# Presenters

Zhiting Hu

Assistant Professor
@UCSD

Hao Zhang

Assistant Professor
@UCSD

Eric P. Xing

Professor @ CMU
President @ MBZUAI
Co-founder @ Petuum

UC San Diego

MOHAMED BIN ZAYED
UNIVERSITY OF
ARTIFICIAL INTELLIGENCE
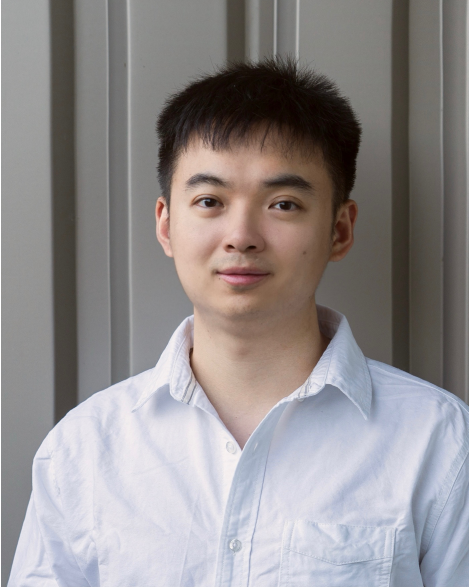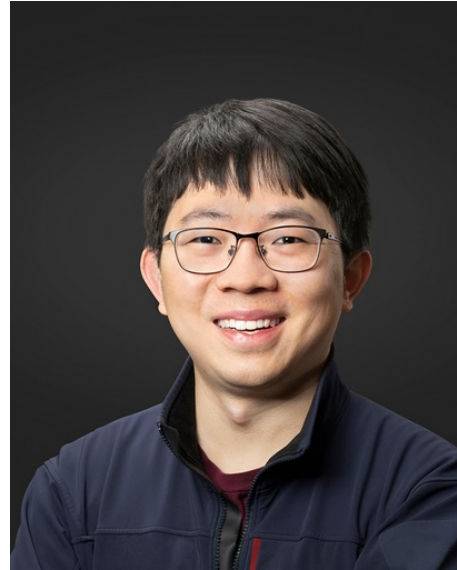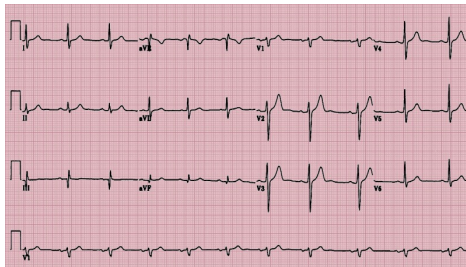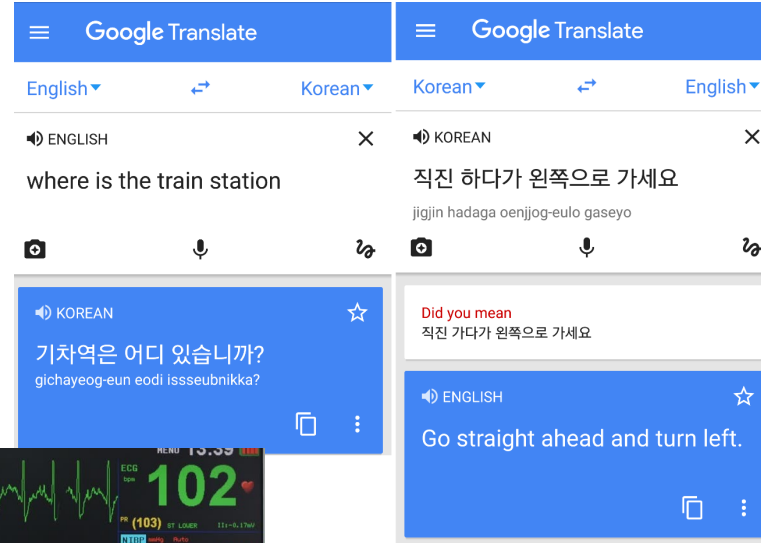
Carnegie Mellon University
School of Computer Science

Petuum

# Real-world Machine Learning Problems

# Data and experience of all kinds



Data examples

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

Rules/Constraints

Knowledge graphs

SCORE: 107

Rewards

Auxiliary agents

Adversaries

should be conceived as a kind of intimate reverie

Master classes

...

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

# An Example: ML for Healthcare



This is where evidence and information start

# A ready-to-use real AI solution is extremely complex, given all these experiences to train on

## Use Case: Automatic Medical (or other) Report Generation



**Findings:**
There are no focal areas of consolidation.
No suspicious pulmonary opacities.
Heart size within normal limits.
No pleural effusions.
There is no evidence of pneumothorax.
Degenerative changes of the thoracic spine.

**Impression:**
No acute cardiopulmonary abnormality.

- Abnormal regions in medical images are difficult to identify.
- How to localize the image regions and tags that are relevant to a sentence?
- How to distribute topics across sentences
- How to make report readable to humans?

AutoML ?!



Raw Data Enrichment

Model/Algorithm

System/Infra

Pe

# Inter-operability between diverse systems?



Data/ML Process Builder

ML Systems

DL Systems

Data Systems

DL Model Interchange ONNX

Hardware requests

Distributed Communications Backend

Hardware Resource Management

Containers and Storage Volumes

# An AI solution

Data wrangling

Feature engineering

Model compiling

Algorithm designing

Distributed training

Debugging

Resource provisioning

Hardware management

Fault recovery

…etc

# Build versus Craft



- **Modules, Building-blocks**
- **Nuts and Bolts**
- **Interoperability**
- **Process**
- **Soundness**

# Schedule

- **Lecture#1:** Theory: The Standard Model of ML

    A blueprint of ML paradigms for ALL experience

    *(Jan 19 Thursday, 4:45pm-6:15pm UK Time)*



- **Lecture#2:** Tooling: Operationalizing The Standard Model

    Compose your ML solutions like playing Lego

    *(Jan 20 Thursday, 1:00pm-2:30pm)*



- **Lecture#3:** Computing: Modern infrastructure for productive ML

    Automatic tuning, distributing, and scheduling

    *(Jan 20 Thursday, 4:45pm-6:15pm)*

# Theory:

# The Standard Model – A Blueprint for ML

# Experience of all kinds


Data examples

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

Rules/Constraints


Knowledge graphs


Rewards


Auxiliary agents


Adversaries


should be conceived as a kind of intimate reverie

Master classes

...

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

# Experience of all kinds

Type-2

missing

*Data examp...*

*Auxiliary agents*

*...ations thereof*

*Adversaries*

*Master classes*

SCORE: 0

# Experience of all kinds



Data examples

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

Rules/Constraints

Knowledge graphs

SCORE: 107

Rewards

Auxiliary agents

Adversaries

should be conceived as a kind of intimate reverie

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



**Different Model Types**

Graphical Models

Deep Neural Network

Symbolic Knowledge

Prompts

World Model

**?**

**Experiences** $f$

Data

Knowledge

Other Models

Reward

# 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

softmax policy gradient

imitation learning

adversarial domain adaptation

posterior regularization

GANs

constraint-driven learning

knowledge distillation

intrinsic reward

prediction minimization

generalized expectation

regularized Bayes

learning from measurements

energy-based GANs

weak/distant supervision

# Really hard to navigate, and to realize



- Depending on individual's expertise and creativity

- Bespoke, delicate pieces of art

- Like an airport with different runways for every different types of aircrafts

# Physics in the 1800's



- Electricity & magnetism:
  - Coulomb's law, Ampère, Faraday, …

- Theory of light beams:
  - Particle theory: Isaac Newton, Laplace, Plank
  - Wave theory: Grimaldi, Chris Huygens, Thomas Young, Maxwell

- Law of gravity
  - Aristotle, Galileo, Newton, …

# Standard Model in Physics



Diverse electro-magnetic theories

*Maxwell's Eqns:* original form

Simplified w/ rotational symmetry

Further simplified w/ symmetry of special relativity

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

*Unification* of fundamental forces?

$$\nabla \cdot \mathbf{D} = \rho_v$$
$$\nabla \cdot \mathbf{B} = 0$$
$$\nabla \times \mathbf{E} = -\frac{\partial \mathbf{B}}{\partial t}$$
$$\nabla \times \mathbf{H} = \frac{\partial \mathbf{D}}{\partial t} + \mathbf{J}$$

$$\varepsilon^{uvk\lambda}\partial_v F_{k\lambda} = 0$$
$$\partial_v F^{uV} = \frac{4\pi}{c}j^u$$

$$\mathcal{L}_{\text{gf}} = -\frac{1}{2}\text{Tr}(F^2)$$
$$= -\frac{1}{4}F^{a\mu\nu}F^a_{\mu\nu}$$

1861         1910s         1970s

# Quest for more standardized, unified ML principles

Machine Learning 3: 253–259, 1989
© 1989 Kluwer Academic Publishers – Manufactured in The Netherlands

EDITORIAL

Toward a Unified Science of Machine Learning

[P. Langley, 1989]

EARLY ACCESS

Model-Based Machine Learning

Click to open

John Winn and Christopher Bishop
with
Thomas Diethe

"Pedro Domingos demystifies machine learning and shows how wondrous and exciting the future will be."
—Walter Isaacson

THE MASTER ALGORITHM

HOW THE QUEST FOR THE ULTIMATE LEARNING MACHINE WILL REMAKE OUR WORLD

PEDRO DOMINGOS

REVIEW ━━━━━━━━━━━━━━━━━━━ Communicated by Steven Nowlan

A Unifying Review of Linear Gaussian Models

Sam Roweis*
*Computation and Neural Systems, California Institute of Technology, Pasadena, CA 91125, U.S.A.*

Zoubin Ghahramani*
*Department of Computer Science, University of Toronto, Toronto, Canada*

# Toward A "Standard Model" of ML

- Loss

- Experience

- Optimization solver

- Model architecture



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

Experience    Divergence    Uncertainty

# Toward A "Standard Model" of ML

Toward a 'Standard Model' of Machine Learning

Zhiting Hu[†,*], Eric P. Xing[‡,◊,♮,**]

[†] Halıcıoğlu 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

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

Experience     Divergence     Uncertainty

# Maximum likelihood estimation (MLE) at a close look:

- The most classical learning algorithm

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

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

- Unsupervised:
  - Observe $\mathcal{D} = \{(\boldsymbol{x}^*)\}$, $\boldsymbol{y}$ is latent variable
  - Posterior $p_{\theta}(\boldsymbol{y} | \boldsymbol{x})$
  - Solve with EM:
    - E-step imputes latent variable $\boldsymbol{y}$ through expectation on complete likelihood
    - M-step: supervised MLE

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

# MLE as Entropy Maximization

- Duality between supervised MLE and maximum entropy, when $p$ is exponential family

*Shannon entropy $H$*

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

*features $T(\boldsymbol{x},\boldsymbol{y})$*

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

data as constraints

*Solve w/ Lagrangian method* $\Downarrow$

$$p(\boldsymbol{x},\boldsymbol{y}) = \exp\{\boldsymbol{\theta} \cdot T(\boldsymbol{x})\} \ / \ Z(\boldsymbol{\theta})$$

*Lagrangian multiplier $\boldsymbol{\theta}$*

$$\min_{\theta} -\mathbb{E}_{(\boldsymbol{x}^*,\boldsymbol{y}^*)\sim\mathcal{D}}[\boldsymbol{\theta} \cdot T(\boldsymbol{x},\boldsymbol{y})] + \log Z(\boldsymbol{\theta})$$ --> *Negative log-likelihood*

*How to estimate $\theta$ — Close form? SGD?*

# MLE as Entropy Maximization

- Unsupervised MLE can be achieved by maximizing the negative free energy:

  - Introduce an auxiliary variational distribution $q(\boldsymbol{y}|\boldsymbol{x})$ (and then play with its entropy and cross entropy, etc.)

$$\log \int_{\boldsymbol{y}} p_\theta(\boldsymbol{x}^*, \boldsymbol{y}) = \mathbb{E}_{q(\boldsymbol{y}|\boldsymbol{x}^*)}\left[\log \frac{p_\theta(\boldsymbol{x}^*, \boldsymbol{y})}{q(\boldsymbol{y}|\boldsymbol{x}^*)}\right] + \text{KL}\big(q(\boldsymbol{y}|\boldsymbol{x}^*) \,||\, p_\theta(\boldsymbol{y}|\boldsymbol{x}^*)\big)$$

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

# Algorithms for Unsupervised MLE

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



1) Solve with EM

$$\log \int_{\boldsymbol{y}} p_{\theta}(\boldsymbol{x}^*, \boldsymbol{y}) = \mathbb{E}_{q(\boldsymbol{y}|\boldsymbol{x}^*)} \left[ \log \frac{p_{\theta}(\boldsymbol{x}^*, \boldsymbol{y})}{q(\boldsymbol{y}|\boldsymbol{x}^*)} \right] + \mathrm{KL}\big(q(\boldsymbol{y}|\boldsymbol{x}^*) \,||\, p_{\theta}(\boldsymbol{y}|\boldsymbol{x}^*)\big)$$

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

- ❑ E-step:  Maximize $\mathcal{L}(q, \boldsymbol{\theta})$ w.r.t $q$, equivalent to minimizing KL by setting

$$q(\boldsymbol{y}|\boldsymbol{x}^*) = p_{\theta^{old}}(\boldsymbol{y}|\boldsymbol{x}^*)$$

- ❑ M-step: Maximize $\mathcal{L}(q, \boldsymbol{\theta})$ w.r.t $\boldsymbol{\theta}$:  $\max_{\theta} \mathbb{E}_{q(\boldsymbol{y}|\boldsymbol{x}^*)}[\log p_{\theta}(\boldsymbol{x}^*, \boldsymbol{y})]$

# Algorithms for Unsupervised MLE (cont'd)

$$\log \int_{\boldsymbol{y}} p_\theta(\boldsymbol{x}^*, \boldsymbol{y}) = \mathbb{E}_{q(\boldsymbol{y}|\boldsymbol{x}^*)}\left[\log \frac{p_\theta(\boldsymbol{x}^*, \boldsymbol{y})}{q(\boldsymbol{y}|\boldsymbol{x}^*)}\right] + \mathrm{KL}\big(q(\boldsymbol{y}|\boldsymbol{x}^*) \,||\, p_\theta(\boldsymbol{y}|\boldsymbol{x}^*)\big)$$

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

2)  When model $p_{\boldsymbol{\theta}}$ is complex, directly working with the true posterior $p_{\boldsymbol{\theta}}(\boldsymbol{y}|\boldsymbol{x}^*)$ is intractable $\Rightarrow$ Variational EM

- Consider a sufficiently restricted family $Q$ of $q(\boldsymbol{y}|\boldsymbol{x})$ so that minimizing the KL is tractable
  - E.g., parametric distributions, factorized distributions

- E-step: Maximize $\mathcal{L}(q, \boldsymbol{\theta})$ w.r.t $q \in Q$, equivalent to minimizing KL
- M-step: Maximize $\mathcal{L}(q, \boldsymbol{\theta})$ w.r.t $\boldsymbol{\theta}$ : $\max\limits_{\theta} \mathbb{E}_{q(\boldsymbol{y}|\boldsymbol{x}^*)}[\log p_\theta(\boldsymbol{x}^*, \boldsymbol{y})]$

# Algorithms for Unsupervised MLE (cont'd)

$$\log \int_{y} p_\theta(x^*, y) = \mathbb{E}_{q(y|x^*)} \left[ \log \frac{p_\theta(x^*, y)}{q(y|x^*)} \right] + \mathrm{KL}\big(q(y|x^*) \,||\, p_\theta(y|x^*)\big)$$

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

3) When $q$ is complex, e.g., deep NNs, optimizing $q$ in E-step is difficult (e.g., high variance) $\Rightarrow$ **Wake-Sleep algorithm** [Hinton et al., 1995]

- Sleep-phase (E-step):  $\min_\phi \mathrm{KL}(p_\theta(y|x^*)||q_\phi(y|x^*))$  - - - - → *Reverse KL*

- Wake-phase (M-step): Maximize $\mathcal{L}(q, \theta)$ w.r.t $\theta$ : $\max_\theta \mathbb{E}_{q(y|x^*)}[\log p_\theta(x^*, y)]$

*Other tricks: reparameterization in VAE ('2014), control variates in NVIL ('2014)*

# Quick summary of MLE

- Supervised:
  - Duality with MaxEnt
  - Solve with SGD

- Unsupervised:
  - Lower bounded by negative free energy
  - Solve with EM, VEM, Wake-Sleep, …

- Close connections to MaxEnt

- With MaxEnt, algorithms (e.g., EM) arises naturally

# Posterior Regularization (PR)

- Make use of constraints in Bayesian learning
  - An auxiliary posterior distribution $q(\theta)$
  - Slack variable $\xi$, constant weight $\alpha = \beta > 0$

$$\min_{q,\xi} \ -\alpha H(q) - \beta \mathbb{E}_q \Big[ \log p_\theta(x,y) \Big] + \xi$$

$$s.t. \ -\mathbb{E}_q [\ f_\theta(x,y)\ ] \le \xi \qquad \text{[Ganchev et al., 2010]}$$

  - E.g., max-margin constraint for linear regression [Jaakkola et al., 1999] and general models (e.g., LDA, NNs) [Zhu et al., 2014]

- Solution for $q$

$$q(\theta) = \exp\left\{ \frac{\beta \log p_\theta(x,y) + f_\theta(x,y)}{\alpha} \right\} / Z$$

# More general learning leveraging PR

- No need to limit to Bayesian learning
- E.g., Complex rule constraints on general models [Hu et al., 2016], where
  - $q$ can be over arbitrary variables, e.g., $q(x, y)$
  - $p_\theta(x, y)$ is NNs of arbitrary architectures with parameters $\theta$

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

$$s.t. \ \mathbb{E}_{q(x,y)} \left[ 1 - r(x, y) \right] \leq \xi$$

E.g., $r(x, y)$ is a 1st-order logical rule:

If sentence $x$ contains word ``but''

$\Rightarrow$ its sentiment $y$ is the same as the sentiment after "but"

# EM for the general PR

- Rewrite without slack variable:

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

  ○ Solve with EM

  ▪ E-step: $\quad q(x,y) = \exp \left\{ \dfrac{\beta \log p_\theta(x,y) + f(x,y)}{\alpha} \right\} / Z$

  ▪ M-step: $\quad \min_{\theta} \mathbb{E}_q \left[ \log p_\theta(x,y) \right]$

# Reformulating unsupervised MLE with PR

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

- Introduce arbitrary $q(y|x)$

$$\min_{q,\theta,\xi} \ -\alpha H(q) - \beta \mathbb{E}_q\Big[\log p_\theta(x, y)\Big] + \xi$$

$$s.t. \ -\mathbb{E}_q\Big[f(x\,;\,\mathcal{D})\Big] < \xi$$

Data as constraint.
Given $x \sim \mathcal{D}$, this constraint doesn't influence the solution of $q$ and $\boldsymbol{\theta}$

- $f(x\,;\,\mathcal{D}) := \log \mathbb{E}_{x^* \sim \mathcal{D}}[\, \mathbb{1}_{x^*}(x)\,]$
  - A constraint saying $x$ must equal to one of the true data points
  - Or alternatively, the (log) expected similarity of $x$ to dataset $\mathcal{D}$, with $\mathbb{1}(\cdot)$ as the similarity measure (we'll come back to this later)
- $\alpha = \beta = 1$

# A "Standard Model" of Machine Learning

# The Standard Equation (SE)



$$\min_{q,\,\theta,\,\xi \geq 0} \beta \mathbb{D}\Big(q(\boldsymbol{x},\boldsymbol{y}),\, p_\theta(\boldsymbol{x},\boldsymbol{y})\Big) - \alpha \mathbb{H}(q) + \xi$$

$$s.t. \; -\mathbb{E}_{q(\boldsymbol{x},\boldsymbol{y})}\Big[f(\boldsymbol{x},\boldsymbol{y})\Big] < \xi$$

Equivalently:

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

*3 terms:*

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

**Divergence**
(fitness)
e.g., Cross Entropy

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

Textbook
$f(\boldsymbol{x},\boldsymbol{y}|\,.\,)$

Teacher
$q(\boldsymbol{x},\boldsymbol{y})$

Student
$p_\theta(\boldsymbol{x},\boldsymbol{y})$

Uncertainty

42

# The Standard Equation (SE)

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



**Different Model Types**

Graphical Models

Deep Neural Network

Symbolic Knowledge

Prompts

Classify the sentences:

GPT3

Output

Prompts

World Model

**Standard Equation (SE)**

Divergence
*(Fitness)*

$\mathbb{D}(q, p_\theta)$

Cross entropy
KL divergence
*f*-divergence
Wasserstein distances
...

$p_\theta$
Student

Model space

Experiences
*(Exogenous regularization)*

$\mathbb{E}_q[f]$

$q$ Teacher

Uncertainty
*(Self-regularization)* $\mathbb{H}(q)$

$q$ space

**Experiences** $f$

Data

KNOWLEDGE GRAPH

Knowledge

Other Models

4,000,000

Reward

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

# Overview: well-known algorithms/paradigms recovered by SE

| Experience type | Experience function $f$ | Divergence $\mathbb{D}$ | $\alpha$ | $\beta$ | Algorithm |
|---|---|---|---|---|---|
| Data instances | $f_{\text{data}}(\boldsymbol{x};\mathcal{D})$ | CE | 1 | 1 | Unsupervised MLE |
| | $f_{\text{data}}(\boldsymbol{x},\boldsymbol{y};\mathcal{D})$ | CE | 1 | $\epsilon$ | Supervised MLE |
| | $f_{\text{data-self}}(\boldsymbol{x},\boldsymbol{y};\mathcal{D})$ | CE | 1 | $\epsilon$ | Self-supervised MLE |
| | $f_{\text{data-w}}(\boldsymbol{t};\mathcal{D})$ | CE | 1 | $\epsilon$ | Data Re-weighting |
| | $f_{\text{data-aug}}(\boldsymbol{t};\mathcal{D})$ | CE | 1 | $\epsilon$ | Data Augmentation |
| | $f_{\text{active}}(\boldsymbol{x},\boldsymbol{y};\mathcal{D})$ | CE | 1 | $\epsilon$ | Active Learning (Ertekin et al., 2007) |
| Knowledge | $f_{rule}(\boldsymbol{x},\boldsymbol{y})$ | CE | 1 | 1 | Posterior Regularization (Ganchev et al., 2010) |
| | $f_{rule}(\boldsymbol{x},\boldsymbol{y})$ | CE | $\mathbb{R}$ | 1 | Unified EM (Samdani et al., 2012) |
| Reward | $\log Q^{\theta}(\boldsymbol{x},\boldsymbol{y})$ | CE | 1 | 1 | Policy Gradient |
| | $\log Q^{\theta}(\boldsymbol{x},\boldsymbol{y}) + Q^{in,\theta}(\boldsymbol{x},\boldsymbol{y})$ | CE | 1 | 1 | + Intrinsic Reward |
| | $Q^{\theta}(\boldsymbol{x},\boldsymbol{y})$ | CE | $\rho > 0$ | $\rho > 0$ | RL as Inference |
| Model | $f_{\text{model}}^{\text{mimicking}}(\boldsymbol{x},\boldsymbol{y};\mathcal{D})$ | CE | 1 | $\epsilon$ | 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 | $f_{\tau}(\boldsymbol{t})$ | 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:

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

*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 Data Experience -- Supervised MLE

Observe data $\mathcal{D} = \{(\boldsymbol{x}^*, \boldsymbol{y}^*)\}$

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

$$f := f(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D}) = \log \mathbb{E}_{(\boldsymbol{x}^*, \boldsymbol{y}^*) \sim \mathcal{D}}\left[\; \mathbb{1}_{(\boldsymbol{x}^*, \boldsymbol{y}^*)}(\boldsymbol{x}, \boldsymbol{y}) \;\right] \qquad \alpha = 1, \beta = \epsilon$$

Teacher step: $q(\boldsymbol{x}, \boldsymbol{y}) = \exp\left\{\dfrac{\beta \log p_\theta(\boldsymbol{x}, \boldsymbol{y}) + f(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})}{\alpha}\right\}/Z \approx \exp\{f(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})\}/Z = \tilde{p}_d(\boldsymbol{x}, \boldsymbol{y})$

Student step: $\min_\theta \; -\mathbb{E}_q\left[\log p_\theta(\boldsymbol{x}, \boldsymbol{y})\right] \dashrightarrow$ Negative data log-likelihood

# SE with Data Experience -- Unsupervised MLE

Observe data $\mathcal{D} = \{(\boldsymbol{x}^*)\}$

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

$$f := f(\boldsymbol{x}\, ;\, \mathcal{D}) = \log \mathbb{E}_{\boldsymbol{x}^* \sim \mathcal{D}} \left[ \mathbb{1}_{\boldsymbol{x}^*}(\boldsymbol{x}) \right] \qquad \alpha = \beta = 1$$

$$q = q(\boldsymbol{y}|\boldsymbol{x})$$

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

Negative variational lower bound

# SE with "Oracle Data Experience": Active Learning

- Have access to a vast pool of unlabeled data instances
- Can select instances (queries) to be labeled by an oracle (e.g., human)

- Experiences:
  - $u(x)$ measures *informativeness* of an instance $x$
    - e.g., Uncertainty on $x$, measured by Shannon entropy $H(p_\theta(y|x))$
  - Encode instances + oracle labels:

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

# SE and Active Learning

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

$$f := f(\boldsymbol{x},\boldsymbol{y}\,;\,Oracle) + u(\boldsymbol{x}) \qquad\qquad \alpha = 1, \beta = \epsilon$$

○ Teacher step: $q(\boldsymbol{x},\boldsymbol{y}) = \exp\left\{\dfrac{\beta\log p_\theta(\boldsymbol{x},\boldsymbol{y}) + f(\boldsymbol{x}\,,\boldsymbol{y}\,;\,Oracle) + u(\boldsymbol{x})}{\alpha}\right\} / Z$

Equivalent to [e.g., Ertekin et al., 07]:
- Randomly draw a subset $\mathcal{D}_{sub} = \{\boldsymbol{x}^*\}$
- Draw a query $\boldsymbol{x}^*$ from $\mathcal{D}_{sub}$ according to $\exp\{u(\boldsymbol{x})\}$
- Get label $\boldsymbol{y}^*$ for $\boldsymbol{x}^*$ from the oracle
- Maximize log likelihood on $(\boldsymbol{x}^*, \boldsymbol{y}^*)$

○ Student step: $\min_{\theta} \; -\mathbb{E}_q\left[\log p_\theta(\boldsymbol{x},\boldsymbol{y})\right]$

# SE with Reward Experience

Markov Decision
Process (MDP)

**AGENT**

- State $\boldsymbol{x}_t$
- Take action $\boldsymbol{y}_t \sim p_\theta(\boldsymbol{y}_t | \boldsymbol{x}_t)$

**ENVIRONMENT**

- Get reward $r_t = r(\boldsymbol{x}_t, \boldsymbol{y}_t)$
- New state $\boldsymbol{x}_{t+1}$

# SE with Reward Experience: Reinforcement Learning

Markov Decision Process (MDP)

**AGENT**

- State $x_t$
- Take action $y_t \sim p_\theta(y_t | x_t)$

**ENVIRONMENT**

- Get reward $r_t = r(x_t, y_t)$
- New state $x_{t+1}$

- $p_\theta(x, y) = p_\theta(y|x) p_0(x)$, where $p_\theta(y|x)$ is the policy, $p_0(x)$ is the start state distribution

- $Q^\theta(x, y)$ − expected future reward of taking action $y$ in state $x$ and continuing the current policy $p_\theta$

$$Q^\theta(x, y) = \mathbb{E}_{p_\theta} \left[ \sum_{t=0}^{\infty} r_t \mid x_0 = x, y_0 = y \right]$$

- $\mu^\theta(x)$ − state distribution

$$\mu^\theta(x) = \sum_{t=0}^{\infty} p(x_t = x)$$

# SE with **Reward** Experience I: Policy Gradient

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



SCORE: 107

- Policy gradient

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

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

- Student step:

$$\mathbb{E}_{q^{(n)}(\boldsymbol{x}, \boldsymbol{y})} \left[ \nabla_\theta \log p_\theta(\boldsymbol{x}, \boldsymbol{y}) \right] + \mathbb{E}_{q^{(n)}(\boldsymbol{x}, \boldsymbol{y})} \left[ \nabla_\theta f^\theta_{\text{reward},1}(\boldsymbol{x}, \boldsymbol{y}) \right] \Big|_{\boldsymbol{\theta} = \boldsymbol{\theta}^{(n)}}$$

$$= 1/Z \cdot \sum_{\boldsymbol{x}} p_0(\boldsymbol{x}) \nabla_\theta \sum_{\boldsymbol{y}} p_\theta(\boldsymbol{y}|\boldsymbol{x}) Q^\theta(\boldsymbol{x}, \boldsymbol{y}) \Big|_{\boldsymbol{\theta} = \boldsymbol{\theta}^{(n)}} \qquad \textit{(log-derivative trick)}$$

$$= 1/Z \cdot \boxed{\sum_{\boldsymbol{x}} \mu^\theta(\boldsymbol{x}) \sum_{\boldsymbol{y}} Q^\theta(\boldsymbol{x}, \boldsymbol{y}) \nabla_\theta p_\theta(\boldsymbol{y}|\boldsymbol{x}) \Big|_{\boldsymbol{\theta} = \boldsymbol{\theta}^{(n)}}} \qquad \textit{(policy gradient theorem)}$$

*policy gradient*

# SE with **Reward** Experience II: RL as Inference

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



SCORE: 107

- RL-as-inference [Dayan'97; Levine'18, …]

$$f^\theta(x, y) := Q^\theta(x, y) \qquad \alpha = \beta = \rho \ (> 0)$$

$$\min_{q,\theta} -\rho H(q) - \rho \mathbb{E}_q \left[ \log p_\theta(x, y) \right] - \mathbb{E}_{q(x,y)} \left[ Q^\theta(x, y) \right]$$

$$\geq - \log \mathbb{E}_{p_\theta(x,y)} \left[ p(o = 1 \mid x, y) \right]$$



Negative variational lower bound

Define random variable $o \in \{0,1\}$, $p(o = 1) \propto \exp\{ Q^{\theta^t}(x, y)/\rho \}$ (reward excitement fuc. )

# SE with Other Experience

| Experience type | Experience function $f$ | Divergence $\mathbb{D}$ | $\alpha$ | $\beta$ | Algorithm |
|---|---|---|---|---|---|
| Data instances | $f_{\text{data}}(\boldsymbol{x}; \mathcal{D})$ | CE | 1 | 1 | Unsupervised MLE |
| | $f_{\text{data}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})$ | CE | 1 | $\epsilon$ | Supervised MLE |
| | $f_{\text{data-self}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})$ | CE | 1 | $\epsilon$ | Self-supervised MLE |
| | $f_{\text{data-w}}(\boldsymbol{t}; \mathcal{D})$ | CE | 1 | $\epsilon$ | Data Re-weighting |
| | $f_{\text{data-aug}}(\boldsymbol{t}; \mathcal{D})$ | CE | 1 | $\epsilon$ | Data Augmentation |
| | $f_{\text{active}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})$ | CE | 1 | $\epsilon$ | Active Learning (Ertekin et al., 2007) |
| Knowledge | $f_{rule}(\boldsymbol{x}, \boldsymbol{y})$ | CE | 1 | 1 | Posterior Regularization (Ganchev et al., 2010) |
| | $f_{rule}(\boldsymbol{x}, \boldsymbol{y})$ | CE | $\mathbb{R}$ | 1 | Unified EM (Samdani et al., 2012) |
| Reward | $\log Q^{\theta}(\boldsymbol{x}, \boldsymbol{y})$ | CE | 1 | 1 | Policy Gradient |
| | $\log Q^{\theta}(\boldsymbol{x}, \boldsymbol{y}) + Q^{in,\theta}(\boldsymbol{x}, \boldsymbol{y})$ | CE | 1 | 1 | + Intrinsic Reward |
| | $Q^{\theta}(\boldsymbol{x}, \boldsymbol{y})$ | CE | $\rho > 0$ | $\rho > 0$ | RL as Inference |
| Model | $f_{\text{model}}^{\text{mimicking}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})$ | CE | 1 | $\epsilon$ | 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 | $f_{\tau}(\boldsymbol{t})$ | CE | $\rho > 0$ | $\rho > 0$ | Multiplicative Weights (Freund & Schapire, 1997) |

*See paper for more details*

# SE Component: Divergence Function $\mathbb{D}$

We now look at the choices of divergence $\mathbb{D}$:

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

*Divergence*
*(fitness)*
*e.g., Cross Entropy*

# SE with Cross Entropy or KL Divergence

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

$$\min_{q,\theta} - \mathbb{E}_{q(\boldsymbol{x},\boldsymbol{y})}\Big[ f(\boldsymbol{x},\boldsymbol{y}) \Big] + \beta \mathbb{E}_q\Big[ \log p_\theta(\boldsymbol{x},\boldsymbol{y}) \Big] - \alpha \mathbb{H}(q)$$

All the algorithms we've just seen

# SE with Other Divergences

- For notation simplicity, we use $\boldsymbol{x}$ to replace $(\boldsymbol{x}, \boldsymbol{y})$

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

# SE with Other Divergences

- For notation simplicity, we use $\boldsymbol{x}$ to replace $(\boldsymbol{x}, \boldsymbol{y})$

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

- Same as supervised MLE: $f := f(\boldsymbol{x}\,;\,\mathcal{D}), \; \alpha = 1, \; \beta = \epsilon$

- Equivalent to $\min_\theta \mathbb{D}\Big( p_d(\boldsymbol{x}), p_\theta(\boldsymbol{x}) \Big)$

# SE with Other Divergences

- For notation simplicity, we use $\boldsymbol{x}$ to replace $(\boldsymbol{x}, \boldsymbol{y})$

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

- Same as supervised MLE: $f := f(\boldsymbol{x}; \mathcal{D}), \ \alpha = 1, \ \beta = \epsilon$

- Equivalent to $\min_\theta \mathbb{D}\Big( p_d(\boldsymbol{x}), p_\theta(\boldsymbol{x}) \Big)$

- Solve with probability functional descent (PFD) [Chu et al., 2019]
  - $p_\theta(\boldsymbol{x})$ can be optimized by minimizing $\mathbb{E}_{p_\theta}[\Psi(\boldsymbol{x})]$, where $\Psi(\boldsymbol{x})$ is the influence function for $\mathbb{D}$ at $p_{\theta^t}$
  - $\Psi$ is obtained with convex duality
  $$\Psi(\boldsymbol{x}) = \operatorname{argmax}_\psi \mathbb{E}_{p_\theta}[\psi(\boldsymbol{x})] - \mathbb{D}^*(\psi)$$

  Convex conjugate of $\mathbb{D}$

  - So the whole optimization is
  $$\min_\theta \max_\psi \mathbb{E}_{p_\theta}[\psi(\boldsymbol{x})] - \mathbb{D}^*(\psi)$$

# SE with JS Divergence: Generative Adversarial Learning (GANs)

$$\min_{\theta} \mathbb{D}\left(p_d(\boldsymbol{x}), p_{\theta}(\boldsymbol{x})\right)$$

- Solve with probability functional descent (PFD) [Chu et al., 2019]
  - $p_{\theta}(\boldsymbol{x})$ can be optimized by minimizing $\mathbb{E}_{p_{\theta}}[\Psi(\boldsymbol{x})]$, where $\Psi(\boldsymbol{x})$ is the influence function for $\mathbb{D}$ at $p_{\theta^t}$
  - $\Psi$ is obtained with convex duality

    $$\Psi(\boldsymbol{x}) = \mathrm{argmax}_{\psi}\, \mathbb{E}_{p_{\theta}}[\psi(\boldsymbol{x})] - \mathbb{D}^*(\psi)$$

  - So the whole optimization is

    $$\min_{\theta} \max_{\psi} \mathbb{E}_{p_{\theta}}[\psi(\boldsymbol{x})] - \mathbb{D}^*(\psi)$$

Parameterize $\psi$ with an NN $C_{\phi}$.
E.g., when $\mathbb{D}$ is JSD and

$$\psi_{\phi}(x) := 0.5\log\left(1 - C_{\phi}\right) - 0.5\log 2$$

Plugging into the equation recovers vanilla GAN training

Jensen-Shannon Divergence: $\mathrm{JS}(q\|p_{\theta}) = \frac{1}{2}\mathrm{KL}(q\|h) + \frac{1}{2}\mathrm{KL}(p_{\theta}\|h)$

where $h = \frac{1}{2}(q + p_{\theta})$

Petuum

60

# SE with **Wasserstein Distance**: **W-GAN**

$$\min_{\theta} \mathbb{D}\left(p_d(\boldsymbol{x}), p_{\theta}(\boldsymbol{x})\right)$$

- Based on the Kantorovich duality, the 1[st]-order Wasserstein distance between two distributions $q$ and $p$ is written as

$$W_1(q, p) = \sup_{||\psi||_L \leq 1} \mathbb{E}_q[\psi(\boldsymbol{x})] - \mathbb{E}_p(\psi(\boldsymbol{x}))$$

  - where $||\psi||_L \leq 1$ is the constraint of $\psi: \mathcal{X} \to \mathbb{R}$ being a 1-Lipschitz function

- Setting $\mathbb{D}$ to $W_1$ leads to the Wasserstein GAN algorithm [Arjovsky et al., 2017]

$$\min_{\theta} W_1(q, p) = \min_{\theta} \sup_{||\psi||_L \leq 1} \mathbb{E}_{p_d}[\psi(\boldsymbol{x})] - \mathbb{E}_{p_{\theta}}(\psi(\boldsymbol{x}))$$

# Dynamic SE

- So far, we have seen SE as the ultimate learning objective
  - Fully defines the learning problem in an analytical form

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

- In a dynamic or online setting, the learning objective itself may be evolving over time
  - Data instances may follow changing distributions or come from evolving tasks (e.g., lifelong learning)
  - Experience in a strategic game context can involve complex interactions with the target model through co-training or adversarial dynamics

# Dynamic SE

- So far, we have seen SE as the ultimate learning objective
  - Fully defines the learning problem in an analytical form

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

- In a dynamic or online setting, the learning objective itself may be evolving over time
- An extended view of the SE for learning in dynamic contexts
  - SE is a core part of an **outer loop**
  - E.g., consider dynamic experience $f_\tau$ indexed by time $\tau$

for $\tau = 1, 2, \dots$ :

    Acquire experience $f_\tau$,

    Solve SE: $\min_{q,\,\theta} -\mathbb{E}_{q(x,y)}\Big[ f_\tau(x,y) \Big] + \beta \mathbb{D}\Big(q(x,y), p_\theta(x,y)\Big) - \alpha\mathbb{H}(q)$

# Dynamic SE with Adversarial Experience: Variations of GAN

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

- Recall in MLE, $f$ is a fixed function

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

- Intuitively, see $f$ as a similarity metric that measures similarity of sample $\boldsymbol{x}$ against real data $\mathcal{D}$

- Instead of the manually fixed metric, can we learn a metric $f_\phi$?

# Dynamic SE with Adversarial Experience: Variations of GAN

- Augment the standard objective to account for $\phi$:

$$\min_{\theta} \max_{\phi} \min_{q} - \alpha \mathbb{H}(q) + \beta \mathbb{D}\Big( q(\boldsymbol{x}), p_{\theta}(\boldsymbol{x}) \Big) - \mathbb{E}_{q(\boldsymbol{x})}\Big[ f_{\phi}(\boldsymbol{x}) \Big] + \mathbb{E}_{p_d(\boldsymbol{x})}\Big[ f_{\phi}(\boldsymbol{x}) \Big]$$

- Set $\alpha = 0, \beta = 1$. Under mild conditions, the objective recovers:
  - Vanilla GAN [Goodfellow et al., 2014], when $\mathbb{D}$ is JS Divergence and $f_{\phi}$ is a binary classifier
  - $f$-GAN [Nowozin et al., 2016], when $\mathbb{D}$ is $f$-divergence
  - W-GAN [Arjovsky et al., 2017], when $\mathbb{D}$ is Wasserstein distance and $f_{\phi}$ is a 1-Lipschitz function
  - PPO-GAN [Wu et al., 2020], when $\mathbb{D}$ is KL divergence

# Quick recap: well-known algorithms/paradigms recovered by SE

| Experience type | Experience function $f$ | Divergence $\mathbb{D}$ | $\alpha$ | $\beta$ | Algorithm |
|---|---|---|---|---|---|
| Data instances | $f_{\text{data}}(\boldsymbol{x}; \mathcal{D})$ | CE | 1 | 1 | Unsupervised MLE |
| | $f_{\text{data}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})$ | CE | 1 | $\epsilon$ | Supervised MLE |
| | $f_{\text{data-self}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})$ | CE | 1 | $\epsilon$ | Self-supervised MLE |
| | $f_{\text{data-w}}(\boldsymbol{t}; \mathcal{D})$ | CE | 1 | $\epsilon$ | Data Re-weighting |
| | $f_{\text{data-aug}}(\boldsymbol{t}; \mathcal{D})$ | CE | 1 | $\epsilon$ | Data Augmentation |
| | $f_{\text{active}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})$ | CE | 1 | $\epsilon$ | Active Learning (Ertekin et al., 2007) |
| Knowledge | $f_{rule}(\boldsymbol{x}, \boldsymbol{y})$ | CE | 1 | 1 | Posterior Regularization (Ganchev et al., 2010) |
| | $f_{rule}(\boldsymbol{x}, \boldsymbol{y})$ | CE | $\mathbb{R}$ | 1 | Unified EM (Samdani et al., 2012) |
| Reward | $\log Q^\theta(\boldsymbol{x}, \boldsymbol{y})$ | CE | 1 | 1 | Policy Gradient |
| | $\log Q^\theta(\boldsymbol{x}, \boldsymbol{y}) + Q^{in,\theta}(\boldsymbol{x}, \boldsymbol{y})$ | CE | 1 | 1 | + Intrinsic Reward |
| | $Q^\theta(\boldsymbol{x}, \boldsymbol{y})$ | CE | $\rho > 0$ | $\rho > 0$ | RL as Inference |
| Model | $f_{\text{model}}^{\text{mimicking}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})$ | CE | 1 | $\epsilon$ | 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 | $f_\tau(\boldsymbol{t})$ | CE | $\rho > 0$ | $\rho > 0$ | Multiplicative Weights (Freund & Schapire, 1997) |

Paradigms not (yet) covered by SE:
- Meta learning
- Lifelong learning
- …

Interesting future work to study the connections

# Why this is useful?

- Panoramic Learning: learning with ALL experience
  - Experience composition
  - Reuse specialized algorithms -- one runway for different aircrafts

- Complex interaction between experience

- Multi-agent game theoretic learning using all experience

# **Panoramic Learning:** *experience composition*

- Distinct types of experience are all formulated with $f(\boldsymbol{x}, \boldsymbol{y})$
- Combine and plug different $f$ functions into SE to drive learning

$$SE(f, \mathbb{D}, \alpha, \beta)$$

$$f = w_1 \cdot f(\boldsymbol{x} \mid \text{📚}) + w_2 \cdot f(\boldsymbol{x} \mid \text{📖}) + w_3 \cdot f(\boldsymbol{x} \mid \text{💲}) + w_4 \cdot f(\boldsymbol{x} \mid \text{🧑‍🏫}) + \cdots$$

Focus on **what** to use, instead of worrying about **how** to use

# **Panoramic Learning:** *experience composition*
## *Ex.1: Using symbolic knowledge to learn neural networks*

Neural model

Taxt samples   Knowledge bases   Medical KG   ConceptNet (CN)

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

*Hu et al., ACL 2016, "Harnessing Deep Neural Networks with Logic Rules"*
*Hu et al., NeurIPS 2020, "Deep Generative Models with Learnable Knowledge Constraints"*
*Tan et al., EMNLP 2020, "Summarizing Text on Any Aspects: A Knowledge-Informed Weakly-Supervised Approach"*

# **Panoramic Learning:** *experience composition*
## *Ex.2: Using neural networks to "learn" symbolic knowledge*

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

- $\theta$: graph structure to be learned
- $p_\theta$: a simulation model generating medical task samples $(\boldsymbol{x},\boldsymbol{y})$ based on the knowledge graph $\theta$

Measuring likelihood of sample $(\boldsymbol{x},\boldsymbol{y})$ under a trained medical neural model



Commonsense graph



Medical KG

72

# Panoramic Learning: *experience composition*
## *Ex.2: Using neural networks to "learn" symbolic knowledge*

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

| Head entity | Relation | Tail entity | Head entity | Relation | Tail entity |
|---|---|---|---|---|---|
| exercise | prevent | obesity | students | worth celebrating | graduate |
| apple | business | Mac | newborn | can but not good at | sit |
| sleep | prevent | illness | social worker | can help | foster child |
| mall | place for | shopping | honey | ingredient for | honey cake |
| gym | place for | sweat | cabbage | ingredient for | cabbage salad |
| wheat | source of | flour | China | separated by the ocean | Japan |
| oil | source of | fuel | Africa | separated by the ocean | Europe |

Figure 4: Examples of knowledge tuples harvested from ROBERTA-LARGE with MULTI-PROMPTS.

Hao, Tan et al., 2022, "BertNet: Harvesting Knowledge Graphs from Pretrained Language Models"

# **Panoramic Learning:** *experience composition*
## *Ex.2: Using neural networks to "learn" symbolic knowledge*



Figure 2: Precision-recall curve on ConceptNet relations.

Figure 3: Precision-recall curve on LAMA relations.

Hao, Tun et al., 2022, "BertNet: Harvesting Knowledge Graphs from Pretrained Language Models"

# **Panoramic Learning:** *experience composition*
## *Ex.3: Learning prompts to control large pretrained models*

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

**Experiences** *f*

Pretrained LM
(e.g., GPT3)

**SE**

Data instances

$ Reward

Generate a story about cat: once upon a time, ...

prompt *θ*    input    continuation

Deng, Wang, Hsieh et al., EMNLP 2022, " RLPrompt: Optimizing Discrete Text Prompts With Reinforcement Learning"

# **Panoramic Learning:** *experience composition*
## *Ex.4: Learning controllable text generation – more in Lecture#2*

- Combine and plug different $f$ functions into SE to drive learning

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

$$\parallel$$

$$w_1 \cdot f_{data} + w_2 \cdot f_{rules} + w_3 \cdot f_{reward} + \cdots$$

- Enable applications for controllable content generation

Controllable text generation

$f$ = sentiment classifier
        + linguistic rules
        + language model

*Controlling sentiment*

Pos | *The film is full of imagination!*

Neg | *The film is strictly routine!*

[Hu et al., 2017; Yang et al., 2018]

# **Panoramic Learning:** *reusing algorithms*

- Unifying perspective of diverse paradigms (each tailored for a specific type of experience) under SE



- Combining or integrating different experiences
- Re-use or repurpose originally specialized algorithms
  - Systematic idea transfer and solution exchange
  - Solving challenges in one paradigm by applying well-known solutions from another
  - Accelerate innovations across research areas

# **Panoramic Learning:** *reusing algorithms – Ex.1*

- Rules in PR ⇔ Reward in RL
- Empower **reward learning** algo. to **learning rules** [Hu et al., 2018]

| Algorithm | $f$ | $\alpha$ | $\beta$ | $\mathbb{D}$ |
|---|---|---|---|---|
| Unsupervised MLE | $f(\boldsymbol{x}; \mathcal{D})$ | 1 | 1 | CE |
| Supervised MLE | $f(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})$ | 1 | $\epsilon$ | CE |
| Active Learn. | $f(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D}) + u(\boldsymbol{x})$ | temp., $> 0$ | $\epsilon$ | CE |
| Reward-augment MLE | $f_{\text{metric}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D}, r)$ | 1 | $\epsilon$ | CE |
| PG for Seq. Gen. | $f_{\text{metric}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D}, r)$ | 1 | 1 | CE |
| Posterior Reg. | $f_{rule}(\boldsymbol{x}, \boldsymbol{y})$ | weight, $> 0$ | $\alpha$ | CE |
| Unified EM | $f_{rule}(\boldsymbol{x}, \boldsymbol{y})$ | weight, $\in \mathbb{R}$ | 1 | CE |
| Policy Gradient (PG) | $\log Q^{ex}(\boldsymbol{x}, \boldsymbol{y})$ | 1 | 1 | CE |
| + Intrinsic Reward | $\log Q^{ex}(\boldsymbol{x}, \boldsymbol{y}) + Q^{in}(\boldsymbol{x}, \boldsymbol{y})$ | 1 | 1 | CE |
| RL as inference | $Q^{ex}(\boldsymbol{x}, \boldsymbol{y})$ | temp., $> 0$ | $\alpha$ | CE |
| Vanilla GAN | binary classifier | 0 | 1 | JSD |
| $f$-GAN | discriminator | 0 | 1 | $f$-div. |
| WGAN | 1-Lipschitz discriminator | 0 | 1 | W dist. |

# **Panoramic Learning:** *reusing algorithms – Ex.1*

- Rules in PR ⇔ Reward in RL
- Empower **reward learning** algo. to **learning rules** [Hu et al., 2018]

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

MaxEnt inverse RL [Ziebart'08]:
- Parameterize reward $f_\phi(x,y)$ with $\phi$
- Learn $\phi$ with the additional optimization step:

$$\min_{\phi} -\mathbb{E}_{(x^*,y^*)\sim\mathcal{D}}\Big[ \log q_\phi(x^*,y^*) \Big]$$

*Reuse to learn parameterized rules*

Note: $q$ is a function of $f_\phi$, thus $q$ depends on $\phi$

PR with learnable rule constraints $f_\phi(x,y)$:
- E-step to get closed-form $q_\phi$
- M-step to update $p_\theta$
- Reused reward-learning step to update $\phi$

# **Panoramic Learning:** *reusing algorithms – Ex.2*

- Data in supervised MLE ⇔ Reward in RL
- Empower **reward learning** algo. to **learning data augmentation** [Hu et al., 2019]

| Algorithm | $f$ | $\alpha$ | $\beta$ | $\mathbb{D}$ |
|---|---|---|---|---|
| Unsupervised MLE | $f(\boldsymbol{x};\mathcal{D})$ | 1 | 1 | CE |
| Supervised MLE | $f(\boldsymbol{x},\boldsymbol{y};\mathcal{D})$ | 1 | $\epsilon$ | CE |
| Active Learn. | $f(\boldsymbol{x},\boldsymbol{y};\mathcal{D}) + u(\boldsymbol{x})$ | temp., $> 0$ | $\epsilon$ | CE |
| Reward-augment MLE | $f_{\text{metric}}(\boldsymbol{x},\boldsymbol{y};\mathcal{D},r)$ | 1 | $\epsilon$ | CE |
| PG for Seq. Gen. | $f_{\text{metric}}(\boldsymbol{x},\boldsymbol{y};\mathcal{D},r)$ | 1 | 1 | CE |
| Posterior Reg. | $f_{rule}(\boldsymbol{x},\boldsymbol{y})$ | weight, $> 0$ | $\alpha$ | CE |
| Unified EM | $f_{rule}(\boldsymbol{x},\boldsymbol{y})$ | weight, $\in \mathbb{R}$ | 1 | CE |
| Policy Gradient (PG) | $\log Q^{ex}(\boldsymbol{x},\boldsymbol{y})$ | 1 | 1 | CE |
| + Intrinsic Reward | $\log Q^{ex}(\boldsymbol{x},\boldsymbol{y}) + Q^{in}(\boldsymbol{x},\boldsymbol{y})$ | 1 | 1 | CE |
| RL as inference | $Q^{ex}(\boldsymbol{x},\boldsymbol{y})$ | temp., $> 0$ | $\alpha$ | CE |
| Vanilla GAN | binary classifier | 0 | 1 | JSD |
| $f$-GAN | discriminator | 0 | 1 | $f$-div. |
| WGAN | 1-Lipschitz discriminator | 0 | 1 | W dist. |

# **Panoramic Learning:** *reusing algorithms – Ex.2*

- Data in supervised MLE ⇔ Reward in RL

- Empower **reward learning** algo. to **learning data augmentation** [Hu et al., 2019]

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

Intrinsic reward learning [Zheng et al.,08]:

- Reward $f_\phi = f^{ex} + f_\phi^{in}$

- I.e., parameterize (intrinsic) reward $f_\phi^{in}$ with $\phi$

- Learn $\phi$ with the additional optimization step:

$$\min_{\phi} \mathcal{L}_{SE}(\theta^t(\phi))$$

*Reuse to learn parameterized data augmentation model* →

MLE with learnable data augmentation $f_\phi(\boldsymbol{x},\boldsymbol{y})$:

- E-step to get closed-form $q_\phi$
- M-step to update $p_\theta$
- Reused reward-learning step to update $\phi$

Same form as standard equation

Note: updates of $\theta$ depend on experience $f_\phi$, thus the resulting $\theta^t$ is a function of $\phi$

# Panoramic Learning: *reusing algorithms – Ex.3*

- GANs ⇔ RL ⇔ VI
- Empower **RL/VI** algo. (e.g., PPO) to **stabilize GAN training** [Wu et al., 2020]

| Algorithm | $f$ | $\alpha$ | $\beta$ | $\mathbb{D}$ |
|---|---|---|---|---|
| Unsupervised MLE | $f(\boldsymbol{x}; \mathcal{D})$ | 1 | 1 | CE |
| Supervised MLE | $f(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D})$ | 1 | $\epsilon$ | CE |
| Active Learn. | $f(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D}) + u(\boldsymbol{x})$ | temp., $> 0$ | $\epsilon$ | CE |
| Reward-augment MLE | $f_{\text{metric}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D}, r)$ | 1 | $\epsilon$ | CE |
| PG for Seq. Gen. | $f_{\text{metric}}(\boldsymbol{x}, \boldsymbol{y}; \mathcal{D}, r)$ | 1 | 1 | CE |
| Posterior Reg. | $f_{rule}(\boldsymbol{x}, \boldsymbol{y})$ | weight, $> 0$ | $\alpha$ | CE |
| Unified EM | $f_{rule}(\boldsymbol{x}, \boldsymbol{y})$ | weight, $\in \mathbb{R}$ | 1 | CE |
| Policy Gradient (PG) | $\log Q^{ex}(\boldsymbol{x}, \boldsymbol{y})$ | 1 | 1 | CE |
| + Intrinsic Reward | $\log Q^{ex}(\boldsymbol{x}, \boldsymbol{y}) + Q^{in}(\boldsymbol{x}, \boldsymbol{y})$ | 1 | 1 | CE |
| RL as inference | $Q^{ex}(\boldsymbol{x}, \boldsymbol{y})$ | temp., $> 0$ | $\alpha$ | CE |
| Vanilla GAN | binary classifier | 0 | 1 | JSD |
| $f$-GAN | discriminator | 0 | 1 | $f$-divg. |
| WGAN | 1-Lipschitz discriminator | 0 | 1 | W dist. |

# **Panoramic Learning:** *reusing algorithms – Ex.3*

- GANs ⇔ RL ⇔ VI
- Empower **RL/VI** algo. (e.g., PPO) to **stabilize GAN training** [Wu et al., 2020]



(a) Re-use PPO objective for GAN training: discourage excessively large updates by "trapping" the update size around 1

(b) Re-use importance weighting in a VI perspective: greatly reduced variance in both generator and discriminator losses

Improved performance on a range of problems, including image generation, text generation, and text style transfer

# Summary so far …

- The standard equation of objective

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

- Experience function $f$ can encode different types of experience
  - Data instances, constraints, informativeness, reward, adversary models, …

- Enable panoramic learning with ALL experience
  - Re-use or repurpose originally specialized algorithms to other contexts
  - Experience compositonality

# Toward A "Standard Model" of ML

- Loss

- Experience

- Optimization solver

- Model architecture

$$\min_{\theta} \mathcal{L}(\theta, \mathcal{E})$$

Optimization solver  Loss  Model architecture  Experience

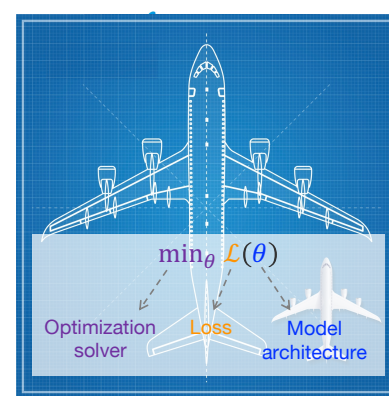# The zoo of optimization solvers



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

> Optimization of the loss, subject to $q \in \mathcal{P}_{\mathrm{prob}}$.
> Convex to $q$ when $\alpha, \beta > 0$ and $\mathbb{D}$ is convex

- Like the Standard Equation as a *master objective* for many paradigms, is there a *master solver* for optimization of loss?

- No (yet) such a general algorithm

- Alternating Projection:
  - Most widely used
  - EM, Variational EM (Variational inference), Wake-Sleep, …

# The Teacher-Student Mechanism

$$\min_{q, \theta} - \alpha \mathbb{H}(q) + \beta \mathbb{D} \Big( q(\boldsymbol{x}), p_\theta(\boldsymbol{x}) \Big) - \mathbb{E}_{q(\boldsymbol{x})} \Big[ f(\boldsymbol{x}\,;\,.) \Big]$$

when $\alpha, \beta > 0$ and $\mathbb{D} = \mathrm{CE}$

(1) Teacher step:

$$q(\boldsymbol{x}) = \exp \left\{ \frac{\beta \log p_\theta(\boldsymbol{x}) + f(\boldsymbol{x}\,;\,.)}{\alpha} \right\} / Z$$

(2) Student step:

$$\min_\theta \mathbb{E}_{q(\boldsymbol{x})} \Big[ \log p_\theta(\boldsymbol{x}) \Big]$$

*Generalization of the classic Variational EM*

- Generalized **E-step**
  Support all types of experience

- *M-step*

# The Teacher-Student Mechanism

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

when $\alpha, \beta > 0$ and $\mathbb{D} = \text{CE}$

(1) Teacher step:

$$q(\boldsymbol{x}) = \exp\left\{\frac{\beta\log p_\theta(\boldsymbol{x}) + f(\boldsymbol{x};\,.\,)}{\alpha}\right\} / Z$$

(2) Student step:

$$\min_\theta \mathbb{E}_{q(\boldsymbol{x})}\left[\log p_\theta(\boldsymbol{x})\right]$$



Model Manifold

Experience Manifold

$p_{\theta^{(1)}}$  $q^{(1)}$

$p_{\theta^{(0)}}$  $q^{(0)}$

EM:
E-step
M-step

Teacher-student Mechanism:
Teacher step
Student step

# Some "advanced" (specialized) techniques

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

- Alternating Projection:
  - EM, Variational EM (Variational inference), Wake-Sleep, …
  - SGD, Back-propagation (BP)

- Convex duality, Lagrangian  --  Kernel Tricks

- Integer linear programming (ILP)

- Probability functional descent (PFD) [Chu et al., 2019] -- Influence function, gives a neat formulation of GAN-like optimization and a few others

89

# I: Duality

- Structured MaxEnt Discrimination (SMED) [Zhu and Xing, 2013]:

$$\min_{q,\,\xi \geq 0} \; -\alpha H(q) - \beta \mathbb{E}_q \left[ \; \log p(\boldsymbol{\theta}) \; \right] + U(\boldsymbol{\xi})$$

$$s.t. \; -\mathbb{E}_q \left[ \Delta F_i(\boldsymbol{y}; \boldsymbol{\theta}) - \Delta \ell_i(\boldsymbol{y}) \right] \leq \xi_i \quad \forall i$$

  ○ Solve the (primal) Lagrangian:

$$q(\boldsymbol{\theta}) = \exp \left\{ \frac{\beta \log p(\boldsymbol{\theta}) + \sum_{i,\boldsymbol{y} \neq \boldsymbol{y}_i^*} \lambda_i(\boldsymbol{y})(\Delta F_i(\boldsymbol{y}; \boldsymbol{\theta}) - \Delta \ell_i(\boldsymbol{y}))}{\alpha} \right\} / Z(\boldsymbol{\lambda})$$

  ○ Solve Lagrangian multipliers $\boldsymbol{\lambda}$ from the **dual problem** (when $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|0, I)$; $U(\xi) = \sum \xi_i,$)

$$\max_{\boldsymbol{\lambda} \geq 0,\, \sum \lambda_i = 1} \sum_{i,\boldsymbol{y} \neq \boldsymbol{y}_i^*} \lambda_i(\boldsymbol{y}) \Delta \ell_i(\boldsymbol{y}) - \frac{1}{2} \left| \sum_{i,\boldsymbol{y} \neq \boldsymbol{y}_i^*} \lambda_i(\boldsymbol{y}) \Delta T_i(\boldsymbol{y}) \right|^2$$

Allows kernel trick for nonlinear interactions b/w experiences

# II: Influence Function and Probability Functional Descent

- Gradient descent in the space of probability measures $\mathcal{P}(X)$

$$\min_{p \in \mathcal{P}(X)} \mathcal{I}(p) \qquad \mathcal{I}: \mathcal{P}(X) \to \mathbb{R} : \text{a probability functional}$$

- Influence function $\Psi_p(x)$:

Gateaux differential of $\mathcal{I}$ at $p$ in the direction $\chi = q - p$

$$d\mathcal{I}_p(\chi) = \int_X \Psi_p(x)\chi(dx)$$
$$= \mathbb{E}_q\big[\Psi_p(x)\big] - \mathbb{E}_p\big[\Psi_p(x)\big]$$

- With a linear approximation $\tilde{\mathcal{I}}(p)$ to $\mathcal{I}(p)$ around $p_0$:

$$\tilde{\mathcal{I}}(p) = \mathcal{I}(p_0) + d\mathcal{I}_{p_t}(p - p_0).$$
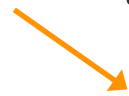$$= \mathbb{E}_{x \sim p}\big[\Psi_{p_0}(x)\big] + const.$$

- Thus, once we obtain the influence function, we can optimize $p$ by decreasing $\mathbb{E}_{x \sim p}\big[\Psi_{p_0}(x)\big]$

[Chu et al., 2019]

# Adversarial learning using PFD

$$\mathcal{I}(p_\theta) = \mathbb{D}\left( p_d(\boldsymbol{x}), p_\theta(\boldsymbol{x}) \right)$$

- Often no closed-form influence function, e.g., when $\mathbb{D}$ is JSD or W-distance

- Approximate with convex duality:

  - Convex conjugate $\mathcal{I}^*(\psi) = \sup_u \int_x \psi(\boldsymbol{x}) u(dx) - \mathcal{I}(u)$

  - Influence function is obtained via $\Psi_{p_\theta}(x) = \text{argmax}_\psi \, \mathbb{E}_{\boldsymbol{x} \sim p_\theta}[\psi(\boldsymbol{x})] - \mathcal{I}^*(\psi)$

  - Parameterize $\psi$ as below to recover optimization of generator and discriminator

    $\psi_{\boldsymbol{\phi}}(\boldsymbol{x}) := 0.5 \log\left(1 - C_\phi\right) - 0.5 \log 2$

    $\Psi_{\boldsymbol{JS}} = \text{argmax}_\phi \, \mathbb{E}_{p_{data}}\big[\log C_\phi\big] - \mathbb{E}_{p_\theta}\big[\log\left(1 - C_\phi\right)\big]$

- The whole optimization of $\mathcal{I}(p)$ is thus

  $$\min_\theta \max_\phi \mathbb{E}_{p_{data}}\big[\log C_\phi\big] - \mathbb{E}_{p_\theta}\big[\log\left(1 - C_\phi\right)\big]$$

[Chu et al., 2019]

# Other popular algorithms in the PFD view

- PFD recovers optimization procedures in some popular algorithms

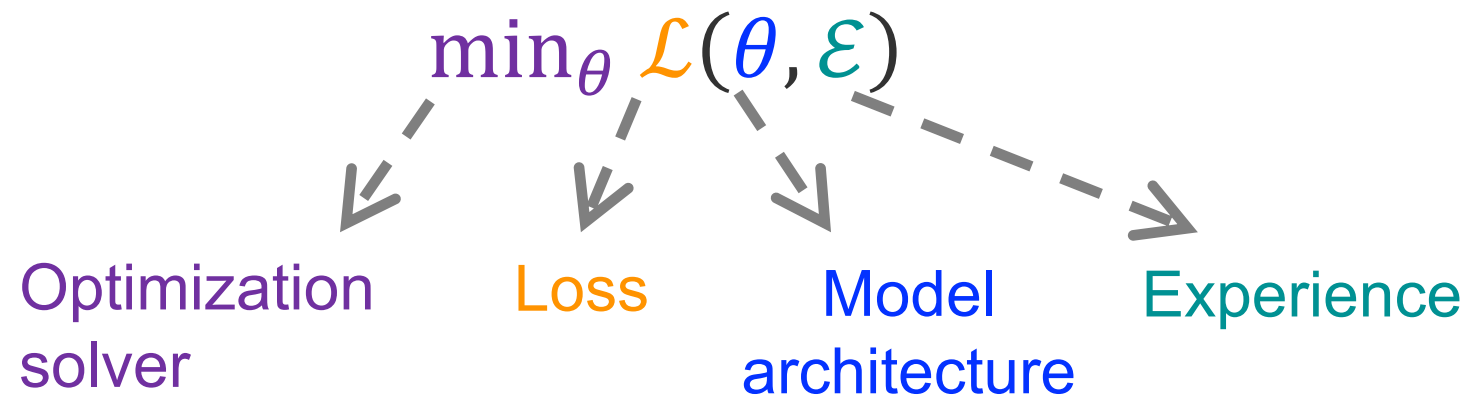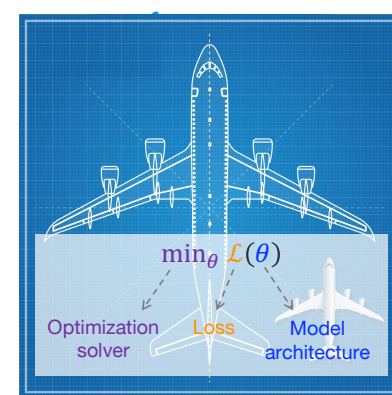| Algorithm | Type of derivative estimator |
|---|---|
| **Generative adversarial networks** | |
| Minimax GAN (Goodfellow et al., 2014) | Convex duality |
| Non-saturating GAN (Goodfellow et al., 2014) | Binary classification |
| Wasserstein GAN (Arjovsky et al., 2017) | Convex duality |
| **Variational inference** | |
| Black-box variational inference (Ranganath et al., 2014) | Exact |
| Adversarial variational Bayes (Mescheder et al., 2017) | Binary classification |
| Adversarial posterior distillation (Wang et al., 2018) | Convex duality |
| **Reinforcement learning** | |
| Policy iteration (Howard, 1960) | Exact |
| Policy gradient (Williams, 1992) | Monte Carlo |
| Actor-critic (Konda & Tsitsiklis, 2000; Sutton et al., 2000) | Least squares |
| Dual actor-critic (Chen & Wang, 2016; Dai et al., 2017b) | Convex duality |

Estimation of the influence function

[Chu et al., 2019]

# Toward A "Standard Model" of ML

- Loss

- Experience

- Optimization solver

- Model architecture

$$\min_\theta \mathcal{L}(\theta, \mathcal{E})$$

Optimization solver     Loss     Model architecture     Experience

# **Model architecture** – *more in Lecture#2*

- Relatively well explored:
  - Neural network design
  - Graphical model design
  - Compositional architectures

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

*Next lecture: a composable catalog of building blocks*
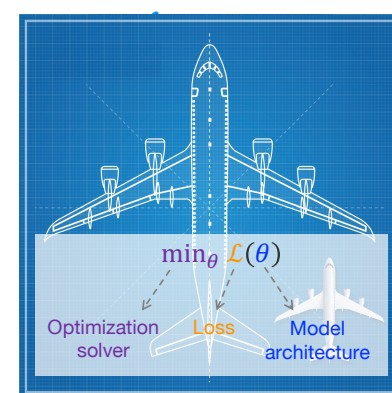
# Summary: A "Standard Model" of ML



- Loss + experience
  - Standard Equation (SE)

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

- Optimization solver

  - The extended EM algorithm gives a general primal solution in many cases

  - PFD gives a neat formulation for some cases (e.g., GANs)

- Model architecture: vast libraries of building blocks → compositionality

*Next: practical implications of the ML "Standard Model"*

# Schedule

- **Lecture#1:** Theory: The Standard Model of ML

  A blueprint of ML paradigms for ALL experience

  *(Jan 19 Thursday, 4:45pm-6:15pm UK Time)*



- **Lecture#2:** Tooling: Operationalizing The Standard Model

  Compose your ML solutions like playing Lego

  *(Jan 20 Thursday, 1:00pm-2:30pm)*



- **Lecture#3:** Computing: Modern infrastructure for productive ML

  Automatic tuning, distributing, and scheduling

  *(Jan 20 Thursday, 4:45pm-6:15pm)*