

Discriminative Embedding of Molecular Structures for Property Prediction

Le Song

College of Computing Georgia Institute of Technology

What are structured input data?



Three key questions



Kernel methods

A modular framework to handle all kinds of complex data Theoretical property well understood



[Shaw-Taylor & Cristianini 04]

Coloumb matrix



Bag of structure (BOS) kernel

Two structured data points χ and χ' , a dictionary of substructures S

$$k(\chi, \chi') = \sum_{s \in S} \#[s \in \chi] \cdot \#[s \in \chi']$$

Eg. RNA sequences (alphabet: A, G, C, U) $S = \{ CUU, UUC, ..., CAG, AGU \}$ ┨╖┝┥╻┝┥┎┝┥ѧ┝┥╔┝╴ χ : -Մχ': -Շ–-Մ–-Մ $\phi(\chi) = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{bmatrix} \begin{array}{c} \mathsf{CUU} \\ \mathsf{UUC} \\ \vdots \\ \mathsf{CAG} \\ \mathsf{AGU} \end{array} \qquad \phi(\chi') = \begin{bmatrix} 2 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \begin{array}{c} \mathsf{CUU} \\ \mathsf{UUC} \\ \vdots \\ \mathsf{CAG} \\ \mathsf{AGU} \end{array} \qquad = \langle \phi(\chi), \phi(\chi') \rangle$

[Leslie et al. 04]



1. Record subtree



2. Re-number (or hash)



3. Record subtree again



4. Re-number (or hash) again



Final features





Big dataset, explosive feature space



Dataset	Harvard clean energy project
Size	2.3 million
Туре	Molecule
Alphabet #	6
Avg node #	28
Avg edge #	33

feature	dimension	MAE
Level 3	1.6 million	0.143
Level 6	1.3 billion	0.096

Kernels based on graphical models (GM)



Define conditional independence structure using structure of the data

Model each data point as a Markov random field

$$p(\chi|\theta) \propto \prod_{i \in \mathcal{V}} \Psi_1(X_i|\theta) \prod_{(i,j) \in \mathcal{E}} \Psi_2(X_i, X_j|\theta)$$



[Jaakkola et al. 99] [Jebara et al. 04]

Limitation of existing two-stage approach



Or learn many graphical models

Not scalable for large datasets

Feature construction not aware of supervised tasks

Represent data point as latent variable model



Model each data point as a Markov random field with latent variables

$$p(\{H_i\}, \{X_i\}|\theta) \propto \prod_{i \in \mathcal{V}} \Psi_1(H_i, X_i|\theta) \prod_{(i,j) \in \mathcal{E}} \Psi_2(H_i, H_j|\theta)$$

Recursive nonlinear feature extraction



Parameter learning

Estimate parameters W and V which minimize empirical loss

$$\min_{V,W} L(V,W) := \sum_{i=1}^{m} (y_i - V^{\top} \mu^a(W,\chi_i))^2$$

Computation	Operation	Similar to
Objective $L(V, W)$	A forward sequence of nonlinear mappings	Graphical model inference
$\frac{\partial L}{\partial W}$	Chain rule of derivatives in reverse order of the mappings	Back propagation in deep learning

Posterior embedding as features



Posterior embedding as features



Operator view of mean field inference



Hilbert (feature) space embedding of distribution



One-to-one mapping: μ_X is a sufficient statistic of p(X)

Operator $\mathcal{T}: \mathcal{P} \mapsto \mathcal{H}$

$$\boldsymbol{\mathcal{T}} \circ \boldsymbol{p}(\boldsymbol{x}) = \widetilde{\boldsymbol{\mathcal{T}}} \circ \boldsymbol{\mu}_{\boldsymbol{X}}$$

[Smola, Gretton, Song and Scholkopf. 2007] [Song et al. 10a,b, 11a,b]

Mean field update using embedding



Parameterize mean field update $\mu_i \leftarrow \widetilde{T}(W) \circ \left(x_i, \{\mu_j\}_{j \in \mathcal{N}(i)}\right)$

One can parameterize it as any flexible nonlinear model

Assume $\mu_i \in \mathcal{R}^d$, $x_i \in \mathcal{R}^n$, parametrize as one-layer neural network



Benefit of the new view: belief propagation

Approximate $p(H_i|\{x_j\}, \theta)$ as

$$q_i(H_i) = \Psi_1(H_i, x_i | \theta)$$
$$\prod_{j \in \mathcal{N}(i)} m_{ji}(H_i)$$

$$H_{1}^{h^{n}} H_{6}^{H_{6}} H_{1}^{H_{1}} H_{5}^{H_{5}} G = (\mathcal{V}, \mathcal{E})$$

$$H_{1}^{h^{n}} X_{6}^{H_{1}} X_{5}^{H_{4}} H_{4}^{H_{4}}$$

$$H_{1}^{h^{n}} X_{6}^{H_{1}} H_{2}^{H_{2}} H_{3}^{H_{4}} X_{4}^{H_{4}}$$

$$H_{1}^{h}, H_{1}^{h}|\theta) \cdot X_{2}^{h^{n}} X_{3}^{h^{n}}$$

with messages updated iteratively:

$$m_{ij}(H_{j}) \leftarrow \int_{\mathcal{H}} \Psi_{1}(H_{i}, x_{i} | \theta) \Psi_{2}(H_{i}, H_{j} | \theta) \cdot X_{2}$$

$$\prod_{\ell \in \mathcal{N}(i) \setminus j} m_{\ell i}(H_{i}) dH_{i}$$

$$q_{i}(H_{i}) \leftarrow \mathcal{T}(\theta) \circ \left(x_{i}, \left\{m_{ji}(H_{j})\right\}_{j \in \mathcal{N}(i)}\right)$$

$$m_{ij}(H_{j}) \leftarrow \mathcal{T}'(\theta) \circ \left(x_{i}, \left\{m_{\ell i}(H_{i})\right\}_{\ell \in \mathcal{N}(i) \setminus j}\right)$$

[Song et al. 11a,b] [Song et al. 10a,b]

Belief propagation using embedding





Experiment II: drug prediction

dataset	MUTAG	NCI1	NCI109	ENZYMES	D&D
type	Chem	Chem	Chem	Protein	Protein
#sample	188	4110	4127	600	1178
#class	2	2	2	6	2
Avg $ V $	18	30	30	33	284
Avg $ E $	20	32	32	62	715
Alphabet size	7	37	38	3	82
Real world problem	Mutagenic / non- mutagenic compounds for <i>Salmonella</i> <i>Typhimurium</i>	Active / inactive compounds in an anti- cancer screen	Active / inactive compounds in an anti- cancer screen	Enzymes / non-enzymes	Enzymes / non-enzymes



Experiment III: graph classification







Predicting efficiency of solar panel materials



Harvard clean energy project
2.3 million
Molecule
6
28
33

Power Conversion Efficiency (PCE) (0 -12 %)

Characteristics of the data



Clusters of data points

Unbalanced output value Resample data during training to make it balance

[Dai et al. NIPS 14] [Xie et al. NIPS 15]



Small model but accurate prediction

10% data for testing

	Test MAE	Test RMSE	# parameters
Mean predictor	1.986	2.406	1
WL level-3	0.143	0.204	1.6 m
WL level-6	0.096	0.137	1378 m
S2V-MF	0.091	0.125	0.1 m
S2V-BP	0.085	0.117	0.1 m

We get ~4% relative error with 10,000 times smaller model!



Effects of inference iterations on errors





Errors across different output range

Prediction quality



Conclusion and future work



Discriminative/generative models

Codes: https://github.com/Hanjun-Dai/graphnn