



UNIVERSIDADE FEDERAL DO PARÁ
NÚCLEO DE DESENVOLVIMENTO AMAZÔNICO EM ENGENHARIA
PROGRAMA DE PÓS-GRADUAÇÃO EM COMPUTAÇÃO APLICADA

GISLENNE DA SILVA MOIA

PredictModelGUI: Ferramenta para classificação de genes essenciais através de técnicas de aprendizado de máquina.

Tucuruí PA

2025

GISLENNE DA SILVA MOIA

PredictModelGUI: Ferramenta para classificação de genes essenciais através de técnicas de aprendizado de máquina.

Dissertação apresentada ao Programa de Pós-Graduação em Computação Aplicada do Núcleo de Desenvolvimento Amazônico em Engenharia, da Universidade Federal do Pará, como requisito para a obtenção do título de Mestre em Computação Aplicada.

Orientador: Prof. Dr. Adonney Allan de Oliveira Veras

Coorientador: Prof. Dr. Cleison Daniel Silva

**Dados Internacionais de Catalogação na Publicação (CIP) de acordo com ISBD
Sistema de Bibliotecas da Universidade Federal do Pará
Gerada automaticamente pelo módulo Ficat, mediante os dados fornecidos pelo(a) autor(a)**

M712p Moia, Gislenne da Silva.
PredictModelGUI: Ferramenta para classificação de genes
essenciais através de técnicas de aprendizado de máquina. /
Gislenne da Silva Moia, . — 2025.
51 f. : il. color.

Orientador(a): Prof. Dr. Adonney Allan de Oliveira Veras
Coorientador(a): Prof. Dr. Cleison Daniel Silva
Dissertação (Mestrado) - Universidade Federal do Pará, Núcleo
de Desenvolvimento Amazônico em Engenharia, Mestrado
Profissional em Computação Aplicada, Tucuruí, 2025.

1. Classificação de Genes Essenciais. 2. Aprendizado de
Máquina. 3. Interface Gráfica (GUI). I. Título.

CDD 005

GISLENNE DA SILVA MOIA

PredictModelGUI: Ferramenta para classificação de genes essenciais através de técnicas de aprendizado de máquina.

Dissertação apresentada ao Programa de Pós-Graduação em Computação Aplicada do Núcleo de Desenvolvimento Amazônico em Engenharia, da Universidade Federal do Pará, como requisito para a obtenção do título de Mestre em Computação Aplicada.

Orientador: Prof. Dr. Adonney Allan de Oliveira Veras

Coorientador: Prof. Dr. Cleison Daniel Silva

Aprovada em 06 de Junho de 2025.

BANCA EXAMINADORA:

Prof. Dr. Adonney Allan de Oliveira Veras - PPCA/UFPA - Orientador

Prof. Dr. Cleison Daniel Silva - PPCA/NDAE/UFPA - Coorientador

Prof. Dr. Otávio Noura Teixeira - PPCA/NDAE/UFPA - Examinador Interno

Prof. Dr. Diego Assis das Graças - PPGBM/UFPA - Examinador Externo

Tucuruí PA
2025

PredictModelGUI: Ferramenta para classificação de genes essenciais através de técnicas de aprendizado de máquina.

Gislenne da Silva Moia¹, Victória Cardoso dos Santos², Saed Silva Sousa³, Walter de Barros Gomes Netto³, Rafael Azevedo Baraúna⁴, Diego Assis das Graças⁴, Artur Silva⁵, Cleison Daniel Silva¹, Adonney Allan de Oliveira Veras^{3,4*}

¹Programa de Pós-Graduação em Computação Aplicada, Universidade Federal do Pará campus Tucuruí (CAMTUC-UFPA), Pará, Brasil

²Faculdade de Engenharia de Computação - Universidade Federal do Pará (UFPA) campus Tucuruí Pará, Brasil

³Faculdade de Computação - FACOMP, UFPA / CCAST.

⁴Laboratório de Engenharia Biológica, Parque de Ciência e Tecnologia, Guamá, Belém Pará, Brasil.

* Autor correspondente

E-mail: allanveras@ufpa.br (AAOV)

RESUMO

As tecnologias de sequenciamento de DNA proporcionaram avanços significativos no conhecimento sobre o conteúdo gênico de inúmeros organismos, desde microrganismos até seres humanos. Dentre as análises realizadas pelas Ciências Ômicas, a Anotação se destaca como uma das mais importantes. Conceitualmente, esse processo consiste na inferência de informações biológicas a partir de sequências genômicas, o que permite aos Pesquisadores compreender a função de produtos genéticos, como os genes — Unidades Básicas da Hereditariedade responsáveis por características físicas e hereditárias de um organismo. Alguns genes desempenham funções vitais, pois codificam proteínas ou RNAs essenciais para processos como o Metabolismo Celular, que participam em vias cruciais como a Glicólise e o Ciclo do Ácido Tricarboxílico. As Plataformas de Sequenciamento passaram a gerar grandes volumes de dados, o que impulsionou avanços nas Áreas Ômicas e fomentou o desenvolvimento de métodos computacionais voltados às mais diversas análises. Mais recentemente, técnicas de Machine Learning e Inteligência Artificial têm sido aplicadas a esses dados, com estudos que demonstram a eficácia de abordagens inspiradas na Biologia. Esses modelos não exigem programação baseada em regras, embora sua criação ainda requeira habilidades avançadas em Programação e Computação. Com o objetivo de contribuir para a solução desse desafio, este estudo apresenta o *PredictModelGUI*, uma interface gráfica desenvolvida em Python que implementa nove modelos para classificar Genes Essenciais. A

interface permite importar conjuntos de dados, re-treinar os modelos e ajustar parâmetros. As informações são armazenadas no banco de dados do *software*, o que assegura rastreabilidade e proporciona uma ferramenta simples e intuitiva para testar diferentes configurações. Disponível em: [Predict Model GUI](#).

Palavras-chave: Classificação de Genes Essenciais, Aprendizado de Máquina, Interface Gráfica (GUI).

ABSTRACT

DNA sequencing technologies have provided significant advances in the understanding of the genetic content of numerous organisms, ranging from microorganisms to humans. Among the analyses performed in the Omics Sciences, Annotation stands out as one of the most important. Conceptually, this process consists of inferring biological information from genomic sequences, which allows researchers to understand the function of genetic products, such as Genes — the Basic Units of Heredity responsible for the physical and hereditary characteristics of an organism. Some Genes perform vital functions by encoding Proteins or RNAs essential for processes such as Cellular Metabolism, which participate in crucial pathways like Glycolysis and the Tricarboxylic Acid Cycle. Sequencing Platforms have started to generate large volumes of data, which has driven advances in the Omics fields and fostered the development of computational methods aimed at diverse analyses. More recently, Machine Learning and Artificial Intelligence techniques have been applied to these data, with studies demonstrating the effectiveness of biology-inspired approaches. These models do not require rule-based programming, although their creation still demands advanced skills in Programming and Computing. To contribute toward solving this challenge, this study presents PredictModelGUI, a graphical interface developed in Python that implements nine models to classify Essential Genes. The interface allows importing datasets, re-training models, and adjusting parameters. The information is stored in the software database, which ensures traceability and provides a simple and intuitive tool to test different configurations. Available at: [Predict Model GUI](#).

Keywords: Essential Gene Classification. Machine Learning. GUI (Graphical User Interface).

INTRODUÇÃO

Plataformas de Sequenciamento são dispositivos utilizados para identificar as bases nucleotídicas — adenina (A), guanina (G), citosina (C) e timina (T) — em amostras de DNA. Esses dispositivos são divididos em três gerações: 1ª geração, 2ª geração (conhecida como NGS – *Next Generation Sequencing*) e 3ª geração (Satam et al., 2023; Mandlik et al., 2024).

O desenvolvimento tecnológico dessas plataformas permitiu a geração de grandes

volumes de dados. A NovaSeq 6000, da Illumina¹, por exemplo, pode produzir até 3 terabytes de dados brutos, o que corresponde a 10 bilhões de leituras. Já a PromethION2 Solo, da *Oxford Nanopore*², pode gerar até 580 gigabytes de dados. Entre as principais características desses dispositivos estão o baixo custo por base sequenciada, a alta velocidade na geração de informações e a elevada precisão (Hu et al., 2021; Lee, 2023; Lu et al., 2025).

A produção desse volume significativo de dados impulsionou o desenvolvimento de ferramentas computacionais para auxiliar em diversas análises, como o processo de anotação. A anotação tem ampliado o conhecimento sobre o conteúdo gênico de diferentes organismos, por meio da inferência de informações biológicas associadas aos produtos gênicos, como proteínas hipotéticas, pseudogenes, produtos com função descrita, tRNA, rRNA e, por fim, produtos com sigla de gene. O gene é compreendido como a unidade fundamental, física e funcional da hereditariedade, responsável pelas características físicas e hereditárias do fenótipo de um organismo. Os genes são compostos por nucleotídeos, considerados subunidades do DNA (ácido desoxirribonucleico) ou do RNA (ácido ribonucleico), e apresentam três componentes principais: uma pentose (açúcar de cinco carbonos), uma base nitrogenada e um ou mais grupos fosfato (Guigó, 2023).

Alguns desses genes são considerados essenciais para processos vitais, como a replicação do DNA, a transcrição e tradução de proteínas e o metabolismo celular. Genes Essenciais também desempenham um papel crucial no transporte e na absorção de nutrientes, codificação das proteínas responsáveis pela importação de nutrientes vitais e pela exportação de resíduos. Esse mecanismo é fundamental para a sobrevivência do Microrganismo em seu ambiente (Liang et al., 2024).

A análise dessas informações contribui para o enriquecimento de bancos de dados biológicos públicos, como o NCBI³ (*National Center for Biotechnology Information*) e o DEG⁴ (*Database of Essential Genes*) (Liang et al., 2024). Métodos computacionais tradicionais, como ferramentas baseadas em alinhamento por similaridade (por exemplo, *BLAST*⁵) ou regras fixas de classificação, enfrentam dificuldades para processar o grande volume e a complexidade das informações disponíveis nesses bancos. Essas abordagens geralmente não são escaláveis, exigem pré-processamento manual e apresentam limitações ao lidar com dados incompletos ou altamente variados. Além disso, métodos baseados exclusivamente em homologia não conseguem identificar genes sem similaridade conhecida, enquanto algoritmos que utilizam limiares fixos — como conteúdo de GC ou tamanho da sequência — tendem a apresentar baixa precisão. Esse cenário reforça a necessidade de técnicas mais adaptativas,

¹ <https://www.illumina.com>

² <https://nanoporetech.com>

³ <https://www.ncbi.nlm.nih.gov/>

⁴ <http://origin.tubic.org/deg/public/index.php>

⁵ <https://blast.ncbi.nlm.nih.gov/Blast.cgi>

como os modelos de Inteligência Artificial - IA, que permitem maior flexibilidade e desempenho.

Técnicas de Machine Learning ou Deep Learning visam acelerar as análises e fornecer resultados com alta precisão. Essas abordagens envolvem algoritmos capazes de melhorar seu desempenho com base no treinamento, sem depender de regras de programação explícitas (Burkart & Huber, 2021).

A abordagem adotada por esses algoritmos pode ser classificada em três categorias: (i) Aprendizado Supervisionado, no qual os dados são rotulados; (ii) Aprendizado Não Supervisionado, em que os dados não são rotulados e o algoritmo identifica padrões ocultos; e (iii) Aprendizado por Reforço, em que o algoritmo interage repetidamente com o ambiente para atingir um objetivo, como vencer um jogo. Essas interações permitem encontrar soluções mais otimizadas para determinados problemas (Burkart & Huber, 2021; Naidu et al., 2023).

Uma vez treinados, esses modelos eliminam a necessidade de aplicar algoritmos baseados em regras clássicas de programação para obtenção de respostas. Basta fornecer os dados ao modelo treinado para que ele gere os resultados. No entanto, a criação desses modelos ainda requer uma etapa de programação, que pode ser desafiadora para usuários com pouca experiência em Computação (Aromolaran et al., 2021; Mukhamediev et al., 2022).

Para solucionar esse problema, é proposto o *PredictModelGUI*, uma interface gráfica desenvolvida em *Python*⁶ para a classificação de Genes Essenciais com o uso de técnicas de Aprendizado de Máquina. Essa ferramenta oferece nove modelos de classificação previamente treinados que permite ao usuário classificar Genes Essenciais de forma simplificada. Além disso, a interface gráfica do *PredictModelGUI* facilita o processo de treinamento dos modelos e possibilita a personalização dos parâmetros conforme a necessidade do usuário. Todas as configurações e resultados são armazenados em um Banco de Dados, para garantir a rastreabilidade e a consistência das informações.

⁶ <https://www.python.org/>

METODOLOGIA

Como funciona

O arquivo de entrada para o *PredictModelGUI* deve estar no formato CSV (*Comma-Separated Values*), composto por colunas que representam a frequência dos aminoácidos e, na última coluna, o nome do produto, *Product Name*, que representa a classe. Esse arquivo resulta da integração de três fontes de dados: informações de anotação e aminoácidos provenientes do banco de dados DEG e arquivos de anotação no formato *GenBank* (extensão .gb) obtidos do NCBI. Para simplificar esse processo, a própria ferramenta possui um módulo dedicado à geração automática do arquivo de entrada, o que elimina a complexidade para o usuário final.

O *PredictModelGUI* disponibiliza ao usuário três funções principais: Treinamento de Modelo, Predição e *Ensemble* — uma técnica que combina os resultados de múltiplos modelos para melhorar a precisão e robustez das classificações. Os parâmetros dessas funções podem ser ajustados diretamente na interface gráfica, o que proporciona flexibilidade, controle e permite a personalização intuitiva das configurações.

Na tarefa de treinamento de modelo, o usuário pode utilizar seu próprio conjunto de dados para treinar os modelos disponíveis no banco de dados do *software*. Nessa etapa, é necessário fornecer o conjunto de dados no formato CSV e indicar qual percentual dos dados deve servir para testar os modelos. Essa tarefa ocorre em paralelo ao treinamento para avaliar a melhoria nas métricas de avaliação à medida que o processo de aprendizagem dos modelos avança. Além disso, nesta etapa, o usuário pode selecionar quais métricas de avaliação observar: Acurácia (*Accuracy*), Precisão (*Precision*), Revocação (*Recall*), Pontuação F1 (*F1-score*), Índice Kappa, Coeficiente de Correlação de *Matthews* (*Matthews Correlation Coefficient*), que servem para avaliar o desempenho dos modelos.

Para realizar a predição, o usuário pode utilizar um modelo previamente treinado (obtido na etapa anterior) no formato PKL (*Python Pickle*). Além do modelo, deve carregar o conjunto de dados no formato CSV. Para facilitar esse processo, um módulo de preparação de dados está disponível, que padroniza os dados a serem utilizados no *PredictModelGUI*, como já mencionado. Vale destacar que a ferramenta possui um modelo padrão chamado *Ensemble*, uma técnica que combina os resultados de múltiplos modelos individuais para gerar predições mais precisas e robustas. Caso o usuário opte por utilizá-lo, basta fornecer o conjunto de dados para a predição.

Na etapa de *ensemble*, o usuário pode gerar um modelo personalizado baseado em seu próprio conjunto de dados. Para treinar um novo modelo *ensemble*, deve fornecer um arquivo CSV e indicar a pasta onde os modelos previamente treinados (em formato PKL) estão armazenados. Isso torna essencial o treinamento prévio dos modelos individuais, que servem

como entrada para a construção do modelo *ensemble*. Assim como na etapa de treinamento, os parâmetros e métricas de avaliação podem ser ajustados diretamente na interface do *PredictModelGUI*, o que possibilita a personalização da configuração do modelo conforme a necessidade do usuário.

Quanto aos resultados de todas as tarefas, o usuário pode optar por recebê-los localmente, ao informar o diretório onde a pasta de resultados será gerada, ou por *e-mail*, desde que forneça o endereço diretamente na interface do *PredictModelGUI*.

Os principais resultados obtidos incluem uma avaliação gráfica dos modelos, que apresenta todas as métricas selecionadas pelo usuário. Além disso, um arquivo CSV com a lista dos Genes Essenciais identificados ao final do processo de classificação é gerado. No contexto das etapas de treinamento e *ensemble*, os modelos resultantes da análise são disponibilizados no formato PKL, o que possibilita seu reuso em etapas futuras. As etapas de funcionamento estão descritas na Figura 1.

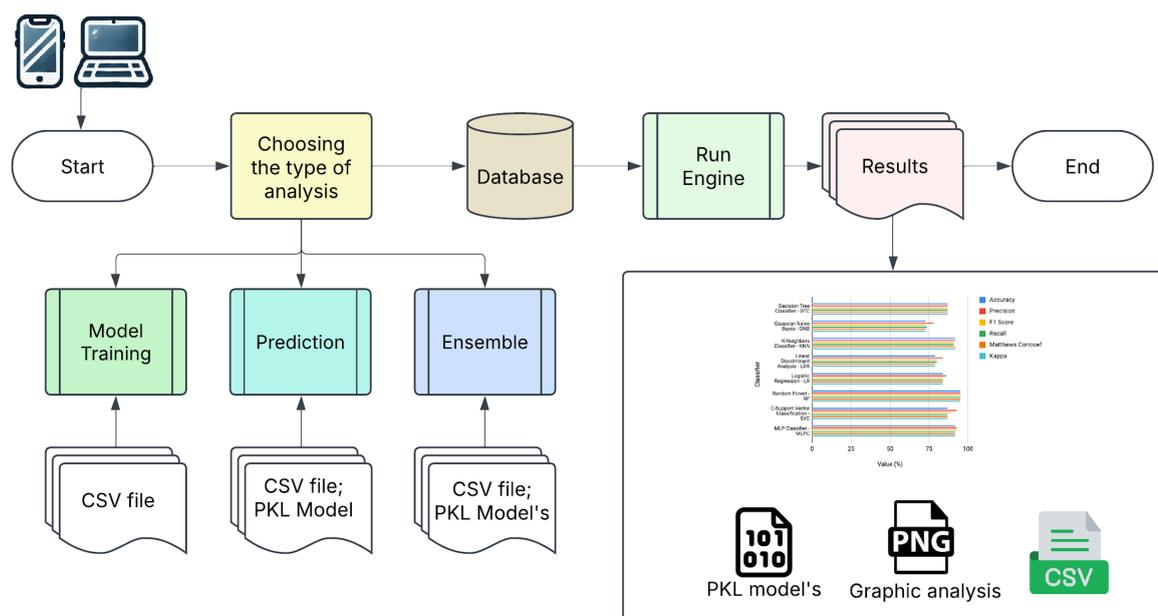


Figura 1. Pipeline do *PredictModelGUI*: A figura mostra o pipeline com as principais etapas do *software*.

Fonte de dados

Os conjuntos de dados utilizados nesta análise estão disponíveis em dois bancos de dados públicos: DEG e NCBI. As informações do DEG incluem sequências de aminoácidos e anotações de produtos gênicos bacterianos, disponíveis no formato FASTA e acompanhadas por um arquivo de anotação dessas sequências no formato CSV.

Os arquivos de anotação no formato *GenBank* foram obtidos do NCBI, que incluem 597 arquivos dos seguintes organismos: *Corynebacterium*, *Escherichia coli*, *Klebsiella*, *Mycobacterium*, *Rhodococcus*, *Burkholderia mallei* e *Salmonella*. Esses arquivos estão listados

no Anexo 01 do Material Suplementar.

Linguagem de Programação e Banco de Dados

A linguagem de programação *Python*, na versão 3.12, foi utilizada tanto para a implementação dos modelos de classificação quanto para o desenvolvimento da ferramenta *PredictModelGUI*. O *Python* é amplamente utilizado devido à sua versatilidade e à grande comunidade de desenvolvedores, especialmente nas áreas de Aprendizado de Máquina e Inteligência Artificial. O ambiente de desenvolvimento adotado foi o *PyCharm*⁷.

A interface gráfica *TkintModelGUI* foi criada a partir do pacote *Tkinter*⁸ na versão 8.5. O banco de dados implementado foi o *SQLite*⁹, na versão 3.42.0, responsável por gerenciar as etapas de processamento e armazenar informações relacionadas aos modelos e parâmetros.

Desenvolvimento da Interface Móvel

A ferramenta móvel foi desenvolvida por meio da linguagem de programação *Kotlin*¹⁰, na versão 2.0.20. Um dos principais motivos para sua escolha foi o suporte ao *Kotlin Multiplatform* (KMP), que facilita o compartilhamento de código entre diferentes plataformas. Para a criação das telas, foi adotado o *Jetpack Compose*¹¹, que permite construir *layouts* de forma simples e eficiente diretamente no código *Kotlin*.

A arquitetura MVVM (*Model-View-ViewModel*) foi selecionada para a aplicação, o que assegura uma separação clara entre a interface do usuário e a lógica da aplicação, conectadas por meio do *ViewModel*, o que promove modularidade e facilidade de manutenção. As bibliotecas utilizadas no desenvolvimento incluem: *Dokka*, *Ktor*, *Koin*, *Navigation Compose*, *Google Fonts*, *Kotlinx Serialization* e *Android Material3*.

Etapas de Pré-processamento

O pré-processamento e a criação do conjunto de dados a ser utilizado na geração dos modelos de classificação são divididos em duas etapas principais: pré-processamento do DEG e pré-processamento do NCBI.

Pré-processamento do DEG

A etapa inicial desse processo consiste no *download* dos arquivos que apresentam as sequências de aminoácidos no formato FASTA e das informações de anotação no formato CSV.

⁷ <https://www.jetbrains.com/pt-br/pycharm>

⁸ <https://docs.python.org/pt-br/3/library/tkinter.html>

⁹ <https://sqlite.org/>

¹⁰ <https://kotlinlang.org/>

¹¹ <https://developer.android.com/compose>

A Figura 2 apresenta um trecho representativo de cada um desses arquivos: (a) referente às sequências de aminoácidos (b) referente às informações de anotação. Os dados presentes nos cabeçalhos das sequências de aminoácidos são usados como referência para a construção do conjunto de dados nesta fase inicial, para assegurar a correta associação entre cada sequência de aminoácidos e sua respectiva informação de anotação.

a)

```

>DEG10010001
MENILDLWNQALAQIEKKLSKPSFETWMKSTKAHSLQGDTLTITAPNEFARDWLESRYLH
LIADTIYELTGEELSIKQVFNQNDVEDFMKPKQVKKAVKEDTSDFPQNMLNPKYTFDTF
VIGSGNRFHAHAASLAVAEAPAKAYNPLFIYGGVGLGKTHLMHAIGHYVIDHNPSAKVVYL
SSEKFTNEFINSDRNKAVDFRNRYRNVVLLIDDIQFLAGKEQTQEEFFHTFNTLHEES
KQIVISSDRPPKEIPTLEDRLRSRFEWGLITDITPPDLETRIAILRKKAKAEGLDIPNEV
MLYIANQIDSNIRELEGALIRVVAYSSLINKDINADLAAEALKDIIIPSSKPKVITIKEIQ
RVVQQQFNKLEDFKAKKRKTSVAFPRQIAMYLSREMTDSSLPKIGEEFGGRDHTTVIHA
HEKISKLLADDEQLQHVKEIKEQLK
>DEG10010002
MKFTIQKDRLVESVQDVLKAVSSRTTIPILTGKIVASDDGVSFTGSDSDISIESFIPKE
EGDKIEVTIEQPGSIVLQARFFSEIVKKLPMATVEIEVQNQYLTIIRSGKAEFNLGLDA
DEYPHLPQIEEHHAIQIPTDLLKNLIRQTVFAVSTSETRPILGWNWKEQSELLCTATD
SHRLALRKAKLDIPEDRSYNNVVPKSLTELSKILDDNQELVDIVITETQVLFKAKNVLF
FSRLLDGNYPDTTSLIPQDSKTEIIVNTKEFLQAIDRASLLAREGRNNVVKLSAKPAESI
EISSNSPEIGKVVEAIVADQIEGEEELNISFSPKYMLDALKVLEGAIEIRVSFTGAMRPFLI
RTPNDETIVQLILPVRTY

```

b)

DEG1001	DEG10010001	<u>dnaA</u>	16077069	COG0593L	DNA replication
DEG1001	DEG10010002	<u>dnaN</u>	16077070	COG0592L	DNA replication
DEG1001	DEG10010003	<u>gyrB</u>	16077074	COG0187L	DNA packaging
DEG1001	DEG10010004	<u>gyrA</u>	16077075	COG0188L	DNA packaging
DEG1001	DEG10010005	<u>guaB</u>	16077077	COG0516F	purine biosynthesis
DEG1001	DEG10010006	<u>serS</u>	16077081	COG0172J	<u>tRNA synthetase</u>
DEG1001	DEG10010007	<u>dnaX</u>	255767015	COG2812L	DNA replication
DEG1001	DEG10010008	<u>tmk/tdk</u>	255767017	COG0125F	pyrimidine biosynthesis
DEG1001	DEG10010009	<u>hoIb</u>	16077099	COG0470L	DNA replication
DEG1001	DEG10010010	<u>metS/metG</u>	16077106	COG0143J	<u>tRNA synthetase</u>
DEG1001	DEG10010011	<u>ispE</u>	16077114	COG1947I	<u>isoprenoid biosynthesis</u>
DEG1001	DEG10010012	<u>gcaD</u>	16077118	COG1207M	<u>aminosugar metabolism</u>
DEG1001	DEG10010013	<u>prs</u>	16077119	COG0462FE	glycolysis
DEG1001	DEG10010014	<u>spoVC</u>	16077121	COG0193J	RNA modification
DEG1001	DEG10010015	<u>divIC</u>	16077130	COG2919D	cell division

Figura 2. Pré-processamento DEG: (a) mostra um exemplo do arquivo DEG10.aa, que contém o cabeçalho de identificação da sequência seguido pelas sequências de aminoácidos. A Figura 1(b) mostra um trecho do arquivo de anotação CSV organizado em 13 colunas, composto por informações como IdDEG, IdDEGFASTA, símbolo de produto, idCOG e a anotação do produto, entre outras.

A etapa seguinte consiste na filtragem dos produtos gênicos que não possuem informação associada ao símbolo gênico, como proteínas hipotéticas ou produtos nomeados exclusivamente pelo identificador “*locus_tag*”. O arquivo resultante desse processo é gerado no formato CSV, composto por um cabeçalho com as letras que representam os aminoácidos, conforme o código genético, e uma coluna chamada “*Product Name*”, na qual são registradas as informações referentes ao símbolo gênico. Um exemplo desse arquivo pode ser visto na

Figura 3.

M	F	L	I	V	S	P	T	A	Y	H	Q	N	K	D	E	C	W	R	G	Product Name
3	0	5	3	7	2	2	8	2	2	2	1	1	6	2	3	0	0	5	6	rpmD
1	0	6	4	7	4	1	5	2	0	2	4	1	3	2	5	0	1	9	3	rpmD
3	1	5	5	5	2	2	6	5	0	2	1	1	5	1	4	0	0	6	5	rpmD
2	0	6	6	7	1	3	5	3	0	3	1	1	5	2	3	0	0	7	6	rpmD
3	0	5	4	6	4	1	4	1	2	3	3	3	6	1	3	1	0	4	4	rpmD
3	0	5	5	7	3	2	5	3	1	2	1	1	5	1	4	0	0	6	5	rpmD
2	0	5	6	7	1	2	4	4	0	3	3	3	6	1	4	0	0	4	6	rpmD
2	0	5	5	9	1	2	4	4	0	4	2	2	6	1	4	0	0	4	6	rpmD
1	0	6	4	9	5	0	4	5	0	1	4	4	5	2	6	0	0	7	3	rpmD
2	0	5	3	9	2	2	5	2	0	1	4	1	6	1	6	0	0	6	6	rpmD
1	0	6	5	6	8	2	5	6	0	2	4	1	4	2	3	0	1	9	6	rpmD
2	0	5	4	8	3	1	5	3	0	2	3	3	6	2	4	0	0	5	3	rpmD

Figura 3. Resultado do pré-processamento DEG: A figura mostra um trecho do conjunto de dados denominado conjunto de dados brutos DEG. Este conjunto de dados contém colunas nomeadas com os aminoácidos, onde a última corresponde ao nome do produto.

Pré-processamento do NCBI

No pré-processamento dos conjuntos de dados no formato *GenBank* (extensão .gb) obtidos do NCBI, são extraídas informações sobre os símbolos gênicos e suas respectivas sequências de aminoácidos. Para orientar essa extração, são utilizadas as *tags* “gene” (/gene) e “translation” (/translation) presentes nos arquivos, conforme mostrado na Figura 4. Em seguida, realiza-se a contagem da frequência de ocorrência dos aminoácidos em cada sequência, o que gera na geração de um arquivo chamado “NCBI raw dataset”, estruturado no mesmo formato do conjunto de dados produzido na etapa de pré-processamento do DEG, conforme ilustrado na Figura 3.

```

CDS      43..1635
         /gene="dnaA"
         /locus_tag="CGLAU_00005"
         /inference="ab initio prediction:Prodigal:2.6"
         /inference="similar to AA sequence:UniProtKB:A0R7K1"
         /codon_start=1
         /transl_table=11
         /product="Chromosomal replication initiator protein DnaA"
         /protein_id="AQQ14009.1"
         /translation="MADQQLEALWRDLVAELLQLSERPNSPVPTLTHQQRAYLQLVKP
VVLVDGYAILSAPHTAAKTVVEESLAPHIAAQLQARLGTPTLAVSIAAPGVQEVGQP
QQEPPIAQPQTEWTATQSSHTLHPPAESTYGFTAPGDQIPMGLDELAQM HARQTEESE
RRQANAHP TVPQIRREKPAHDPDREASLNP KYTFENFVIGSSNRFANGAAVAVAENP
ARAYNPLFIWGGSGLGKTHLLHAAGNYARVLEPNLKIKYVSSEEF TNDYINSVRDDRQ
ESFKRRYRNL DILMVDDIQFLEGEKGTQEEFFHTFNALHQANKQIILSSDRPPRQLTT
LEDRLRTRFEGGLITDIQPPDLETRIAILMKAAADGTQVSHDVLELIASQFESSIRE
LEGALIRVSAYSSLIDEPITMDVAQVALRDILPDENDVTVTATI IKDAAA EYFQVGVE
QLTGAGKTRHVAHARQIAMYLCRELTDLSLPKIGDEFGGKDHTTVMYADRKIRKEMTE
NRGTYDEIQELTQLIKNRARTR"
    
```

Figura 4. Arquivo de anotação do NCBI Genbank: A figura mostra um trecho do arquivo *Genbank* no qual as tags “gene” e “translation” são selecionadas, que exemplifica quais informações são extraídas desses arquivos.

Os arquivos CSV gerados nas etapas de pré-processamento do DEG e do NCBI são integrados, o que resulta em um conjunto de dados enriquecido. Essa integração baseia-se nos símbolos dos Genes Essenciais identificados na etapa inicial (pré-processamento do DEG). A construção desse conjunto combinado (DEGvsNCBI) segue os seguintes critérios: (i) o símbolo gênico deve estar presente em pelo menos seis ocorrências distintas, e (ii) a frequência de aminoácidos nessas ocorrências apresenta variação em pelo menos um aminoácido, que assegura a unicidade da entrada no conjunto de dados. O produto final consiste em um arquivo CSV semelhante ao ilustrado na Figura 3, porém expandido para cobrir um maior número de observações (linhas).

Modelos de Aprendizado

Neste estudo, foram utilizados nove modelos para classificar Genes Essenciais: *Logistic Regression* (LR), *Random Forest* (RF), *Support Vector Classification* (SVC), *MLP Classifier* (MLPC), *Linear Discriminant Analysis* (LDA), *Decision Tree* (DTC), *Gaussian Naive Bayes* (NB), *K-Neighbors Classifier* (KNN), *Ensemble Learning – VotingClassifier* (EL). Cada modelo foi gerado com os valores de parâmetros descritos na Tabela 1. Ressalta-se que todos esses valores podem ser ajustados diretamente na interface gráfica do *PredictModelGUI*.

Tabela 1. Modelos e Parâmetros de Classificação: Lista de modelos, parâmetros e seus respectivos valores padrão usados na *PredictModelGUI*.

Modelos	Parâmetros
<i>Decision Tree Classifier - DTC</i>	"criterion": "gini", "splitter": "best", "max_depth": None, "min_samples_split": 2, "min_samples_leaf": 1, "min_weight_fraction_leaf": 0.0, "max_features": None, "random_state": None, "max_leaf_nodes": None, "min_impurity_decrease": 0.0, "class_weight": None, "ccp_alpha": 0.0
<i>Gaussian Naive Bayes - GNB</i>	"priors": None, "var_smoothing": 1e-9
<i>K-Neighbors Classifier - KNN</i>	"n_neighbors": 5, "weights": "uniform", "algorithm": "auto", "leaf_size": 30, "p": 2, "metric": "minkowski", "metric_params": None, "n_jobs": None
<i>Linear Discriminant Analysis - LDA</i>	"solver": "svd", "shrinkage": None, "priors": None, "n_components": None, "store_covariance": False, "tol": 1e-4

<i>Logistic Regression - LR</i>	<i>"penalty": "l2", "dual": False, "tol": 1e-4, "C": 1.0, "fit_intercept": True, "intercept_scaling": 1, "class_weight": None, "random_state": None, "solver": "lbfgs", "max_iter": 100, "multi_class": "auto", "verbose": 0, "warm_start": False, "n_jobs": None, "l1_ratio": None</i>
<i>Random Forest - RF</i>	<i>"n_estimators": 100, "criterion": "gini", "max_depth": None, "min_samples_split": 2, "min_samples_leaf": 1, "min_weight_fraction_leaf": 0.0, "max_features": "sqrt", "max_leaf_nodes": None, "min_impurity_decrease": 0.0, "bootstrap": True, "oob_score": False, "n_jobs": None, "random_state": None, "verbose": 0, "warm_start": False, "class_weight": None, "ccp_alpha": 0.0, "max_samples": None</i>
<i>C-Support Vector Classification - SVC</i>	<i>"C": 1.0, "kernel": "rbf", "degree": 3, "gamma": "scale", "coef0": 0.0, "shrinking": True, "probability": False, "tol": 1e-3, "cache_size": 200, "class_weight": None, "verbose": False, "max_iter": -1, "decision_function_shape": "ovr", "break_ties": False, "random_state": None</i>
<i>MLP Classifier - MLPC</i>	<i>"hidden_layer_sizes": (100,), "activation": "relu", "solver": "adam", "alpha": 0.0001, "batch_size": "auto", "learning_rate": "constant", "learning_rate_init": 0.001, "power_t": 0.5, "max_iter": 200, "shuffle": True, "random_state": None, "tol": 1e-4, "verbose": False, "warm_start": False, "momentum": 0.9, "nesterovs_momentum": True, "early_stopping": False, "validation_fraction": 0.1, "beta_1": 0.9, "beta_2": 0.999, "epsilon": 1e-8, "n_iter_no_change": 10, "max_fun": 15000</i>
<i>Ensemble Learning - VotingClassifier (EL)</i>	<i>estimators, *, voting='hard', weights=None, n_jobs=None, flatten_transform=True, verbose=False</i>

Padronização e Métricas de Avaliação

A avaliação dos modelos implementados no *PredictModelGUI* foi realizada com base nas seguintes métricas: Acurácia, Precisão, Revocação, *F1-score*, Índice Kappa, Coeficiente de Correlação de *Matthews*. Essas métricas são essenciais para fornecer uma visão abrangente do desempenho dos modelos, que ajudam a entender sua efetividade e confiabilidade, e garantem uma avaliação detalhada e precisa.

Para o processo de normalização, foi utilizado o *Standard Scaler*, que padroniza os dados para que tenham média zero e desvio padrão igual a um. Essa técnica é menos sensível a valores extremos (*outliers*) em comparação com outros métodos de escalonamento. A escolha foi feita com o objetivo de aumentar a robustez e consistência do desempenho dos modelos.

Estação de Trabalho

A estação de trabalho utilizada para gerar os modelos e realizar os testes gerais da ferramenta foi uma Penguin Computing Tempest 4400, equipada com 1024 GB de RAM, Processador AMD Opteron 6376 (64 núcleos), 4 TB de armazenamento. O sistema operacional utilizado foi o Debian 12.

RESULTADOS E DISCUSSÃO

Os resultados deste estudo estão organizados em uma análise detalhada da etapa de pré-processamento do conjunto de dados, seguida pela avaliação dos modelos preditivos por meio das métricas selecionadas. Além disso, examina-se a estratégia metodológica adotada e apresenta-se a ferramenta desenvolvida, que termina com a exibição do aplicativo móvel.

Análise da etapa de pré-processamento

As etapas de pré-processamento e os arquivos produzidos foram organizados nas seguintes fases: (1) Pré-processamento DEG, conjuntos de dados DEG10.aa e deg_annotation_p.csv: esses conjuntos de dados apresentam informações relacionadas a sequências de aminoácidos e anotações, cada um com inicialmente 26.619 observações (linhas). Após a aplicação das regras de pré-processamento, foi produzido um arquivo CSV com 12.607 observações. (2) Pré-processamento NCBI: o pré-processamento inicial do conjunto de dados NCBI gerou um arquivo CSV com 830.153 observações. (3) Cruzamento dos conjuntos de dados DEG e NCBI: esta etapa envolveu o cruzamento dos arquivos CSV DEG e NCBI, o que resultou em um arquivo com 99.818 observações.

Nesta etapa, o arquivo DEG que inclui os Genes Essenciais é usado como referência para guiar o processo, de modo que as classes (nome do produto) com a mesma frequência de aminoácidos sejam removidas. Em seguida, as classes com frequência inferior a seis representações são removidas da análise, o que gera um arquivo com 99.818 observações, equivalente a 1.720 classes distintas. A Figura 5 mostra a distribuição de frequência das classes.

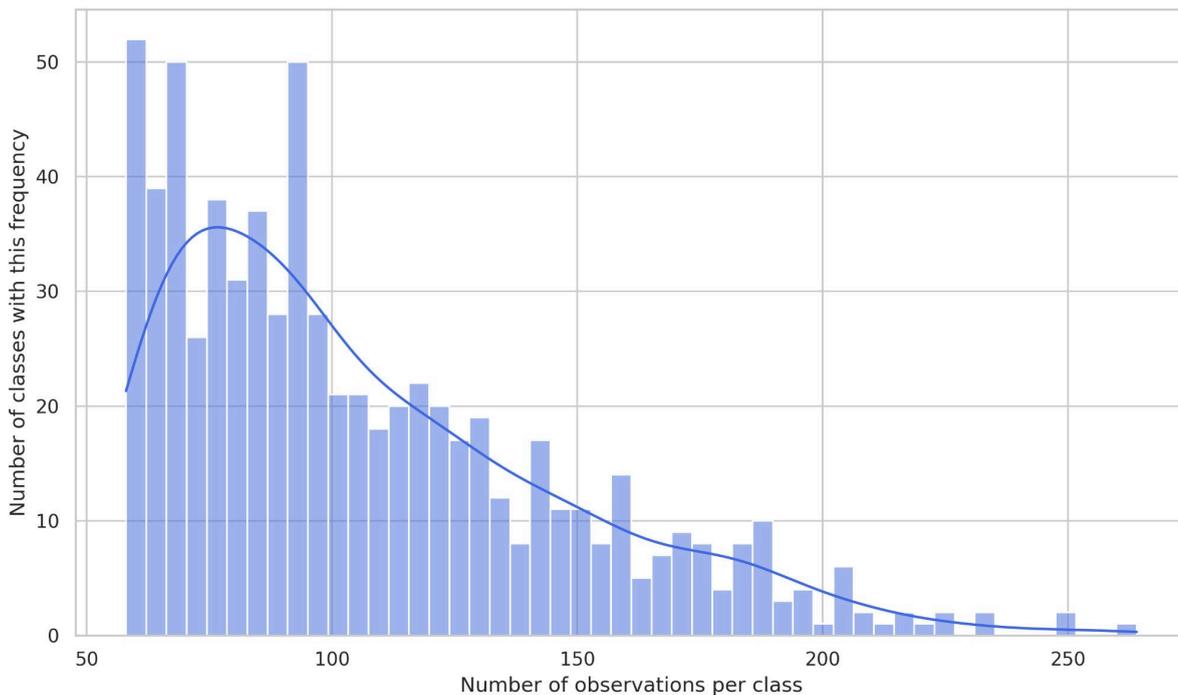


Figura 5. Distribuição da frequência de observações por classe (símbolo do gene): O eixo X representa o número de vezes que cada classe (produto com símbolo do gene na coluna "Nome do Produto") aparece no conjunto de dados. O eixo Y indica o número de classes que ocorrem com uma determinada frequência. A curva KDE (*Kernel Density Estimation*) sobreposta fornece uma estimativa suave da densidade da distribuição de frequência.

Dada a natureza desbalanceada do conjunto de dados, foi adotada uma estratégia baseada na definição da proporcionalidade das classes, que é determinada pela razão entre o número total de observações no conjunto de dados e o número de representações por classe. Como resultado dessa análise, foi estabelecido que cada classe deveria conter pelo menos 58 observações. Com base nesse critério, foi gerado um conjunto de dados composto por 74.132 observações e 696 classes, denominado *Dataset1*. As classes cuja proporcionalidade ficou abaixo desse limite foram alocadas a um segundo conjunto de dados, denominado *Dataset2*, composto por 25.686 observações e 1.024 classes.

Ao final, é realizada uma amostragem estratificada, onde se separou 1% de cada conjunto de dados para criar um conjunto de dados de validação que não participa do processo de construção do modelo. Nesse ponto, os conjuntos de dados são denominados *dataset1_toBuildModel.csv* e *dataset2_toBuildModel.csv*; para simplificar, serão denominados no manuscrito como *dataset1* e *dataset2*.

Construção de modelos e análise de métricas de avaliação

O próximo passo foi criar os modelos de classificação. Ambos os conjuntos de dados foram submetidos a uma amostragem aleatória estratificada, dividido em 80% para treinamento e

20% para teste dos modelos, seguidos pela criação dos oito modelos de classificação. Oito modelos de classificação iniciais foram desenvolvidos e, em seguida, combinados em um modelo final por meio da abordagem de Aprendizado por Conjunto. Para essa integração, foi aplicado o método *VotingClassifier*, que combina as previsões dos modelos base para melhorar a precisão da classificação.

Ao final desse processo, dezessete modelos foram criados: oito com base em cada conjunto de dados (*datasets 1 e 2*) e, por fim, o modelo resultante do agrupamento por Aprendizado de Conjunto. Os modelos foram avaliados por meio das métricas de desempenho para todos os modelos criados a partir dos dois conjuntos de dados. A análise das métricas de avaliação — Acurácia, Precisão, Recall, Pontuação F1, MCC e Kappa — é apresentada nas Figuras 6 e 7. Essas figuras mostram o desempenho dos modelos de classificação.

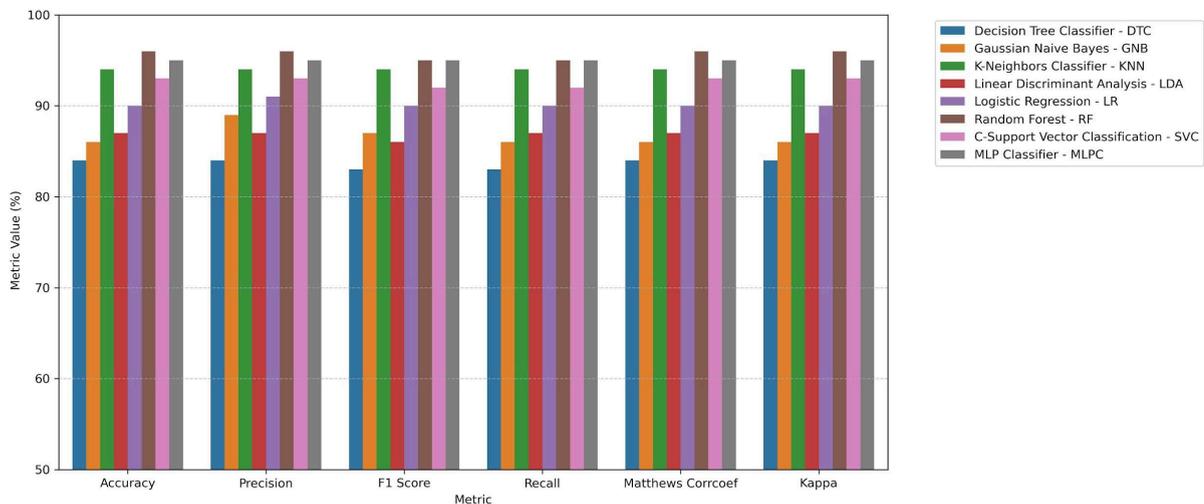


Figura 6. Avaliação dos modelos de classificação para o Conjunto de Dados 1: Análise das métricas de avaliação dos modelos criados com o Dataset1.

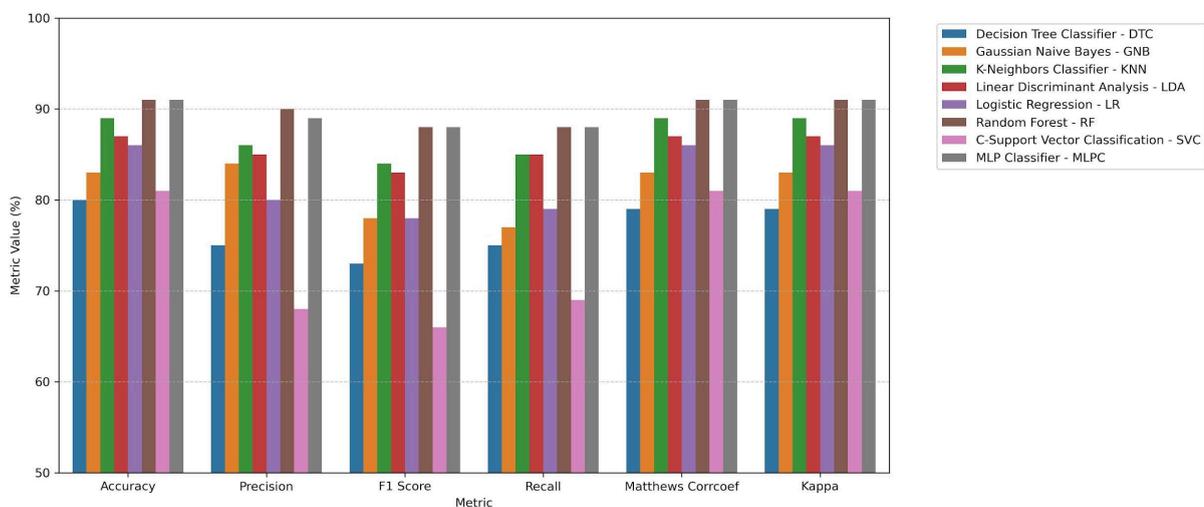


Figura 7. Avaliação dos modelos de classificação para o Conjunto de Dados 2: Análise das métricas de avaliação dos modelos criados com o Dataset2.

A análise das Figuras 6 e 7 revela um comportamento semelhante entre os modelos avaliados com diferentes conjuntos de dados. O modelo que apresentou os melhores resultados nas métricas de avaliação foi o *Random Forest*, seguido pelo *MLP Classifier*, em ambos os conjuntos. Essa consistência indica que o *Random Forest* possui desempenho superior e maior robustez nesta análise, independentemente das variações nos dados utilizados. Esse desempenho pode estar relacionado às características do próprio algoritmo, que combina múltiplas árvores de decisão e utiliza uma votação majoritária para determinar a predição final. Essa estratégia reduz o risco de sobreajuste (*overfitting*) e melhora a capacidade de generalização do modelo. Além disso, o *Random Forest* é capaz de lidar adequadamente com ruído nos dados e com o desbalanceamento entre classes, características frequentemente presentes em dados biológicos.

Embora o *MLP Classifier* também tenha apresentado desempenho elevado, sua performance tende a ser mais sensível à escolha dos hiperparâmetros e à normalização dos dados de entrada. Esses fatores podem afetar negativamente sua capacidade de generalização caso não sejam ajustados adequadamente, o que justifica sua performance ligeiramente inferior em relação ao *Random Forest*.

A boa performance de ambos os modelos, baseados em fundamentos distintos — *Random Forest* por meio de árvores de decisão e *MLP Classifier* por Redes Neurais Artificiais —, demonstra que diferentes abordagens podem ser eficazes na tarefa de classificação de Genes Essenciais. No entanto, isso não implica que a escolha do modelo seja indiferente. Cada algoritmo responde de forma diferente às variações nos dados e apresenta vantagens e limitações próprias. Dessa forma, a adoção de estratégias que exploram múltiplas técnicas, como o uso de métodos de Aprendizado de Conjunto (*ensemble*), pode contribuir para aumentar a robustez da análise, e resultar em classificações mais confiáveis e com melhor desempenho geral.

Para analisar a evolução das métricas de avaliação dos modelos de classificação criados neste estudo, foram desenvolvidos três modelos de agrupamento baseados em Aprendizado de Conjunto. O primeiro, chamado *Ensemble Step01*, utilizou o agrupamento de modelos criados com o *dataset1*. O segundo, chamado *Ensemble Step02*, foi formado a partir de modelos baseados exclusivamente no *dataset2*. Finalmente, os dezesseis modelos foram combinados em um modelo final chamado *Ensemble Final*. É importante notar que as métricas de desempenho dos modelos de conjunto (Etapa 1, Etapa 2 e Conjunto Final) foram avaliadas com um conjunto de validação composto por uma amostra estratificada dos *dataset 1* e *2*. Essa amostra não foi usada em nenhuma etapa do processo de construção do modelo, a fim de se garantir uma avaliação imparcial e representativa do desempenho preditivo. Essa abordagem tornou possível avaliar o impacto dos diferentes agrupamentos na melhoria do desempenho do

modelo Figura 8.

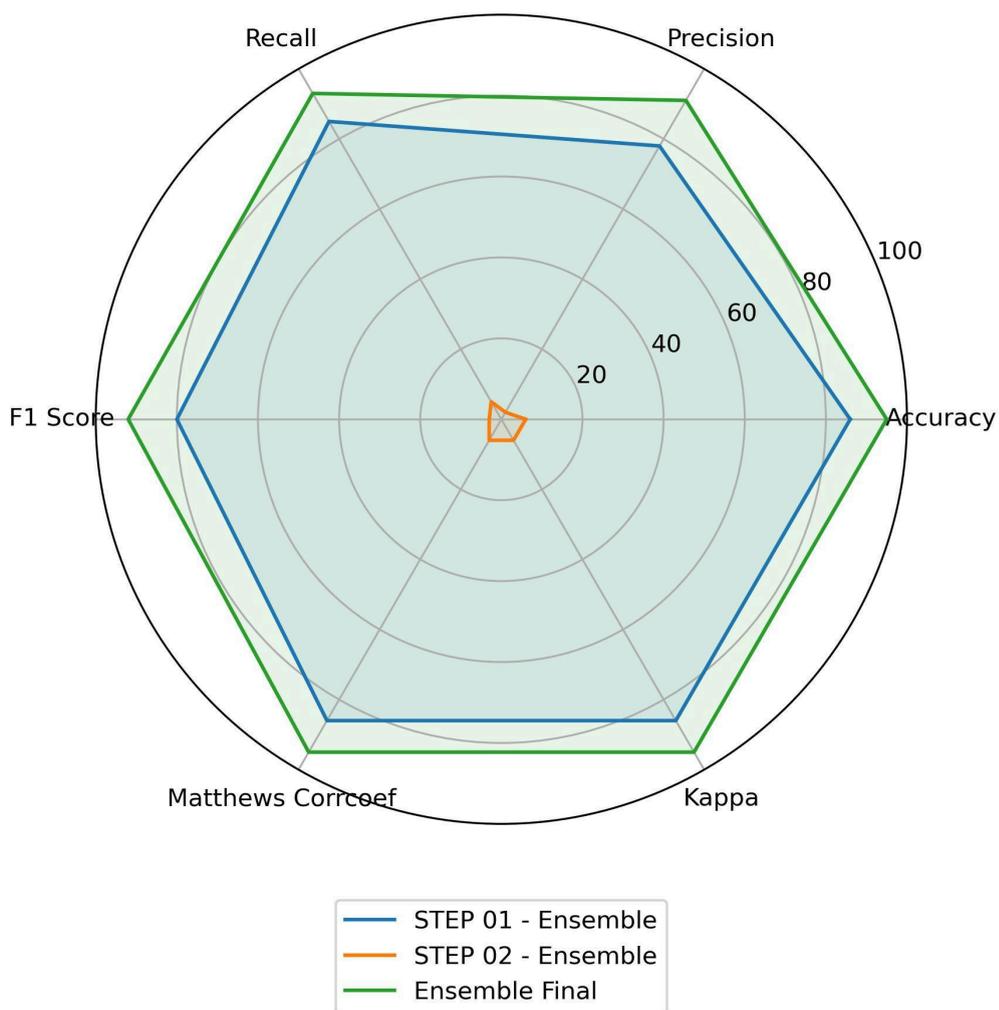


Figura 8. Avaliação do modelo de conjunto de classificação: Esta figura mostra a avaliação dos modelos gerados por meio do método de Aprendizado de Conjunto para o conjunto de oito modelos com os *dataset 1* e *2*, os quais resultaram em *ensemble_step01* e *ensemble_step02*, respectivamente.

O gráfico de radar na Figura 8 permite visualizar e comparar rapidamente o desempenho de diferentes métricas em cada estágio do conjunto. Isso facilita a análise de tendências e flutuações no desempenho do modelo em diferentes estágios de refinamento do conjunto. As métricas de avaliação que obtiveram as maiores pontuações na análise foram Precisão, MCC e Kappa, todas com 95%.

É importante enfatizar que o uso do método de aprendizagem por conjunto, que combina os modelos gerados a partir dos *dataset 1* e *2*, permite uma representação abrangente de todas as classes de Genes Essenciais do banco de dados DEG. Isso sugere que o modelo resultante é capaz de identificar um gene essencial em qualquer conjunto de dados, independentemente de sua alta ou baixa representação. Além disso, os valores das métricas de avaliação demonstram a eficácia do modelo gerado, o que reforça sua capacidade de predição precisa e confiável.

Versões para desktop e aplicativo do *PredictModelGUI*

O *PredictModelGUI* possui uma interface intuitiva e simplificada (Figura 9), desenvolvida em *Python*, com o objetivo de proporcionar uma experiência eficiente ao usuário. Ao iniciar a aplicação, uma janela introdutória é exibida, o que permite acesso imediato às funcionalidades da ferramenta. Após a inicialização, o usuário pode explorar os principais módulos de análise, conforme ilustrado na Figura 9: *Training*, que possibilita o treinamento de modelos previamente definidos e salvos no banco de dados; *Prediction*, que realiza a predição sobre um conjunto de dados e identifica automaticamente as classes (Genes Essenciais) presentes; e *Ensemble*, que possibilita a combinação de múltiplos modelos para melhorar a precisão dos resultados preditivos. Essa estrutura organizacional torna a ferramenta mais fácil de navegar e utilizar, o que promove maior eficiência e praticidade no processo de modelagem preditiva.

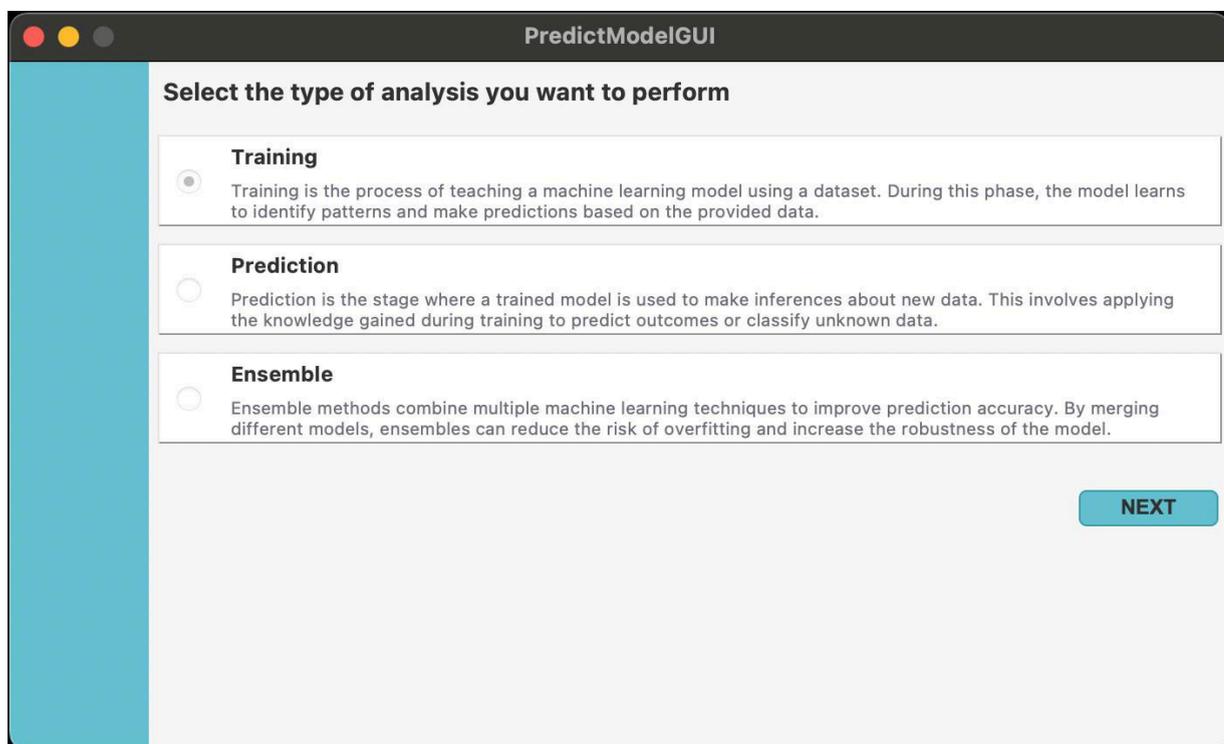


Figura 9. Janela Tipos de Análise: Esta figura exhibe as opções de análise disponíveis no *PredictModelGUI*.

A seguir, são detalhados os recursos disponíveis na interface *PredictModelGUI*, com início na etapa de treinamento. Após selecionar essa opção, o usuário deve escolher entre criar um novo projeto ou carregar um existente. Caso opte por criar um novo projeto, ao inserir seu conjunto de dados, será perguntado a porcentagem do conjunto de dados a ser utilizada no processo de teste do modelo, que é de 30% por padrão na ferramenta e pode ser alterada de acordo com as necessidades do usuário. Em seguida, será exibida uma janela com todos os modelos implementados na ferramenta, conforme mostrado na Figura 10.

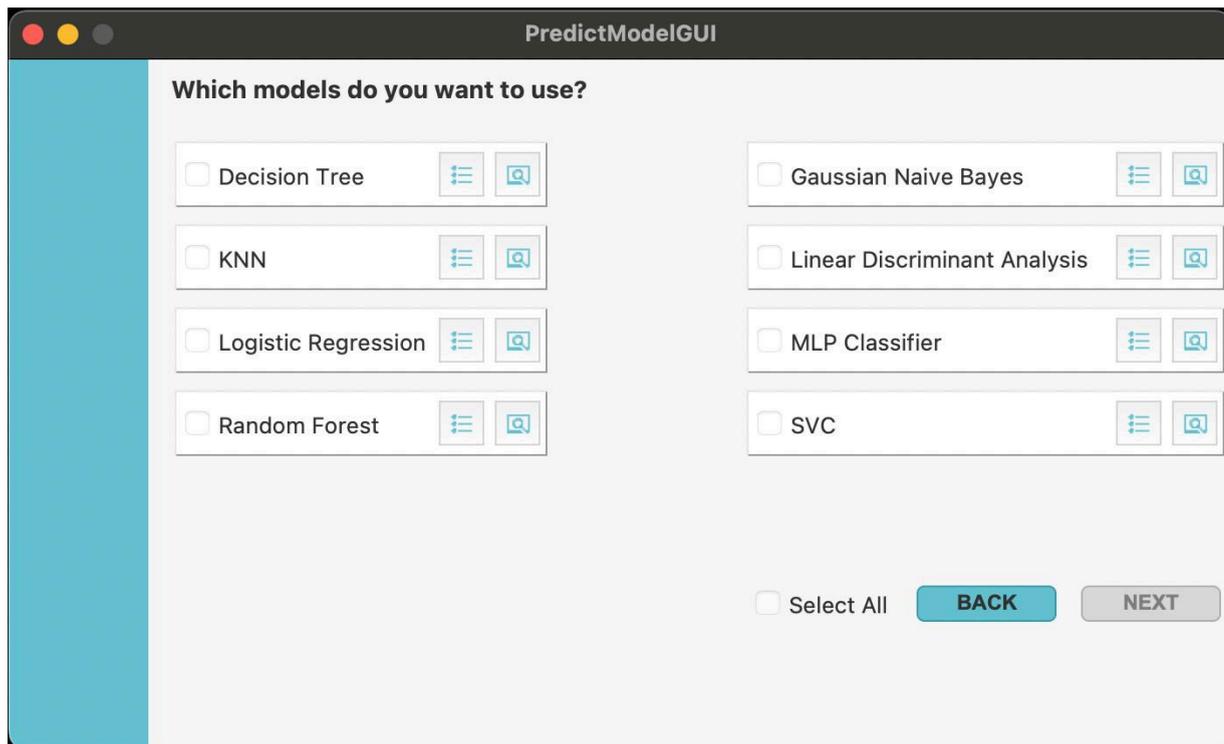


Figura 10. Lista de Modelos: A figura mostra a lista de modelos disponíveis na ferramenta *PredictModelGUI*, permite ao usuário visualizar as opções implementadas para análise.

Cabe ressaltar que todos os parâmetros podem ser modificados diretamente na interface do *software*, o que permite maior flexibilidade e personalização do processo de treinamento do modelo preditivo.

Na próxima etapa, mostrada na Figura 11, a interface exibe os parâmetros correspondentes a cada modelo selecionado na fase anterior. Nesse momento, o usuário pode optar por manter os valores padrão ou ajustá-los de acordo com as necessidades específicas de sua análise. Para facilitar essa configuração, a interface fornece acesso direto à documentação, que contém informações detalhadas sobre os parâmetros.

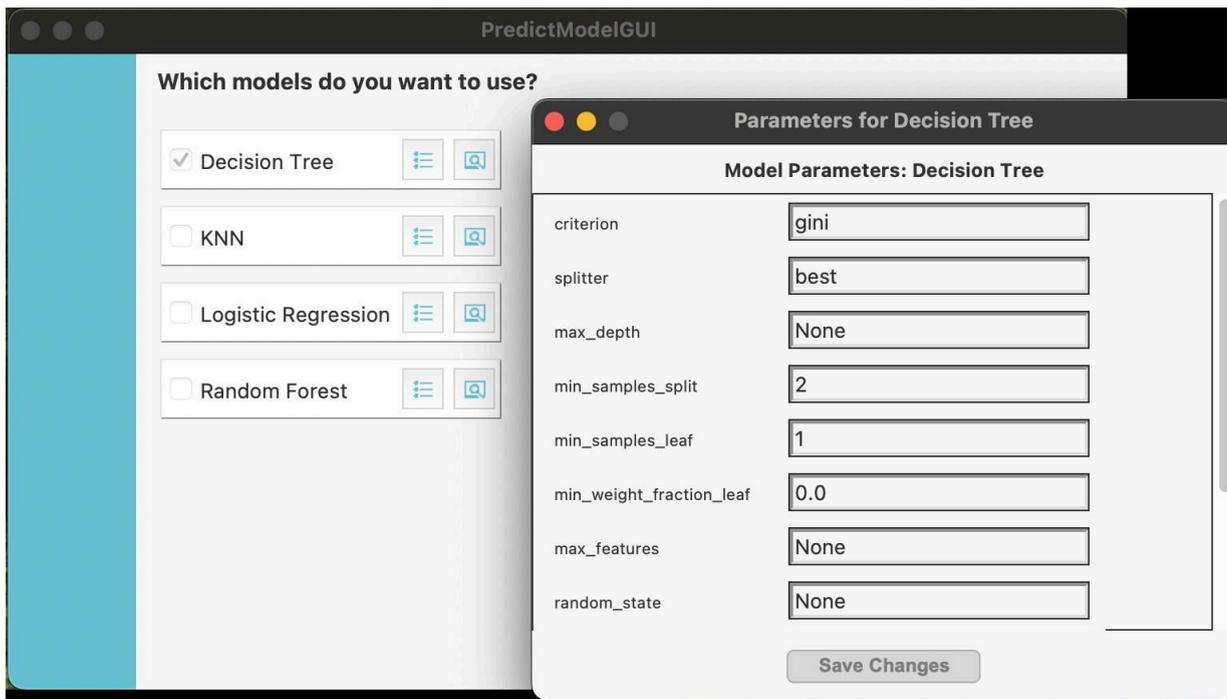


Figura 11. Lista de Parâmetros: Janela para edição de parâmetros, no exemplo são mostrados os parâmetros do modelo *Decision Tree*.

Na próxima etapa, o usuário pode selecionar as métricas que serão utilizadas para avaliar os modelos treinados. A Figura 12 mostra todas as métricas implementadas na *PredictModelGUI*, o que possibilita que o usuário defina quais serão utilizadas no processo de avaliação. Dessa forma, é possível escolher um subconjunto específico de métricas ou utilizar todas elas, a partir dos critérios adotados.

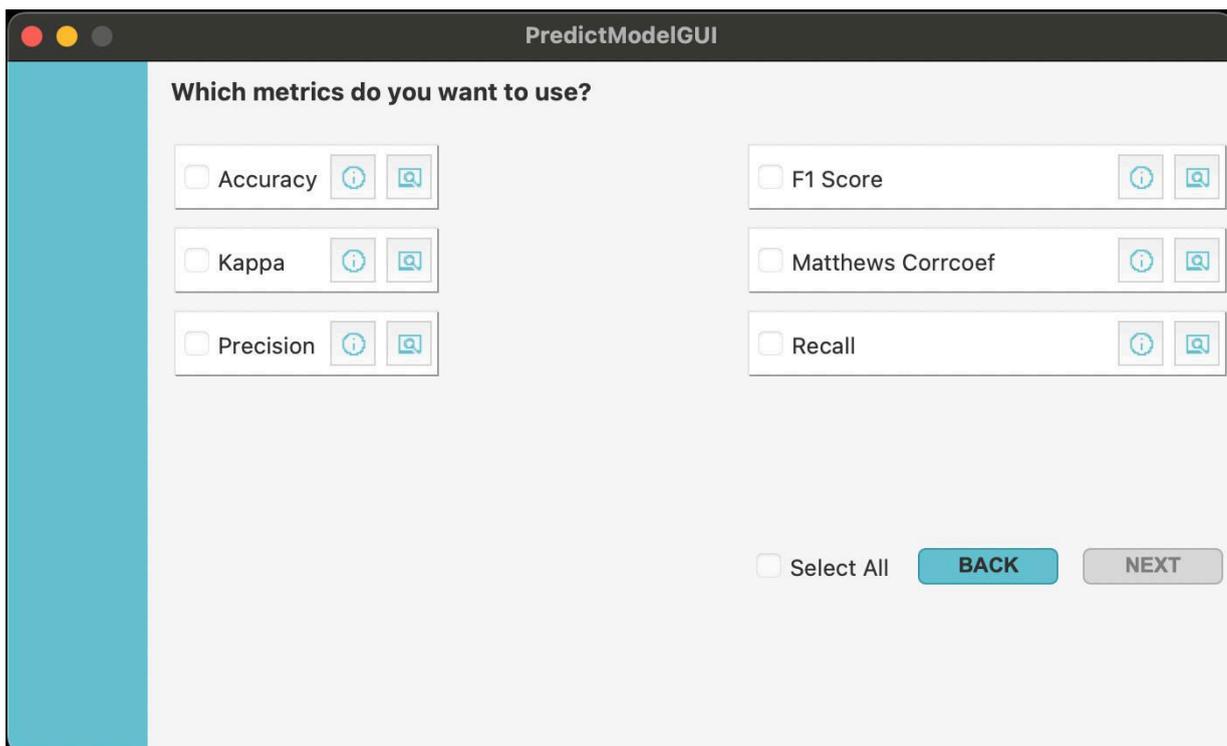


Figura 12. Lista de métricas de avaliação do modelo: Esta janela exibe todas as métricas de avaliação

disponíveis na ferramenta *PredictModelGUI*, permite ao usuário escolher as opções mais adequadas para analisar o desempenho do modelo.

Por fim, o usuário pode especificar o método de recebimento dos resultados. Se optar pelo armazenamento local, deverá indicar o diretório onde será gerada a pasta com os resultados. Alternativamente, o usuário pode optar por enviar os resultados por *e-mail*, ao informar o endereço de e-mail diretamente na interface do *software*.

Na etapa de análise, chamada Predição, o usuário pode fazer previsões a partir de um conjunto de dados em formato CSV como entrada. O *PredictModelGUI* possui um modelo padrão chamado *Ensemble*, que é compactado na primeira utilização. Portanto, antes da utilização, o *software* descompacta automaticamente o modelo e, em seguida, carrega-o na ferramenta. Esse modelo é então armazenado localmente em uma pasta específica, criada automaticamente durante o processo de descompactação.

Além disso, o equipamento utilizado deve atender aos requisitos mínimos de *hardware* para garantir o funcionamento adequado do modelo. Para facilitar essa verificação, o *PredictModelGUI* realiza automaticamente um teste de compatibilidade e informa o usuário sobre a viabilidade de utilização do modelo no dispositivo em questão.

O usuário tem a opção de utilizar um modelo previamente treinado, por meio dos procedimentos descritos na etapa de treinamento. Para isso, é necessário carregar o modelo em formato PKL, juntamente com o conjunto de dados em formato CSV. As etapas subsequentes, são compostas por edição de parâmetros, definição de métricas de avaliação e escolha do método de recebimento dos resultados, e seguem o mesmo fluxo descrito acima.

Para facilitar a padronização do arquivo de entrada na ferramenta *PredictModelGUI*, um módulo específico é fornecido para a conversão e preparação do conjunto de dados em formato CSV a partir de arquivos de anotação no formato *GenBank*. Este módulo visa garantir que os dados estejam em conformidade com os requisitos da ferramenta.

Na última etapa, correspondente à tarefa *Ensemble*, o usuário tem a opção de recriar o modelo de conjunto padrão da ferramenta através do seu conjunto de dados. Para isso, é necessário primeiro treinar os modelos individualmente, com os procedimentos descritos na etapa de treinamento. Após a conclusão desta etapa, os modelos treinados devem ser organizados no formato PKL e armazenados em um único diretório.

Na fase de *ensemble*, o usuário deve fornecer o conjunto de dados em formato CSV e indicar o diretório onde os modelos previamente treinados estão armazenados. Com base nessas informações, o modelo de conjunto será treinado.

As janelas já descritas para parametrização e definição de métricas de avaliação e resultados são as mesmas para esta etapa. É importante observar que, uma vez treinado, o modelo *ensemble* deve ser carregado sempre que o usuário desejar utilizá-lo para fazer

previsões.

Para otimizar o uso da ferramenta *PredictModelGUI*, foi desenvolvida uma [versão mobile](#), mostrada na Figura 13, disponível na *Google Play Store*. Essa versão permite que o usuário acesse as funcionalidades do *PredictModelGUI* por meio da comunicação com uma versão *Server*. Dessa forma, foi implementada uma arquitetura cliente/servidor. Após a conclusão do processo selecionado pelo usuário, uma notificação de encerramento é enviada para o endereço de *e-mail* cadastrado.

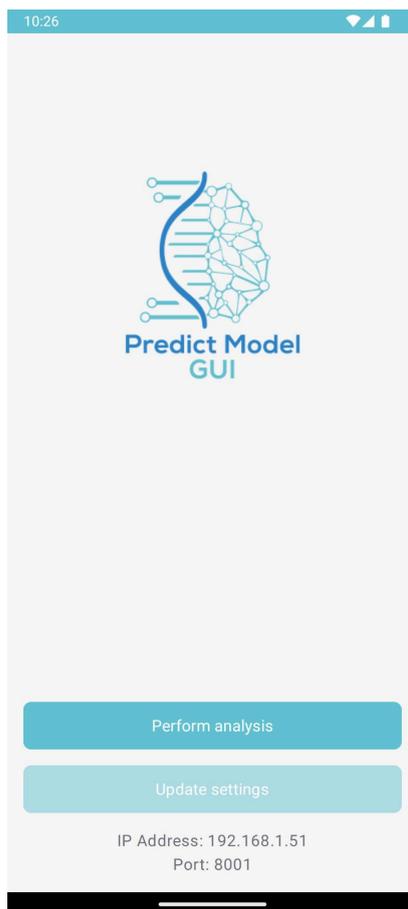


Figura 13. Versão móvel do *PredictModelGUI*: Janela principal do aplicativo *PredictModelGUI*.

O padrão visual das janelas foi mantido consistente entre as versões *desktop* e *mobile*, o que oferece uma experiência de usuário homogênea e intuitiva em ambas as plataformas.

A análise dos resultados demonstra três principais contribuições científicas do desenvolvimento da *PredictModelGUI*. Primeiramente, a metodologia inovadora adotada para a geração e o processamento do conjunto de dados melhorou a qualidade das informações, o que gerou em modelos preditivos mais precisos e confiáveis. Em segundo lugar, a implementação de uma abordagem diferenciada para a construção dos modelos de classificação permitiu uma cobertura abrangente de todas as classes genéticas essenciais, para promover melhorias significativas nas métricas de avaliação. Por fim, o desenvolvimento da ferramenta computacional permitiu a integração dos modelos em uma interface intuitiva e amigável, para reduzir a complexidade operacional para o usuário final.

Além disso, o *PredictModelGUI* possui uma versão compatível com arquiteturas cliente-servidor, o que permite a comunicação via dispositivos móveis e aumenta sua acessibilidade e aplicabilidade em diferentes contextos. Informações detalhadas sobre como executar cada etapa descrita neste manuscrito são fornecidas no manual do usuário do *PredictModelGUI*.

Por fim, é importante destacar que o desenvolvimento da ferramenta *PredictModelGUI* proporcionou uma maneira eficiente de abstrair a complexidade envolvida na construção, no treinamento e no uso de modelos de Aprendizado de Máquina. Além disso, a ferramenta amplia o acesso a essa tecnologia, especialmente para usuários sem conhecimento aprofundado de computação, e torna a aplicação mais acessível e intuitiva.

CONCLUSÃO

Neste estudo, investigou-se o uso de técnicas de Aprendizado de Máquina na classificação de Genes Essenciais. Para tanto, adotou-se uma abordagem inovadora na criação do conjunto de dados que utiliza informações de dois bancos de dados: DEG e NCBI. As informações sobre Genes Essenciais são registradas no DEG e comparadas com os dados de anotação do NCBI.

Vale ressaltar que o conjunto de dados criado exibe, ao final, a frequência de aminoácidos associada a cada produto, com um símbolo de gene que corresponde às observações do conjunto de dados. Ao criar o conjunto de dados, foi estabelecida uma estratégia para determinar o valor mínimo de proporcionalidade de representação que cada classe deveria ter. Assim, foram criados dois conjuntos de dados: um com valores iguais ou superiores ao estabelecido pela estratégia e outro com valores menores.

A próxima etapa envolveu o treinamento dos modelos com base nos conjuntos de dados criados anteriormente. No entanto, cada modelo gerado a partir de um determinado conjunto de dados contém apenas uma fração do conhecimento sobre Genes Essenciais, influenciado pela proporcionalidade das classes utilizadas para construir esses conjuntos de dados.

Essa limitação fica clara ao analisar os resultados das métricas de avaliação. A comparação dos valores obtidos para os dois conjuntos de dados mostra uma disparidade significativa entre as métricas. A título de comparação, os resultados para os *datasets* 1 e 2 foram, respectivamente: Acurácia de 86% e 6%; Precisão de 78% e 2%; F1-Score de 85% e 5%; Recall de 80% e 3%; as métricas Coeficiente de Correlação de *Matthews* (MCC) e Kappa obtiveram os mesmos valores de 86% e 6%.

Para mitigar essa discrepância, os dois modelos gerados anteriormente a partir dos conjuntos de dados foram combinados para formar um modelo final, denominado *ensemble* final. Essa técnica, amplamente utilizada em Inteligência Artificial e Aprendizado de Máquina, permite a agregação de diferentes modelos, o que proporciona um melhor desempenho

preditivo do que modelos individuais.

Essa abordagem se mostrou eficaz, como demonstrado pelas métricas de desempenho: Acurácia, MCC e Kappa atingiram 95%; Precisão atingiu 91%; F1-Score atingiu 93%; e Recall atingiu 92%. É importante observar que esses valores foram obtidos por meio do conjunto de dados de validação, não a parte de teste utilizada durante o treinamento do modelo. Esses resultados demonstram a eficiência do modelo na identificação de Genes Essenciais e reforçam a aplicação inovadora do método de Aprendizado por Conjunto de dados biológicos.

A análise e a avaliação dos modelos implementados no *PredictModelGUI* demonstraram alta eficiência na identificação de Genes Essenciais. Essa funcionalidade é aprimorada por uma interface gráfica intuitiva, que simplifica o processo de treinamento e utilização dos modelos, o que elimina a necessidade de conhecimentos avançados de computação por parte do usuário. Além disso, foi desenvolvida uma versão *mobile* da ferramenta, que envia os dados para processamento em um dispositivo remoto e, em seguida, recebe os resultados diretamente no *e-mail*, o que possibilita o uso da ferramenta de forma mais acessível, dinâmica e versátil.

É importante observar que, à medida que o usuário expande o conjunto de dados e ajusta os valores dos parâmetros, os modelos podem apresentar desempenho aprimorado, refletido em melhores métricas de avaliação. Isso indica um aumento em sua eficácia preditiva. O principal objetivo deste trabalho foi desenvolver uma ferramenta que simplifique a complexidade inerente ao desenvolvimento dos *scripts* necessários para a criação dos modelos e facilite tanto o processo de treinamento quanto a sua utilização.

AGRADECIMENTOS

Agradecimentos ao Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq) e à Universidade Federal do Pará. Este trabalho foi financiado pela PROPESP/PAPQ. Este trabalho faz parte da pesquisa realizada pelo BIOD (Grupo de Pesquisa em Bioinformática, Ômicas e Desenvolvimento) e pelo EngBio (Laboratório de Engenharia Biológica).

CONFLITOS DE INTERESSE

Os autores declaram que nenhum conflito de interesse pode ser percebido como prejudicial à imparcialidade da pesquisa relatada.

CONTRIBUIÇÕES DOS AUTORES

GSM desenvolvimento de *software* e aplicativo móvel, redação do manuscrito, redação – revisão e edição, análise formal e redação de manuais técnicos. VCS desenvolvimento de *software*, redação de manuais técnicos e testes de *software*. SSS desenvolvimento de aplicativo

móvel. WBGN teste de *software* e aplicativo móvel. RAB, DAG e AS redação – revisão e edição. CDS redação – revisão e edição e análise formal. AAOV investigação, metodologia, papéis/redação – rascunho original e desenvolvimento de *software*. Todos os autores leram e aprovaram a versão final.

REFERÊNCIAS

Burkart, N. and Huber, M.F. (2021) ‘A survey on the explainability of supervised machine learning’, *Journal of Artificial Intelligence Research*, 70, pp. 245–317. doi:10.1613/jair.1.12228.

Guigó, R. (2023) ‘Genome annotation: From human genetics to biodiversity genomics’, *Cell Genomics*, 3(8), p. 100375. doi:10.1016/j.xgen.2023.100375.

Lee JY. The Principles and Applications of High-Throughput Sequencing Technologies. *Dev Reprod.* 2023 Apr;27(1):9-24. doi: 10.12717/DR.2023.27.1.9. Epub 2023 Mar 31. PMID: 38075439; PMCID: PMC10703097.

Liang, Y. et al. (2024) ‘Recent advances in the characterization of essential genes and development of a database of essential genes’, *iMeta*, 3(1). doi:10.1002/imt2.157.

Lu, Y.; Li, M.; Gao, Z.; Ma, H.; Chong, Y.; Hong, J.; Wu, J.; Wu, D.; Xi, D.; Deng, W. Advances in Whole Genome Sequencing: Methods, Tools, and Applications in Population Genomics. *Int. J. Mol. Sci.* 2025, 26, 372. <https://doi.org/10.3390/ijms26010372>.

Mandlik, J.S., Patil, A.S. and Singh, S. (2024) ‘Next-generation sequencing (NGS): Platforms and applications’, *Journal of Pharmacy and Bioallied Sciences*, 16(Suppl 1). doi:10.4103/jpbs.jpbs_838_23.

Mukhamediev RI, Popova Y, Kuchin Y, Zaitseva E, Kalimoldayev A, Symagulov A, Levashenko V, Abdoldina F, Gopejenko V, Yakunin K, et al. Review of Artificial Intelligence and Machine Learning Technologies: Classification, Restrictions, Opportunities and Challenges. *Mathematics.* 2022; 10(15):2552. <https://doi.org/10.3390/math10152552>

Naidu, G., Zuva, T., Sibanda, E.M. (2023). A Review of Evaluation Metrics in Machine Learning Algorithms. In: Silhavy, R., Silhavy, P. (eds) Artificial Intelligence Application in Networks and Systems. CSOC 2023. Lecture Notes in Networks and Systems, vol 724. Springer, Cham. https://doi.org/10.1007/978-3-031-35314-7_2

Olufemi Aromolaran, Damilare Aromolaran, Itunuoluwa Isewon, Jelili Oyelade, Machine learning approach to gene essentiality prediction: a review, Briefings in Bioinformatics, Volume 22, Issue 5, September 2021, bbab128, <https://doi.org/10.1093/bib/bbab128>

Satam H, Joshi K, Mangrolia U, Waghoo S, Zaidi G, Rawool S, Thakare RP, Banday S, Mishra AK, Das G, et al. Next-Generation Sequencing Technology: Current Trends and Advancements. *Biology*. 2023; 12(7):997. <https://doi.org/10.3390/biology12070997>

Taishan Hu, Nilesh Chitnis, Dimitri Monos, Anh Dinh, Next-generation sequencing technologies: An overview, *Human Immunology*, Volume 82, Issue 11, 2021, Pages 801-811, ISSN 0198-8859, <https://doi.org/10.1016/j.humimm.2021.02.012>.

MATERIAL SUPLEMENTAR

Anexo 01. Lista de organismos usados para construir os conjuntos de dados: A coluna 1 lista os nomes dos organismos e a coluna 2 lista o número de acesso.

Nome do organismo	Número de acesso
<i>Corynebacterium Ammoniogenes Dsm 20306 9.6</i>	CP009244
<i>Corynebacterium Ammoniogenes Kccm 40472</i>	CP019705
<i>Corynebacterium Amycolatum Fdaargos 1107</i>	CP068169
<i>Corynebacterium Amycolatum Fdaargos 1108</i>	CP068168
<i>Corynebacterium Anserum 23h37-10</i>	CP046883
<i>Corynebacterium Aquatimens Cnctc</i>	CP068278
<i>Corynebacterium Argentoratense Dsm 44202</i>	CP006365
<i>Corynebacterium Aurimucosum Atcc 700975</i>	CP001601
<i>Corynebacterium Aurimucosum Fdaargos 1109</i>	CP068166
<i>Corynebacterium Aurimucosum Fdaargos 1110</i>	CP068165
<i>Corynebacterium Bovis Fdaargos 1052</i>	CP066067
<i>Corynebacterium Callunae Dsm 20147</i>	CP004354
<i>Corynebacterium Camporealensis Cip 105508</i>	CP027001
<i>Corynebacterium Camporealensis Dsm 44610</i>	CP011311
<i>Corynebacterium Casei Lmg S-19264</i>	CP004350
<i>Corynebacterium Coyleae Fdaargos 1425</i>	CP077302
<i>Corynebacterium Coyleae Fdaargos 1492</i>	CP083648
<i>Corynebacterium Diphtheriae 241</i>	CP003207

<i>Corynebacterium Diphtheriae 31a</i>	CP003206
<i>Corynebacterium Diphtheriae B-D-16-78</i>	CP018331
<i>Corynebacterium Diphtheriae Bh8</i>	CP003209
<i>Corynebacterium Diphtheriae Bq11</i>	CP029644
<i>Corynebacterium Diphtheriae C7(Beta)</i>	CP003210
<i>Corynebacterium Diphtheriae Cdce 8392</i>	CP003211
<i>Corynebacterium Diphtheriae Cn2000</i>	CP039522
<i>Corynebacterium Diphtheriae Dongyang</i>	CP074413
<i>Corynebacterium Diphtheriae Fdaargos 1552</i>	CP085960
<i>Corynebacterium Diphtheriae Fdaargos 1553</i>	CP085959
<i>Corynebacterium Diphtheriae Fdaargos 1592</i>	CP085994
<i>Corynebacterium Diphtheriae Fdaargos 197</i>	CP020410
<i>Corynebacterium Diphtheriae Hc01</i>	CP003212
<i>Corynebacterium Diphtheriae Th1526</i>	CP038504
<i>Corynebacterium Doosanense Cau 212 Dsm 45436</i>	CP006764
<i>Corynebacterium Falsenii Fdaargos 1493</i>	CP083647
<i>Corynebacterium Falsenii Fdaargos 1494</i>	CP083646
<i>Corynebacterium Flavescens Oj8</i>	CP009246
<i>Corynebacterium Freneyi Fdaargos 1426</i>	CP077258
<i>Corynebacterium Glaucum Dsm 30827</i>	CP019688
<i>Corynebacterium Glucuronolyticum Fdaargos 1053</i>	CP066007
<i>Corynebacterium Glucuronolyticum Fdaargos 1111</i>	CP068162
<i>Corynebacterium Glucuronolyticum Fdaargos 1190</i>	CP069485
<i>Corynebacterium Glutamicum Ar1</i>	CP007724
<i>Corynebacterium Glutamicum Atcc 13032</i>	CP025533
<i>Corynebacterium Glutamicum Atcc 13869</i>	CP016335
<i>Corynebacterium Glutamicum Atcc 14067</i>	CP022614

<i>Corynebacterium Glutamicum Atcc 21573</i>	CP068290
<i>Corynebacterium Glutamicum Atcc 21831</i>	CP007722
<i>Corynebacterium Glutamicum B253</i>	CP010451
<i>Corynebacterium Glutamicum B414</i>	CP012297
<i>Corynebacterium Glutamicum Bca</i>	CP059382
<i>Corynebacterium Glutamicum Be</i>	CP053188
<i>Corynebacterium Glutamicum C1</i>	CP017995
<i>Corynebacterium Glutamicum Cgmcc1.15647</i>	CP073911
<i>Corynebacterium Glutamicum Cicc10064</i>	CP012298
<i>Corynebacterium Glutamicum Cp</i>	CP012194
<i>Corynebacterium Glutamicum Cr101</i>	CP080542
<i>Corynebacterium Glutamicum Ha</i>	CP025534
<i>Corynebacterium Glutamicum Jh41</i>	CP041729
<i>Corynebacterium Glutamicum Mb001</i>	CP005959
<i>Corynebacterium Glutamicum Scgg1</i>	CP004047
<i>Corynebacterium Glutamicum Scgg2</i>	CP004048
<i>Corynebacterium Glutamicum Tccc11822</i>	CP020033
<i>Corynebacterium Glutamicum Tq2223</i>	CP020658
<i>Corynebacterium Glutamicum Usda-Ars-USmarc-56828</i>	CP013991
<i>Corynebacterium Glutamicum Wm001</i>	CP022394
<i>Corynebacterium Glutamicum Xv</i>	CP018175
<i>Corynebacterium Glutamicum Yi</i>	CP014984
<i>Corynebacterium Glutamicum Zl-6</i>	CP004062
<i>Corynebacterium Glyciniphilum Aj 3170</i>	CP006842
<i>Corynebacterium Halotolerans Yim 70093 Dsm 44683</i>	CP003697
<i>Corynebacterium Humireducens Nbrc 106098 Dsm 45392</i>	CP005286
<i>Corynebacterium Imitans Dsm 44264</i>	CP009211

<i>Corynebacterium Jeikeium Fdaargos 328</i>	CP022054
<i>Corynebacterium Jeikeium Fdaargos 574</i>	CP033784
<i>Corynebacterium Kefirresidentii Fdaargos 1055</i>	CP067012
<i>Corynebacterium Kroppenstedtii Dsm 44385</i>	CP001620
<i>Corynebacterium Kroppenstedtii Fdaargos 1192</i>	CP069509
<i>Corynebacterium Kutscheri Dsm 20755</i>	CP011312
<i>Corynebacterium Lactis Rw2-5</i>	CP006841
<i>Corynebacterium Liangguodongii 2184</i>	CP026948
<i>Corynebacterium Lizhenjunii Zj-599</i>	CP064954
<i>Corynebacterium Lujinxingii Zg-917</i>	CP061032
<i>Corynebacterium Macginleyi 160811</i>	CP068292
<i>Corynebacterium Macginleyi 180208</i>	CP068291
<i>Corynebacterium Macginleyi Fdaargos 1195</i>	CP069516
<i>Corynebacterium Marinum Dsm 44953</i>	CP007790
<i>Corynebacterium Matruchotii Atcc 14266</i>	CP050134
<i>Corynebacterium Minutissimum Fdaargos 1196</i>	CP069533
<i>Corynebacterium Minutissimum Fdaargos 894</i>	CP065689
<i>Corynebacterium Minutissimum Fdaargos 992</i>	CP065964
<i>Corynebacterium Mustelae Dsm 45274</i>	CP011542
<i>Corynebacterium Nuruki S6-4</i>	CP042429
<i>Corynebacterium Pelargi 136-3</i>	CP035299
<i>Corynebacterium Phocae M408-89-1</i>	CP009249
<i>Corynebacterium Propinquum Fdaargos 1112</i>	CP068161
<i>Corynebacterium Pseudotuberculosis 39</i>	CP015188
<i>Corynebacterium Pseudotuberculosis Atcc 19410</i>	CP021251
<i>Corynebacterium Pseudotuberculosis Mb66</i>	CP013263
<i>Corynebacterium Striatum 216</i>	CP024932

<i>Corynebacterium Striatum Kc-Na-01</i>	CP021252
<i>Escherichia Coli Isolate Eccnb12-2</i>	NZ_CP033635
<i>Escherichia Coli Strain 06-00048</i>	NZ_CP015229
<i>Escherichia Coli Strain 06-3462</i>	NZ_CP034794
<i>Escherichia Coli Strain 08-3914</i>	NZ_CP034808
<i>Escherichia Coli Strain 105</i>	NZ_CP028700
<i>Escherichia Coli Strain 106</i>	NZ_CP028698
<i>Escherichia Coli Strain 107</i>	NZ_CP028695
<i>Escherichia Coli Strain 108</i>	NZ_CP028693
<i>Escherichia Coli Strain 109</i>	NZ_CP028690
<i>Escherichia Coli Strain 110</i>	NZ_CP028687
<i>Escherichia Coli Strain 111</i>	NZ_CP028685
<i>Escherichia Coli Strain 112</i>	NZ_CP028683
<i>Escherichia Coli Strain 113</i>	NZ_CP028680
<i>Escherichia Coli Strain 114</i>	NZ_CP028677
<i>Escherichia Coli Strain 116</i>	NZ_CP028671
<i>Escherichia Coli Strain 117</i>	NZ_CP028668
<i>Escherichia Coli Strain 118</i>	NZ_CP028665
<i>Escherichia Coli Strain 119</i>	NZ_CP028662
<i>Escherichia Coli Strain 120</i>	CP028659
<i>Escherichia Coli Strain 121</i>	NZ_CP028656
<i>Escherichia Coli Strain 122</i>	NZ_CP028654
<i>Escherichia Coli Strain 123</i>	NZ_CP028652
<i>Escherichia Coli Strain 124</i>	NZ_CP028650
<i>Escherichia Coli Strain 130</i>	NZ_CP028647
<i>Escherichia Coli Strain 131</i>	NZ_CP028644
<i>Escherichia Coli Strain 132</i>	NZ_CP028641

<i>Escherichia Coli Strain 133</i>	NZ_CP028638
<i>Escherichia Coli Strain 134</i>	NZ_CP028635
<i>Escherichia Coli Strain 135</i>	NZ_CP028632
<i>Escherichia Coli Strain 136</i>	NZ_CP028629
<i>Escherichia Coli Strain 137</i>	NZ_CP028626
<i>Escherichia Coli Strain 138</i>	NZ_CP028623
<i>Escherichia Coli Strain 139</i>	NZ_CP028620
<i>Escherichia Coli Strain 140</i>	NZ_CP028617
<i>Escherichia Coli Strain 141</i>	NZ_CP028614
<i>Escherichia Coli Strain 143</i>	NZ_CP028607
<i>Escherichia Coli Strain 144</i>	NZ_CP028603
<i>Escherichia Coli Strain 150</i>	NZ_CP028592
<i>Escherichia Coli Strain 155</i>	NZ_CP018237
<i>Escherichia Coli Strain 2009c-3554</i>	CP034803
<i>Escherichia Coli Strain 2009c-4687</i>	NZ_CP034799
<i>Escherichia Coli Strain 2010c-3142</i>	NZ_CP034801
<i>Escherichia Coli Strain 2010c-3347</i>	NZ_CP034806
<i>Escherichia Coli Strain 2011c-4251</i>	NZ_CP027388
<i>Escherichia Coli Strain 2013c-3264</i>	NZ_CP027544
<i>Escherichia Coli Strain 2013c-3277</i>	NZ_CP027331
<i>Escherichia Coli Strain 2013c-3342</i>	NZ_CP027766
<i>Escherichia Coli Strain 2013c-3513</i>	NZ_CP027555
<i>Escherichia Coli Strain 2013c-4187</i>	NZ_CP027546
<i>Escherichia Coli Strain 2014c-3051</i>	NZ_CP027338
<i>Escherichia Coli Strain 2014c-3655</i>	NZ_CP027351
<i>Escherichia Coli Strain 2015c-3163</i>	NZ_CP027219
<i>Escherichia Coli Strain 20r2r</i>	CP062160

<i>Escherichia Coli Strain 266917 2</i>	NZ_CP026723
<i>Escherichia Coli Strain 272</i>	NZ_CP018239
<i>Escherichia Coli Strain 28rc1</i>	NZ_CP015020
<i>Escherichia Coli Strain 319</i>	NZ_CP018241
<i>Escherichia Coli Strain 472</i>	NZ_CP018245
<i>Escherichia Coli Strain 9000</i>	NZ_CP018252
<i>Escherichia Coli Strain 95nr1</i>	NZ_CP021339
<i>Escherichia Coli Strain Atcc 43889</i>	NZ_CP015853
<i>Escherichia Coli Strain Ausmdu00002545</i>	NZ_CP045975
<i>Escherichia Coli Strain Ausmdu00014361</i>	NZ_CP045827
<i>Escherichia Coli Strain Cau16175</i>	NZ_CP047378
<i>Escherichia Coli Strain Cfsan027346</i>	NZ_CP037945
<i>Escherichia Coli Strain Cfsan027350</i>	NZ_CP037941
<i>Escherichia Coli Strain Cv261</i>	NZ_CP040316
<i>Escherichia Coli Strain Dsm 103246</i>	NZ_CP019944
<i>Escherichia Coli Strain E2855</i>	NZ_AP018796
<i>Escherichia Coli Strain Erl04-3476</i>	NZ_CP032808
<i>Escherichia Coli Strain Erl05-1306</i>	NZ_CP032803
<i>Escherichia Coli Strain Erl06-2442</i>	NZ_CP032801
<i>Escherichia Coli Strain Erl06-2497</i>	NZ_CP032797
<i>Escherichia Coli Strain Erl06-2503</i>	NZ_CP032795
<i>Escherichia Coli Strain F1 E4</i>	NZ_CP040307
<i>Escherichia Coli Strain Forc 044</i>	NZ_CP016755
<i>Escherichia Coli Strain Fsis11705876</i>	NZ_CP035545
<i>Escherichia Coli Strain Hb6</i>	NZ_CP040305
<i>Escherichia Coli Strain Ma11</i>	NZ_CP040314

<i>Escherichia Coli Strain Nadc 5570-86-24-6564 Isolate Wild Type</i>	NZ_CP017251
<i>Escherichia Coli Strain Nadc 5570-86-24-6565 Isolate Mutant</i>	NZ_CP017249
<i>Escherichia Coli Strain Nctc9112</i>	NZ_LR134079
<i>Escherichia Coli Strain Nzrm4169</i>	NZ_CP032789
<i>Escherichia Coli Strain Pa20</i>	NZ_CP017669
<i>Escherichia Coli Strain Pk8241</i>	NZ_CP080139
<i>Escherichia Coli Strain Rm13322</i>	NZ_CP050498
<i>Escherichia Coli Strain Rm19259</i>	NZ_CP046527
<i>Escherichia Coli Strain Scaid Wnd1-2021 (1-128)</i>	NZ_CP082831
<i>Escherichia Coli Strain Sj7</i>	NZ_CP044315
<i>Escherichia Coli Strain Srcc 1675</i>	NZ_CP015023
<i>Escherichia Coli Strain Stec2018-553</i>	NZ_CP075665
<i>Klebsiella Aerogenes 035</i>	CP050067
<i>Klebsiella Aerogenes 18-2341</i>	CP049600
<i>Klebsiella Aerogenes Ar 0007</i>	CP024883
<i>Klebsiella Aerogenes Ar 0018</i>	CP024880
<i>Klebsiella Aerogenes Ar 0062</i>	CP026756
<i>Klebsiella Aerogenes Ar0009</i>	CP024885
<i>Klebsiella Aerogenes Auh-Kam-9</i>	CP048598
<i>Klebsiella Aerogenes C9</i>	CP042530
<i>Klebsiella Aerogenes Cav1320</i>	CP011574
<i>Klebsiella Aerogenes Ea46506</i>	CP070520
<i>Klebsiella Aerogenes Fdaargos 1442</i>	CP077293
<i>Klebsiella Aerogenes Fdaargos 327</i>	CP031756
<i>Klebsiella Aerogenes Fdaargos 363</i>	CP023963
<i>Klebsiella Aerogenes G7</i>	CP011539

<i>Klebsiella Aerogenes Hnhf1</i>	CP047669
<i>Klebsiella Aerogenes Ka P10 L5 03.19</i>	CP044214
<i>Klebsiella Aerogenes Ka37751</i>	CP041925
<i>Klebsiella Aerogenes Kae3sp</i>	CP082898
<i>Klebsiella Aerogenes Lu2</i>	CP035466
<i>Klebsiella Aerogenes Ntt31xs Chromosome Ntt31xs-1</i>	CP077429
<i>Klebsiella Aerogenes Rhbstw-00898</i>	CP056260
<i>Klebsiella Aerogenes Rhbstw-00938</i>	CP055904
<i>Klebsiella Aerogenes Y1</i>	CP045870
<i>Klebsiella Aerogenes Y3</i>	CP045869
<i>Klebsiella Aerogenes Y6</i>	CP045868
<i>Klebsiella Africana 200023</i>	CP084874
<i>Klebsiella Africana Ff1003</i>	CP059391
<i>Klebsiella Grimontii 2481359</i>	CP067380
<i>Klebsiella Grimontii 5512.56</i>	CP067391
<i>Klebsiella Grimontii Kox 60</i>	CP067433
<i>Klebsiella Grimontii M5a1</i>	CP063301
<i>Klebsiella Grimontii Ny3</i>	CP079754
<i>Klebsiella Grimontii Rhbstw-00039</i>	CP058180
<i>Klebsiella Grimontii Rhbstw-00165</i>	CP055364
<i>Klebsiella Grimontii Rhbstw-00853</i>	CP056150
<i>Klebsiella Grimontii Rhbstw-00866</i>	CP055309
<i>Klebsiella Grimontii Ss141</i>	CP044527
<i>Klebsiella Huaxiensis Wchki090001</i>	CP036175
<i>Klebsiella Michiganensis 12084</i>	CP072119
<i>Klebsiella Michiganensis 53828cz</i>	CP085764
<i>Klebsiella Michiganensis 7525</i>	CP065474

<i>Klebsiella Michiganensis 8-1</i>	CP089448
<i>Klebsiella Michiganensis Akkl-001</i>	CP060111
<i>Klebsiella Michiganensis Ar375</i>	CP029141
<i>Klebsiella Michiganensis B106</i>	CP067093
<i>Klebsiella Michiganensis Bd-50-Km</i>	CP061930
<i>Klebsiella Michiganensis Bd177</i>	CP048108
<i>Klebsiella Michiganensis C52</i>	CP042545
<i>Klebsiella Michiganensis Ccri-24235</i>	CP081351
<i>Klebsiella Michiganensis Cz598</i>	CP073305
<i>Klebsiella Michiganensis Fdaargos 647</i>	CP044109
<i>Klebsiella Michiganensis Jnqh491</i>	CP075881
<i>Klebsiella Michiganensis K516</i>	CP022348
<i>Klebsiella Michiganensis K518</i>	CP023185
<i>Klebsiella Michiganensis K92</i>	CP089315
<i>Klebsiella Michiganensis Km41</i>	CP090078
<i>Klebsiella Michiganensis Kmfe267</i>	CP071393
<i>Klebsiella Michiganensis Knu07</i>	CP041515
<i>Klebsiella Michiganensis Ko 408</i>	CP091470
<i>Klebsiella Michiganensis Kox101</i>	CP089407
<i>Klebsiella Michiganensis Kox58</i>	CP089395
<i>Klebsiella Michiganensis Lds17</i>	CP065338
<i>Klebsiella Michiganensis M1</i>	CP008841
<i>Klebsiella Michiganensis M82255</i>	CP035214
<i>Klebsiella Michiganensis Rhb20-C02</i>	CP058212
<i>Klebsiella Michiganensis Rhbstw-00409</i>	CP055325
<i>Klebsiella Michiganensis Sb-24</i>	CP073236
<i>Klebsiella Michiganensis X2-1</i>	CP051427

<i>Klebsiella Oxytoca Ar 0147</i>	CP020358
<i>Klebsiella Oxytoca Cav1335</i>	CP011618
<i>Klebsiella Oxytoca Fdaargos 335</i>	CP027426
<i>Klebsiella Oxytoca Konih4</i>	CP026269
<i>Klebsiella Oxytoca Rhbstw-00493</i>	CP056453
<i>Klebsiella Pneumoniae Ba2275</i>	CP053364
<i>Klebsiella Pneumoniae F1</i>	CP026130
<i>Klebsiella Pneumoniae F138</i>	CP026149
<i>Klebsiella Pneumoniae F5</i>	CP026132
<i>Klebsiella Pneumoniae F77</i>	CP026136
<i>Klebsiella Pneumoniae Fdaargos 775</i>	CP040993
<i>Klebsiella Pneumoniae Frpdr</i>	CP063759
<i>Klebsiella Pneumoniae Kp1517</i>	CP072463
<i>Klebsiella Pneumoniae Kp18-1</i>	CP082001
<i>Klebsiella Pneumoniae Kp18-2079</i>	CP048933
<i>Klebsiella Pneumoniae Kp19-2029</i>	CP047160
<i>Klebsiella Pneumoniae Kp36</i>	CP047192
<i>Klebsiella Pneumoniae Kpc-2</i>	CP078122
<i>Klebsiella Pneumoniae P1428</i>	CP017994
<i>Klebsiella Pneumoniae Qd23</i>	CP042858
<i>Klebsiella Quasipneumoniae A708</i>	CP026368
<i>Klebsiella Quasipneumoniae G4584</i>	CP034129
<i>Klebsiella Quasipneumoniae Kqpf26</i>	CP065838
<i>Klebsiella Quasipneumoniae Ww-14a</i>	CP080099
<i>Klebsiella Quasipneumoniae Ydkl-002</i>	CP068237
<i>Klebsiella Quasivariicola 08a119</i>	CP084768
<i>Klebsiella Quasivariicola Kpn1705</i>	CP022823

<i>Klebsiella Variicola</i> 342	CP000964
<i>Klebsiella Variicola</i> Dsm 15968	CP010523
<i>Klebsiella Variicola</i> F2r9	CP072130
<i>Klebsiella Variicola</i> Fh-1	CP054254
<i>Klebsiella Variicola</i> X39	CP018307
<i>Mycobacterium Avium</i> 104	CP000479
<i>Mycobacterium Avium</i> Dsm44156	CP046507
<i>Mycobacterium Avium</i> Fdaargos 1606	CP085978
<i>Mycobacterium Avium</i> Fdaargos 1607	CP085977
<i>Mycobacterium Avium</i> Hjw	CP028731
<i>Mycobacterium Avium</i> Rcad0278	CP016396
<i>Mycobacterium Avium</i> Subsp <i>Hominissuis</i> 101034	CP040247
<i>Mycobacterium Avium</i> Subsp <i>Hominissuis</i> 101115	CP040255
<i>Mycobacterium Avium</i> Subsp <i>Hominissuis</i> 101174	CP040250
<i>Mycobacterium Avium</i> Subsp <i>Hominissuis</i> Cam177	CP076851
<i>Mycobacterium Avium</i> Subsp <i>Hominissuis</i> Hp17	CP016818
<i>Mycobacterium Avium</i> Subsp <i>Hominissuis</i> Mac109	CP029332
<i>Mycobacterium Avium</i> Subsp <i>Hominissuis</i> Mah11	CP035744
<i>Mycobacterium Avium</i> Subsp <i>Hominissuis</i> Mc2 2500	CP036220
<i>Mycobacterium Avium</i> Subsp <i>Hominissuis</i> W14	CP060407
<i>Mycobacterium Avium</i> Subsp <i>Hominissuis</i> W9	CP060405
<i>Mycobacterium Avium</i> Subsp <i>Homonissuis</i> H87	CP018363
<i>Mycobacterium Avium</i> Subsp <i>Paratuberculosis</i> Dsm 44135	CP053068
<i>Mycobacterium Avium</i> Subsp <i>Paratuberculosis</i> Telford	CP033688
<i>Mycobacterium Colombiense</i> Cect 3035	CP020821
<i>Mycobacterium Diernhoferi</i> Atcc 19340	CP080332
<i>Mycobacterium Frederiksbergense</i> Lb 501t	CP038799

<i>Mycobacterium Goodii X7b</i>	CP012150
<i>Mycobacterium Gordonae X7091</i>	CP070973
<i>Mycobacterium Grossiae Dsm 104744</i>	CP043474
<i>Mycobacterium Haemophilum Dsm 44634 Atcc 29548</i>	CP011883
<i>Mycobacterium Heraklionense Jcm 30995</i>	CP080997
<i>Mycobacterium Holsaticum Jcm 12374</i>	CP080998
<i>Mycobacterium Intracellulare Atcc 13950</i>	CP003322
<i>Mycobacterium Intracellulare Fdaargos 1564</i>	CP085945
<i>Mycobacterium Intracellulare Fdaargos 1611</i>	CP089221
<i>Mycobacterium Intracellulare Fdaargos 1612</i>	CP089220
<i>Mycobacterium Intracellulare Fdaargos 1616</i>	CP089231
<i>Mycobacterium Intracellulare Flac0133</i>	CP023146
<i>Mycobacterium Intracellulare Flac0181</i>	CP023149
<i>Mycobacterium Intracellulare Mott-02</i>	CP003323
<i>Mycobacterium Kansasii Fdaargos 1614</i>	CP089218
<i>Mycobacterium Kansasii Fdaargos 1615</i>	CP089216
<i>Mycobacterium Kansasii 10mk</i>	CP019886
<i>Mycobacterium Kansasii 11mk</i>	CP019887
<i>Mycobacterium Kansasii 4mk</i>	CP019884
<i>Mycobacterium Kansasii Atcc 12478</i>	CP006835
<i>Mycobacterium Kansasii Fdaargos 1534</i>	CP084364
<i>Mycobacterium Kubicae Jcm 13573</i>	CP065047
<i>Mycobacterium Kubicae Njh Mkub1</i>	CP045081
<i>Mycobacterium Kubicae Njh Mkub2</i>	CP045075
<i>Mycobacterium Leprae Mrhru-235-G</i>	CP029543
<i>Mycobacterium Malmoense Atcc 29571</i>	CP080999
<i>Mycobacterium Marinum 1218r</i>	CP025779

<i>Mycobacterium Marinum Ccug20998</i>	CP024190
<i>Mycobacterium Marinum Mma1</i>	CP058277
<i>Mycobacterium Marseillense Flac0026</i>	CP023147
<i>Mycobacterium Pallens Jcm 16370</i>	CP080333
<i>Mycobacterium Paragordoniae 49061</i>	CP025546
<i>Mycobacterium Riyadhense Ntm</i>	CP045092
<i>Mycobacterium Senegalense Atcc 35796</i>	CP081000
<i>Mycobacterium Shigaense Scy</i>	CP022927
<i>Mycobacterium Shottsii M175</i>	CP014860
<i>Mycobacterium Tuberculosis 2-0029p6c4</i>	CP041837
<i>Mycobacterium Tuberculosis 2-0046p6c4</i>	CP041833
<i>Mycobacterium Tuberculosis 060827</i>	CP058236
<i>Mycobacterium Tuberculosis 1-0107p6c4</i>	CP041853
<i>Mycobacterium Tuberculosis 11502</i>	CP070338
<i>Mycobacterium Tuberculosis 120-26cao</i>	CP071127
<i>Mycobacterium Tuberculosis 2-0022p6c4</i>	CP041840
<i>Mycobacterium Tuberculosis 2-0023p6c4</i>	CP041839
<i>Mycobacterium Tuberculosis 267-47w148</i>	CP071128
<i>Mycobacterium Tuberculosis 3-0096p6c4</i>	CP041827
<i>Mycobacterium Tuberculosis 3-0124p6c4</i>	CP041826
<i>Mycobacterium Tuberculosis Beijing-Like-36918</i>	CP017594
<i>Mycobacterium Tuberculosis Beijing-Like-50148</i>	CP017597
<i>Mycobacterium Tuberculosis Cg20</i>	CP072764
<i>Mycobacterium Tuberculosis Cg21</i>	CP072763
<i>Mycobacterium Tuberculosis Cg23</i>	CP072762
<i>Mycobacterium Tuberculosis Cg24</i>	CP072761
<i>Mycobacterium Tuberculosis Dk9897</i>	CP018778

<i>Mycobacterium Tuberculosis F28</i>	CP010330
<i>Mycobacterium Tuberculosis Fdaargos 751</i>	CP046308
<i>Mycobacterium Tuberculosis Gg-111-10</i>	CP025593
<i>Mycobacterium Tuberculosis Gg-37-11</i>	CP025598
<i>Mycobacterium Tuberculosis Gg-90-10</i>	CP025601
<i>Mycobacterium Tuberculosis H107</i>	CP019612
<i>Mycobacterium Tuberculosis H37rv Cg</i>	CP072765
<i>Mycobacterium Tuberculosis H54</i>	CP019610
<i>Mycobacterium Tuberculosis H83</i>	CP019611
<i>Mycobacterium Tuberculosis Mt-0080</i>	CP041207
<i>Mycobacterium Tuberculosis Mtb1</i>	CP020381
<i>Mycobacterium Tuberculosis Mtb2</i>	CP022014
<i>Mycobacterium Tuberculosis Ocu901s S2 2s</i>	CP018014
<i>Mycobacterium Tuberculosis Rus B0</i>	CP030093
<i>Mycobacterium Tuberculosis Scaid 187.0</i>	CP012506
<i>Mycobacterium Tuberculosis Scaid 252-0</i>	CP016888
<i>Mycobacterium Tuberculosis Sea07010354p6c4</i>	CP041813
<i>Mycobacterium Tuberculosis Sit745-Eai1-Mys</i>	CP046529
<i>Mycobacterium Tuberculosis Tb282</i>	CP017920
<i>Mycobacterium Tuberculosis Tcdc11</i>	CP046728
<i>Mycobacterium Tuberculosis Wc059</i>	CP022578
<i>Mycobacterium Tuberculosis Wc078</i>	CP022577
<i>Mycobacterium Tuberculosis Zmc13-88</i>	CP009101
<i>Rhodococcus Aetherivorans Isolate Psbb011</i>	NZ_CP069306
<i>Rhodococcus Aetherivorans Strain Cbo21-1</i>	NZ_CP088969
<i>Rhodococcus Erythropolis Ccm2595</i>	CP003761
<i>Rhodococcus Erythropolis Pr4</i>	AP008957

<i>Rhodococcus Erythropolis R138</i>	CP007255
<i>Rhodococcus Erythropolis Strain Bg43</i>	CP011295
<i>Rhodococcus Erythropolis Strain Kb1</i>	NZ_CP050124
<i>Rhodococcus Erythropolis Strain R85</i>	NZ_CP070870
<i>Rhodococcus Erythropolis Strain X5</i>	NZ_CP044284
<i>Rhodococcus Fascians A21d2</i>	NZ_CP049748
<i>Rhodococcus Fascians A25f</i>	NZ_CP049744
<i>Rhodococcus Fascians D188</i>	CP015235
<i>Rhodococcus Fascians Strain Pbts 2</i>	NZ_CP015220
<i>Rhodococcus Globerulus Strain D757</i>	CP079698
<i>Rhodococcus Hoagii Jcm94-14</i>	AP024181
<i>Rhodococcus Hoagii Jcm94-16</i>	AP024183
<i>Rhodococcus Hoagii Jcm94-25</i>	AP024185
<i>Rhodococcus Hoagii Jcm94-27</i>	AP024187
<i>Rhodococcus Hoagii Jid03-27</i>	AP024194
<i>Rhodococcus Hoagii Jid03-46</i>	AP024196
<i>Rhodococcus Hoagii Jid03-56</i>	AP024198
<i>Rhodococcus Hoagii Strain Dsskp-R-001</i>	CP027793
<i>Rhodococcus Hoagii Strain Fdaargos 952</i>	CP065594
<i>Rhodococcus Hoagii Strain Jcm94-3</i>	NZ_AP024192
<i>Rhodococcus Hoagii Strain Jcm94-31</i>	NZ_AP024189
<i>Rhodococcus Hoagii Strain Wy</i>	NZ_CP041647
<i>Rhodococcus Koreensis Strain R85</i>	NZ_CP070609
<i>Rhodococcus Opacus Pd630</i>	NZ_CP080954
<i>Rhodococcus Opacus Strain Icp</i>	CP009111
<i>Rhodococcus Opacus Strain Dsm 44186</i>	CP082160
<i>Rhodococcus Pyridinivorans Sb3094</i>	CP006996

<i>Rhodococcus Pyridinivorans Strain 5ap</i>	NZ_CP063450
<i>Rhodococcus Pyridinivorans Strain B403</i>	NZ_CP066853
<i>Rhodococcus Pyridinivorans Strain Gf3</i>	NZ_CP022915
<i>Rhodococcus Pyridinivorans Strain Tg9</i>	CP022208
<i>Rhodococcus Pyridinivorans Strain Yc-Jh2</i>	NZ_CP050178
<i>Rhodococcus Pyridinivorans Strain Yf3</i>	NZ_CP040719
<i>Rhodococcus Qingshengii Cs98</i>	AP023172
<i>Rhodococcus Qingshengii Strain 7b</i>	CP063234
<i>Rhodococcus Qingshengii Strain Cl-05</i>	CP072108
<i>Rhodococcus Qingshengii Strain Cx-1</i>	CP054207
<i>Rhodococcus Qingshengii Strain R11</i>	CP042917
<i>Rhodococcus Qingshengii Strain Tg-1</i>	CP077417
<i>Rhodococcus Rhodochrous Strain Atcc Baa870</i>	NZ_CP032675
<i>Rhodococcus Rhodochrous Strain Bx2</i>	NZ_CP027557
<i>Rhodococcus Ruber Strain C1</i>	NZ_CP044211
<i>Rhodococcus Ruber Strain P14</i>	NZ_CP024315
<i>Rhodococcus Ruber Strain R1</i>	NZ_CP038030
<i>Rhodococcus Ruber Strain Sd3</i>	NZ_CP029146
<i>Rhodococcus Ruber Strain Yc-Yt1</i>	NZ_CP023714
<i>Rhodococcus Ruber Strain Yyl</i>	NZ_CP024890
<i>Rhodococcus Sp. Djl-6-2</i>	CP025959
<i>Rhodococcus Triatomae Strain Dsm 44892</i>	NZ_CP048814
<i>Rhodococcus Triatomae Strain Dsm 44893</i>	NZ_CP048813
<i>Salmonella Bongori Cfsan000510</i>	CP074233
<i>Salmonella Bongori N268-08</i>	CP006608
<i>Salmonella Bongori Nctc 12419</i>	FR877557
<i>Salmonella Bongori Serovar 40 Z35 - Strain Cfsan001045</i>	NZ_CP074592

<i>Salmonella Bongori</i> Serovar 48 Z41 -- Str. Rks3044	CP006692
<i>Salmonella Bongori</i> Serovar 48 Z81 - Strain 08-0158	CP053336
<i>Salmonella Bongori</i> Serovar 66 Z41 - Str. Sa19983605	CP022120
<i>Salmonella Bongori</i> Strain 04-0440	CP035676
<i>Salmonella Bongori</i> Strain 85-0051	CP053416
<i>Salmonella Bongori</i> Strain 92-0238	CP053417
<i>Salmonella Bongori</i> Strain Se40	CP067369
<i>Salmonella Enterica</i> Sehaa3795	AP020330
<i>Salmonella Enterica</i> Sesen3709	AP020332
<i>Salmonella Enterica</i> Strain 2011k-1440	NZ_CP053585
<i>Salmonella Enterica</i> Strain 2012k-0845	NZ_CP053579
<i>Salmonella Enterica</i> Strain 2014k-1020	NZ_CP053584
<i>Salmonella Enterica</i> Strain 85-0120	CP054715
<i>Salmonella Enterica</i> Strain 94-0093	NZ_CP053581
<i>Salmonella Enterica</i> Strain C629	NZ_CP015724
<i>Salmonella Enterica</i> Strain Cfsan033950	NZ_CP075141
<i>Salmonella Enterica</i> Strain Cfsan044865	NZ_CP075140
<i>Salmonella Enterica</i> Strain Cfsan044875	NZ_CP075138
<i>Salmonella Enterica</i> Strain Cfsan044885	NZ_CP075137
<i>Salmonella Enterica</i> Strain Cfsan044888	NZ_CP075135
<i>Salmonella Enterica</i> Strain Cfsan044925	NZ_CP075127
<i>Salmonella Enterica</i> Strain Cfsan044945	NZ_CP075125
<i>Salmonella Enterica</i> Strain Cfsan060804	CP075110
<i>Salmonella Enterica</i> Strain Cfsan060807	NZ_CP075109
<i>Salmonella Enterica</i> Strain Cfsan060808	NZ_CP075108
<i>Salmonella Enterica</i> Strain Cfsan060809	NZ_CP075107
<i>Salmonella Enterica</i> Strain Cfsan064033	NZ_CP028172

<i>Salmonella Enterica Strain Cfsan064034</i>	NZ_CP028169
<i>Salmonella Enterica Strain Cfsan064276</i>	CP075106
<i>Salmonella Enterica Strain Cfsan096147</i>	NZ_CP044257
<i>Salmonella Enterica Strain Da34821</i>	NZ_CP029567
<i>Salmonella Enterica Strain Fdaargos 1066</i>	NZ_CP066047
<i>Salmonella Enterica Strain Fdaargos 1067</i>	NZ_CP066009
<i>Salmonella Enterica Strain Fdaargos 1271</i>	NZ_CP069518
<i>Salmonella Enterica Strain Fdaargos 687</i>	NZ_CP046283
<i>Salmonella Enterica Strain Fdaargos 688</i>	NZ_CP046280
<i>Salmonella Enterica Strain Fdaargos 70</i>	NZ_CP026052
<i>Salmonella Enterica Strain Fdaargos 707</i>	NZ_CP046279
<i>Salmonella Enterica Strain Fdaargos 708</i>	NZ_CP046278
<i>Salmonella Enterica Strain Fdaargos 710</i>	NZ_CP046277
<i>Salmonella Enterica Strain Fdaargos 712</i>	NZ_CP046291
<i>Salmonella Enterica Strain Fdaargos 765</i>	NZ_CP041011
<i>Salmonella Enterica Strain Fdaargos 768</i>	NZ_CP041005
<i>Salmonella Enterica Strain Fdaargos 878</i>	CP065718
<i>Salmonella Enterica Strain Fdaargos 928</i>	NZ_CP065639
<i>Salmonella Enterica Strain Forc 019</i>	NZ_CP012396
<i>Salmonella Enterica Strain Forc 038</i>	NZ_CP015574
<i>Salmonella Enterica Strain Forc 051</i>	NZ_CP017232
<i>Salmonella Enterica Strain Forc 056</i>	NZ_CP017177
<i>Salmonella Enterica Strain Forc 074</i>	NZ_CP023436
<i>Salmonella Enterica Strain Forc 078</i>	NZ_CP026713
<i>Salmonella Enterica Strain Fsw0104</i>	NZ_CP037894
<i>Salmonella Enterica Strain Gsj 2017-Sal.-009</i>	NZ_CP050833
<i>Salmonella Enterica Strain Gsj2016-Sal.-018</i>	NZ_CP069166

<i>Salmonella Enterica Strain Gx1006</i>	NZ_CP060585
<i>Salmonella Enterica Strain K Sa184</i>	NZ_CP061159
<i>Salmonella Enterica Strain Lhica E3</i>	NZ_CP079839
<i>Salmonella Enterica Strain Lt2</i>	CP014051
<i>Salmonella Enterica Strain Mac15</i>	NZ_CP030749
<i>Salmonella Enterica Strain Mfds1004024</i>	NZ_CP025745
<i>Salmonella Enterica Strain Mfds1004839</i>	NZ_CP026569
<i>Salmonella Enterica Strain No75</i>	NZ_CP075372
<i>Salmonella Enterica Strain Osf005645</i>	NZ_CP040380
<i>Salmonella Enterica Strain Qh</i>	NZ_CP043773
<i>Salmonella Enterica Strain S146</i>	NZ_CP077662
<i>Salmonella Enterica Strain S44712</i>	NZ_CP035917
<i>Salmonella Enterica Strain S61394</i>	NZ_CP035915
<i>Salmonella Enterica Strain S639</i>	NZ_CP089207
<i>Salmonella Enterica Strain S90</i>	NZ_CP077670
<i>Salmonella Enterica Strain Sa19992307</i>	NZ_CP030207
<i>Salmonella Enterica Strain Sa20021456</i>	NZ_CP030219
<i>Salmonella Enterica Strain Sa20025921</i>	NZ_CP030214
<i>Salmonella Enterica Strain Sa20030575</i>	NZ_CP030181
<i>Salmonella Enterica Strain Sa20031245</i>	NZ_CP030235
<i>Salmonella Enterica Strain Sa20041605</i>	NZ_CP030225
<i>Salmonella Enterica Strain Sa20043041</i>	NZ_CP030231
<i>Salmonella Enterica Strain Sa20044414</i>	CP030209
<i>Salmonella Enterica Strain Sa20051401</i>	NZ_CP030196
<i>Salmonella Enterica Strain Sa20051528</i>	NZ_CP030211
<i>Salmonella Enterica Strain Sa20052327</i>	NZ_CP030202
<i>Salmonella Enterica Strain Sa20055162</i>	NZ_CP030238

<i>Salmonella Enterica Strain Sa20075157</i>	NZ_CP030217
<i>Salmonella Enterica Strain Sa20080453</i>	CP030194
<i>Salmonella Enterica Strain Sa20083039</i>	NZ_CP030223
<i>Salmonella Enterica Strain Sa20083530</i>	NZ_CP030203
<i>Salmonella Enterica Strain Sa20094620</i>	NZ_CP030185
<i>Salmonella Enterica Strain Sa20100201</i>	NZ_CP030180
<i>Salmonella Enterica Strain Sa20101045</i>	NZ_CP030233
<i>Salmonella Enterica Strain Sa20104250</i>	NZ_CP030190
<i>Salmonella Enterica Strain Scsw714</i>	NZ_CP051213
<i>Salmonella Enterica Strain Sg1722-1</i>	NZ_CP081187
<i>Salmonella Enterica Strain Slr1 7627</i>	NZ_CP060517
<i>Salmonella Enterica Strain Slr1 7697</i>	NZ_CP060508
<i>Salmonella Enterica Strain Slr1 7966</i>	NZ_CP060512
<i>Salmonella Enterica Strain Slr1 8094</i>	NZ_CP060515
<i>Salmonella Enterica Strain Slr1 8245</i>	NZ_CP080091
<i>Salmonella Enterica Strain Slr1 8250</i>	NZ_CP060522
<i>Salmonella Enterica Strain Sp</i>	NZ_CP077668
<i>Salmonella Enterica Strain Src27</i>	NZ_CP058807
<i>Salmonella Enterica Strain Szl 30</i>	NZ_CP085981
<i>Salmonella Enterica Strain Szl 31</i>	NZ_CP085983
<i>Salmonella Enterica Strain Szl 38</i>	NZ_CP085987
<i>Yersinia Pestis 1412</i>	CP006783
<i>Yersinia Pestis 1413</i>	CP006762
<i>Yersinia Pestis 1522</i>	CP006758
<i>Yersinia Pestis 2944</i>	CP006792
<i>Yersinia Pestis 3067</i>	CP006754
<i>Yersinia Pestis 3770</i>	NZ_CP006751

<i>Yersinia Pestis 8787</i>	NZ_CP006748
<i>Yersinia Pestis A1122</i>	CP002956
<i>Yersinia Pestis Angola</i>	CP000901
<i>Yersinia Pestis Antiqua</i>	CP000308
<i>Yersinia Pestis Co92</i>	AL590842
<i>Yersinia Pestis D182038</i>	CP001589
<i>Yersinia Pestis Ev Niieg</i>	CP065918
<i>Yersinia Pestis Kim10+</i>	AE009952
<i>Yersinia Pestis Pestoides F</i>	NZ_CP009715
<i>Yersinia Pestis Str. Pestoides B</i>	NZ_CP010023
<i>Yersinia Pestis Strain 14d</i>	NZ_CP063303
<i>Yersinia Pestis Strain 195-P</i>	CP019708
<i>Yersinia Pestis Strain C-781</i>	NZ_CP064117
<i>Yersinia Pestis Strain C-783</i>	NZ_CP064118
<i>Yersinia Pestis Strain C-792</i>	NZ_CP064119
<i>Yersinia Pestis Strain C-830</i>	CP064127
<i>Yersinia Pestis Strain Cadman</i>	NZ_CP016273
<i>Yersinia Pestis Strain El Dorado</i>	NZ_CP009785
<i>Yersinia Pestis Strain Fdaargos 601</i>	NZ_CP033699
<i>Yersinia Pestis Strain Fdaargos 602</i>	NZ_CP033696
<i>Yersinia Pestis Strain Fdaargos 603</i>	NZ_CP033690
<i>Yersinia Pestis Strain Harbin35</i>	NZ_CP009704
<i>Yersinia Pestis Strain I-1252</i>	NZ_CP064126
<i>Yersinia Pestis Strain Java9</i>	NZ_CP009996
<i>Yersinia Pestis Strain Km 567</i>	CP064120
<i>Yersinia Pestis Strain M-1482</i>	NZ_CP064121
<i>Yersinia Pestis Strain M-1974</i>	NZ_CP064124

<i>Yersinia Pestis Strain M2029</i>	NZ_CP064123
<i>Yersinia Pestis Strain M2085</i>	NZ_CP064125
<i>Yersinia Pestis Strain M2086</i>	NZ_CP064128
<i>Yersinia Pestis Strain Nicholisk 41</i>	NZ_CP009991
<i>Yersinia Pestis Strain Pbm19</i>	NZ_CP009492
<i>Yersinia Pestis Strain S19960127</i>	NZ_CP045640
<i>Yersinia Pestis Strain Scpm-O-B-5935 (I-1996)</i>	NZ_CP045154
<i>Yersinia Pestis Strain Scpm-O-B-5942 (I-2638)</i>	NZ_CP045258
<i>Yersinia Pestis Strain Scpm-O-B-6291 (C-25)</i>	NZ_CP045163
<i>Yersinia Pestis Strain Scpm-O-B-6899 (231)</i>	NZ_CP045145
<i>Yersinia Pestis Strain Scpm-O-Dna-18 (I-3113)</i>	NZ_CP045149
<i>Yersinia Pestis Strain Shasta</i>	NZ_CP009723
<i>Yersinia Pestis Subsp. Pestis Bv. Medievalis</i>	CP045158
<i>Yersinia Pestis Z176003</i>	CP001593