



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: GIORDANO BRUNO SANTOS SOUZA

No. fl.
1

Questão 1) → Ferramentas estatísticas em biologia computacional

A estatística é um campo do conhecimento orientado a descrever populações (ou amostras) e tentar inferir a probabilidade de os eventos. Mais é surpreendente que esteja intimamente atrelada à história evolutiva dos organismos e à descrição dessas populações. Grandes avanços no campo da estatística foram realizados por geneticistas como Galton e Fisher que ~~faziam~~ criavam novos modelos e experimentos para descrever as populações e tentar prever alterações nestas.

De maneira geral, a estatística pede ser dividida em uma parte descriptiva que busca caracterizar as amostras. As características de cada elemento dessa população podem ser denominadas variáveis. Variáveis assumem valores qualitativos (como as categóricas - associadas à nominais; e as ordinais - que representam a ordenação da amostra). As variáveis podem assumir ainda valores discretos ou contínuos.

Em populações biológicas, exemplos de variáveis contínuas são ~~podem ser~~ massa e altura dos indivíduos, ~~variáveis discritas são~~ a contagem de alelos em um indivíduo assim como a ~~é~~ contagem de alelos de um locus pode assumir caráter discreto em um ~~individuo~~. Alelos são abstrações para representar a herança genética de indivíduos e, dessa forma, nos apre-

claram a interpretar e inferir as características e a história evolutiva das populações. O ato de tentar prever essas características e as mudanças temporais constitui a segunda parte da estatística, denominada induativa ou preditiva.

Como no caso de Galton que tentava inferir o impacto da genética na altura dos indivíduos, outras abordagens são possíveis. Pode-se querer inferir quais fatores genéticos e socioeconómicos estão associados a uma doença como o câncer gástrico. Neste caso, proponho um teste de hipótese onde a hipótese nula representa que não há associação entre um genótipo ou a renda com a doença e a alternativa que indica a associação. Para testar essas hipóteses recorremos aos testes estatísticos, estes podem ser paramétricos quando a distribuição de variável é conhecida e normal, ou não-paramétrico para os demais casos. Entre os testes estatísticos mais utilizados, podemos usar o teste de Kolmogorov-Smirnov para identificar o status socioeconómico está associado à doença ou a regressão para inferir a relação entre a variável dependente e a resposta. Esses testes citados dependem de que as variáveis testadas sejam independentes, o que ~~num~~ sempre é o caso. Um teste de associação onde as variáveis são tomadas ao longo do tempo, essa associação é violada. Em um estudo para associar a cor da pele (constitutiva e a responsável ao UV) a dados genéticos e socioeconómicos usando uma coorte familiar, a assumção de independência é violada uma vez que os indivíduos são aparentados e compartilham grande parte do genoma. Nestes casos onde a colinearidade está presente, é conveniente usar outras ferramentas como os modelos lineares mistos, que estes são capazes



NOME: GIORGANO BRUNO SOARES SOUZA

No. fl.
3

de criar diferentes interceptos para amostragens temporais - um intercepto para cada unidade de tempo ou famílias - um intercepto por família.

Fica claro neste exemplo que é necessário conhecer o grau de compartilhamento genético para utilizar os métodos estatísticos apropriados. De uma forma conveniente de se analisar a diversidade genética das populações e identificar estruturas cripticas é realizar o cálculo das frequências aleáticas das populações e aplicar o modelo de Equilíbrio de Hardy-Weinberg. Neste modelo, a probabilidade conjunta das frequências aleáticas de duas populações é estimada e comparada a observada, diferenças significativas indicam a estrutura populacional. Esse mesmo princípio pode ser utilizado para inferir a divergência genética de grupos populacionais a partir da análise de recombinação. Esse cálculo denominado Índice de fração (FST) na genética de populações nos permite identificar o grau de divergência entre os loci das populações com grande serventia nas áreas evolutiva e biomedicina.

Em um estudo realizado em populações nativas do Peru, conseguimos infiar polymorfismos possivelmente sob seleção natural devido ao grau de divergência das frequências locais específicas e identificamos um conjunto de marcadores com estruturação e associados à doença em outros estudos na literatura científica. Isso pode indicar que alguns desses achados são falso-positivos e a associação não é causal devido a fatores genéticos mas a algum fator associado. Por exemplo, para um estudo de associação com o câncer gástrico, havia associação entre a ancestralidade ameríndia e fatores sócioeconômicos e desfecho clínico, mas a menor loci foi



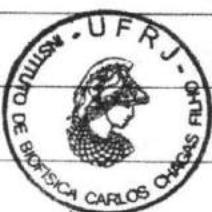
individuais associadas devido à incorporação da ancestralidade. Em alguns casos, indivíduos muito diferentes geneticalemente em um grupo case ou controle podem ser excluídos das análises. Para exclusão em para incorporação ao modelo, as análises ^{de} reduções de dimensionalidade, tais como, análise de componentes principais são usualmente empregadas.

A identificação das bases moleculares das doenças pode levar a melhorias no prognóstico, diagnóstico e tratamento, mas devido ao caráter complexo da etiologia das desfechas é conveniente incorporar a incerteza nos modelos de previsão e decisão. Em geral, isso pode ser alcançado através da estatística Bayesiana. Esta, se assenta sobre o teorema de Bayes e estima a probabilidade de um evento a partir da incorporação do conhecimento prévio ou de incerteza ao modelo.

Nesta figura, temos uma distribuição de probabilidades a priori - que saímos, uma distribuição de probabilidades a posteriori - que desejamos conhecer e uma verossimilhança associada que nos indica quais valores são mais prováveis. Análises Bayesianas são utilizadas em modelos de decisão, como ajustar a dose de um medicamento de acordo com o perfil genético do indivíduo ou até mesmo serem incorporadas um modelos mais complexos como o ABC (Computação Bayesiana Aproximada) que facilita e permite inferências evolutivas. ⁴ nos estudos evolutivos, essa ferramenta é utilíssima pois descrevemos os valores dos parâmetros do modelo e, por isso, é conveniente representar essa incerteza.



O uso das ferramentas estatísticas pela biologia computacional é essencial para descrever as populações, inferir sua história evolutiva e as implicações destas na expressão de fenótipos. Durante minha trajetória acadêmica, este tem sido o meu principal objetivo.através da descrição da diversidade de populações e loci, da tentativa de identificar as bases moleculares de fenótipos como o câncer gástrico e a cor da pele, assim como a associação destes a fatores sócioeconômicos e no desenvolvimento de ferramentas para summarizar dados biológicos tenho continuamente aplicado vários dos conceitos e ferramentas necessárias para aumentar o ^{meu} conhecimento sobre a evolução e o impacto da diversidade.



Q2) modelagem de redes biológicas

Redes são abstrações que nos permitem identificar os componentes, entender o comportamento dinâmico e a interação entre os elementos. Constituem um método ^{conveniente} ~~apropriado~~ para modelar a relações entre entidades por não requerer regularidade e por serem objetos combinatorios flexíveis. ~~Não~~ constituídos ^e ~~constituídos~~ São representados como grafos ($G = \{V, E\}$) constituídos por vértices que representam os elementos e arestas que representam os relacionamentos. Estas podem ser direcionadas (sentido) ou não-direcionadas, e ~~e~~ ^e capazes de expressar confiança ou força através dos pesos, ou atração e repressão através de sinais.

Outras aplicações de redes envolvem a visualização de relações entre variáveis causais de diferentes tipos e que nos permitem identificar elementos comuns às ciências e também os específicos. Nesta aplicação refer-se à dissolução de redes familiares para cálculo de ancestralidade. Neste caso, uma rede é modelada tendo as arestas representando o grau de relacionamento entre os indivíduos e, recursivamente, grafos são formados ~~com~~ ^{com} o intuito de agrupar a maioria quantitativa de indivíduos mais aparentados.

Outro tipo de rede bastante utilizada são as redes de georreferências onde as arestas representam probabilidade e os grafos são direcionados. Equações diferenciais podem ser usadas para modelar sistemas dinâmicos no espaço.



~~individuos e para estimar o grau de dissimilaridade genômica entre
individuos de espécies próximas.~~

~~Di teoria de sistemas está intimamente conectada à biologia
de sistemas por ilustrar modelar e prever os resultados de perturbações.
Um sistema é descrito como um grupo caos, interacionando, interdependentemente,
situado no espaço e no tempo que relaciona-se com o ambiente e possui uma
estrutura e ~~um~~ propósito. Di modelagem dos sistemas nos permitem
compreender o seu funcionamento e identificar sinergia e propriedades
emergentes. Uma das aplicações desta~~

~~e no tempo. Di representação das interações e topologias podem
auxiliar a entender a dinâmica do sistema imodelável. Algumas
medidas importantes para caracterizar reais existem em:~~

~~Conectividade - ~~o~~ número de arestas necessárias para juntar e
separar um gráfo sub-gráfo; central centralidade que atesta a importân-
cia do nó. homogeneidade que representa a quantidade média
de ligações; distância que é calculada pela soma das arestas
e intermediárias que representa o nó pelo qual passam a maior
parte dos caminhos mais curtos.~~



A3) Fundamentos teóricos de algoritmos da filosofia computacional
e filosofia lógica

É filosofia computacional é um campo interdisciplinar que envolve as áreas da ciência da computação, ciências filosóficas e sociais. Os fundamentos teóricos sobre os quais se assenta a filosofia da computacional é a filosofia formalística:

- 1) Teoria da informática; 2) Teoria dos sistemas; 3) Teoria do controle;
- 4) Matemática discreta e estatística; 5) Genética de populações e evolução;
- 6) Filosofia e sistemática.

Entre os fatores que podemos modelar problemas e aplicar ~~técnicas~~ de algoritmos para solucionar as questões postas pelos pesquisadores é a teoria da informática estabelece a base técnica para se entender como as mensagens são emitidas e de codificadas, além de estabelecer a unívocia como uma medida de unicidade associada a um experimento probabilístico. O autor pode ser medida entre duas mensagens ~~separadas~~ indicando independência ou assimilação. Essa também relaciona-se à complexidade de uma sequência, por exemplo, e pode indicar o grau de compressão possível. Um conceito de univocidade tem sido aplicado a diferentes ~~regras~~ algoritmos para estabelecer o grau de similaridade entre sequências, como a identificação de assimilações genéticas ou comparacões de sequências com alinhamentos. Em trabalho recente, ~~utilizando~~ este ferramentas associadas à este conceito de similaridade para confirmar que ~~sequências~~ o mesmo tipo de resultados em tempos e espacos distintos eram provenientes das 

indivíduos e para estimar o grau de dissimilaridade genômica entre indivíduos de espécies próximas.

A teoria de sistemas ~~de~~ busca modelar e prever o impacto das ~~de~~ perturbações em um sistema. Este pode ser descrito como um grupo de elementos coeso, interrelacionado, interdependente, situado no tempo e no espaço, que relaciona-se com o ambiente e possui estrutura e propósito. Essas características que descrevem o sistema são amplamente encontradas na biologia e, por isso, a teoria de sistemas fornece subsídios para a modelagem de sistemas biológicos afim de prever suas funcionalidades. Um exemplo de aplicação ocorre na modelagem natural de seleção péligenica onde ~~uma~~ ^{a presença} ~~outro~~ de vários genes com indicativo de seleção em subgrafos de uma rede metabólica podem indicar a ocorrência de um processo de adaptação. Em trabalho com nativos americanos, conseguimos identificar agrupamentos de genes relacionados à metabolização de fármacos.

~~BB~~

Quando se ~~se~~ aplica seus próprios algoritmos, há diferentes ~~abordagens~~ fundamentos e abordagens que podem ser aplicadas para a resolução de uma tarefa. Em geral, os algoritmos podem ser divididos em: I) programação dinâmica; II) programação linear; III) algoritmos quentes; IV) dividir e conquistar; V) grafos. A essência de cada categoria e exemplos de utilização na minha carreira acadêmica são descritos ~~nos~~ nos parágrafos a seguir.



Os algoritmos baseados em Hidden Markov ~~mo~~ ^{contêm} modelos probabilísticos gerativos onde se representa ~~o~~ uma cadeia de estados com os seguintes parâmetros: ~~o~~ inicial, transição e emissão. O sentido e o estado ^{só} dependem do próprio estado. Uma distribuição de probabilidades é associada a cada comunidade e o algoritmo de Viterbi escolhe a melhor comunidade aplicadas para fixar parâmetros iniciais em programas de infecção populacional.

Programações dinâmica ~~também~~ consiste em dividir os problemas em instâncias menores. São utilizadas em algoritmos de ~~otimizar~~ sequências.

Algoritmos gulosos resolvem problemas peça-a-peça e só se importam com a melhor solução local e imediata.

Algoritmos de Dividir e Conquistar quebram os problemas em instâncias menores e resolvem recursivamente. Basta aplicá-los em ordenações e com aplicações na seleção de primos para PCL multiplex.

Programação ~~dinâmica~~ ^{linear} é um processo de otimização onde as restrições e otimizações são lineares.

Algoritmos de grafos podem auxiliar na identificação de subconjuntos através da decomposição modular e podem ser aplicados em árvores de decisões.

