

DSC291: Machine Learning with Few Labels

“Standard Model” of ML

Zhiting Hu

Lecture 24, May 31, 2024

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

This Lecture

- “Standard Model” (20mins)
- Presentation #1 (10mins):
 - Shanglin Zeng, What is important about the No Free Lunch theorems?
- Presentation #2 (10mins):
 - Jiayu Li, TBD
- Presentation #3 (10mins):
 - Jiaxian Xiang, ChineseBERT: Chinese Pretraining Enhanced by Glyph and Pinyin Information

Google form for presentation questions and feedback:

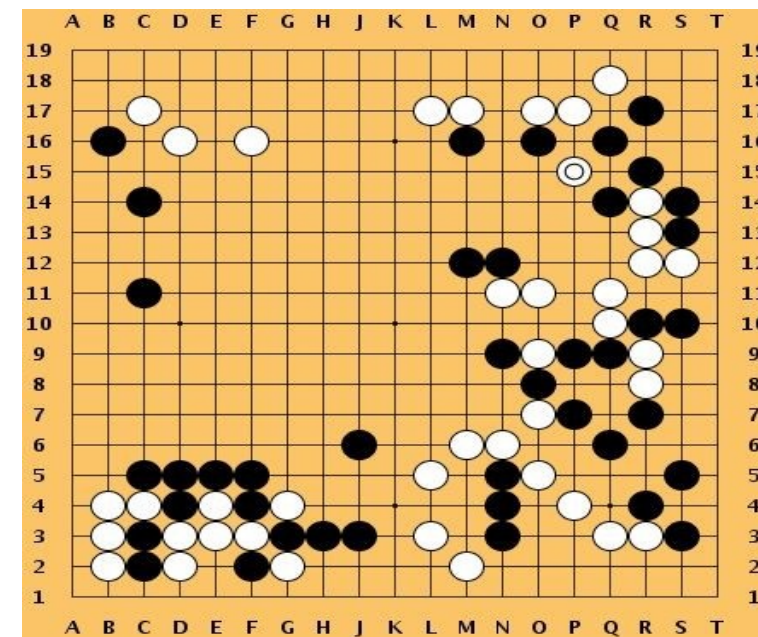
Policy gradients: AlphaGo

Overview:

- Mix of supervised learning and reinforcement learning
- Mix of old methods (Monte Carlo Tree Search) and recent ones (deep RL)

How to beat the Go world champion:

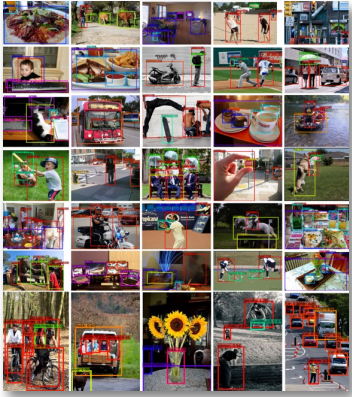
- Featurize the board (stone color, move legality, bias, ...)
- Initialize policy network with supervised training from professional go games, then continue training using policy gradient (play against itself from random previous iterations, +1 / -1 reward for winning / losing)
- Also learn value network (critic)
- Finally, combine policy and value networks in a Monte Carlo Tree Search algorithm to select actions by lookahead search



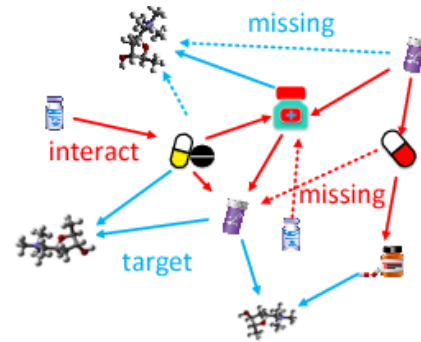
[Silver et al.,
Nature 2016]

This image is [CC0 public domain](#)

Experience of all kinds



Type-2 diabetes is 90% more common than type-1



Data examples

Rules/Constraints

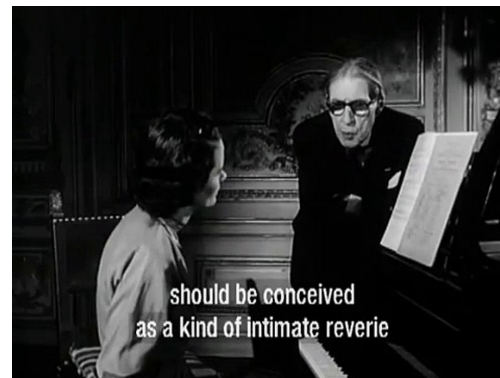
Knowledge graphs

Rewards

Auxiliary agents



Adversaries

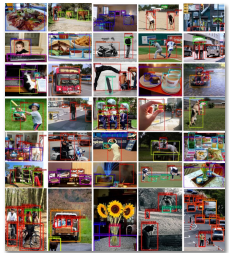


Master classes

...

- *And all combinations of such*
- *Interpolations between such*
- ...

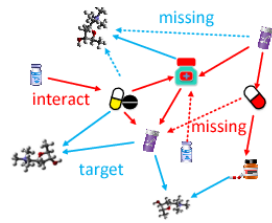
Human learning vs machine learning



Data examples

Type-2 diabetes is 90% more common than type-1

Rules/Constraints



Knowledge graphs



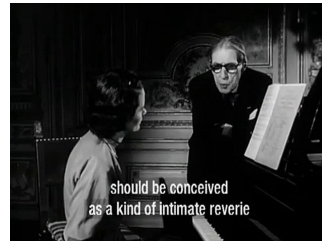
Rewards



Auxiliary agents



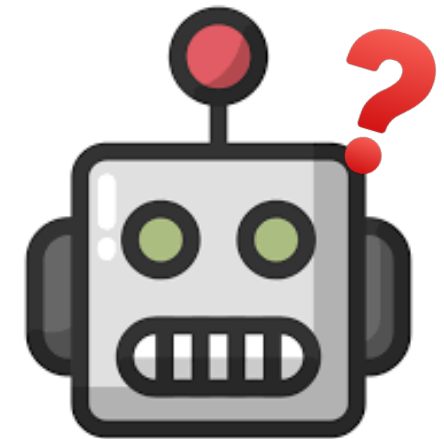
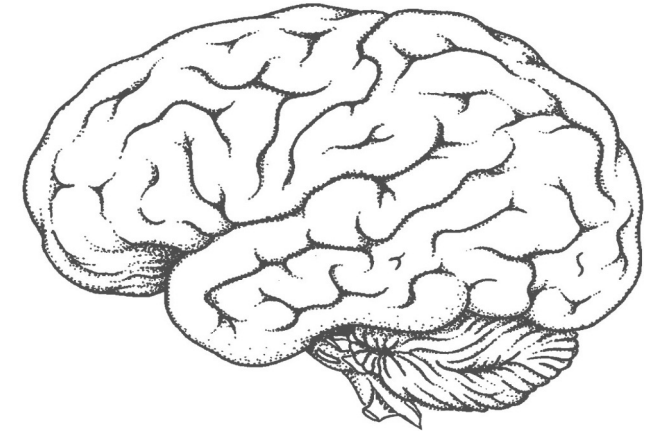
Adversaries



Master classes

...

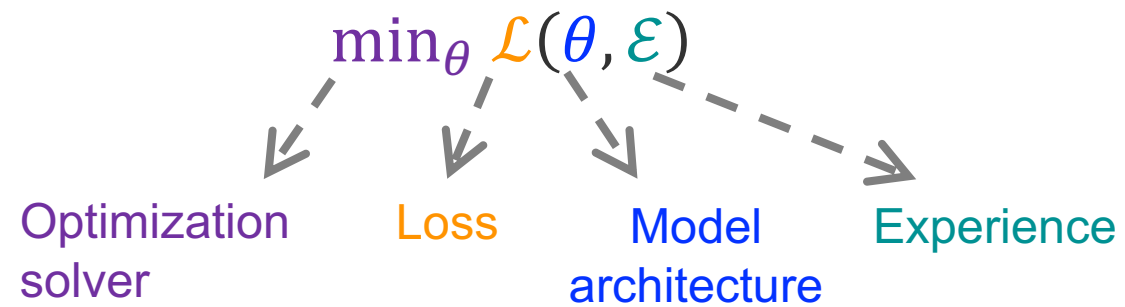
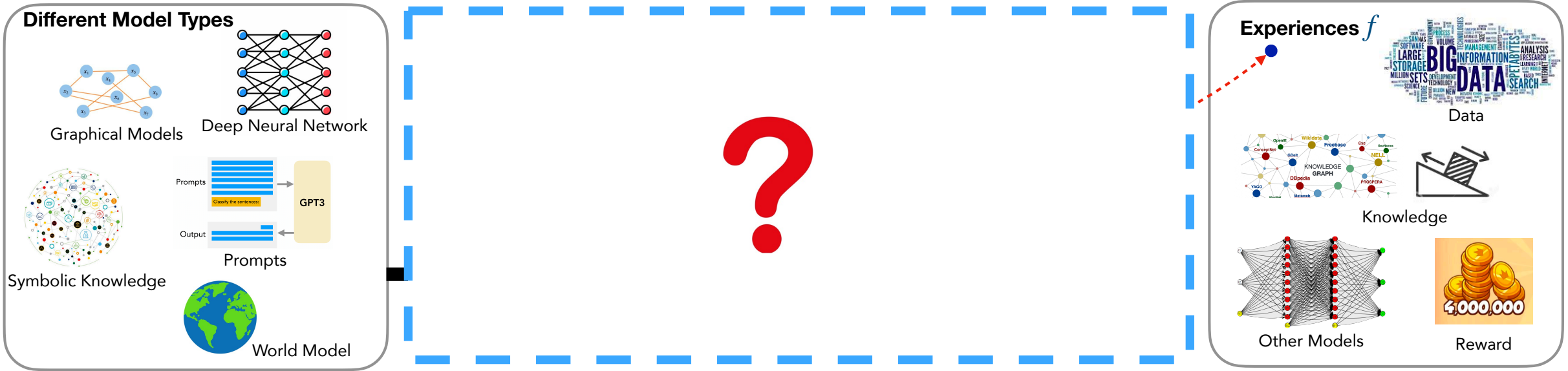
- And all combinations of such
- Interpolations between such
- ...



The zoo of ML/AI models

- Neural networks
 - Convolutional networks
 - AlexNet, GoogleNet, ResNet
 - Recurrent networks, LSTM
 - Transformers
 - BERT, GPTs
- Graphical models
 - Bayesian networks
 - Markov Random fields
 - Topic models, LDA
 - HMM, CRF
- Kernel machines
 - Radial Basis Function Networks
 - Gaussian processes
 - Deep kernel learning
 - Maximum margin
 - SVMs
- Decision trees
- PCA, Probabilistic PCA, Kernel PCA, ICA
- Boosting

The zoo of ML/AI algorithms



The zoo of ML/AI algorithms

maximum likelihood estimation reinforcement learning as inference
data re-weighting inverse RL policy optimization active learning
data augmentation actor-critic reward-augmented maximum likelihood
label smoothing imitation learning softmax policy gradient
adversarial domain adaptation posterior regularization
GANs constraint-driven learning
knowledge distillation intrinsic reward
prediction minimization generalized expectation
regularized Bayes
energy-based GANs learning from measurements
weak/distant supervision

Standard Model in Physics

*Maxwell's Eqns:
original form*

$e + \frac{df}{dx} + \frac{dg}{dy} + \frac{dh}{dz} = 0$	(1) Gauss' Law
$\mu\alpha = \frac{dH}{dy} - \frac{dG}{dz}$ $\mu\beta = \frac{dF}{dz} - \frac{dH}{dx}$ $\mu\gamma = \frac{dG}{dx} - \frac{dF}{dy}$	(2) Equivalent to Gauss' Law for magnetism
$P = \mu \left(\gamma \frac{dy}{dt} - \beta \frac{dz}{dt} \right) - \frac{dF}{dt} - \frac{d\Psi}{dz}$ $Q = \mu \left(\alpha \frac{dz}{dt} - \gamma \frac{dx}{dt} \right) - \frac{dG}{dt} - \frac{d\Psi}{dy}$ $R = \mu \left(\beta \frac{dx}{dt} - \alpha \frac{dy}{dt} \right) - \frac{dH}{dt} - \frac{d\Psi}{dz}$	(3) Faraday's Law (with the Lorentz Force and Poisson's Law)
$\frac{d\gamma}{dy} - \frac{d\beta}{dz} = 4\pi p'$ $\frac{d\alpha}{dz} - \frac{d\gamma}{dx} = 4\pi q'$ $\frac{d\beta}{dx} - \frac{d\alpha}{dy} = 4\pi r'$	(4) Ampère-Maxwell Law
$P = -\xi p \quad Q = -\xi q \quad R = -\xi r$	Ohm's Law
$P = kf \quad Q = kg \quad R = kh$	The electric elasticity equation ($\mathbf{E} = \mathbf{D}/\epsilon$)
$\frac{de}{dt} + \frac{dp}{dx} + \frac{dq}{dy} + \frac{dr}{dz} = 0$	Continuity of charge

*Simplified w/
rotational
symmetry*

$$\nabla \cdot \mathbf{D} = \rho_V$$

$$\nabla \cdot \mathbf{B} = 0$$

$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}$$

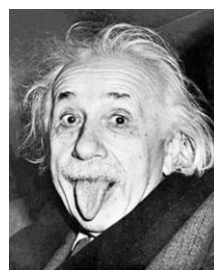
$$\nabla \times \mathbf{H} = \frac{\partial \mathbf{D}}{\partial t} + \mathbf{J}$$



*Further
simplified w/
symmetry of
special relativity*

$$\epsilon^{uvk\lambda} \partial_v F_{k\lambda} = 0$$

$$\partial_v F^{uv} = \frac{4\pi}{c} j^u$$



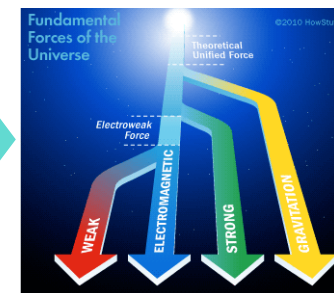
*Standard Model
w/ Yang-Mills
theory and US(3)
symmetry*

$$\mathcal{L}_{gf} = -\frac{1}{2} \text{Tr}(F^2)$$

$$= -\frac{1}{4} F^{a\mu\nu} F_{\mu\nu}^a$$



*Unification of
fundamental
forces?*



*Diverse
electro-
magnetic
theories*



1861

1910s

1970s



Toward a ‘Standard Model’ of Machine Learning

Zhiting Hu^{†,*}, Eric P. Xing^{‡,◇,‡,**}

[†] Halicioğlu Data Science Institute, University of California San Diego, San Diego, USA

[‡] Machine Learning Department, Carnegie Mellon University, Pittsburgh, USA

[‡] Mohamed bin Zayed University of Artificial Intelligence, Abu Dhabi, UAE

[◇] Petuum Inc., Pittsburgh, USA



[Hu & Xing, Harvard Data Science Review, 2022]: <https://arxiv.org/abs/2108.07783>

$$\min_{q, \theta} - \mathbb{E} + \mathbb{D} - \mathbb{H}$$

Experience Divergence Uncertainty

Maximum likelihood estimation (MLE) at a close look:

- The most classical learning algorithm

- Supervised:

- Observe data $\mathcal{D} = \{(\mathbf{x}^*, \mathbf{y}^*)\}$
- Solve with SGD

$$\min_{\theta} - \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} \left[\log p_{\theta}(\mathbf{y}^* | \mathbf{x}^*) \right]$$

- Unsupervised:

- Observe $\mathcal{D} = \{(\mathbf{x}^*)\}$, \mathbf{y} is latent variable
- Posterior $p_{\theta}(\mathbf{y} | \mathbf{x})$
- Solve with EM:

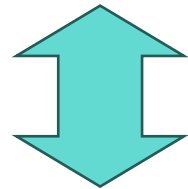
$$\min_{\theta} - \mathbb{E}_{\mathbf{x}^* \sim \mathcal{D}} \left[\log \int_{\mathbf{y}} p_{\theta}(\mathbf{x}^*, \mathbf{y}) \right]$$

- E-step imputes latent variable \mathbf{y} through expectation on complete likelihood
- M-step: supervised MLE

MLE as Entropy Maximization

- Equivalence between supervised MLE and maximum entropy (when p_θ is an exponential family distribution)

$$\min_{\theta} - \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} \left[\log p_{\theta}(\mathbf{y}^* | \mathbf{x}^*) \right]$$



Shannon entropy H

$$\min_{p(\mathbf{x}, \mathbf{y})} H(p)$$

features $T(\mathbf{x}, \mathbf{y})$

$$s.t. \mathbb{E}_p[T(\mathbf{x}, \mathbf{y})] = \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}}[T(\mathbf{x}, \mathbf{y})]$$

data as constraints

MLE as Entropy Maximization

- **Unsupervised MLE** can be achieved by maximizing the negative free energy:
 - Introduce an **auxiliary** variational distribution $q(\mathbf{y}|\mathbf{x})$ (and then play with its entropy and cross entropy, etc.)

$$\log \int_{\mathbf{y}} p_{\theta}(\mathbf{x}^*, \mathbf{y}) \geq H(q(\mathbf{y}|\mathbf{x}^*)) + \mathbb{E}_{q(\mathbf{y}|\mathbf{x}^*)}[\log p_{\theta}(\mathbf{x}^*, \mathbf{y})]$$

The general expression as a constrained optimization:

MaxEnt perspective

$$\begin{aligned} & \min_{q, \theta} \mathcal{L}(q, \theta) \\ & \text{s.t. } q \in \mathcal{Q}. \end{aligned}$$

(auxiliary) distribution q *loss*
constrained set

- Supervised MLE and maximum entropy
- Unsupervised MLE and maximum entropy
- Bayesian inference and maximum entropy
 - Bayesian inference as optimization

The general expression as a constrained optimization:

MaxEnt perspective

$$\begin{aligned} & \min_{q, \theta} \mathcal{L}(q, \theta) \\ & \text{s.t. } q \in \mathcal{Q}. \end{aligned}$$

(auxiliary) distribution q $\mathcal{L}(q, \theta)$ *loss*
constrained set

- Supervised MLE and maximum entropy
- Unsupervised MLE and maximum entropy
- Bayesian inference and maximum entropy

$$\begin{aligned} & \min_{q(\mathbf{z})} -\mathbf{H}(q(\mathbf{z})) + \log p(\mathcal{D}) - \mathbb{E}_{q(\mathbf{z})} \left[\log \pi(\mathbf{z}) + \sum_{\mathbf{x}^* \in \mathcal{D}} \log p(\mathbf{x}^* | \mathbf{z}) \right] \\ & \text{s.t. } q(\mathbf{z}) \in \mathcal{P} \end{aligned}$$

The Standard Equation (SE)

- Let \mathbf{t} be the variable of interest
 - E.g., the input-output pair $\mathbf{t} = (\mathbf{x}, \mathbf{y})$ in a prediction task
 - or $\mathbf{t} = \mathbf{x}$ in generative modeling
- $p_{\theta}(\mathbf{t})$: the target model to be learned
- $q(\mathbf{t})$: auxiliary distribution
- The SE:
$$\min_{q, \theta, \xi} -\alpha \mathbb{H}(q) + \beta \mathbb{D}\left(q(\mathbf{t}), p_{\theta}(\mathbf{t})\right) + U(\xi)$$
$$s. t. -\mathbb{E}_{q(\mathbf{t})} \left[f_k(\mathbf{t}) \right] < \xi_k, \quad k = 1, \dots, K$$
 - Experience function f represents external experiences of different kinds for training the model
 - $f_k(\mathbf{t}) \in \mathbb{R}$: measures the goodness of a configuration \mathbf{t} in light of any given experiences
 - Data, constraints, reward, adversarial discriminators, etc., can all be formulated as an experience function (later)
 - Maximizing $\mathbb{E}_{q(\mathbf{t})} [f_k(\mathbf{t})]$ \rightarrow q is encouraged to produce samples receiving high scores

The Standard Equation (SE)

- Let \mathbf{t} be the variable of interest
 - E.g., the input-output pair $\mathbf{t} = (\mathbf{x}, \mathbf{y})$ in a prediction task
 - or $\mathbf{t} = \mathbf{x}$ in generative modeling

- $p_{\theta}(\mathbf{t})$: the target model to be learned

- $q(\mathbf{t})$: auxiliary distribution

- The SE:
$$\min_{q, \theta, \xi} -\alpha \mathbb{H}(q) + \beta \mathbb{D} \left(q(\mathbf{t}), p_{\theta}(\mathbf{t}) \right) + U(\xi)$$

$$s. t. -\mathbb{E}_{q(\mathbf{t})} \left[f_k(\mathbf{t}) \right] < \xi_k, \quad k = 1, \dots, K$$

- Divergence \mathbb{D} : measures the distance between the target model p_{θ} to be trained and the auxiliary model q
 - E.g., cross entropy

The Standard Equation (SE)

- Let \mathbf{t} be the variable of interest
 - E.g., the input-output pair $\mathbf{t} = (\mathbf{x}, \mathbf{y})$ in a prediction task
 - or $\mathbf{t} = \mathbf{x}$ in generative modeling
- $p_{\theta}(\mathbf{t})$: the target model to be learned
- $q(\mathbf{t})$: auxiliary distribution
- The SE:
$$\min_{q, \theta, \xi} -\alpha \mathbb{H}(q) + \beta \mathbb{D}\left(q(\mathbf{t}), p_{\theta}(\mathbf{t})\right) + U(\xi)$$
$$s. t. -\mathbb{E}_{q(\mathbf{t})} \left[f_k(\mathbf{t}) \right] < \xi_k, \quad k = 1, \dots, K$$
- Uncertainty \mathbb{H} : controls the compactness of the model
 - E.g., Shannon entropy

The Standard Equation (SE)

$$\min_{q, \theta, \xi} -\alpha \mathbb{H}(q) + \beta \mathbb{D}\left(q(\mathbf{t}), p_{\theta}(\mathbf{t})\right) + U(\xi)$$

$$s. t. -\mathbb{E}_{q(\mathbf{t})}\left[f_k(\mathbf{t})\right] < \xi_k, \quad k = 1, \dots, K$$

Assuming penalty $U = \sum_k \xi_k$, and $f = \sum_k f_k$:

$$\min_{q, \theta} -\alpha \mathbb{H}(q) + \beta \mathbb{D}\left(q(\mathbf{t}), p_{\theta}(\mathbf{t})\right) - \mathbb{E}_{q(\mathbf{t})}\left[f(\mathbf{t})\right]$$

3 terms:

Uncertainty
(self-regularization)
e.g., Shannon entropy



Uncertainty

Divergence
(fitness)
e.g., Cross Entropy

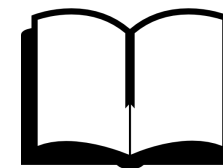


Teacher
 $q(\mathbf{t})$

Student
 $p_{\theta}(\mathbf{t})$

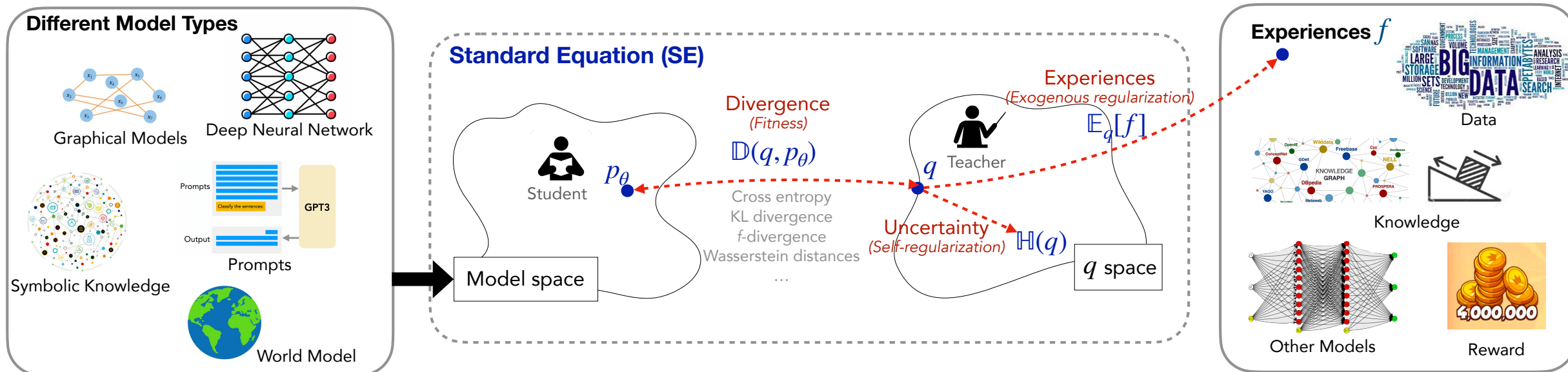
Experiences
(exogenous regularizations)
e.g., data examples, rules

Textbook
 $f(\mathbf{t})$



The Standard Equation (SE)

$$\min_{q, \theta} -\alpha \mathbb{H}(q) + \beta \mathbb{D}\left(q(t), p_{\theta}(t)\right) - \mathbb{E}_{q(t)} \left[f(t) \right]$$



[Note: in SE, experience function f can also depends on θ . See the paper for mor details]

The Standard Equation (SE)

$$\min_{q, \theta} -\alpha \mathbb{H}(q) + \beta \mathbb{D} \left(q(\mathbf{t}), p_{\theta}(\mathbf{t}) \right) - \mathbb{E}_{q(\mathbf{t})} \left[f(\mathbf{t}) \right]$$

- Formulates a large space of learning algorithms, which encompasses many well-known algorithms

SE encompasses many well-known algorithms (more later)

Experience type	Experience function f	Divergence \mathbb{D}	α	β	Algorithm
Data instances	$f_{\text{data}}(\mathbf{x}; \mathcal{D})$	CE	1	1	Unsupervised MLE
	$f_{\text{data}}(\mathbf{x}, \mathbf{y}; \mathcal{D})$	CE	1	ϵ	Supervised MLE
	$f_{\text{data-self}}(\mathbf{x}, \mathbf{y}; \mathcal{D})$	CE	1	ϵ	Self-supervised MLE
	$f_{\text{data-w}}(\mathbf{t}; \mathcal{D})$	CE	1	ϵ	Data Re-weighting
	$f_{\text{data-aug}}(\mathbf{t}; \mathcal{D})$	CE	1	ϵ	Data Augmentation
	$f_{\text{active}}(\mathbf{x}, \mathbf{y}; \mathcal{D})$	CE	1	ϵ	Active Learning (Ertekin et al., 2007)
Knowledge	$f_{\text{rule}}(\mathbf{x}, \mathbf{y})$	CE	1	1	Posterior Regularization (Ganchev et al., 2010)
	$f_{\text{rule}}(\mathbf{x}, \mathbf{y})$	CE	\mathbb{R}	1	Unified EM (Samdani et al., 2012)
Reward	$\log Q^\theta(\mathbf{x}, \mathbf{y})$	CE	1	1	Policy Gradient
	$\log Q^\theta(\mathbf{x}, \mathbf{y}) + Q^{\text{in}, \theta}(\mathbf{x}, \mathbf{y})$	CE	1	1	+ Intrinsic Reward
	$Q^\theta(\mathbf{x}, \mathbf{y})$	CE	$\rho > 0$	$\rho > 0$	RL as Inference
Model	$f_{\text{model}}^{\text{mimicking}}(\mathbf{x}, \mathbf{y}; \mathcal{D})$	CE	1	ϵ	Knowledge Distillation (G. Hinton et al., 2015)
Variational	binary classifier	JSD	0	1	Vanilla GAN (Goodfellow et al., 2014)
	discriminator	f -divergence	0	1	f-GAN (Nowozin et al., 2016)
	1-Lipschitz discriminator	W_1 distance	0	1	WGAN (Arjovsky et al., 2017)
	1-Lipschitz discriminator	KL	0	1	PPO-GAN (Y. Wu et al., 2020)
Online	$f_\tau(\mathbf{t})$	CE	$\rho > 0$	$\rho > 0$	Multiplicative Weights (Freund & Schapire, 1997)

SE Component: Experience Function f

Different choices of experience function f lead to different algorithms:

$$\min_{q, \theta} - \mathbb{E}_{q(x, y)} \left[f(x, y) \right] + \beta \mathbb{D} \left(q(x, y), p_{\theta}(x, y) \right) - \alpha \mathbb{H}(q)$$

Experience
(exogenous regularizations)
e.g., data examples, rules

Set Divergence to Cross
Entropy
 $\mathbb{D}(q, p_{\theta}) = -\mathbb{E}_q[\log p_{\theta}]$

Set Uncertainty to
Shannon Entropy
 $\mathbb{H}(q) = H(q) := -\mathbb{E}_q[\log q]$

SE with supervised data experience

$$\min_{q, \theta} -\alpha \mathbb{H}(q) + \beta \mathbb{D} \left(q(\mathbf{t}), p_{\theta}(\mathbf{t}) \right) - \mathbb{E}_{q(\mathbf{t})} \left[f(\mathbf{t}) \right]$$

- Input-output variables $\mathbf{t} = (\mathbf{x}, \mathbf{y})$
- Experience: dataset $\mathcal{D} = \{(\mathbf{x}^*, \mathbf{y}^*)\}$ of size N
 - defines the empirical distribution

$$\tilde{p}(\mathbf{x}, \mathbf{y}) = \frac{m(\mathbf{x}, \mathbf{y})}{N} = \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} [\mathbb{1}_{(\mathbf{x}^*, \mathbf{y}^*)}(\mathbf{x}, \mathbf{y})]$$

- Define the experience function

$$f := f_{data}(\mathbf{x}, \mathbf{y}; \mathcal{D}) = \log \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} [\mathbb{1}_{(\mathbf{x}^*, \mathbf{y}^*)}(\mathbf{x}, \mathbf{y})]$$

- Let \mathbb{D} cross entropy, \mathbb{H} Shannon entropy, $\alpha = 1, \beta = \epsilon$ (a very small value)

$$\min_{q, \theta} -H(q) - \epsilon \mathbb{E}_q \left[\log p_{\theta}(\mathbf{x}, \mathbf{y}) \right] - \mathbb{E}_q \left[f_{data}(\mathbf{x}, \mathbf{y}; \mathcal{D}) \right]$$

SE with supervised data experience

$$f := f_{data}(\mathbf{x}, \mathbf{y}; \mathcal{D}) = \log \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} \left[\mathbb{1}_{(\mathbf{x}^*, \mathbf{y}^*)}(\mathbf{x}, \mathbf{y}) \right]$$

$$\min_{q, \theta} -H(q) - \epsilon \mathbb{E}_q \left[\log p_{\theta}(\mathbf{x}, \mathbf{y}) \right] - \mathbb{E}_q \left[f_{data}(\mathbf{x}, \mathbf{y}; \mathcal{D}) \right]$$

- At each iteration n :

Teacher: $q^{(n+1)}(\mathbf{t}) = \exp \left\{ \frac{\beta \log p_{\theta^{(n)}}(\mathbf{t}) + f(\mathbf{t})}{\alpha} \right\} / Z \approx \tilde{p}(\mathbf{x}, \mathbf{y})$

Student: $\theta^{(n+1)} = \operatorname{argmax}_{\theta} \mathbb{E}_{q^{(n+1)}(\mathbf{t})} \left[\log p_{\theta}(\mathbf{t}) \right],$

Maximizes data log-likelihood

q reduces to the empirical distribution

- Recovers supervised MLE!

SE with unsupervised data experience

$$\min_{q, \theta} -\alpha \mathbb{H}(q) + \beta \mathbb{D}\left(q(\mathbf{t}), p_{\theta}(\mathbf{t})\right) - \mathbb{E}_{q(\mathbf{t})} \left[f(\mathbf{t}) \right]$$

- Input-output variables $\mathbf{t} = (\mathbf{x}, \mathbf{y})$
- Experience: dataset $\mathcal{D} = \{(\mathbf{x}^*)\}$ of size N , i.e., we only observe the \mathbf{x} part
 - defines the empirical distribution

$$\tilde{p}(\mathbf{x}) = \frac{m(\mathbf{x})}{N} = \mathbb{E}_{\mathbf{x}^* \sim \mathcal{D}} [\mathbb{1}_{\mathbf{x}^*}(\mathbf{x})]$$

- Define the experience function

$$f := f_{data}(\mathbf{x}; \mathcal{D}) = \log \mathbb{E}_{\mathbf{x}^* \sim \mathcal{D}} [\mathbb{1}_{\mathbf{x}^*}(\mathbf{x})]$$

- Let \mathbb{D} cross entropy, \mathbb{H} Shannon entropy, $\alpha = 1, \beta = 1$

$$\min_{q, \theta} -H(q) - \mathbb{E}_q \left[\log p_{\theta}(\mathbf{x}, \mathbf{y}) \right] - \mathbb{E}_q \left[f_{data}(\mathbf{x}; \mathcal{D}) \right]$$

- Assume $q(\mathbf{x}, \mathbf{y}) = \tilde{p}(\mathbf{x})q(\mathbf{y}|\mathbf{x})$

Recovers unsupervised
MLE (EM)!

SE with manipulated data experience

- Input-output variables $\mathbf{t} = (\mathbf{x}, \mathbf{y})$
- Experience: dataset $\mathcal{D} = \{(\mathbf{x}^*, \mathbf{y}^*)\}$ of size N
 - defines the empirical distribution

$$\tilde{p}(\mathbf{x}, \mathbf{y}) = \frac{m(\mathbf{x}, \mathbf{y})}{N} = \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} [\mathbb{1}_{(\mathbf{x}^*, \mathbf{y}^*)}(\mathbf{x}, \mathbf{y})]$$

- Define the experience function

$$f := f_{data}(\mathbf{x}, \mathbf{y}; \mathcal{D}) = \log \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} [\mathbb{1}_{(\mathbf{x}^*, \mathbf{y}^*)}(\mathbf{x}, \mathbf{y})]$$

- The similarity measure $\mathbb{1}_a(b)$ is too restrictive. Let's enrich it:

- Don't have to be 0/1, we can scale it

$$f := f_{data-w}(\mathbf{x}, \mathbf{y}; \mathcal{D}) = \log \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} [w(\mathbf{x}^*, \mathbf{y}^*) \cdot \mathbb{1}_{(\mathbf{x}^*, \mathbf{y}^*)}(\mathbf{x}, \mathbf{y})]$$

- Plug f_{data-w} into SE, keep all other configurations the same as supervised MLE, we recover **data re-weighting** in the "student" step

$$\max_{\theta} \mathbb{E}_{\mathbf{t}^* \sim \mathcal{D}} [w(\mathbf{t}^*) \cdot \log p_{\theta}(\mathbf{t}^*)]$$

SE with manipulated data experience

- Input-output variables $\mathbf{t} = (\mathbf{x}, \mathbf{y})$
- Experience: dataset $\mathcal{D} = \{(\mathbf{x}^*, \mathbf{y}^*)\}$ of size N
 - defines the empirical distribution

$$\tilde{p}(\mathbf{x}, \mathbf{y}) = \frac{m(\mathbf{x}, \mathbf{y})}{N} = \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} [\mathbb{1}_{(\mathbf{x}^*, \mathbf{y}^*)}(\mathbf{x}, \mathbf{y})]$$

- Define the experience function

$$f := f_{data}(\mathbf{x}, \mathbf{y}; \mathcal{D}) = \log \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} [\mathbb{1}_{(\mathbf{x}^*, \mathbf{y}^*)}(\mathbf{x}, \mathbf{y})]$$

- The similarity measure $\mathbb{1}_a(b)$ is too restrictive. Let's enrich it:

- Don't have to match exactly, we can relax it

$$f := f_{data-aug}(\mathbf{x}, \mathbf{y}; \mathcal{D}) = \log \mathbb{E}_{(\mathbf{x}^*, \mathbf{y}^*) \sim \mathcal{D}} [a_{(\mathbf{x}^*, \mathbf{y}^*)}(\mathbf{x}, \mathbf{y})]$$

- $a_{(\mathbf{x}^*, \mathbf{y}^*)}(\mathbf{x}, \mathbf{y})$: assigns non-zero probability to not only the exact $(\mathbf{x}^*, \mathbf{y}^*)$ but also other (\mathbf{x}, \mathbf{y}) configurations
- Plug $f_{data-aug}$ into SE, keep all other configurations the same as supervised MLE, we recover **data augmentation** in the "student" step $\max_{\theta} \mathbb{E}_{\mathbf{t}^* \sim \mathcal{D}, \mathbf{t} \sim a_{\mathbf{t}^*}(\mathbf{t})} [\log p_{\theta}(\mathbf{t})]$.

SE with **reward** experience -- Policy Gradient

$$\min_{q, \theta} -\alpha H(q) - \beta \mathbb{E}_q \left[\log p_{\theta}(\mathbf{x}, \mathbf{y}) \right] - \mathbb{E}_q \left[f(\mathbf{x}, \mathbf{y}) \right]$$

- Policy gradient

$$f^{\theta}(\mathbf{x}, \mathbf{y}) := \log Q^{\theta}(\mathbf{x}, \mathbf{y}) \quad \alpha = \beta = 1$$

- Teacher step: $q^{(n)}(\mathbf{x}, \mathbf{y}) = p_{\theta^{(n)}}(\mathbf{x}, \mathbf{y}) Q^{\theta^{(n)}}(\mathbf{x}, \mathbf{y}) / Z$
- Student step:

$$\begin{aligned} & \mathbb{E}_{q^{(n)}(\mathbf{x}, \mathbf{y})} \left[\nabla_{\theta} \log p_{\theta}(\mathbf{x}, \mathbf{y}) \right] + \mathbb{E}_{q^{(n)}(\mathbf{x}, \mathbf{y})} \left[\nabla_{\theta} f_{\text{reward}, 1}^{\theta}(\mathbf{x}, \mathbf{y}) \right] \Big|_{\theta = \theta^{(n)}} \\ &= 1/Z \cdot \sum_{\mathbf{x}} p_0(\mathbf{x}) \nabla_{\theta} \sum_{\mathbf{y}} p_{\theta}(\mathbf{y}|\mathbf{x}) Q^{\theta}(\mathbf{x}, \mathbf{y}) \Big|_{\theta = \theta^{(n)}} \quad (\text{log-derivative trick}) \\ &= 1/Z \cdot \sum_{\mathbf{x}} \mu^{\theta}(\mathbf{x}) \sum_{\mathbf{y}} Q^{\theta}(\mathbf{x}, \mathbf{y}) \nabla_{\theta} p_{\theta}(\mathbf{y}|\mathbf{x}) \Big|_{\theta = \theta^{(n)}} \quad (\text{policy gradient theorem}) \end{aligned}$$

policy gradient



$$\begin{aligned} & \text{(auxiliary) distribution } q \leftarrow \min_{q, \theta} \mathcal{L}(q, \theta) \rightarrow \text{loss} \\ & \text{s.t. } q \in \mathcal{Q}. \rightarrow \text{constrained set} \end{aligned}$$

Key Takeaways

- The MaxEnt perspective converts learning into a constrained optimization problem
- The standard equation (SE):

$$\min_{q, \theta} -\alpha \mathbb{H}(q) + \beta \mathbb{D}\left(q(t), p_{\theta}(t)\right) - \mathbb{E}_{q(t)} \left[f(t) \right]$$

3 terms:

Uncertainty
(self-regularization)
e.g., Shannon entropy



Divergence
(fitness)
e.g., Cross Entropy



Experiences
(exogenous regularizations)
e.g., data examples, rules



Questions?