

AI for S/W Eng Assignment #1 Report

Recommending code tokens via N-gram models

Nishat Sultana, Danny Otten, Joseph Call

Dataset Creation Process

Our initial plan was to create a new corpus of at least 25,000 Java methods, but due to DataHub creating overly large datasets or failing to generate them entirely, we had no access to data early on in the model development process. Therefore, we decided to use the example dataset generated in class on DataHub to test our N-Gram model code. This was enough to get us started in terms of:

- Data format
- Data preprocessing

We then created our own dataset on DataHub for **Java** methods using the following parameters:

The screenshot shows the DataHub dataset creation interface. It is divided into three main sections: Repository Sample Characteristics, Dataset Characteristics, and Code Filters & Processing.

Repository Sample Characteristics: A dropdown menu shows 'Java'. Below it, four filters are set: Commits (100), Issues (100), Contributors (10), and Stars (100). There are checkboxes for 'Has Open-source License' (unchecked) and 'Exclude Forks' (checked).

Dataset Characteristics: A section for 'Granularity' with a 'Metadata' tab selected. It contains four checkboxes: 'File' (unchecked), 'Function' (checked), 'Pair each instance with its Symbolic Expression representation' (unchecked), 'Pair each instance with its AST-based representation' (unchecked), and 'Pair each instance with tree-sitter parser metadata' (unchecked).

Code Filters & Processing: A section for 'Characters' with three filters: 'Characters' (200000), 'Tokens' (10000), and 'Lines' (100000). All are set to 'max'. There is an 'Exclude' section with checkboxes for 'Test code' (checked), 'Boilerplate code' (checked), 'Instances with syntax errors' (checked), and 'Instances with non-ASCII characters' (unchecked). There is also an 'Ignore' section with checkboxes for 'Duplicates' (checked), 'Near-clones' (unchecked), 'Regular comments' (unchecked), and 'Documentation comments' (unchecked). A 'Submit' button is at the bottom right.

This dataset was about 6 gigabytes when compressed and 40 gigabytes when uncompressed, which is too large to use for a full training run. Because of this, we created a subset of this file which we used to actually train our models. Our final dataset was 134 megabytes and included 30,000 Java methods.

Model Training Methodology

For simplicity, we decided to write our model code in Python. Our control file input parameters such as:

- Number of models to train and compare
- Number of grams
- Number of classes to use in the training set

- Number of classes to use in the testing set

The model training process takes on average <training minutes> and we found our model was **overfit**. When a model is finished training, its performance will be evaluated by calculating the perplexity, and the model with lowest perplexity out of all those created will be returned as the best.

From there, the user is able to enter code snippets for code generation. The user can then interact with the model.

Code

We used a GitHub repository to manage our code during this project, which can be found here: <https://github.com/dsotten/plm>