of *words* [1] often have structures with analogies like "king - man + woman = queen". Is there evidence of similar analogies in our code embedding?

```
lst.insert(0, v)
for i in range(len(lst)-1):
    if lst[i] > lst[i+1]:
        lst[i], lst[i+1] = lst[i+1], lst[i]

lst.insert(0, v)
for i in range(len(lst)):
    if lst[i] > lst[i+1]:
        lst[i], lst[i+1] = lst[i+1], lst[i]
```

+

```
i = len(lst)
while lst[i] > v:
    if v < lst[i]:
        lst.insert(i - 1, v)
    if v > lst[i]:
        lst.insert(0, v)
```

≈

```
i = len(lst)-1
while lst[i] > v:
    if v < lst[i]:
        lst.insert(i - 1, v)
    if v > lst[i]:
        lst.insert(0, v)
```

In this example, we:

1. Embed the first three pieces of code
2. Compute the addition and subtraction of the vectors
3. Find the code submission closest to resulting vector

This result suggests that direction vectors in the embedding space are meaningful, and can represent complex programming concepts.