

Nanopore automata

Ian Holmes^{1,2,*}

1 Lawrence Berkeley National Laboratory, Berkeley, CA, USA

2 Department of Bioengineering, University of California, Berkeley, CA, USA

Contents

1	Abstract	2
2	Specification	2
2.1	Parameterization algorithm	2
2.2	Reference search algorithm	2
2.3	Implementation	3
2.4	Evaluation	3
3	Methods	3
3.1	Model	3
3.2	Baum-Welch algorithm	5
3.3	Viterbi algorithm	5
4	Results	6
5	Discussion	6
6	Acknowledgments	7
7	Figure Legends	8
8	Appendix	9
8.1	Gamma distribution	9
8.2	Normal distribution	9

1 Abstract

State machine algorithms for aligning Nanopore reads. Initial goal is simple reusable code for aligning a nanopore read to a reference sequence. No attempt at optimization yet.

2 Specification

2.1 Parameterization algorithm

Given the following inputs

- Reference genome (FASTA)
- Segment-called reads (FAST5/HDF5)

Perform the following steps

- Perform Baum-Welch to fit a rich model

Rich model incorporates segment statistics.

2.2 Reference search algorithm

Given the following inputs

- Reference genome
- Segment-called reads (FAST5/HDF5)
- Parameterized rich model

Perform the following steps

- Perform Viterbi alignment

2.3 Implementation

Libraries etc.

HDF5...

2.4 Evaluation

Strategy...

Data sets...

3 Methods

Model & inference algorithms.

3.1 Model

- Order- N transducer.
- Input: nucleotide
- Output: nucleotide, segment mean, duration
- Emissions:
 - categorical (base k -mer)
 - mixture of Normal/gamma (mean/duration)
- Transitions:
 - *Match*: emit single segment, absorb 1 base
 - *Insert*: affine gap insertion of bases: emits segments, absorbs no bases
 - *Delete*: affine gap deletion of bases: emits no segments, absorbs bases
 - *Merge*: emit single segment, absorb 2 or 3 bases

- *Split*: emit single segment, absorb 0 bases
- *Skip*: emit single segment, absorb $2 \dots K$ bases (large K , low extension penalty)

This can be achieved by a Mealy transducer with 3×4^N states. The factor of 4^N accounts for the order- N context. For each such context, the three states are MAT, INS and DEL.

Parameters:

- Gap opening & extension probabilities θ_{go}, θ_{gx}
- Merge probability θ_{mo} , probability that it's a 3-merge is θ_{mx}
- Split probability θ_s
- Skip probability θ_{ko} , skip extension probability θ_{kx}

The transition table is as follows:

Transition	From	To	Weight	Input x	Output (y, m, d)
Match	MAT	MAT	$(1 - \theta_{go})(1 - \theta_{mo})(1 - \theta_s)(1 - \theta_{ko})$	$x \in \Omega$	$y \sim \text{Categorical}(\mathbf{p}_{x,c}^m),$ $m \sim \text{Normal}(\mu_{x,c}^m, \tau_{x,c}^m),$ $d \sim \text{Gamma}(\alpha_{x,c}^m, \beta_{x,c}^m)$
Insert	MAT	INS			
	INS	INS			
	INS	MAT			
Delete	MAT	DEL			
	DEL	DEL			
	DEL	MAT			
Merge	MAT	MAT			
Split	MAT	MAT			
Skip	MAT	MAT			

Here $y \in \Omega$ where Ω is the nucleotide alphabet and $c \in \Omega^N$ is the context.

3.2 Baum-Welch algorithm

3.3 Viterbi algorithm

4 Results

5 Discussion

6 Acknowledgments

7 Figure Legends

8 Appendix

8.1 Gamma distribution

$$x \sim \text{Gamma}(\alpha, \beta)$$

$$\mathbb{E}[x] = \alpha/\beta$$

$$\text{Var}[x] = \alpha/\beta^2$$

Shape parameter α , rate parameter β .

$$P(x|\alpha, \beta) = \frac{x^{\alpha-1} \beta^\alpha \exp(-x\beta)}{\Gamma(\alpha)}$$

where Γ is the gamma function

$$\Gamma(\alpha) = \int_0^\infty z^{\alpha-1} \exp(-z) dz$$

Note $\Gamma(n) = (n-1)!$ for positive integer n .

8.2 Normal distribution

$$x \sim \text{Normal}(\mu, \tau)$$

Mean μ , precision τ (precision is reciprocal of variance).

$$P(x|\mu, \tau) = \sqrt{\frac{\tau}{2\pi}} \exp\left(-\frac{\tau}{2}(x-\mu)^2\right)$$