# EARLY TEXT CLASSIFICATION FOR EARLY DETECTION OF SIGNS OF DEPRESSION

Juan Martín Loyola[1,2*] and Marcelo L. Errecalde[2]

1 IMASL-CONICET, Italia 1556, San Luis, Argentina.
2 LIDIC, Ejército de los Andes 950, San Luis, Argentina.
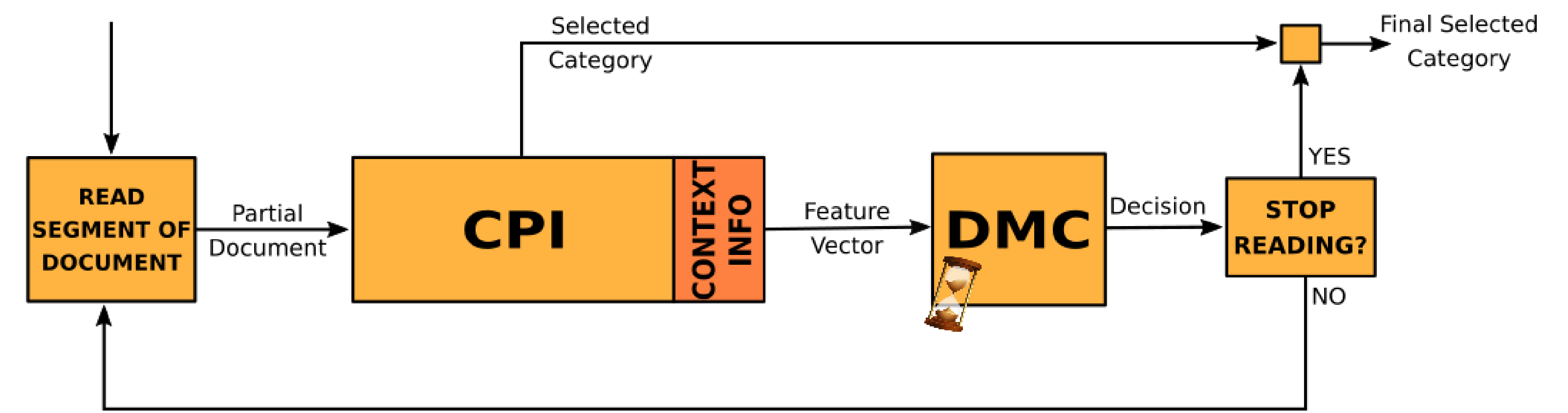
* jmloyola@unsl.edu.ar

## INTRODUCTION

The problem of classification in supervised learning is a widely studied one. Nonetheless, there are scenarios that received little attention despite its applicability. One of such scenarios is **early text classification**, which deals with the development of predictive models that can determine the class a document belongs to as soon as possible. Here a document is assumed to be processed sequentially, starting at the beginning and reading its containing parts one by one. In this context, it is desired to make predictions with as little information (as soon) as possible. The importance of this variant of the classification problem is evident in tasks like sexual predator detection, where one wants to identify an offender as early as possible. [1]

It is important to note that the early text classification problem consists of two related and complementary tasks. On the one hand, the task of **classification with partial information** (CPI), which consists of obtaining an efficient predictive model when only partial information is available that has been read sequentially up to a certain point in time. Here, the emphasis is to determine which classification methods are more likely to achieve performance comparable to that obtained when classified using the entire document. On the other hand, we have the task of **decision of the moment of classification** (DMC), that is, in which point in time one can stop reading and classify with some degree of confidence that the prediction is going to be correct. [2]

In this work, we apply this framework to the early detection of signs of depression in users in an online forum [3].
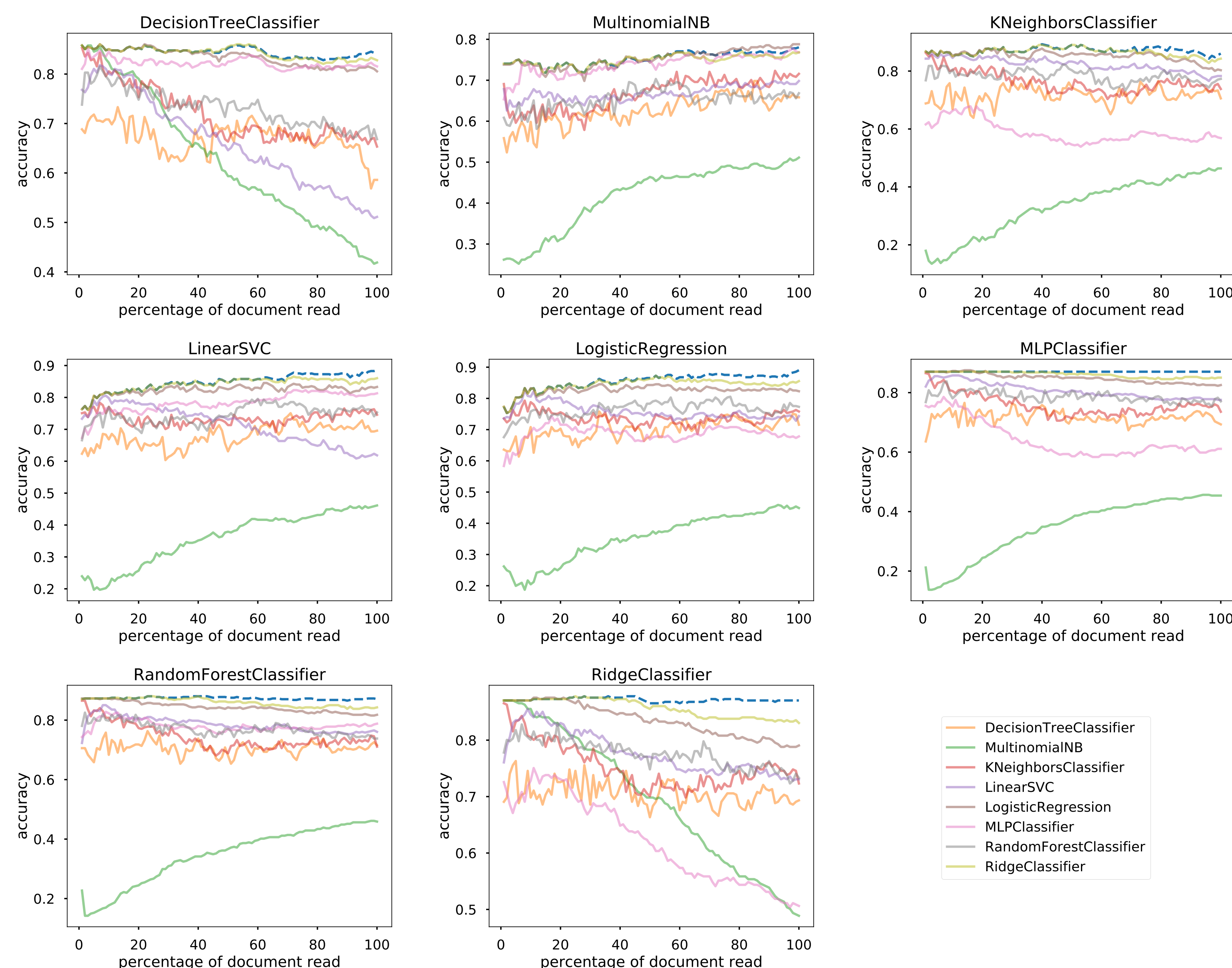
## METHOD



## EVALUATION METRIC

$$\text{EDE}_o(d, k) = \begin{cases} lc_o(k) \cdot c_{\text{tp}}^i & \text{if the decision } d_i \text{ is correctly positive} \\ c_{\text{fn}}^j + c_{\text{fp}}^i & \text{if the decision } d_j \text{ is incorrectly negative} \\ & \text{and if the decision } d_i \text{ is incorrectly positive} \end{cases}$$

where $d$ represents the decision made for all the categories, $d_i$ the decision on category $i$ and $k$ the time when the decision is made. Constants $c_{\text{fp}}^i$, $c_{\text{fn}}^i$ and $c_{\text{tp}}^i$ indicate the cost associated with the decision on the category being false positive, false negative or true positive, respectively. The values given to these constants depend on the particular addressed problem. The factor $lc_o(k) \in [0, 1]$ encodes the cost associated to the delay in detecting true positives. [2]

## RESULTS

### Model comparison



### Results of the model comparison

| CPI Model | DMC Model | Precision | Recall | F1 Score | Accuracy | EDE $o = 5$ Proportional | EDE $o = 50$ Proportional |
|---|---|---|---|---|---|---|---|
| MultinomialNB | DecisionTreeClassifier | 0.608 | 0.693 | 0.617 | 0.751 | 0.087 | 0.076 |
| MultinomialNB | KNeighborsClassifier | 0.601 | 0.677 | 0.610 | 0.751 | 0.087 | 0.080 |
| MultinomialNB | RandomForestClassifier | 0.601 | 0.671 | 0.610 | 0.756 | 0.090 | 0.082 |
| MultinomialNB | RidgeClassifier | 0.589 | 0.655 | 0.594 | 0.741 | 0.087 | 0.086 |
| MultinomialNB | LogisticRegression | 0.589 | 0.655 | 0.594 | 0.741 | 0.089 | 0.088 |
| MultinomialNB | MLPClassifier | 0.585 | 0.646 | 0.590 | 0.741 | 0.107 | 0.092 |
| LogisticRegression | DecisionTreeClassifier | 0.567 | 0.582 | 0.573 | 0.786 | 0.108 | 0.106 |
| MultinomialNB | LinearSVC | 0.607 | 0.682 | 0.618 | 0.761 | 0.111 | 0.106 |
| LogisticRegression | RandomForestClassifier | 0.535 | 0.538 | 0.536 | 0.781 | 0.118 | 0.117 |
| LogisticRegression | LogisticRegression | 0.530 | 0.534 | 0.531 | 0.773 | 0.119 | 0.118 |
| LogisticRegression | RidgeClassifier | 0.530 | 0.534 | 0.531 | 0.773 | 0.119 | 0.118 |
| LogisticRegression | MLPClassifier | 0.553 | 0.571 | 0.557 | 0.766 | 0.129 | 0.119 |
| LinearSVC | RandomForestClassifier | 0.523 | 0.526 | 0.524 | 0.773 | 0.121 | 0.121 |
| LogisticRegression | KNeighborsClassifier | 0.517 | 0.520 | 0.518 | 0.763 | 0.122 | 0.122 |
| KNeighborsClassifier | DecisionTreeClassifier | 0.738 | 0.526 | 0.518 | 0.873 | 0.124 | 0.123 |

### Temporal model against linear model

| Type of Model | Precision | Recall | F1 Score | Accuracy | EDE $o = 50$ Proportional |
|---|---|---|---|---|---|
| Temporal | 0.608 | 0.693 | 0.617 | 0.751 | 0.076 |
| Linear | 0.747 | 0.770 | 0.758 | 0.885 | 0.138 |

## FUTURE WORK

1. Adapt the framework to read chunks of posts so we can compare our results with those reported in the *erisk* task [4].
2. Use a different document representation for the CPI model, for example the TVT [5].
3. Augment the contextual information of the DMC model with more informative features, for instance use the words with highest information gain as relevant words or use external information like a depression lexicon.

## REFERENCES

[1] H. J. Escalante, M. Montes-y-Gómez, L. V. Pineda, and M. L. Errecalde, "Early text classification: a naïve solution," in *Proceedings of the 7th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis, WASSA@NAACL-HLT 2016, June 16, 2016, San Diego, California, USA*, pp. 91–99, 2016.

[2] J. M. Loyola, M. L. Errecalde, H. J. Escalante, and M. Montes y Gomez, "Learning when to classify for early text classification," in *Computer Science – CACIC 2017* (A. E. De Giusti, ed.), (Cham), pp. 24–34, Springer International Publishing, 2018.

[3] D. E. Losada and F. Crestani, "A test collection for research on depression and language use," in *Experimental IR Meets Multilinguality, Multimodality, and Interaction* (N. Fuhr, P. Quaresma, T. Gonçalves, B. Larsen, K. Balog, C. Macdonald, L. Cappellato, and N. Ferro, eds.), (Cham), pp. 28–39, Springer International Publishing, 2016.

[4] D. E. Losada, F. Crestani, and J. Parapar, "erisk 2017: Clef lab on early risk prediction on the internet: Experimental foundations," in *Experimental IR Meets Multilinguality, Multimodality, and Interaction* (G. J. Jones, S. Lawless, J. Gonzalo, L. Kelly, L. Goeuriot, T. Mandl, L. Cappellato, and N. Ferro, eds.), (Cham), pp. 346–360, Springer International Publishing, 2017.

[5] M. L. Errecalde, M. P. Villegas, D. G. Funez, M. J. G. Ucelay, and L. C. Cagnina, "Temporal variation of terms as concept space for early risk prediction," in *Working Notes of CLEF 2017 - Conference and Labs of the Evaluation Forum, Dublin, Ireland, September 11-14, 2017.*, 2017.

## SOURCE CODE

The source code of the early text classification framework and the jupyter notebook that produces this results are available at:
https://github.com/jmloyola/early-classification

## ACKNOWLEDGEMENTS