



UNIVERSIDADE FEDERAL DO RIO DE JANEIRO

Instituto de Biofísica Carlos Chagas Filho

Concurso para Professor Adjunto

MS-020 – Biologia Computacional

Edital nº 953, de 20 de dezembro de 2019, publicado no DOU nº 248, de 24 de dezembro de 2019 - consolidado com as alterações dos Editais nº 9, de 9 de janeiro de 2020, nº 31, de 03 de fevereiro de 2020, nº 48, de 11 de fevereiro de 2020 e nº 116 de 25 de março de 2020



NOME: Maria Fernanda Ribeiro Dias

No. fl.  
1/14

### Tema 3 - Ferramentas estatísticas em Biologia Computacional

A Biologia Computacional se caracteriza pela elaboração de algoritmos matemáticos/computacionais para solucionar um problema biológico. A implementação desses algoritmos pode ser escrita em diferentes linguagens e para cada algoritmo é necessário métricas e programos estatísticos que visam avaliar o desempenho do mesmo.

Dentro da modelagem comparativa podemos citar o RMSD que avalia a diferença na posição dos átomos entre duas moléculas (modelo e molde), sendo útil para saber se o modelo criado segue parâmetros coerentes e pode ser utilizado em etapas seguintes na geração de uma proteína por exemplo.

Podemos ainda acrescentar à essa análise o gráfico de Ramachandran, que avalia os ângulos ( $\phi$  e  $\psi$ ) de um polipeptídeo, permitindo assim avaliar sua este

quissimetria, por meio de todas as combinações possíveis desses ângulos. As duas métricas citadas não utilizadas em Biologia Estrutural e limitam-se a avaliação de estruturas 3D. Essas análises podem ser feitas usando a ferramenta ~~Pymol~~ PyMOL, por exemplo, um programa para visualização de moléculas 3D, ou ainda programas para análise de ancoramento molecular como o Docking (Autodock VINA).

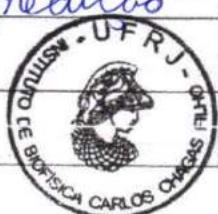
Uma outra análise seria a análise de correlação entre duas variáveis, essa pode ser utilizada quando temos um conjunto de dados sendo representados por diferentes atributos, por exemplo, propriedades físico-químicas e assim, se faz necessário analisar se esses atributos estão diretamente ou inversamente correlacionados. Para esta análise podemos usar a correlação de Pearson, por exemplo, onde os atributos são a par não comparados, e gerados valores entre -1 e 1, valores próximos a -1 indica que os atributos se correlacionam de forma negativa; valores próximos a 1 indica que estes se correlacionam de forma positiva, logo um deles pode ser eliminado da representação do conjunto de dados, sem gerar danos para os mesmos; valores próximos a 0 indica que os atributos não independentes. Essas análises não utilizadas em problemas cujo



conjunto de dados apresenta alta dimensão e é necessário uma redução dimensional em sua representação. Problemas biológicos aplicados à análise por Inteligência Artificial necessitam passar por essa etapa de pré-processamento. Para fazer essas análises é possível utilizar a ferramenta R, um pacote estatístico com interface gráfica, que permite ao usuário utilizar bibliotecas prontas ou implementar algoritmos direcionados para o conjunto de dados em questão.

Ainda dentro da Inteligência Artificial, especificamente na área de Aprendizagem de Máquina, diversas ferramentas estatísticas são usadas. Em Aprendizagem de máquina não supervisionada, quando não se tem informações prévias sobre o conjunto de dados, pode-se usar métrica de distância, como análise de distância Euclidiana para separar os dados e extrair informações desse mesmo. Já a técnica de Aprendizagem supervisionada apresenta uma gama de ferramentas, implementadas em C, ai para o PYTHON, para avaliar o desempenho de um classificador ou modelos de classificação que utilizaram diferentes parâmetros ou algoritmos.

Suponhamos que temos um algoritmo de predição cujo foco é classificar possíveis ligantes com potencial para inibição enzimática.



As análises estatísticas para avaliação de desempenho desse classificador iniciam com um processo de validação cruzada, que utiliza o que chamamos de conjunto teste (ligantes e decoys, conhecidos e retirados do banco de dados). O processo de validação cruzada pode ocorrer de diferentes formas, uma delas é o uso K-fold-cross-validation onde utiliza-se todos os dados divididos em K partes, usando K-1 como conjunto de treinamento e 1 parte separada como teste. Aplicando o classificador nesse conjunto de dados temos a formação da que chamamos de matriz de confusão, constando de: VP → verdadeiros positivos; VN → verdadeiros negativos; FP → falsos positivos; FN → falsos negativos; Esses constam de dados cujo classificador faz a previsão de forma correta e dados que são preditos de forma erronea pelo classificador, respectivamente. Diante dessa classificação, podemos extrair métricas estatísticas como Recall, Especificidade e acurácia.

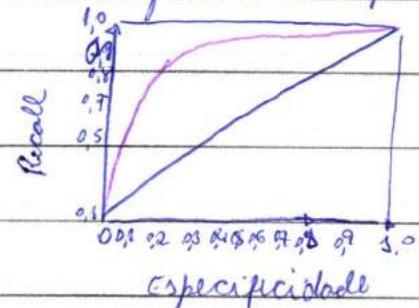
Recall  $\rightarrow \frac{VP}{VP + FN}$  → Avalia a quantidade de dados preditos como positivos diante de todos os positivos.

Especificidade  $\rightarrow \frac{VN}{VN + FP}$  → Avalia a quantidade de dados preditos como negativos diante todos os dados negativos.

Acurácia  $\rightarrow \frac{VP + VN}{VP + FP + VN + FN}$  → Avalia a quantidade de dados preditos de forma correta diante o conjunto de todos os dados.



As métricas de Recall e especificidade permitem a elaboração da análise do desempenho global de desempenho da curva ROC. Esta análise traz um comparativo entre os classificações corretas de dados positivos e negativos, assim digo que a área sob a curva (AUC) demonstra o desempenho de classificação em um determinado problema, dado o exemplo do gráfico a seguir; podemos dizer que o classificador apresenta área AUC = 0,90 na classificação de dados positivos.



Além da curva ROC é necessário ainda citar o fator de Enriquecimento que avalia a quantidade de dados classificados, com um determinado score a partir de um ponto de corte.

Embora essas métricas sejam muito utilizadas, é necessário que estas não sejam avaliadas sozinhas, mas em parceria com análise acurada do conjunto de dados. Em problemas biológicos é comum que tenhamos muitos dados relacionados a uma única classe, o que constitui o que chamamos de dados desbalanceados. Esse tipo de dado gera o que chamamos de overfitting (quando o preditor se adapta ao conjunto de dados específico, sendo muito eficiente no conjunto teste e apresentando valores próximos a 0,5 ao serem aplicados em conjuntos de





componibile (guscio, peristoma, testa), se possiamo altro  
far quei neuromotori extra para expressar sua  
aqui nesse fases nesse o exequente quando  
de um organismo ou complexo biológico. Esse  
que eu entendo que humanos é fundamental  
de uso, formadas por massas químicas (de hidrata  
ou da Regulação da sustância e contra-a de matéria  
de aquela de fluidos metabólicos contínuo num adulto

#### Tema 4 - Modelagem de Redes Metabólicas

as estruturas, utilizando métodos matemáticos.  
e figura shows como implementamos e como elas podem  
ampliar com outras funções e facilitar a geração  
deles utilizados para representar fluxo-químico e  
O foco principal era focar na tipologia de redes.  
as proteínas de HU, assim como proteínas inibidoras.  
evidências de como os de degradação e reguladoras metabólicas.  
como conseguem de degradação de forma muito elegante  
mecanismo de feedbacks inibidores de processos, utilizados  
formam dinâmica entre de degradação, suas metas  
também mudam face de degradação, suas metas  
quando de degradação tem a técnica SMTG, por exemplo  
que de degradação tem certa vez de certa vez temos  
concluídas. As regras que são feitas tem a que

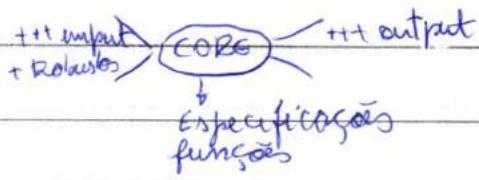
6/19

No. II.

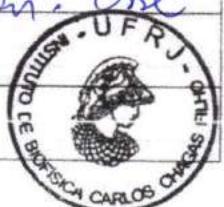
NOME: Afonso

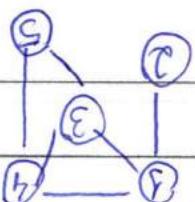
para o Todo?" Essa pergunta permite que as pesquisas biológicas itinem novos rumos que se contrapõem à biologia reducionista e abrange analisar um problema em um organismo como um todo. Assim podemos citar duas metodologias; Top-down e bottom-up. <sup>Metodologias</sup> Análises do tipo top down geram um fluxo descendente de informações, onde dados de genómica ou transcriptómica são utilizados por exemplo para nutrir uma rede. já nas metodologias do tipo bottom-up, há um fluxo ascendente de informações, onde dados de metabólitos são utilizados para nutrir e montar a rede metabólica.

~~Análises~~ na montagem de redes metabólicas ~~base~~ com necessita de <sup>modelos</sup> ~~modelos~~ matemáticos/computacionais que sejam robustos (avalia a estabilidade do método, mesmo tendo condições de perturbação) e degradativos. Uma arquitetura viável é do tipo bow-tie que se assemelha a a gravata borboleta, onde as entradas (várias) passam por um core <sup>especializado</sup> e menos degenerado



O protocolo para modelagem de redes metabólicas foi estabelecido em 2010 por Thiele e Palsson. Esse protocolo conta com 4 etapas e surgiu a partir de um genoma. São elas:





$$\{4,5\}, \{5,3\}, \{3,4\} \}$$
$$G = \{(1,2), (1,3), (1,4)\}$$

grado G, formas: V = {1, 2, 3, 4, 5}

afé e escravizá tu' outas. Dánum para um  
and u' no Voltais que pediu' re' gom, perdiu'as  
Aplicando lema das gráficas tâmas que G = {1, 2, 3}  
ca u' vez mais de gráficas u' mdu'as us' ação.  
uma forma de transformar uma redi multabel.  
gab.

u' , por exemplo a mensage da informaçao em  
atômas e subatômas auxiliu' su' memória, pedindo fa'  
Tâmas an 4 tâmas no' multâmas, pedindo asu' mdu'  
mdu' que se' amde fomada com rru'as q' fomadas  
cuidar de gols' jumentos ma' redi e operações dea  
4) Aprendendo des' redi + Consulta ma' aulação  
permeiso' des' redi, filhos' e escopo e organização da redi.  
as fumâmas e atômas que mdu'as se' aplicam ma'  
3) Aprendendo de modulos matemáticos + Consulta ma' bude  
a mdu'as que se' p'ra' medidau', assim a' confunção des' atômas.  
uma e' matemática contumânta Técnica e biológico des'  
us' de fumâmas e atômas que mdu'as se' aplicam ma'  
2) Cuidar'as mdu'as + Consulta mo' mo' as opereções  
us' de fumâmas des' as mdu'as p'ra' mdu' + p'ra'  
a mdu'as que se' p'ra' medidau', assim a' confunção des' atômas.  
us' de fumâmas e atômas que mdu'as se' aplicam ma'

3) Draft - Constitui' a formação da redi multâmbia  
dâmodas des' operações u' lâmbos das dâmodas globo's, come

8/14

No. II.

Miguelina

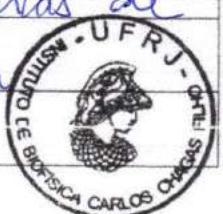
NOME:

Vamos, cada vértice corresponde a um nó e cada linha corresponde a uma reação ou interação entre um nó e outro. Supondo que cada nó seja uma proteína, ou uma reação vírica, podemos ter uma análise de fluxo direcional, avaliando os substratos gerados em determinadas reações.

A utilização de matrizes estequiométricas conta com análises de restrição onde cada reação forma uma coluna com índices que correspondem a valores positivos e negativos, sendo produtores e reagentes (-) e produtores (+), respectivamente.

Além dos gaps, um outro componente que se deve avaliar nas redes metabólicas são fatores externos, por exemplo fatores ambientais. Como o processo de produção é itatitivo, a cada informação ou atualização da rede, um novo nó ou informação pode ser acrescentado e assim a rede metabólica é atualizada.

Uma das ferramentas muito utilizadas na implementação de redes é a COBRA (constraint-based Reconstruction and Analysis), uma ferramenta que foi implementada em MATLAB, mas atualmente já apresenta implementação em Python. COBRA oferece ainda ferramentas de avaliação da rede, por exemplo ao modelar uma rede metabólica de prokariotos, mon-



NOME: Wyllis

No. fl.

10/14

Na avaliação o grau de 1 nó ( $K$ ) indica quantas conexões a ele mesmo possui com outros nós, sendo  $p(K)$  a distribuição probabilística desses nós em um grafo.

Por ser um protocolo interativo, temos os modelos criados <sup>baseados</sup> em restrições, por meio de imposições de restrições físico-químicas e identificação de pressões setitivas, por exemplo fatores ambientais e o fluxo de reações pode ser representado por escalas estequiométricas, utilizando ainda otimização linear, baseada em estado estacionário, sem levar em conta parâmetros cinéticos.

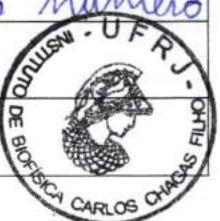
Para análise da rede, podemos levar em conta ainda análise de modularidade, ontologia genética e centralidade que constituem respectivamente <sup>em</sup> módulos (como subconjunto de nodos) padronizados que juntamente com localizações e funcionalidade celular constituem ontologia genética, baseando ainda a análise de centralidade em métricas locais de grafos.

## Tema 9. Fundamentos Teóricos de Algoritmos de Biologia Computacional e Bioinformática

Para embasar esse tema gostaria de citar Margarit Dayoff, que no final do



século XX, mais especificamente na década de 70, iniciou a partir de análise comparativa de sequências de proteínas curadas manualmente, a produção de uma matriz de ~~pontuações~~<sup>substituição</sup> para mutações, a matriz PAM. Margaret era matemática e utilizou para isso o conjunto de 5660 proteínas que compunham o 1º atlas proteico, todas curadas manualmente pela própria Margaret. O que Margaret fez foi criar matrizes em que comparava duas sequências proteicas  $S_1$  e  $S_2$  e estabelecia pontuações para cada mutação encontrada, assim tínhamos os valores de 0 na diagonal das matrizes e de acordo com a taxa de mutação dos aminoácidos, esses iam ganhando valores maiores mais positivos quando apresentavam altas taxas de mutação e negativos quando apresentavam valores baixo. Partindo das sequências  $S_1$  e  $S_2$ , podemos dizer que ambas as sequências são 1 PAM divergentes se  $S_1$  puder ser substituída por  $S_2$ , sem danos, ocorrendo uma mutação a cada 200 resíduos de aminoácidos. Ao fazer essas comparações, Margaret ~~até~~ gerou árvores de substituição em cada sequência e ia anotando as mutações ocorridas. Ela gerou duas métricas i) mutabilidade relativa; ii) frequência relativa; que corresponde ao número de vezes que um aminoácido foi substituído por outro e na posição do mesmo; e a



NOME: Málias

No. fl.

12/19

freqüência de substituições, na qual foi gerado uma tabela com valores para padrão para cada aminoácido. Nessa tabela é possível observar que o triptofano sofre pouquíssima mutação, e essa informação é importante até hoje quando avaliamos sitios de interação entre proteínas de forma computacional pois a substituição desse aminoácido geralmente leva à inativação da proteína, por ser um aminoácido grande e interferir na estrutura proteica. A matriz PAM foi calculada em diferentes níveis tendo hoje PAM 30, PAM 70 e até 250 PAM.

20 anos depois a matriz PAM formação dessas matrizes, surgiu a matriz BLOSUM, gerada na década de 90 pelo casal Henaff. As matrizes BLOSUM também eram calculadas a partir da comparação de sequências biológicas de proteínas utilizando pontuação, porém essa comparação era feita em blocos e utilizava sequências mais distantes e dissimilares. Tanto a matriz PAM, quanto a BLOSUM não utilizadas em algoritmos de alinhamento como o BLAST, essas constituem métricas importantes para análise comparativa de genoma, por exemplo uma comparação do genoma do SARS-CoV-2 com outros coronavírus, inferindo informações acerca da origem.



de vírus, assim como similaridades com outros organismos.

Outra base que fundamenta metodologias computacionais para aplicações em biologia é o uso de ferramentas que simulam questões biológicas como o uso de algoritmos genéticos, utilizados em algoritmos de docking molecular, esses não são considerados algoritmos evolutivos por apresentarem abordagens diversificadas. Os algoritmos evolutivos surgiram por John Holland e são algoritmos estocásticos que no docking, permitiam o enriquecimento da técnica por permitir a inclusão de mais graus de liberdade relacionados à flexibilidade molecular. Esses algoritmos são otimizados por meio do que chamamos de mutações e recombinações e assim, inicialmente temos um cromossoma com possíveis soluções para o problema; esse vai sofrendo cross over e sendo apresentadas novas soluções. Essas soluções e a otimização do algoritmo podem ser avaliadas por meio dos valores de energia do docking ou da análise de curva ROC, já descrito no item 3.

Além das ~~anota~~ algoritmos que dão inicio à análise de sequências lineares e análises estruturais, gostaria ainda de citar os algoritmos baseados em Inteligência Artificial (IA) área da computação, que desde o inicio do



NOME: Thales

No. fl.

19/19

século XXI com a intensificação de dados geração de dados vem sendo aplicada em problemas biológicos. Esta aplicação se intensifica conforme surgem mais procedimentos e melhorias nas técnicas "ômicas". Um dos primeiros trabalhos a utilizar <sup>IA</sup> técnicas de aprendizagem de máquina foi do Golub, em 1999 que utilizou análise de cluster como K-means e técnicas baseadas em distância como a euclidiana para separar dados de células de origem linfóide e mielóide. As análises feitas por Golub utilizaram uma plataforma chamada genepattern. Esta plataforma permite análise de genes utilizando técnicas de aprendizagem supervisionada e não supervisionada.

Nesta área de aprendizagem supervisionada, algoritmos como árvore de decisão foram aplicados para classificar sitios de ubiquitinação, utilizando informações físico-químicas de proteínas. Este projeto ocorreu na primeira década do século XXI e se constituiu em um dos primeiros trabalhos desenvolvidos com aplicação de metodologia de inteligência Artificial em problemas biológicos. Hoje, temos uma gama de algoritmos como SVM, redes neurais e sua evolução para aplicação e resolução de problemas biológicos, sendo usados na avaliação de redes metabólicas e até mesmo na elucidação da estrutura ~~3D~~ terciária de uma proteína.

