

Predicting Research Paper Citation Counts with Gradient Boosting Regression

Question 3: Model Pipeline Explanation

The model pipeline uses a gradient boosting regressor model that is highly parameterized with multi-stage data preprocessing. The feature transformation starts with performing TF-IDF vectorization of the combined text features based on bigram and with 1500 maximum features. Author count, reference count, year of publication and length of the texts have their median values imputed and then scaled to standard values. Two other descriptive features: categorical venue features are one-hot encoded with the maximum number of categories of 10. GradientBoostingRegressor uses 250 estimators, 0.07 learning rate, maximum depth of 4 and 'sqrt' for features sampling. This pipeline consists of several feature engineering to produce precise academic paper citation prediction along with strong regression techniques to handle complex feature dependencies and intelligent feature transformation.

Question 4: Model Development Journey

The development of models that are being described in this paper proceeded through the refinement of the initial ideas. Firstly, we considered a basic application of linear regression and found its application unsuitable to explain intricate citation behavior. Switching to the random forest algorithm then improved the performance but at the same time reduced the explanatory power of the results. Surprisingly, an enhanced approach known as Gradient Boosting appeared to provide the greatest level of predictive accuracy. In a similar manner where machine learning is done in a black box, we systematically began to do feature engineering incrementally by adding more metadata features including author count, reference networks, and text based characteristics. Thus, the preprocessing methods changed from the simplest scaling option to median imputation and TF-IDF vectors. To fine-tune hyperparameters, cross-validation was employed used to run the model, and to balance between model complexity and the risks of overfitting in order to minimize mean absolute error.

Question 5: Additional Solution Insights

It asserts generalization as the solution with increased intelligent feature representation. Validation procedures guarantee that a model would be accurate across different academic paper datasets. To augment the initial feature space, we built a textual representation of the papers by processing the abstract, title, and venue of each paper, even though the information in each can often be quantified by simple numeric features. The preprocessing strategy deals with missing values appropriately, replacing them by median for the numeric variables and by a constant if the original variables were

categorical. The flexibility is offered by the unknown category handling while performing the one-hot encoding. The proposed model successfully avoids overfitting by introducing model complexity constraint, implementing moderate learning rate, and by performing thorough feature engineering of citations where needed and none where not thus the developed citation prediction model is flexible and effective.