



# CLASSIFICAÇÃO DE SUICIDALIDADE EM UMA VASTA COORTE OCUPACIONAL: UMA ANÁLISE DE ALGORITMOS DE APRENDIZADO DE MÁQUINA APLICADOS AO ESTUDO ELSA-BRASIL

Gabriel de Souza Seibel

Trabalho de Graduação — Dezembro de 2020

## AGENDA

1. Introdução
2. Trabalhos Relacionados
3. Metodologia
4. Experimentos
5. Resultados
6. Conclusões

## Introdução

---

- Mais de 800.000 por ano <sup>1</sup>
- 1 a cada 40 segundos <sup>1</sup>
- Segunda maior causa de mortes entre pessoas de 15 a 29 anos <sup>2</sup>
- 80% acontecem em países em desenvolvimento <sup>2</sup>

1: REID (2010), 2: WHO (2017)

- Ideação como vulnerabilidade
- Suicidalidade e intenção

Suicidalidade = Ideação<sub>e/ou</sub> Taedium Vitae<sub>e/ou</sub> Desesperança

Nos últimos 7 dias:

- Se sentiu completamente sem esperança, por exemplo, em relação ao seu futuro?
- Sentiu que não vale a pena viver?
- Pensou em se matar?

Em um população de  
**brasileiros adultos com transtorno mental comum,**

criar modelos de  
**classificação de suicidalidade** com alto desempenho

e extrair **padrões e fatores** socioeconômicos, biológicos e  
comportamentais indicadores de suicidalidade.

- Mitigar desbalanço de classes
- Obter conjunto pequeno de atributos
- Testar diferentes classificadores e abordagem ensemble



**Trabalhos Relacionados**

---

Figura: Artigo de revisão sistemática da literatura (Burke, 2019)

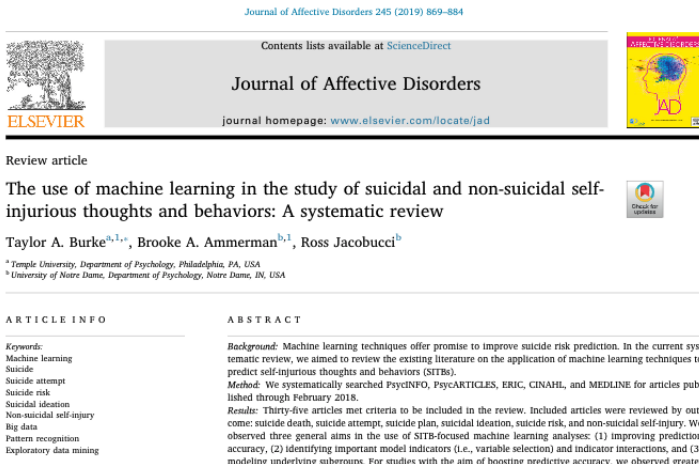
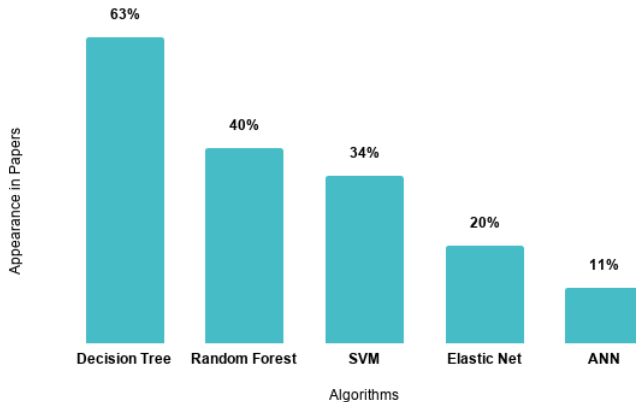


Figura: Resumo de prevalência de algoritmos na revisão de Burke (2019)



Fonte: Autor

- Modelos simples - interpretabilidade
- Tratamento de desbalanço de classes
- Métricas insuficientes

| <i>Paper</i> | <i>Algorithm</i> | <i>F<sub>2</sub>-Score</i> | <i>AUCROC</i> | <i>Sensitivity</i> | <i>Specificity</i> |
|--------------|------------------|----------------------------|---------------|--------------------|--------------------|
| A            | XGB              | 0.84                       | 0.86          | 0.79               | 0.79               |
| B            | ANNs/RF          | 0.71                       | 0.88          | 0.80               | 0.79               |
| C            | ANN              | 0.48                       | 0.88          | 0.81               | 0.77               |
| D            | EN               | 0.45                       | 0.79          | 0.67               | 0.78               |
| E            | RFs              |                            | 0.98          |                    |                    |
| F            | RF               |                            | 0.92          |                    |                    |
| G            | SVM              |                            |               | 0.77               | 0.79               |

A: JUNG et al. (2019);

B: ROY et al. (2020);

C: OH et al. (2020);

D: LIBRENZA-GARCIA et al. (2020);

E: SCHUBACH et al. (2017);

F: GRADUS et al. (2017);

G: BARROS et al. (2017).

## Metodologia

---

Figura: Modelagem de dados e classificação

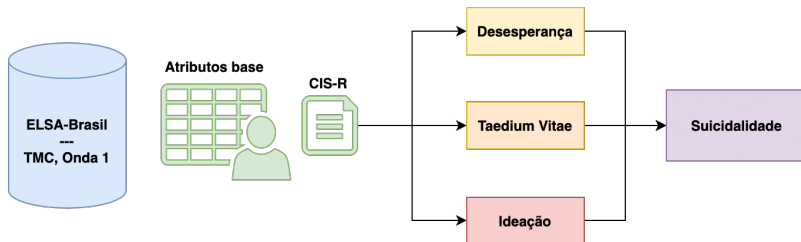


Tabela: Número de variáveis envolvidas em etapas da limpeza dos dados

| <i>Conjunto de Atributos</i>        | <i>Tamanho do Conjunto</i> |
|-------------------------------------|----------------------------|
| Total (dados brutos)                | 2463 (100%)                |
| Removidos (vazamento de informação) | 13 (0.69%)                 |
| Removidos (valores faltantes)       | 773 (31.38%)               |
| Removidos (texto livre)             | 47 (1.91%)                 |
| <b>Restante (dados limpos)</b>      | <b>1626 (66.02%)</b>       |



Tabela: Principais características do conjunto de dados limpo

| <i>Característica</i> | <i>Valor</i>  |
|-----------------------|---------------|
| #Instâncias           | 4039          |
| #Atributos            | 1626          |
| #Positivos            | 1120 (27.73%) |
| #Negativos            | 2919 (72.27%) |

- Downsampling
- Atribuição de valores faltantes
- Corte por variância quase nula
- Filtragem de correlações altas
- SMOTE - Synthetic Minority Oversampling Technique

- Elastic Nets
- Redes Neurais
- Florestas Aleatórias

Figura: Pseudocódigo do algoritmo RFE

---

**Algorithm 2:** Recursive feature elimination incorporating resampling

---

```
2.1 for Each Resampling Iteration do
2.2   Partition data into training and test/hold-back set via resampling
2.3   Tune/train the model on the training set using all predictors
2.4   Predict the held-back samples
2.5   Calculate variable importance or rankings
2.6   for Each subset size  $S_i$ ,  $i = 1 \dots S$  do
2.7     Keep the  $S_i$  most important variables
2.8     [Optional] Pre-process the data
2.9     Tune/train the model on the training set using  $S_i$  predictors
2.10    Predict the held-back samples
2.11    [Optional] Recalculate the rankings for each predictor
2.12  end
2.13 end
2.14 Calculate the performance profile over the  $S_i$  using the held-back samples
2.15 Determine the appropriate number of predictors
2.16 Estimate the final list of predictors to keep in the final model
2.17 Fit the final model based on the optimal  $S_i$  using the original training set
```

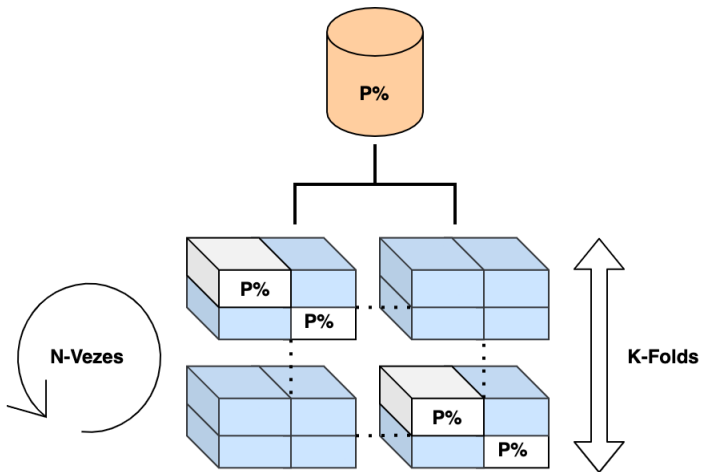
---

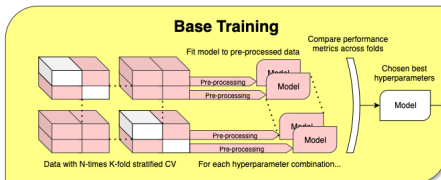
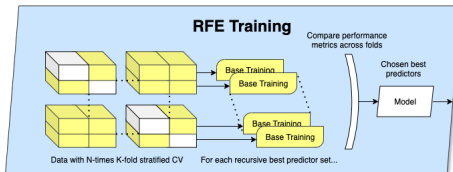
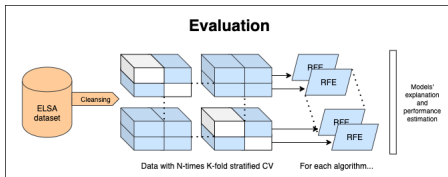
- Precisão ( $Pr$ ) e Recall ( $Re$ )
- Sensibilidade e Especificidade
- Área sob a curva ROC ( $AUCROC$ )
- $F_2$ -Score

$$F_2 = \frac{5 * TP}{5 * TP + 4 * FN + FP} \quad (1)$$

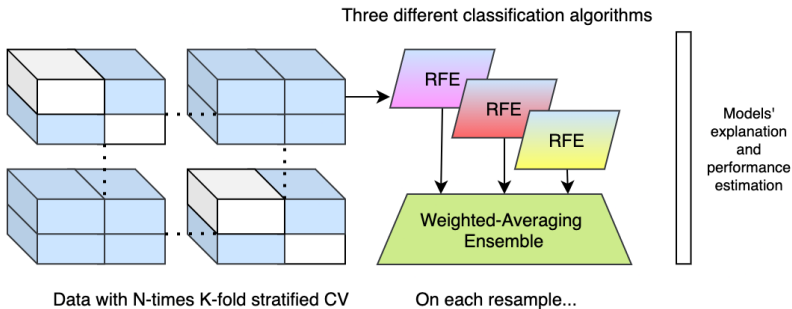
$$F_2 = \frac{5 * Pr * Re}{4 * Pr + Re} \quad (2)$$

Figura: Validação cruzada repetida e estratificada





## Evaluation Layer - Ensemble Constitution





## Experimentos

---

| <i>Parameter</i>              | <i>Value</i>        |
|-------------------------------|---------------------|
| CV Avaliação - K (folds)      | 10                  |
| CV Avaliação - N (vezes)      | 3                   |
| CV Treino RFE - K (folds)     | 5                   |
| CV Treino RFE - N (vezes)     | 2                   |
| CV Treino Base - K (folds)    | 5                   |
| CV Treino Base - N (vezes)    | 2                   |
| Downsampling - Taxa positivos | 33.3%               |
| SMOTE - Taxa positivos        | 50%                 |
| RFE - Número de atributos     | $(2^k)_{k=3}^{k=9}$ |

- Busca em grade

| <i>Parâmetro</i>                | <i>Valores</i>                         |
|---------------------------------|--|
| Elastic Net - Alpha             | 0.1 , 0.325 , 0.550 , 0.775 , 1        |
| Elastic Net - Lambda            | 2e-4 , 9.2e-4 , 4.3e-3 , 2e-2 , 9.2e-2 |
| Rede Neural - Camada 1          | 1 , 2 , 3 , 4 , 5                      |
| Rede Neural - Camada 2          | 0 , 1 , 2 , 3 , 4                      |
| Floresta Aleatória - Atributos  | 2, 17, 33, 48, 64                      |
| Floresta Aleatória - Node-split | <i>gini, extratrees</i>                |
| Ensemble - Pesos dos modelos    | 1/3                                    |

## Linguagem *R*

Pacotes principais:

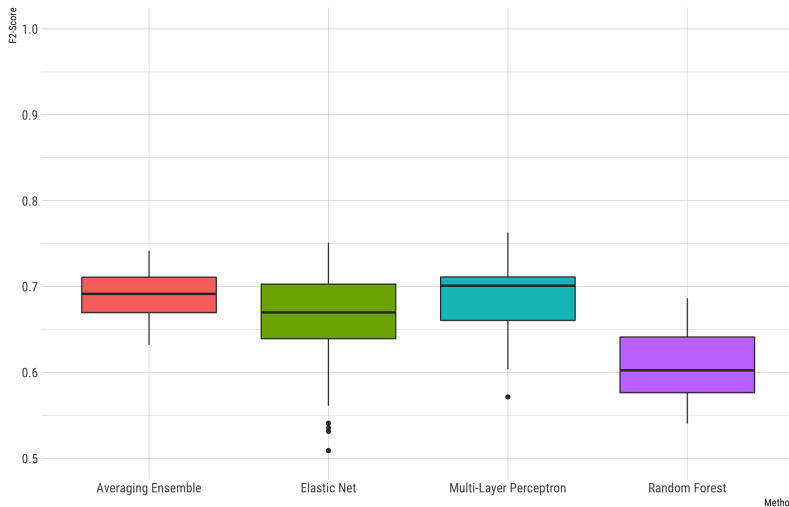
- *caret*
- *recipes*
- *dplyr*
- *purrr*
- *ggplot2*

## Resultados

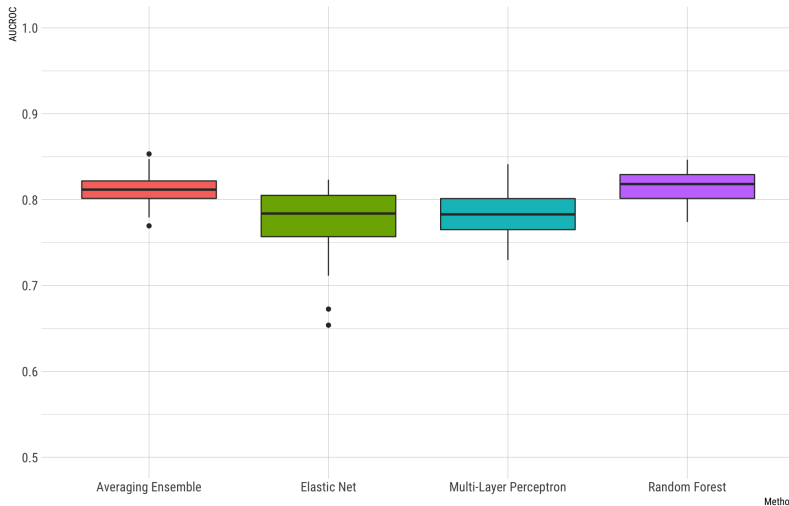
---

Tabela: Médias e desvios padrão de estimativas de desempenho

| <i>Algoritmo</i> | <i>F<sub>2</sub>-Score</i> | <i>AUCROC</i> | <i>Sens.</i> | <i>Espe.</i> |
|------------------|----------------------------|---------------|--------------|--------------|
| Ensemble         | 0.69 ± 0.03                | 0.81 ± 0.02   | 0.78 ± 0.05  | 0.67 ± 0.05  |
| R. Neurais       | 0.69 ± 0.04                | 0.76 ± 0.08   | 0.81 ± 0.09  | 0.59 ± 0.17  |
| Elastic N.       | 0.66 ± 0.07                | 0.77 ± 0.04   | 0.75 ± 0.11  | 0.66 ± 0.09  |
| Florestas A.     | 0.61 ± 0.04                | 0.81 ± 0.02   | 0.63 ± 0.05  | 0.79 ± 0.03  |

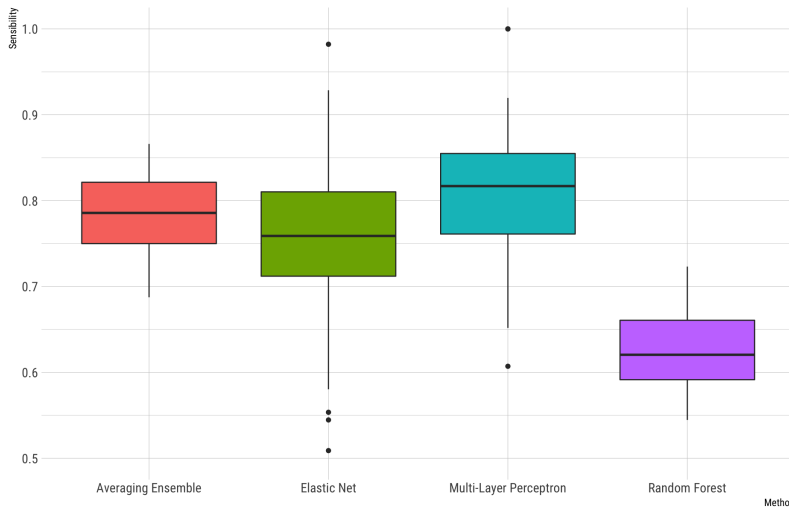
**F2-Score per method**

AUCROC per method

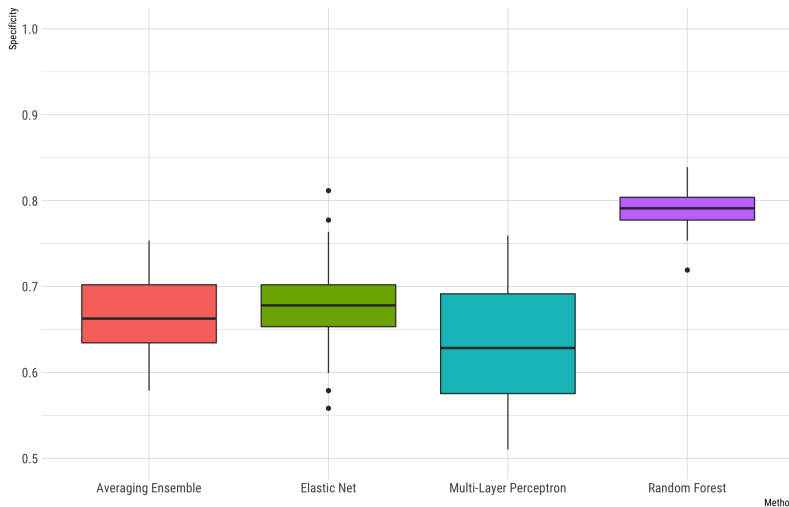




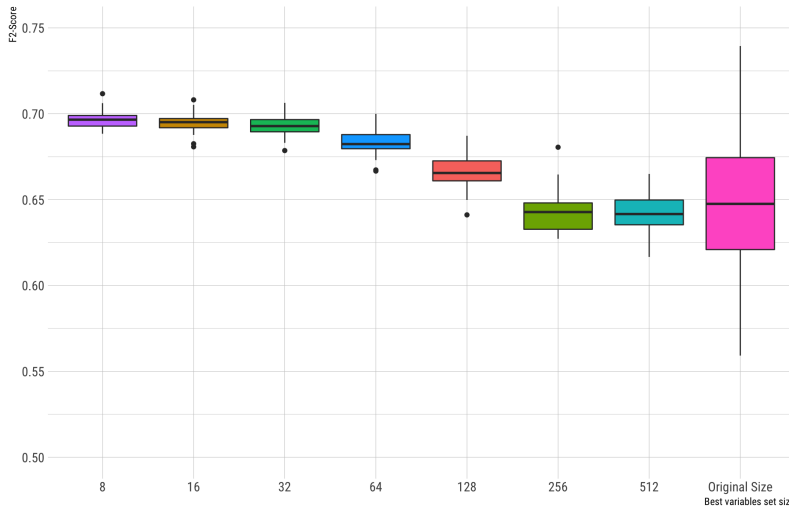
Sensibility per method

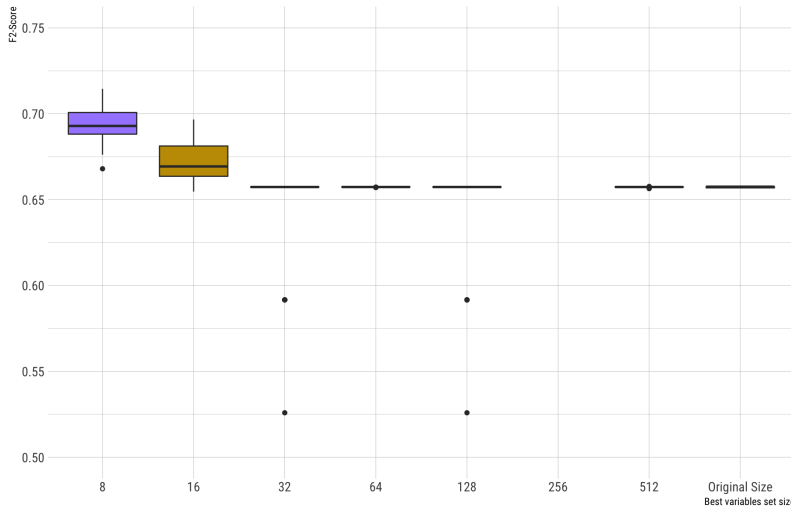


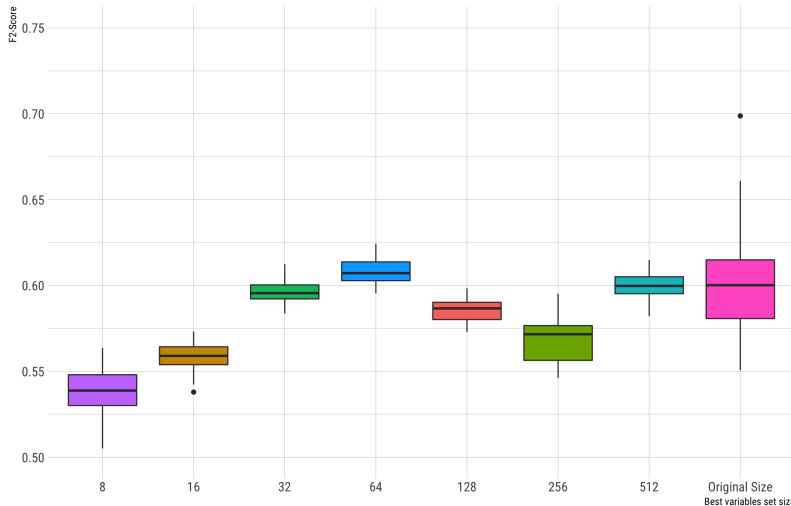
Specificity per method



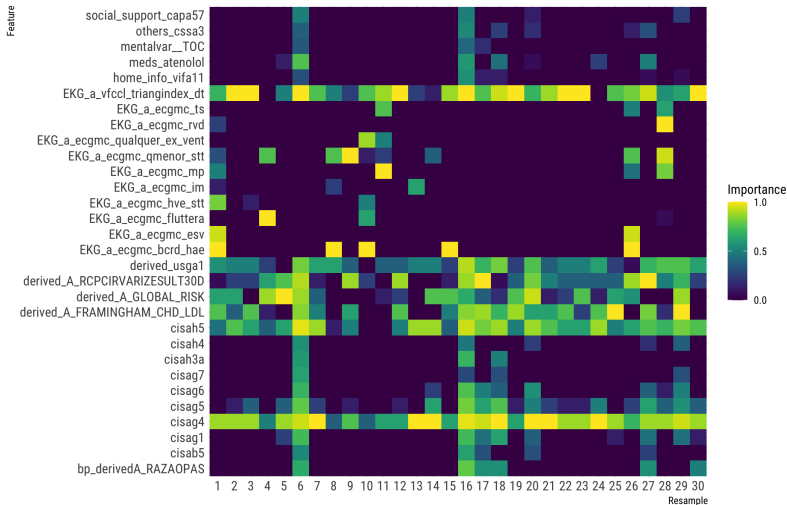
F2-Score per RFE size - Elastic Net



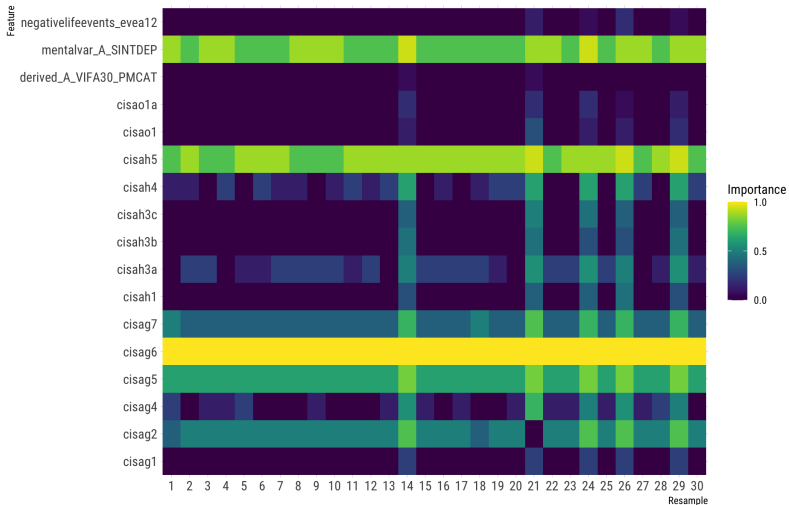
**F2-Score per RFE size - Multi-Layer Perceptron**

**F2-Score per RFE size - Random Forest**

Feature importance per resample - Elastic Net

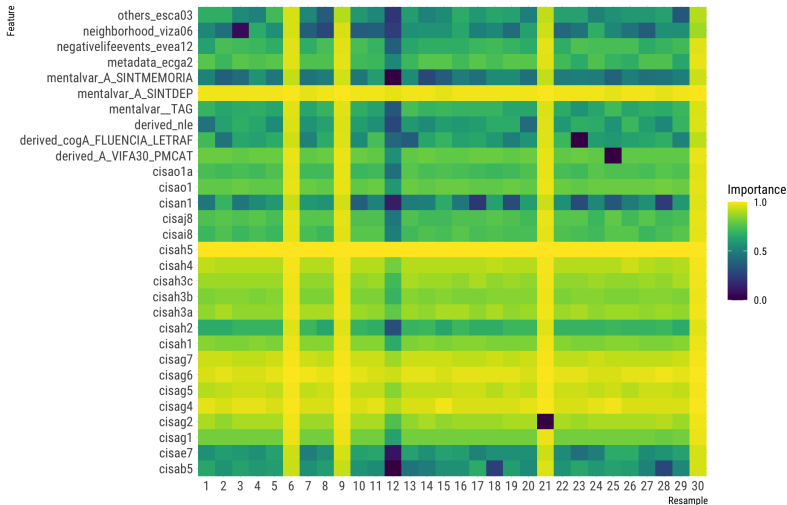


Feature importance per resample - Multi-Layer Perceptron



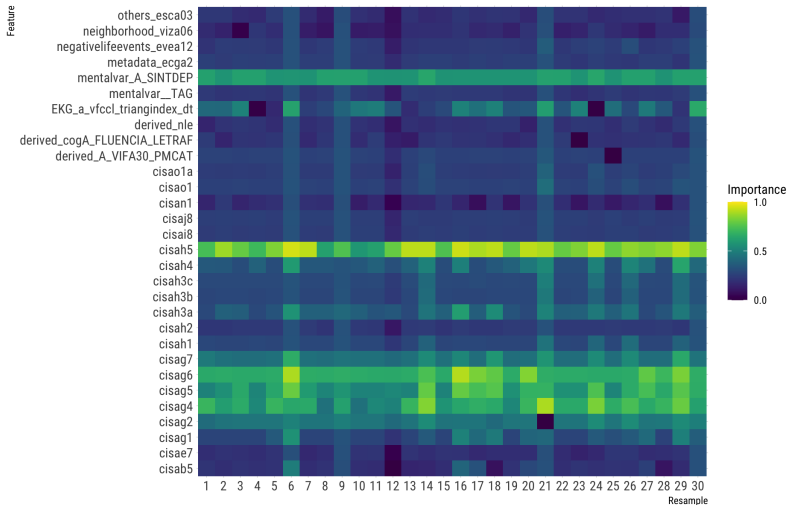
MAPA DE ATRIBUTOS - FLORESTAS  
ALEATÓRIASCLASSIFICAÇÃO DE  
SUICIDALIDADE EM  
ADULTOS BRASILEIROS

Feature importance per resample - Random Forest

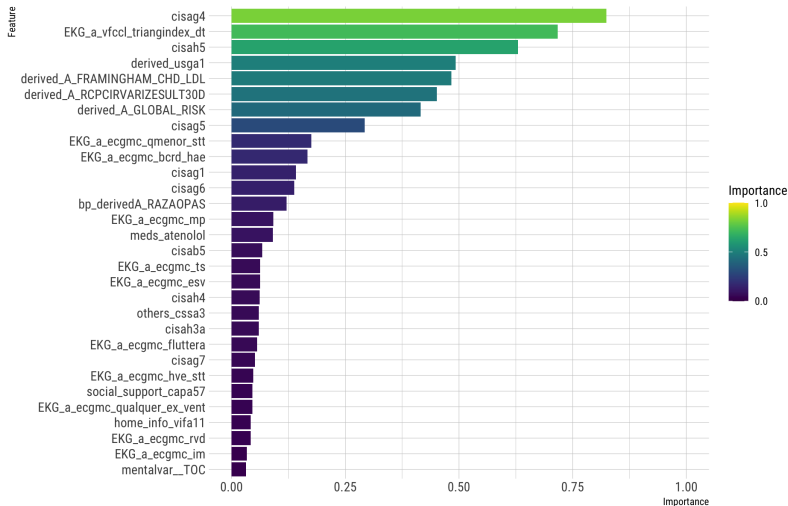




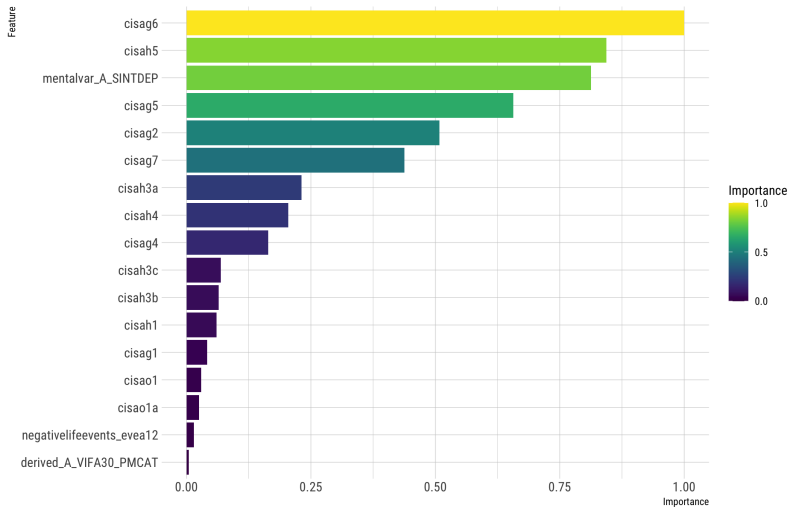
Feature importance per resample - Averaging Ensemble



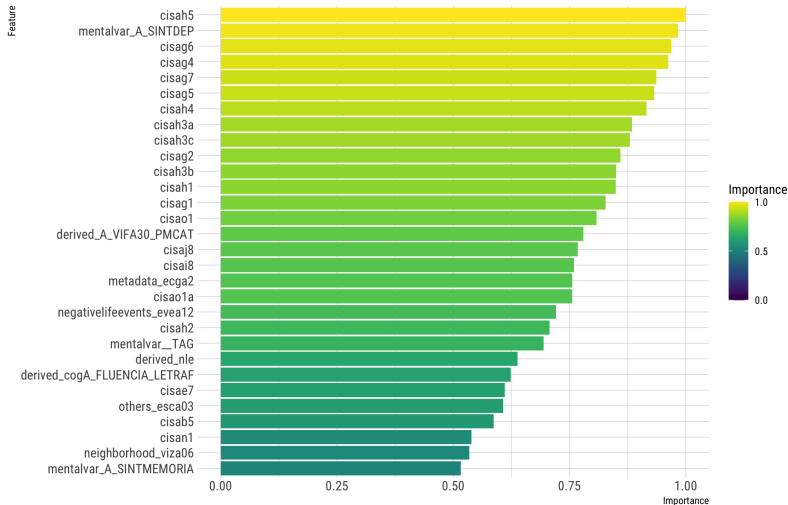
Feature importance overall - Elastic Net



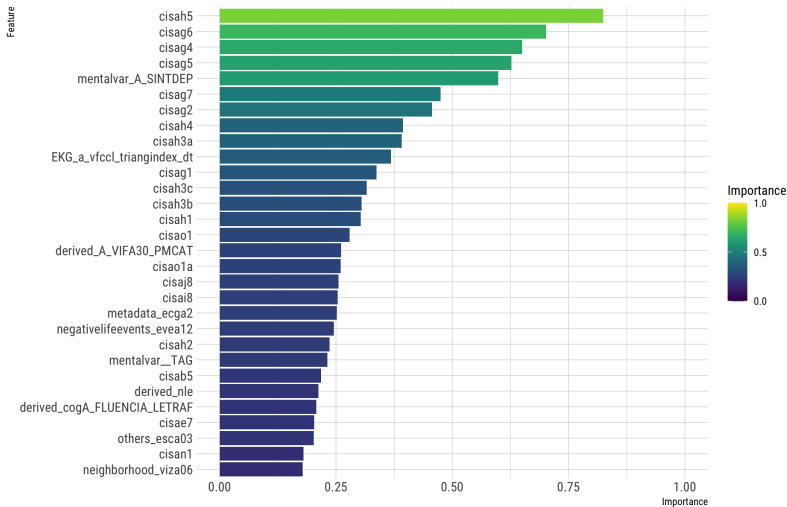
Feature importance overall - Multi-Layer Perceptron



Feature importance overall - Random Forest



Feature importance overall - Averaging Ensemble



1. Sentimento de inferioridade
2. Tristeza
3. Desaparecimento de interesses
4. Auto-culpa desnecessária
5. Energia (disposição)
6. Incapacidade de realizar atividades
7. Renda
8. Ansiedade
9. Preocupação
10. Libido
11. Irritabilidade
12. Obsessão
13. Atividades físicas

## Conclusões

---

Classificação de suicidalidade com o ELSA-Brasil

Relevância de variáveis para classificação

Metodologia (desbalanço de classes, RFE, etc.)



Tabela: Estimativas de desempenho - comparação com trabalhos similares

| <i>Paper</i> | <i>Algorithm</i> | <i>F<sub>2</sub>-Score</i> | <i>AUCROC</i> | <i>Sens.</i> | <i>Espec.</i> |
|--------------|------------------|----------------------------|---------------|--------------|---------------|
| A            | XGB              | 0.84                       | 0.86          | 0.79         | 0.79          |
| B            | ANNs+RF          | 0.71                       | 0.88          | 0.80         | 0.79          |
| <b>Ours</b>  | <b>EN/ANN/RF</b> | <b>0.69</b>                | <b>0.81</b>   | <b>0.78</b>  | <b>0.67</b>   |
| C            | ANN              | 0.48                       | 0.88          | 0.81         | 0.77          |
| D            | EN               | 0.45                       | 0.79          | 0.67         | 0.78          |

A: JUNG et al. (2019);

B: ROY et al. (2020);

C: OH et al. (2020);

D: LIBRENZA-GARCIA et al. (2020).

Análise de variação de valores dos atributos

Explorar mais os dados do ELSA-Brasil

Estudo com foco clínico

Aplicações de assistência e suporte clínicos

## AGREDECIMENTOS ESPECIAIS

André Russowsky Brunoni (USP)

Ives Cavalcante Passos (UFRGS)

Mariana Recamonde Mendoza (UFRGS)

OBRIGADO!

CLASSIFICAÇÃO DE  
SUICIDALIDADE EM  
ADULTOS BRASILEIROS

**Gabriel de Souza Seibel**

Instituto de Informática — UFRGS

`inf.ufrgs.br/~gsseibel`

