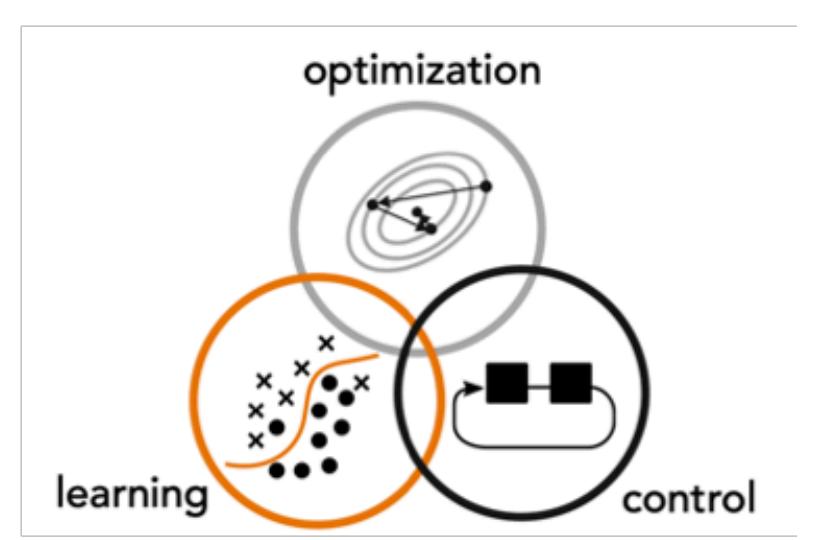


# Regularization for Adversarial Robust Learning

Jie Wang †, Rui Gao ‡, Yao Xie †

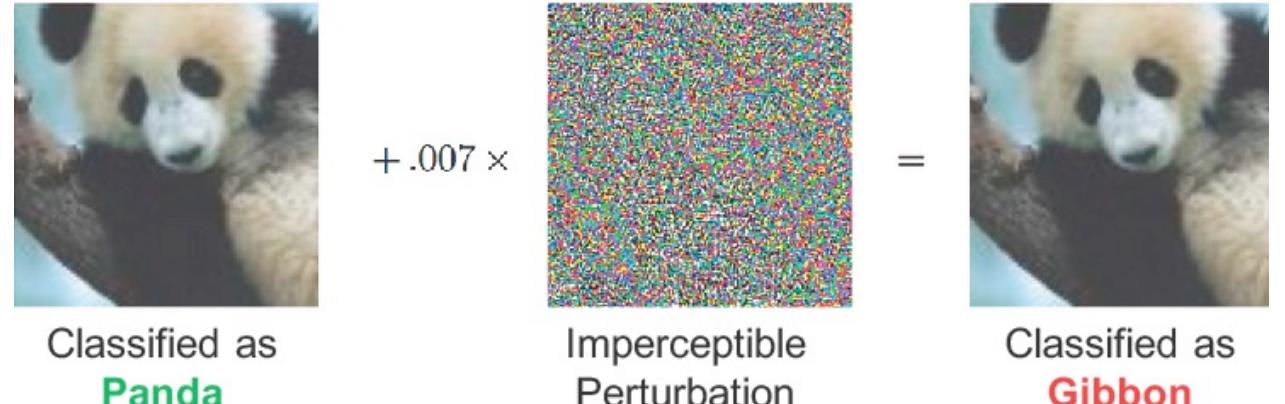
† H. Milton Stewart School of Industrial and Systems Engineering, Georgia Institute of Technology



‡ Department of Information, Risk, and Operations Management, McCombs School of Business, University of Texas at Austin

## Motivation and Background

In Picture



In Text

A **inspire** movie. It is entire of feelings and **wonderful** behaving. I could have sat through it a seconds period.  
A **touch** movie. It is entire of feelings and **good** behaving. I could have sat through it a seconds period.

**Positive**

**Negative**

## Adversarial Training (Goodfellow et al, 2014)

$$\min_{\theta \in \Theta} \left\{ \mathbb{E}_{z \sim \mathbb{P}_n} \left[ \sup_{z' : \|z - z'\| \leq \rho} \ell(z'; \theta) \right] \right\}$$

- $z \sim \mathbb{P}_n$ : Data (e.g., feature-label pair)
- $\|z - z'\| \leq \rho$ : Perturbation Constraints
- $\ell(z'; \theta)$ : Loss Function

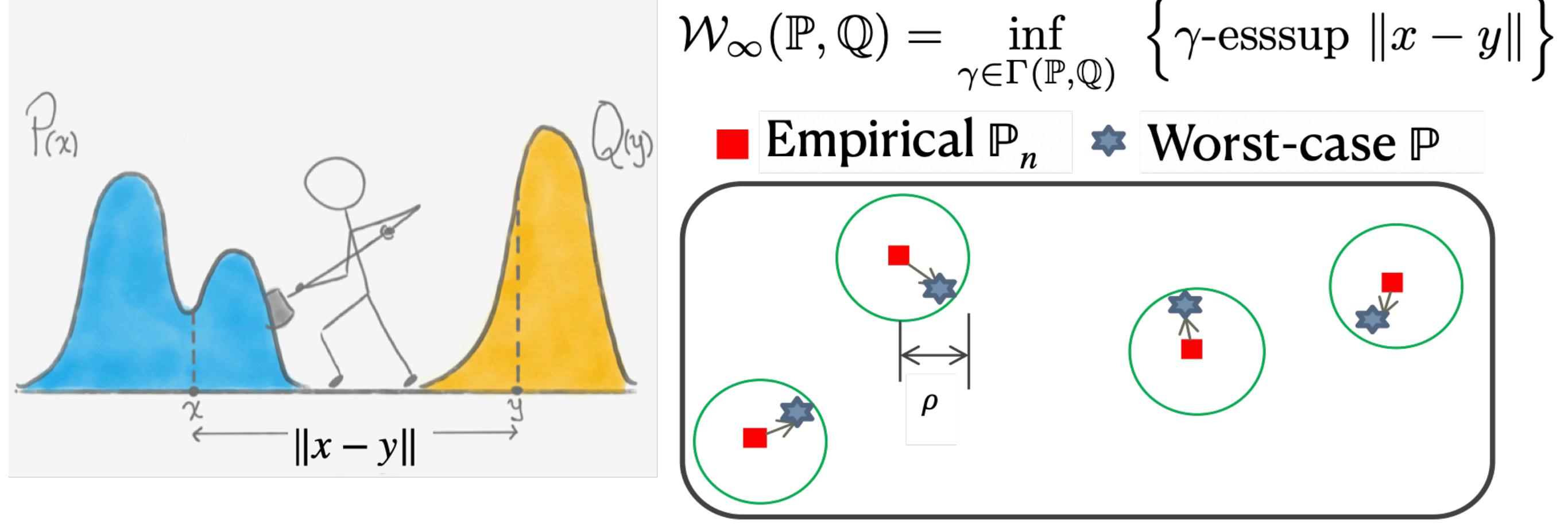
• Cons: Inner supremum is generally **nonconcave** in  $z$

• Heuristics: Fast but **no theoretical guarantees of robustness**

$$\Delta(\theta; z) := \arg \max_{\|\Delta\| \leq \rho} \left\{ \ell(z; \theta) + \nabla_z \ell(z; \theta)^\top \Delta \right\}, \quad \text{then perturb } z \leftarrow z + \Delta(\theta; z).$$

## Wasserstein Distributionally Robust Optimization (DRO)

$$\begin{aligned} (\text{AT}) &= \min_{\theta \in \Theta} \left\{ \sup_{\mathbb{P}: \mathcal{W}_\infty(\mathbb{P}, \mathbb{P}_n) \leq \rho} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] \right\} \\ &= \min_{\theta \in \Theta} \left\{ \sup_{\mathbb{P}, \gamma} \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] : \begin{array}{l} \gamma \in \Gamma(\mathbb{P}_n, \mathbb{P}) \\ \gamma\text{-esssup } \|\zeta_1 - \zeta_2\| \leq \rho \end{array} \right\} \end{aligned}$$



## Proposed Formulation

$$\min_{\theta \in \Theta} \sup_{\mathbb{P}, \gamma} \left\{ \mathbb{E}_{z \sim \mathbb{P}} [\ell(z; \theta)] - \eta \mathbb{D}_f(\gamma, \gamma_0) : \begin{array}{l} \gamma \in \Gamma(\mathbb{P}_n, \mathbb{P}) \\ \gamma\text{-esssup } \|x - y\| \leq \rho \end{array} \right\} \triangleq V_{\text{Primal}}$$

•  $f$ -divergence (e.g., KL or  $\chi^2$ ):  $\mathbb{D}_f(\gamma, \gamma_0) = \int f\left(\frac{d\gamma}{d\gamma_0}\right) d\gamma_0$

• Reference transport  $\gamma_0$  is uniform: For each  $z \in \text{supp}(\mathbb{P}_n)$ ,

$\mathbb{Q}_z(\cdot) \triangleq \gamma_0(\cdot | z)$  is an uniform measure on  $\mathbb{B}_\rho(z)$ .

• Under mild conditions,  $V_{\text{Primal}} = V_{\text{dual}}$ :

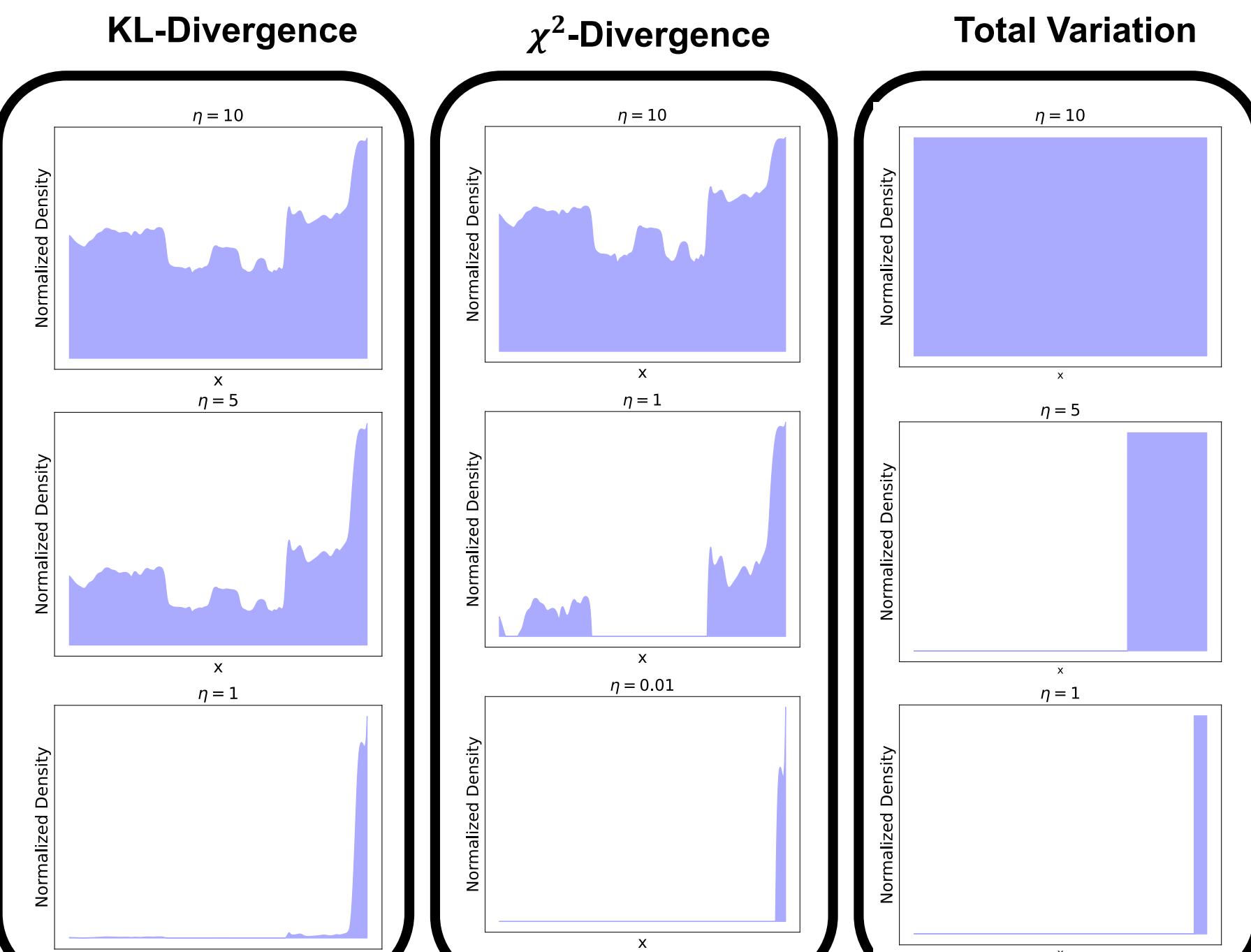
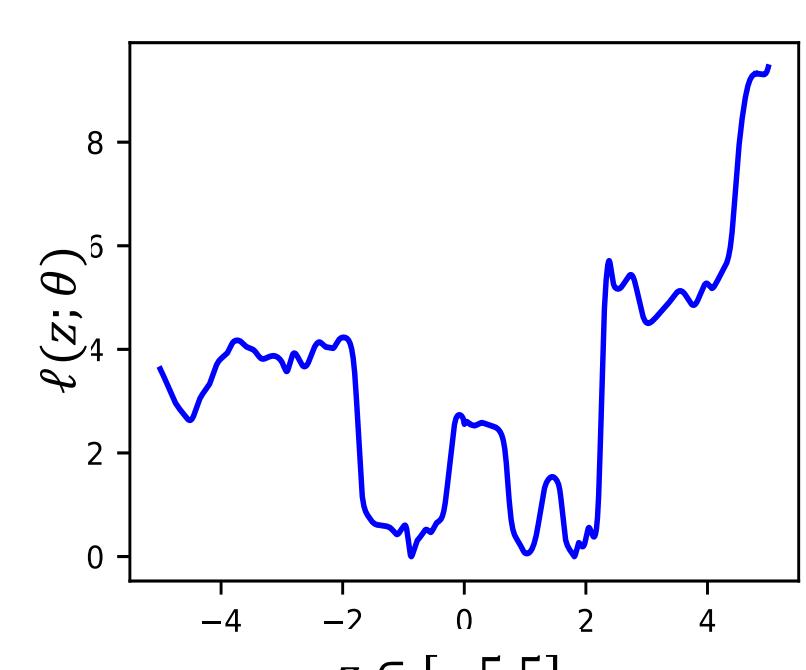
$$V_{\text{Dual}} = \mathbb{E}_{\mathbb{P}_n} \left[ \inf_{\mu \in \mathbb{R}} \left\{ \mu + \mathbb{E}_{z' \sim \mathbb{Q}_z} [(\eta f)^*(\ell(z'; \theta) - \mu)] \right\} \right]$$

Penalized  $f$ -divergence DRO with reference measure  $\mathbb{Q}_z$

$$\frac{d\mathbb{P}^*(\omega)}{d\omega} = \mathbb{E}_{z \sim \mathbb{P}_n} \left[ \alpha_z \cdot \mathbf{1}\{\|\omega - z\| \leq \rho\} \cdot (\eta f)^*(\ell(\omega; \theta) - \mu_z^*) \right]$$

- Normalizing Constant
- Support Constraint
- Density contributed by  $z$

- Visualization of worst-case distribution
- $\ell(z; \theta)$  is a feed-forward neural network



## Optimization Algorithm

• (Regularized-AT) as a bilevel optimization with many lower level constraints:

$$\min_{\theta \in \Theta} \left\{ F(\theta) = \mathbb{E}_{z \sim \mathbb{P}_n} [R(\theta; z)] \right\}$$

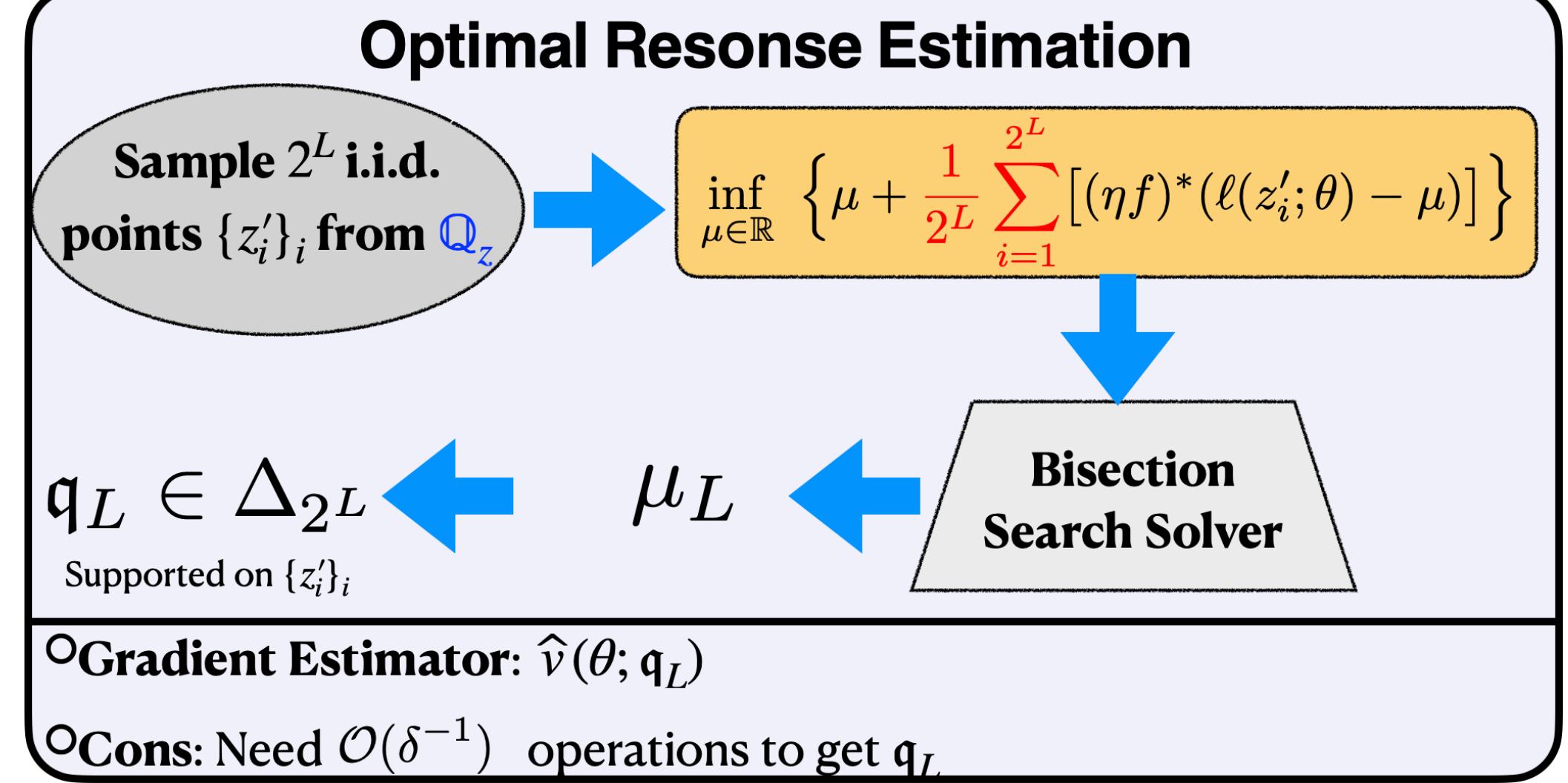
$$\text{where } R(\theta; z) = \inf_{\mu \in \mathbb{R}} \left\{ \mu + \mathbb{E}_{z' \sim \mathbb{Q}_z} [(\eta f)^*(\ell(z'; \theta) - \mu)] \right\}, \quad \forall z \\ = \sup_{q \in \mathcal{P}} \left[ \mathbb{E}_{z' \sim q} [\ell(z'; \theta)] - \eta \mathbb{E}_{z' \sim \mathbb{Q}_z} \left[ \phi \left( \frac{dq(z')}{d\mathbb{Q}_z(z')} \right) \right] \right], \quad \forall z$$

• For iterations  $t = 0, 1, \dots, T-1$ :

- (i) Estimate biased (sub-)gradient; (ii) Perform projected gradient descent.

Algorithm	Naive Estimator		Random Sampling Estimator	
Loss $\ell(z, \cdot)$	Convex	Nonconvex Smooth	Convex	Nonconvex Smooth
Choice of $f$ -divergence	Arbitrary	KL-divergence	Arbitrary	KL-divergence
Complexity	$\tilde{O}(\delta^{-3})$	$\tilde{O}(\delta^{-6})$	$\tilde{O}(\delta^{-2})$	$\tilde{O}(\delta^{-4})$

•  $\nabla F(\theta) = \mathbb{E}_{z \sim \mathbb{P}_n} \mathbb{E}_{z' \sim \mathbb{Q}(z; \theta)} [\nabla_\theta \ell(z'; \theta)]$ , where optimal response  $q^*(z; \theta)$  solves  $R(\theta; z)$ .



## Random Sampling Gradient Estimator

$$\begin{aligned} \hat{v}(\theta; q_{L+1}) &= \hat{v}(\theta; q_1) + \sum_{\ell=1}^L \left[ \hat{v}(\theta; q_{\ell+1}) - \hat{v}(\theta; q_\ell) \right] \\ &= \hat{v}(\theta; q_1) + \sum_{\ell=1}^L p_\ell \frac{\hat{v}(\theta; q_{\ell+1}) - \hat{v}(\theta; q_\ell)}{p_\ell} = \mathbb{E}_{\ell \sim \mathbb{P}_\ell} \left[ \hat{v}(\theta; q_1) + \frac{\hat{v}(\theta; q_{\ell+1}) - \hat{v}(\theta; q_\ell)}{p_\ell} \right]. \end{aligned}$$

• Sample  $\ell$  according to pmf  $p_\ell \propto 2^{-\ell}$ ,  $\sum_\ell p_\ell = 1$ , then construct

$$\hat{v}(\theta) = \hat{v}(\theta; q_1) + \frac{\hat{v}(\theta; q_{\ell+1}) - \hat{v}(\theta; q_\ell)}{p_\ell}$$

• High probability: generate small  $\ell$ ; Low probability: generate large  $\ell$

• Per-iteration cost reduction: From  $\mathcal{O}(\delta^{-1})$  to  $\tilde{O}(1)$

• Variance reduction effect as  $\hat{v}(\theta; q_{\ell+1}) - \hat{v}(\theta; q_\ell) \rightarrow 0$

## Regularization Effect

$$( \text{Regularized-AT} ) \approx \begin{cases} \min_{\theta} \mathbb{E}_{\mathbb{P}_n} [\ell(z; \theta)] + \rho \mathbb{E}_{\mathbb{P}_n} [\|\nabla \ell(z; \theta)\|_*], & \text{if } \rho/\eta \rightarrow \infty \\ \min_{\theta} \mathbb{E}_{\mathbb{P}_n} [\ell(z; \theta)] + \frac{\rho^2}{2f''(1) \cdot \eta} \text{Var}(\nabla \ell(z; \theta)^\top \Delta), & \text{if } \rho/\eta \rightarrow 0 \\ \min_{\theta} \mathbb{E}_{\mathbb{P}_n} [\ell(z; \theta)] + \frac{C}{\rho} \mathbb{E}_{\mathbb{P}_n} \left[ \inf_{\mu \in \mathbb{R}} \{ \mu + \mathbb{E}[f^*(C \cdot \nabla \ell(z; \theta)^\top \Delta - \mu)] \} \right], & \text{if } \rho/\eta \rightarrow C \end{cases}$$

## Numerical Study on Learning and Control

