

DSC190: Machine Learning with Few Labels

Knowledge driven learning

Zhiting Hu

Lecture 12, November 2, 2021

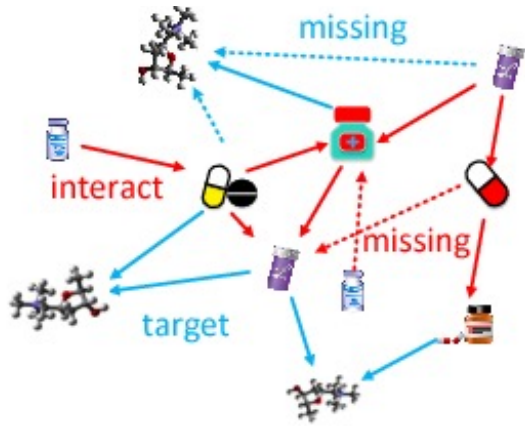
UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Outline

- Knowledge-driven learning (65mins)
- 1 Paper presentations (15 mins)
 - **Kejin Wu:** Variational Inference with Normalizing Flows

Structured knowledge



Knowledge graphs

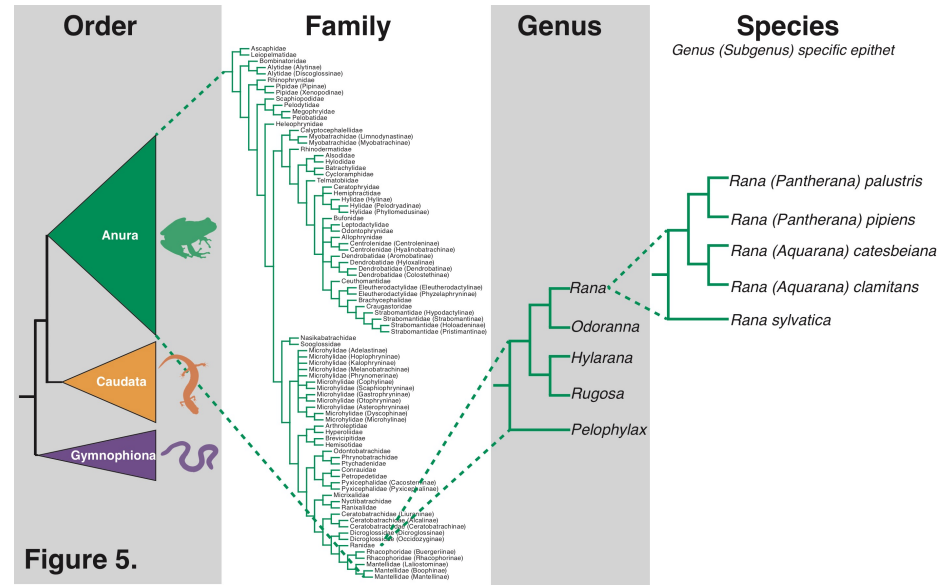


Figure 5.

Taxonomy



WIKIPEDIA
The Free Encyclopedia

- Main page
- Contents
- Current events
- Random article
- About Wikipedia
- Contact us
- Donate

Article [Talk](#)

Machine learning

From Wikipedia, the free encyclopedia

For the journal, see [Machine Learning \(journal\)](#).

"Statistical learning" redirects here. For statistical learning in linguistics, see [Statistical learning theory](#).

Machine learning (ML) is the study of computer algorithms that can automatically learn and improve from experience without being explicitly programmed. Machine learning algorithms build a model based on training data, which the model can then use to make predictions or decisions without being explicitly programmed to do so. Machine learning is a subset of artificial intelligence. Machine learning algorithms are used in a wide variety of applications, such as in medicine, email filtering, speech recognition, and recommendation systems.

Encyclopaedia

Rules:

- “Every part of speech sequence should have a verb”
- “Type-2 diabetes is 90% more common than type-1”

Machine Learning, esp., deep learning

- Heavily rely on massive labeled data
- Uninterpretable
- Hard to encode human intention and domain knowledge

How Humans Learn

- Learn from **concrete** examples (as deep neural networks do)
- Learn from **abstract** knowledge (definitions, logic rules, etc) [Minsky 1980; Lake et al., 2015]

How Humans Learn

- Learn from **concrete** examples (as deep neural networks do)
- Learn from **abstract** knowledge (definitions, logic rules, etc) [Minsky 1980; Lake et al., 2015]

Past tense of verb

Examples:

add → added
accept → accepted
ignore → ignored
end → ended
block → blocked
love → loved

...

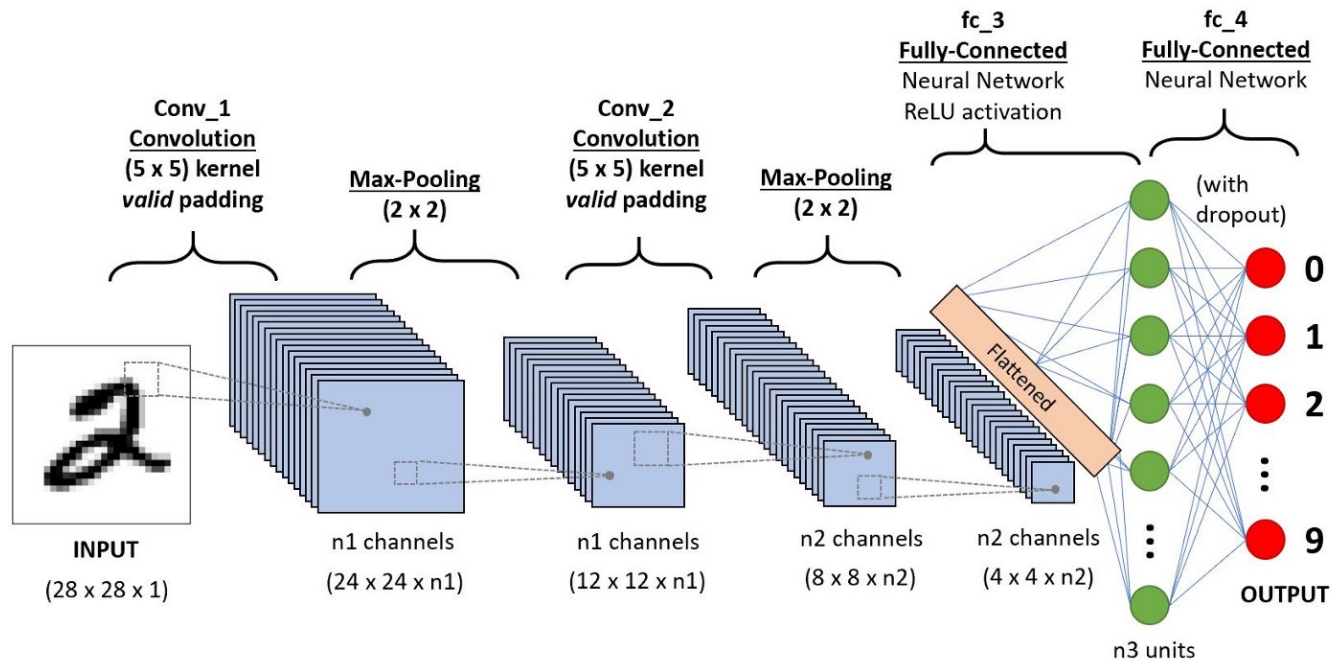
V.S.

Rule:

regular verbs -d/-ed

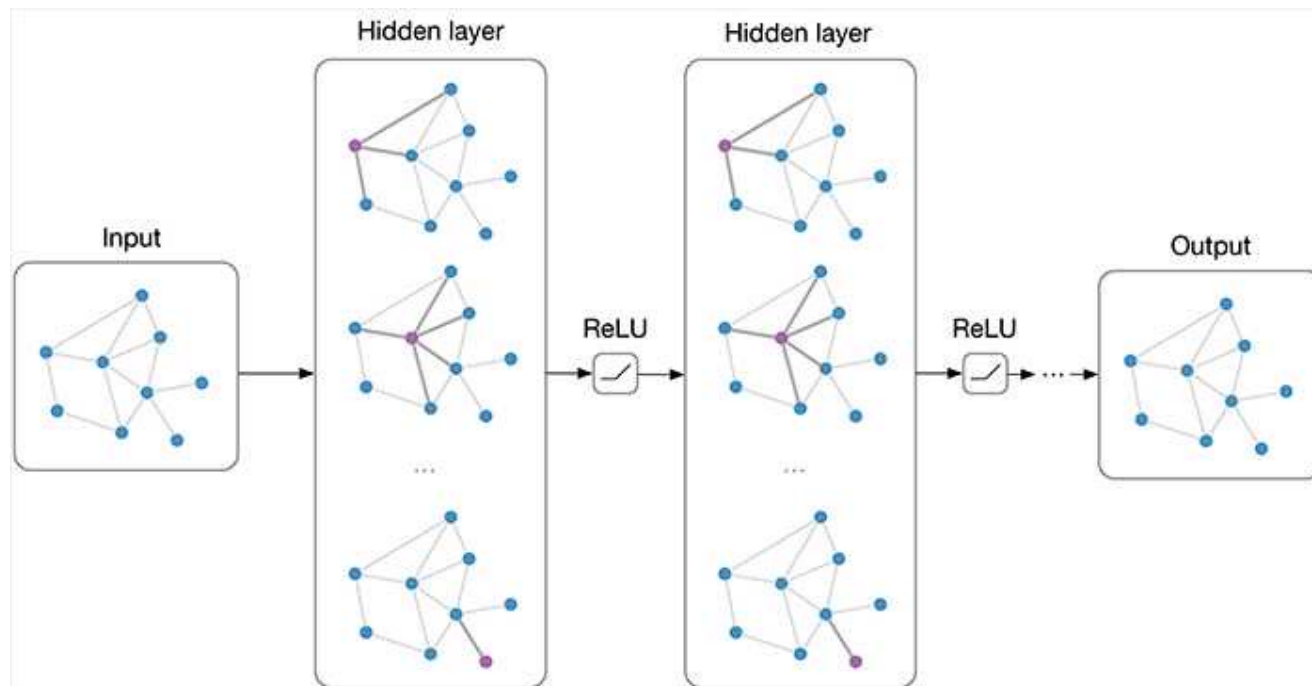
Integrating structured knowledge with machine learning

- Two general ways:
 - Bake structured knowledge into specifically-designed neural architecture (“inductive bias”)
 - E.g., Convolutional networks (ConvNets): translation-invariance of image



Integrating structured knowledge with machine learning

- Two general ways:
 - **Bake structured knowledge into specifically-designed neural architecture** (“inductive bias”)
 - E.g., Convolutional networks (ConvNets): translation-invariance of image
 - E.g., graph neural networks



Integrating structured knowledge with machine learning

- Two general ways:
 - Bake structured knowledge into specifically-designed neural architecture (“inductive bias”)
 - E.g., Convolutional networks (ConvNets): translation-invariance of image
 - E.g., graph neural networks

The Bitter Lesson

Rich Sutton

March 13, 2019

The biggest lesson that can be read from 70 years of AI research is that generalization is the key. The ultimate reason for this is Moore's law, or rather its generalization of computing power as if the computation available to the agent were constant (in which case level of performance grows slightly longer time than a typical research project, massively more computation, researchers seek to leverage their human knowledge of the domain, but this runs counter to each other, but in practice they tend to. Time spent on one is time spent on the other. And the human-knowledge approach tends to complicate methods of computation. There were many examples of AI researchers' belated learning.

<http://www.incompleteideas.net/Incldeas/BitterLesson.html>



RODNEY BROOKS *Robots, AI, and other stuff* **BLOG** MIT ROBUST.AI

POST: A BETTER LESSON

MARCH 19, 2019 — REVIEWS

A Better Lesson  
rodneybrooks.com/a-better-lesson/

Just last week Rich Sutton published a very short blog post titled [The Bitter Lesson](#). I'm going to try to keep this review shorter than his post. Sutton is well known for his long and sustained contributions to reinforcement learning.

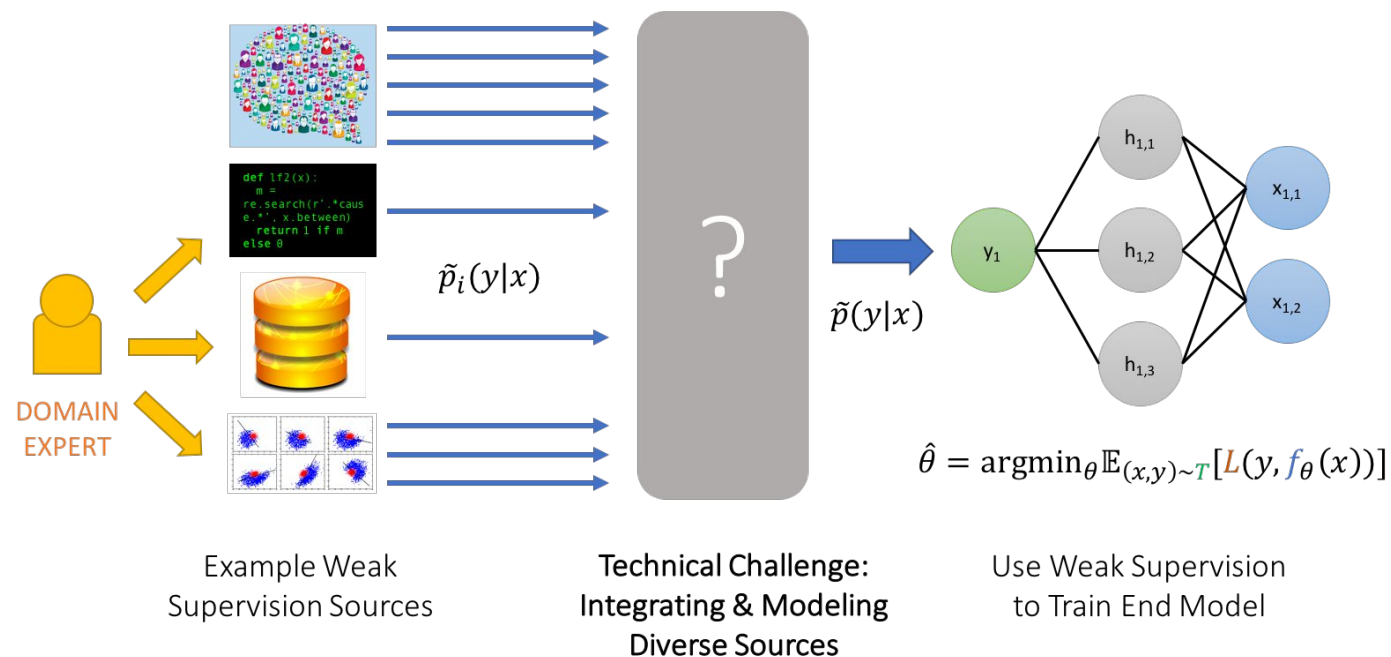
<https://rodneybrooks.com/a-better-lesson/>

Integrating structured knowledge with machine learning

- Two general ways:
 - **Bake structured knowledge into specifically-designed neural architecture** (“inductive bias”)
 - E.g., Convolutional networks (ConvNets): translation-invariance of image
 - E.g., graph neural networks
 - **Integrating knowledge through learning:**
 - Loss and/or constraints defined by the structured knowledge
 - Model-agnostic
 - E.g., Weak supervision
 - E.g., Posterior regularization
 - E.g., Integer linear programming (ILP)

Recap: Weakly supervised learning

- Converts knowledge into weak-supervision labels
- Learn with supervised learning methods



Source: A. Ratner et. al <https://dawn.cs.stanford.edu/2017/07/16/weak-supervision/>

Recap: Posterior regularization

- Knowledge as constraints

$$\min_{q, \xi} \quad -H(q(\mathbf{z})) - \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{x}^* | \mathbf{z}) \pi(\mathbf{z})] + \sum_i \xi_i$$

$$s.t. \quad \mathbb{E}_q [T(\mathbf{x}^*, \mathbf{z})] \leq \xi$$

$$\xi \geq 0,$$

Knowledge as constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
 - Conditional model, $p_{\theta}(\mathbf{x} | \textit{inputs})$
 - Generative model, e.g., \mathbf{x} is an image
 - Discriminative model, e.g., x is a sentence label

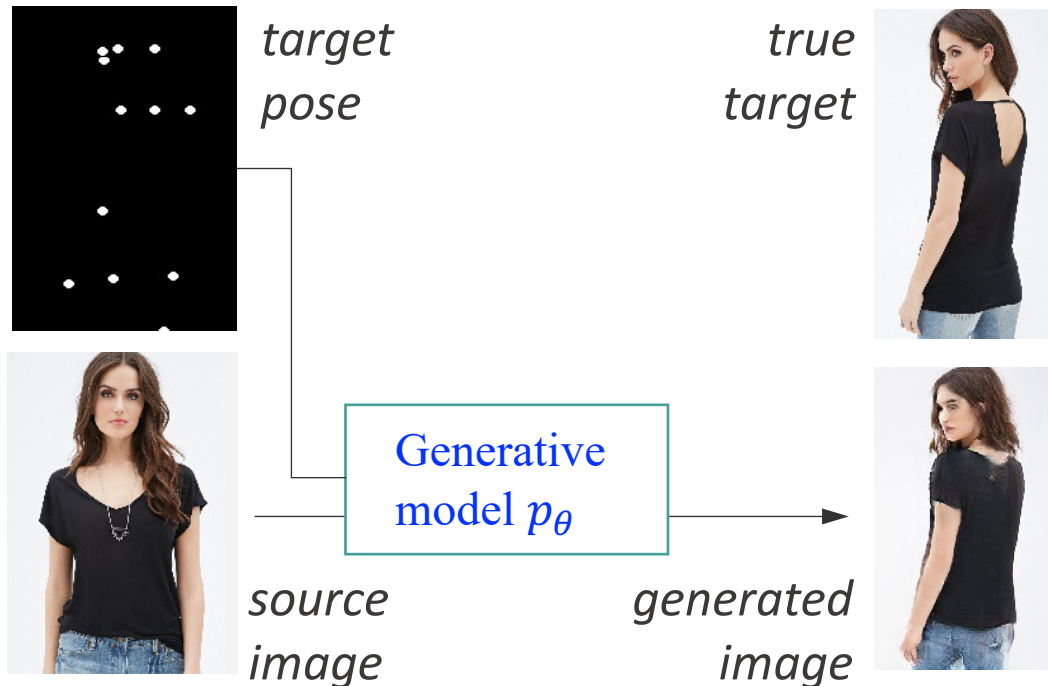
Knowledge as constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f(\mathbf{x}) \in \mathbb{R}$
 - Higher f value, better \mathbf{x} w.r.t. the knowledge

Knowledge as constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f(\mathbf{x}) \in \mathbb{R}$
 - Higher f value, better \mathbf{x} w.r.t. the knowledge

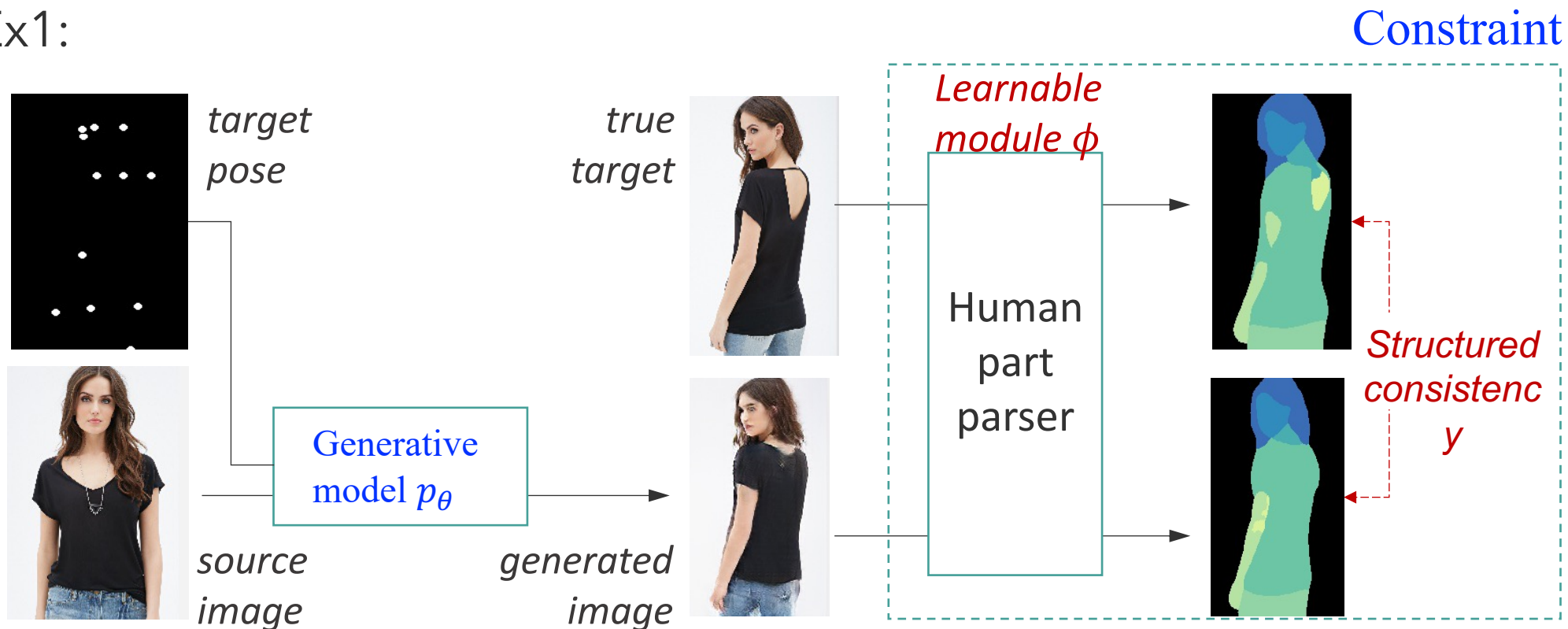
Ex1:



Knowledge as constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f(\mathbf{x}) \in \mathbb{R}$
 - Higher f value, better \mathbf{x} w.r.t. the knowledge

Ex1:



Knowledge as constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f(\mathbf{x}) \in \mathbb{R}$
 - Higher f value, better \mathbf{x} w.r.t. the knowledge

Ex2:

- Sentiment classification
 - “This was a terrific movie, but the director could have done better”
- Logical Rules:
 - Sentence S with structure A -but- B \Rightarrow sentiment of B dominates

Knowledge as constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f(\mathbf{x}) \in \mathbb{R}$
 - Higher f value, better \mathbf{x} w.r.t. the knowledge
- One way to impose the constraint is to maximize: $\mathbb{E}_{p_{\theta}}[f(\mathbf{x})]$

Knowledge as constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f(\mathbf{x}) \in \mathbb{R}$
 - Higher f value, better \mathbf{x} w.r.t. the knowledge
- One way to impose the constraint is to maximize: $\mathbb{E}_{p_{\theta}}[f(\mathbf{x})]$
- Objective:

$$\min_{\theta} \mathcal{L}(\theta) - \alpha \mathbb{E}_{p_{\theta}}[f(\mathbf{x})]$$

Regular objective (e.g.,
cross-entropy loss, etc.)

Regularization:
imposing constraints
(difficult to compute)

Knowledge as constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f(\mathbf{x}) \in \mathbb{R}$

$$\min_{\theta} \mathcal{L}(\theta) - \alpha \mathbb{E}_{p_{\theta}}[f(\mathbf{x})]$$

Knowledge as constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f(\mathbf{x}) \in \mathbb{R}$

$$\min_{\theta} \mathcal{L}(\theta) - \alpha \mathbb{E}_{p_{\theta}}[f(\mathbf{x})]$$

$$\mathcal{L}(\theta, q) = \text{KL}(q(\mathbf{x}) || p_{\theta}(\mathbf{x})) - \lambda \mathbb{E}_q[f(\mathbf{x})]$$

Knowledge as constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f(\mathbf{x}) \in \mathbb{R}$

$$\min_{\theta} \mathcal{L}(\theta) - \alpha \mathbb{E}_{p_{\theta}}[f(\mathbf{x})]$$

$$\mathcal{L}(\theta, q) = \text{KL}(q(\mathbf{x}) || p_{\theta}(\mathbf{x})) - \lambda \mathbb{E}_q[f(\mathbf{x})]$$

Posterior Regularization
[Ganchev et al., 2010]

- Introduce variational distribution q
 - Impose constraint on q
 - Encourage q to stay close to p

Knowledge as constraints

- Consider a statistical model $\mathbf{x} \sim p_{\theta}(\mathbf{x})$
- Consider a constraint function $f(\mathbf{x}) \in \mathbb{R}$

$$\min_{\theta} \mathcal{L}(\theta) - \alpha \mathbb{E}_{p_{\theta}}[f(\mathbf{x})]$$

$$\mathcal{L}(\theta, q) = \text{KL}(q(\mathbf{x}) || p_{\theta}(\mathbf{x})) - \lambda \mathbb{E}_q[f(\mathbf{x})]$$

Posterior Regularization
[Ganchev et al., 2010]

- Introduce variational distribution q
 - Impose constraint on q
 - Encourage q to stay close to p
- Objective

$$\min_{\theta, q} \mathcal{L}(\theta) + \alpha \mathcal{L}(\theta, q)$$

Knowledge as constraints

$$\min_{\theta, q} \mathcal{L}(\boldsymbol{\theta}) + \alpha \mathcal{L}(\boldsymbol{\theta}, q)$$

$$\mathcal{L}(\boldsymbol{\theta}, q) = \text{KL}(q(\mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{x})) - \lambda \mathbb{E}_q[f(\mathbf{x})]$$

- EM-style procedure for solving the problem

Knowledge as constraints

$$\min_{\theta, q} \mathcal{L}(\theta) + \alpha \mathcal{L}(\theta, q)$$

$$\mathcal{L}(\theta, q) = \text{KL}(q(\mathbf{x}) || p_{\theta}(\mathbf{x})) - \lambda \mathbb{E}_q[f(\mathbf{x})]$$

- EM-style procedure for solving the problem
 - E-step

$$q^*(\mathbf{x}) = p_{\theta}(\mathbf{x}) \exp\{ \lambda f(\mathbf{x}) \} / Z$$

Knowledge as constraints

$$\min_{\theta, q} \mathcal{L}(\theta) + \alpha \mathcal{L}(\theta, q)$$

$$\mathcal{L}(\theta, q) = \text{KL}(q(\mathbf{x}) || p_{\theta}(\mathbf{x})) - \lambda \mathbb{E}_q[f(\mathbf{x})]$$

- EM-style procedure for solving the problem
 - E-step

$$q^*(\mathbf{x}) = p_{\theta}(\mathbf{x}) \exp\{\lambda f(\mathbf{x})\} / Z$$



Higher value -- higher probability under q – “soft constraint”

Knowledge as constraints

$$\min_{\theta, q} \mathcal{L}(\theta) + \alpha \mathcal{L}(\theta, q)$$

$$\mathcal{L}(\theta, q) = \text{KL}(q(\mathbf{x}) || p_{\theta}(\mathbf{x})) - \lambda \mathbb{E}_q[f(\mathbf{x})]$$

- EM-style procedure for solving the problem
 - E-step

$$q^*(\mathbf{x}) = p_{\theta}(\mathbf{x}) \exp\{\lambda f(\mathbf{x})\} / Z$$

- M-step

$$\min_{\theta} \mathcal{L}(\theta) - \mathbb{E}_{q^*}[\log p_{\theta}(\mathbf{x})]$$

Higher value -- higher probability under q -- “soft constraint”



Logical Rule Constraints

- Consider a supervised learning: $p_{\theta}(y|\mathbf{x})$
- Input-Target space (X, Y)

Logical Rule Constraints

- Consider a supervised learning: $p_{\theta}(y|\mathbf{x})$
- Input-Target space (X, Y)
- First-order logic rules: (r, λ)

Logical Rule Constraints

- Consider a supervised learning: $p_{\theta}(y|\mathbf{x})$
- Input-Target space (X, Y)
- First-order logic rules: (r, λ)
 - $r(X, Y) \in [0, 1]$, could be soft
 - λ is the confidence level of the rule

Logical Rule Constraints

- Consider a supervised learning: $p_{\theta}(y|\mathbf{x})$
- Input-Target space (X, Y)
- First-order logic rules: (r, λ)
 - $r(X, Y) \in [0, 1]$, could be soft
 - λ is the confidence level of the rule

- Given l rules:

- E-step:
$$q^*(y|\mathbf{x}) = p_{\theta}(y|\mathbf{x}) \exp \left\{ \sum_l \lambda_l r_l(y, \mathbf{x}) \right\} / Z$$

- M-step:

$$\min_{\theta} \mathcal{L}(\theta) - \mathbb{E}_{q^*} [\log p_{\theta}(y|\mathbf{x})]$$

Logical Rule Constraints

- Consider a supervised learning: $p_{\theta}(y|\mathbf{x})$
- Input-Target space (X, Y)
- First-order logic rules: (r, λ)
 - $r(X, Y) \in [0, 1]$, could be soft
 - λ is the confidence level of the rule

- Given l rules:

- E-step:
$$q^*(y|\mathbf{x}) = p_{\theta}(y|\mathbf{x}) \exp \left\{ \sum_l \lambda_l r_l(y, \mathbf{x}) \right\} / Z$$

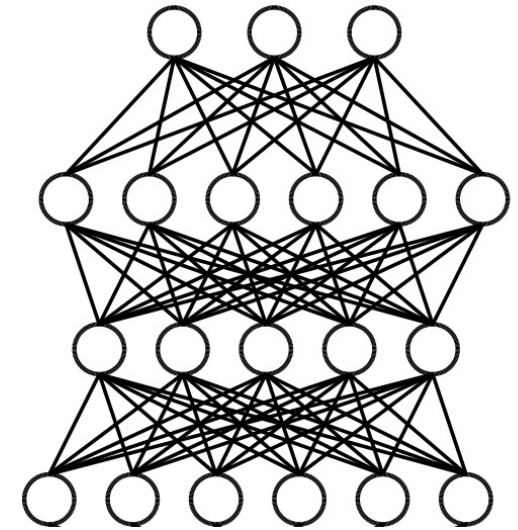
- M-step:

$$\min_{\theta} \mathcal{L}(\theta) - \mathbb{E}_{q^*} [\log p_{\theta}(y|\mathbf{x})]$$

Knowledge distillation [Hinton et al., 2015; Bucilu et al., 2006]

Knowledge Distillation

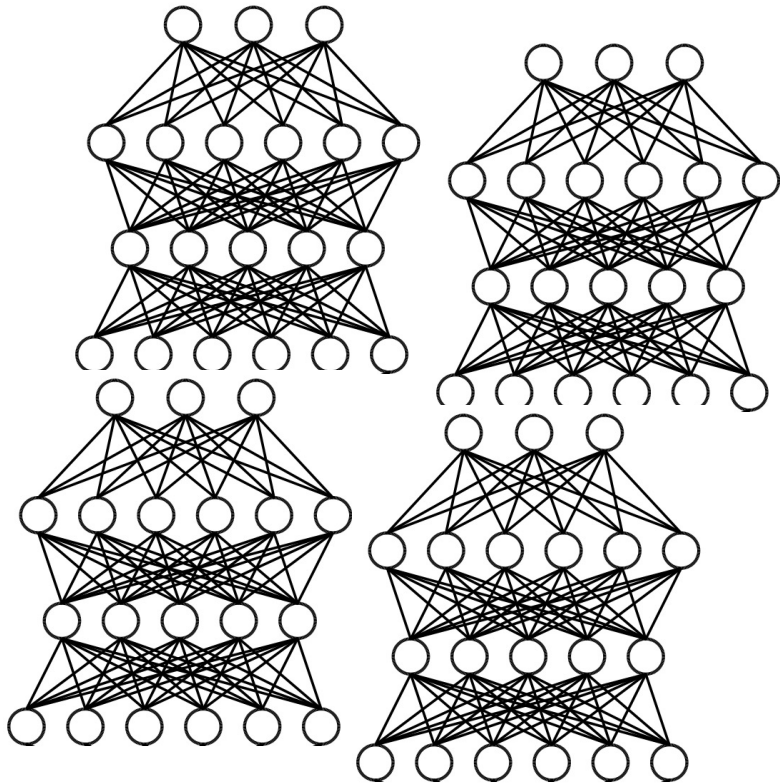
$$p_{\theta}(y|\mathbf{x})$$



Student

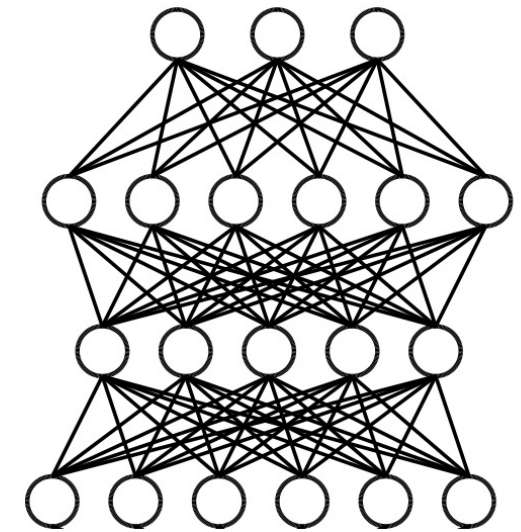
Knowledge Distillation

$$q(y|\mathbf{x})$$



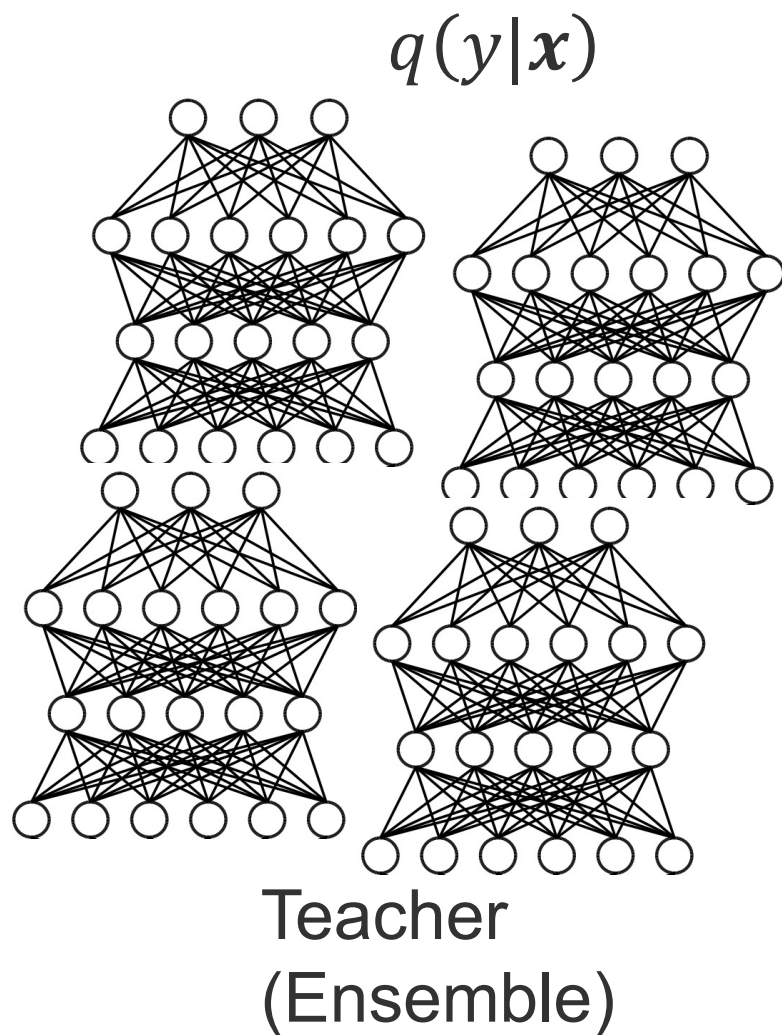
Teacher
(Ensemble)

$$p_{\theta}(y|\mathbf{x})$$

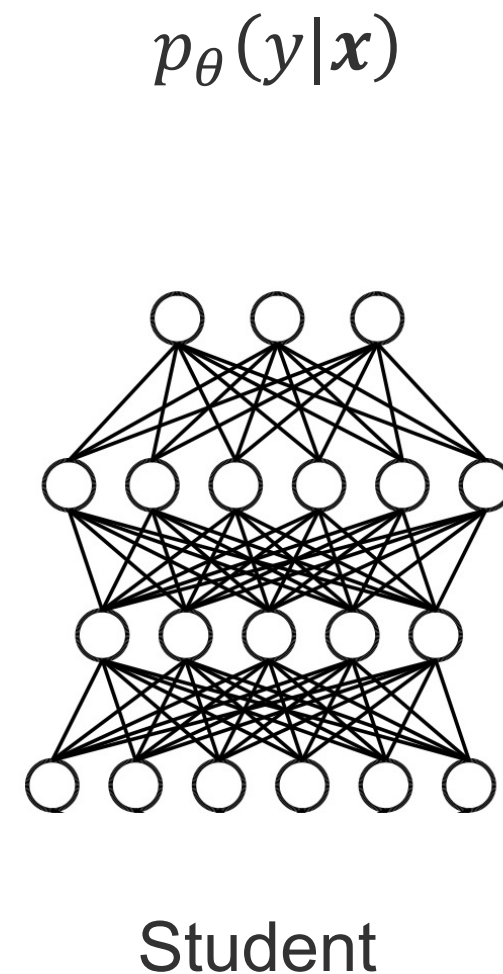
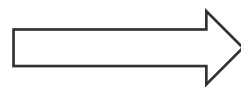


Student

Knowledge Distillation



Match soft predictions of the teacher network and student network



Rule Knowledge Distillation

$$\min_{\theta} \mathcal{L}(\theta) - \mathbb{E}_{q^*} [\log p_{\theta}(y|\mathbf{x})]$$

- Neural network $p_{\theta}(y|\mathbf{x})$
- Train to imitate the outputs of the rule-regularized teacher network

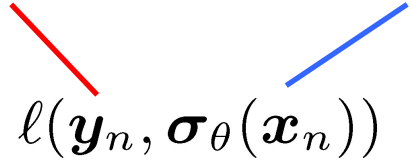
Rule Knowledge Distillation

$$\min_{\theta} \mathcal{L}(\theta) - \mathbb{E}_{q^*} [\log p_{\theta}(y|\mathbf{x})]$$

- Neural network $p_{\theta}(y|\mathbf{x})$
- Train to imitate the outputs of the rule-regularized teacher network
- At iteration t :

$$\theta^{(t+1)} = \arg \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{y}_n, \sigma_{\theta}(\mathbf{x}_n))$$

true hard label soft prediction of $p_{\theta}(y|x)$



Rule Knowledge Distillation

$$\min_{\theta} \mathcal{L}(\theta) - \mathbb{E}_{q^*} [\log p_{\theta}(y|\mathbf{x})]$$

- Neural network $p_{\theta}(y|\mathbf{x})$
- Train to imitate the outputs of the rule-regularized teacher network
- At iteration t :

$$\theta^{(t+1)} = \arg \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \left(\begin{array}{l} \text{true hard label} \\ \ell(\mathbf{y}_n, \sigma_{\theta}(\mathbf{x}_n)) \\ \text{soft prediction of } p_{\theta}(y|\mathbf{x}) \end{array} \right) + \left(\begin{array}{l} \ell(\mathbf{s}_n^{(t)}, \sigma_{\theta}(\mathbf{x}_n)), \\ \text{soft prediction of the} \\ \text{teacher network} \end{array} \right)$$

$$q^*(y|\mathbf{x}) = p_{\theta}(y|\mathbf{x}) \exp \left\{ \sum_l \lambda_l r_l(y, \mathbf{x}) \right\} / Z$$

Rule Knowledge Distillation

$$\min_{\theta} \mathcal{L}(\theta) - \mathbb{E}_{q^*} [\log p_{\theta}(y|\mathbf{x})]$$

- Neural network $p_{\theta}(y|\mathbf{x})$
- Train to imitate the outputs of the rule-regularized teacher network
- At iteration t :

$$\theta^{(t+1)} = \arg \min_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N (1 - \pi) \ell(\mathbf{y}_n, \sigma_{\theta}(\mathbf{x}_n)) + \pi \ell(\mathbf{s}_n^{(t)}, \sigma_{\theta}(\mathbf{x}_n)),$$

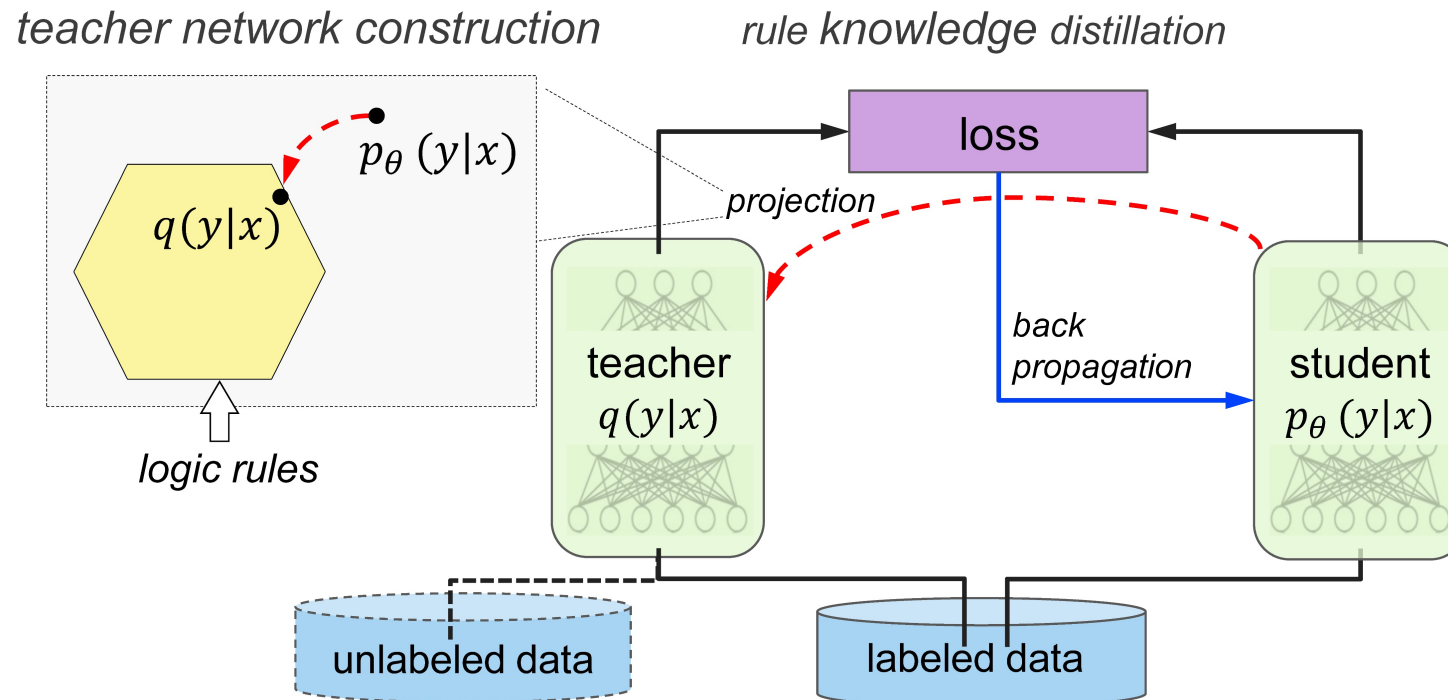
Annotations for the equation above:

- Red line: true hard label
- Blue line: soft prediction of $p_{\theta}(y|\mathbf{x})$
- Blue line: balancing parameter
- Green line: soft prediction of the teacher network

$$q^*(y|\mathbf{x}) = p_{\theta}(y|\mathbf{x}) \exp \left\{ \sum_l \lambda_l r_l(y, \mathbf{x}) \right\} / Z$$

Rule Knowledge Distillation

- Neural network $p_{\theta}(y|x)$
- At each iteration
 - Construct a teacher network with “soft constraint”
 - Train DNN to emulate the teacher network



Learning Rules / Constraints

$$q^*(y|\mathbf{x}) = p_{\theta}(y|\mathbf{x}) \exp \left\{ \sum_l \lambda_l r_l(y, \mathbf{x}) \right\} / Z$$

- Learn the confidence value λ_l for each logical rule [Hu et al., 2016b]

Learning Rules / Constraints

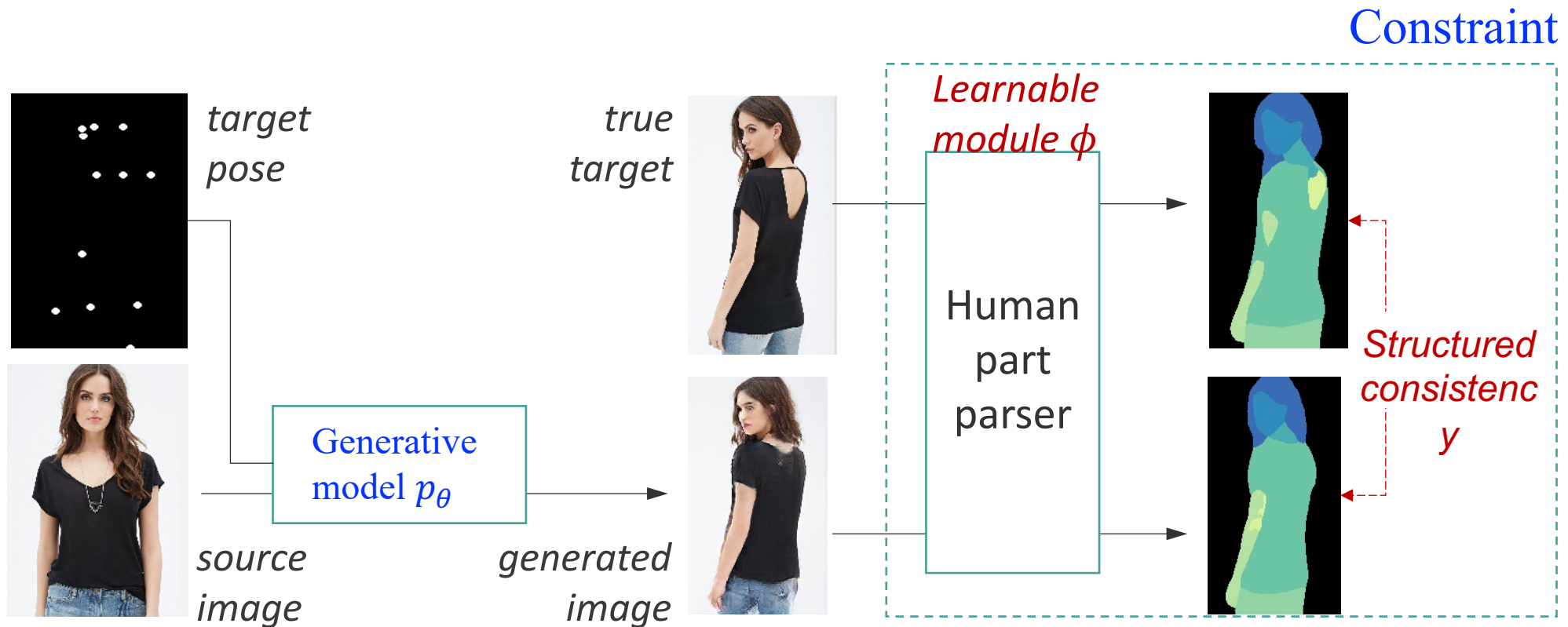
$$q^*(y|\mathbf{x}) = p_\theta(y|\mathbf{x}) \exp \left\{ \sum_l \lambda_l r_l(y, \mathbf{x}) \right\} / Z$$

- Learn the confidence value λ_l for each logical rule [Hu et al., 2016b]

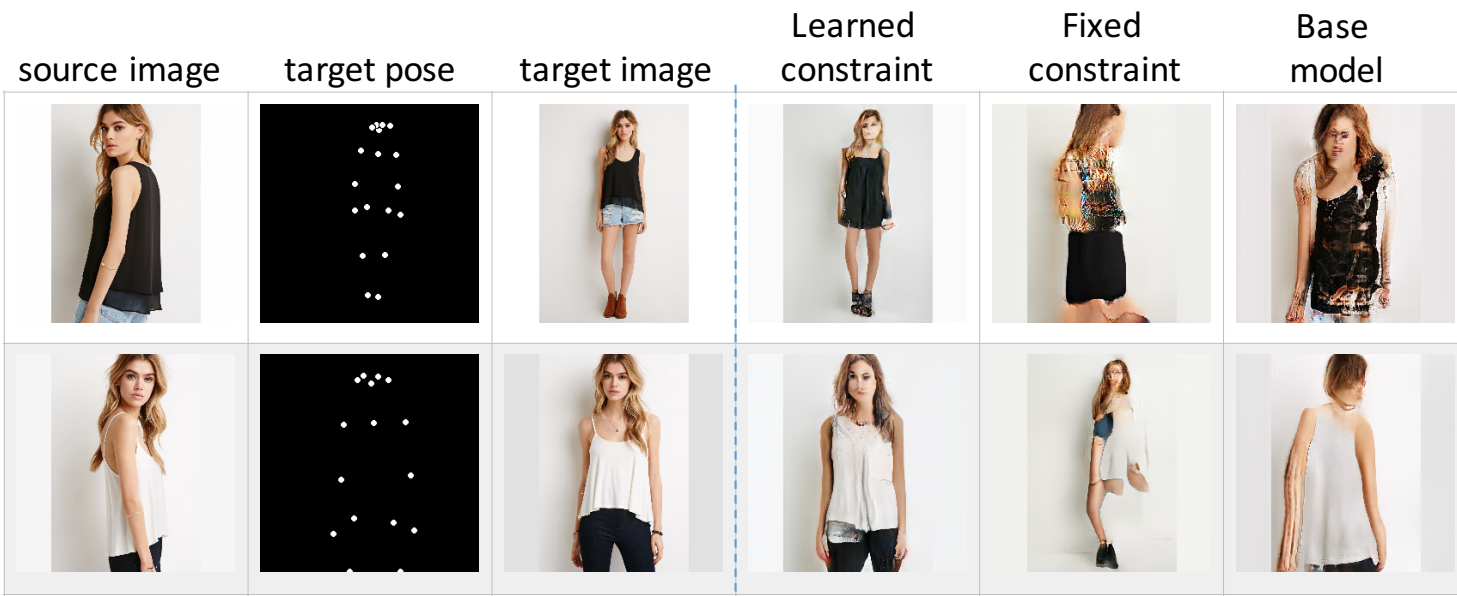
$$q^*(\mathbf{x}) = p_\theta(\mathbf{x}) \exp \{ \lambda f_\phi(\mathbf{x}) \} / Z$$

- More generally, optimize parameters of the constraint $f_\phi(\mathbf{x})$ [Hu et al., 2018]
 - Treat $f_\phi(\mathbf{x})$ as an extrinsic reward function
 - Use MaxEnt Inverse Reinforcement Learning to learn the “reward”

Pose-conditional Human Image Generation



Pose-conditional Human Image Generation



Samples generated by the models

	Method	SSIM	Human
1	Ma et al. [38]	0.614	—
2	Pumarola et al. [44]	0.747	—
3	Ma et al. [37]	0.762	—
4	Base model	0.676	0.03
5	With fixed constraint	0.679	0.12
6	With learned constraint	0.727	0.77

Quantitative and Human Evaluation

Knowledge as constraints: Integer Linear Programming (ILP)

- An integer linear program (ILP) is an optimization problem of the form

$$\max_z \quad a \cdot z \quad \text{subj. to} \quad \text{linear constraints on } z$$

$z \in \mathbb{Z}^n$ (integers)

- For a fixed vector a
- Well-engineered solvers exist
 - e.g, Gurobi
 - Useful for prototyping
 - But general not as efficient as dynamic programming

Ex: Casting sequence labeling as an ILP

- Sequence labeling, e.g., named entity recognition
 - names of **people**, **organizations**, **locations**

Ex: "Brendan Iribe, a co-founder of Oculus VR and a prominent University of Maryland donor, is leaving Facebook four years after it purchased his company."

- BIO labeling scheme for NER

x = [Brendan, Iribe, ",", a, co-founder, of, Oculus, VR, and, a, prominent, University, of, Maryland, donor, ",", is, leaving, Facebook, four, years, after, it, purchased, his, company, "."]

y = [B-PER, I-PER, O, O, O, O, B-ORG, I-ORG, O, O, O, B-ORG, I-ORG, I-ORG, O, O, O, B-ORG, O, O, O, O, O, O, O, O]

Ex: Casting sequence labeling as an ILP

Many NLP tasks can be framed as sequence labeling

x = [Brendan, Iribe, ",", a, co-founder, of, Oculus, VR, and, a, prominent, University, of, Maryland, donor, ",", is, leaving, Facebook, four, years, after, it, purchased, his, company, "."]

y = [B-PER, I-PER, O, O, O, O, B-ORG, I-ORG, O, O, O, B-ORG, I-ORG, I-ORG, O, O, O, B-ORG, O, O, O, O, O, O, O]

"BIO" labeling scheme for named entity recognition

Ex: Casting sequence labeling as an ILP

- Step 1: Define variables z as binary indicator variables which encode an output sequence y

$$z_{l,k',k} = \mathbf{1}[\text{label } l \text{ is } k \text{ and label } l - 1 \text{ is } k']$$

- Step 2: Construct the linear objective function

$$a_{l,k',k} = \boldsymbol{w} \cdot \phi_l(\boldsymbol{x}, \langle \dots, k', k \rangle)$$

Ex: Casting sequence labeling as an ILP

- Step 3: Define constraints to ensure a well-formed solution
 - Z's should be binary: for all l, k', k

$$z_{l,k',k} \in \{0, 1\}$$

- For a given position l , there is exactly one active z

$$\sum_k \sum_{k'} z_{l,k',k} = 1 \text{ for all } l$$

- The z 's are internally consistent

$$\sum_{k'} z_{l,k',k} = \sum_{k''} z_{l+1,k,k''} \text{ for all } l, k$$

Key Takeaways

- Two general ways of integrating structured knowledge with ML:
 - Model architecture (inductive bias)
 - Integrating knowledge through learning (loss, constraints)
 - Weak supervision
 - Posterior regularization
 - Integer linear programming (ILP)\
 - Others ..

Questions?