

# DSC190: Machine Learning with Few Labels

A “standardized” view of ML

**Zhiting Hu**

Lecture 9, October 21, 2021

**UC San Diego**

**HALICIOĞLU DATA SCIENCE INSTITUTE**

# Logistics

- HW1 due extended to Sunday (10/24)
- HW2 out: (much) easier than HW1 !
- Mid-term survey

# Outline

- Functional derivative
- A “standardized” view of ML

# Functional derivative

- $\nabla_q - \mathbb{H}(q) = \log q + 1$
- Functional  $F(y)$ : an operator that takes a function  $y(x)$  and returns an output value  $F$
- Functional derivative (aka, variational derivative) relates a change in a Functional  $F(y)$  to a change in the function  $y$

# Functional derivative

- Recall the conventional derivative  $\frac{dy}{dx}$ 
  - Taylor expansion

$$y(x + \epsilon) = y(x) + \frac{dy}{dx}\epsilon + O(\epsilon^2)$$

- Functional derivative
  - How much a functional  $F[y]$  changes when we make a small change  $\epsilon\eta(x)$  to the function  $y(x)$

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2)$$

- A function  $y(x)$  that maximizes (or minimizes) a functional  $F[y]$  must satisfy

$$\frac{\delta F}{\delta y(x)} = 0 \text{ for all } x$$

# Functional derivative

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2)$$

- Consider a functional that is defined by an integral over a function  $G(y, x)$

$$F[y] = \int G(y, x) dx$$

- Consider variations in the function  $y(x)$ ,

$$F[y + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\partial G}{\partial y} \eta(x) dx + O(\epsilon^2)$$

# Functional derivative

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2)$$

- Consider a functional that is defined by an integral over a function  $G(y, x)$

$$F[y] = \int G(y, x) dx$$

- Consider variations in the function  $y(x)$ ,

$$F[y + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\partial G}{\partial y} \eta(x) dx + O(\epsilon^2)$$

- Ex.1,  $-\mathbb{H}(q) = \int q(x) \log q(x) dx$ 
  - $G = q(x) \log q(x)$
- EX.2, posterior regularization

## Ex.2: Posterior Regularization

assume single data point  $\mathbf{x}^*$

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

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

$$\xi \geq 0,$$

- *Lagrangian*

$$\begin{aligned} \max_{\mu \geq 0, \eta \geq 0, \alpha \geq 0} \min_{q, \xi} & -\mathbf{H}(q(\mathbf{z})) - \mathbb{E}_{q(\mathbf{z})} [\log p(\mathbf{x}^* | \mathbf{z}) \pi(\mathbf{z})] \\ & + \sum_i (1 - \mu_i) \xi_i + \sum_i \eta_i (\mathbb{E}_q [T_i(\mathbf{x}^*; \mathbf{z})] - \xi_i) + \alpha \left( \sum_z q(\mathbf{z}) - 1 \right) \end{aligned}$$



# A “Standardized” View of ML

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
- Posterior regularization:
  - Constrained Bayesian inference => constrained optimization

The general expression as a constrained optimization:

## MaxEnt perspective

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

*(auxiliary) distribution  $q$*

- Supervised MLE and maximum entropy

$$\min_{q(x, y)} H(q)$$

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

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

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

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

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

The 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  $q$*  (arrow pointing to  $q$ )  
*loss* (arrow pointing to  $\mathcal{L}(q, \theta)$ )  
*constrained set* (arrow pointing to  $\mathcal{Q}$ )

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

$$\begin{aligned} \min_{q, \xi} \quad & -\mathbf{H}(q(\mathbf{z})) - \mathbb{E}_{q(\mathbf{z})} \left[ \sum_{\mathbf{x}^* \in \mathcal{D}} \log p(\mathbf{x}^* | \mathbf{z}) \pi(\mathbf{z}) \right] + U(\xi) \\ \text{s.t.} \quad & q(\mathbf{z}) \in \mathcal{Q}(\xi) \\ & \xi \geq 0, \end{aligned}$$

# The Standard Equation (SE)

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

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

- $q(t)$ : auxiliary distribution

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

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

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



# The Standard Equation (SE)

- Let  $t$  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)

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

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

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

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

3 terms:

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



Uncertainty

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

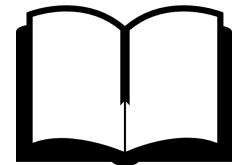


Teacher  
 $q(\mathbf{t})$

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

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

Textbook  
 $f(\mathbf{t})$



# The Standard Equation (SE)

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

- The introduction of the auxiliary distribution  $q$  relaxes the learning problem of  $p_{\theta}$ , originally only over  $\theta$ , to be now alternating between  $q$  and  $\theta$ 
  - Recall in EM, we introduced  $q$  to deal with the intractable marginal log-likelihood
- $q$  acts as a conduit between the exogenous experience and the target model
  - subsumes the experience, by maximizing the expected  $f$  value
  - passes it incrementally to the target model, by minimizing the divergence  $\mathbb{D}$
- E.g., assume  $\mathbb{D}$  is cross entropy, and  $\mathbb{H}$  is Shannon entropy
  - The above optimization, at each iteration  $n$ :

$$q^{(n+1)}(\mathbf{t}) = \exp\left\{\frac{\beta \log p_{\theta^{(n)}}(\mathbf{t}) + f(\mathbf{t})}{\alpha}\right\} / Z$$

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

# The Standard Equation (SE)

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

- The introduction of the auxiliary distribution  $q$  relaxes the learning problem of  $p_{\theta}$ , originally only over  $\theta$ , to be now alternating between  $q$  and  $\theta$ 
  - Recall in EM, we introduced  $q$  to deal with the intractable marginal log-likelihood
- $q$  acts as a conduit between the exogenous experience and the target model
  - subsumes the experience, by maximizing the expected  $f$  value
  - passes it incrementally to the target model, by minimizing the divergence  $\mathbb{D}$
- E.g., assume  $\mathbb{D}$  is cross entropy, and  $\mathbb{H}$  is Shannon entropy
  - The above optimization, at each iteration  $n$ :

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

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

# The Standard Equation (SE)

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

- 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 $\mathbb{D}$ | $\alpha$     | $\beta$    | 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            | $\epsilon$ | Supervised MLE                                  |
|                 | $f_{\text{data-self}}(\mathbf{x}, \mathbf{y}; \mathcal{D})$                             | CE                      | 1            | $\epsilon$ | Self-supervised MLE                             |
|                 | $f_{\text{data-w}}(\mathbf{t}; \mathcal{D})$  | CE                      | 1            | $\epsilon$ | Data re-weighting                               |
|                 | $f_{\text{data-aug}}(\mathbf{t}; \mathcal{D})$  | CE                      | 1            | $\epsilon$ | Data Augmentation                               |
|                 | $f_{\text{active}}(\mathbf{x}, \mathbf{y}; \mathcal{D})$                                | CE                      | 1            | $\epsilon$ | 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                      | $\tau > 0$   | $\tau > 0$ | RL as Inference                                 |
| Other advanced  | binary classifier   | JSD                     | 0            | 1          | Vanilla GAN (Goodfellow et al., 2014)           |
|                 | discriminator   | f-divg.                 | 0            | 1          | f-GAN (Nowozin et al., 2016)                    |
|                 | 1-Lipschitz discriminator   | $W_1$ dist.             | 0            | 1          | WGAN (Arjovsky et al., 2017)                    |
|                 | 1-Lipschitz discriminator   | KL                      | 0            | 1          | PPO-GAN (Wu et al., 2020)                       |

# 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

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

The expected similarity between  $(\mathbf{x}, \mathbf{y})$  and observed data  $(\mathbf{x}^*, \mathbf{y}^*)$ , with similarity measure  $\mathbb{1}_a(b)$ , i.e., an indicator function (1 if  $a=b$ , 0 otherwise)

# 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

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

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

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

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 actively supervised experience

- Have access to a vast pool of unlabeled data instances
- Can select instances (queries) to be labeled by an oracle (e.g., human)
  
- Experiences:
  - $u(\mathbf{x})$  measures *informativeness* of an instance  $\mathbf{x}$ 
    - e.g., Uncertainty on  $\mathbf{x}$ , measured by predictive entropy
  - Instances + oracle labels:

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

# SE with actively supervised experience

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

$$f := f(\mathbf{x}, \mathbf{y}; \text{Oracle}) + u(\mathbf{x})$$

$$\alpha = 1, \beta = \epsilon$$



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

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

Equivalent to **active learning** [e.g., Ertekin et al., 07]:

- Randomly draw a subset  $\mathcal{D}_{sub} = \{\mathbf{x}^*\}$
- Draw a query  $\mathbf{x}^*$  from  $\mathcal{D}_{sub}$  according to  $\exp\{u(\mathbf{x})\}$
- Get label  $\mathbf{y}^*$  for  $\mathbf{x}^*$  from the oracle
- Maximize log likelihood on  $(\mathbf{x}^*, \mathbf{y}^*)$

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

# 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



- Functional derivative

$$F[y(x) + \epsilon \eta(x)] = F[y(x)] + \epsilon \int \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2)$$

Questions?