

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	3
2.3	Implementation	3
2.4	Evaluation	3
3	Methods	3
3.1	Model	4
3.2	Baum-Welch algorithm	4
3.3	Viterbi algorithm	5
4	Results	6
5	Discussion	6
6	Acknowledgments	7
7	Figure Legends	8
8	Appendix	9
8.1	Exponential distribution	9
8.2	Gamma distribution	9

1 Abstract

State machine algorithms for aligning Nanopore reads.

2 Specification

Initial goal (Preliminary Results) is simple reusable code for aligning a segmented nanopore read (with segment currents summarized) to a reference sequence.

Longer-term goals (Specific Aims) include

- quasi-hierarchical series of models for processed→raw data (raw, FAST5, FASTQ, FASTA)
- transducer intersection-style models for read-pair alignment
- systematic strategies for approximation/optimization algorithms, climbing the hierarchy (starting with k-mer or FM-index approaches)
- transducer intersection models for aligning reads from different sequencing technologies

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: current levels
- States: Start, Match $_{x_1 \dots x_N}$, Delete $_{x_1 \dots x_N}$, End
- Transitions
 - Start \rightarrow Start: weight $p^{\text{StartEmit}}$, emits current $y \sim \text{Normal}(\mu^{\text{Start}}, \tau^{\text{Start}})$
 - Start \rightarrow Match $_{x_1 \dots x_N}$: weight $1 - p^{\text{StartEmit}}$, absorbs $\geq N$ bases ending in $x_1 \dots x_N$
 - Match $_{x_1 \dots x_N} \rightarrow$ Match $_{x_1 \dots x_N}$: weight $p^{\text{MatchEmit}}$, emits current $y \sim \text{Normal}(\mu^{\text{Match}}_{x_1 \dots x_N}, \tau^{\text{Match}}_{x_1 \dots x_N})$
 - Match $_{x_1 \dots x_N} \rightarrow$ Match $_{x_2 \dots x_{N+1}}$: weight $(1 - p^{\text{MatchEmit}}_{x_1 \dots x_N})(1 - p^{\text{BeginDelete}})$, absorbs 1 base x_{N+1}
 - Match $_{x_1 \dots x_N} \rightarrow$ Delete $_{x_2 \dots x_{N+1}}$: weight $(1 - p^{\text{MatchEmit}}_{x_1 \dots x_N})p^{\text{BeginDelete}}$, absorbs 1 base x_{N+1}
 - Match $_{x_1 \dots x_N} \rightarrow$ End: weight $1 - p^{\text{MatchEmit}}_{x_1 \dots x_N}$, absorbs any string of bases
 - Delete $_{x_1 \dots x_N} \rightarrow$ Delete $_{x_2 \dots x_{N+1}}$: weight $p^{\text{ExtendDelete}}$, absorbs 1 base x_{N+1}
 - Delete $_{x_1 \dots x_N} \rightarrow$ Match $_{x_1 \dots x_N}$: weight $1 - p^{\text{ExtendDelete}}$

3.2 Baum-Welch algorithm

As usual.

3.3 Viterbi algorithm

As usual.

4 Results

5 Discussion

6 Acknowledgments

7 Figure Legends

8 Appendix

8.1 Exponential distribution

$$\begin{aligned}
 x &\sim \text{Exponential}(\kappa) \\
 P(x|\kappa) &= \kappa \exp(-\kappa x) \\
 \mathbb{E}[x] &= \kappa^{-1} \\
 \text{Var}[x] &= \kappa^{-2}
 \end{aligned}$$

Rate parameter κ .

8.2 Gamma distribution

$$\begin{aligned}
 x &\sim \text{Gamma}(\alpha, \beta) \\
 P(x|\alpha, \beta) &= \frac{x^{\alpha-1} \beta^\alpha \exp(-x\beta)}{\Gamma(\alpha)} \\
 \mathbb{E}[x] &= \alpha/\beta \\
 \text{Var}[x] &= \alpha/\beta^2
 \end{aligned}$$

Shape parameter α , rate parameter β . $\Gamma()$ is the gamma function

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

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

8.3 Normal distribution

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

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

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