

Cell2Doc Supplementary Material

August 18, 2023

1 Cell2Doc Outputs on Sample Examples

Figures 1, 2, 3 and 4 show output documentation on sample code examples for PLBART [1], CodeBERT [2], UnixCoder [3], and GraphCodeBERT [4] respectively. We can see from these examples that with Cell2Doc pipeline, all these pre-trained models have provided longer and more informative documentation compared to their base stand-alone versions (shown in the previous columns in the table). These observations highlight that our Cell2Doc pipeline is model agnostic and can improve the effectiveness of other pre-trained models.

2 Cell2Doc Human Evaluation with CodeBERT

We have done the human evaluation of Cell2Doc using CodeT5 and CodeBERT as the CoDoc model. As CodeBERT is an encoder-only model, we have used a 6-layer transformer decoder with it to build the CoDoc model. We have only used CodeT5 in the main paper because it has performed better than CodeBERT overall. Figure 5 summarizes the human evaluation results using CodeBERT in CoDoc of Cell2Doc. Like CodeT5 [5], we can see that Cell2Doc also improves CodeBERT’s effectiveness across all three dimensions.

3 Other Input Representations

Inline code comments can provide additional information about the code, which can be utilized to create better documentation. We tested this input representation idea with the CodeBERT model in CoDoc, which can take both code and natural language as input. Code and comments present in a code cell can be represented as $[code_snippet_1, comment_1, code_snippet_2, comment_2, \dots, code_snippet_n]$ where $code_snippet_i$ are code segments (present between two consecutive comments) and $comment_i$ are the comments. We use it as input to CodeBERT by using a delimiter between code and comment. As target documentation, we use only summarized markdown in these cases. The BLEU [6] and ROUGE [7] (only F1 for ROUGE-1, ROUGE-2 and ROUGE-L) scores for this experiment are **20.37** and (**25.29**, **11.60**, and **28.24**) respectively. Compared with the automated evaluation where CodeBERT is used in CoDoc, we can see that the numbers are better when comments are included in the input with respect to other baselines. It does not beat the SCSCM representation, though, and the reason for it might be that SCSCM has more data points in the training set as it uses comments as separate labels, and comments serve as good quality labels as well.

Code cell	PLBART	Cell2Doc (PLBART)
<pre> n_classes = 12 data, labels = make_classification(n_samples =2000, n_features=100, n_informative=50, n_classes=n_classes, random_state= random_state) X, X_test, y, y_test = train_test_split(data , labels, test_size=0.2, random_state= random_state) X_train, X_valid, y_train, y_valid = train_test_split(X, y, test_size=0.25, random_state=random_state) print('Data shape:') print('X_train: %s, X_valid: %s, X_test: %s \n' %(X_train.shape, X_valid.shape, X_test.shape)) </pre>	Splitting the data into training and validation data	<p>STEP 1: Number of classes</p> <p>STEP 2: Create classification object</p> <p>STEP 3: taking holdout set for validating with random state</p> <p>STEP 4: Splitting the dataset into training and validation</p> <p>STEP 5: Check data shape</p>
<pre> def cel(y_true, y_pred): y_true = torch.argmax(y_true, axis=-1) return nn.CrossEntropyLoss()(y_pred, y_true.squeeze()) def accuracy(y_true, y_pred): y_true = torch.argmax(y_true, axis=-1). squeeze() y_pred = torch.argmax(y_pred, axis=-1). squeeze() return (y_true == y_pred).float().sum()/ len(y_true) </pre>	Evaluation of the Model	<p>STEP 1: The function for evaluation is borrowed from</p> <p>STEP 2: Cross entropy loss</p>
<pre> train_images = [] image_dirs = np.take(os.listdir('../input/ train'), select_rows) for image_dir in tqdm(sorted(image_dirs)): image = imread('../input/train/'+ image_dir) train_images.append(image) del image gc.collect() train_images = np.array(train_images) </pre>	Read in the train images	<p>STEP 1: Read in the training images</p> <p>STEP 2: Read in all the images</p> <p>STEP 2.1: Read in image and save in list</p> <p>STEP 2.2: Delete to reduce memory usage</p>

Figure 1: Code and generated documentation examples using PLBART in CoDoc of Cell2Doc

Code cell	CodeBERT	Cell2Doc (CodeBERT)
<pre> n_classes = 12 data, labels = make_classification(n_samples =2000, n_features=100, n_informative=50, n_classes=n_classes, random_state= random_state) X, X_test, y, y_test = train_test_split(data , labels, test_size=0.2, random_state= random_state) X_train, X_valid, y_train, y_valid = train_test_split(X, y, test_size=0.25, random_state=random_state) print('Data shape:') print('X_train: %s, X_valid: %s, X_test: %s \n' %(X_train.shape, X_valid.shape, X_test.shape)) </pre>	Train a Classifier	<p>STEP 1: Number of classes per image</p> <p>STEP 2: Create a Classification object</p> <p>STEP 3: Train Validation Split</p> <p>STEP 4: Split the dataset into training and validation</p> <p>STEP 5: Checking the shape of the data</p>
<pre> def cel(y_true, y_pred): y_true = torch.argmax(y_true, axis=-1) return nn.CrossEntropyLoss()(y_pred, y_true.squeeze()) def accuracy(y_true, y_pred): y_true = torch.argmax(y_true, axis=-1). squeeze() y_pred = torch.argmax(y_pred, axis=-1). squeeze() return (y_true == y_pred).float().sum()/ len(y_true) </pre>	Cross entropy loss	<p>STEP 1: Cross entropy loss</p> <p>STEP 2: The accuracy of the classifier</p>
<pre> train_images = [] image_dirs = np.take(os.listdir('../input/ train'), select_rows) for image_dir in tqdm(sorted(image_dirs)): image = imread('../input/train/'+ image_dir) train_images.append(image) del image gc.collect() train_images = np.array(train_images) </pre>	load train images	<p>STEP 1: Get train images</p> <p>STEP 2: Read in each image</p> <p>STEP 2.1: Read in the image</p> <p>STEP 2.2: Delete image to reduce memory usage</p>

Figure 2: Code and generated documentation examples using CodeBERT in CoDoc of Cell2Doc

Code cell	UnixCoder	Cell2Doc (UnixCoder)
<pre> n_classes = 12 data, labels = make_classification(n_samples =2000, n_features=100, n_informative=50, n_classes=n_classes, random_state= random_state) X, X_test, y, y_test = train_test_split(data , labels, test_size=0.2, random_state= random_state) X_train, X_valid, y_train, y_valid = train_test_split(X, y, test_size=0.25, random_state=random_state) print('Data shape:') print('X_train: %s, X_valid: %s, X_test: %s \n' %(X_train.shape, X_valid.shape, X_test.shape)) </pre>	Prepare data for modeling	<p>STEP 1: Number of classes used for visualization</p> <p>STEP 2: Create data and labels for model</p> <p>STEP 3: Split the data into train and test data</p> <p>STEP 4: Split the dataset into training and validation sets</p> <p>STEP 5: Check data types</p>
<pre> def cel(y_true, y_pred): y_true = torch.argmax(y_true, axis=-1) return nn.CrossEntropyLoss()(y_pred, y_true.squeeze()) def accuracy(y_true, y_pred): y_true = torch.argmax(y_true, axis=-1). squeeze() y_pred = torch.argmax(y_pred, axis=-1). squeeze() return (y_true == y_pred).float().sum() / len(y_true) </pre>	Cross entropy loss	<p>STEP 1: Cross entropy loss</p> <p>STEP 2: Defining the accuracy</p>
<pre> train_images = [] image_dirs = np.take(os.listdir('../input/ train'), select_rows) for image_dir in tqdm(sorted(image_dirs)): image = imread('../input/train/' + image_dir) train_images.append(image) del image gc.collect() train_images = np.array(train_images) </pre>	Read and shuffle file-names	<p>STEP 1: Read in the training dataset</p> <p>STEP 2: Read in images</p> <p>STEP 2.1: Read in the training images</p> <p>STEP 2.2: Free up some memory</p>

Figure 3: Code and generated documentation examples UnixCoder in CoDoc of Cell2Doc

Code cell	GraphCodeBERT	Cell2Doc (GraphCodeBERT)
<pre> n_classes = 12 data, labels = make_classification(n_samples =2000, n_features=100, n_informative=50, n_classes=n_classes, random_state= random_state) X, X_test, y, y_test = train_test_split(data , labels, test_size=0.2, random_state= random_state) X_train, X_valid, y_train, y_valid = train_test_split(X, y, test_size=0.25, random_state=random_state) print('Data shape:') print('X_train: %s, X_valid: %s, X_test: %s \n' %(X_train.shape, X_valid.shape, X_test.shape)) </pre>	Splitting the dataset into train and validation set	<p>STEP 1: Define some constants</p> <p>STEP 2: Use data augmentation</p> <p>STEP 3: Splitting the into train and test data</p> <p>STEP 4: Split datas in train and validation set</p> <p>STEP 5: Overview of the data</p>
<pre> def cel(y_true, y_pred): y_true = torch.argmax(y_true, axis=-1) return nn.CrossEntropyLoss()(y_pred, y_true.squeeze()) def accuracy(y_true, y_pred): y_true = torch.argmax(y_true, axis=-1). squeeze() y_pred = torch.argmax(y_pred, axis=-1). squeeze() return (y_true == y_pred).float().sum() / len(y_true) </pre>	NumtaDB Classification Report	<p>STEP 1: Define loss function</p> <p>STEP 2: define eval metrics</p>
<pre> train_images = [] image_dirs = np.take(os.listdir('../input/ train'), select_rows) for image_dir in tqdm(sorted(image_dirs)): image = imread('../input/train/' + image_dir) train_images.append(image) del image gc.collect() train_images = np.array(train_images) </pre>	Load the data	<p>STEP 1: List of images</p> <p>STEP 2: Iterate over all images</p> <p>STEP 2.1: Read in the training data</p> <p>STEP 2.2: Delete unnecessary images</p>

Figure 4: Code and generated documentation examples using GraphCodeBERT in CoDoc of Cell2Doc

Index	Model	Correctness	Informativeness	Readability
1	CodeBERT (CM)	$\mu = 3.00, \sigma = 1.34$	$\mu = 2.8, \sigma = 1.24$	$\mu = 3.65, \sigma = 1.18$
2	CodeBERT (CSM)	$\mu = 2.74, \sigma = 1.39$	$\mu = 2.73, \sigma = 1.29$	$\mu = 3.62, \sigma = 1.14$
3	CodeBERT (ECSM)	$\mu = 2.57, \sigma = 1.31$	$\mu = 2.57, \sigma = 1.28$	$\mu = 3.62, \sigma = 1.11$
4	CodeBERT (SCSCM)	$\mu = 3.10, \sigma = 1.30$	$\mu = 3.04, \sigma = 1.19$	$\mu = 3.72, \sigma = 1.10$
5	CodeBERT (Cell2Doc)	$\mu = \mathbf{3.81}, \sigma = 1.13$	$\mu = \mathbf{4.00}, \sigma = 1.11$	$\mu = \mathbf{4.22}, \sigma = 0.90$

Figure 5: Results of the human evaluation using CodeBERT in CoDoc of Cell2Doc

References

- [1] W. Ahmad, S. Chakraborty, B. Ray, and K.-W. Chang, “Unified pre-training for program understanding and generation,” in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Online: Association for Computational Linguistics, Jun. 2021, pp. 2655–2668. [Online]. Available: <https://aclanthology.org/2021.naacl-main.211>
- [2] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang, and M. Zhou, “CodeBERT: A pre-trained model for programming and natural languages,” in *Findings of the Association for Computational Linguistics: EMNLP 2020*. Online: Association for Computational Linguistics, Nov. 2020, pp. 1536–1547. [Online]. Available: <https://aclanthology.org/2020.findings-emnlp.139>
- [3] D. Guo, S. Lu, N. Duan, Y. Wang, M. Zhou, and J. Yin, “Unixcoder: Unified cross-modal pre-training for code representation,” 2022. [Online]. Available: <https://arxiv.org/abs/2203.03850>
- [4] D. Guo, S. Ren, S. Lu, Z. Feng, D. Tang, S. Liu, L. Zhou, N. Duan, A. Svyatkovskiy, S. Fu *et al.*, “Graph-codebert: Pre-training code representations with data flow,” *arXiv preprint arXiv:2009.08366*, 2020.
- [5] Y. Wang, W. Wang, S. Joty, and S. C. H. Hoi, “Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation,” 2021. [Online]. Available: <https://arxiv.org/abs/2109.00859>
- [6] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*. Philadelphia, Pennsylvania, USA: Association for Computational Linguistics, Jul. 2002, pp. 311–318. [Online]. Available: <https://aclanthology.org/P02-1040>
- [7] C.-Y. Lin, “ROUGE: A package for automatic evaluation of summaries,” in *Text Summarization Branches Out*. Barcelona, Spain: Association for Computational Linguistics, Jul. 2004, pp. 74–81. [Online]. Available: <https://aclanthology.org/W04-1013>