Knowledge Distillation via Constrained Variational Inference Bridging Discriminative Power and Interpretability

Ardavan Saeedi¹, Yuria Utsumi², Li Sun³, Kayhan Batmanghelich³, Li-wei H. Lehman²

¹Hyperfine ²MIT ³University of Pittsburgh

AAAI 2022

《曰》 《聞》 《臣》 《臣》 三臣

Roadmap

Motivation & Problem Statement

- 2 Methodology
- Experimental Validation
- Technical Contributions
- 5 Discussion & Future Work

6 Conclusion

The Interpretability-Performance Trade-off

Discriminative Models

- High predictive performance
- Black-box nature
- Limited interpretability
- Dense, complex representations

Probabilistic Graphical Models

- Interpretable structure
- Principled uncertainty quantification
- Poor predictive performance
- Two-stage feature extraction

Core Challenge

Can we distill the discriminative power of neural networks into interpretable graphical models while preserving their generative capabilities?

Teacher Model

(Discriminative)

Knowledge Distillation

Student Model (Graphical Model)

Desired Properties

- Interpretability
- Generative modeling
- Predictive power
- Uncertainty quantification

Framework Overview



Key Innovation

Similarity-preserving knowledge distillation constraint integrated into black-box variational inference

Saeedi et al.	

KD via Constrained VI

AAAI 2022

5/18

Constrained Variational Objective

 \mathcal{L}

$$\begin{aligned} (\phi) &= \mathbb{E}_{q_{\phi}}[\log p(\boldsymbol{x}, \boldsymbol{z}, \boldsymbol{\theta}) - \log q_{\phi}(\boldsymbol{z}, \boldsymbol{\theta})] \\ &+ \gamma_{\eta} \cdot \frac{1}{N^{2}} \| \bar{\boldsymbol{F}}^{s} - \bar{\boldsymbol{F}}^{t} \|_{F}^{2} \end{aligned} \tag{2}$$

Similarity Matrices		Variational Approximation
$egin{array}{ll} ilde{m{ extbf{F}}}^s = m{m{ extbf{F}}}^s \cdot (m{m{ extbf{F}}}^s)^{ op} \ ilde{m{ extbf{F}}}^t = m{m{ extbf{F}}}^t \cdot (m{m{ extbf{F}}}^t)^{ op} \end{array}$	(3) (4)	$egin{aligned} q(m{ heta}) &= \mathcal{N}(m{\mu}, ext{diag}(ext{exp}(m{\omega})^2)) \ (5) \ q(m{z} m{x}) &= \mathcal{N}(m{\mu}_{\phi_z}(m{x}), ext{diag}(ext{exp}(m{\omega}_{\phi_z}(m{x}))^2)) \end{aligned}$
$\bar{F}^{s}, \bar{F}^{t}: \ell_{2}$ -normalized versions		(6) <ロ> < () < () < () < () < () < () < () < (
Speedi et al	KD via Car	

Black-Box Variational Inference Integration

ADVI for Global Variable	es $ heta$	Advantages	
 Transform to unconstrai 	ned space: $oldsymbol{\Theta} = \mathcal{T}(oldsymbol{ heta})$	 Model-agnostic 	
 Factorized Gaussian app 	 Scalable inference 		
 Automatic differentiatio 	 Easy implementation 		
		 Broad applicability 	
AEVB for Local Variable	es z		
• Recognition network: <i>x</i>	$\mapsto q(\boldsymbol{z} \boldsymbol{x})$	Technical Detail	
 Amortized inference 		Combined ADVI + AEVB enables	
 Reparameterization trick 	<	handling both global and local	
		latent variables in a unified	
		framework	
		framework	

Application 1: COPD Disease Subtyping

Dataset & Task

- COPDGene: 7,292 subjects
- CT lung images
- Predict clinical severity (FEV1pp, FEV1/FVC, etc.)
- Identify interpretable disease subtypes

Results: Coefficient of Determination (R^2)

Model	FEV1pp	FEV1/FVC	FVCpp	Distance
G-LDA	0.16	0.30	0.03	0.05
Sup. G-LDA	0.29	0.49	-0.47	0.06
KD-LDA	0.49	0.61	0.15	0.14
Teacher	0.65	0.70	0.28	0.16

Models

- Teacher: Subject2Vec (deep sets)
- Student: Gaussian-LDA variant
- Features: Topic proportions π_s

Key Finding

Substantial improvement in predictive performance while maintaining generative capabilities (ELBO preserved)

COPD Results: Clinical Interpretability



Clinical Insights

- Subtype 1 proportion increases with disease severity (GOLD 0 ${
 ightarrow}4)$
- Subtype 2 decreases with severity
- Learned subtypes correlate with anatomical regions
- Interpretable disease progression patterns

Application 2: Sepsis Dynamics Modeling

Dataset & Task	Results
• MIMIC-III: 11,648 sepsis patients	Model AUROC ELBO
 29 physiological variables 	G-ARHMM 0.56 -55.24
 Predict in-hospital mortality 	Sup. G-ARHMM 0.56 -55.24
 Model clinical state transitions 	KD-ARHMM 0.65 -3.94
	Teacher (LSTM) 0.71 N/A
Models	
• Teacher: LSTM	
• Student: AR-HMM	Remarkable Result
• Features: State marginals	Improved both predictive performance AND generative modeling (ELBO

improvement)

Sepsis Results: Clinical State Analysis



Clinical Validation

- KD-ARHMM states better correlate with SOFA scores
- More clinically meaningful state transitions
- Improved end-organ dysfunction tracking

Similarity-Preserving Knowledge Distillation

- Matches pairwise similarity matrices instead of direct feature matching
- Handles heterogeneous feature spaces (neural net i-i graphical model)

2 Unified ADVI + AEVB Framework

- Enables black-box inference for both global and local variables
- Model-agnostic approach

Onstrained Variational Inference

- Integrates knowledge distillation as Lagrangian constraint
- Balances generative modeling and discriminative power

Broader Impact

Framework applicable to any probabilistic graphical model with local latent variables

Algorithmic Complexity & Scalability

Computational Considerations

- Similarity Matrix: $O(N^2)$ space and computation per batch
- Mini-batch Strategy: Reduces to $O(B^2)$ where $B \ll N$
- Gradient Computation: Automatic differentiation through constraint

Hyperparameter Sensitivity

- Regularization Weight γ_{η} : Controls knowledge distillation strength
- **Batch Size** *B*: Affects similarity matrix approximation quality
- Architecture: Recognition network capacity impacts local variable inference

Convergence Properties

Inherits convergence guarantees from ADVI with additional constraint satisfaction

Saeedi et al

Limitations & Future Directions

Current Limitations

- BBVI weaknesses:
 - Posterior variance underestimation
 - Initialization sensitivity
 - Amortization gap
- Quadratic scaling in similarity computation
- Hyperparameter tuning complexity

Future Research Directions

- Discrete latent variable support
- Alternative similarity measures
- Multi-teacher distillation
- Theoretical analysis of constraint effects
- Hierarchical model extensions

Open Questions

- How does constraint strength affect posterior approximation quality?
- Can we develop adaptive weighting schemes for γ_{η} ?
- What are the theoretical guarantees for knowledge preservation?

Broader Implications

Methodological Impact

- Interpretable AI: Bridge between black-box and interpretable models
- Transfer Learning: Novel approach to cross-model knowledge transfer
- Probabilistic ML: Enhanced inference for structured models

Application Domains

- Healthcare: Disease subtyping, patient monitoring, clinical decision support
- Natural Language: Topic modeling with neural guidance
- Computer Vision: Interpretable scene understanding
- Time Series: Dynamics modeling with neural priors

Paradigm Shift

From choosing between interpretability and performance to achieving both through principled knowledge distillation

Saeedi et a

Summary

Problem Addressed

How to distill discriminative knowledge into interpretable probabilistic models without sacrificing generative capabilities

Solution Approach

Similarity-preserving knowledge distillation integrated into black-box variational inference via constrained optimization

Key Results

- Substantial predictive performance improvements (COPD: 0.16→0.49 R², Sepsis: 0.56→0.65 AUROC)
- Preserved or improved generative modeling capabilities
- Clinically meaningful and interpretable learned representations
- General framework applicable to diverse graphical models

Take-Home Message

Principled knowledge distillation can bridge the interpretability-performance gap in probabilistic machine learning

AAAI 2022

Questions & Discussion

Thank you for your attention

av.saeedi@gmail.com

<□▶ <⊡▶ <글▶ <글▶