
Semantic Code Classification for Automated Machine Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Generating a complete machine learning pipeline from a short description in natural
2 language is a notoriously difficult task because the code is generally much longer
3 than the description and uses a lot of information from different parts of the
4 description. Using intermediate representation might help with this task since
5 specific actions require specific information. In this work, we present a semantic
6 code classification task and a way to represent the machine learning pipeline
7 as a sequence of such semantic classes. Finally, we discuss methods for solving
8 semantic code classification problem on the Natural Language to Machine Learning
9 (NL2ML) dataset.

10 **1 Introduction**

11 In recent years, machine learning was successfully applied to many tasks such as image classification
12 (He et al. (2015)) and even playing the game of Go (Silver et al. (2016)). Such models are able to
13 show superhuman performance (Russakovsky et al. (2014)). However, they were manually designed
14 by machine learning experts through trial and error. Such a process requires substantial computational,
15 temporary, and monetary resources.

16 The idea of automated machine learning (AutoML) emerged to reduce the costs of building a machine
17 learning pipeline and decrease the demand for data scientists (Zöller & Huber (2021)). Despite the
18 high interest in AutoML, so far, there is not much progress regarding a complete automated pipeline.
19 Most of the work in this area is focused on the specific steps of AutoML pipeline or severely limit the
20 set of models that can be used, so a complete machine learning pipeline remains an open problem
21 (He et al. (2021); Yang et al. (2020); Elsken et al. (2017)).

22 One of the standard formulations of AutoML task is building a system that can surpass humans on
23 Kaggle, a popular platform for open data science competitions. Its description defines the competition
24 in natural language and a dataset. However, generating a complete pipeline from a short description
25 is typically tricky because many generative models struggle with long sequences because of either
26 loss of context (Hochreiter (1998)) or high computational complexity (Vaswani et al. (2017)).

27 In this work, we propose an intermediate representation that might help with generating complete
28 machine learning pipelines. The proposed representation splits code into snippets so that each snippet
29 is responsible for a single small action in the pipeline. These small actions are called semantic code
30 classes. Such representation allows to better understand the structure of the ML code and the contents
31 of each step of possible pipelines, which can contribute to the generation task. We also present a
32 baseline for the semantic code classification task.

33 **2 Dataset**

34 Natural Language to Machine Learning (NL2ML) corpus is a collection of code blocks from Kaggle
35 competition notebooks. It includes 101 071 different notebooks comprising 2 587 074 cells (or code
36 block), which participants submitted to 266 Kaggle competitions. 4 748 of those code blocks have
37 passed through manual assessment of its semantics. Each block is attributed to an upper-level category,
38 which corresponds to a coarse meaning, and a second-level category, representing the finer semantic
39 of that block. Publication of the corpus, along with the detailed description of the annotation process,
40 is delayed due to legal considerations with Kaggle.

41 The upper level categories attribute the general kind of code block. One block can be-
42 long to a single category. Here is the list of possible values: Hypothesis, Environment,
43 Data extraction, Exploratory data analysis, Data transform, Model train, Model
44 evaluation, Model interpretation, Hyperparameter tuning, Visualization. The lower-
45 level category represents a more specific type, for example, upper-level class Data transform
46 contains such items as `create dataframe`, `remove duplicates`, `correct missing values`,
47 `feature engineering`, `filter` and others.

48 The dataset is very unbalanced (appendix 1). Some simple classes (`show table`) are overrepresented
49 because data contains many similar examples. At the same time, many classes (`load data from`
50 `zip`, `remove duplicates`) are represented by less than 10 examples. Another difficulty is that data
51 is heavily skewed because a large portion of the snippets came from the same competition; hence the
52 snippets in the same class tend to be similar in general.

53 **3 Code classification**

54 **3.1 Baseline solution: SVM**

55 Support vector machines classifier was used as baseline solution because the method does not require
56 a large amount of samples. The multiclass support is handled according to a one-vs-one scheme
57 (Chang & Lin (2011)). Baseline F1-score is 0.676 ± 0.029 .

58 Due to the imbalance in the data, the F1-score was chosen as the target metric.

59 **3.2 Preprocessing**

60 All of the code blocks were modified according to PEP8 using a standard Python utility `autopep8` to
61 eliminate the problems with code-style differences.

62 Comments have been removed from the code. Despite the fact that they can contain a lot of
63 information about the blocks semantic type, their presence, as a rule, slightly degrades the quality of
64 predictions. This is most likely because the content of the comments can be both useful and confusing.
65 The text may contain information unrelated to the code block or unnecessary code that does not affect
66 the semantic type.

67 For tokenization, the Byte-pair encoding method was used (Sennrich et al. (2015)). The tokenizer
68 was trained using unlabeled data. The size of the dictionary used is 30,000 tokens; during the training,
69 the BPE-Dropout regularization technique was used (Provlkov et al. (2019)).

70 This tokenization has proven to be quite helpful for many models. The names of entities in the
71 program code can consist of several words, written together or through an underscore "_". Splitting
72 by the spaces cannot distinguish different words included in the names of entities, and breaking into
73 n-grams may not be flexible since words in names can have different lengths. The BPE method was
74 optimal because it can extract information about the donkeys that are often found in the names of
75 entities.

76 Term frequency-inverse document frequency (TF-IDF) was used as features for models. These
77 statistics show the importance of the token in the document, given the context of the corpus.

78 3.3 Naive Bayesian Classifier and NB-SVM

79 The Naive Bayesian Classifier (NB) is a probabilistic model based on Bayes' theorem. Let y be
80 the target variable, and x_1, \dots, x_d be features. Then the probabilistic model approximates the
81 conditional density:

$$p(y|x_1, \dots, x_d) = \frac{p(x_1, \dots, x_d|y)}{p(x_1, \dots, x_d)}$$

82 Naive Bayesian Classifier is a common pre-neural network method for classifying text data. As a
83 rule, it is assumed that the data has either a multinomial distribution or a Bernoulli distribution.

84 Since the dataset is unbalanced, estimates for the parameters of the multinomial distribution may be
85 unstable, which negatively affects the quality of predictions. To combat this problem, a complemen-
86 tary NB was proposed in Rennie et al. (2003).

87 The model with the Bernoulli distribution showed the best result among naive Bayesian classifiers.
88 This is because this model works best with short texts (Mccallum & Nigam (2001)). Short texts
89 predominate in the training dataset.

90 Metrics for NB classifier are in rows Multinomial NB, Complementary NB, Bernoulli NB of
91 table 2.

92 NB classifier and SVM are often used as basic solutions in text classification problems. NB tends
93 to perform better on shorter passages, while SVM works best on longer texts. The article Wang
94 & Manning (2012) discusses the composition of an NB classifier and SVM, which combines the
95 advantages of both methods.

96 3.4 Neural networks

97 Neural networks are an important class of machine learning algorithms. It is due to neural networks
98 that the quality of text data processing has been improved. Unfortunately, many advanced neural
99 network architectures require a lot of training data. There was not too much data in this case, but it
100 was possible to obtain relatively good results.

101 There is not enough data to train vector representations, so frozen pretrained CodeBERT (Feng et al.
102 (2020)) was used. CodeBERT has several drawbacks that are worth noting. Firstly, the important thing
103 is that this model uses its own tokenizer. The tokens used do not quite match the machine learning
104 code as they were derived from more general code, which means that machine learning-specific
105 tokens may be too rare to enter the dictionary. Secondly, CodeBERT was trained on the program
106 code not only in Python3 but also in several other languages: JavaScript, Java, PHP, Ruby, Go. Due
107 to this, the connection of tokens with machine learning is becoming even less.

108 Another important point when training a neural network is the loss function, on which the success
109 of the model depends a lot. The standard loss function in a classification problem is cross-entropy.
110 Conventional cross-entropy is unstable to unbalanced data; the situation can be improved by adding
111 weights for examples from different classes. Since there are a lot of classes, a vast number of
112 hyperparameters appear, which is a big drawback since neural networks take a long time to learn.

113 For unbalanced data, there are alternative loss functions:

- 114 1. Soft-F1. This loss function directly optimizes the target metric in the task at hand, which is a
115 significant advantage. The critical drawback of this loss function is that it requires a sizeable
116 mini-batch size. If the batch size is too small, rare classes will almost not be included in the
117 batch, which is why the completeness and accuracy for such classes will often take large
118 values. If there are many such rare classes, the loss function will be extremely close to the
119 optimum, leading to the fact that learning does not occur. In the dataset, quite a few classes
120 are containing less than 10 objects, which means that large batch size is required.
- 121 2. Focal Loss. This loss function comes from computer vision (Lin et al. (2017)). The main
122 idea is to adaptively reduce the weight of objects for which the model determines the class
123 correctly and confidently. Due to this, it is possible to learn patterns even in rare classes.

124 The following architectures were used:

- 125 1. Recurrent neural network (RNN). GRU and LSTM were tried as a recurrent layer; the result
126 is approximately the same. The recurrent network is bi-directional since this configuration
127 usually gives the best result, allowing more context for the tokens. Row RNN in results table
128 (2).
129 The token sequence is fed into the inputs of the recurrent layer. The output goes through a
130 dense layer, GELU activation, dropout, and outputs probabilities through dense layer and
131 Softmax.
132 Layer sizes are optimized with random search. The models with best scores uses a single
133 LSTM layer with hidden state size of 180 and a hidden fully connected layer of size 269.
134 The scheme is in appendix 2.
- 135 2. RNN with additive attention (Bahdanau et al. (2014)). The model differs from the usual
136 RNN in only one layer: the attention layer added after the recurrent layer.
137 The sizes of all layers were optimized using random search. The best model uses a single
138 LSTM layer with hidden state size of 100. After attention, there is a hidden fully connected
139 layer of size 135. The scheme of the model can be found in appendix 3.
140 A layer of attention helps to identify the most important values in a sequence, which leads
141 to a dramatic improvement in the result. Row RNN+Attention in results table (2).

142 3.5 Augmentation

143 Since the marked-up data is insufficient, augmentation was used. However, it is necessary that the
144 code remains compilable and belongs to the same semantic type after modifications.

145 Many code blocks contain similar variable names (eg. "df", "train", "test"). They do not
146 always provide additional information about the block type, but due to the frequency they are used by
147 the classifier. Therefore, augmentation was applied, which masks some fraction of the variables in
148 the block. The variables in code were found using AST trees and replaced by masks, each variable in
149 a block with a different one.

150 The augmentation was used with SVM and Attention RNN. Results are in rows SVM+Augmentation,
151 RNN+Attention+Augmentation in table (2).

152 3.6 Hierarchy

153 The structure by which the code snippets are classified is a two-level graph: the first is the general
154 class of action, the second is what exactly is done in the block. Therefore, it was decided to implement
155 this structure as a two-level classifier, where the upper vertex is predicted first, and then the lower by
156 a separate classifier, already trained on specific objects belonging to the child vertices.

157 The advantage of considering the upper classes separately is that in contrast to the original dataset,
158 each class is represented by at least several dozen objects. At the same time, most of the very small
159 original classes go into such a classification together with other small classes. Thus, the determination
160 of their upper class is a more realistic task. If we correctly guessed the upper class, then small objects
161 are likely to have no large competitors. However, if the class was identified incorrectly in the first
162 step, the subsequent classification is guaranteed to be inaccurate.

163 SVM was chosen as the lower classifier since the model demonstrated the best results in previous
164 experiments, while the upper classifiers were

- 165 1. RNN containing a bidirectional LSTM, the vector representations were taken from the
166 frozen pre-trained CodeBERT (Feng et al. (2020)).
- 167 2. SVM

168 In table 1 the metrics of the upper level classification are presented. SVM achieved sufficient
169 results. The overall results of hierarchy classifier are in 2 - rows RNN + SVM Hierarchy, SVM +
170 SVM Hierarchy - this method allowed to significantly improve metrics.

171 3.7 Semi-supervised, pseudo labels

172 Since the amount of labeled data is extremely small, it was decided to experiment with semi-
173 supervised models that allow to use unlabelled data (Xie et al. (2019)).

Table 1: Upper level prediction results.

Model	Metrics	
	F1-score	Accuracy
RNN	0.759	0.773
SVM	0.937	0.938

174 First, we trained SVM model on the marked-up dataset. It’s predictions on unlabeled data were used
 175 as pseudo labels. Then another SVM model with separately selected hyperparameters was trained on
 176 pseudo-labeled data and a part of marked-up data. It was tested on the rest of the labeled data.

177 We tested this approach using 20%, 40%, and 100% of unlabeled data. The results are in rows
 178 Pseudo labels 20%, Pseudo labels 40%, Pseudo labels 100% in table (2).

179 4 Results

180 We evaluated models on marked-up code snippets of the NL2ML dataset. Several approaches
 181 demonstrated some improvements in metrics from the baseline solution. All of the resulting models
 182 are compared in the table (2) by F1-score and accuracy. The best results are achieved by hierarchy
 183 and pseudo-labels models. The first one utilizes the knowledge of the classes structure, and the
 184 second one uses a great amount of unlabeled data.

185 The target metrics were computed on a test dataset that contains 20 % of all labeled data. For each
 186 model, we optimized hyperparameters with random search. For neural networks, we used a validation
 187 set that contains 20 % of labeled data. For other methods, we used cross-validation with 10 folds.

Table 2: Summary table of all the models and their results

Model	Metrics	
	F1-score	Accuracy
SVM+Linear (Baseline)	0.676 ± 0.029	0.678 ± 0.031
SVM + Poly	0.620 ± 0.033	0.627 ± 0.030
SVM + RBF	0.672 ± 0.031	0.676 ± 0.031
Multinomial NB	0.561 ± 0.031	0.583 ± 0.031
Complementary NB	0.576 ± 0.024	0.587 ± 0.023
Bernoulli NB	0.609 ± 0.021	0.623 ± 0.026
NBSVM	0.694 ± 0.013	0.704 ± 0.015
NBSVM (binarization)	0.684 ± 0.016	0.694 ± 0.015
NBSVM (bigrams)	0.689 ± 0.14	0.697 ± 0.14
RNN	0.541	0.556
RNN + Attention	0.639	0.648
SVM + Augmentation	0.686 ± 0.031	0.696 ± 0.029
RNN + Attention + Augmentation	0.638	0.652
RNN + SVM Hierarchy	0.291	0.282
SVM+SVM Hierarchy	0.741	0.744
Pseudo labels 20 %	0.713 ± 0.014	0.724 ± 0.014
Pseudo labels 40 %	0.727 ± 0.016	0.738 ± 0.014
Pseudo labels 100 %	0.755 ± 0.03	0.764 ± 0.028

188 **5 Conclusion**

189 This paper compares methods for semantic classification of code snippets from NL2ML corpus, a
190 large-scale dataset of data science-specific code harvested from Kaggle, the largest platform for data
191 science competitions. The dataset includes only 4 748 manually labeled snippets.

192 We provided detailed qualitative and quantitative comparisons of the used classification methods. The
193 final proposed model demonstrated better results than the baseline SVM solution, which has an
194 F1-score of 0.67. The method with the best F1-score of 0.755 is an SVM trained on pseudo-labels.
195 The demonstrated results are significantly better than those obtained by solutions without Machine
196 Learning, which only achieve an F1-score of ≈ 0.4 .

197 The proposed method allows to automatically markup a program code corpus with better accuracy
198 and performance as well as effectively assess the quality of data science code in Python. We also
199 genuinely hope that this work can contribute to the development of the next generation of robust
200 AutoML systems as semantic classification provides an interlingua that can be used for translating
201 between natural and programming languages.

202 **References**

203 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly
204 learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.

205 Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. *ACM*
206 *transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.

207 Thomas Elsken, Jan-Hendrik Metzen, and Frank Hutter. Simple and efficient architecture search for
208 convolutional neural networks, 2017.

209 Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing
210 Qin, Ting Liu, Daxin Jiang, and Ming Zhou. Codebert: A pre-trained model for programming and
211 natural languages. *CoRR*, abs/2002.08155, 2020. URL <https://arxiv.org/abs/2002.08155>.

212 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
213 recognition. *CoRR*, abs/1512.03385, 2015. URL <http://arxiv.org/abs/1512.03385>.

214 Xin He, Kaiyong Zhao, and Xiaowen Chu. Automl: A survey of the state-of-the-art. *Knowledge-*
215 *Based Systems*, 212:106622, 2021.

216 Sepp Hochreiter. The vanishing gradient problem during learning recurrent neural nets and problem
217 solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 6:
218 107–116, 04 1998. doi: 10.1142/S0218488598000094.

219 Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense
220 object detection. *CoRR*, abs/1708.02002, 2017. URL <http://arxiv.org/abs/1708.02002>.

221 Andrew McCallum and Kamal Nigam. A comparison of event models for naive bayes text classifica-
222 tion. *Work Learn Text Categ*, 752, 05 2001.

223 Ivan Provilkov, Dmitrii Emelianenko, and Elena Voita. Bpe-dropout: Simple and effective subword
224 regularization. *CoRR*, abs/1910.13267, 2019. URL <http://arxiv.org/abs/1910.13267>.

225 Jason D. M. Rennie, Lawrence Shih, Jaime Teevan, and David R. Karger. Tackling the poor assump-
226 tions of naive bayes text classifiers. In *Proceedings of the Twentieth International Conference on*
227 *International Conference on Machine Learning, ICML'03*, pp. 616–623. AAAI Press, 2003. ISBN
228 1577351894.

229 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
230 Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-
231 Fei. Imagenet large scale visual recognition challenge. *CoRR*, abs/1409.0575, 2014. URL
232 <http://arxiv.org/abs/1409.0575>.

233 Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with
234 subword units. *CoRR*, abs/1508.07909, 2015. URL <http://arxiv.org/abs/1508.07909>.

- 235 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
236 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
237 the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- 238 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
239 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information*
240 *processing systems*, pp. 5998–6008, 2017.
- 241 Sida Wang and Christopher Manning. Baselines and bigrams: Simple, good sentiment and topic
242 classification. In *Proceedings of the 50th Annual Meeting of the Association for Computational*
243 *Linguistics (Volume 2: Short Papers)*, pp. 90–94, Jeju Island, Korea, July 2012. Association for
244 Computational Linguistics. URL <https://aclanthology.org/P12-2018>.
- 245 Qizhe Xie, Eduard H. Hovy, Minh-Thang Luong, and Quoc V. Le. Self-training with noisy student
246 improves imagenet classification. *CoRR*, abs/1911.04252, 2019. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1911.04252)
247 [1911.04252](http://arxiv.org/abs/1911.04252).
- 248 Chengrun Yang, Jicong Fan, Ziyang Wu, and Madeleine Udell. Automl pipeline selection: Efficiently
249 navigating the combinatorial space. In *Proceedings of the 26th ACM SIGKDD International*
250 *Conference on Knowledge Discovery & Data Mining, KDD '20*, pp. 1446–1456, New York, NY,
251 USA, 2020. Association for Computing Machinery. ISBN 9781450379984. doi: 10.1145/3394486.
252 3403197. URL <https://doi.org/10.1145/3394486.3403197>.
- 253 Marc-André Zöller and Marco F Huber. Benchmark and survey of automated machine learning
254 frameworks. *Journal of Artificial Intelligence Research*, 70:409–472, 2021.

255 **A** Distribution of snippets by semantic types

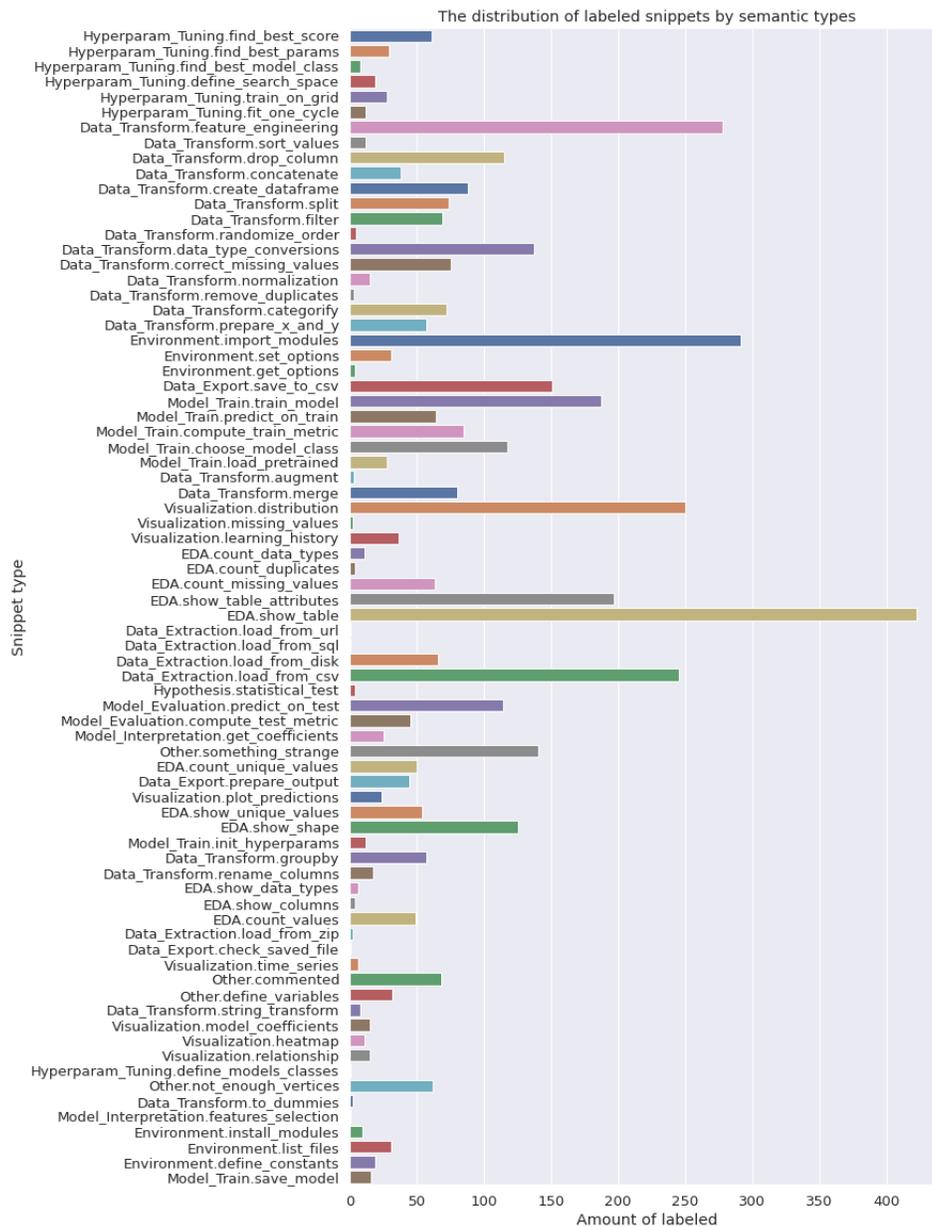


Figure 1: The distribution of labeled snippets by semantic types.

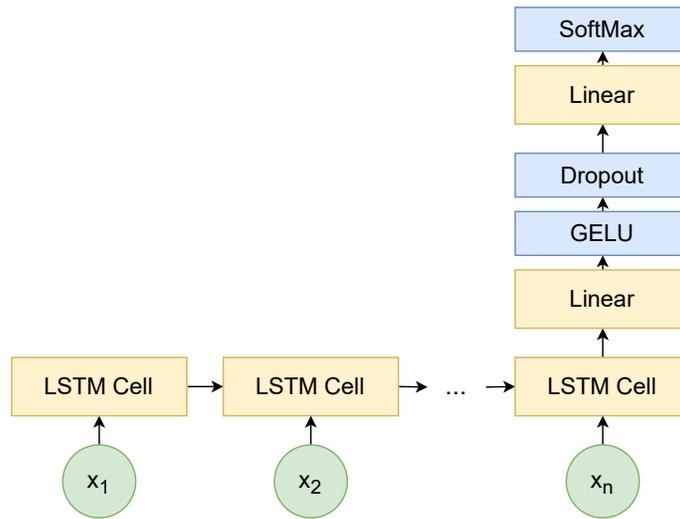


Figure 2: Architecture of recurrent neural network

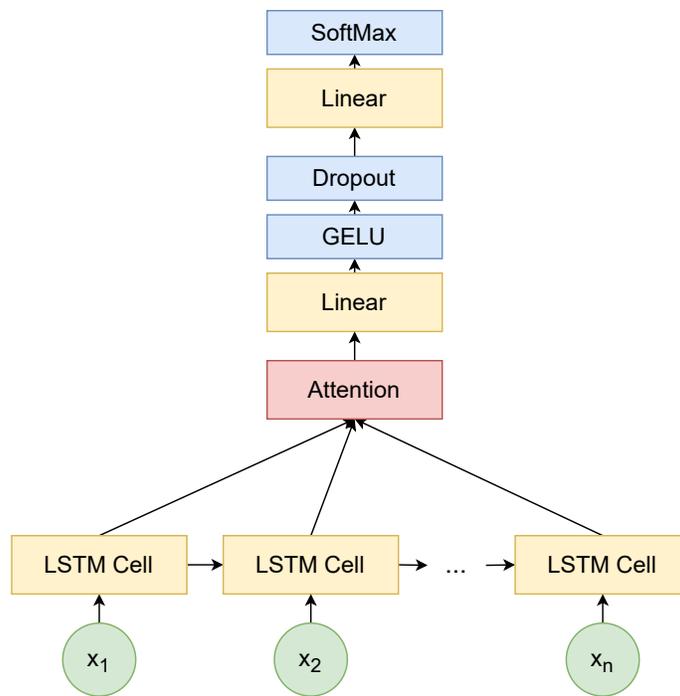


Figure 3: Architecture of recurrent neural network with attention

Table 3: Examples of code snippets from NL2ML dataset

Semantic class	Example
Data_Transform.drop_column	<pre>train_df.drop("Date", inplace=True, axis=1) test_df.drop("Date", inplace=True, axis=1)</pre>
Model_Train.choose_model_class	<pre>from sklearn import linear_model reg_CC = linear_model.Lasso(alpha=0.1) reg_Fat = linear_model.Lasso(alpha=0.1)</pre>
Hyperparams.define_search_space	<pre>parameters = {'lstm_nodes': [14,16,20], 'nb_epoch': [50], 'batch_size': [32], 'optimizer': ['adam']}</pre>
Visualization.distribution	<pre>fig = plt.figure() fig.suptitle("Algorithm_Comparison") ax = fig.add_subplot(111) plt.boxplot(results) ax.set_xticklabels(names) plt.show()</pre>
Data_Transform.normalization	<pre>scaler = MinMaxScaler() df['revenue'] = scaler.fit_transform(df[['revenue']])</pre>
Model_Train.train_model	<pre>rfr = RandomForestRegressor(n_estimators=200, max_depth=5, max_features=0.5, random_state=449, n_jobs=-1) rfr.fit(x_train, y_train)</pre>

D Selected hyperparameters for each model

Table 4: Hyperparameters for models without hierarchy and first level of hierarchical models

Model	Hyperparameter	Type	Value
SVM+Linear (Baseline)	C	numeric	37.17
	min_df for TF-IDF	integer	2
	max_df for TF-IDF	numeric	0.31
SVM + Poly	C	numeric	1.43
	Degree of poly kernel	integer	3
	min_df for TF-IDF	integer	6
	max_df for TF-IDF	numeric	0.30
SVM + RBF	C	numeric	8.71
	min_df for TF-IDF	integer	7
	max_df for TF-IDF	numeric	0.39
Multinomial NB	Alpha	numeric	0.01
	min_df for TF-IDF	integer	2
	max_df for TF-IDF	numeric	0.57
Complementary NB	Alpha	numeric	0.29
	min_df for TF-IDF	integer	2
	max_df for TF-IDF	numeric	0.30
Bernoulli NB	Alpha	numeric	0.01
	min_df for TF-IDF	integer	5
	max_df for TF-IDF	numeric	0.30
NBSVM	C for SVM	numeric	0.50
	Kernel type	categorical	Linear
	Alpha for Naive Bayes	numerical	4.22
	min_df for TF-IDF	integer	1
	max_df for TF-IDF	numeric	0.94
NBSVM (binarization)	C for SVM	numeric	0.10
	Kernel type	categorical	Linear
	Alpha for Naive Bayes	numerical	0.28
	min_df for TF-IDF	integer	3
	max_df for TF-IDF	numeric	0.94
NBSVM (bigrams)	C for SVM	numeric	0.10
	Kernel type	categorical	Linear
	Alpha for Naive Bayes	numerical	2.51
	min_df for TF-IDF	integer	2
	max_df for TF-IDF	numeric	0.95
RNN	RNN cell type	categorical	LSTM
	Size of hidden state	integer	180
	Size of linear layer	integer	269
RNN + Attention	RNN cell type	categorical	LSTM
	Size of hidden state	integer	100
	Size of linear layer	integer	135
SVM + Augmentation	C	numeric	5.82
	Kernel type	categorical	Linear
	min_df for TF-IDF	integer	2
	max_df for TF-IDF	numeric	0.30
	Percent of masked variables	numeric	0.93
RNN + Attention + Augmentation	RNN cell type	categorical	LSTM
	Size of hidden state	integer	169

	Size of linear layer	integer	255
	Percent of masked variables	numeric	99.7%
SVM+SVM Hierarchy (1st level)	C	numeric	149.65
	Kernel type	categorical	Poly
	Degree of poly kernel	integer	2
RNN+SVM Hierarchy (1st level)	RNN cell type	categorical	LSTM
	Size of hidden state	integer	143
	Size of linear layer	integer	185
Pseudo labels 20 %	C	numeric	98.37
	Kernel type	categorical	linear
	min_df for TF-IDF	integer	3
	max_df for TF-IDF	numeric	0.53
Pseudo labels 40 %	C	numeric	121.59
	Kernel type	categorical	linear
	min_df for TF-IDF	integer	3
	max_df for TF-IDF	numeric	0.41
Pseudo labels 100%	C	numeric	145.56
	Kernel type	categorical	linear
	min_df for TF-IDF	integer	2
	max_df for TF-IDF	numeric	0.26

Table 5: Hyperparameters for SVM+SVM Hierarchy (2nd level). Separate parameters for each upper level class. If upper level class is not present in this table, then it has one subclass and hence does not require a second level classifier.

Upper level class	Hyperparameter	Type	Value
Data Extraction	C	numeric	135.49
	Kernel type	categorical	Linear
	min_df for TF-IDF	integer	50
	max_df for TF-IDF	numeric	0.89
EDA	C	numeric	4.70
	Kernel type	categorical	Linear
	min_df for TF-IDF	integer	4
	max_df for TF-IDF	numeric	0.61
Environment	C	numeric	232.06
	Kernel type	categorical	Linear
	min_df for TF-IDF	integer	6
	max_df for TF-IDF	numeric	0.81
Model_Train	C	numeric	9.63
	Kernel type	categorical	RBF
	min_df for TF-IDF	integer	17
	max_df for TF-IDF	numeric	0.47
Data_Transform	C	numeric	298.83
	Kernel type	categorical	Linear
	min_df for TF-IDF	integer	2
	max_df for TF-IDF	numeric	0.20
Model_Evaluation	C	numeric	0.84
	Kernel type	categorical	Poly
	Degree of poly kernel	integer	3
	min_df for TF-IDF	integer	16
	max_df for TF-IDF	numeric	0.50

Data_Export	C	numeric	5.35
	Kernel type	categorical	RBF
	min_df for TF-IDF	integer	23
	max_df for TF-IDF	numeric	0.83
Hyperparam_Tuning	C	numeric	20.03
	Kernel type	categorical	RBF
	min_df for TF-IDF	integer	6
	max_df for TF-IDF	numeric	0.26
Visualization	C	numeric	865.58
	Kernel type	categorical	Linear
	min_df for TF-IDF	integer	2
	max_df for TF-IDF	numeric	0.66

Table 6: Hyperparameters for RNN+SVM Hierarchy (2nd level). Separate parameters for each upper level class. If upper level class is not present in this table, then it has one subclass and hence does not require a second level classifier.

Upper level class	Hyperparameter	Type	Value
Data Extraction	C	numeric	7.66
	Kernel type	categorical	RBF
	min_df for TF-IDF	integer	17
	max_df for TF-IDF	numeric	0.88
EDA	C	numeric	1.71
	Kernel type	categorical	Linear
	min_df for TF-IDF	integer	6
	max_df for TF-IDF	numeric	0.32
Environment	C	numeric	462.65
	Kernel type	categorical	RBF
	min_df for TF-IDF	integer	16
	max_df for TF-IDF	numeric	0.87
Model_Train	C	numeric	4.58
	Kernel type	categorical	Linear
	min_df for TF-IDF	integer	1
	max_df for TF-IDF	numeric	0.90
Data_Transform	C	numeric	204.62
	Kernel type	categorical	Linear
	min_df for TF-IDF	integer	2
	max_df for TF-IDF	numeric	0.45
Model_Evaluation	C	numeric	2.03
	Kernel type	categorical	Linear
	min_df for TF-IDF	integer	3
	max_df for TF-IDF	numeric	0.63
Data_Export	C	numeric	26.49
	Kernel type	categorical	Poly
	Degree of poly kernel	numeric	2
	min_df for TF-IDF	integer	23
	max_df for TF-IDF	numeric	0.93
Hyperparam_Tuning	C	numeric	14.58
	Kernel type	categorical	Linear
	min_df for TF-IDF	integer	10
	max_df for TF-IDF	numeric	0.65
Visualization	C	numeric	2.48

Kernel type	categorical	Linear
min_df for TF-IDF	integer	2
max_df for TF-IDF	numeric	0.21
