

Minicurso: Ferramenta rSWeeP – método para representação vetorial de sequências biológicas.

Camila P Perico
Escola Paranaense de Bioinformática – 2024

Universidade Federal do Paraná

6–8 de Novembro de 2024



Organização do Minicurso

- 1 O que são métodos livres de alinhamento (AF)
 - problemas dos métodos baseados em alinhamento
 - tipos de métodos AF
- 2 SWeeP
 - descrição do método
 - pacote rSWeeP
- 3 **Aplicações e Prática com rSWeeP**
 - Aplicação em proteoma mitocondrial e quickStart
 - Aplicação em SARS-CoV-2 e Prática 1
 - Parametrização
 - Aplicação em proteomas bacterianos
- 4 Conclusão

Métodos livres de Alinhamento

Qual o grande desafio atual da bioinformática (e outras áreas)?

Grande volume de dados (**big data**)

Novas tecnologias → mais dados coletados

- Grande número de sequências biológicas
- Sequências biológicas longas (genomas, metagenomas etc)
- Transcriptomas de milhares de células (scRNAseq)
- entre outros

Desafios na análise de sequências biológicas

Desafios:

- Dificuldade atual em analisar grandes quantidades de dados moleculares
- **Alto custo computacional** para realizar **alinhamento** (método tradicional)

Soluções:

- Métodos **livres de alinhamento** vêm se mostrando vantajosos para Big data¹
- Particularmente métodos de **representação vetorial** de sequências biológicas possibilitam usar técnicas de **ML** e aplicações mais efetivas em bioinformática²

¹ Asgari and Mofrad (2015) Continuous distributed representation of biological sequences for deep proteomics and genomics. PloS one, 10(11).

² Randhawa, et al. (2020) Machine learning using intrinsic genomic signatures for rapid classification of novel pathogens: Covid-19 case study. Plos one, 15(4)
Leimeister, et al. (2019). Prot-spam: Fast alignment-free phylogeny reconstruction based on whole-proteome sequences. GigaScience, 8(3)
Zhang, et al. (2017) Viral phylogenomics using an alignment-free method: A three-step approach to determine optimal length of k-mer. Scientific reports, 7(1)

Porque usar métodos livres de alinhamento (AF)?

5 casos em que métodos de alinhamento falham (e AF não)³:

- 1 Alinhamentos pressupõem segmentos homólogos linearmente arranjados.
Mundo real: genomas são **mosaicos**.
Alinhamentos são sensíveis a alta taxa mutacional, recombinação, duplicações, transferência horizontal de genes, ganho/perda de genes
- 2 acurácia dos alinhamentos cai rapidamente conforme **diminui a identidade**
- 3 alinhamentos são muito **custosos em RAM** e tempo de processamento
- 4 limitação de escala
múltiplos genomas → impossível de resolver em tempo realista
- 5 Pressupostos na **parametrização**: escolha da matriz de substituição, penalidade de gaps, limites estatísticos etc. Diferentes parâmetros → **resultados diferentes**

³ Zieleszinski, et al. (2017). Alignment-free sequence comparison: benefits, applications, and tools. Genome biology, 18

Métodos livres de alinhamento (AF)

Definição *Métodos AF podem ser definidos como qualquer método para quantificar a similaridade/dissimilaridade de sequências que não use ou produza alinhamento em nenhuma etapa da aplicação do algoritmo*⁴

Métodos AF **não** enfrentam os mesmos desafios dos métodos baseados em alinhamentos

Há 2 classes principais:

- 1 *word-based* ou frequência de palavras (k-mers): calculam similaridade/dissimilaridade baseada na frequência de subsequências
- 2 baseados em teoria da informação: avaliam o conteúdo informacional completo entre duas sequências

****** *word-based* se mostram superiores na detecção de padrões complexos

⁴ Zieleszinski, et al. (2017). Alignment-free sequence comparison: benefits, applications, and tools. Genome biology, 18

Métodos livres de alinhamento (AF)

Os métodos AF estão se tornando cada vez mais presentes na literatura⁵

Query sequences	x	ATGTGTG	y	CATGTG
Word size: 3	W_3^x	ATG TGT GTG TGT GTG	W_3^y	CAT ATG TGT GTG
3-mer				
Union of two sets	$W_3 = W_3^x \cup W_3^y$	CAT ATG TGT GTG		
Word counts	c_3^x	0 1 2 2	c_3^y	1 1 1 1
Euclidean distance	$\ c_3^x - c_3^y\ $	$\sqrt{(0-1)^2 + (1-1)^2 + (2-1)^2 + (2-1)^2} = \sqrt{3} = 1.73$		

word-based ou frequência de palavras (k-mers)

Considere duas sequências de DNA:
ATGTGTG e CATGTG

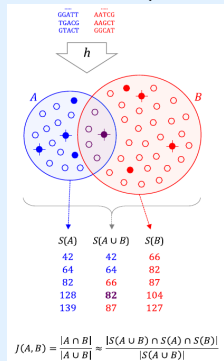
A distância euclidiana entre elas é 1.73

Adaptado de: Zielezinski, et al (2017)

⁵ Forsdyke (2019). Success of alignment-free oligonucleotide (k-mer) analysis confirms relative importance of genomes not genes in speciation and phylogeny. Biological Journal of the Linnean Society, 128(2)
Zielezinski, et al. (2017). Alignment-free sequence comparison: benefits, applications, and tools. Genome biology, 18

Alguns importantes métodos AF

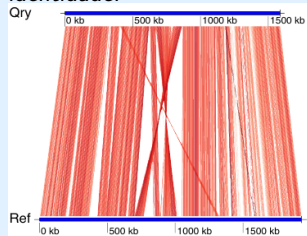
Mash → técnica baseada em estatística de conjuntos (índice de Jaccard) de k-mers



Ondov, et al (2016). Mash: fast genome and metagenome distance estimation using minhash. Genome biology, 17:1–14.

fastANI (Average Nucleotide Identity) → **NCBI**

fragmentação de sequências de genoma, seguida de pesquisa de sequência de nucleotídeos, alinhamento e cálculo de identidade.



Jain, et al (2018). High throughput ani analysis of 90k prokaryotic genomes reveals clear species boundaries. Nature communications,9(1):5114.

Alguns importantes métodos AF

Mash → técnica baseada em estatística de conjuntos (índice de Jaccard) de k-mers
Ondov, et al (2016). Mash: fast genome and metagenome distance estimation using minhash. *Genome biology*, 17:1–14.

fastANI (Average Nucleotide Identity)
→ **NCBI**

fragmentação de sequências de genoma, seguida de pesquisa de sequência de nucleotídeos, alinhamento e cálculo de identidade.

Jain, et al (2018). High throughput ani analysis of 90k prokaryotic genomes reveals clear species boundaries. *Nature communications*, 9(1):5114.

** ambos baseados em nucleotídeos

** ambos retornam um índice de similaridade (matriz de distâncias), a partir de FASTAS fornecidos

** ambos podem analisar genomas completos

O Método SWeeP

O que é SWeeP?

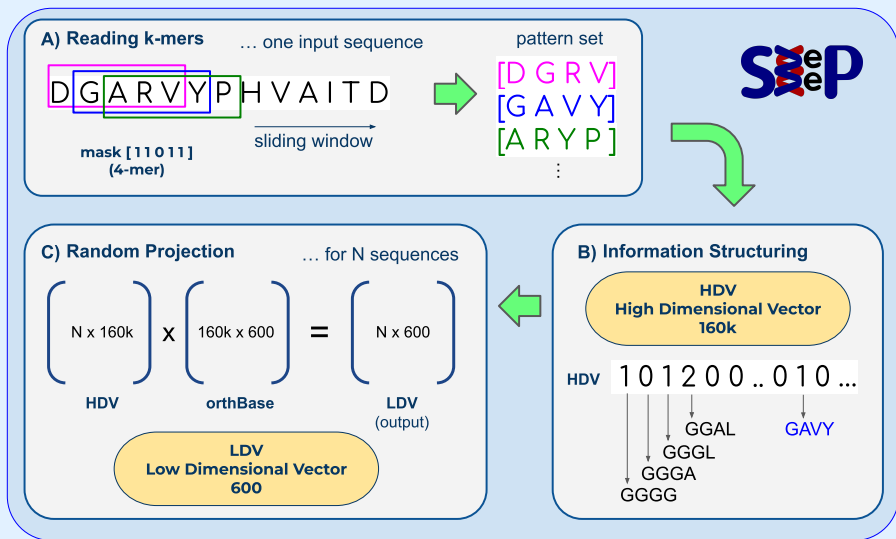
SWeeP (Spaced Words Projection)

método de representação vetorial de sequências biológicas⁶

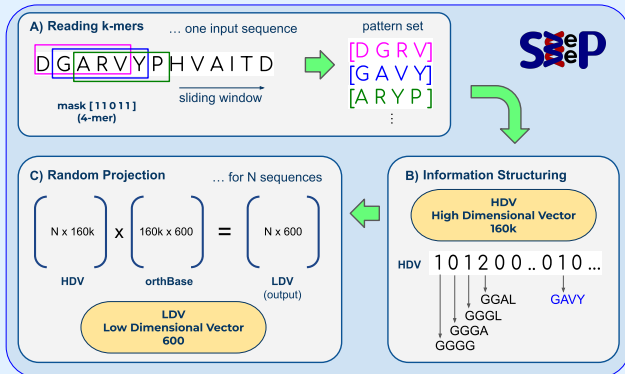
Ferramenta **AF** *word-based* que representa as sequências (nucleotídicas ou proteicas) em vetores compactos por meio de **Random Projections** (apoiada pelo lema de Johnson-Linderstrauss JL⁷).

⁶ De Pierri, et al. (2020) Sweep: representing large biological sequences datasets in compact vectors. Scientific reports, 10(1)

⁷ Johnson, Lindenstrauss (1984) Extensions of Lipschitz mappings into a Hilbert space, Contemp. Matemática



Processamento: embedding por projeção randômica da frequência de k-mers de sequências biológicas.



Propriedades:

- preserva distâncias entre as amostras no espaço de menor dimensão
- vetor numérico estruturado e diretamente comparável
- vetor de dimensões reduzidas

Terminologia

embedding conversão de um vetor de alta dimensionalidade em outro de baixa dimensionalidade (ex. projeção randômica, RP, ou PCA)

alta dimensão o número de características é maior que o número de observações (amostras)

HDV vetor estruturado de alta dimensão (frequência de k-mers, vetor contagem RNAseq)

LDV projeção do HDV em vetor compacto

k-mer oligonucleotídeo

máscara [1111] e [11011] são duas máscaras possíveis para 4-mer

k-mers e significado biológico

O uso de *k-mers* é válido? Traz informações relevantes da evolução das espécies?

Segundo Forsdyke (2019), *k-mers* possuem **justificativa biológica**⁸:

- Variações na frequência de *k-mers* está fortemente relacionado à **especiação** (redução na capacidade de hibridização das fitas de DNA)
espécies possuem frequência de *k-mers* bem determinado
- **códons preferenciais** → tendem a ser preservados na espécie
- pode revelar relações que métodos de alinhamento não conseguiriam

Em DNA:

- 1-mer - (%GC) permite diferenciar principais grupos de eubacterias
- 3-mer - genomas específicos
- 4-mer - padrão ouro para vírus e procariotos
- N-mer - maior capacidade de discriminação até nível de indivíduo

⁸ Forsdyke, D. R. (2019). Success of alignment-free oligonucleotide (*k*-mer) analysis confirms relative importance of genomes not genes in speciation and phylogeny. *Biological Journal of the Linnean Society*, 128(2)

Desafios na alta dimensionalidade

Quanto maior o k-mer, maior o vetor HDV (vetor contagem)

$$1\text{-mer} = 4$$

$$2\text{-mer} = 16$$

$$4\text{-mer} = 256$$

$$7\text{-mer} = 16384$$

$$15\text{-mer} = 268\,435\,456$$

$$25\text{-mer} = 1.1259 \times 10^{15} \text{ (ref}^9\text{)}$$

maior k	→	maior o vetor HDV	→	maior consumo de RAM	→	maior a dificuldade em analisar os vetores
-----------	---	----------------------	---	-------------------------	---	---

Métodos de ML têm dificuldade de lidar com a alta dimensionalidade (tempo computacional)

Necessário o uso de métodos para a **redução da dimensionalidade**

⁹ Chan CX, Bernard G, Poirion O, Hogan JM, Ragan MA. Inferring phylogenies of evolving sequences without multiple sequence alignment. Sci Rep. 2014;4:6504.

Redução de Dimensionalidade

3 técnicas básicas de redução de dimensão linear: **RP**, **PCA**, **FS**

RP – projeções randômicas ; **PCA** – Análise de Componentes Principais; **FS** – seleção de características

¹⁰ Grellmann, et al. (2016). Random projection for fast and efficient multivariate correlation analysis of high-dimensional data: A new approach. *Frontiers in Genetics*, 7:102.

Ray, et al. (2021) Various dimension reduction techniques for high dimensional data analysis: a review. *Artificial Intelligence Review*, 54

Redução de Dimensionalidade

3 técnicas básicas de redução de dimensão linear: **RP, PCA, FS**

Limitações:¹⁰

- **FS, PCA** – depende explicitamente do conjunto de dados de entrada
- **PCA** – computacionalmente custoso com tempo de processamento quadrático em relação ao número de dimensões (características)
- **FS** – perda de informação potencialmente relevante
- **FS** – problema dependente
- **RP, PCA** – dificuldade em alocar altíssima dimensionalidade em RAM

Não foi encontrado na literatura sintetize do gap na redução de dimensionalidade.

Então, introduzimos **nosso conceito de generalidade** e propomos uma **solução** que satisfaz as limitações dos métodos atuais.

¹⁰ Grellmann, et al. (2016). Random projection for fast and efficient multivariate correlation analysis of high-dimensional data: A new approach. *Frontiers in Genetics*, 7:102.

Ray, et al. (2021) Various dimension reduction techniques for high dimensional data analysis: a review. *Artificial Intelligence Review*, 54

GAPS nos métodos AF e Perspectiva de Generalidade

Na ausência de um termo específico, vamos utilizar a designação de “**generalidade**” como um método que incorpora os seguintes conceitos:

- a) **usabilidade** → capaz de realizar *embedding* de alta dimensionalidade¹¹.
- b) **conservação das distâncias** (local e global) → preserva informação¹²
- c) método **conjunto independente** → *embedding* **reprodutível** a qualquer conjunto subsequente, dada a mesma parametrização

Impomos tais exigências para obter de um método:

geral, implementável, representação vetorial de sequências universal

**Nenhum método da literatura satisfaz tais propriedades
(exceto SWeeP)**

¹¹ Grellmann, et al. (2016). Random projection for fast and efficient multivariate correlation analysis of high-dimensional data: A new approach. *Frontiers in Genetics*, 7:102.

¹² Heiser e Lau (2020). A quantitative framework for evaluating single-cell data structure preservation by dimensionality reduction techniques. *Cell reports*, 31(5)

Motivação

Principal motivação generalização da ferramenta → nova versão: SWeePlite

Não se limita a sequências biológicas (genes, genomas, proteínas e proteomas)¹³ mas possui **potencial** para:

metagenomas, espectrometria de massa, expressão gênica, análise multiômica, textos etc

SWeeP já foi implementado em **mineração de textos**¹⁴ de forma bem sucedida.

¹³ da Silva Filho, et al. (2021). Prediction and analysis in silico of genomic islands in aeromonas hydrophila. Frontiers in microbiology, 12.

Raittz, et al. (2021). Comparative genomics provides insights into the taxonomy of azoarcus and reveals separate origins of nif genes in the proposed azoarcus and aromatoleum genera. Genes, 12(1)

¹⁴ leger Raittz, et al. (2022). What are we learning with yoga: a text mining approach to literature. medRxiv

Consumo memória do vetor HDV

quero projetar uma sequência de **aminoácidos** em um vetor de comprimento **1000**

SWeeP Original

máscara	k-mer	# combinações	tamanho da matriz	memória (float)
11011	4	20^4	$1.6 \cdot 10^5 \times 1000$	640 Mb
110111	5	20^5	$3.2 \cdot 10^6 \times 1000$	12.8 Gb
1110111	6	20^6	$6.4 \cdot 10^7 \times 1000$	256 Gb

SWeePlite

uso de memória: 80 - 800 Mb (parametrizável)

**permite alocar a matriz completa se quiser*

O Pacote rSWeeP

<https://bioconductor.org/packages/release/bioc/html/rSWeeP.html>

<https://bioconductor.org/packages/devel/bioc/html/rSWeeP.html>

Package ‘rSWeeP’

October 5, 2024

Type Package

Title Spaced Words Projection (SWeeP)

Version 1.17.3

Date 2024-04-28

Maintainer Camila P Perico <camilapp94@gmail.com>

Description ``Spaced Words Projection (SWeeP)'' is a method for representing biological sequences using vectors preserving inter-sequence comparability.

Depends foreach, doParallel, parallel, Biostings, methods, utils

Imports tools, stringi,

Suggests Rtsne, ape, Seurat, knitr, rmarkdown, tictoc, BiocStyle, testthat (>= 3.0.0)

Encoding UTF-8

RoxygenNote 7.3.2

VignetteBuilder knitr

Config/testthat/edition 3

O Pacote rSWeeP

<https://bioconductor.org/packages/release/bioc/html/rSWeeP.html>

<https://bioconductor.org/packages/devel/bioc/html/rSWeeP.html>

Contents

rSWeeP-package	2
extractHDV	3
orthBase	5
PCCI	6
PMPG	7
SWeeP	8
SWeePlite	13

2 funções principais e 4 auxiliares:

- SWeeP – função SWeeP original, necessita que forneça a matriz de projeção
- SWeePlite – versão lite com função FGOP para gerar internamente matriz proj.
- orthBase – gera matriz de projeção
- extractHDV – extrai o HDV das sequências biológicas
- PMPG – Percentage of Mono or Paraphyletic Groups
- PCCI – PhyloTaxonomic Consistency Cophenetic Index. Estimativa do grau de agrupamento das amostras do mesmo taxon na árvore filogenética.

Função original:

- SWeeP – função SWeeP original, necessita que forneça a matriz de projeção
- orthBase – gera matriz de projeção

```
library(rSWeeP)

# mascara
mask = c(2,1,2)
# tamanho da projecao
psz = 600
# sequencia tipo aminoacido
seqtype="AA"
# endereco da pasta com FASTAS
path = "pasta/FASTAS/"

# crio a base ortonormal de projecao - 160000 x 600
base160k = orthBase(mask=mask,col=psz,seqtype='AA')

sw = SWeeP(path,base160k,seqtype=seqtype,mask=mask)
```

Função lite:

- SWeePlite – versão lite com função FGOP para gerar internamente matriz proj.

```
library(rSWeeP)

# mascara
mask = c(2,1,2)
# tamanho da projecao
psz = 600
# sequencia tipo aminoacido
seqtype="AA"
# endereco da pasta com FASTAS
path = "pasta/FASTAS/"

sw = SWeePlite(path,seqtype=seqtype,mask=mask,psz=psz)
```

Saída das funções SWeeP e SWeePlite

- `output$proj` uma matriz numérica com `psz` colunas e uma linha por sequência, correspondendo cada linha a um vetor compacto
- `output$info` informação adicional do processo. Este objeto subdivide-se em:
 - `$ProjectionSize`: um número inteiro correspondente a `psz` (o comprimento do vetor de saída)
 - `$bin`: um booleano que contém se é binário (TRUE) ou de contagem (FALSE)
 - `$mask`: um vetor que contém a máscara utilizada
 - `$SequenceType`: um carácter que contém o tipo da sequência (aminoácido: 'AA', ou nucleótido: 'NT')
 - `$concatenate`: um booleano que corresponde à concatenação de sequências
 - `$version`: um carácter correspondente à versão do pacote
 - `$norm`: um carácter que contém a normalização utilizada
 - `$extension`: um carácter que contém a lista de extensões consideradas
 - `$timeElapsed`: um double que contém o tempo decorrido em segundos
 - `$headers`: lista de cabeçalhos para cada sequência analisada

Funções adicionais

- `extractHDV` – extrai o HDV das sequências biológicas

```
library(rSWeeP)

# mascara
mask = c(2,1,2)
# sequencia tipo aminoacido
seqtype="AA"
# endereco da pasta com FASTAS
path = "pasta/FASTAS/"

output = extractHDV(path,mask=mask,seqtype=seqtype)
```

- `output$HDV` High Dimensional Vectors (vetores de alta dimensão) do FASTAS fornecido. A dimensão do vetor corresponde a 20^k para os aminoácidos e 4^k para os nucleótidos. k corresponde ao k-mer da máscara utilizada.
- `output$info` informação adicional do processo.

Funções adicionais

- PMPG – Percentage of Mono or Paraphyletic Groups
- PCCI – PhyloTaxonomic Consistency Cophenetic Index. Estimativa do grau de agrupamento das amostras do mesmo taxon na árvore filogenética.

```
library(rSWeeP)
[...]
```

```
mdist = dist(sw$proj,method='euclidean')
tree = ape::nj(mdist)
```

```
PMPG(tree,taxons)
PCCI(tree,taxons)
```

- PCCI\$mean retorna o valor médio do agrupamento
- PMPG\$percMono retorna o percentual de táxons monofiléticos
- PMPG\$percPara retorna o percentual de táxons parafiléticos
- PMPG\$mean retorna o percentual somado entre mono e paraifilético



Prática

PRÁTICA

Tutoriais

<https://aibialab.github.io/rSWeeP>

rSWeeP

Artificial Intelligence Applied to Bioinformatics Laboratory (AIBIA) - UFPR

rSWeeP

HOME

The *rSWeeP* package is an R implementation of the Spaced Words Projection (SWeeP) method (De Pierri, 2019). The main function of this package is to provide a vector representation of biological sequences (nucleotides or amino acids), and thus favor alignment-free phylogenetic studies. Each sequence provided is represented by a compact numerical vector which is easier to analyze. SWeeP uses k-mers counting for representing the sequences in high dimensional vector (HDV) and then projected into a low dimensional vector (LDV) through random projection using an orthonormal base. The LDV represents the biological sequence and is handable for comparative analysis and machine learning implements. In addition, the package allows general dimensionality reduction of RNAseq data and generic matrices.

INSTALAÇÃO PACOTES R NECESSÁRIOS

Instalação do rSWeeP:

```
# instalo o manager do Bioconductor
install.packages("BiocManager")

# instalo o rSWeeP depositado no Bioconductor
BiocManager::install("rSWeeP")
```

Pacotes necessários para a prática

```
install.packages("ape")
install.packages("umap")
install.packages("Rtsne")
BiocManager::install("ggtree")
```

Manual do pacote

```
?rSWeeP
```


Dúvidas

Dúvidas e dificuldades instalação R