

# RESSALVA

Atendendo solicitação do(a) autor(a), o texto completo desta dissertação será disponibilizado somente a partir de 21/10/2024.



UNIVERSIDADE ESTADUAL PAULISTA  
"JÚLIO DE MESQUITA FILHO"  
Campus de Botucatu



UNIVERSIDADE ESTADUAL PAULISTA

"Júlio de Mesquita Filho"

INSTITUTO DE BIOCÊNCIAS DE BOTUCATU

OTIMIZAÇÃO DO SEQUENCE SLIDER: UM MÉTODO DE  
ELUCIDAÇÃO DE ESTRUTURAS CRISTALOGRAFICAS  
PROVENIENTES DE FONTES NATURAIS

**Aluno: João Paulo Ballerini Bruno**

**Orientador: Prof. Titular Marcos Roberto de Mattos Fontes**

**Coorientador: Dr. Rafael Junqueira Borges**

Dissertação apresentada ao Instituto de Biociências,  
Campus de Botucatu, UNESP, para obtenção do título  
de Mestre no Programa de Pós-Graduação em  
Biologia Geral e Aplicada, Área de concentração  
Biomoléculas: estrutura e função (BEF)

**BOTUCATU - SP  
2022**

JOÃO PAULO BALLERINI BRUNO

Otimização do SEQUENCE SLIDER: um método de elucidação de estruturas cristalográficas  
provenientes de fontes naturais

Dissertação apresentada ao Instituto de Biociências,  
Campus de Botucatu, UNESP, para obtenção do título  
de Mestre no Programa de Pós-Graduação em  
Biologia Geral e Aplicada, Área de concentração  
Biomoléculas: estrutura e função (BEF)

Orientador: Prof. Titular Marcos Roberto de  
Mattos Fontes

Coorientador: Dr. Rafael Junqueira Borges

FICHA CATALOGRÁFICA ELABORADA PELA SEÇÃO TÉC. AQUIS. TRATAMENTO DA INFORM.  
DIVISÃO TÉCNICA DE BIBLIOTECA E DOCUMENTAÇÃO - CÂMPUS DE BOTUCATU - UNESP  
BIBLIOTECÁRIA RESPONSÁVEL: ROSEMEIRE APARECIDA VICENTE-CRE 8/5651

Bruno, João Paulo Ballerini.

Otimização do Sequence Slider : um método de elucidação de estruturas cristalográficas provenientes de fontes naturais / João Paulo Ballerini Bruno. - Botucatu, 2022

Dissertação (mestrado) - Universidade Estadual Paulista (UNESP), Instituto de Biociências, Botucatu

Orientador: Marcos Roberto de Mattos Fontes

Coorientador: Rafael Junqueira Borges

Capes: 20804008

1. Envenenamento - Tratamento. 2. Aprendizado de máquina. 3. Toxicologia. 4. Sequence Slider. 5. Elucidação de estruturas.

Palavras-chave: Aprendizado de máquinas; Elucidação de estruturas; Sequence Slider; Toxinologia; XGBoost.

**ATA DA DEFESA PÚBLICA DA DISSERTAÇÃO DE MESTRADO DE JOÃO PAULO BALLERINI BRUNO, DISCENTE DO PROGRAMA DE PÓS-GRADUAÇÃO EM BIOLOGIA GERAL E APLICADA, DO INSTITUTO DE BIOCIÊNCIAS - CÂMPUS DE BOTUCATU.**

Aos 21 dias do mês de outubro do ano de 2022, às 14:30 horas, no(a) Anfiteatro Minoru Sakate, realizou-se a defesa de DISSERTAÇÃO DE MESTRADO de JOÃO PAULO BALLERINI BRUNO, intitulada **Otimização do SEQUENCE SLIDER: um método de elucidação de estruturas cristalográficas provenientes de fontes naturais**. A Comissão Examinadora foi constituída pelos seguintes membros: Prof. Tit. MARCOS ROBERTO DE MATTOS FONTES (Orientador(a) - Participação Presencial) do(a) Departamento de Biofísica e Farmacologia / Instituto de Biociências de Botucatu - UNESP, Prof. Assoc. JOEL MESA HORMAZA (Participação Presencial) do(a) Departamento de Biofísica e Farmacologia / Instituto de Biociências de Botucatu - UNESP, Prof. Dr. ANGELO JOSÉ MAGRO (Participação Presencial) do(a) Departamento de Bioprocessos e Biotecnologia / Faculdade de Ciências Agrônômicas de Botucatu - UNESP. Após a exposição pelo mestrando e arguição pelos membros da Comissão Examinadora que participaram do ato, de forma presencial e/ou virtual, o discente recebeu o conceito final APROVADO. Nada mais havendo, foi lavrada a presente ata, que após lida e aprovada, foi assinada pelo(a) Presidente(a) da Comissão Examinadora.



Prof. Dr. MARCOS ROBERTO DE MATTOS FONTES

---

## AGRADECIMENTOS

- Ao meu pai **Paulo Sérgio Giannelli Bruno** e a minha mãe **Trézia Ieda Ballerini Bruno** por todo o suporte e apoio;
- Ao meu orientador, **Prof. Dr. Marcos Roberto de Mattos Fontes**, pela oportunidade e orientação;
- Ao meu coorientador, **Dr. Rafael Junqueira Borges**, pela atenção, paciência e pelos ensinamentos ministrados.
- Ao programa de Pós-graduação em Biologia Geral Aplicada, por possibilitar a realização deste trabalho.
- O presente trabalho foi realizado com apoio da Coordenação de Aperfeiçoamento de Pessoal de Nível Superior – Brasil (CAPES) – Código de Financiamento 001.

---

## RESUMO

A cristalografia desempenha papel essencial na elucidação dos mecanismos de ação de proteínas, por oferecer dados em nível atômico. Para elucidar a estrutura de uma macromolécula é fundamental o conhecimento da exata composição do seu cristal, o que geralmente é o caso de proteínas obtidas de forma recombinante. Porém, em diversas áreas de estudos, como na toxinologia, as amostras são geralmente obtidas através da purificação direta de fontes naturais, como por exemplo veneno de serpentes, onde propriedades físico-químicas semelhantes de isoformas podem dificultar seu isolamento. Na incapacidade de determinar uma única sequência em um cristal e na ausência de dados cristalográficos à resolução atômica, não existem métodos que auxiliem na elucidação destas estruturas *ab initio*. O método SEQUENCE SLIDER foi desenvolvido para avaliar diferentes possibilidades de cadeias laterais em um modelo cristalográfico no âmbito do faseamento no software ARCIMBOLDO e da incerteza da sequência na toxinologia. Nesta última finalidade, SLIDER integra dados de cristalografia, espectrometria de massa e análises filogenéticas. Assim, o objetivo deste trabalho foi otimizar SLIDER através da técnica de aprendizado de máquinas supervisionado *eXtreme Gradient Boosting* (XGBoost) sobre dados de análise de densidade eletrônica e do ambiente físico-químico de cada resíduo para estimar a atribuição do amino ácido correto. Foram utilizadas 41 estruturas cristalográficas de fosfolipases A<sub>2</sub>, 15 de receptores de porina e 149 metaloproteases, obtidas de fonte recombinante cuja sequência é conhecida para treinamento e teste da metodologia. Resultados obtidos apresentam acurácia de 94.3% a 98.4% para 16.919 resíduos. É esperado que a aplicação deste método a dados inéditos provenientes de proteínas purificadas a partir de fontes naturais com sequência desconhecida possa melhor caracterizar seus componentes e, conseqüentemente, auxiliar na compreensão de seus mecanismos de ação e estratégias de inibição. SLIDER ainda poderá auxiliar outros cristalógrafos e biólogos estruturais ao ser disponibilizado à comunidade científica e, utilizado em diferentes sistemas biológicos obtidos de fontes naturais.

**Palavras-chave:** Aprendizado de Máquinas. SEQUENCE SLIDER. Toxinologia. XGBoost. Elucidação de estruturas.

---

## ABSTRACT

Crystallography plays an essential role for the understanding of the action mechanisms of proteins, as it offers atomic resolution data. In order to elucidate the structure of a macromolecule, it is fundamental to know its exact crystal composition, which is usually the case for recombinant proteins. However, in several areas of study, such as toxinology, samples are usually obtained through direct purification from natural source, such as snake venom, where similar physico-chemical properties of the toxins can cause its isolation to be a challenge. Thus, in case of the inability to determine a single sequence in a crystal and in the absence of crystallographic data at atomic resolution, there are no methods for aiding *ab initio* elucidation of structures. The SEQUENCE SLIDER software was developed to evaluate different side chains possibilities for a crystallographic model in the scope of the ARCIMBOLDO phasing method and the sequences uncertainty in toxinology. In this last aim, SLIDER integrates crystallographic, mass spectrometry and phylogenetic data. Therefore, the goal of this work was to optimize SLIDER through application of the supervised machine learning *eXtreme Gradient Boosting* (XGBoost) with data from electron density and to physico-chemical environment analysis of each residue to estimate the correct amino acid assignment. Train and test data are composed of 41 crystallographic structures of phospholipases A<sub>2</sub>, 15 porine receptors and 149 metaloproteases, obtained from recombinant source, whose sequence is known. Obtained results show accuracy ranging from 94.3% to 98.4% for 16.919 residues. It is expected that the application of the method to elucidate novel data from proteins purified from natural source with unknown sequence can better characterize their components and, consequently, aid action mechanisms comprehension and inhibition strategies developments. SLIDER may be able to assist other crystallographers and structural biologists as it will be available to the scientific community and, used for different biological systems whose source are natural.

**Keywords:** Machine Learning. SEQUENCE SLIDER. Toxinology. XGBoost. Structure Elucidation.

---

## LISTA DE ILUSTRAÇÕES

Figura 1.1 – Metodologia de deslizamento ( <i>sliding</i> ) vertical e horizontal do método SEQUENCE SLIDER.	16
Figura 1.2 – Aminoácidos semelhantes.....	22
Figura 1.3 – Exemplos e diferenças ilustrativas entre algoritmos supervisionados e não-supervisionados .....	29
Figura 1.4 – Representação da ramificação de uma árvore de decisão .....	31
Figura 1.5 – Fluxograma de um <i>Boosting Ensemble</i> .....	34
Figura 1.6 – Cálculo da qualidade de uma árvore de decisão .....	37
Figura 1.7 – Recursos e vantagens de XGBoost .....	38
Figura 2.1 – Metodologia proposta .....	40
Figura 3.1 – Curva de aprendizagem do modelo.....	47
Figura 3.2 – Ganho médio das variáveis do modelo .....	49
Figura 3.3 – Curvas de avaliação de desempenho.....	50
Figura 3.4 – Matrizes de confusão do modelo para dois diferentes limiares de atribuição.....	52
Figura 3.5 – Representação por resíduo para limiar igual a 0,3987 .....	54
Figura 3.6 – Representação por resíduo para limiar igual a 0,9240 .....	55
Figura 3.7 – Representação da densidade eletrônica para o resíduo 3ELO-6 .....	57
Figura 3.8 – Representação da densidade eletrônica para o resíduo 3ELO-31 .....	58
Figura 3.9 – Representação da densidade eletrônica para o resíduo 3ELO-1 .....	59

---

## LISTA DE TABELAS

Tabela 1.1	– Os 20 aminoácidos principais .....	24
Tabela 1.2	– Representação do arquivo final_model.log.....	25
Tabela 1.3	– Representação do arquivo clashes.log.....	26
Tabela 1.4	– Representação do arquivo resdepth.log.....	26
Tabela 1.5	– Representação do arquivo interactions.log.....	27
Tabela 1.6	– Representação do arquivo rotamers.log .....	27
Tabela 3.1	– Resultado do teste de validação cruzada.....	47
Tabela 3.2	– 3ELO-A-6.....	57
Tabela 3.3	– 3ELO-A-31 .....	57
Tabela 3.4	– 3ELO-A-1 .....	57

---

# SUMÁRIO

PREÂMBULO .....	12
INTRODUÇÃO.....	14
1.1 Cristalografia de proteínas e o método SEQUENCE SLIDER.....	14
1.1.1 O método ARCIMBOLDO .....	15
1.1.2 O método SEQUENCE SLIDER.....	16
1.2 Desafios na elucidação de proteínas purificadas de venenos.....	17
1.3 Toxinologia .....	18
1.3.1 A importância de estudar compostos obtidos de fontes naturais .....	18
1.4 Aplicação de SEQUENCE SLIDER à toxinologia.....	20
1.4.1 Estruturas utilizadas.....	20
1.4.2 Variáveis obtidas .....	20
1.4.2.1 Coeficiente de correlação no espaço real - RSCC .....	20
1.4.2.2 Número de interações de contato.....	21
1.4.2.3 Energia de ligação .....	22
1.4.2.4 Clashes.....	23
1.4.2.5 Rotâmeros.....	23
1.4.2.6 Profundidade do resíduo.....	24
1.5 Dados obtidos .....	24
1.6 SLIDER e aprendizado de máquinas.....	28
1.7 Introdução ao aprendizado de máquinas.....	28
1.7.1 Aprendizado supervisionado e não-supervisionado.....	29
1.7.2 Árvores de decisão .....	30
1.7.3 Ensemble Learning.....	31
1.7.3.1 Sabedoria das multidões .....	32
1.7.3.2 <i>Boosting</i> .....	33
1.7.4 <i>eXtreme Gradient Boosting</i> – XGBoost .....	35
1.7.4.1 Aspectos teóricos.....	35
1.7.4.2 Recursos de XGBoost.....	38

MATERIAIS E MÉTODOS.....	40
2.1 Ferramentas Utilizadas .....	41
2.1.1 <i>Numpy</i> .....	41
2.1.2 <i>Pandas</i> .....	41
2.1.3 <i>Matplotlib</i> e <i>Seaborn</i> .....	42
2.1.4 <i>Scikit-Learn</i> .....	42
2.2 Estruturação dos dados .....	42
2.3 Aplicação de XGBoost .....	43
2.4 Métricas de Avaliação .....	44
RESULTADOS E DISCUSSÃO.....	46
3.1 Validação Cruzada.....	46
3.2 Curvas ROC e recall-precisão .....	50
3.3 Matrizes amostrais e de confusão .....	52
3.4 Saída do algoritmo.....	56
CONCLUSÃO.....	61
REFERÊNCIAS .....	62
APÊNDICE A – ESTRUTURAS UTILIZADAS .....	67

---

## PREÂMBULO

A cristalografia desempenha papel essencial na elucidação dos mecanismos de ação de proteínas ao revelar os detalhes de estruturas ao nível atômico, fator chave no entendimento de função e estratégias de inibição. Para se resolver a estrutura cristalográfica de uma macromolécula, é necessário conhecer sua composição, o que é usual ao considerar que a maior parte das proteínas são produzidas de maneira recombinante. Na toxinologia, as amostras geralmente são obtidas a partir de origem natural, isto é, do veneno bruto. Nos passos finais da purificação de venenos, amostras podem ser compostas por misturas de isoformas e/ou diferentes proteínas, já que as propriedades físico-químicas semelhantes dessas moléculas dificultam a separação. É possível cristalizar essas amostras, porém determinar a sequência na estrutura cristalográfica é desafiador na ausência de dados a resolução atômica. O desenvolvimento de métodos que auxiliem na elucidação de estruturas de maneira *ab initio* podem auxiliar cientistas da área de toxinologia a compreenderem melhor os mecanismos de ação das proteínas. Idealizado pelo coorientador desta dissertação, Dr. Rafael Junqueira Borges, SEQUENCE SLIDER aplicado a compostos naturais integra dados de cristalografia, espectrometria de massa e análises filogenéticas para avaliar diferentes hipóteses de cadeias laterais para cada posição de resíduo e assignar sequência (BORGES *et al.* 2022).

Esta dissertação almeja melhorar e complementar o método SEQUENCE SLIDER utilizando um algoritmo de aprendizado de máquinas denominado XGBoost, que calculará a probabilidade de atribuição correta dentre 20 hipóteses de aminoácidos naturais para cada resíduo de uma proteína de interesse a partir de análise de densidade eletrônica e do entorno físico-químico. Esta implementação facilitará a distinção de diferentes hipóteses de cadeias laterais e conclusão das sequências mais prováveis.

A hipótese deste trabalho é que o algoritmo XGBoost pode conseguir diferenciar possibilidades de aminoácidos naturais corretas de erradas para resíduos de uma proteína de interesse utilizando informações acerca do ambiente físico-químico local e da densidade eletrônica de dados cristalográficos. A construção do modelo utilizando o algoritmo XGBoost foi otimizada utilizando métodos de validação cruzada para evitar gastos desnecessários de poder computacional sem prejudicar o desempenho de predição. Para analisar o desempenho do algoritmo frente a diferentes variáveis da densidade eletrônica e do ambiente físico-químico

local, foram feitas análises utilizando o índice de Youden e que maximizam a métrica f-score para propor diferentes limiares de atribuição correta.

A Introdução desta dissertação apresenta as principais motivações teórico-práticas para a criação do método SEQUENCE SLIDER, utilizando como exemplo de possível aplicação, a área da toxicologia. Discorre sobre o método de cristalografia e como esse está relacionado diretamente ao método SLIDER. Exemplifica desafios recorrentes encontradas na área da toxicologia, além de destacar a importância da pesquisa científica nesta área. Explica a metodologia utilizada no cálculo das variáveis de origem físico-química por SLIDER para utilização no desenvolvimento do modelo de aprendizado de máquinas. Posteriormente, a Introdução apresenta a teoria acerca de técnicas de aprendizado de máquinas. Explica e exemplifica métodos supervisionados e não-supervisionados brevemente e o conceito de árvores de decisão, aplicando esses na apresentação de métodos de *ensemble learning*, com enfoque em técnicas de *boosting*. Apresenta XGBoost, o algoritmo utilizado para complementar a saída de SLIDER ao discorrer acerca de seus aspectos particulares, formulações matemáticas e funcionamento, e a metodologia utilizada em sua aplicação, fazendo uso de técnicas de validação cruzada e utilização de métricas na construção de modelos mais robustos iterativamente. Os resultados são analisados utilizando matrizes de confusão para limiares distintos de atribuição e curvas ROC (*receiver operating characteristic*) e recall-precisão, com finalidade de avaliar o desempenho do modelo final. Exemplifica, por fim, a saída do algoritmo, discorrendo acerca de possíveis contribuições na elucidação de estruturas.

---

# INTRODUÇÃO

## 1.1 Cristalografia de proteínas e o método SEQUENCE SLIDER

A cristalografia desempenha papel essencial na elucidação de mecanismos de ação de proteínas, por oferecer dados em nível atômico destas moléculas biológicas na forma nativa ou com mutações e também complexadas com cofatores, inibidores, ou com outras proteínas, com RNA ou DNA. A relevância desta técnica é ilustrada pelos avanços em várias áreas da bioquímica que repercutiram em dezenas de cientistas laureados com prêmios Nobel em virtude do desenvolvimento de pesquisas relacionadas ou baseadas na cristalografia de pequenos compostos inorgânicos ou de origem biológica (JASKOLSKI; DAUTER; WLODAWER, 2014).

O processo de elucidar uma estrutura cristalográfica de uma molécula se inicia em sua cristalização, cujo sucesso está diretamente relacionado à pureza e monodispersidade da amostra. O resultado final é um modelo atômico construído sobre um mapa de densidade eletrônica, sendo este último calculado com base nas intensidades e fases das reflexões provenientes da difração de raios X de cristais homogêneos. Contudo, neste experimento, apenas as intensidades são observadas, de maneira que recuperar as fases é o segundo maior desafio depois da cristalização. Experimentalmente, é possível obter as fases de uma estrutura por meio da introdução de átomos pesados nos cristais (Substituição Isomórfica Múltipla) ou alterando-se o comprimento de onda dos raios X para explorar o espalhamento anômalo de átomos específicos. Ao se obter as fases, que representam apenas parcela da estrutura total, é possível revelar o restante da estrutura utilizando algoritmos de modificação e interpretação de densidade eletrônica.

O conhecimento das estruturas de diversas proteínas permitiu o advento da Substituição Molecular, método mais utilizado para recuperar as fases. Este método consiste em utilizar as fases de um modelo de proteína homóloga que deveria ser semelhante ao da estrutura proteica desconhecida, já que a conservação de estrutura terciária é maior que a primária em proteínas. Nessa técnica, são realizadas duas buscas, uma rotação e outra translação. A amostragem do espaço tridimensional em conjunto com uma função pontuação baseada em máxima parcimônia encontra a solução correta. Com a otimização desta função de pontuação (considerando

resolução e erros dos dados e a fração de espalhamento e divergência esperada do modelo), com melhoria da precisão das medidas de difração de raios X e do poder computacional, surgiu a Substituição Molecular a partir de fragmentos. Nesta, fragmentos são localizados e, por se tratar de porção pequena da estrutura final, a função de pontuação não é suficiente para discriminar verdadeiro de falso positivos. A estratégia é submeter múltiplas possibilidades dos fragmentos aos algoritmos de modificação e interpretação de densidade eletrônica e, caso a subestrutura (fragmentos) esteja acuradamente posicionada, suas fases seriam suficientes para revelar o restante da proteína.

### 1.1.1 O método ARCIMBOLDO

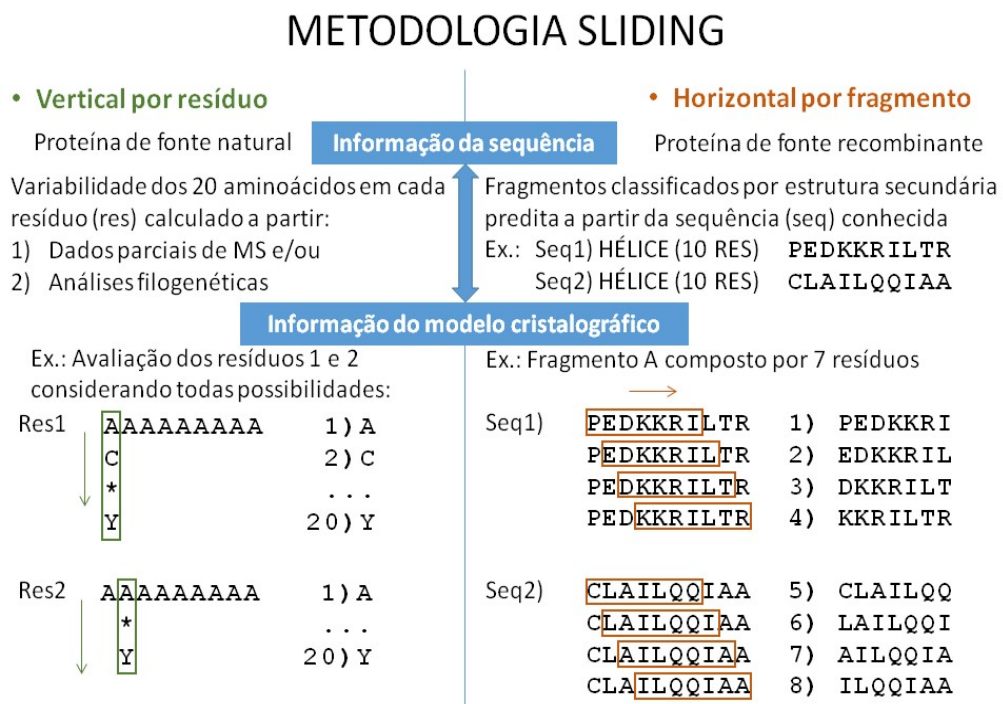
O método ARCIMBOLDO (SAMMITO et al., 2015) foi originalmente concebido como método de faseamento *ab initio*, utilizando a premissa que toda proteína tem em sua composição estrutural fragmentos de estruturas secundárias que se repetem, geralmente hélices  $\alpha$  e folhas  $\beta$  constituídas por 8 a 20 resíduos de aminoácidos (MEDINA et al., 2020). A quantidade e as características dos elementos de estruturas secundárias podem ser preditas utilizando ferramentas de bioinformática a partir da sequência. ARCIMBOLDO reúne um conjunto de hipóteses através da localização do(s) fragmento(s) escolhido(s) com o programa de substituição molecular PHASER (MCCOY et al., 2007). Caso o posicionamento dos fragmentos tenha sido correto (desvio quadrático médio (RMSD) abaixo de 0,5), a modificação e interpretação de densidade eletrônica é realizado pelo programa SHELXE (USÓN; SHELDRIK, 2018), construindo-se, então, o restante do modelo.

O caráter *ab initio* do ARCIMBOLDO destaca sua importância como alternativa às demais técnicas de faseamento que requerem ou dados experimentais adicionais (Espalhamento anômalo ou Substituição Isomórfica Múltipla) ou disponibilidade de um modelo de uma proteína homóloga (Substituição Molecular). Porém, um pré-requisito para ARCIMBOLDO funcionar é que sejam localizados corretamente fragmentos totalizando pelo menos 10% dos átomos proteicos. Entretanto, estender o faseamento a dados de difração com resolução inferior a 2,5 Å e/ou mais de 300 resíduos na unidade assimétrica é desafiador, particularmente na interpretação da densidade.

## 1.1.2 O método SEQUENCE SLIDER

Incorporar cadeias laterais aos fragmentos inicialmente encontrados permite estender ARCIMBOLDO a resoluções mais baixas, já que os átomos presentes nos modelos estariam aumentando. A inspeção manual da densidade eletrônica não permite distinguir a que região pertence o fragmento da sequência, restando comprovar massivamente as combinações de cadeias laterais que a sequência conhecida permite. Para reduzir o número de combinações, foi proposto o método nomeado SEQUENCE SLIDER. A estratégia é restringir as hipóteses de sequências no fragmento coincidindo sua estrutura secundária com a predita a partir da sequência com o algoritmo PSIPRED (BUCHAN et al., 2013; MCGUFFIN; BRYSON; JONES, 2000). Trechos da sequência são deslizados sobre os fragmentos encontrados por PHASER ou traçados por SHELXE contendo mesma estrutura secundária (quadro direito da Figura 1.1). Com essa estratégia foram capazes de elucidar estruturas desconhecidas com resolução a 2.7 Å contendo 600 resíduos na unidade assimétrica (BORGES et al., 2020).

**Figura 1.1** - Metodologia de deslizamento (*sliding*) vertical e horizontal do método SEQUENCE SLIDER



**Fonte:** Autor - Informações complementares provenientes de espectrometria de massa ou análise filogenética são relacionadas ao modelo cristalográfico para gerar hipóteses das diferentes sequências possíveis. Neste, a janela do *sliding* é vertical e cada resíduo é avaliado por vez. No caso ilustrado, 20 possibilidades são avaliadas para os resíduos 1 (Res1) e 2 (Res) de uma cadeia de polialanina. No caso de faseamento, a janela do *sliding* é horizontal e tem comprimento igual ao tamanho do fragmento avaliado (FragA), hipóteses são avaliadas por fragmento e são

compostas pelas sequências que foram preditas com mesma estrutura secundária do fragmento. No caso ilustrado, as hipóteses permitidas pelas sequências preditas como hélice (Seq 1 e 2) e pelo tamanho do fragmento presente no modelo cristalográfico (FragA) são 8.

## 1.2 Desafios na elucidação de proteínas purificadas de venenos

Na toxilogia as amostras de toxinas de serpentes são geralmente obtidas através da purificação direta do veneno, uma vez que a extração do mesmo é prática comum pela acessibilidade da glândula e pela grande quantidade armazenada. Caracterizadas por uma das mais rápidas divergências e variabilidades evolutivas em qualquer categoria de proteínas, as toxinas e suas diferentes isoformas compartilham propriedades físico-químicas entre si (CALVETE et al., 2009). Entre as PLA<sub>2</sub>s por exemplo, foram identificadas e caracterizadas pelo menos 16 isoformas de  $\beta$ -bungarotoxina no veneno de *Bungarus multicinctus* e 15 isoformas de crotoxina no veneno de *Crotalus durissus terrificus* (DOLEY; KINI, 2009). A amoditoxina de *Vipera ammodytes* tiveram duas isoformas elucidadas por cristalografia e apenas suas duas mutações naturais (F113I e K117E) são suficientes para alterar sua toxicidade e atividade anticoagulante (SAUL et al., 2010).

Em virtude dessa complexidade, as amostras obtidas nos passos finais da purificação podem ser misturas de isoformas e/ou diferentes proteínas. Nesses casos, resultados funcionais e estruturais, incluindo espectrometria de massa, são provenientes de combinação de tais compostos. Na ausência de uma sequência conhecida e de dados cristalográficos a resolução atômica (valores nominais abaixo de 1 Å, que se trata de menos de 1% dos dados disponíveis no PDB ([https://www.rcsb.org/stats/distribution\\_resolution](https://www.rcsb.org/stats/distribution_resolution))), poucos são os métodos cristalográficos que auxiliam a elucidação de estruturas com característica *ab initio*, e mesmo estes, apresentam baixa eficácia a resoluções comumente encontradas no PDB.

Modelos cristalográficos com alta confiabilidade são essenciais para servir de suporte à elucidação de mecanismos de ação e desenvolvimento de inibidores específicos que auxiliem na suplementação do soro antiofídico. Geralmente, dados provenientes de espectrometria de massa não alcançam coberturas de 100%, sendo o restante da sequência completado com o observado no banco de dados. Ao desprezar a complexidade destas isoformas, um viés é produzido à informação já disponível no banco de dados evitando introdução de novas sequências. Algumas PLA<sub>2</sub>s, cujos dados cristalográficos são provenientes de amostras, não possuem sequência 100% conhecida, o que torna o método SLIDER, que propõe a avaliação de

diferentes possibilidades de cadeias laterais em fragmentos, imediatamente aplicável à incerteza da sequência nos casos de proteínas purificadas de fontes naturais.

### **1.3 Toxinologia**

A toxinologia, o estudo dos venenos, é uma área de notoriedade recente. Descrito como coquetel de moléculas potentes e estáveis de fácil extração a partir de glândulas externas, o veneno apresenta um grande potencial farmacêutico. O exemplo clássico é o famoso Captopril®, medicamento para tratar a hipertensão, que foi sintetizado em laboratório a partir das estruturas de peptídeos purificados do veneno de serpentes *Bothrops jararaca* (E SILVA; BERALDO; ROSENFELD, 1949). Além da utilidade na medicina, a toxinologia auxilia descrição de processos fisiológicos, como a elucidação do mecanismo catalítico de fosfolipases A<sub>2</sub> (PLA<sub>2</sub>s) a partir de toxinas de serpente e abelha (SCOTT et al., 1990a, 1990b), grupo de proteínas que também está envolvido na inflamação, apoptose e digestão. A capacidade de analisar amostras complexas quali e quantitativamente abrangeu a toxinologia com a proposição da venômica, conjunto de ferramentas ômicas especializadas em caracterizar venenos (CALVETE; JUÁREZ; SANZ, 2007). Contudo, a área ainda carece do sequenciamento do DNA de animais venenosos (KERKKAMP et al., 2016) e de antídotos capazes de neutralizar os efeitos decorrentes de envenenamento de animais peçonhentos. É desejável, portanto, o desenvolvimento de ferramentas que auxiliem pesquisadores a elucidar novas estruturas de toxinas, permitindo um melhor entendimento de processos fisiológicos nos quais essas participam.

#### **1.3.1 A importância de estudar compostos obtidos de fontes naturais**

Os acidentes ofídicos, foram classificados pela Organização Mundial de Saúde como uma doença tropical negligenciada (WORLD HEALTH ORGANIZATION, 2007). É estimado que ocorram aproximadamente 421.000 a 1.800.000 acidentes envolvendo envenenamento por serpentes em todo mundo, em que a taxa de letalidade é cerca de 5% e chances de sequelas físicas permanentes de até 15%. As regiões tropicais e subtropicais, por apresentarem clima ideal para serpentes, registram números de acidentes ofídicos mais elevados, com destaque para o sul e sudeste asiático, África subsaariana e centro e sul do continente americano (KASTURIRATNE et al., 2008).

No Brasil, a taxa de mortalidade é de apenas 0,36% (BRASIL, 2019b) em aproximadamente 28.000 mil casos por ano considerando os dados da última década (BRASIL, 2019a), o que pode ser atribuído à eficiência na produção e distribuição do único tratamento disponível, o soro antiofídico (KASTURIRATNE et al., 2008). Para que os soros antiofídicos tenham um amplo espectro de ação, sua produção é realizada a partir dos venenos das espécies de serpentes mais comuns de uma determinada região (KALIL; FAN, 2016). Fatores como distribuição geográfica, gênero, idade, temporada do ano e variação genética interferem nas características do veneno de serpentes de mesma espécie e, portanto, na complexidade para produção de soros de menor especificidade (KALIL; FAN, 2016).

O gênero *Bothrops* é responsável por cerca de 70% dos acidentes ofídicos no Brasil. Estes eventos apresentam complicações como distúrbios na coagulação, necrose do tecido muscular, hipotensão e disfunções renais, com risco proporcional ao atraso na aplicação do soro específico (BRASIL, 2010). Os danos locais musculares ainda não são neutralizados com eficácia pelo soro antiofídico e podem levar à amputação do membro. O quadro é agravado também pelo fato da grande maioria dos acidentes ofídicos ocorrerem na zona rural durante a execução de trabalhos manuais, fator que dificulta e retarda o acesso aos hospitais (BRASIL, 2010; GUTIÉRREZ; THEAKSTON; WARRELL, 2006).

Dois grupos de toxinas são responsáveis pela mionecrose nos acidentes botrópicos, as metaloproteases e as fosfolipases A<sub>2</sub> (PLA<sub>2</sub>s – E.C. 3.1.1.4). A quantidade do último grupo pode variar em até 80% do total proteico do veneno dependendo da espécie (CALVETE et al., 2010; CALVETE; JUÁREZ; SANZ, 2007). As PLA<sub>2</sub>s apresentam um amplo espectro de ação. *In vivo*, observa-se miotoxicidade local, edema, liberação de citoquinina, recrutamento de leucócitos, hiperalgesia, analgesia e inibição de crescimento tumoral (LOMONTE et al., 2009). Já nos experimentos *in vitro*, os efeitos tóxicos são vastos, tais como citotoxicidade em diferentes tipos de células humanas, ação bactericida, fungicida e antiparasitária, entre outras (LOMONTE et al., 2009). Os estudos dos efeitos farmacológicos *in vivo* são relevantes no envenenamento ofídico para desenvolvimento de estratégias de complementação do soro antiofídico. Outros estudos na área têm como objetivo auxiliar na compreensão da ação das toxinas e no desenvolvimento de medicamentos, já que os efeitos ocorrem em ambiente controlado diferente do natural (LOMONTE et al., 2009). Essa variada gama de efeitos farmacológicos atraíram a atenção dos pesquisadores devido ao potencial das PLA<sub>2</sub>s contra doenças que acometem o ser humano (RODRIGUES et al., 2009).

## 1.4 Aplicação de SEQUENCE SLIDER à toxilogia

O método SEQUENCE SLIDER foi desenvolvido pelo Dr. Rafael Junqueira Borges, coorientador desta dissertação, com o objetivo de avaliar diferentes hipóteses de sequência a dados cristalográficos no âmbito do faseamento foi adaptado para toxilogia. Novas funções foram criadas para avaliar o ambiente físico-químico local e densidade eletrônica. Essas informações podem ajudar na formulação de hipóteses mais robustas de atribuição de sequência, contribuindo na elucidação de uma estrutura de interesse.

### 1.4.1 Estruturas utilizadas

Foram utilizadas 205 estruturas de proteínas de fonte recombinante, retiradas do banco de dados: *Research Collaboratory for Structural Bioinformatics Protein Data Bank* - RCSB PDB (BERMAN, 2000), divididas em 41 PLA<sub>2</sub>s, majoritariamente compostas por  $\alpha$ -hélices, 15 receptores porina compostos por folhas- $\beta$  e 149 metaloproteases, com diferentes estruturas secundárias, totalizando 67.683 resíduos. Os dados apresentam resoluções de 0,85 a 3,0 Å. As estruturas utilizadas estão elencadas no apêndice A.

### 1.4.2 Variáveis obtidas

Para o método SLIDER aplicado à toxilogia, o algoritmo avaliou as diferentes possibilidades de aminoácidos nas densidades eletrônicas e o entorno físico-químico. A presença de dados de espectrometria de massa, peptídeos sequenciados podem servir para restringir as possibilidades avaliadas por SLIDER (BORGES et al., 2022). Além de dados físico-químicos, análises filogenéticas em conjunto com cálculo da taxa de mutação pontual dos resíduos são outras estratégias possíveis para restringir as múltiplas hipóteses (BORGES et al., 2022).

#### 1.4.2.1 Coeficiente de correlação no espaço real - RSCC

SLIDER avalia as diferentes hipóteses geradas através do cálculo do coeficiente de correlação no espaço real (RSCC – real space correlation coefficient) da densidade observada  $\rho(r)_{obs}$  (dados experimentais) com a esperada  $\rho(r)_{calc}$  (teórica a partir do modelo) conforme a Equação 1.1. Existem diferentes estratégias para calcular esse RSCC; a atualmente implementada no SLIDER é o POLDER MAP (LIEBSCHNER et al., 2017), que intensifica o

sinal omitindo determinado(s) átomo(s) escolhido(s) pelo usuário no cálculo das fases e evitando que essa região seja alvo da modelagem de solvente.

$$RSCC = \frac{\sum_r (\rho(r)_{obs} - \overline{\rho(r)_{obs}}) \cdot (\rho(r)_{calc} - \overline{\rho(r)_{calc}})}{(\sum_r (\rho(r)_{obs} - \overline{\rho(r)_{obs}})^2 \cdot \sum_r (\rho(r)_{calc} - \overline{\rho(r)_{calc}})^2)^{1/2}} \quad (1.1)$$

O RSCC é uma variável que fornece informação da densidade eletrônica. Foi calculado para a cadeia principal do aminoácido selecionando os mesmos átomos para diferentes aminoácidos, e para a lateral, que varia de acordo com a hipótese analisada. Se o valor do RSCC é considerado baixo, ou seja, a densidade eletrônica observada e esperada são distintas, as informações obtidas para um determinado resíduo podem ser ruidosas. Uma das razões que ocasionam baixa densidade eletrônica é a flexibilidade de resíduos em regiões da proteína expostas ao solvente.

#### 1.4.2.2 Número de interações de contato

Para cada posição de um resíduo, foi utilizado o *software* PYMOL para gerar as demais moléculas por simetria (SCHRÖDINGER, LLC, 2015), excluindo átomos que estivessem a uma distância maior que 5.5 Å, para reduzir custo computacional e solvente a menos de 2.5 Å para remover impedimento estérico. Foram utilizadas, para cada resíduo, variáveis que representam o número de interações de contato entre átomos. Essas interações podem ocorrer a curtas distâncias; como ligações de hidrogênio (2.7 – 3.3 Å) estabelecendo ligação com átomos hidrofílicos, ou a longas distâncias; como interações hidrofóbicas (3.3 – 4.0 Å). Os números das seguintes interações de contato foram utilizados como variável:

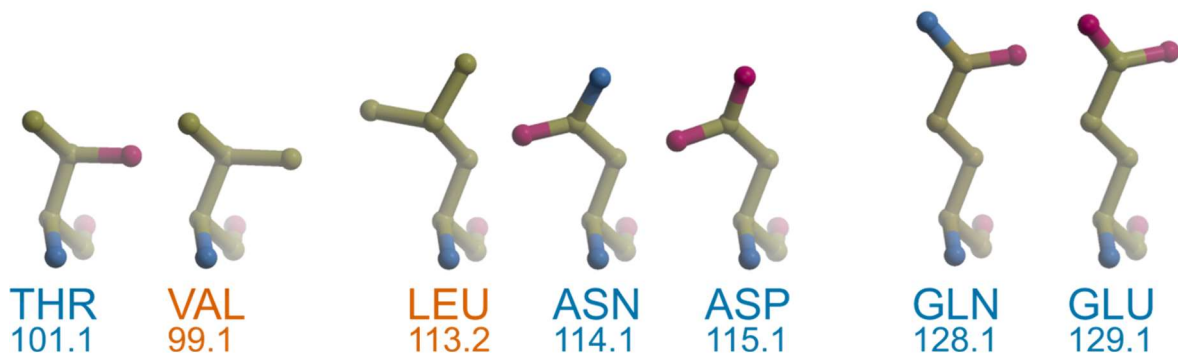
- Ligações de dissulfeto – nSSb – ligação covalente com importante função na formação da estrutura terciária proteica que ocorre entre aminoácidos cisteína.
- Ligações de hidrogênio – nHb – interação polar que ocorre entre átomos de hidrogênio eletropositivos e átomos eletronegativos como O<sup>-</sup> e N<sup>+</sup> com função de estabilização das estruturas secundárias das proteínas. Esta medida foi realizada com o *software* HBPLUS (MCDONALD; THORNTON, 1994), e apenas as ligações com átomos da cadeia lateral foram consideradas.
- Pontes salinas – nSalt – interação iônica e forte. Ocorre principalmente entre aminoácidos com cadeia lateral ácida (átomo eletronegativo) e básica (grupo amino de cadeia lateral). Foram contadas exclusivamente para interações entre

os últimos átomos de interações entre lisina/arginina e ácido aspártico/ácido glutâmico.

- Interações hidrofóbicas – nHydInt – interações que surgem entre moléculas de água e moléculas hidrofóbicas apolares. Quando moléculas hidrofóbicas estão em meio aquoso, as ligações de hidrogênio das moléculas de água de seu entorno serão quebradas para conceder espaço para as moléculas hidrofóbicas, gerando calor ao sistema. As moléculas de água então distorcidas, vão se rearranjar, formando novas ligações de hidrogênio, em estruturas envelopadas com menor entropia. Foram contados os átomos hidrofóbicos da proteína que interagiriam com os hidrofóbicos da cadeia lateral (HAGEMANS et al., 2015).

Comparar número de interações hidrofílicas e hidrofóbicas pode possibilitar a distinção entre aminoácidos aproximadamente isostéricos com densidade eletrônica semelhante, mais especificamente treonina da valina; e leucina do(a) ácido aspártico/asparagina. A utilização da presença de pontes salinas pode diferenciar a glutamina do ácido glutâmico. A figura 1.2 ilustra aminoácidos com densidade eletrônica semelhante.

**Figura 1.2** – Aminoácidos semelhantes.



**Fonte:** Autor – A figura ilustra três casos de aminoácidos que possuem densidade eletrônica semelhante, possuindo apenas átomos diferentes. Podem ser diferenciados por sua massa atômica (valor indicado) e propriedades físico-químicas distintas da cadeia lateral.

### 1.4.2.3 Energia de ligação

Define-se energia de ligação como a quantidade de energia necessária para remover uma partícula de um determinado sistema. A energia de ligação de átomos da cadeia lateral com

demais átomos da proteína, ligantes e solventes próximos foi calculada segundo a Equação 1.2, onde  $nHb$ ,  $nSalt$ ,  $nHydInt$  e  $nSSb$  são o número de ligações de hidrogênio, número de pontes salinas, número de interações hidrofóbicas e número de ligações de dissulfeto. Os valores de energia por ligação foram extraídos do livro ‘Biomolecular Crystallography’ do autor Bernhard Rupp (RUPP, 2010, p. 50).

$$Energia = (nHb + nSalt) \times 3 + (nHydInt) \times 0.7 + (nSSb) \times 62 \text{ kcal/mol} \quad (1.2)$$

#### 1.4.2.4 Clashes

Para analisar colisões, foi avaliado apenas os átomos com distância inferior a 5.5 Å da cadeia lateral escolhida incluindo pares simétricos. Foram utilizadas duas variáveis extraídas do software PHENIX (ADAMS et al., 2010) e uma mensurando distâncias. Quanto maior o valor, maior são as colisões.

- Score de colisões – PhClSc – número obtido através de PHENIX utilizando a função `phenix.clashscore`.
- Superfície de Van der Waals – PhClSum – obtido da somatória da intersecção da superfície de Van der Waals de átomos da cadeia lateral com outros átomos.
- Somatória de scores de distância – SCldist – consiste na somatória da diferença entre 2.5 Å e a distância atômica entre átomos externos, representada pela Equação 1.3.

$$SCldist = \sum_{i=1}^n (2.5 - AtomicDist_i) \quad (1.3)$$

#### 1.4.2.5 Rotâmeros

Foi realizada uma análise utilizando o software PHENIX (`phenix.rotalyze`) que valida os rotâmeros da cadeia lateral de aminoácidos (ADAMS et al., 2010). Os rotâmeros de cada hipótese de aminoácido de um resíduo de interesse foram classificados em favoráveis, permitidos ou não permitidos (*outliers*).

### 1.4.2.6 Profundidade do resíduo

A profundidade do resíduo é medida através de sua distância às regiões externas da proteína, com contato com solvente. Quanto maior o valor, mais profunda na proteína é a localização do resíduo de interesse, o que pode ajudar a diferenciar aminoácidos hidrofílicos de hidrofóbicos. Para este cálculo, foi utilizada a função *ResidueDepth* da biblioteca PDB de BioPython (COCK et al., 2009).

## 1.5 Dados obtidos

Os dados obtidos por esta modalidade de SLIDER são em formato LOG. Para cada proteína descrita na seção 1.3.1, cinco arquivos foram obtidos, contendo dados acerca do ambiente físico-químico do resíduo. Para facilitar, as abreviações de letra única dos 20 diferentes aminoácidos estão dispostas na tabela 1.1.

*Tabela 1.1 – Os 20 aminoácidos principais*

---

A	Ala	Alanina
C	Cis ou Cys	Cisteína
D	Asp	Aspartato (Ácido aspartico)
E	Glu	Glutamato (Ácido glutâmico)
F	Fen ou Phe	Fenilalanina
G	Gli ou Gly	Glicina
H	His	Histidina
I	Ile	Isoleucina
K	Lis ou Lys	Lisina
L	Leu	Leucina
M	Met	Metionina
N	Asn	Asparagina
P	Pro	Prolina
Q	Gln	Glutamina (Glutamida)
R	Arg	Arginina
S	Ser	Serina
T	Tre ou Thr	Treonina
V	Val	Valina
W	Trp	Triptofano (Triptofana)
Y	Tir ou Tyr	Tirosina

---

Apenas 20 aminoácidos encontrados em proteínas de fonte natural estão representados.

As tabelas 1.2, 1.3, 1.4, 1.5 e 1.6 ilustram a saída de SLIDER para cada resíduo para as diferentes variáveis que serão analisadas para cada hipótese de um mesmo resíduo.

Tabela 1.2 – Representação do arquivo *final\_model.log*

Chain	ResN	PCCres	PCC	PCCdif%	
A	1	M!	87.7	6.6	No
A	1	Q	81.1	1.4	No
A	1	A	79.7	1.4	No
A	1	E	78.3	0.3	No
A	1	S	78.0	1.2	No
A	1	K	76.8	0.9	No
A	1	L	75.9	0.1	No
A	1	H	75.8	4.2	No
A	1	D	71.6	1.4	No
A	1	I	70.2	0.4	No
A	1	N	69.8	0.5	No
A	1	F	69.3	0.0	No
A	1	T	69.3	0.2	No
A	1	V	69.1	1.9	No
A	1	R	67.2	0.5	No
A	1	G	66.7	3.7	No
A	1	C	63.0	9.5	No
A	1	Y	53.5	6.4	No
A	1	W	47.1	5.4	No
A	1	P	41.7	last	No
A	1		96.2	BaselineMainChain	

Saída de SLIDER para variáveis relacionadas ao coeficiente de correlação de espaço real (RSCC), onde, Chain, ResN, PCCres, PCC e PCCdif% dizem respeito à cadeia, número do resíduo, hipótese de aminoácido, RSCC e diferença de contraste de RSCC, respectivamente. A última coluna da direita indica se a cadeia lateral é ou não excluída. Na última linha é mostrado o RSCC da cadeia principal, igual para todas as hipóteses de aminoácidos analisadas. Os dados estão arranjados de maneira decrescente utilizando o RSCC, onde a hipótese da primeira linha (Metionina - M) possui maior similaridade de densidade eletrônica com o resíduo analisado. A hipótese correta é mostrada pelo símbolo “!”, dado o conhecimento da sequência correta dos dados utilizados (fonte recombinante).

Tabela 1.3 – Representação do arquivo *clashes.log*

Ch	ResN	ResT	SClDist	PhClSum	PhClSc
P	16	A	0.0	0.0	10.5
P	16	C	0.0	0.0	6.9
P	16	D	0.0	0.0	0.0
P	16	E	0.2	0.5	12.9
P	16	F	1.3	0.0	6.2
P	16	G	0.0	0.0	10.9
P	16	H	0.7	2.2	31.9
P	16	I	0.0	0.9	10.6
P	16	K	0.0	1.6	25.8
P	16	L	0.0	0.5	5.3
P	16	M	0.1	0.8	20.0
P	16	N	0.0	0.0	5.4
P	16	P	0.0	0.0	10.1
P	16	Q	0.0	0.6	4.4
P	16	R	7.2	0.0	5.9
P	16	S	0.0	0.0	0.0
P	16	T	0.0	0.0	0.0
P	16	V	0.0	0.0	0.0
P	16	W	1.7	0.0	6.1
P	16	Y	1.9	1.4	24.8

Saída de SLIDER para as variáveis relacionadas às colisões, Ch ResN, ResT, SClDist, PhClSum e PhClSc são abreviações para cadeia, número de resíduo, abreviação de uma letra do aminoácido, somatória de distância em colisão, somatória da intersecção da superfície de Van der Waals e score de colisões, respectivamente.

Tabela 1.4 – Representação do arquivo *resdepth.log*

Ch	ResN	ResDepth
A	1	3.4
A	2	4.3
A	3	2.0
A	4	2.0
A	5	4.6
A	6	3.3
A	7	2.0
A	8	3.5
A	9	5.2

Saída de SLIDER para a variável de profundidade de cada resíduo. A abreviação Ch, ResN e ResDepth representam a cadeia, o número do resíduo e a profundidade do resíduo em Å, respectivamente.

Tabela 1.5 – Representação do arquivo *interactions.log*

Ch	ResN	ResT	nSSb	nHb	nSalt	nHydInt	Energy
A	1	A	0	0	0	2	1.4
A	1	C	0	1	0	0	3.0
A	1	D	0	1	0	1	3.7
A	1	E	0	1	0	1	3.7
A	1	F	0	0	0	0	0.0
A	1	G	0	0	0	0	0.0
A	1	H	0	2	0	2	7.4
A	1	I	0	0	0	4	2.8
A	1	K	0	0	0	2	1.4
A	1	L	0	0	0	3	2.1
A	1	M	0	0	0	2	1.4
A	1	N	0	2	0	2	7.4
A	1	P	0	0	0	3	2.1
A	1	Q	0	1	0	2	4.4
A	1	R	0	0	0	2	1.4
A	1	S	0	1	0	0	3.0
A	1	T	0	1	0	0	3.0
A	1	V	0	0	0	3	2.1
A	1	W	0	0	0	0	0.0
A	1	Y	0	0	0	5	3.5

Saída de SLIDER para a variáveis relacionadas às interações, Ch ResN, ResT, nSSb, nHb, nSalt, nHydInt são abreviações para cadeia, número de resíduo, abreviação de uma letra do aminoácido, número de pontes de dissulfeto, número de ligações de hidrogênio, número de pontes salina, número de interações hidrofóbicas, respectivamente.

Tabela 1.6 – Representação do arquivo *rotamers.log*

Ch	ResN	C	D	E
A	1	Favored	Favored	OUTLIER
A	2	OUTLIER	Allowed	OUTLIER
A	3	Favored	OUTLIER	Allowed
A	4	Favored	Allowed	OUTLIER
A	5	Favored	Favored	OUTLIER
A	6	Favored	Favored	Allowed
A	7	Favored	Favored	Allowed
A	8	OUTLIER	Favored	OUTLIER
A	9	Favored	OUTLIER	OUTLIER

Saída de SLIDER para a análise de rotâmeros. Ch e ResN são abreviações para cadeia e número de resíduo. Apenas três hipóteses representadas pela abreviação de letra única dos vinte aminoácidos foram exemplificadas.

## 1.6 SLIDER e aprendizado de máquinas

Considerando que um grande volume de dados de estruturas de proteínas depositadas no banco de dados (PDB) possuem sequência conhecida, dado que foram obtidas de fonte recombinante, é direta a possibilidade de predizer tipo do aminoácido para cada resíduo utilizando técnicas de aprendizado de máquinas. Esses algoritmos são capazes de encontrar padrões e correlações em um grande volume de dados para criar um modelo capaz de predizer informações novas. Portanto, obter resultados satisfatórios utilizando um modelo de aprendizado de máquinas em dados de proteínas recombinantes disponíveis no PDB, a partir de dados obtidos por SLIDER, pode auxiliar na elucidação de estruturas provenientes de fontes naturais.

Os algoritmos de aprendizado de máquinas são eficientes na identificação de imagens e voz, funções de material genético não-codificante, de estruturas cristalinas a partir de imagens eletrônicas e padrões de difração sem orientação definida (AGUIAR et al., 2019).

## 1.7 Introdução ao aprendizado de máquinas

Técnicas de aprendizado de máquinas, que fazem parte do extenso ramo da inteligência artificial, estão cada vez mais sendo utilizadas para solucionar problemas encontrados no meio acadêmico e empresarial que, embora disponham de uma grande quantidade de dados, são dificilmente resolvidos por seres humanos. Alguns exemplos mais comuns de utilização dessas técnicas são no reconhecimento de padrões, na interpretação de textos e na detecção de objetos. (MOHAMMED, et al., 2016)

Modelos preditivos buscam encontrar, por meio de aprendizado utilizando um grande número observacional de dados, uma função  $f$  que melhor mapeia variáveis de entrada  $x$  em variáveis de saída  $y$ , de acordo com a Equação 2.1.

$$y = f(x) + e \quad (2.1)$$

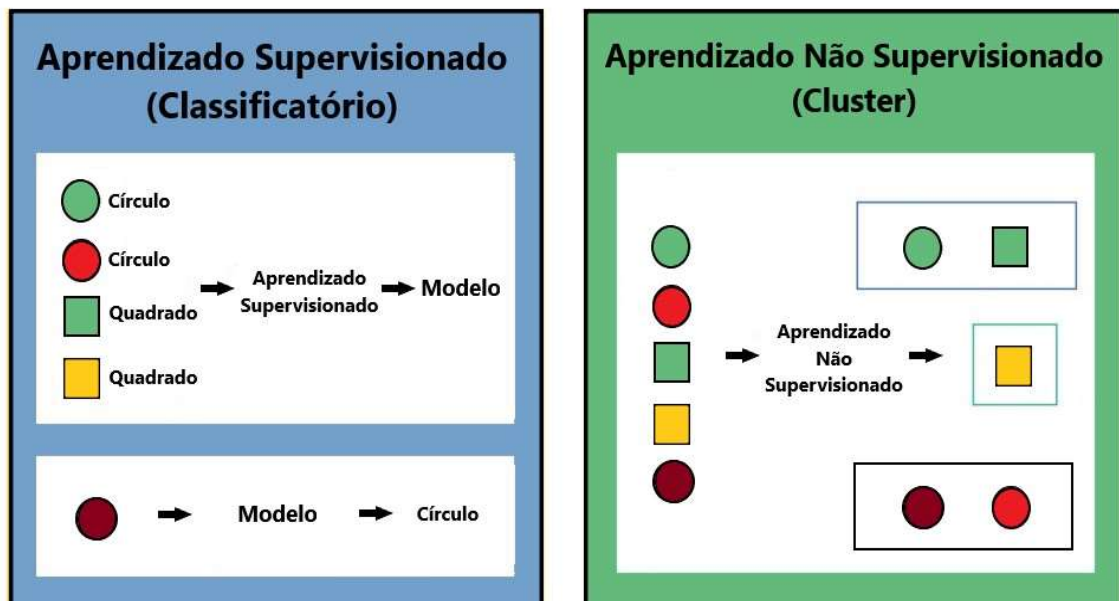
Embora seja desejável que a acurácia do modelo seja perfeita, é provável que exista um erro  $e$  que depende da escolha das variáveis de entrada, pois essas podem não ser suficientes para caracterizar um mapeamento perfeito da função  $f$ . Logo, grande parte das aplicações de técnicas de aprendizado de máquina consistem na busca de melhores escolhas e refinamento de variáveis de entrada  $x$ , com intuito de facilitar a busca pela melhor função de mapeamento

possível, que minimiza o erro  $e$  de um determinado modelo preditivo utilizado na resolução do problema de interesse.

### 1.7.1 Aprendizado supervisionado e não-supervisionado

Existem dois tipos principais de algoritmos de aprendizado de máquina: supervisionados e não-supervisionados. A Figura 1.3 diferencia ambos métodos com um exemplo com formas geométricas.

**Figura 1.3** – Exemplos e diferenças ilustrativas entre algoritmos supervisionados e não-supervisionados.



**Fonte:** Autor - O modelo supervisionado no exemplo tem como variável de saída a forma geométrica; no não supervisionado, a forma geométrica pode ter sido considerada, mas as cores também foram utilizadas no processo de agrupamento dos dados.

Em algoritmos de aprendizado de máquina supervisionados, as variáveis de saída ou, informalmente, respostas procuradas, devem ser bem definidas. Isso possibilita que dados prévios sejam utilizados juntamente com suas respostas referentes já conhecidas para construir um modelo. Para aplicar um método supervisionado, deve-se realizar a separação dos dados do problema entre dois principais grupos: treino e teste. No primeiro, as respostas são evidenciadas juntamente com os dados, e no segundo, são ocultadas, reservando ao algoritmo, com base em

análises realizadas nos dados de treino, prever a resposta correspondente a cada linha de informação. É importante pontuar que observações devem ser independentes, ou seja, não devem ser repetidas em ambos grupos na mesma análise. Um exemplo de algoritmo supervisionado pode ser uma regressão logística, que realiza previsões sobre uma variável discreta binária.

A principal característica de algoritmos não-supervisionados, ao contrário dos supervisionados, é que a saída não é uma variável definida. O modelo é criado sem a necessidade de dados de treino previamente categorizados. É recomendado quando há necessidade de se evidenciar padrões em um grupo de dados de interesse. Para exemplificar; é possível aplicar um *k-means*, método não-supervisionado, que agrupa os dados conforme suas características em comum.

Ao escolher um algoritmo de aprendizado de máquina, deve-se atentar às suas características, pois essas vão implicar na separação e organização dos dados do problema a ser resolvido.

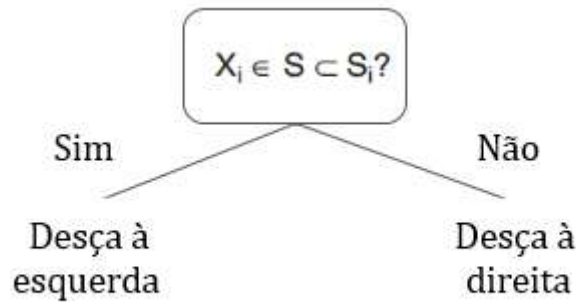
### 1.7.2 Árvores de decisão

Algoritmos de árvores de decisão são utilizados rotineiramente na construção de modelos de aprendizado de máquina, sendo a base de muitas técnicas poderosas de predição como florestas de decisão aleatória (*random forests*), e aceleração (*boosting*).

Atualmente, utiliza-se árvores de regressão (CART), termo introduzido por Leo Breiman em 1984 (BREIMAN, et al., 2017) que, considerando métodos supervisionados, fornece relações entre um grupo de atributos (discretos ou contínuos) e uma classe contínua, ao contrário de árvores de decisão de classificação, onde a classe analisada é discreta.

A representação de um modelo construído utilizando uma CART condiz com uma estrutura de dados denominada árvore binária. Os nós representam a variável de entrada, e também pontos de separação com uma determinada condição. Quando não há ramificações adicionais, os nós terminais são chamados de folhas, e contêm variáveis de saída que são utilizadas para prever algo acerca dos dados analisados. Na figura 2.2, é mostrado um ponto de separação (CUTLER et al., 2012).

**Figura 1.4** – Representação da ramificação de uma árvore de decisão



**Fonte:** Adaptado de (CUTLER et al., 2012) – Considerando uma variável de predição categórica  $X_i$ , que possui valores de um conjunto finito de categorias  $S_i = \{S_{i,1}, \dots, S_{i,m}\}$ , a ramificação ocorre dada a condição  $X_i \in S \subset S_i$ . Para valores discretos ou categóricos, a árvore é denominada de classificação; para valores contínuos, tipicamente números reais, regressão.

As ramificações em árvores de regressão são definidas utilizando critérios de minimização. Considerando valores resultantes de um nó  $y_1, \dots, y_k$ , pode se usar como critério o erro quadrático médio residual

$$Q = \frac{1}{k} \sum_{i=1}^k (y_i - \bar{y})^2, \quad (2.2)$$

onde  $\bar{y}$ , é o valor predito no nó

$$\bar{y} = \frac{1}{k} \sum_{i=1}^k y_i. \quad (2.3)$$

A escolha é dada minimizando  $Q_{sep} = k_E Q_E + k_D Q_D$ , onde  $Q_E$  e  $Q_D$  são os resultados numéricos em um nó resultante da esquerda e direita, respectivamente, e  $k_E$  e  $k_D$  o tamanho das amostras, utilizando algoritmos com alta taxa de atualização (CUTLER et al., 2012).

### 1.7.3 Ensemble Learning

Algoritmos de *ensemble learning* (aprendizagem em comitês) utilizam da combinação de múltiplos modelos construídos previamente para melhorar o desempenho geral na predição de dados de um determinado problema de interesse. Combinar a informação presente de

modelos individuais tem como premissa o conceito de ‘*wisdom of crowds*’ (sabedoria das multidões). Existem diversas técnicas utilizadas para agregar informação de múltiplos modelos.

### 1.7.3.1 Sabedoria das multidões

Uma tomada de decisão é algo recorrente na vida de um ser humano. Para que esta seja a melhor possível, um indivíduo racional poderia se basear nas opiniões distintas de outras pessoas previamente e optar pela mais observada. Para exemplificar: um indivíduo poderia, antes de realizar uma reserva em um hotel, avaliar os comentários de múltiplos clientes anteriores sobre as possíveis qualidades e defeitos do estabelecimento antes de consumir tal escolha. Essa abordagem, a uma tomada de decisão, com base em múltiplas decisões de nível inferior, é definida como sabedoria das multidões ou ‘*wisdom of crowds*’, que utiliza da opinião resultante de um agregado de grupos de pessoas ao invés da opinião de um único indivíduo. (SUROWIECKY, 2005)

A aplicação deste conceito é direta a algoritmos de *ensemble learning*. Construir um modelo robusto utilizando da informação agregada de múltiplos modelos menos eficientes, usualmente descritos como ‘*weak learners*’, é uma estratégia que apresenta melhores resultados em diversas áreas como segurança e detecção de fraudes (VANERIO, CASAS, 2017). As propriedades que fazem esse sistema ser mais eficiente, segundo o livro ‘*Pattern Classification Using Ensemble Methods*’ são elencadas a seguir (ROKACH, 2010):

- Independência: as predições de modelos menores não podem ser influenciadas por outros modelos de mesma categoria, ou seja, as análises devem ser independentes;
- Descentralização: garantir que os modelos possam se especializar e, portanto, resolver um problema de maneira grosseira através de conclusões locais.
- Diversidade de opinião: a interpretação dos dados por um modelo menor, mesmo que crua, deve ser preservada;
- Mecanismos de junção: transformação da informação, anteriormente de caráter privado de um modelo, em uma decisão que pondera a individualidade de todos os modelos menores.

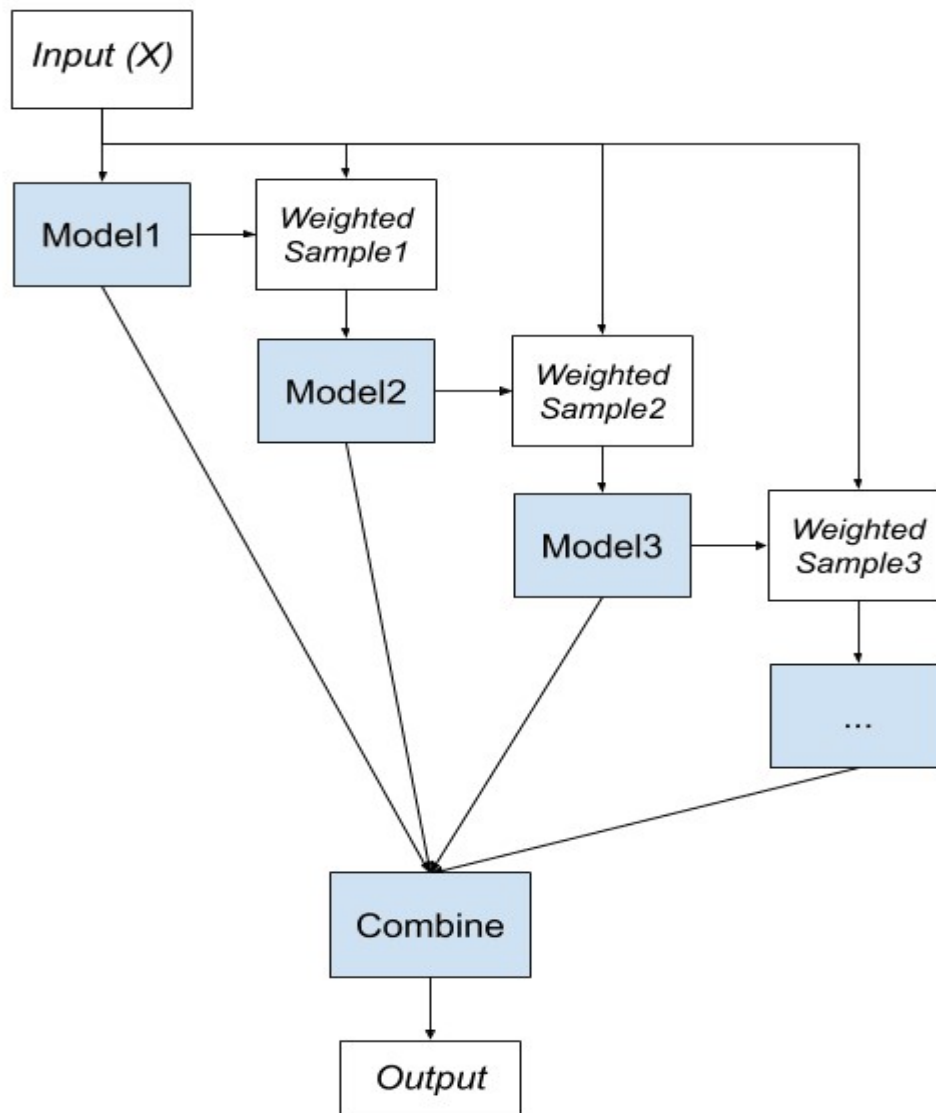
É possível utilizar de variações dessas propriedades para criar abordagens que condizem melhor com o problema a ser solucionado utilizando aprendizado de máquinas.

### 1.7.3.2 *Boosting*

A criação de novas abordagens na solução de problemas envolvendo métodos de *ensemble learning* é limitada apenas na criatividade do desenvolvedor em definir como os modelos menores vão compor o modelo principal (ZHANG, MA, 2021). Existem, portanto, várias estratégias para aplicar um algoritmo de *ensemble learning* que são utilizadas usualmente, como *bagging* (bolsas), *stacking* (empilhar), *boosting*, *bucket of models*, entre outras (ZHOU, 2019).

O *boosting*, é uma técnica de *ensemble learning* que utiliza da manipulação analítica do grupo de dados de treino em um modelo supervisionado para focalizar nos erros cometidos por *weak learners* com o propósito de corrigir essas imprecisões. As iterações são feitas sequencialmente, onde há enfoque na correção dos erros cometidos do primeiro modelo ao construir um segundo através de métodos como votação e utilização de um cálculo de média ponderada, com repetição deste processo até a construção de um modelo com melhor acurácia e desempenho. O enfoque pode ser dado utilizando diferentes pesos para observações (linhas de dados) que são mais escassas na predição de um resultado correto e em exemplos que podem ser mais difíceis de se prever. É mostrado, no fluxograma representado na Figura 1.5, a lógica por trás de um algoritmo de *boosting*.

**Figura 1.5** – Fluxograma de um *Boosting Ensemble*.



**Fonte:** Adaptado de: <https://machinelearningmastery.com/tour-of-ensemble-learning-algorithms/>

O primeiro algoritmo de *boosting* que se destacou em meio a outras técnicas de *ensemble learning* foi o *Adaptive Boosting*, abreviado em *AdaBoost*, onde os modelos são criados utilizando árvores de decisão (FREUND et al., 1997), demonstrando que a estratégia teórica do ‘*wisdom of crowds*’ é utilizada de maneira eficiente na prática para resolver muitos problemas atuais. Ao utilizar algoritmos de *boosting*, deve-se haver precauções quanto à facilidade de ocorrência de *overfitting* (sobreajuste), que prejudica o desempenho do modelo na predição de novos dados, mesmo que os dados de treino tenham sido muito bem preditos. Métodos mais recentes, porém, possuem ferramentas mais eficientes para evitar esse problema recorrente, as

mais usuais são: utilização de parâmetros como ‘*learning rate*’ (taxa de aprendizagem); controlar diretamente o número de iterações para obter o modelo final; e utilização de métodos de regularização.

Utilizando do conceito inicial de *AdaBoost*, Friedman (2001) generalizou o algoritmo, o denominando *Gradient Boosting Machine*, onde modificou a etapa de junção ao utilizar uma estratégia de estágios ao invés de uma sequencial pura; exemplificando: significa que todos os *weak learners* são preservados, reforçando a ideia de independência dos modelos ao invés de interagirem entre si a cada interação na construção do modelo final. Com a utilização de uma otimização numérica, com objetivo de minimizar a perda de informação desses modelos, por método do gradiente, os algoritmos de *boosting* mais modernos podem ser utilizados em problemas de classificação de classes múltiplas, com possibilidade de suporte a regressão, entre outras vantagens, em relação a seus predecessores (TANHA, 2020).

Os conceitos principais do método de *gradient boosting* são, portanto: a criação de modelos do tipo *weak learners*, construídos a partir de árvores de decisão; uma função de perda genérica e, por fim; a construção de um modelo de junção, que possibilita a adição de árvores com independência, garantindo a individualidade da informação dos modelos, possibilitada pela aplicação do método do gradiente. Para garantir que um determinado modelo minimize o erro a cada interação, este é parametrizado e modificado utilizando parâmetros usuais de árvores de decisão, sempre em direção ao mínimo local.

#### **1.7.4 *eXtreme Gradient Boosting* – XGBoost**

O método *eXtreme Gradient Boosting* (XGBoost) (CHEN, GUESTRIN, 2016) é um algoritmo *open source* (código aberto, disponível em <https://github.com/dmlc/xgboost/>) de aprendizado de máquina construído com base em árvores de decisão, utilizando técnicas de aumento de gradiente – *gradient boosting*. Desenvolvido para ser altamente flexível e eficiente, possibilita resolução de diversos problemas encontrados atualmente na ciência e no meio privado, apresentando resultados de última geração em competições atuais de resolução de problemas por métodos de aprendizado de máquina (NIELSEN, 2016).

##### **1.7.4.1 Aspectos teóricos**

O algoritmo XGBoost (CHEN, GUESTRIN, 2016), é considerado um *gradient boosting* regularizado por utilizar dos conceitos de otimização por funções de custo, além de diversas

técnicas de regularização. A função objetiva do algoritmo, que se deseja minimizar, em uma interação  $k$ , é dada por

$$\mathcal{L}^{(k)} = \sum_{i=1}^n l\left(y_i, \hat{y}_i^{(k-1)} + f_k(\chi_i)\right) + \Omega(f_k) \quad (2.4)$$

onde  $l$  é uma função de árvores de regressão,  $\hat{y}_i^{(k)}$  é a predição de uma instância  $i$ -nésima e  $\Omega(f_k)$  é um termo de regularização. Ao constatar que a função objetiva de XGBoost é uma função com parâmetros de funções, deve-se destacar a impossibilidade de utilização de métodos tradicionais de otimização no espaço euclidiano de maneira imediata. Uma aproximação por séries de Taylor de segunda ordem para transformar a função objetiva original foi utilizada para possibilitar uso de técnicas de otimização;

$$f(x) \approx f(a) + f'(a)(x - a) + \frac{1}{2}f''(a)(x - a)^2 \quad (2.5)$$

$$\mathcal{L}^{(k)} \cong \sum_{i=1}^n \left[ l(y_i, \hat{y}_i^{(k-1)}) + g_i f_k(\chi_i) + \frac{1}{2} h_i f_k^2(\chi_i) \right] + \Omega(f_k) \quad (2.6)$$

$$\tilde{\mathcal{L}}^{(k)} = \sum_{i=1}^n \left[ g_i f_k(\chi_i) + \frac{1}{2} h_i f_k^2(\chi_i) \right] + \Omega(f_k) \quad (2.7)$$

onde (2.7) é a função objetiva em uma interação  $k$  e  $g_i$  e  $h_i$  são estatísticas do gradiente de primeira e segunda ordem da função de custo. É possível expandir  $\Omega(f_k)$  ao definir uma coleção ordenada de instâncias  $I_j = \{i | q(\chi_i) = j\}$  de uma folha  $j$ ;

$$\Omega(f_k) = \gamma K + \frac{1}{2} \lambda \sum_{j=1}^K w_j^2 \quad (2.8)$$

$$\tilde{\mathcal{L}}^{(k)} = \sum_{j=1}^K \left[ \left( \sum_{i \in I_j} g_i \right) w_j + \frac{1}{2} \left( \sum_{i \in I_j} h_i + \lambda \right) w_j^2 \right] + \gamma K \quad (2.9)$$

Para uma estrutura fixa de uma árvore  $q(\chi)$ , pode-se computar o peso ideal de uma folha  $j$ :  $w_j^*$ ;

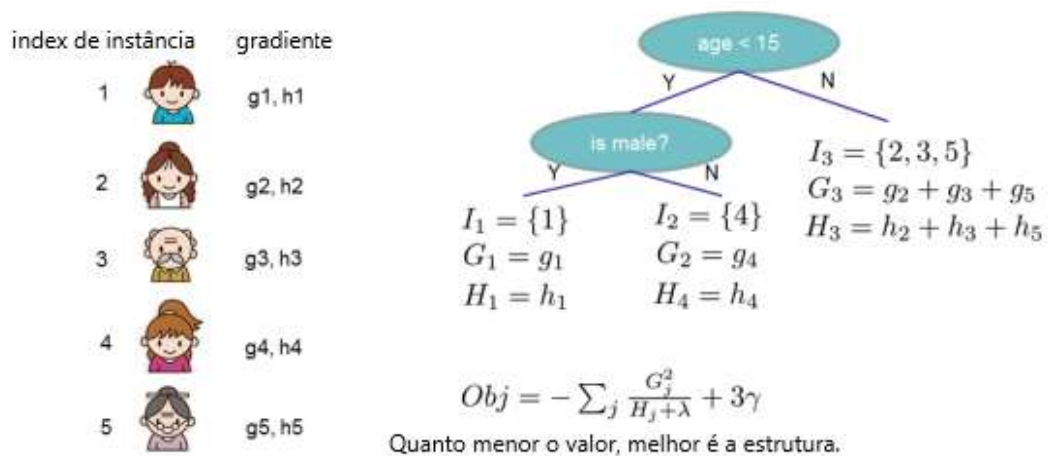
$$w_j^* = -\frac{\sum_{i \in I_j} g_i}{\sum_{i \in I_j} h_i + \lambda}, \quad (2.10)$$

utilizado para calcular o valor ideal correspondente pela função objetiva;

$$\tilde{\mathcal{L}}^{(k)}(q) = -\frac{1}{2} \sum_{j=1}^K \frac{\left(\sum_{i \in I_j} g_i\right)^2}{\sum_{i \in I_j} h_i + \lambda} + \gamma K \quad (2.11)$$

que por sua vez é utilizada como avaliadora da qualidade de uma árvore  $q$ ; avaliação exemplificada na Figura 1.6.

**Figura 1.6** – Cálculo da qualidade de uma árvore de decisão



**Fonte:** Adaptado/Traduzido de (CHEN, GUESTRIN, 2016) - Utilizar o valor da somatória das estatísticas dos gradientes de primeira e segunda ordem em cada folha na fórmula de avaliação representada resulta em um avaliador de qualidade da estrutura.

É utilizado um algoritmo ambicioso exato, que itera sobre todas as possíveis condições de ramificação, escolhendo, de maneira eficiente, melhores escolhas em múltiplos estágios menores para se obter um melhor resultado global. O algoritmo utiliza, dado que  $I_E$  e  $I_D$  são

instâncias dos nós esquerdo e direito após a ramificação, respectivamente, e considerando  $I = I_E \cup I_D$ , a fórmula;

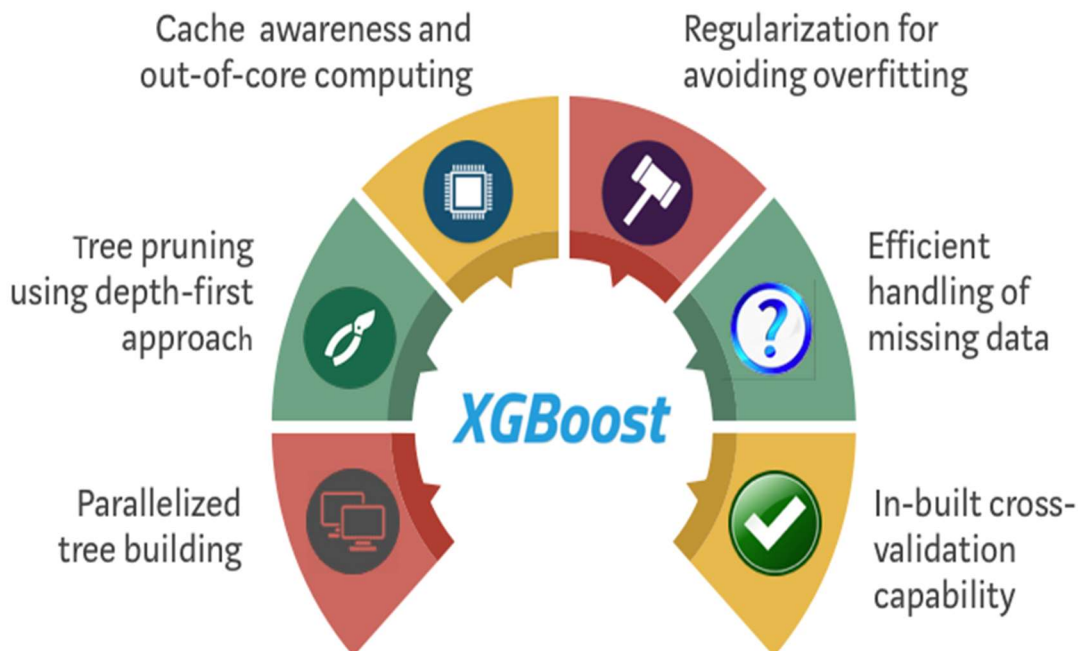
$$\mathcal{L}_{ram} = \frac{1}{2} \left[ \frac{(\sum_{i \in I_E} g_i)^2}{\sum_{i \in I_E} h_i + \lambda} + \frac{(\sum_{i \in I_D} g_i)^2}{\sum_{i \in I_D} h_i + \lambda} - \frac{(\sum_{i \in I} g_i)^2}{\sum_{i \in I} h_i + \lambda} \right] - \gamma \quad (2.12)$$

para encontrar a melhor ramificação de maneira ambiciosa.

#### 1.7.4.2 Recursos de XGBoost

Nos últimos anos, diversas bibliotecas foram construídas para várias linguagens de programação como Python, R, Julia, entre outras. XGBoost apresenta inúmeros recursos desenvolvidos para auxiliar nas dificuldades rotineiramente encontradas na aplicação de métodos de aprendizado de máquina. Alguns deles são exemplificados na figura 1.7 e sintetizados a seguir:

**Figura 1.7** – Recursos e vantagens de XGBoost.



**Fonte:** <http://dataanalyticsedge.com/2019/11/23/xgboost-using-python/> - acesso mais recente em 12/11/2021

1. Construção de árvores de maneira paralelizada ou aprendizagem paralelizada: Uso de todos os núcleos da CPU durante a fase de treinamento dos dados em virtude da

utilização de uma estrutura em forma de blocos organizados com base nas diferentes variáveis do problema em questão.

2. Poda da árvore por profundidade: Há opção de construir as árvores de decisão com profundidade máxima definida com base na escolha do hiperparâmetro '*max\_depth*', o que concede velocidade na fase de treinamento. Após a construção das árvores de decisão, o algoritmo realiza exclusão de nós com parâmetros de regularização não-críticos ou redundantes para a classificação correta. A abordagem de poda de XGBoost consiste em começar pelo nó de maior profundidade, dependendo de hiperparâmetros definidos pelo desenvolvedor, que podem deixar o algoritmo mais ou menos conservador ao realizar podas.
3. Cache-consciente e cálculo *out-of-core*: Implementação de um algoritmo para lidar com problemas relacionados à alocação de memória de maneira indireta resultante da estrutura de blocos característica do algoritmo. Essa estrutura permite um sistema de alocação de uso do disco ao invés de memória temporariamente, minimizando utilização do segundo.
4. Controle de *overfitting*: Observado na Equação 2.4, o termo adicional de regularização auxilia a penalizar a complexidade de modelos (*weak learners*), ponderando melhor o peso de dados na construção de um modelo definitivo. Além disso, vários hiperparâmetros de controle são disponibilizados para uso que contribuem para controlar o *overfitting*, preocupação recorrente em técnicas de *boosting*.
5. Alocação consciente de dados faltantes: O algoritmo é capaz de interpretar facilmente agrupamentos de dados nulos nos dados de maneira otimizada.
6. Validação cruzada embutida: Funções disponíveis que concedem suporte à validação cruzada como `xgboost.cv()`.

---

## CONCLUSÃO

Proteínas purificadas a partir de fontes naturais, como hemoglobinas, lisozimas ou toxinas, podem apresentar incerteza de sequência e sua purificação pode ser desafiador pela coexistência de isoformas. O método SEQUENCE SLIDER aplicado a toxinologia propõe assignar sequência destas proteínas integrando com dados disponíveis, sendo utilizado para a elucidação de uma metaloproteases (SALVADOR et al., 2020) e uma fosfolipase A<sub>2</sub> homóloga de *Bothrops moojeni* (SALVADOR et al., 2021) e duas isoformas de uma PLA<sub>2</sub> de *Bothrops jararacussu* (BORGES et al., 2021). Nestes artigos, dados de cristalografia, espectrometria de massa e análise filogenética foram integrados (BORGES et al., 2022). A análise das densidades eletrônicas e do entorno físico-químico foi realizada manualmente, desta maneira a implementação de aprendizado de máquinas para atribuir probabilidades automatiza SLIDER e auxiliará a elucidação de estruturas futuras.

A implementação do algoritmo de aprendizado de máquinas XGBoost aos dados cristalográficos do SLIDER cria uma nova métrica que integra caracterização de densidade eletrônica e do entorno físico-químico: a probabilidade de assignação de hipóteses de aminoácidos para cada resíduo. Esta implementação foi possível utilizando diversos dados cristalográficos de proteínas com sequência conhecida, já que sua produção foi de forma recombinante. Com os resultados do modelo, é proposto dois limiares de assignação, o primeiro com alta especificidade onde um aminoácido se discrimina dos demais; e o segundo com alta sensibilidade considera mais de um aminoácido como possível. Neste sentido, a metodologia agora otimizada, com uma etapa adicional automática envolvendo aprendizado de máquinas, se torna aplicável a elucidar novas estruturas cristalográficas de proteínas purificadas a partir de fontes naturais e podem ser utilizados por outros biólogos estruturais ao ser disponibilizado à comunidade. Os resultados gerados podem melhor caracterizar os componentes de venenos de serpente, e auxiliar no desenvolvimento de estratégias mais eficientes no tratamento contra envenenamento ofídico.

---

## REFERÊNCIAS

ADAMS, P. D. ET AL. PHENIX: A COMPREHENSIVE PYTHON-BASED SYSTEM FOR MACROMOLECULAR STRUCTURE SOLUTION. **Acta Crystallographica Section D Biological Crystallography**, v. 66, n. 2, p. 213–221, 22 jan. 2010.

AGUIAR, J. A. ET AL. DECODING CRYSTALLOGRAPHY FROM HIGH-RESOLUTION ELECTRON IMAGING AND DIFFRACTION DATASETS WITH DEEP LEARNING. **Science Advances**, v. 5, n. 10, p. eaaw1949, out. 2019.

BERGSTRA, J., YAMINS, D., COX, D. D. MAKING A SCIENCE OF MODEL SEARCH: HYPERPARAMETER OPTIMIZATION IN HUNDREDS OF DIMENSIONS FOR VISION ARCHITECTURES. **30th International Conference on Machine Learning**, 2013.

BERMAN, H. M. THE PROTEIN DATA BANK. **Nucleic Acids Research**, v. 28, n. 1, p. 235–242, 1 jan. 2000.

BORGES, R. J. ET AL. SEQUENCE SLIDER : EXPANDING POLYALANINE FRAGMENTS FOR PHASING WITH MULTIPLE SIDE-CHAIN HYPOTHESES. **Acta Crystallographica Section D Structural Biology**, v. 76, n. 3, 1 mar. 2020.

BORGES, R. J. ET AL. BTHTX-II FROM BOTHROPS JARARACUSSU VENOM HAS VARIANTS WITH DIFFERENT OLIGOMERIC ASSEMBLIES: AN EXAMPLE OF SNAKE VENOM PHOSPHOLIPASES A2 VERSATILITY. **International Journal of Biological Macromolecules**, v. 191, p. 255–266, nov. 2021.

BORGES, R. J. ET AL. SEQUENCE SLIDER: INTEGRATION OF STRUCTURAL AND GENETIC DATA TO CHARACTERIZE ISOFORMS FROM NATURAL SOURCES. **Nucleic Acids Research**, gkac029, 2022.

BRASIL. MINISTÉRIO DA SAÚDE. SECRETARIA DE VIGILÂNCIA EM SAÚDE. DEPARTAMENTO DE VIGILÂNCIA EPIDEMIOLÓGICA. In: **Doenças Infeciosas e Parasitárias: Guia de Bolso**. 7. ed. Brasília, DF: Ministério da Saúde, 2010. p. 431–437.

BRASIL. MINISTÉRIO DA SAÚDE. SECRETARIA DE VIGILÂNCIA EM SAÚDE. DEPARTAMENTO DE VIGILÂNCIA EPIDEMIOLÓGICA. Casos de acidentes por serpentes. Brasil, Grandes Regiões e Unidades Federadas 2000 a 2018. [s.l: s.n.].

BRASIL. MINISTÉRIO DA SAÚDE. SECRETARIA DE VIGILÂNCIA EM SAÚDE. DEPARTAMENTO DE VIGILÂNCIA EPIDEMIOLÓGICA. Óbitos por serpentes. Brasil, Grandes Regiões e Unidades Federadas 2000 a 2018. [s.l: s.n.].

BREIMAN, LEO ET AL. **Classification and regression trees**. Routledge, 2017.

BUCHAN, D. W. A. ET AL. SCALABLE WEB SERVICES FOR THE PSIPRED PROTEIN ANALYSIS WORKBENCH. **Nucleic Acids Research**, v. 41, n. W1, p. W349–W357, 1 jul. 2013.

CALVETE, J. J.; JUÁREZ, P.; SANZ, L. SNAKE VENOMICS. STRATEGY AND APPLICATIONS. **Journal of mass spectrometry: JMS**, v. 42, n. 11, p. 1405–1414, nov. 2007.

CALVETE, J. J. et al. SNAKE VENOMICS OF THE CENTRAL AMERICAN RATTLESNAKE CROTALUS SIMUS AND THE SOUTH AMERICAN CROTALUS DURISSUS COMPLEX POINTS TO NEUROTOXICITY AS AN ADAPTIVE PAEDOMORPHIC TREND ALONG CROTALUS DISPERSAL IN SOUTH AMERICA. **Journal of proteome research**, v. 9, n. 1, p. 528–544, jan. 2010.

CHAP T. LE. A SOLUTION FOR THE MOST BASIC OPTIMIZATION PROBLEM ASSOCIATED WITH AN ROC CURVE. **Statistical Methods in Medical Research**. 15: 571-584, 2006.

CHEN, Tianqi; GUESTRIN, Carlos. Xgboost: A scalable tree boosting system. In: **Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining**. 2016. p. 785-794.

COCK P.A., ET AL., BIOPYTHON: FREELY AVAILABLE PYTHON TOOLS FOR COMPUTATIONAL MOLECULAR BIOLOGY AND BIOINFORMATICS. **Bioinformatics**, 25, 1422-1423

CUTLER, ADELE; CUTLER, D. RICHARD; STEVENS, JOHN R. RANDOM FORESTS. In: **Ensemble machine learning**. Springer, Boston, MA, 2012. p. 157-175.

DOLEY, R.; KINI, R. M. PROTEIN COMPLEXES IN SNAKE VENOM. **Cellular and molecular life sciences: CMLS**, v. 66, n. 17, p. 2851–2871, set. 2009.

FERNÁNDEZ, A. ET AL. **Learning from Imbalanced Data Sets**. Cham: Springer International Publishing, 2018.

FLUSS, RONEN; FARAGGI, DAVID; REISER, BENJAMIN. ESTIMATION OF THE YOUDEN INDEX AND ITS ASSOCIATED CUTOFF POINT. **Biometrical Journal: Journal of Mathematical Methods in Biosciences**, v. 47, n. 4, p. 458-472, 2005.

FREUND, Yoav; SCHAPIRE, Robert E. A decision-theoretic generalization of on-line learning and an application to boosting. **Journal of computer and system sciences**, v. 55, n. 1, p. 119-139, 1997.

FRIEDMAN, Jerome H. Greedy function approximation: a gradient boosting machine. **Annals of statistics**, p. 1189-1232, 2001.

GUTIÉRREZ, J. M.; THEAKSTON, R. D. G.; WARRELL, D. A. CONFRONTING THE NEGLECTED PROBLEM OF SNAKE BITE ENVENOMING: THE NEED FOR A GLOBAL PARTNERSHIP. **PLoS Medicine**, v. 3, n. 6, jun. 2006.

HAGEMANS, D. ET AL., A SCRIPT TO HIGHLIGHT HYDROPHOBICITY AND CHARGE ON PROTEIN SURFACES. **Frontiers in Molecular Biosciences**, v. 2, 13 out. 2015.

HARRIS, C.R., MILLMAN, K.J., VAN DER WALT, S.J. ET AL. ARRAY PROGRAMMING WITH NUMPY. **Nature** **585**, 357–362 (2020).

HUNTER, J. D. MATPLOTLIB: A 2D GRAPHICS ENVIRONMENT. **Computing in Science & Engineering**, v. 9, n. 3, p. 90–95, 2007.

JASKOLSKI, M.; DAUTER, Z.; WLODAWER, A. A BRIEF HISTORY OF MACROMOLECULAR CRYSTALLOGRAPHY, ILLUSTRATED BY A FAMILY TREE AND ITS NOBEL FRUITS. **The FEBS journal**, v. 281, n. 18, p. 3985–4009, set. 2014.

KALIL, J.; FAN, H. W. PRODUCTION AND UTILIZATION OF SNAKE ANTIVENOMS IN SOUTH AMERICA. In: GOPALAKRISHNAKONE, P. (Ed.). **Toxins and Drug Discovery**. Dordrecht: Springer Netherlands, 2016. p. 1–22.

KASTURIRATNE, A. ET AL. THE GLOBAL BURDEN OF SNAKEBITE: A LITERATURE ANALYSIS AND MODELLING BASED ON REGIONAL ESTIMATES OF ENVENOMING AND DEATHS. **PLoS medicine**, v. 5, n. 11, p. e218, 4 nov. 2008.

LIEBSCHNER, D. ET AL. POLDER MAPS: IMPROVING OMIT MAPS BY EXCLUDING BULK SOLVENT. **Acta Crystallographica Section D Structural Biology**, v. 73, n. 2, p. 148–157, 1 fev. 2017.

LOMONTE, B. ET AL. THE PHOSPHOLIPASE A2 HOMOLOGUES OF SNAKE VENOMS: BIOLOGICAL ACTIVITIES AND THEIR POSSIBLE ADAPTIVE ROLES. **Protein and peptide letters**, v. 16, n. 8, p. 860–876, 2009.

MCCOY, A. J. ET AL. PHASER CRYSTALLOGRAPHIC SOFTWARE. **Journal of Applied Crystallography**, v. 40, n. 4, p. 658–674, 13 jul. 2007.

MCDONALD, I. K.; THORNTON, J. M. SATISFYING HYDROGEN BONDING POTENTIAL IN PROTEINS. **Journal of Molecular Biology**, v. 238, n. 5, p. 777–793, 19 mai. 1994.

MCGUFFIN, L. J.; BRYSON, K.; JONES, D. T. THE PSIPRED PROTEIN STRUCTURE PREDICTION SERVER. **Bioinformatics (Oxford, England)**, v. 16, n. 4, p. 404–405, abr. 2000.

MCKINNEY, DATA STRUCTURES FOR STATISTICAL COMPUTING IN PYTHON. **Proceedings of the 9th Python in Science Conference**, Volume 445, p. 56-61, 2010.

MEDINA, A. ET AL. ALEPH: A NETWORK-ORIENTED APPROACH FOR THE GENERATION OF FRAGMENT-BASED LIBRARIES AND FOR STRUCTURE INTERPRETATION. **Acta Crystallographica Section D Structural Biology**, v. 76, n. 3, 1 mar. 2020.

MOHAMMED, M.; KHAN, M. B.; BASHIER, E. B.. **Machine learning: algorithms and applications**. Crc Press, 2016.

NIELSEN, Didrik. **Tree boosting with xgboost-why does xgboost win" every" machine learning competition?**. 2016. Dissertação de Mestrado. NTNU.

PEDREGOSA, F. ET AL. SCIKIT-LEARN: MACHINE LEARNING IN PYTHON. **The Journal of Machine Learning Research**, v. 12, n. null, p. 2825–2830, 1 nov. 2011.

RODRIGUES, R. S. ET AL. SNAKE VENOM PHOSPHOLIPASES A2: A NEW CLASS OF ANTITUMOR AGENTS. **Protein and Peptide Letters**, v. 16, n. 8, p. 894–898, 2009.

RUPP, B. **Biomolecular crystallography: principles, practice, and application to structural biology**. New York: Garland Science, 2010.

SAITO, TAKAYA; REHMSMEIER, MARC. THE PRECISION-RECALL PLOT IS MORE INFORMATIVE THAN THE ROC PLOT WHEN EVALUATING BINARY CLASSIFIERS ON IMBALANCED DATASETS. **PLoS one**, v. 10, n. 3, p. e0118432, 2015.

SALVADOR, G. H. M. ET AL. BIOCHEMICAL, PHARMACOLOGICAL AND STRUCTURAL CHARACTERIZATION OF BMOOMP-I, A NEW P-I METALLOPROTEINASE FROM BOTHROPS MOOJENI VENOM. **Biochimie**, v. 179, p. 54–64, dez. 2020.

SALVADOR, G. H. M. ET AL. THE SYNTHETIC VARESPLADIB MOLECULE IS A MULTI-FUNCTIONAL INHIBITOR FOR PLA2 AND PLA2-LIKE OPHIDIC TOXINS. **Biochimica et Biophysica Acta (BBA) - General Subjects**, v. 1865, n. 7, p. 129913, jul. 2021.

SAMMITO, M. ET AL. ARCIMBOLDO\_LITE: SINGLE-WORKSTATION IMPLEMENTATION AND USE. **Acta Crystallographica. Section D, Biological Crystallography**, v. 71, n. Pt 9, p. 1921–1930, 1 set. 2015.

SAUL, F. A. ET AL. COMPARATIVE STRUCTURAL STUDIES OF TWO NATURAL ISOFORMS OF AMMODYTOXIN, PHOSPHOLIPASES A2 FROM VIPERA AMMODYTES AMMODYTES WHICH DIFFER IN NEUROTOXICITY AND ANTICOAGULANT ACTIVITY. **Journal of Structural Biology**, v. 169, n. 3, p. 360–369, mar. 2010.

SCHRÖDINGER, LLC. **The PyMOL Molecular Graphics System, Version 1.8**. nov. 2015.

SUROWIECKI, JAMES. **The wisdom of crowds**. Anchor, 2005.

ROKACH, LIOR. **Pattern classification using ensemble methods**. World Scientific, 2010.

TANHA, Jafar et al. Boosting methods for multi-class imbalanced data classification: an experimental review. **Journal of Big Data**, v. 7, n. 1, p. 1-47, 2020.

USÓN, I.; SHELDRIK, G. M. AN INTRODUCTION TO EXPERIMENTAL PHASING OF MACROMOLECULES ILLUSTRATED BY SHELX ; NEW AUTOTRACING FEATURES. **Acta Crystallographica Section D Structural Biology**, v. 74, n. 2, p. 106–116, 1 fev. 2018.

VANERIO, Juan; CASAS, Pedro. Ensemble-learning approaches for network security and anomaly detection. In: **Proceedings of the Workshop on Big Data Analytics and Machine Learning for Data Communication Networks**. 2017. p. 1-6.

VAN ROSSUM, G., & DRAKE, F. L. **Python 3 Reference Manual**. Scotts Valley, CA: CreateSpace, 2009.

WASKOM, M. SEABORN: STATISTICAL DATA VISUALIZATION. **Journal of Open Source Software**, v. 6, n. 60, p. 3021, 6 abr. 2021.

WORLD HEALTH ORGANIZATION. RABIES AND ENVENOMINGS: A NEGLECTED PUBLIC HEALTH ISSUE: REPORT OF A CONSULTATIVE MEETING, World Health Organization, Geneva, 10 January 2007. p. 32, 2007.

ZHANG, CHA; MA, YUNQIAN (ED.). **Ensemble machine learning: methods and applications**. Springer Science & Business Media, 2012.

ZHOU, ZHI-HUA. **Ensemble methods: foundations and algorithms**. Chapman and Hall/CRC, 2019.