



EpiSegMix

A Flexible Distribution Hidden Markov Model with Duration Modeling for Chromatin State Discovery

Johanna E. Schmitz, Nihit Aggarwal, Lukas Laufer, Jörn Walter, Abdulrahman Salhab, Sven Rahmann

DSB 2024

Algorithmic Bioinformatics







Biological Background









Biological Background



Modification	Function	Genomic Context
H3K9me3	Transcriptional repression	Heterochromatin
H3K27me3	Transcriptional repression	CpG-rich promoters, intergenic regions
H3K36me3	Transcriptional elongation	Gene bodies
H3K27ac	Transcriptional activation	Promoters, enhancers
H3K4me1	Transcriptional activation	Promoters, enhancers, intergenic regions
H3K4me3	Transcriptional activation	Promoters





Chromatin Immunoprecipitation with DNA sequencing (ChIP-seq)



Crosslinking, fragmentation and immunoprecipitation

Adapter III J Sequencing library

DNA purification and sequencing libary preparation



Mapping reads to reference genome







Chromatin Segmentation



Algorithmic Bioinformatics



SIC Saarland Informatics Campus



Introduction to Hidden Markov Models









Introduction to Hidden Markov Models









Introduction to Hidden Markov Models



Algorithmic Bioinformatics







Hidden Markov Models for Chromatin Segmentation



2. Probabilistic Model

3. Segmentation



¹ Ernst, Kellis (2010) 2 Mammana, Chung (2015)





Hidden Markov Models for Chromatin Segmentation

1. Binned read counts

2. Probabilistic Model

3. Segmentation

chr	start	end	H3K9me3	H3K27me3	
1	0	200	0	0	
1	200	400	5	3	
1	400	600	20	4	
1	600	800	9	0	
1	800	1000	8	1	
1	1000	1200	10	1	
1		1	1	1	



ſ	chr	start	end	State
[1	0	200	1
[1	200	400	2
[1	400	600	3
[1	600	800	2
I	1	800	1000	2
[1	1000	1200	2
ſ	1			





Flexible Modeling with EpiSegMix









Flexible Emission Modeling

Table: Overview of available discrete distributions	(Johnson	et al.,	1993)
---	----------	---------	-------

Name	Parameters	Flexibility
Poisson	$\lambda \in \mathbb{R}^+$	+
Binomial	$n\in\mathbb{N},$ $oldsymbol{ ho}\in[0,1]$	+
Negative Binomial	$r \in \mathbb{R}^+, oldsymbol{p} \in [0,1]$	++
Beta Binomial	$\mathbf{n} \in \mathbb{N}, \alpha \in \mathbb{R}^+, \beta \in \mathbb{R}^+$	+ + +
Beta Negative Binomial	$\mathbf{r} \in \mathbb{R}^+, \alpha \in \mathbb{R}^+, \beta \in \mathbb{R}^+$	+ + +
Sichel	$\mu \in \mathbb{R}^+, \sigma \in \mathbb{R}^+, \textit{v} \in \mathbb{R}$	+ + +





Comparison Poisson and Beta Negative Binomial Distribution









Flexible Duration Modeling



1 Russell, Cook (1987)

Algorithmic Bioinformatics







HMM parameters

- Transition probabilities
- Emission probabilities
 - \rightarrow parameters of discrete distributions







HMM parameters

- Transition probabilities
- Emission probabilities
 - \rightarrow parameters of discrete distributions

Problem

 No mapping between observations and states available









Solution

- Baum-Welch algorithm (expectation-maximization):
 - E-step: estimate state membership probabilities
 - M-step: estimate model parameters
- Emission parameters are updated by maximum likelihood







13











chr	start	end	H3K27me3		H3K36me3
1	0	200	0		0
1	200	400	5		10
1	400	600	20		1
:	:	:	:	:	:































 $\mathbb{N} \cdots \mathbb{N}$

Results - Advantage of Flexible Modeling



Results - Increased Segment Length



Algorithmic Bioinformatics







Results - EpiSegMix is Predictive of Cell Biology



Figure: Prediction of (a) gene expression (b) ATAC-seq counts using state labels as categorical predictors.







Conclusion

- Automated chromatin state discovery enables the annotation of both coding and non-coding regions of the genome.
- The assumptions of the probabilitic model have an impact on the segmentation and biological interpretability of states.
- EpiSegMix provides a flexible framework that can be applied to read count data with different underlying distribution types.



https://gitlab.com/rahmannlab/episegmix







- Ernst Jason, Kellis Manolis. Discovery and characterization of chromatin states for systematic annotation of the human genome // Nature Biotechnology. VIII 2010. 28, 8.
- Johnson Norman L., Kotz Samuel, Kemp Adrienne W., others . Univariate discrete distributions. New York, NY [u.a.]: Wiley, 1993. 2nd ed.
- Mammana Alessandro, Chung Ho-Ryun. Chromatin segmentation based on a probabilistic model for read counts explains a large portion of the epigenome // Genome Biology. VII 2015. 16, 1. 151.
- Russell M., Cook A. Experimental evaluation of duration modelling techniques for automatic speech recognition // ICASSP '87. IEEE International Conference on Acoustics, Speech, and Signal Processing. 12. IV 1987. 2376–2379.









Algorithmic Bioinformatics



SIC Saarland Informatics Campus











A: Negative binomial emissions



Algorithmic Bioinformatics





Algorithmic Bioinformatics



SIC Saarland Informatics Campus

