

# PySS3: A new interpretable and simple machine learning model for text classification

Sergio G. Burdisso<sup>\*1, 2</sup>, Marcelo Errecalde<sup>1</sup>, and Manuel Montes-y-Gómez<sup>3</sup>

<sup>1</sup> Universidad Nacional de San Luis (UNSL), Argentina <sup>2</sup> Consejo Nacional de Investigaciones Científicas y Técnicas (CONICET), Argentina <sup>3</sup> Instituto Nacional de Astrofísica, Óptica y Electrónica (INAOE), México

DOI: [10.21105/joss.03934](https://doi.org/10.21105/joss.03934)

## Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Brian McFee](#) ↗

## Reviewers:

- [@hbaniecki](#)
- [@kmichael08](#)

Submitted: 08 November 2021

Published: 19 November 2021

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

## Summary

In this paper, we briefly introduce PySS3<sup>1</sup> and share it with the community. PySS3 is an open-source Python package that implements a novel machine learning model for text classification, called SS3. PySS3 comes with two useful tools that allow working with SS3 in an effortless, interactive, and visual way. For instance, one of these tools provides post hoc explanations using visualizations that directly highlight relevant portions of the raw input document, allowing researchers to better understand the models being deployed. Therefore, PySS3 could be especially useful for those working with sensitive text classification problems by which people's lives could be affected since it allows researchers and practitioners to deploy interpretable (i.e. self-explainable) and more reliable models for text classification.

## Statement of need

A challenging scenario in the machine learning field is the one referred to as *early classification*. Early classification deals with the problem of classifying data streams as early as possible without having a significant loss in performance. The reasons behind this requirement of *earliness* could be diverse, but the most important and interesting case is when the classification delay has negative or risky implications. This scenario, known as *Early Risk Detection* (ERD), has gained increasing interest in recent years with potential applications in rumor detection ([Kwon et al., 2017](#); [Ma et al., 2015, 2016](#)), sexual predator detection, aggressive text identification ([Escalante et al., 2017](#)), depression detection ([Losada et al., 2017](#); [Losada & Crestani, 2016](#)), and terrorism detection ([Iskandar, 2017](#)), among others.

A recently introduced machine learning model for text classification ([Sergio G. Burdisso et al., 2019a, 2020](#)), called SS3, has shown to be well suited to deal with ERD problems on social media streams. It obtained state-of-the-art performance on early depression, anorexia and self-harm detection on the latest CLEF's eRisk lab challenges ([Sergio G. Burdisso et al., 2019b, 2019a, 2020](#); [Loyola et al., 2021](#)). Unlike standard classifiers, this new classification model was specially designed to deal with ERD problems since: (a) it is interpretable and therefore can naturally self-explain its rationale, and (b) it naturally supports incremental training and classification over text streams.

However, little attention has been paid to the potential use of SS3 as a general-purpose classifier for other text classification tasks. One of the main reasons could be the fact that there is

<sup>\*</sup>corresponding author

<sup>1</sup><https://github.com/sergioburdisso/pyss3>

no open-source implementation of SS3 available yet. We believe that the availability of open-source implementations is of critical importance to foster the use of new tools, methods, and algorithms. On the other hand, Python has come to be the most popular programming language in the machine learning community thanks to its simple syntax and a rich ecosystem of efficient open-source implementations of popular algorithms (Oliphant, 2007). Therefore, we decided to develop an open-source Python package to provide the first official implementation of this new classifier.

## Implementation

PySS3 was coded to be compatible with Python 2 and Python 3 as well as with different operating systems, such as Linux, macOS, and Microsoft Windows. To ensure this compatibility holds whenever the source code is updated, we have configured and linked the github repository with the Travis CI service. This service automatically runs the test scripts using different operating systems and versions of Python whenever new code is pushed to the repository.

The package is composed of one main module and three submodules.<sup>2</sup> The main module is called `pyss3` and contains the classifier's implementation *per se* in a class called `SS3`. This class not only implements the “plain-vanilla” version of the classifier (Sergio G. Burdisso et al., 2019a) but also different variations, such as the one introduced later by the same authors (Sergio G. Burdisso et al., 2020), which allows `SS3` to recognize important word n-grams on the fly. As the reader will notice in the example shown in the next section, this class exposes a clear and simple API that is similar to that of *Scikit-learn* models. For instance, it provides standard methods like `fit()` and `predict()` for training and classification, respectively.<sup>3</sup> Finally, the three submodules, `pyss3.server`, `pyss3.cmd_line`, and `pyss3.util`, provide a collection of useful tools and utility functions such as, for instance, those related to data loading, evaluation or, as will be shown in the next section, “live” testing the models.

## Illustrative examples

In this section, we will only introduce two simple illustrative examples. For full and real working examples, please refer to the [tutorials](#) in the documentation. Additionally, for readers interested in trying PySS3 out “right away,” we have created Jupyter Notebooks for the tutorials,<sup>4</sup> which can be executed online, for instance, using [MyBinder](#).

Before introducing the examples, we will assume the user has already loaded the training and test documents and category labels, as usual, in the `x_train`, `y_train`, `x_test`, `y_test` lists, respectively. For instance, this could be done using the `Dataset` class from the `pyss3.util` submodule, as follows:

```
from pyss3.util import Dataset

x_train, y_train = Dataset.load_from_files("path/to/train")
x_test, y_test = Dataset.load_from_files("path/to/test")
```

## Training and test example

This simple example shows how to train and test an `SS3` model using default values. Since `SS3` creates a language model for each category, we do not need to create a document-term

<sup>2</sup>A more detailed description of the package is given in the official documentation (<https://pyss3.rtdf.io>). Additionally, an extended version of this paper can be found in ArXiv (Sergio G. Burdisso et al., 2019).

<sup>3</sup><https://pyss3.rtdf.io/en/latest/api>

<sup>4</sup><https://github.com/sergioburdisso/pyss3/tree/master/examples>

74 matrix, we can simply use the raw `x_train` and `x_test` documents for training and test,  
75 respectively, as follows:

```
from pyss3 import SS3

clf = SS3()
clf.fit(x_train, y_train)
y_pred = clf.predict(x_test)

print("Accuracy:", accuracy(y_pred, y_test))
```

## 76 Training and live test example

77 This example is similar to the previous one. However, instead of simply using `predict` and  
78 `accuracy` to measure our model's performance, here we are using the PySS3's "Live Test"  
79 tool. The "Live Test" is an interactive visualization tool that allows users to test the models  
80 "on the fly" with a single line of extra code, as follows:

```
from pyss3 import SS3
from pyss3.server import Live_Test

clf = SS3()
clf.fit(x_train, y_train)

# The following line will open up the Live Test tool, shown in Figure 1
Live_Test.run(clf, x_test, y_test)
```

81 As shown in [Figure 1](#), the tool provides a user interface by which users can manually and  
82 actively test their model using either the documents in the test set or just typing in their own.  
83 More precisely, this tool allows researchers to analyze and understand what their models are  
84 learning by providing an interactive and visual explanation of the classification process at three  
85 different levels (word n-grams, sentences, and paragraphs).<sup>5</sup>

86 Alternatively to the Live Test tool, we could also use the `clf.extract_insight()` function.  
87 When this function is applied to a document, it returns the list of text fragments involved  
88 in the classification decision, ordered by a confidence value. For instance, suppose we need  
89 to obtain the text fragment with the highest confidence value that was used to classify the  
90 document shown in [Figure 1](#) as "sports," we could use `extract_insight()` as follows:

```
# Assign to "doc" the same document that is shown in Figure 1
doc = "Last year, Moore became Liverpool's CEO. This season, his club ..."

fragments = clf.extract_insight(doc, "sports")

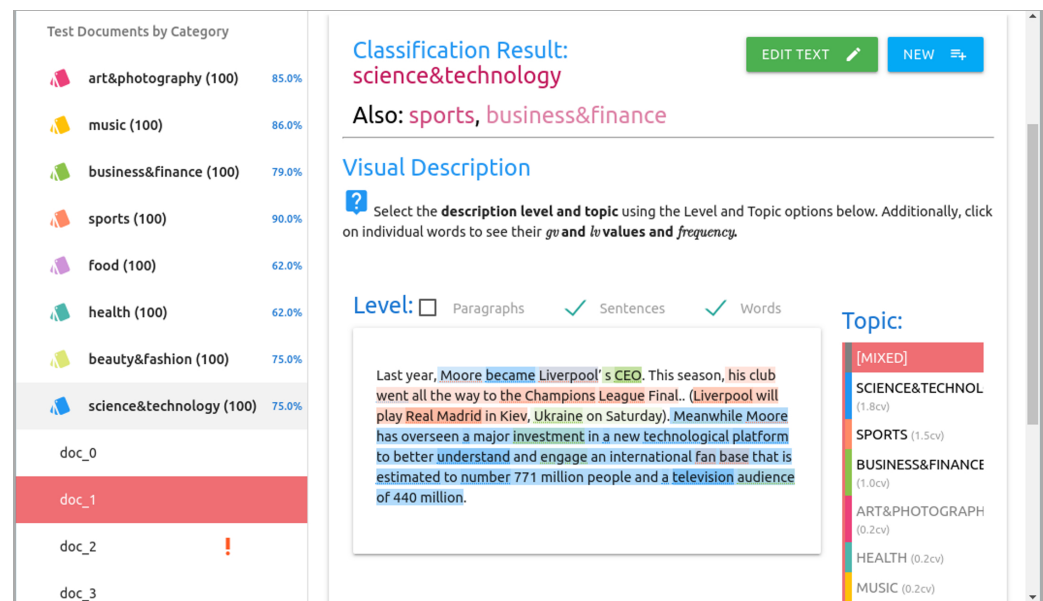
print(fragments[0])
```

91 Which will print the following (fragment, confidence value) pair:

```
('to the Champions League Final. (Liverpool will play Real Madrid in
Kiev, Ukraine on Saturday). Meanwhile Moore has', 0.97)
```

<sup>5</sup>We recommend trying out the "Topic Categorization" or the "Sentiment Analysis on Movie Reviews" online live demo available at <http://tworld.io/ss3>

92 This pair indicates the document was classified as “sports,” mainly due to our model finding  
93 this fragment as belonging to “sports” with a confidence value of 0.97.<sup>6</sup>



**Figure 1:** Live Test screenshot. On the left side, the list of test documents grouped by category is shown along with the percentage of success. Note the doc\_2 document is marked with an exclamation mark (!); this mark indicates it was misclassified, which eases error analysis. The user has selected doc\_1, the “classification result” is shown above the visual description. In this figure, the user has chosen to display the visual explanation at sentence-and-word level, using mixed topics. For instance, the user can confirm that, apparently, the model has learned to recognize important words and that it has correctly classified the document. Also, by using the colors, the user could notice that the first sentence belonging to multiple topics, the second sentence shifted the topic to *sports*, and finally, from “Meanwhile” on, the topic is shifted to *technology* (and a little bit of *business* given by the words “investment” or “engage” colored in green). Note that the user can also edit the document or even create new ones using the two buttons on the upper-right corner.

## 94 Conclusions

95 We have briefly introduced PySS3, an open-source Python package that implements a novel  
96 machine learning model for text classification that comes with useful visualization tools. This  
97 software could be especially useful for researchers and practitioners interested in deploying  
98 interpretable and more reliable models for text classification. Finally, we hope to continue  
99 the advancement and improvement of PySS3 through the direct or indirect help of those in  
100 the mathematical and computer science disciplines who want to be part of this open-source  
101 project.

## 102 References

103 Burdisso, Sergio G., Errecalde, M., & Montes-y-Gómez, M. (2019a). A text classification  
104 framework for simple and effective early depression detection over social media streams.  
105 *Expert Systems with Applications*, 133, 182–197. <https://doi.org/10.1016/j.eswa.2019.05.023>  
106

<sup>6</sup>Interested readers may refer to the “getting the text fragments involved in the classification decision” tutorial [available online](#)

- 107 Burdisso, Sergio G., Errecalde, M., & Montes-y-Gómez, M. (2019). PySS3: A python pack-  
108 age implementing a novel text classifier with visualization tools for explainable AI. *arXiv*  
109 *Preprint arXiv:1912.09322*.
- 110 Burdisso, Sergio G., Errecalde, M., & Montes-y-Gómez, M. (2019b). UNSL at eRisk 2019: A  
111 unified approach for anorexia, self-harm and depression detection in social media. *Working*  
112 *Notes of CLEF 2019*.
- 113 Burdisso, Sergio G., Errecalde, M., & Montes-y-Gómez, M. (2020).  $\tau$ -SS3: A text classifier  
114 with dynamic n-grams for early risk detection over text streams. *Pattern Recognition*  
115 *Letters*, 138, 130–137. <https://doi.org/10.1016/j.patrec.2020.07.001>
- 116 Escalante, H. J., Villatoro-Tello, E., Garza, S. E., López-Monroy, A. P., Montes-y-Gómez,  
117 M., & Villaseñor-Pineda, L. (2017). Early detection of deception and aggressiveness using  
118 profile-based representations. *Expert Systems with Applications*, 89, 99–111. <https://doi.org/10.1016/j.eswa.2017.07.040>
- 120 Iskandar, B. S. (2017). Terrorism detection based on sentiment analysis using machine learn-  
121 ing. *Journal of Engineering and Applied Sciences*, 12(3), 691–698.
- 122 Kwon, S., Cha, M., & Jung, K. (2017). Rumor detection over varying time windows. *PloS*  
123 *One*, 12(1), e0168344. <https://doi.org/10.1371/journal.pone.0168344>
- 124 Losada, D. E., & Crestani, F. (2016). A test collection for research on depression and language  
125 use. *International Conference of the Cross-Language Evaluation Forum for European*  
126 *Languages*, 28–39. [https://doi.org/10.1007/978-3-319-44564-9\\_3](https://doi.org/10.1007/978-3-319-44564-9_3)
- 127 Losada, D. E., Crestani, F., & Parapar, J. (2017). eRisk 2017: CLEF lab on early risk  
128 prediction on the internet: Experimental foundations. *International Conference of the*  
129 *Cross-Language Evaluation Forum for European Languages*, 346–360. [https://doi.org/10.1007/978-3-319-65813-1\\_30](https://doi.org/10.1007/978-3-319-65813-1_30)
- 131 Loyola, J. M., Burdisso, S. G., Thompson, H., Cagnina, L., & Errecalde, M. (2021). UNSL  
132 at eRisk 2021: A comparison of three early alert policies for early risk detection. *Working*  
133 *Notes of CLEF 2021*. <http://ceur-ws.org/Vol-2936/paper-81.pdf>
- 134 Ma, J., Gao, W., Mitra, P., Kwon, S., Jansen, B. J., Wong, K.-F., & Cha, M. (2016).  
135 Detecting rumors from microblogs with recurrent neural networks. *IJCAI*, 3818–3824.
- 136 Ma, J., Gao, W., Wei, Z., Lu, Y., & Wong, K.-F. (2015). Detect rumors using time series  
137 of social context information on microblogging websites. *Proceedings of the 24th ACM*  
138 *International on Conference on Information and Knowledge Management*, 1751–1754.
- 139 Oliphant, T. E. (2007). Python for scientific computing. *Computing in Science Engineering*,  
140 9(3), 10–20.