

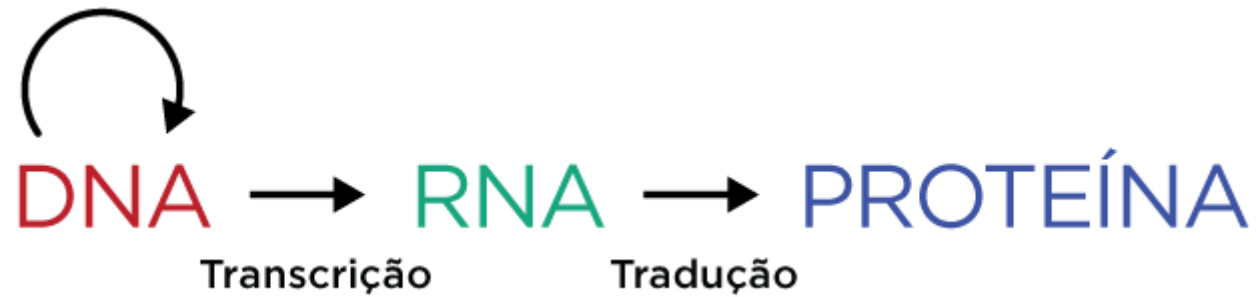
Projeto III – Introdução à modelagem de Big Data

Classificação de essencialidade gênica em organismos modelo

Iara Souza

Introdução

- Dogma central da biologia



Genes —————> Função

Genes essenciais

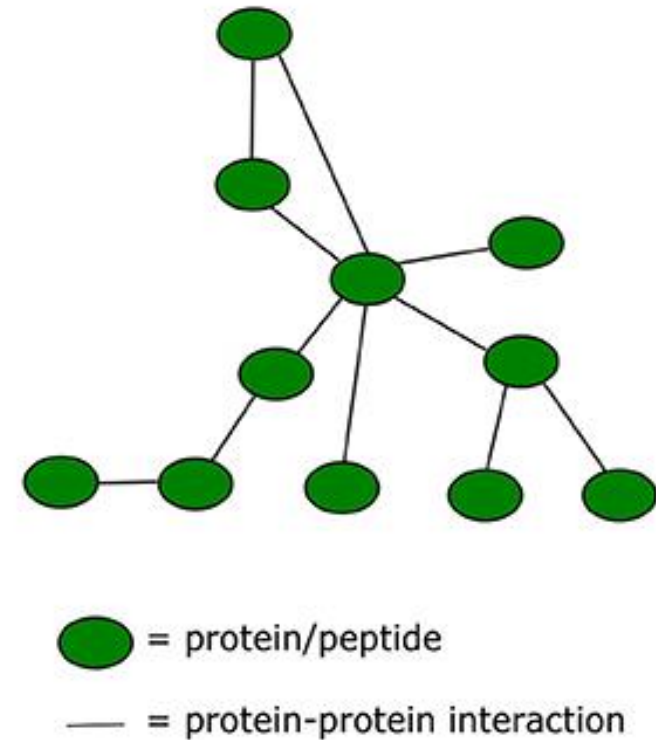
Genes essenciais $\xrightarrow{\text{Alteração deletéria}}$ Morte do organismo

Genes não-essenciais $\xrightarrow{\text{Alteração deletéria}}$ Viabilidade do organismo

Genes essenciais

- Genes essenciais são, em geral:
 - Mais antigos;
 - Interagem mais com outros genes;
 - Participam de funções cruciais na célula

Protein-protein interaction network



Objetivo

- Classificar os genes de dois organismos modelo (camundongo e levedura) em essenciais ou não-essenciais com base em dados evolutivos e de propriedades de redes biológicas.

Coleta de dados

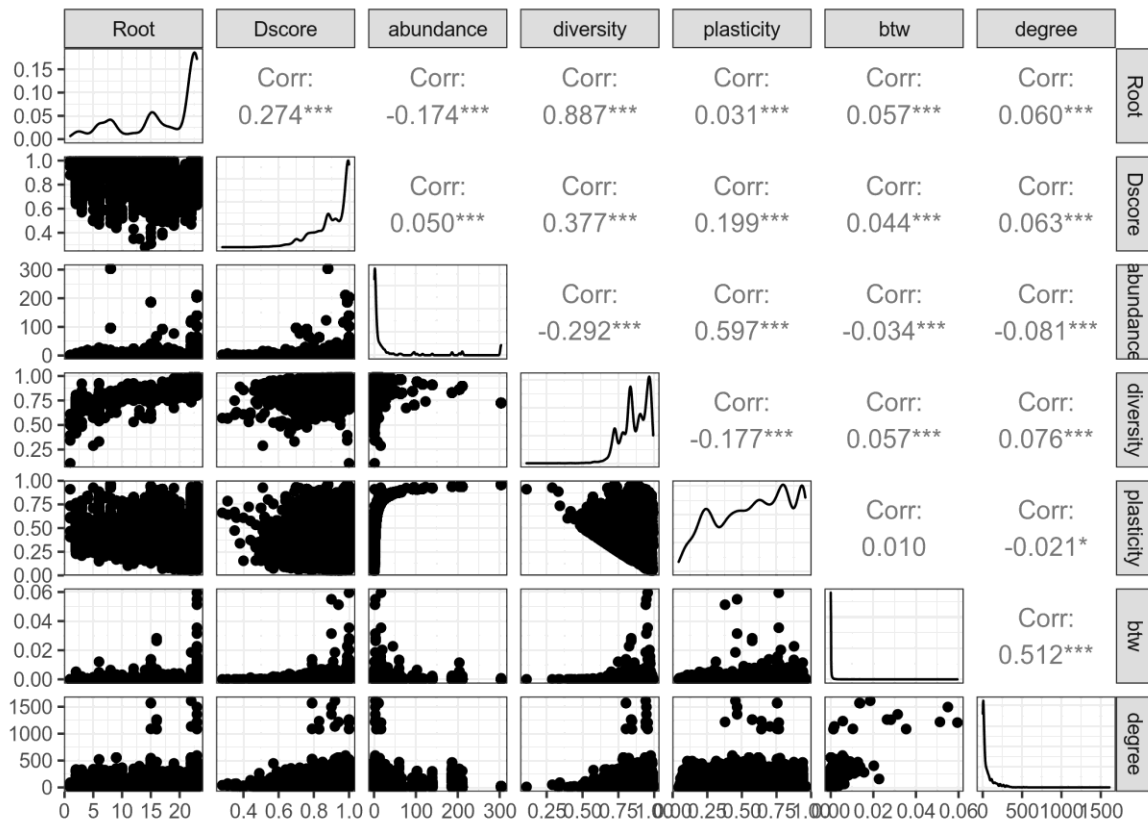
- A obtenção dos dados foi feita em múltiplas etapas:
 - 1. Anotação funcional dos genes
 - 2. Inferência das raízes evolutivas
 - 3. Construção das redes de interação proteína-proteína para cada organismo

Variáveis do conjunto de dados

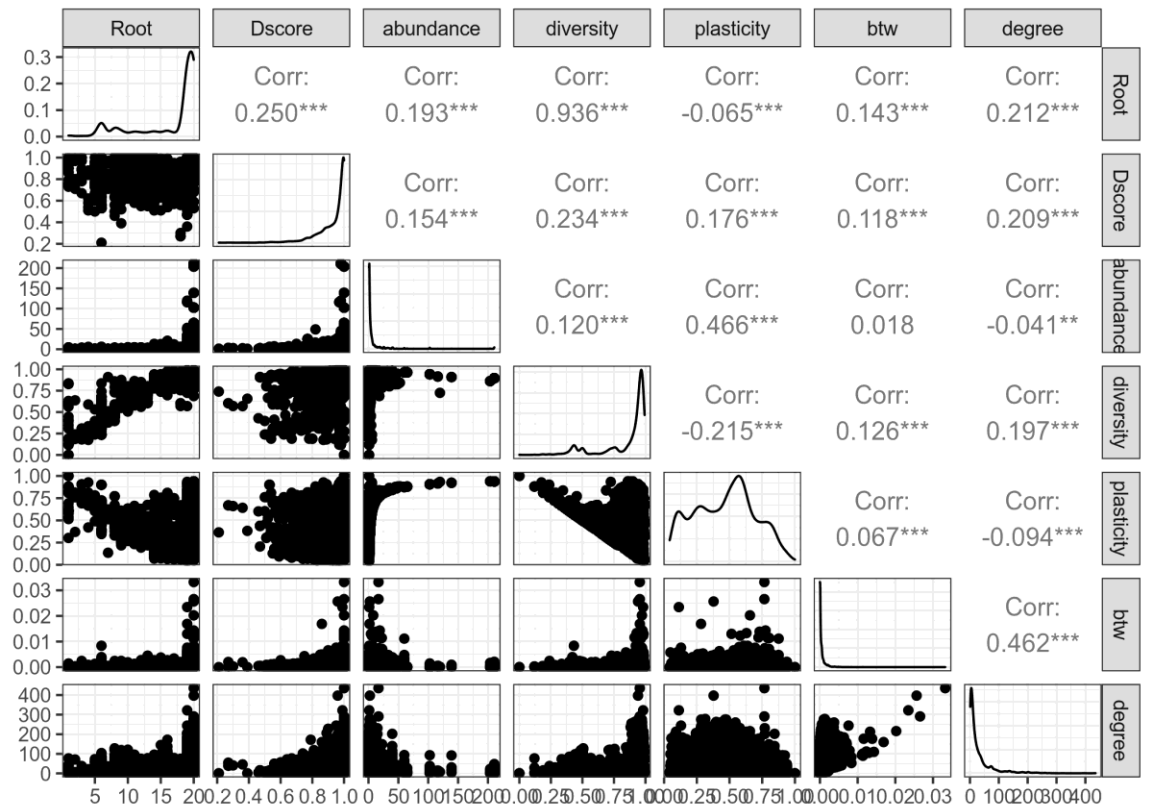
- `ensembl_peptide_id`: identificador da proteína (gene);
- `cog_id`: identificador do grupo de ortólogos ao qual o gene pertence;
- `Root`: raiz evolutiva inferida. Quanto maior, mais antigo;
- `Dscore`: escore de consistência para a inferência da raiz;
- `Pvalue` e `AdjPvalue`: p-valor e p-valor ajustado para a inferência da raiz evolutiva;
- `abundance`: medida de abundância dos grupos de ortólogos;
- `diversity`: medida de diversidade dos grupos de ortólogos;
- `plasticity`: índice de plasticidade evolutiva;
- `ancestry`: ancestralidade do gene. Transformação da variável `Root` (escala de 0 - mais novo - a 1 - mais antigo);
- `btw`: centralidade do nó na rede;
- `degree`: grau do nó na rede.

Análise exploratória

Camundongo

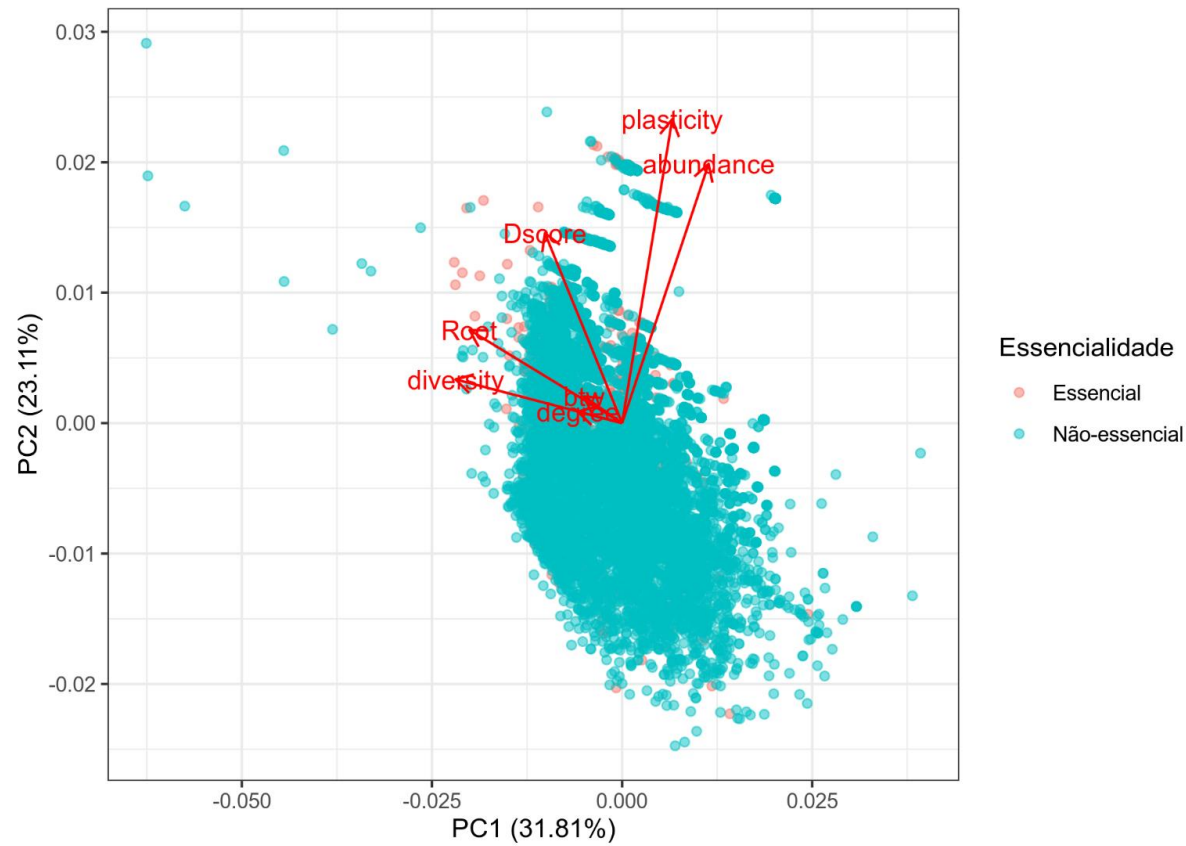


Levedura

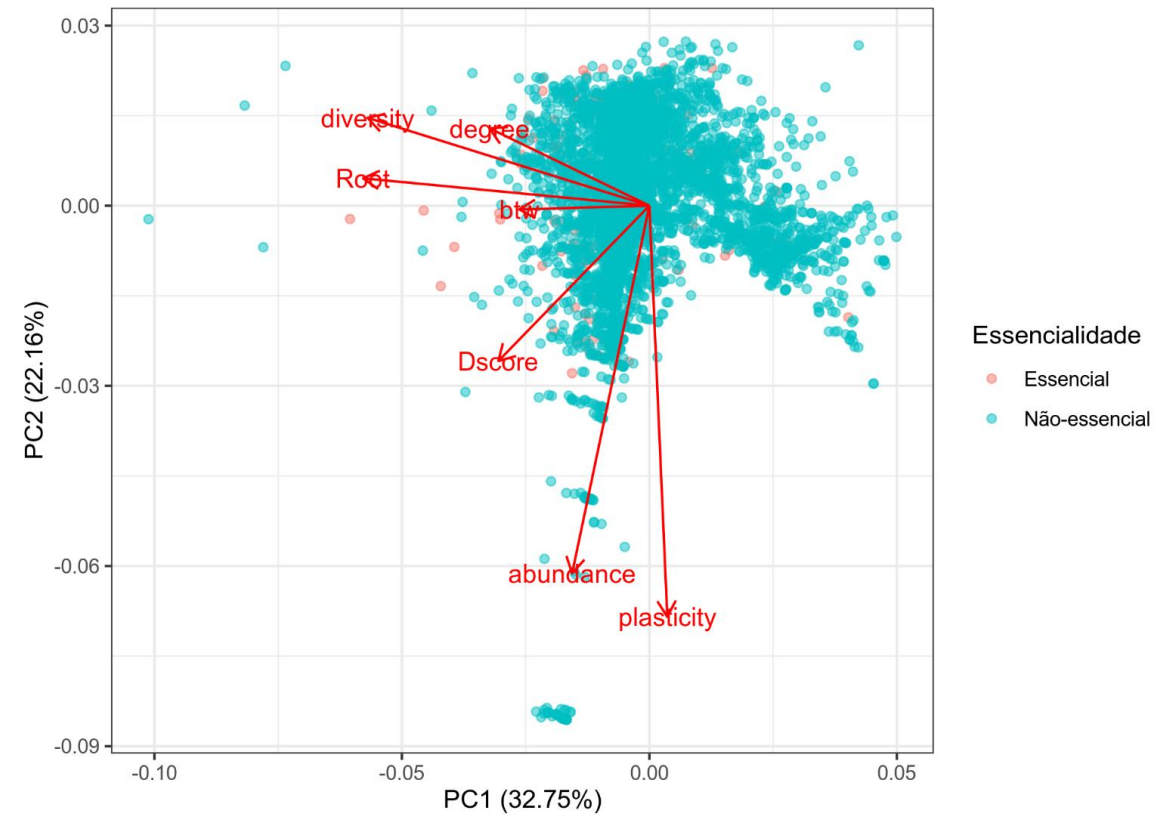


Análise exploratória

Camundongo



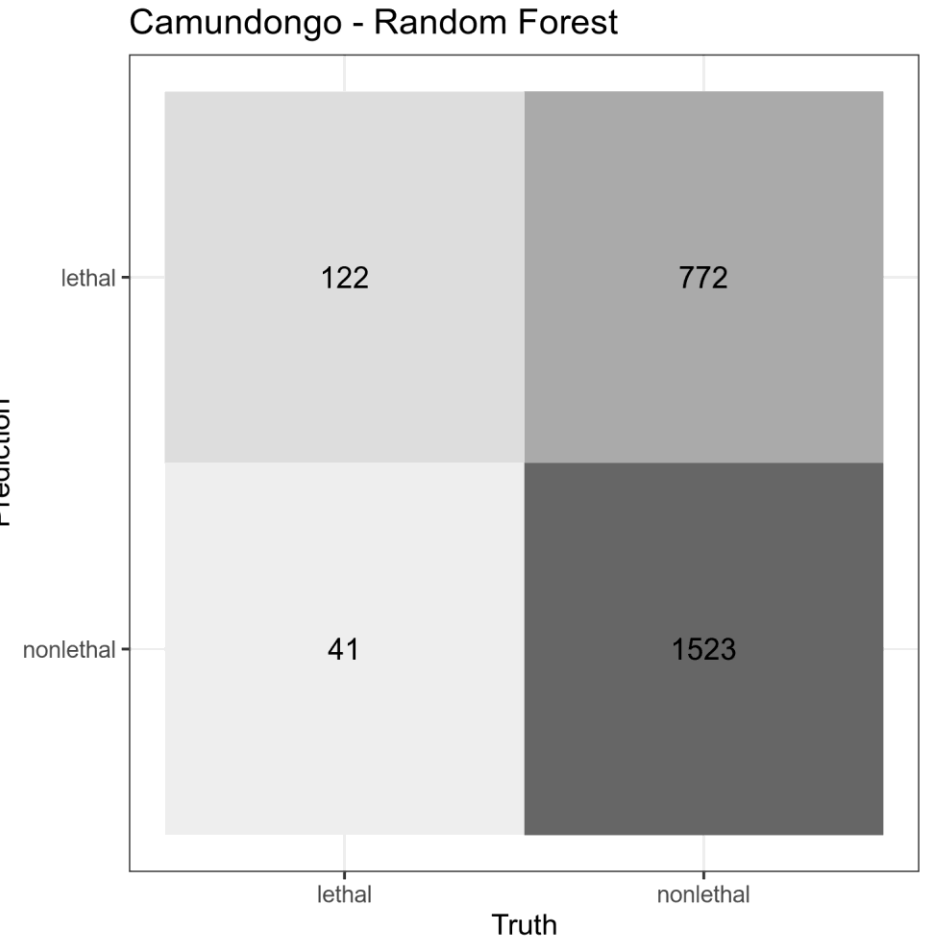
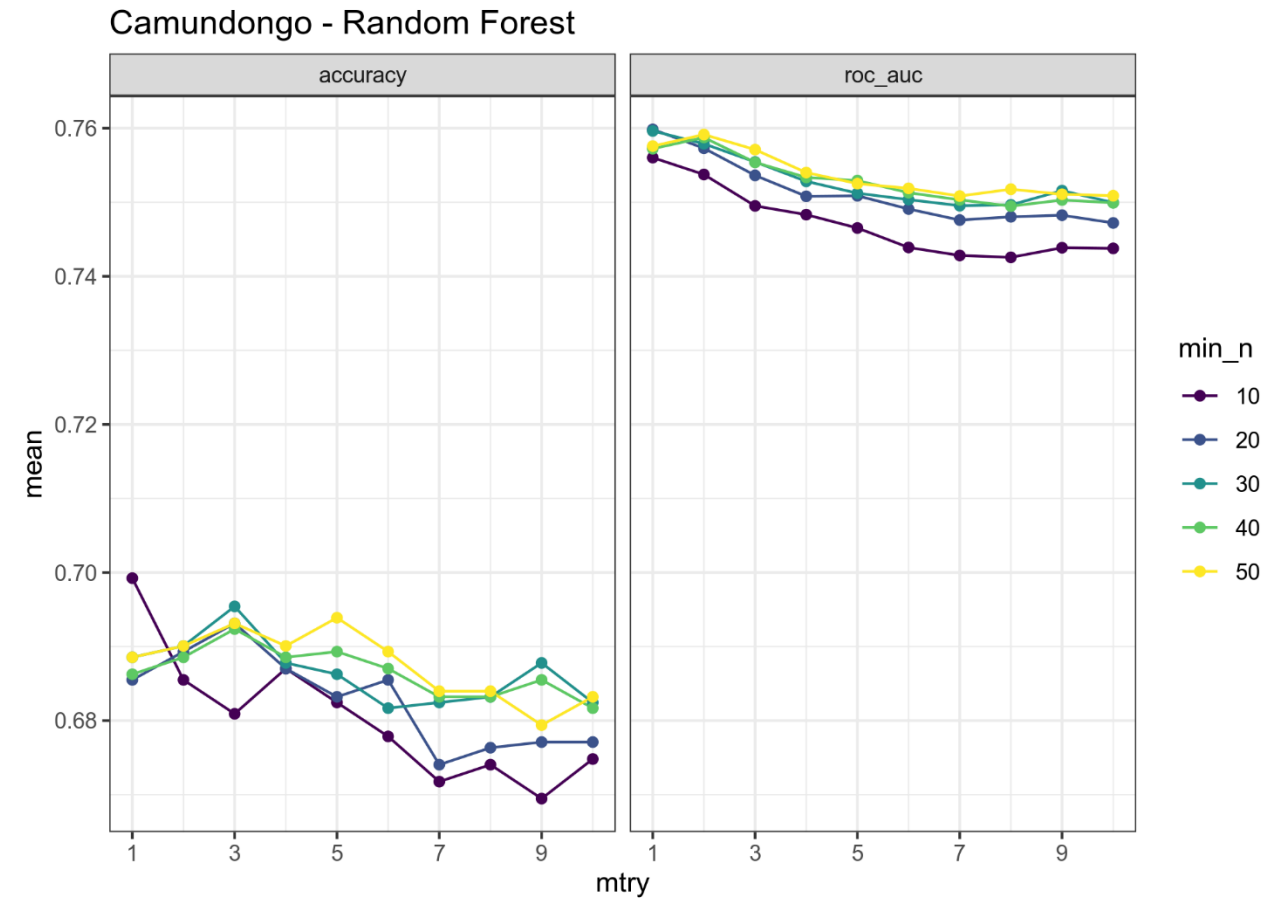
Levedura



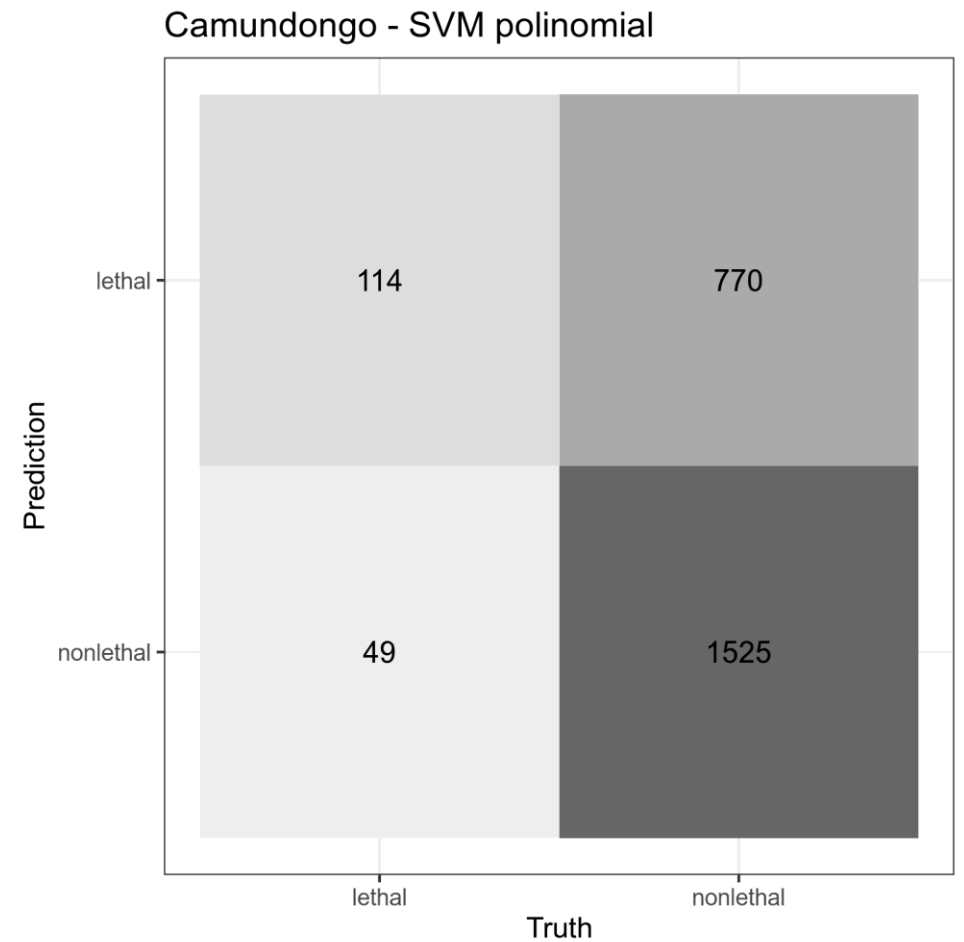
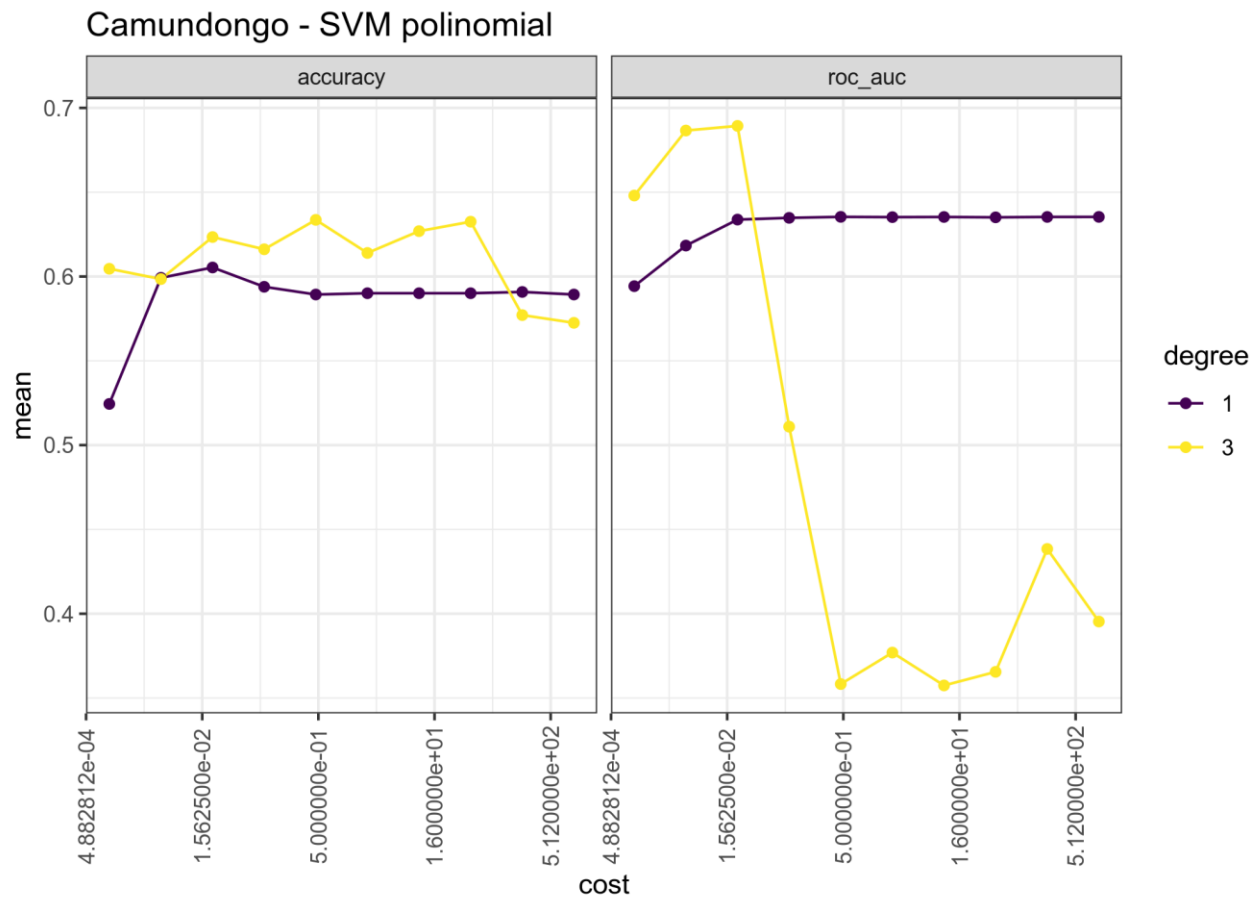
Resultados – Algoritmos usados

- Random Forest:
 - `mtry`: 10 níveis, de 1 a 10.
 - `min_n`: 5 níveis, de 10 a 50.
- SVM polinomial:
 - `cost`: 10 níveis, de -10 a 10
 - `degree`: 2 níveis (1 e 3)
- SVM radial:
 - `cost`: 10 níveis, de -10 a 10.
 - `rbf_sigma`: 5 níveis

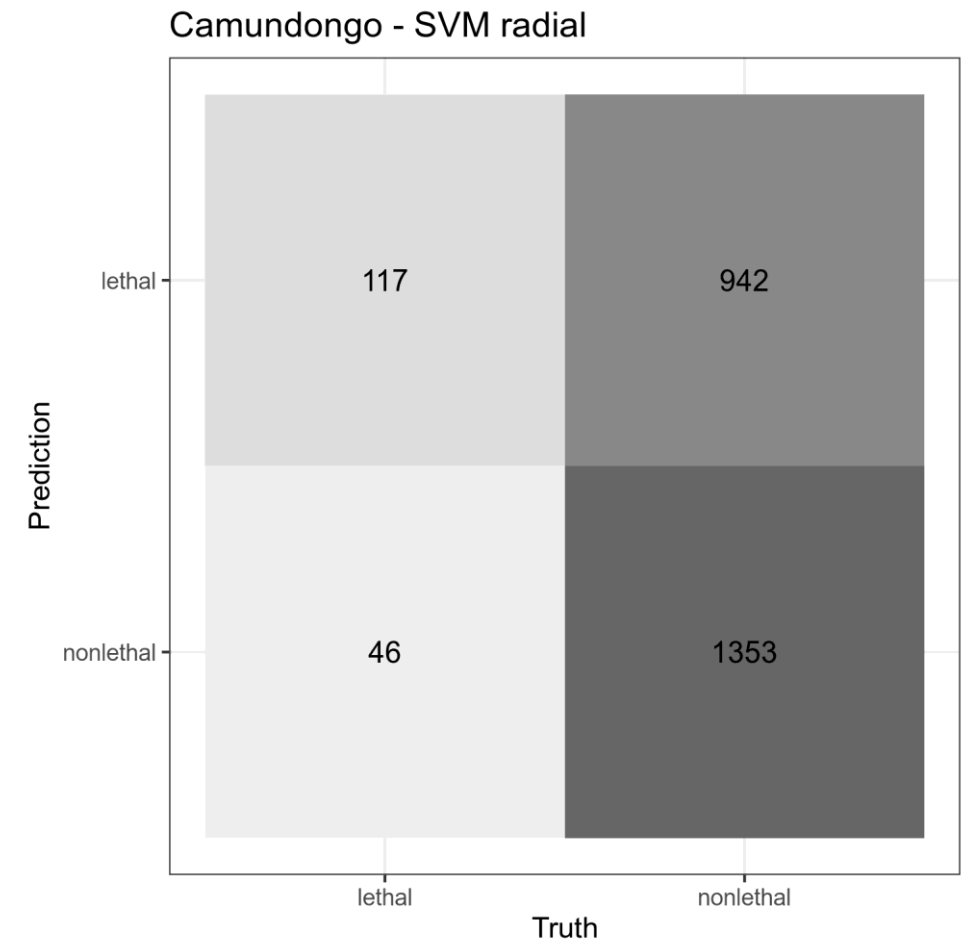
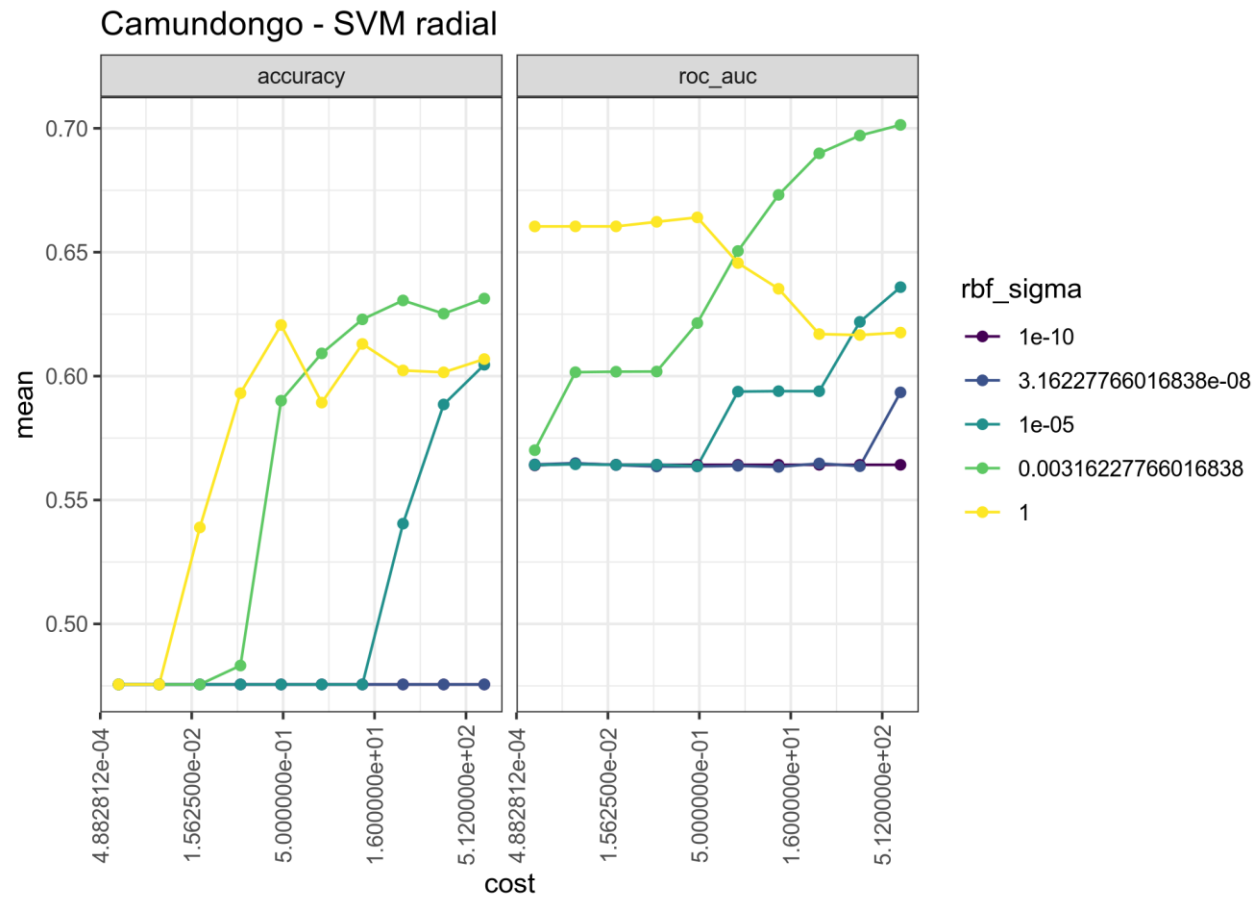
Resultados



Resultados

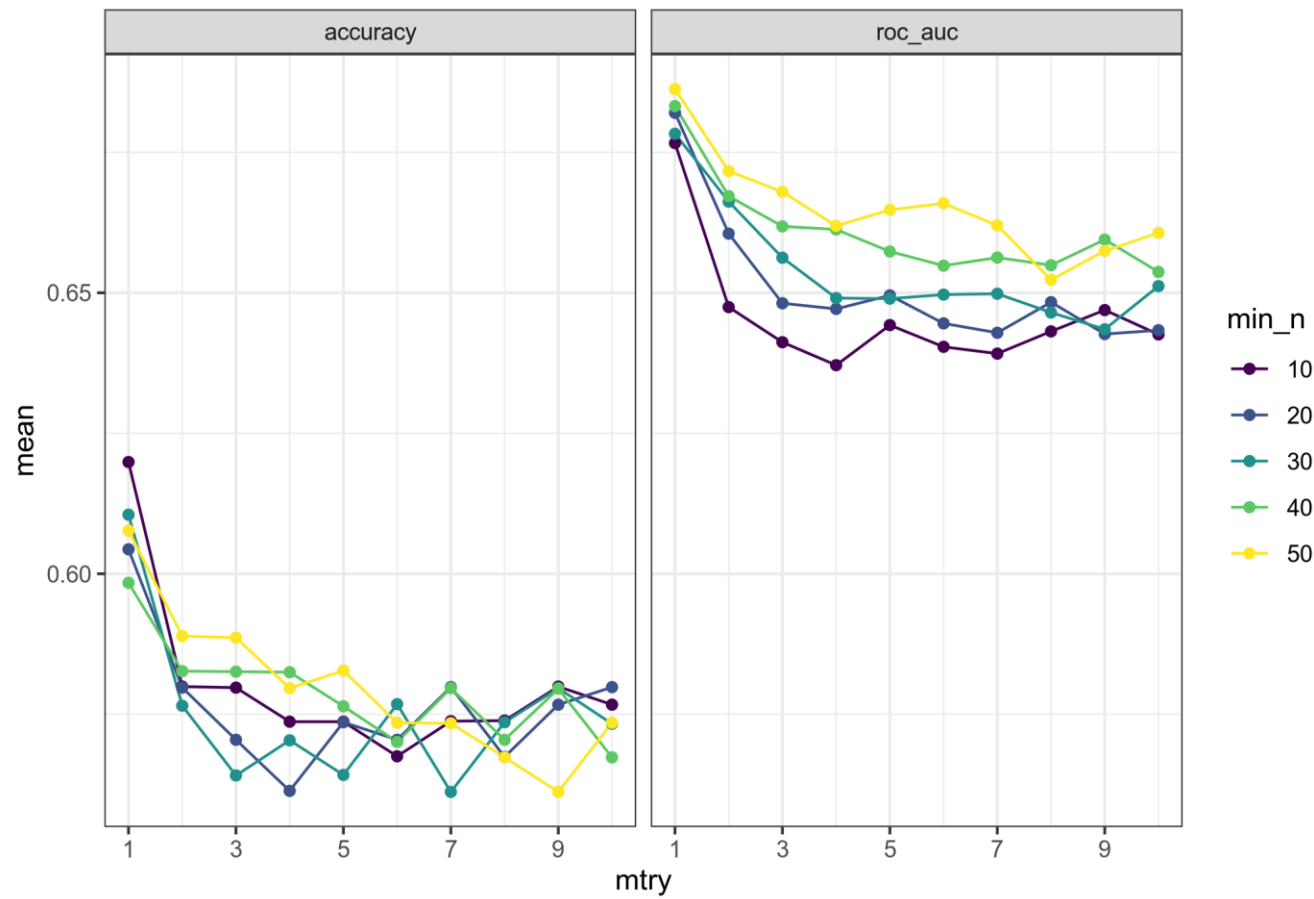


Resultados

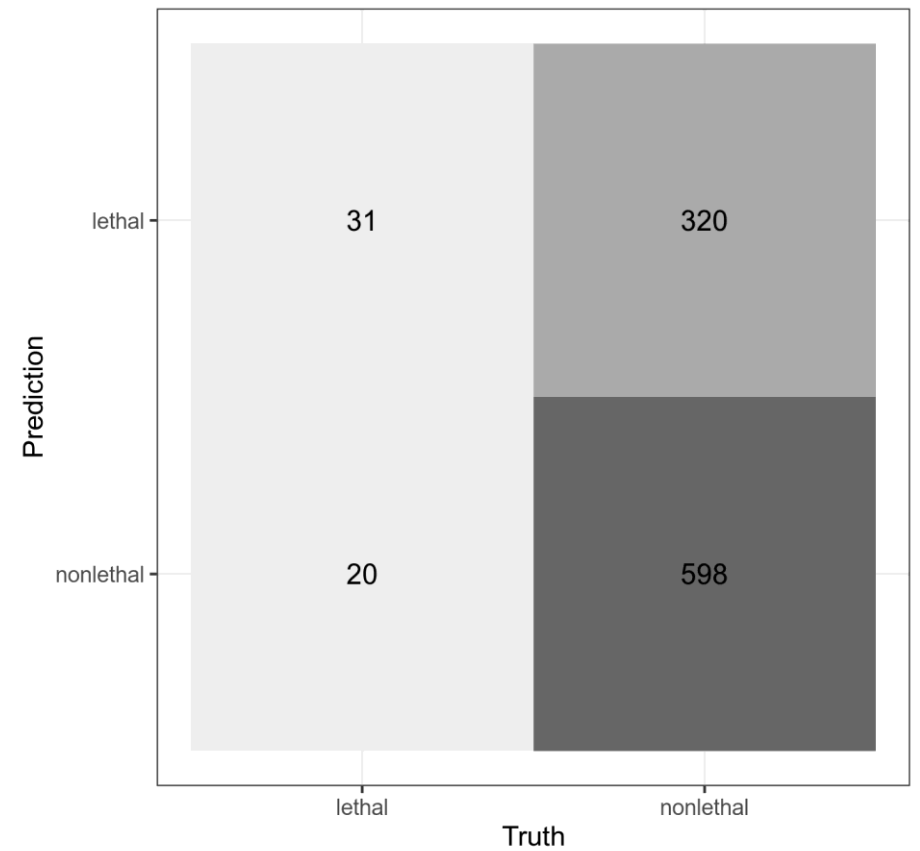


Resultados

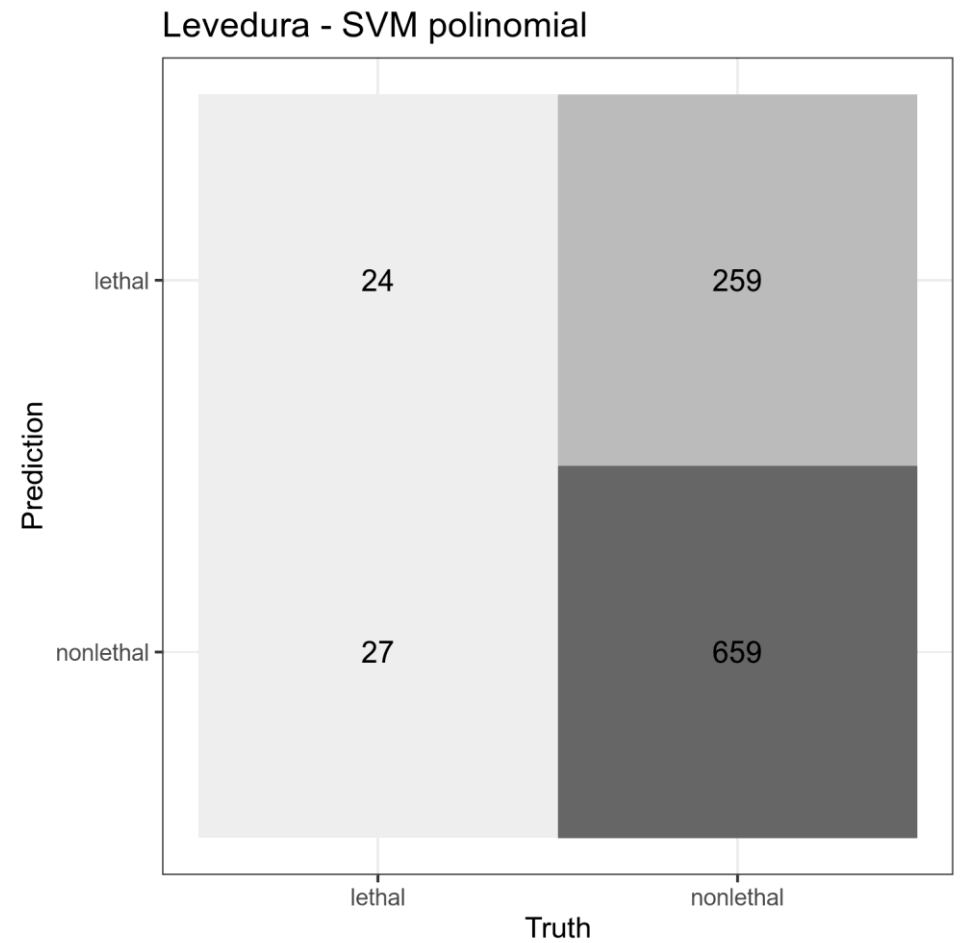
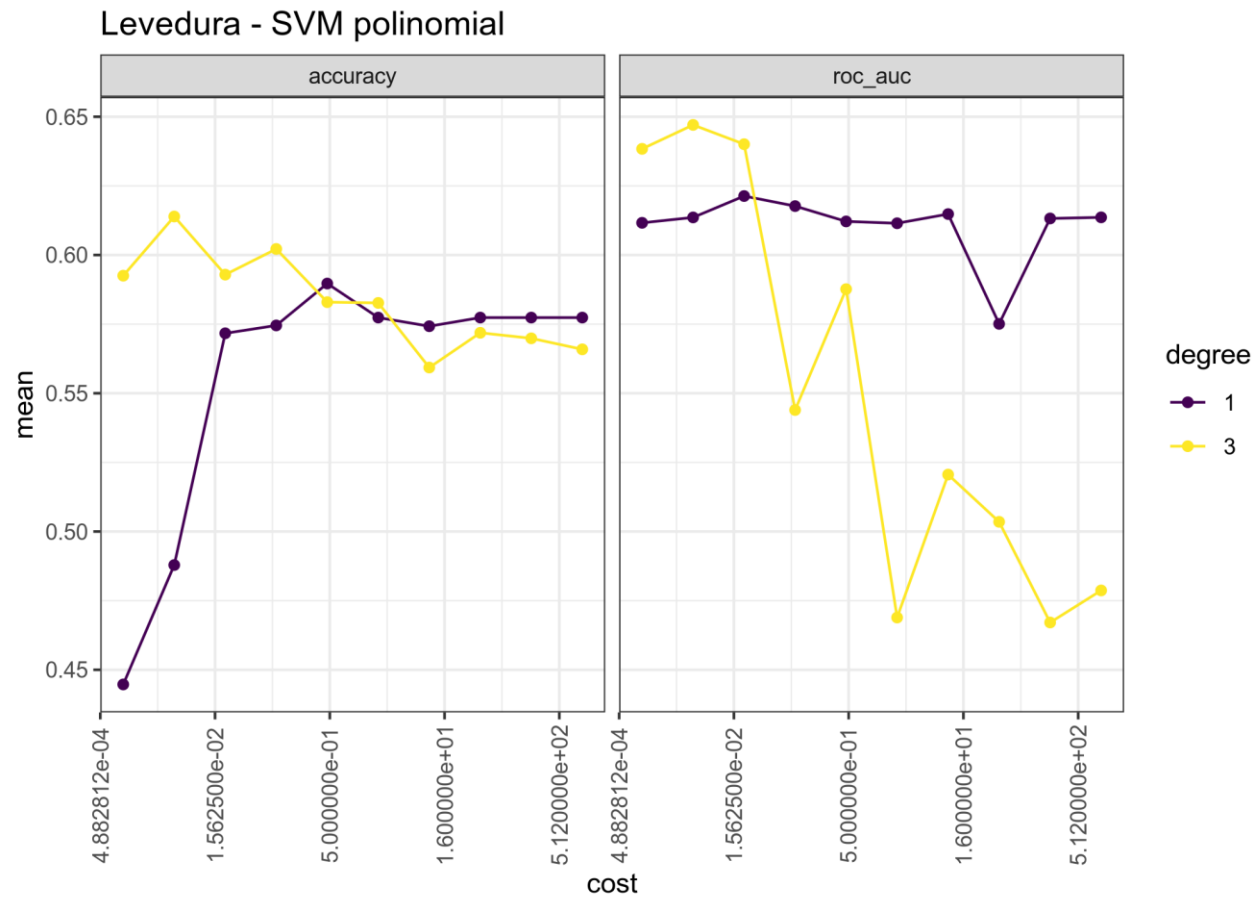
Levedura - Random Forest



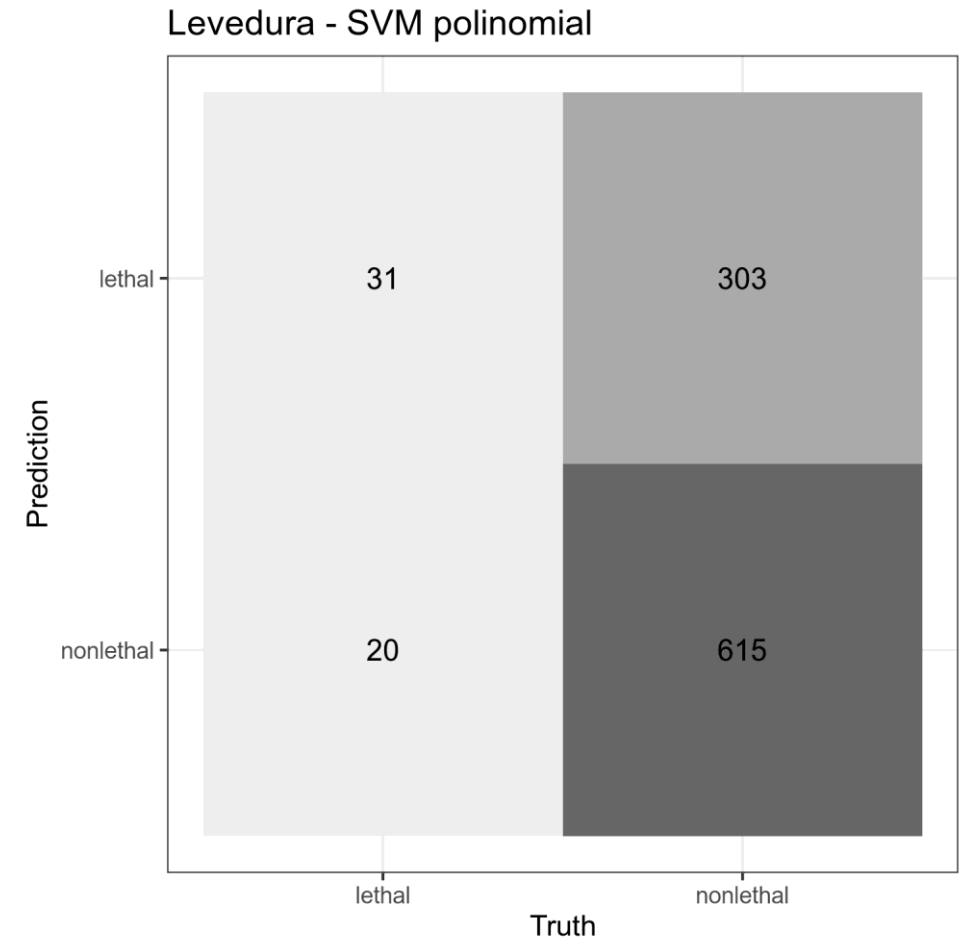
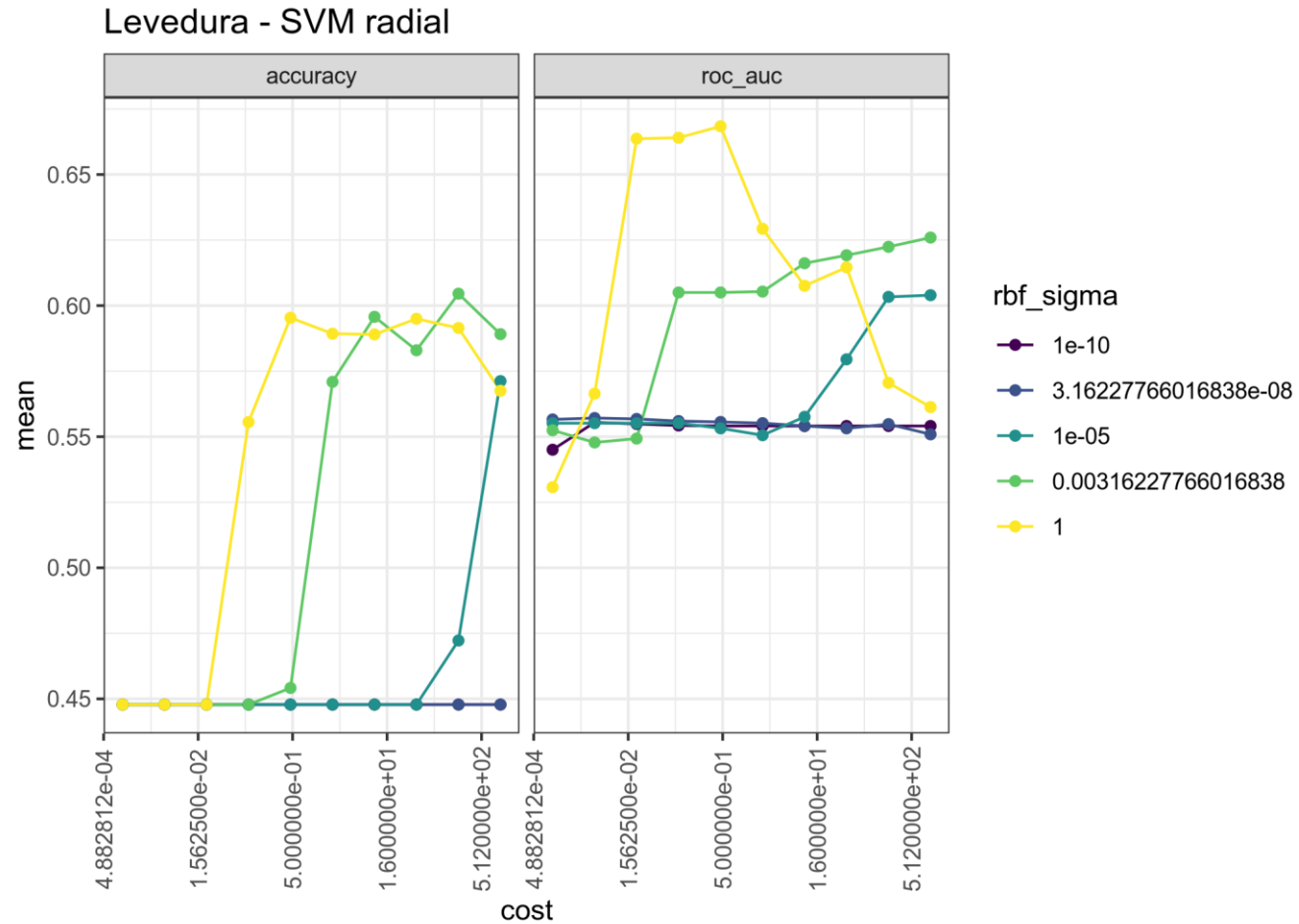
Levedura - Random Forest



Resultados



Resultados



Considerações finais

- Nenhum dos algoritmos produziu um bom modelo;
 - Todos abaixo de 70% de acurácia;
- Isto era esperado, devido a grande variabilidade das classes e ao fato de a essencialidade ser uma característica complexa;
- Os modelos podem ser melhorados com mais variáveis preditoras