



CODIGO DA PROVA: RP-001/0001



1
↓ de 13

UNIVERSIDADE FEDERAL DO RIO DE JANEIRO
INSTITUTO DE CIÊNCIAS BIOMÉDICAS
CONCURSO:

FOLHA DE RESPOSTA

Importante: O código da prova só será colocado na entrega da prova ao fiscal. As provas serão escaneadas e enviadas aos membros da banca avaliadora sem o nome do candidato.

Item a - Controle da expressão gênica com ênfase em métodos computacionais.

- O controle da expressão é uma das formas que buscamos para entender uma das perguntas mais fundamentais dentro de áreas de biologia celular e do desenvolvimento ~~para~~ que é como o mesmo DNA de um organismo é capaz de dar origem a múltiplos ~~tipos~~ tipos celulares. Isto é, é a partir do controle da expressão gênica que genes são ativados ou desativados ao longo da diferenciação celular e desenvolvimento.

O nosso entendimento deste processo está historicamente relacionado aos avanços em tecnologias para analisar a própria expressão gênica e fatores fundamentais relacionados a como mensuramos o quanto, onde e como genes são expressos.

Inicialmente, métodos computacionais ~~de~~ que nos ajudam a entender expressão gênica evoluem de acordo com as tecnologias que possuímos para extrair a informação biológica e quantificá-la. Apesar desta forma, a abordagem inicial usava métodos de quantificação de expressão gênica como Northern blot e quantitative PCR, passaram a serem substituídos por métodos como microarray entre o final de 1990 e 2000 e então por métodos baseados em síntese com RNAseq, em especial bulk-RNAseq e a partir de meados 2010-2013 sendo em grande medida métodos de transcriptômica de célula única. Portanto, houve a necessidade de se



desenvolver métodos que atendem estas demandas. Para métodos ~~como~~ de alto rendimento onde é possível interrogar a expressão gênica de forma genome-wide é necessário extrair a informação, ~~dele~~ normalizá-la e quantificá-la, ~~de~~ buscando posteriormente buscar ~~seus~~ métodos estatísticos para analisá-los.

Apesar de estar em desuso atualmente Microarray deu o primeiro impulso em método com putações voltados para expressão gênica, e um dos métodos desenvolvidos para análise de expressão diferencial é o LIMMA, baseado em regressão linear, que é suficiente para lidar com a variância de dados derivado do sinal fluorescente derivado dos chips. Ao passar para uso de bulk-RNA seq é possível olhar de forma mais ampla para o conjunto de transcritos devido ao fato de não ser mais dependente de sondas com as sequências previamente conhecidas são usados métodos por síntese, especificamente sequencing-by-synthesis; com por exemplo illumina. Permitindo encontrar e quantificar não apenas regiões previamente notadas mas também novos isoformas e regiões por exemplo não-codificante de proteínas previamente desconhecidas.

A partir do uso de microarray e bulk-RNA seq começaram a obter ~~dados~~ uma quantidade de dados massivos. E a partir desse momento foi possível construir grandes redes gênicas reguladas (GRN - gene regulatory network).

A construção dessas é baseada por exemplo em níveis de co-expressão, ~~grupos~~ set-genes (grupos de genes de são conhecidos por estarem em um viz em comum) e ~~construção~~ modelos de correlação e regressão. Combinando essas informações com informações prévias de ~~notas~~ é possível reconhecer genes regulados por exemplo pelo mesmo fator de transcrição.

Um grande desafio encontrado pela bioinformática na análise de dados envolvidos na regulação gênica é a abundância dos dados, e com isso quer dizer a quantidade



de dados que é necessário gerar e preparar. Dessa maneira mesmo se tratando de ~~de~~ microarray e bulk-RNA seq é necessário organizar os dados de forma que possa extrair insights com relevância biológica e ao mesmo tempo ter ferramentas para verificar a qualidade e dispersão dos dados. Entre essas formas usando clusterização (clustering), isso é agrupar ~~em~~ genes com perfil de expressão semelhantes. Para isso usando ~~para~~ abordagens de bioinformática como clusterização hierárquica, sendo possível aplicar diferentes métodos, como inferência euclidiana ou então usando K-means. A partir disso é possível obter uma representação de genes ^{com up ou down regulados ao comparar diferentes grupos experimentais} sendo essa representação ainda mais clara ao visualizarmos usando heatmaps, ~~usando~~ ^{usando} (gerados com ~~heatmaps~~ ^{heatmap} ~~em~~ ^{em} R). Ao termos usando visão olística, em seguida por reduzir a dimensionalidade destes dados, uma vez que se trata de milhares de genes com alta variabilidade de abundância de transcritos, isso é feito ao usar análise de Componentes Principais (PCA), normalizado para log (facilitando a visualização em escala) de TPM (transcritos por milhão). Ao usar essa abordagem é possível ver quais componente ~~tem~~ ^{possuem} maior variabilidade ou explicar a variabilidade e distância entre grupos, por exemplo ao comparar dois tipos celulares distintos, como neurônio motor e fibroblasto.

Para encontrar agora genes diferencialmente expressos aplicam-se métodos estatísticos com teste-t ao se comparar dois grupos experimentais ou ANOVA ao se comparar mais de dois grupos experimentais entre si.

Uma vez que temos uma visão geral do perfil de expressão gênica utilizando os métodos bioinformáticos descritos acima para detectar ~~expressões~~ ^{expressões} e quantificar a expressão gênica, pode-se buscar entender por exemplo como ou quais são os mecanismos que



levaram estes genes a serem diferenciatamente expressos.

Assim, precisamos agora usar métodos em bioinformática que nos ajudem a identificar dentro desse rede de genes o GW possui em suas estruturas moleculares semelhantes ou distintas. É possível fazer isso olhando para além do gene body (exon-intron) do gene, isso é para regiões regulatórias, como TSS (transcription start site), promotor proximal, promotor distal, enhancers e silencer regions, além de sítios de ligação de fatores de transcrição.

Métodos de bioinformática são capazes de prever/identificar motivos específicos baseado em dados experimentais previamente identificados ao buscar sequências consenso de fatores de transcrição, sequências conservadas baseadas em dados de alinhamentos múltiplos como Clustal e aplicando modelos como Hidden-Markovchain. Estes abordagens usam softwares como JASPAR, HOMER, MEME / TRANSFAC para a identificação de regiões de putativos binding-sites de fatores de transcrição e de sequência para início de transcrição. Um grande desafio neste campo é o grande número de falso positivo encontrados ao usar estes métodos. Portanto, é necessário combinar e alimentar métodos de Machine Learning supervisionado, Forest Trees e neural network para reconhecer padrões vindos de abordagens ~~como~~ que visam estabelecer por exemplo a acessibilidade de cromatina como ATAC-seq e métodos de buscar identificar regiões de interação Proteína-Proteína, como Chip-seq e CUT-and-run.

É fundamental ressaltar que métodos de bioinformática e métodos experimentais devem ser usados em conjunto. Assim como um único método de bioinformática por exemplo predição de fatores de transcrição não é suficiente para explicar a ~~regulação~~ mudança de expressão por exemplo.





Desta maneira vem a abordagem de bioinformática
 desenvolvida para integrar as diferentes camadas
 de mecanismos de regulação gênica é a multi-ômica,
 onde é possível integrar dados de expressão gênica
 derivados de RNA-seq, junto com dados de acessibilidade da
 cromatina disponíveis para transposase (ATAC-seq), mais
 dados de chip-seq para analisar a ligação de um fator de
 transcrição ou modificação específica de histona ex
 H3K9me3. Ao combinarmos estes dados temos
 a possibilidade de encontrar fatores em comum
 correlacionados entre si. Isso é possível ao usar
 modelos de variáveis latentes como uma forma
 de redução de dimensionalidade para em seguida
 combinar estes dados em matriz de fatorização.
 seguindo o workflow para clusterização e identificação
 de fatores que correlacionam e convergem em
 mostrar que mecanismos em comum atuam para
 fazer com que um gene ou um grupo de genes
 sejam ~~regulados~~ ativamente expressos.

Até o momento olhamos e buscamos entender
 redes regulatórias com métodos que foram gerados
 uma medida de expressão gênica derivada de um grupo de
 células, usando por exemplo bulk-RNA-seq. No entanto,
 como mencionado anteriormente entre métodos de
 2010-2013 tentativas de se obter dados de expressão
 gênica derivados de células únicas cresceram e tiveram
 um boom a partir de meados de 2017, com abordagens
 com drop-seq exemplo 10x chromium e smt-seq.

~~estes dados~~ a bioinformática aplicada a redes de regulação
 usando dados derivados de transcriptoma e ATAC-seq
 de células únicas nos levam a ter maior resolução,
 a nível de quais redes regulatórias estão sendo utilizadas
 por exemplo durante a diferenciação de uma
 célula tronco germinativa ~~de~~ humana até estágios
 de gestulação, sendo possível fazer a mudança de



redes regulatórias presentes em mudanças de expressão gênica e mudança de isoformas de de-transcritos derivados de splicing diferencial, usando ferramentas como Monocle e RNA velocity.

É claro que a combinação de dados multi-ômicos deve ser integrada com análises como identificação de tipos celulares identificados assim como estudos celulares vindos de scRNAseq usando por exemplo Cluster profiler em R.

E a partir disso combinar com informações de gene ontology, usando DAVID por exemplo ou IPA (Qiagen). A análise de gene ontology se faz crucial ao usar tanto bulk RNA seq quanto sc-RNAseq. Ou mesmo a combinação paralela de sc-RNAseq mais sc-ATAC-seq.

Para concluir se faz necessário pontuar que a integração de múltiplos ômicos e a predição e validação são de modelos e ferramentas de bioinformática por métodos experimentais é a chave para entendermos redes de regulação gênica, assim como a integração de dados de interação proteína-proteína derivados de ~~yeast two hybrid assay~~ ensaio duplo híbrido de levedura e predições derivadas deste dados, assim como de co-immunoprecipitação seguida de mass espec. (espectrometria de massa). Uma vez combinada a informação vinda destes métodos softwares de produção e visualização de redes como cytoscape podem ser usados para visualizar redes de regulatórias envolvendo múltiplos caminhos de regulação molecular heveladas ao utilizar abordagens de bioinformática.



- **topico 7** - Bioinformática aplicada a análise epigenômica.

- A pergunta de como ~~uma~~ a informação contida em uma célula ~~para múltiplos tempos~~ leva a multitude de tipos celulares é chave para ~~o~~ o campo da epigenômica. Porque a epigenética por si é o conjunto de modificações no DNA possuindo informação biológica sem alterar a sequência de nucleotídeos que compõe o DNA.

Desta forma ao analisar o conjunto de modificações no DNA desde metilação em citosina até mesmo pós-modificação em ~~casos~~ de histonas como acetilação ^{tradicional} metilação, ubiquitinação e fosforilação traz informações de acessibilidade de informação contida ~~no~~ no DNA e portanto fundamental para que circuitos genéticos e programas celulares sejam utilizados. Para por exemplo ~~obtidos de~~ a diferenciação ou reprogramação (iPS) de tipos celulares seja obtida.

Métodos de bioinformática para acessar por exemplo citosinas metiladas, em especial dinucleotídeos CpG ~~podem~~ são aplicados ao utilizar métodos de sequenciamento de alto rendimento por exemplo Bisulfite sequencing usando illumina paired-end após tratar o DNA com bissulfite; levando a ter no final do ^{processo} conversão citosina para timina. É a partir disso, é possível identificar citosinas metiladas, uma vez que citosinas metiladas não são convertidas para timina. Ao produzir bibliotecas de BS-seq ou RR-BSseq (reduced representation bisulfite sequencing) é possível ter uma visão geral do metiloma. Isso graças ao uso de softwares como ~~o~~ ^{padrão} para análise qualidade e limpeza dos reads como FASTQC, seguida de normalização, análise de ~~com~~ ^{com} ~~por~~ ^{por} ~~centos~~ ^{centos} principais e análises específicas para metilação



Como Methkit e Bismark para ~~essa~~ análise de reads ~~o~~ metilados e não metilados. A partir disso, busca-se ~~identificar~~ mapear onde estão os dinucleotídeos CpG ao longo do genoma, assim com CpG islands contendo de 200-3000 repetições. Uma vez que se identifica sítios metilados e não metilados o mapeamento usando por exemplo com UCSC genome Browser ou IGV reads metilados ao redor de por exemplo gene body (exon-intron). Uma vez que ~~em~~ a metilação de citosina está relacionada com regiões não acessíveis a maquinaria que controla a transcrição, esta abundância de bioinformática pode ser associada a inativação da expressão gênica ou na ausência de metilação de citosina com região após a setem transcriçionalmente ativa.

É fundamental lembrar que um único método de bioinformática sozinho não leva a uma compreensão abrangente de epigenômica. Por isso é necessário combinar métodos como ATAC-seq (chromatina acessível ao ataque de transposase TNS - seguida de sequenciamento). Com esta abordagem é possível detectar regiões de cromatina em estado de eucromatina (aberta) e disponível a transposição por TNS, e a partir disso com a inclusão de adaptadores illumina, por exemplo, é possível amplificar essa fragmento ^{por PCR} e em seguida sequenciá-los. Ao mapear os reads derivados da análise seguida de análise de qualidade e remoção de adaptadores e normalização, ~~com o~~ ~~se~~ é possível identificar regiões de cromatina acessível. É portanto, a partir desse momento identificar ~~presença~~ ~~de~~ ~~motivos~~ e informações prévias se esta região é ~~em~~ uma região regulatória (promotor/enhancer) ^{em} que pode conter sítios putativos para ligação de fatores de transcrição. É possível fazer uma



busca por estes sítios utilizando preditores de fatores de transcrição e elementos regulatórios como MEME, TRANSFAC e JASPAR por exemplo.

No entanto, graças a outras abordagens de bioinformática podemos buscar validar experimentalmente estes sítios putativos. A abordagem funciona da seguinte forma, ao usarmos preditores podemos encontrar o enriquecimento de sítios putativos binding sites na região upstream do gene body, dentro de uma janela genômica de 3Kb (3 mil pares de bases), um exemplo é Kappa B sites, desta forma o pesquisador pode direcionar o experimento de Chip-seq para usar o anticorpo anti-NFKB (p65). Ao usar Chip-seq como sugerido pelo nome imunoprecipitação de cromatina seguida de sequenciamento. Desta maneira no precipitado imunoprecipitado o DNA genômico com anti-NFKB, o pesquisador ~~deve~~ ~~obter~~ fragmenta o DNA, por exemplo por sonicação e precipita fragmentos que são enriquecidos ~~na~~ de proteínas de interesse.

A partir desse momento, ~~o~~ se ~~o~~ ~~amplifica~~ ~~o~~ por PCR os fragmentos ~~de~~ foram enriquecidos na imunoprecipitação, ~~os~~ prepara os fragmentos para bibliotecas illumina e analisa ~~o~~ as qualidades dos reads com FastQC e em seguida normaliza ~~o~~ pelo background do amostra controle sem imunoprecipitação. Ao final 1990, agora é possível validar ~~o~~ se os putativos sítios de K_B sites contem a interação DNA-Proteína predita.

Ao combinar metiloma, ATAC-seq e Chip-seq é possível ter uma visão com perspectiva de mediadores, epigenômicos controlando por exemplo a resposta a inflamação ou processos de diferenciação celular ou mesmo fatores chave que ~~o~~ são suficientes

10 de 13



UNIVERSIDADE FEDERAL DO RIO DE JANEIRO
INSTITUTO DE CIÊNCIAS BIOMÉDICAS



para por exemplo reprogramar uma célula
diferenciada e se torna pluripotente novamente
como ocorre com células tronco de pluripotência
indolida, com



tópico 9 - Aplicação de bioinformática no estudo das patologias humanas.

A bioinformática tem um papel fundamental no entendimento e busca de terapias para patologias humanas.

Com o avanço e barateamento de sequenciamento ~~de~~ de genomas inteiros ~~de~~ (WGS) assim como de WES, hoje é possível sequenciar o genoma ou exoma de milhares de indivíduos saudáveis ou com patologias, sejam elas familiares ou esporádicas. A partir do sequenciamento desses genomas e exomas inteiros é possível montar usando o método de re-sequencing usando short reads Illumina ou long-reads (PACBio HiFi, Oxford Nanopore) para encontrar variantes de nucleotídeos únicos, indels ou variantes estruturais.

Com estes dados é possível usar PLINK para conduzir estudos de GWAS (genome-wide association studies). Desta forma encontram o odds-ratio de uma determinada variante genética ser benéfica ou patogênica dentro de uma população ou coorte.

A partir disso é possível buscar ~~variação~~ variantes associadas desde doenças de alta frequência comum até doenças raras. Portanto ferramentas como GATK e SAMtools são essenciais para encontrar variantes no processo de variant calling.

Uma vez que variantes foram associadas com uma determinada patologia, é possível estratificar estes pacientes, o que vem se tornando fundamental para a utilização de novas terapias que são desenvolvidas de forma personalizada. Um exemplo é a estratificação de pacientes



com a doença neurodegenerativa para esboçar
intervalo a miotônica. devido a natureza de ter
múltiplos genes e variantes relacionadas,
descobriu usando métodos de bioinformática
que terapias com drogas específicas são
efetivas para genes e variantes específicas.

É importante ressaltar que a bioinformática
pode ser usada muito antes dos pacientes
manifestarem seus sintomas/fenótipos. Um exemplo é Huntington,
onde é possível detectar a expansão somática
de polyQ sendo aumentada antes do paciente
manifestarem sintomas ao usar amplicon seq
seguido de mapeamento de reads detectando
a expansão.

Para o câncer, hoje se aplica exome-seq e
WGS para detectar variantes potencialmente
cruciais como driver mutations, mas vai
além ao empregar RNA-seq (RNA-seq) combinado
com a transcriptoma spatial.

Ao usar singletons com Seurat(R) é possível
revelar a heterogeneidade celular desde
tipos e estados celulares e por comparação usando
MAST(R), revelando população complexas
sendo visualizadas com redução de dimensionalidade
com UMAP ou t-SNE.

Porém, ao condicionar RNA-seq perdemos
uma informação fundamental que é a espacial
onde estas células estão, e por exemplo se
tratando de câncer onde está o infiltrado de
células imunes. Ao usar Visium 10x para gene
wide transcriptome com resolução de spots
ou ao usar MERFISH com áreas direcionadas,
baseado em métodos únicos derivados de
RNA-seq é possível do pacote SPARCS identificar



13 de 13

transcritos únicos para MER fish. MER fish,
no contrario de Visum resolve transcriptoma
especial com resolução sub-celular.

Ao encontrar cluster únicos residuais
especialmente e com essa resolução celular
única, é possível selecionar marcadores
únicos para drug discovery/ desenho em
estrutura.