DegSampler: Collapsed Gibbs sampler for detecting E3 binding sites

Osamu Maruyama Faculty of Design Kyushu University

Collaborator

Dr. Fumiko Matsuzaki Medical Institute of Bioregulation Kyushu University

- 1. Motivation
- 2. Data sets
- 3. Problem formulation
- 4. Methods
- 5. Results
- 6. Future works

Motivation

Proteasome (Wikipedia)

Unneeded and damaged proteins are degraded and recycled



E3 ubiquitin ligase selectively recognizes a substrate protein



- 1. Motivation
- 2. Data sets
- 3. Problem formulation
- 4. Methods
- 5. Results
- 6. Future works

Available data sets

E3NET:

E3-substrate interactions



ELM: DEG class motifs: E3 binding sites of protein sequences



- 1. Motivation
- 2. Data sets
- 3. Problem formulation
- 4. Methods
- 5. Results
- 6. Future works

Problem: Estimate the motif of binding site of an E3



- Performance measure =#position covered by known and predicted motifs simultaneously / W (motif width)
- 36 E3-specific sets of substrate proteins, some of which have known motifs.

- 1. Motivation
- 2. Data sets
- 3. Problem formulation

4. Methods

- 5. Results
- 6. Future works

Posterior probability distribution for motif identification

- $X = (X_1, ..., X_N)$: N given sequences.
- $\mathbf{Z} = (z_1, ..., z_N)$: starting positions of motif occurrences.
- $\theta_{0:W} = (\theta_0, \theta_1, ..., \theta_W)$: W + 1 categorical distributions of letters.
- Posterior: $p(\mathbf{Z}, \boldsymbol{\theta}_{0:W} | \mathbf{X}) \propto p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\theta}_{0:W}) \cdot p(\mathbf{Z}) \cdot p(\boldsymbol{\theta}_{0:W})$

MEME: meme-suite.org Gibbs Motif Sampler:

Likelihood function $L_{pssm}(\mathbf{Z}, \boldsymbol{\theta}_{0:W} | \mathbf{X})$



Categorical distribution

Prior distributions and collapsing

p(**Z**): Uniform distribution

 $p(\boldsymbol{\theta}_{0:W}):$

Dirichlet distribution

 \cdot conjugate prior of the categorical distribution.

• Recall posterior:

 $p(\mathbf{Z}, \dot{\boldsymbol{\theta}}_{0:W} | \mathbf{X}) \propto p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\theta}_{0:W}) \cdot p(\mathbf{Z}) \cdot p(\boldsymbol{\theta}_{0:W})$

• Collapsed posterior:

$$p(\mathbf{Z} | \mathbf{X}) = \int p(\mathbf{Z}, \boldsymbol{\theta}_{0:W} | \mathbf{X}) d\boldsymbol{\theta}_{0:W}$$

Modeling: Prior knowledge



E3 binding sites are often located in disordered regions of substrates [Guharoy *et al.*, 2016]

https://www.quantamagazine.org/howdisordered-proteins-are-upendingmolecular-biology-20170118/

Prior of starting position z_i

 $disorder_{i,j}$: disorderness of $X_{i,j}$

$$p(\mathbf{z_i}) \propto \left(\prod_{w=1}^{W} disorder_{i,z_i+w-1}\right)^{\frac{1}{W}}$$

$$p(\mathbf{Z}) \propto \prod_{i=1}^{N} p(\mathbf{z}_i)$$

Modeling: 2^{nd} likelihood based on amino acid indexing $L_{aai}(\mathbb{Z} | \mathbb{X})$

Resulting likelihood: $p(X|Z, \theta_{0:W}) \propto L_{pssm}(Z, \theta_{0:W}|X) \cdot L_{aai}(Z|X)$



- 1. Motivation
- 2. Data sets
- 3. Problem formulation
- 4. Methods
- 5. Results
- 6. Future works

Result: Effectiveness of disorder prior



- c_d : coefficient of disorder prior. (W,R)=(9,6)
 - W: motif width
 - R: #selected columns as motif part)

Result: Effectiveness of (W,R)



Example: E3 SIAH2_HUMAN

			•	Р		Α		V	•	Р
Substrate	Pos	*	*	*		*		*		*
O14974	390	А	Α	Ρ	V	Α	V	Т	Т	Р
075376	124	Α	Α	V	L	Р	L	V	Η	Р
075925	562	Α	А	Α	А	Α	А	V	S	D
P29590	437	Α	Q	Ρ	Μ	Α	V	V	Q	S
P37840	89	А	А	Α	Т	G	F	V	Κ	K
P43146+	1330	Т	Ι	Ρ	Т	Α	С	V	R	Р
Q16633+	45	Р	А	Р	Т	Α	V	V	L	Р
Q92667	284	А	А	Р	А	Р	Р	V	Α	D
Q9GZT9	77	Р	Α	Ρ	Р	Α	Α	V	Р	Р
Q9H6Z9	62	Α	G	Р	R	Α	G	V	S	K

- 1. Motivation
- 2. Data sets
- 3. Problem formulation
- 4. Methods
- 5. Results
- 6. Future works

Future works

- Improvement of the likelihood function based on amino acid indexing.
- Simulated annealing as a postprocessing.