



Distributed representations of code: a tutorial

¹E Krishna Assistant Professor, krishna.cseit@gmail.com,

²Bandam Naresh Assistant Professor, nareshbandam4@gmail.com,

³Banothu Usha Assistant Professor, banothuusha@gmail.com,

⁴Dr. J Rajaram Assistant Professor, drrajram81@gmail.com,

Department of CSE Engineering,
Nagole , Institute of Engineering and Technology collage in Hyderabad.

INTRODUCTION

Using distributed representations of words, sentences, paragraphs, and documents (such as doc2vec) has been critical in unlocking the potential of neural networks for natural language processing (NLP) tasks. An object's low-dimensional vector representation, known as an embedding, may be learned using techniques for learning distributed representation. Because the meaning of an element is spread out over several vector components in these vectors, things with semantically comparable meanings are translated to nearby vectors.

```
String[] f(final String[] array) {
    final String[] newArray = new String[array.length];
    for (int index = 0; index < array.length; index++) {
        newArray[array.length - index - 1] = array[index];
    }
    return newArray;
}
```

Predictions



Fig. 1. A code snippet and its predicted labels as computed by our model.

Purpose: The goal of this study is to learn code embeddings, continuous vectors for expressing snippets of code. By learning code embeddings, our long-term objective is to allow the application of neural approaches to a broad variety of programming-language problems. In this study, we employ the motivating aim of semantic labelling of

code fragments. Motivating task: semantic labelling of code fragments. Consider the procedure in Figure 1. The procedure comprises just low-level assignments to arrays, yet a person reading the code may (correctly) describe it as doing the opposite action. Our objective is to forecast such labels automatically. The right-hand side of Figure 1 displays the labels predicted automatically using our technique. The most probable guess (77.34 percent) is reverseArray. Section 6 includes other instances. This challenge is complex since it needs learning a relationship between the complete content of a technique and a semantic label. That is, it involves aggregating maybe hundreds of phrases and assertions from the method body into a single, descriptive name. Our approach. We provide an unique approach for predicting programme attributes using neural networks. Our fundamental contribution is a neural network that learns code embeddings \mathbb{D} continuous distributed vector representations for code. The code embeddings provide us a natural and efficient way to represent relationships between code snippets and labels. Because of the organised nature of source code, our neural network design can learn to combine several syntactic pathways into a single vector. As in NLP, word embeddings are important to the use of deep learning for NLP tasks, and this capacity is crucial to the application of deep learning in programming languages. It is necessary to provide the model with an appropriate tag, caption or name for the code snippet. Using this label, we specify the semantic attribute we want the network to represent, for example, a tag applied to the snippet, or the name of the method, class or project from whence the sample was extracted. The code fragment is C, and the label or tag is L. It is our primary assumption that the distribution of labels may be derived from C's syntactic routes. That is why the label distribution $P(L|C)$ is part of our model's learning process. For the goal of guessing the name of a technique based on its body, we show the usefulness of our methodology. This is an essential issue since clear method names



make it simpler to write and maintain software. Choosing a method's name carefully allows you to communicate the method's primary goal in a concise and memorable way. Ideally, "you don't need to look at the body if you have a suitable method name." A few years later [Fowler and Beck 1999]. Choosing the right names for public API methods is extremely important, since bad method names may condemn a project to obscurity. It has been reported [Allamanis et al, 2015a; Hùst, 2009]; [Allamanis et al, 2015b; Hùst, 2009]. Semantic equivalence between names is captured. Learning code vectors requires the acquisition of a vocabulary of vectors for the labels. To forecast method names using our approach, method-name vectors reveal unexpectedly close semantic and parallels. $\text{vector}(\text{equals}) + \text{vector}(\text{toLowerCase})$ produces a vector that is very similar to $\text{vector}(\text{equalsIgnoreCase})$.

Table 1. Semantic similarities between method names.

1.00	0.9999999999999999	0.9999999999999999	0.9999999999999999
0.9999999999999999	1.00	0.9999999999999999	0.9999999999999999
0.9999999999999999	0.9999999999999999	1.00	0.9999999999999999
0.9999999999999999	0.9999999999999999	0.9999999999999999	1.00

A well-known NLP analogy is that the model learns analogies that are relevant to source code, such as: receive is to transmit as download is to upload [Mikolov and colleagues 2013c]. There are more instances in Table 1, and a more in-depth explanation may be found in Section 6.4.

Applications

It's possible to use machine learning methods by embedding code snippets as vectors, since machine learning techniques often use vectors as inputs. The following direct applications are examined in this paper: Automatic code review suggests better method names when the name supplied by the developer does not correspond to method functionality. To avoid naming issues, increase readability and maintenance of code, and promote the usage of public APIs, better method names are needed. Allamanis et al. 2015a; Fowler and Beck 1999; Hùst and stvold 2009] have already demonstrated that this application is of major value. Semantic similarities allow search in "the issue domain" rather than "the solution domain" for retrieval and API discovery. It is common for developers to search for a serialise function when the corresponding method of the class is called toJson

since serialisation is done through json. Finding toJson is easy with the use of a computer algorithm that scans all of the techniques for the one that is most close to the required name (Table 1). Without our method, it's impossible to detect semantic similarities like this. Our vectors may also be used by an automated tool to tell if someone is using equals right after toLowerCase and propose that they switch to equalsIgnoreCase in its place (Table 6). To help with activities like code retrieval, captioning, categorization, and tagging, or as a measure for comparing the similarity of code snippets in order to aid in ranking and clone detection, we generate vectors of code. As a result of our method's unique ability to generate a unique vector for each and every line in a snippet of code, comparable lines of code may be given similar vectors. Machine learning algorithms may be used in a wide range of contexts thanks to this capacity. For our assessment benchmark, we chose the tough job of method name prediction, which has previously shown disappointing results [Allamanis et al. 2015a, 2016; Alon et al. 2018]. Predicting whether a programme does or does not conduct I/O, discovering its dependencies, and determining whether it is suspected of being a malicious software are all activities that are necessary for success in this assignment. We demonstrate that our method significantly outperforms the findings of prior studies on this difficult benchmark.

Challenges: Representation and Attention

Giving a name to a method (or assigning a semantic label to it) is an example of an issue that calls for a concise way to describe the meaning of the code in question. When it comes to code snippets, the challenge is how to encode them in such a manner that they may be reused across different applications and used to anticipate attributes such as their labels. Two issues arise as a result of this: Incorporating an example from one software into another in order to make it easier to understand. It's important to discover which components of the representation are significant to predicting the desired feature, and the order in which they are important. Representation. Text is often treated as a linear series of tokens by NLP techniques. There are a number of current techniques that treat source code as a token stream [Allamanis et al. 2014, 2016; Allamanis and Sutton 2013; Hindle et al. 2012; Movshovitz-Attias and Cohen 2013; White et al. 2015]. In contrast to this, programming languages may considerably benefit from representations that use the organised character



of their grammar [Alon et al. 2018; Bielik et al. 2016; Raychev et al. 2015]. A trade-off exists between the amount of analysis necessary to extract the representation and the amount of learning effort required to use it. ' In many cases, learning about a programme just reading the program's surface language requires a substantial amount of time and effort. So much data is needed for this learning endeavour because the learning model must re-learn programming syntax and semantics from the input. However, if the representation is extracted by a thorough study of the programme code, the learnt model may become language-specific (and even task-specific). For our representation, we follow prior research [Alon et al. 2018; Raychev et al. 2015] by using paths in the program's abstract syntax tree (AST). The regularities that reflect frequent coding patterns may be captured by describing a sample of code using its syntactic pathways. For vast volumes of code and a broad variety of issues, we show that this approach greatly reduces the learning effort (relative to learning from programme text). A code snippet is represented by a collection of its extracted pathways, which we call a bag. The difficulties Attention. To put it another way, the issue is that you need to figure out how to connect a bag to a computer. with a label and a set of path contexts. Even comparable procedures will not have the exact same bag of path-contexts if they represent each one monolithically. As a result, we require a compositional technique that can aggregate a bag of route contexts such that bags that generate the same label are mapped to nearby vectors. Compositional mechanisms of this kind might generalise and describe previously unknown bags by exploiting previously learned path-contexts and their constituent components (paths, values, etc.). A new neural attention network design is used to tackle this problem. Attention models have lately attracted a lot of attention, particularly in the fields of neural machine translation (NMT), reading comprehension, voice recognition, and computer vision. Neural mechanisms learn how much attention to pay to each item in a bag of path-contexts (the attention"). Each path-context is aggregated into a single vector that contains all the information about the whole code snippet. This may be shown in Section 6.4, where we explain how the weights assigned by our attention mechanism can be represented to comprehend the relative relevance of each path-context in a prediction. Both the atomic representations of pathways and the capacity to assemble many contexts into a single code vector are optimised concurrently while learning both the attention mechanism and the

embeddings. Disciplined and unguided focus By Xu et al. [2015], soft and hard attention were used to describe the work of creating picture captions. While hard attention refers to selecting one path-context to concentrate on, soft attention indicates that all path-contexts are weighted equally in our context. The enhanced outcomes may be attributed to the usage of soft focus on syntactic routes. It is shown that our approach is more efficient for modelling code when compared to an analogous model based on hard focus.

Existing Techniques

Recent years have seen a lot of attention and development in the challenge of predicting programme attributes by learning from huge code [Allamanis et al. 2014; Allamanis and Sutton 2013; Bielik et al. 2016; Hindle et al. 2012; Raychev et al. 2016a]. For a variety of applications, it is essential to be able to predict semantic properties of a programme without running it and with minimal or no semantic analysis: predicting names for programme entities [Allamanis et al. 2015a; Alon et al. 2018; Raychev et al. 2015], code completion [Mishne et al. 2012; Raychev et al. 2014], code summarization [Allamanis et al. 2016], code generation [Amod (see [Allamanis et al. 2017; Vechev and Yahav 2016] for a survey).

Contributions

It uses a path-based attention model for vector learning in arbitrary-sized chunks of code. Using this technique, we may feed a programme, which is a discrete object, into a deep learning pipeline for a variety of tasks.

As a benchmark for our methodology, we conduct a quantitative assessment for predicting the names of cross-project methods, trained on more than 12 million methods of real-world data, and compared to earlier research. • This new technique outperforms prior efforts that employed LSTMs, CNNs, and CRF-based networks. For example, a qualitative assessment based on how much attention the model has learnt to pay to the diverse path circumstances while generating predictions Method name embeddings, which commonly assign similar names to comparable vectors, and even make it easy to calculate analogies using basic vector arithmetic, are included. A comparison of our model to prior non-neural efforts, such as Alon et al. [2018] and Raychev et al. [2015],

to highlight the major benefits in terms of generalisation and spatial complexity of our model.

OVERVIEW

It's in this part that we show how our model is able to identify tiny changes between comparable pieces of code. A prediction may be made about each snippet even if it has not been seen in its whole in the training data because of the vectors. In order to generate a single code snippet, we employ an attention method to learn the weighted average of the route vectors and extract syntactic pathways from the snippet. These path vectors are then represented as a bag of distributed vector representations. As a final benefit, this code vector may be used to make educated guesses about the snippet's name.



Figure 2 shows three strategies that, despite their similar syntactic form, may be clearly identified by our model: Our approach accurately predicts meaningful names by capturing the minor distinctions between them. Each technique depicts the model's top four most important routes. Colored pathways' widths are proportionate to the amount of attention they received.

Motivating Example

We show how to learn code vectors for method bodies and predict the method name given the body using our technique since method names are often descriptive and accurate labels for code snippets. The same method may be used to apply to any code fragment that contains a label. The three Java methods shown in Figure 2 are good examples. In terms of syntax, these methods all contain a single argument named target, (ii) iterate through a field called elements, and (iii) have an if condition within the loop body. They all follow a similar pattern: As can be seen in Figures 2a and 2b, the former returns true if elements include target and the latter returns false if it does not; Figure 2b returns the element from elements whose hashCode matches target, while Figure 2c provides the index of target inside elements. Despite their overlapping qualities, our model is able to accurately predict the descriptive method names: contains, get, and indexOf, all of which have a distinct meaning. Extraction of a route. An AST is constructed for each query method in the training corpus. Syntactic analysis is then carried out

by traversing the ASTPaths between the AST leaves are extracted. Paths are shown as a series of AST nodes connected by arrows pointing up or down in the tree, respectively. We refer to this tuple as a path-context since it contains the values of the AST leaves to which it is connected. Section 3 clearly defines these concepts. Using the AST of the technique in Figure 2a, Figure 3 shows the top four path contexts that were given the greatest attention by the model during this prediction, with the width of each route corresponding to the attention it received from the AST. Contexts are represented in a distributed manner. It is mapped to the real-valued vector representation, or embedding, of each of the route and leaf values of a path-context. For each path-context, a single vector is concatenated from the three contexts. It is possible to learn the embedding values as well as other network parameters during training. A network of paths and attention. An entire method's path-context embeddings are combined into one vector by the Path-Attention network. It's the attention mechanism that learns to score each path-context, such that the more the attention, the better the score.

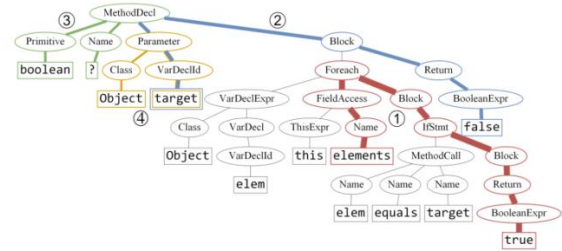


Figure 3 shows the model's top-four attended pathways from Figure 2a on the AST of the same sample. In each hue, the amount of attention it receives is reflected in the path's width (red 1: 0.23, blue 2: 0.14, green 3: 0.09, orange 4: 0.07).

The attention ratings are used to combine these many embeddings into a single code vector. The network then estimates the likelihood of each target method name given the code vector. Section 4 explains the network's structure. Interpretation of the path of travel. The attention scores that each path-context received from the network may be seen, notwithstanding the difficulty of interpreting particular values of vector components in neural networks. Figures 2 and 3 show snippets of code that indicate the top four route contexts in each case, as determined by the model. Depending on how much attention these path-contexts get, the pathways' widths vary. As a result of training on millions of samples, the model has learnt how much weight to



assign each feasible route. For example, in Figure 3, the red 1 path-context, which runs from the field items to the return value true, was given the most attention. In contrast, less attention was paid to the blue 2 path-context, which extends from the argument target to the return value false. Think about the red 1 path-context shown in Figure 2a and Figure 3a. It's explained in Section 3 as: (elements, Name'FieldAccess'Foreach'Block'IfSmt'Block'Return' BooleanExpress>true') It is clear that this single route contains the method's core functionality, since it iterates over a field named elements and verifies an if condition for each value; if the condition is true, the method returns true. It is easy to see why this route was given the most attention by the model since we employ "soft attention," which takes into consideration other pathways such as those that explain the "if condition" itself." In addition, the model's top-five recommendations for each approach are shown in Figure 2. It is clear from the samples that the top proposals are highly similar to each other, and all of them are descriptive of the process. For example, a method named matches is likely to include an if condition within a for loop and to return true if a condition is true, as shown in Figure 2a's top-5 choices (contains and containsExact are two of the most correct ones). Figure 2a's orange 4 path-context, which runs from Object to target, received less attention than other path-contexts in the same procedure but more attention than the orange 4 path-context in Figure 2c. Although attention is not constant, it is provided to the various path-contexts in the code. Comparative study of n-gram structures. Figure 2a displays the four path-contexts that received the greatest attention during the prediction of the method name. The orange 4 path-context, for example, connects the tokens "object" and "target" in a chain. This may give the idea that a bag-of-bigrams representation of this approach may be as expressive as a syntactic route representation. While the AST node of type Parameter differentiates it from, for example, a local variable declaration of the same name and type, the orange 4 path does not. An object model uses the same representation regardless of whether an object is sent as a method argument or stored locally. Using a syntactic representation of a code sample, the model can discriminate between two snippets of code that other models cannot. The model may take advantage of small changes across snippets to provide a more precise prediction by combining all contexts with attention. The essentials. Highlights of our methodology may be seen in the examples provided. • A collection of path-contexts may be used to represent a code snippet. • Making an

accurate forecast requires more than just one context. To create a forecast, an attention-based neural network takes into account different route contexts. While code samples with identical syntactic structure contain many of the same n-grams, our model can quickly discern subtle variations across code snippets. This model may be used to forecast method names across large datasets and projects. In spite of its neural network foundation, our model is human-interpretable and generates intriguing insights.

CONCLUSION

New attention-based neural networks for encoding arbitrary-sized code chunks using fixed-length continuous vectors were described. Using a soft-attention method, the snippet's Abstract Syntax Tree (AST) vector representations are aggregated to form a single vector representation. Predicting method names using a model trained on over 12 million methods was one way we showed off our technique. With our model, we are able to forecast file names across many projects, which is a huge improvement over prior methods. We hypothesise that our model's simplicity and dispersed nature allow it to be generalised. The prediction findings are understandable and engaging because of the attention mechanism. As a foundation for a broad variety of programming language processing activities, we think the attention-based approach that leverages a structural representation of code may be used. All of our code and our trained model may be found on <https://github.com/tech-srl/code2vec> for this reason.

REFERENCES

- (1) MiltiadisAllamanis, Earl T. Barr, Christian Bird, and Charles Sutton. 2014. Learning Natural Coding Conventions. In Proceedings of the 22Nd ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE 2014). ACM, New York, NY, USA, 2816293.
- (2) <https://doi.org/10.1145/2635868.2635883> MiltiadisAllamanis, Earl T. Barr, Christian Bird, and Charles Sutton. 2015a. Suggesting Accurate Method and Class Names. In Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering (ESEC/FSE 2015).
- (3) ACM, New York, NY, USA, 38649. <https://doi.org/10.1145/2786805.2786849> MiltiadisAllamanis, Earl T. Barr, PremkumarDevanbu, and Charles Sutton. 2017. A Survey of Machine Learning for Big Code and Naturalness. arXiv preprint arXiv:1709.06182 (2017).
- (4) MiltiadisAllamanis, Marc Brockschmidt, and Mahmoud Khademi. 2018. Learning to Represent Programs with Graphs. In ICLR. MiltiadisAllamanis, Hao Peng, and Charles A.



- Sutton. 2016. A Convolutional Attention Network for Extreme Summarization of Source Code.
- (5) In Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016. 2091&2100. <http://jmlr.org/proceedings/papers/v48/allamanis16.html> MiltiadisAllamanis and Charles Sutton.
- (6) 2013. Mining Source Code Repositories at Massive Scale Using Language Modeling. In Proceedings of the 10th Working Conference on Mining Software Repositories (MSR '13). IEEE Press, Piscataway, NJ, USA, 207&216.
- (7) <http://dl.acm.org/citation.cfm?id=2487085.2487127> MiltiadisAllamanis and Charles Sutton. 2014. Mining Idioms from Source Code. In Proceedings of the 22Nd ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE 2014). ACM, New York, NY, USA, 472&483. <https://doi.org/10.1145/2635868.2635901> MiltiadisAllamanis, Daniel Tarlow, Andrew D. Gordon, and Yi Wei. 2015b. Bimodal Modelling of Source Code and Natural Language.
- (8) In Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37 (ICML'15). JMLR.org, 2123&2132. <http://dl.acm.org/citation.cfm?id=3045118.3045344> Uri Alon, MeitalZilberstein, Omer Levy, and Eran Yahav. 2018. A General Path-based Representation for Predicting Program Properties.
- (9) In Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2018). ACM, New York, NY, USA, 404&419.
- (10) Matthew Amodio, Swarat Chaudhuri, and Thomas W. Reps. 2017. Neural Attribute Machines for Program Generation. CoRR abs/1705.09231 (2017). arXiv:1705.09231 <http://arxiv.org/abs/1705.09231> Thierry Artieres et al.
- (11) 2010. Neural conditional random fields. In Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics. 177&184. Jimmy Ba, Volodymyr Mnih, and KorayKavukcuoglu. 2014. Multiple object recognition with visual attention.
- (12) arXiv preprint arXiv:1412.7755 (2014). DzmitryBahdanau, Kyunghyun Cho, and YoshuaBengio. 2014. Neural Machine Translation by Jointly Learning to Align and Translate. CoRR abs/1409.0473 (2014). <http://arxiv.org/abs/1409.0473> DzmitryBahdanau, Jan Chorowski, Dmitriy Serdyuk, Philemon Brakel, and YoshuaBengio. 2016. End-to-end attentionbased large vocabulary speech recognition.
- (13) In Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE International Conference on. IEEE, 4945&4949. YoshuaBengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. 2003. A Neural Probabilistic Language Model. J. Mach. Learn. Res. 3 (March 2003), 1137&1155.
- (14) <http://dl.acm.org/citation.cfm?id=944919.944966> PavolBielik, VeselinRaychev, and Martin T. Vechev. 2016. PHOG: Probabilistic Model for Code. In Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016. 2933&2942. <http://jmlr.org/proceedings/papers/v48/bielik16.html> Chris Callison-Burch, Miles Osborne, and Philipp Koehn. 2006. Re-evaluation the role of bleu in machine translation research.
- (15) In 11th Conference of the European Chapter of the Association for Computational Linguistics. Jan KChorowski, DzmitryBahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and YoshuaBengio. 2015. Attention-based models for speech recognition. In Advances in Neural Information Processing Systems. 577&585. Ronan Collobert and Jason Weston.
- (16) 2008. A Unified Architecture for Natural Language Processing: Deep Neural Networks with Multitask Learning. In Proceedings of the 25th International Conference on Machine Learning (ICML '08). ACM, New York, NY, USA, 160&167. <https://doi.org/10.1145/1390156.1390177> Yaniv David, Nimrod Partush, and Eran Yahav. 2016. Statistical Similarity in Binaries. In PLDI'16: Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation. Yaniv David, Nimrod Partush, and Eran Yahav.
- (17) 2017. Similarity of Binaries through re-optimization. In PLDI'17: Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation