

Distributionally & Adversarially Robust Logistic Regression Intersecting Wasserstein Balls

Aras Selvi | Eleonora Kreacic | Mohsen Ghassemi | Vamsi Potluru | Tucker Balch | Manuela Veloso

IMPERIAL
BUSINESS SCHOOL



JPMC

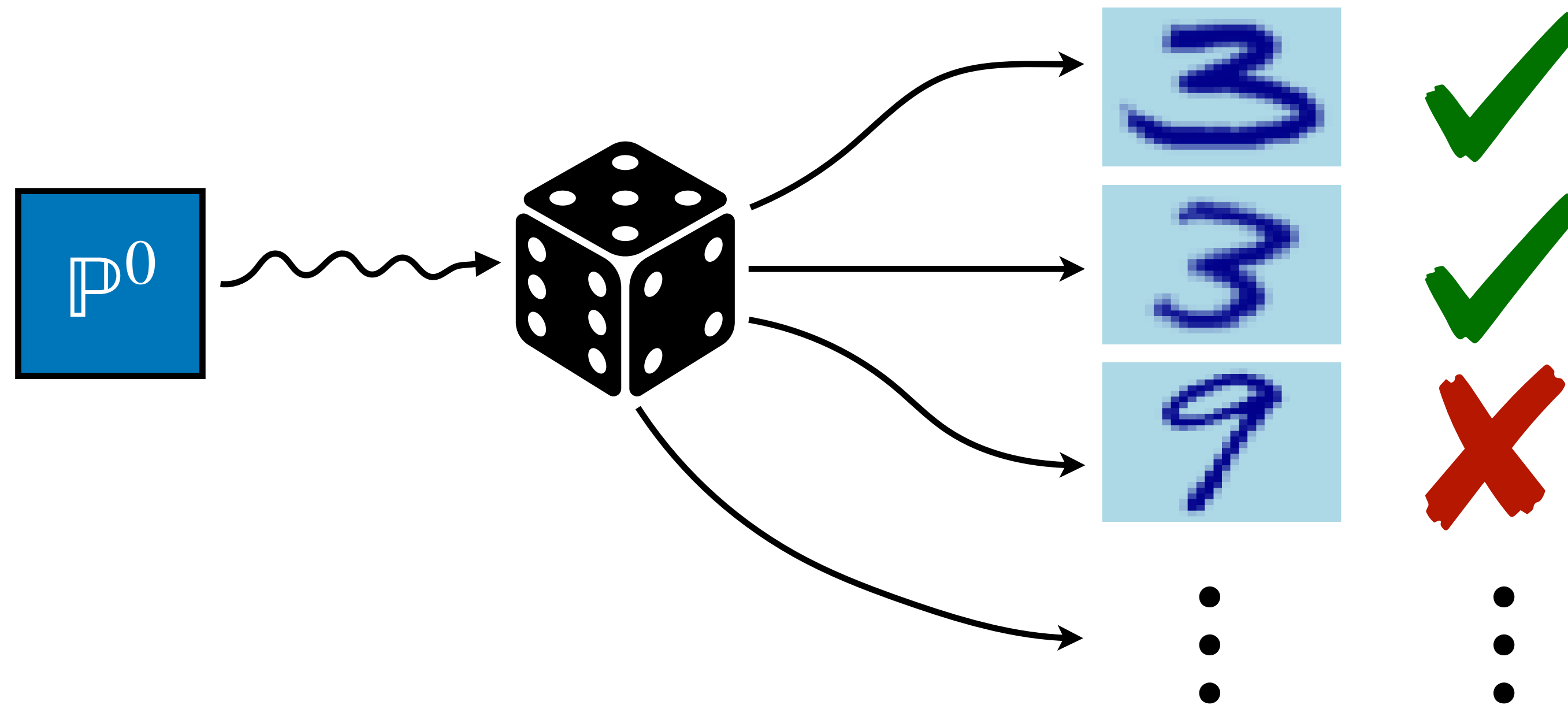
Sargent Centre
for Process Systems
Engineering

True Risk Minimization

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$$

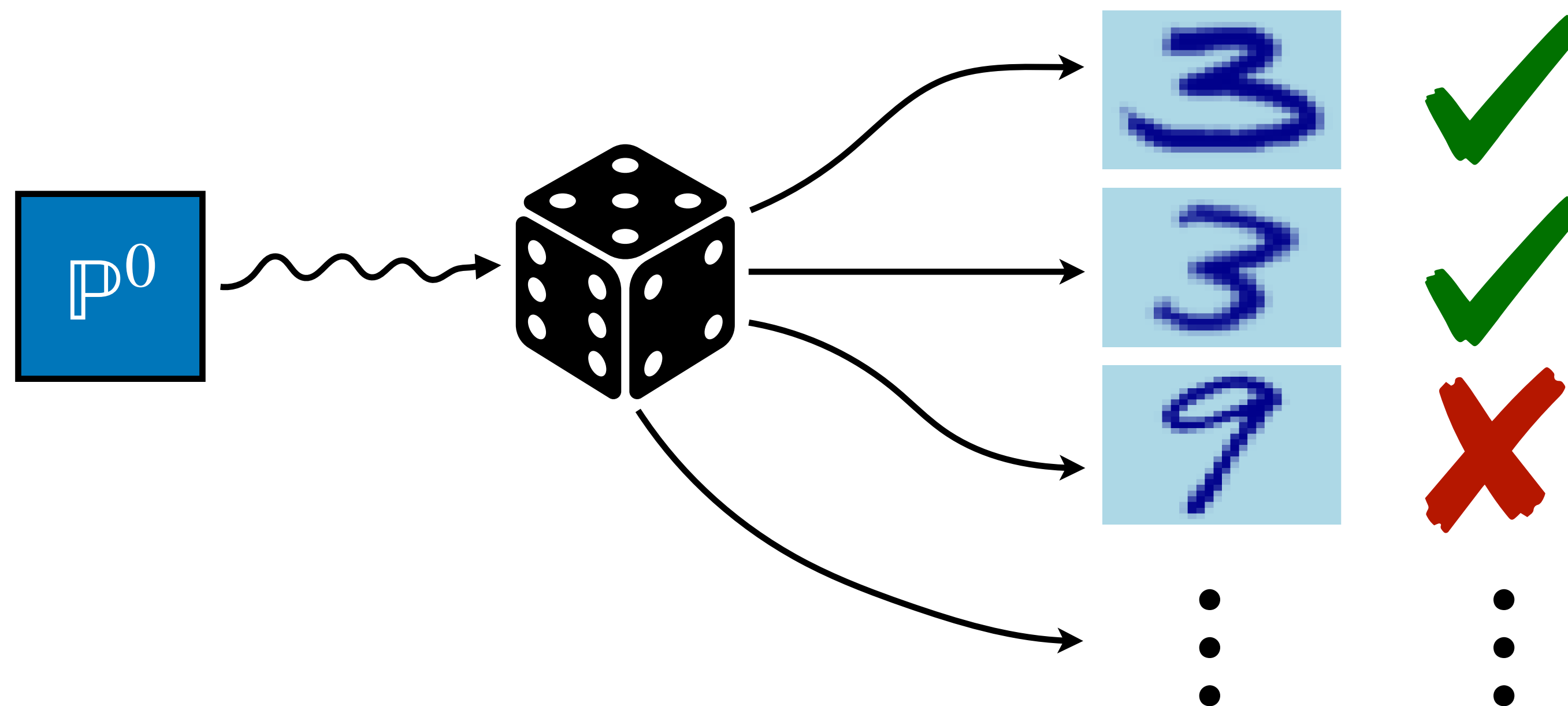
True Risk Minimization

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$$



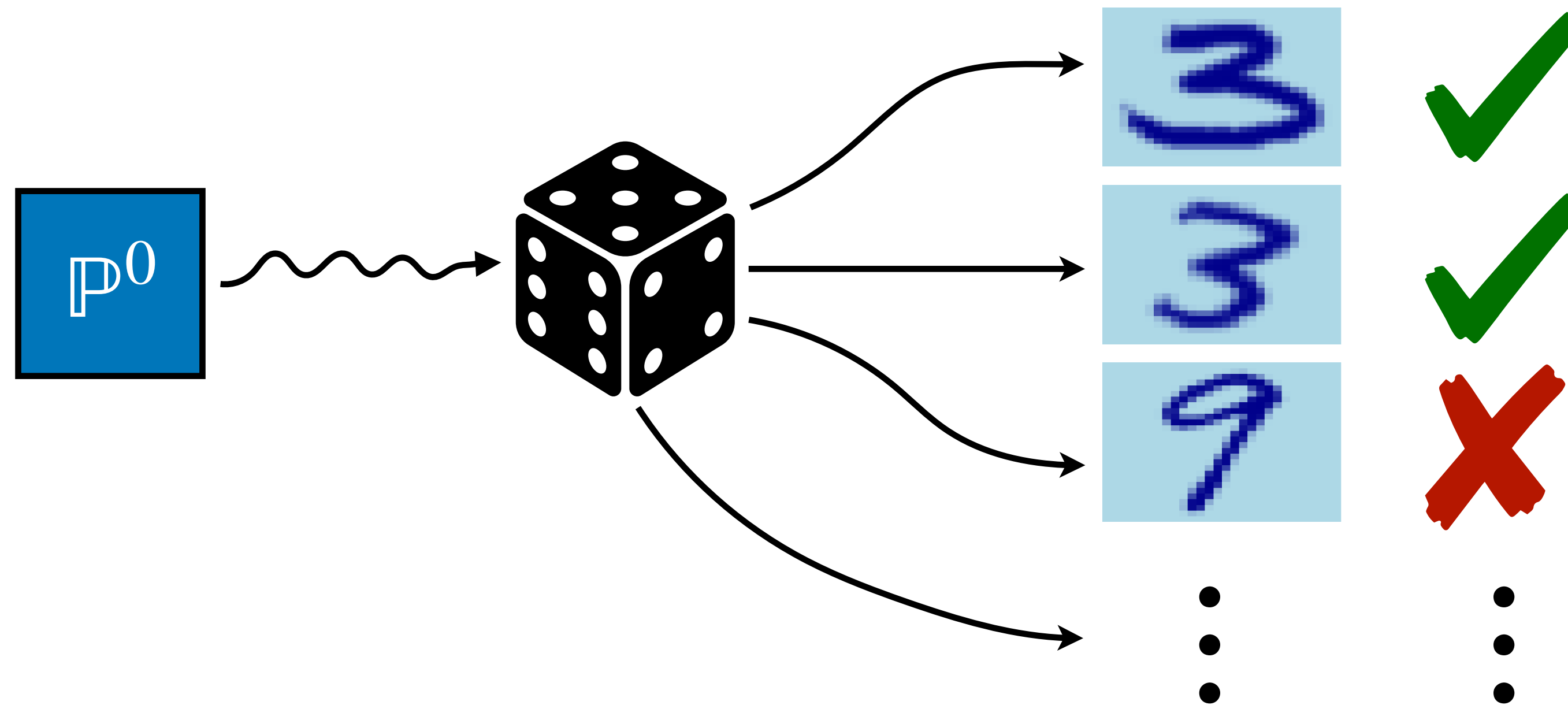
True Risk Minimization

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$$



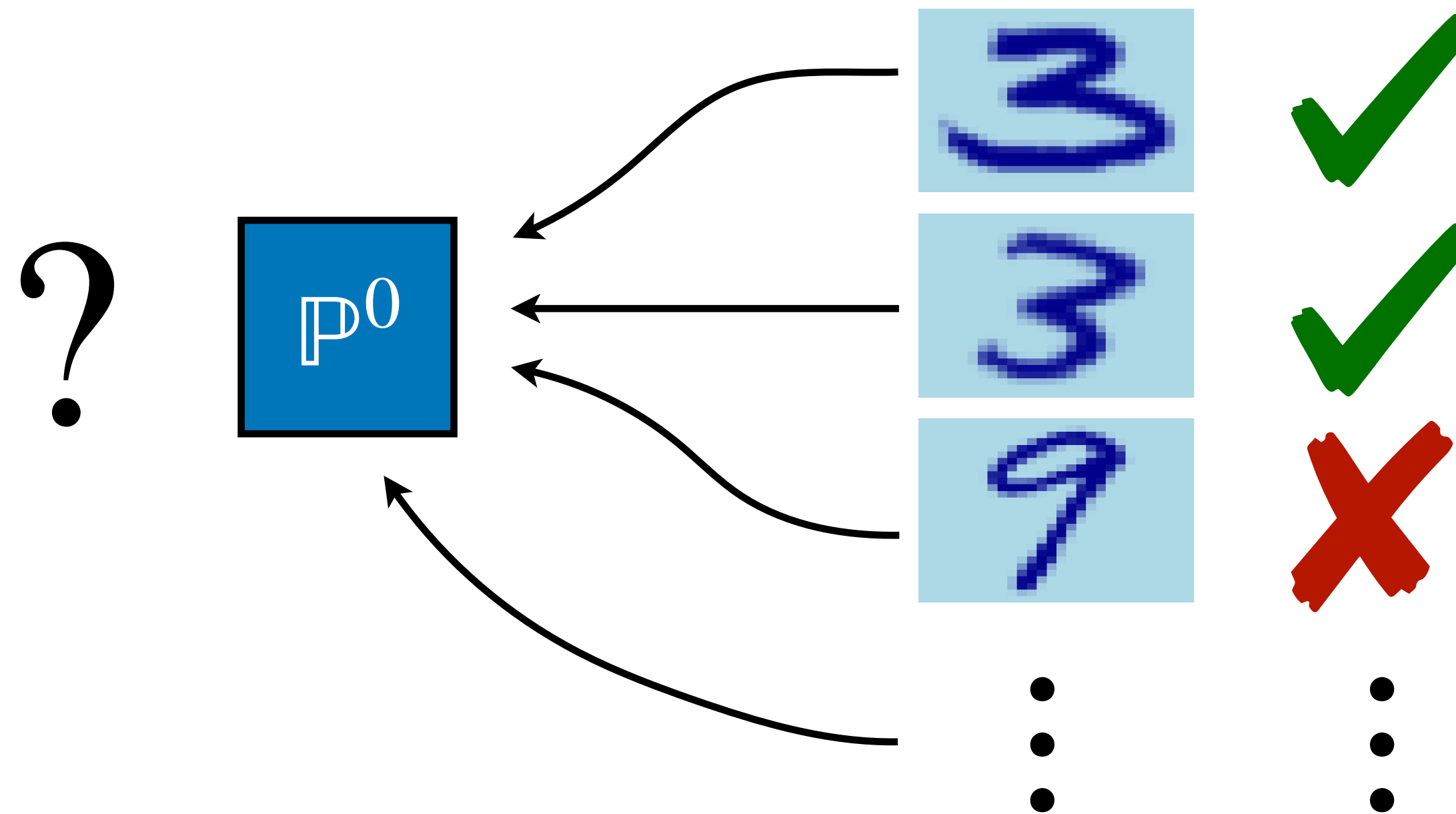
True Risk Minimization

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$$



True Risk Minimization

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$$



True Risk Minimization

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$$

Data

$$\{\xi^i = (\mathbf{x}^i, y^i)\}_{i \in [N]} \xrightarrow{\sim \text{iid } \mathbb{P}^0}$$

Empirical Risk Minimization

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}(\mathbf{x}, y)]$$

$$\mathbb{P}_N := \frac{1}{N} \sum_{i \in [N]} \delta_{\xi^i}$$

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N}[\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0}[\ell_{\beta}(\mathbf{x}, y)]$

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N}[\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0}[\ell_{\beta}(\mathbf{x}, y)]$

Overfitting

We tend to **underestimate** the true risk with β^{ERM}

$$\mathbb{E}_{\mathbb{P}^0}[\ell_{\beta^{\text{ERM}}}(\mathbf{x}, y)] - \mathbb{E}_{\mathbb{P}_N}[\ell_{\beta^{\text{ERM}}}(\mathbf{x}, y)]$$

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N}[\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0}[\ell_{\beta}(\mathbf{x}, y)]$

Overfitting

We tend to **underestimate** the true risk with β^{ERM}

$$\mathbb{E}_{\mathbb{P}^0}[\ell_{\beta^{\text{ERM}}}(\mathbf{x}, y)] - \mathbb{E}_{\mathbb{P}_N}[\ell_{\beta^{\text{ERM}}}(\mathbf{x}, y)]$$

DRO philosophy: statistical error of estimating \mathbb{P}^0 via \mathbb{P}_N

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N}[\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0}[\ell_{\beta}(\mathbf{x}, y)]$

Overfitting

We tend to **underestimate** the true risk with β^{ERM}

$$\mathbb{E}_{\mathbb{P}^0}[\ell_{\beta^{\text{ERM}}}(\mathbf{x}, y)] - \mathbb{E}_{\mathbb{P}_N}[\ell_{\beta^{\text{ERM}}}(\mathbf{x}, y)]$$

DRO philosophy: statistical error of estimating \mathbb{P}^0 via \mathbb{P}_N

$$W(\mathbb{P}_N, \mathbb{P}^0) > 0$$

Metric on the feature-label space Ξ

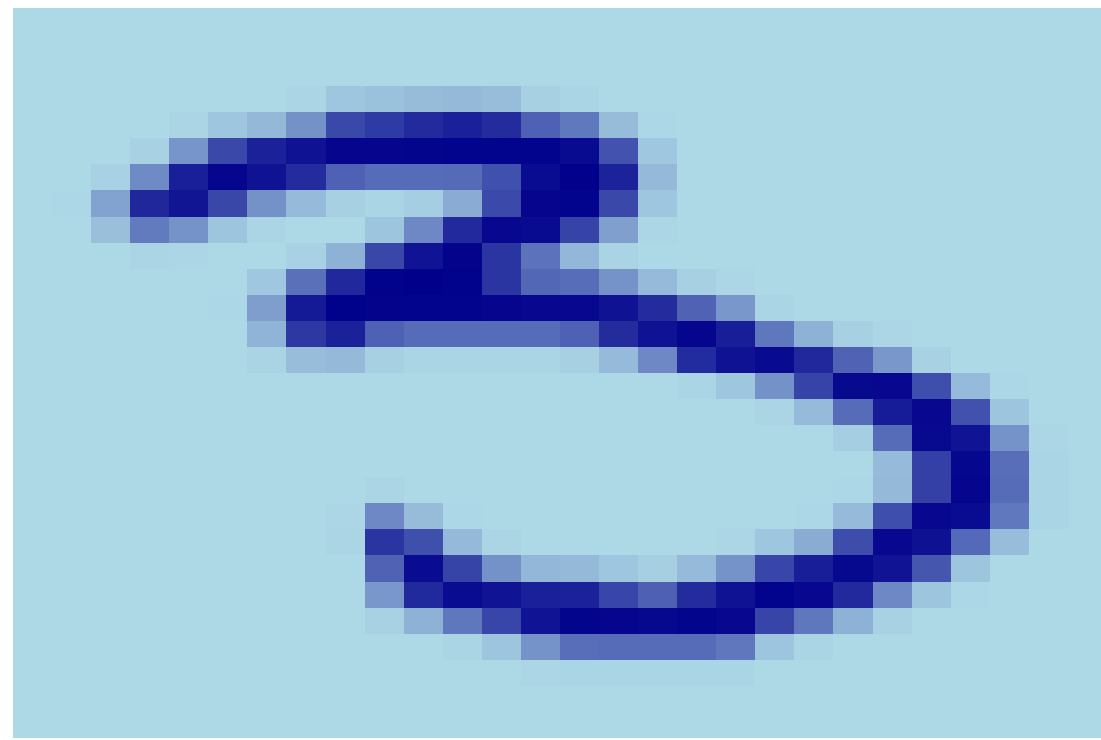
Distance between $\xi = (x, y) \in \Xi$ and $\xi' = (x', y') \in \Xi$ is

$$d(\xi, \xi') = \|x - x'\|_q + \kappa \cdot \mathbf{1}[y \neq y']$$

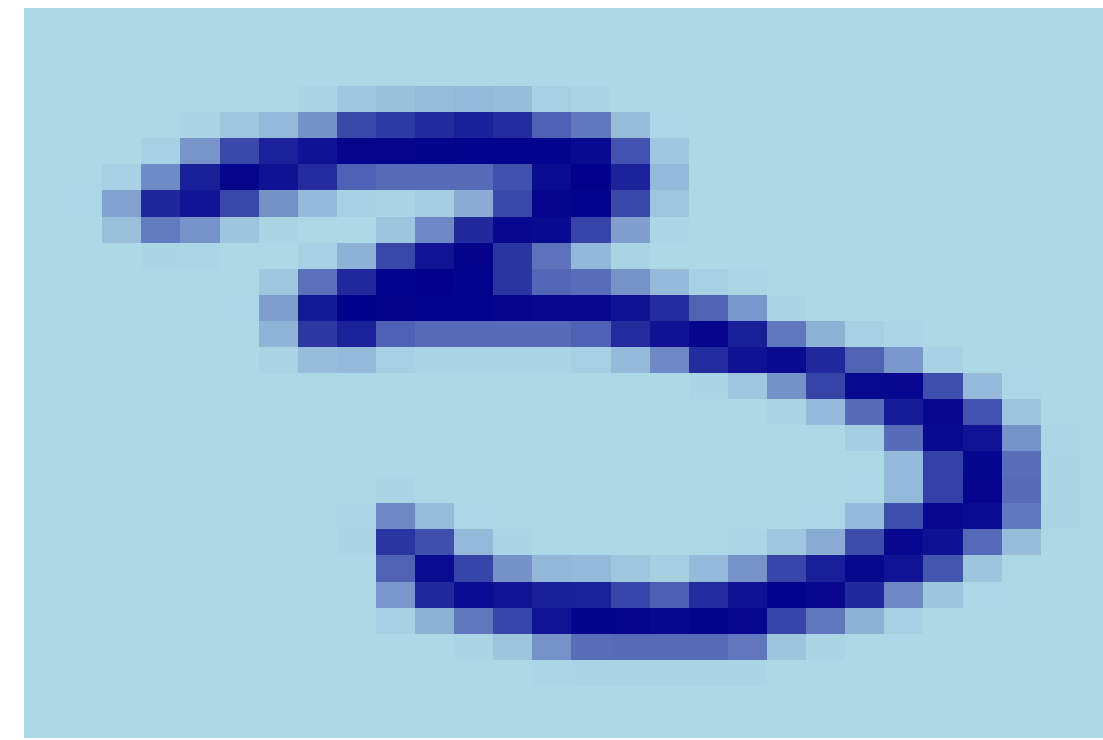
The Wasserstein Distance

Metric on the feature-label space Ξ

Example $\kappa = 28, q = 2$



ξ



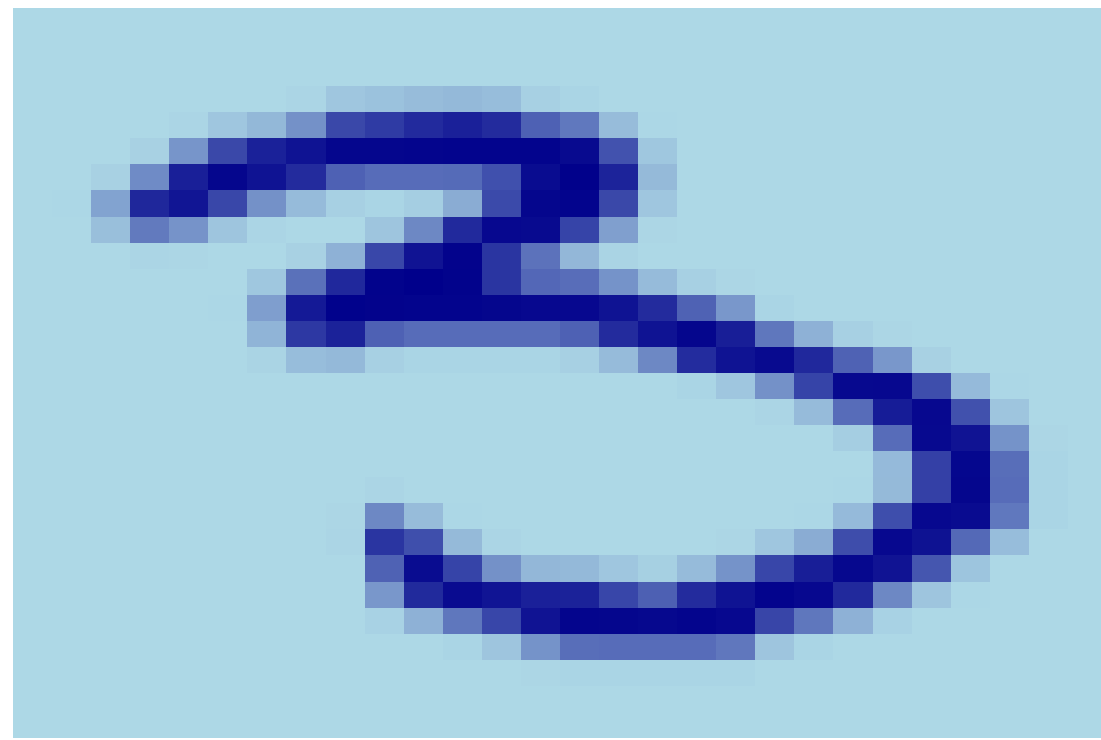
ξ'

$$d(\xi, \xi') = 0$$

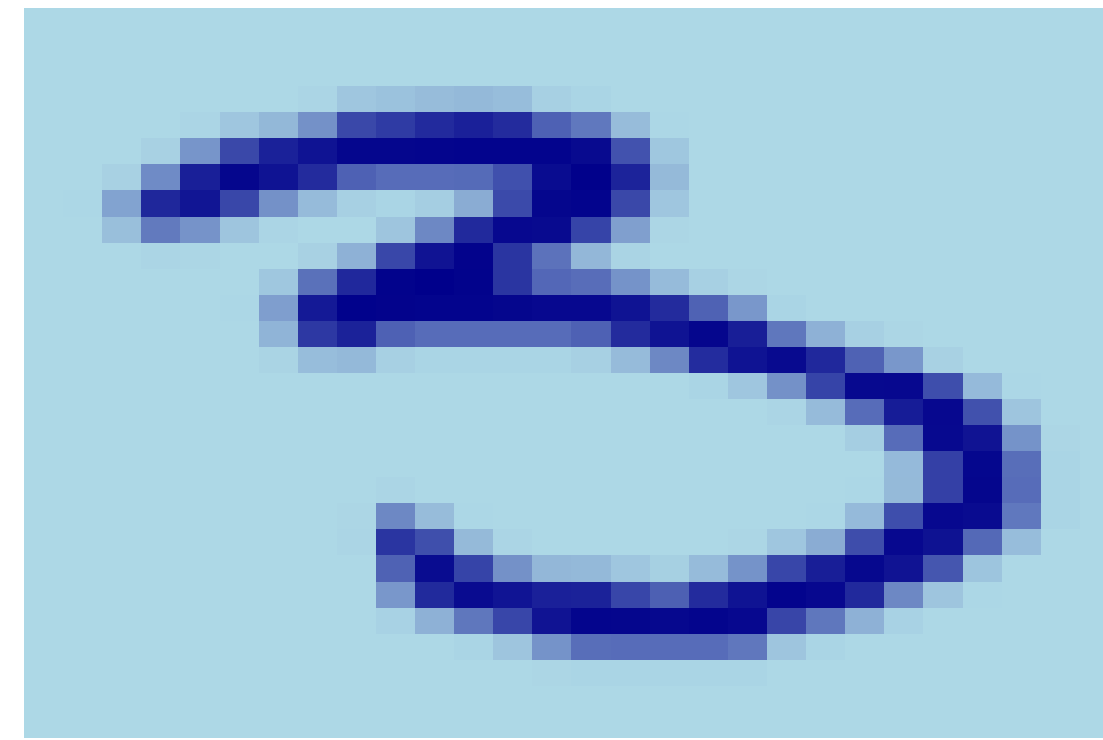
The Wasserstein Distance

Metric on the feature-label space Ξ

Example $\kappa = 28, q = 2$



ξ



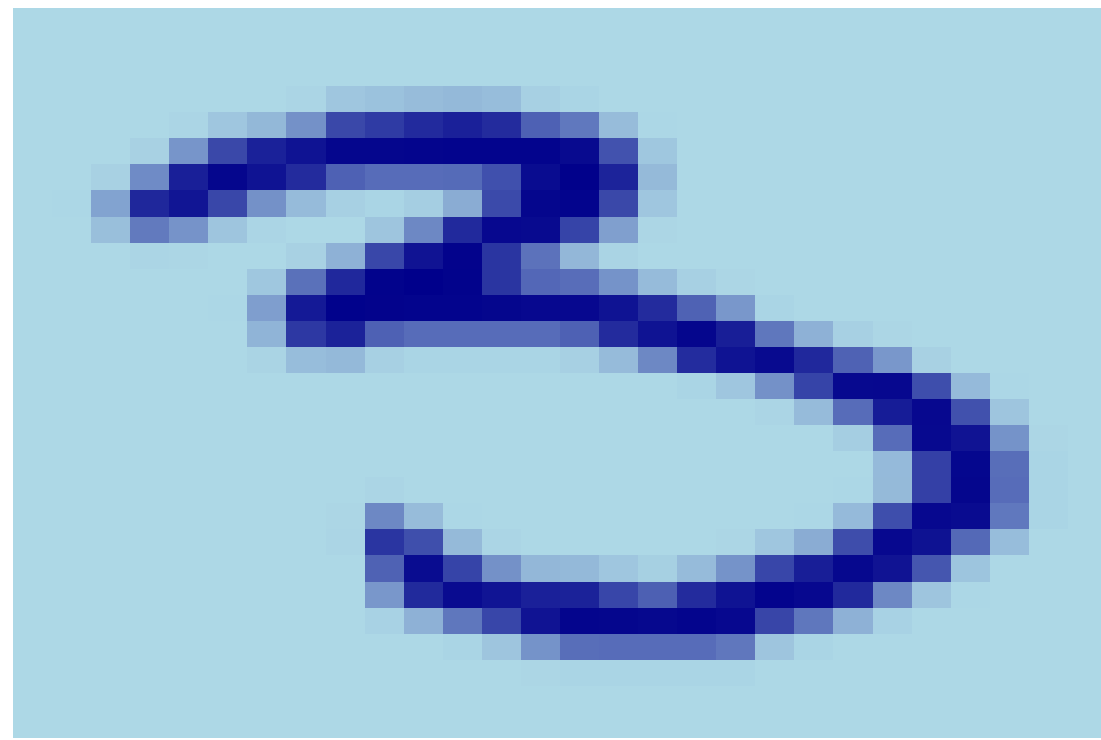
ξ'

$$d(\xi, \xi') = 0.5$$

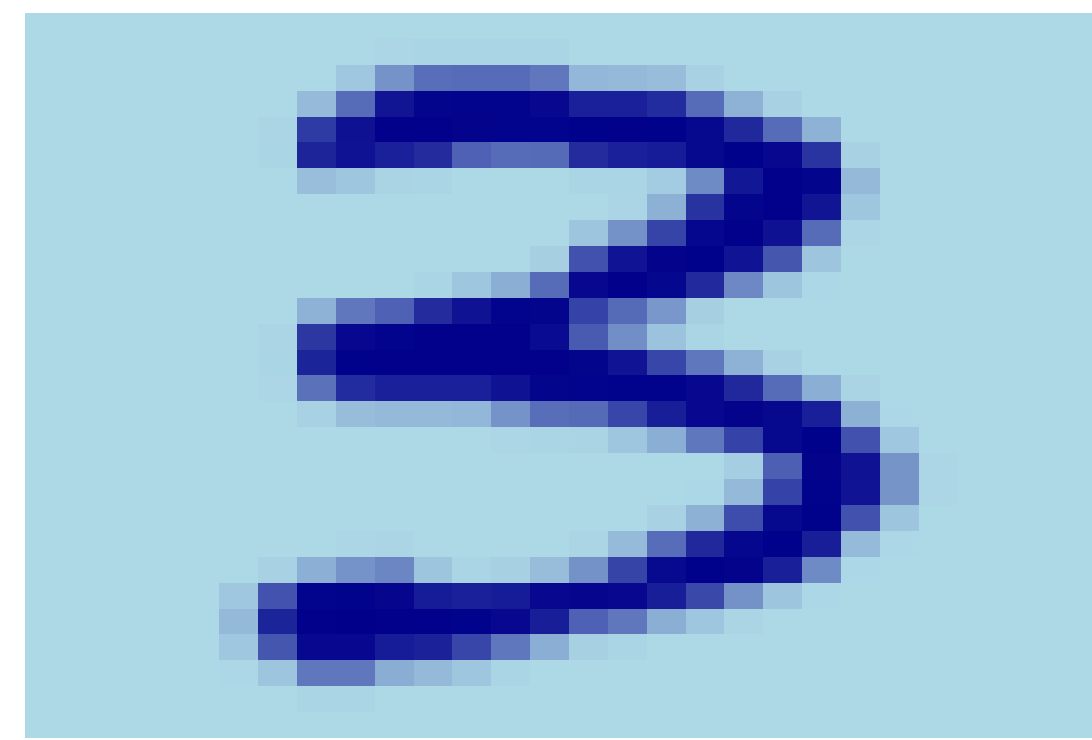
The Wasserstein Distance

Metric on the feature-label space Ξ

Example $\kappa = 28, q = 2$



ξ



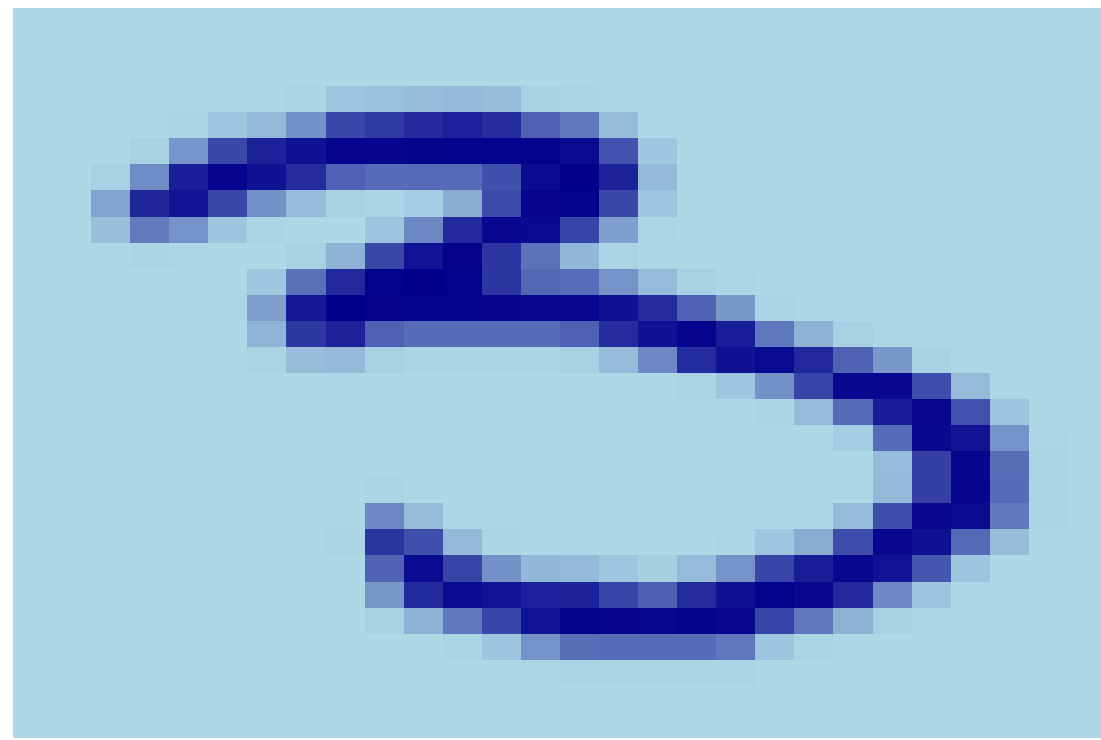
ξ'

$$d(\xi, \xi') = 0.21$$

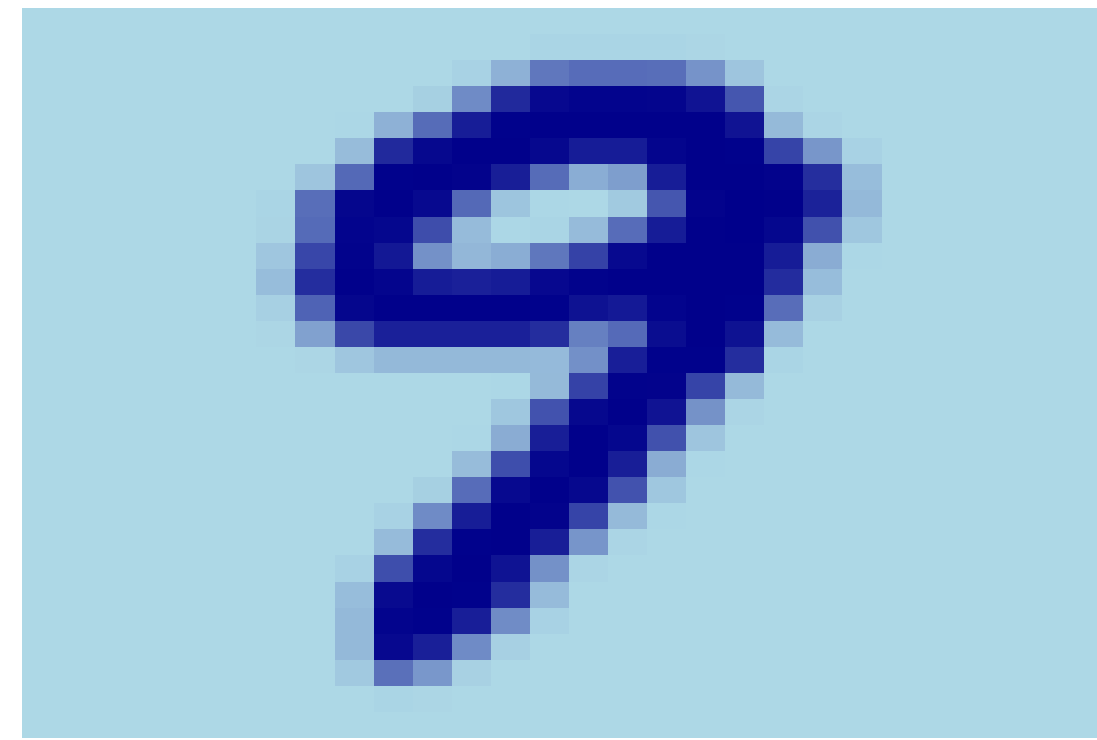
The Wasserstein Distance

Metric on the feature-label space Ξ

Example $\kappa = 28, q = 2$



ξ



ξ'

$$d(\xi, \xi') = 0.70$$

The Wasserstein Distance

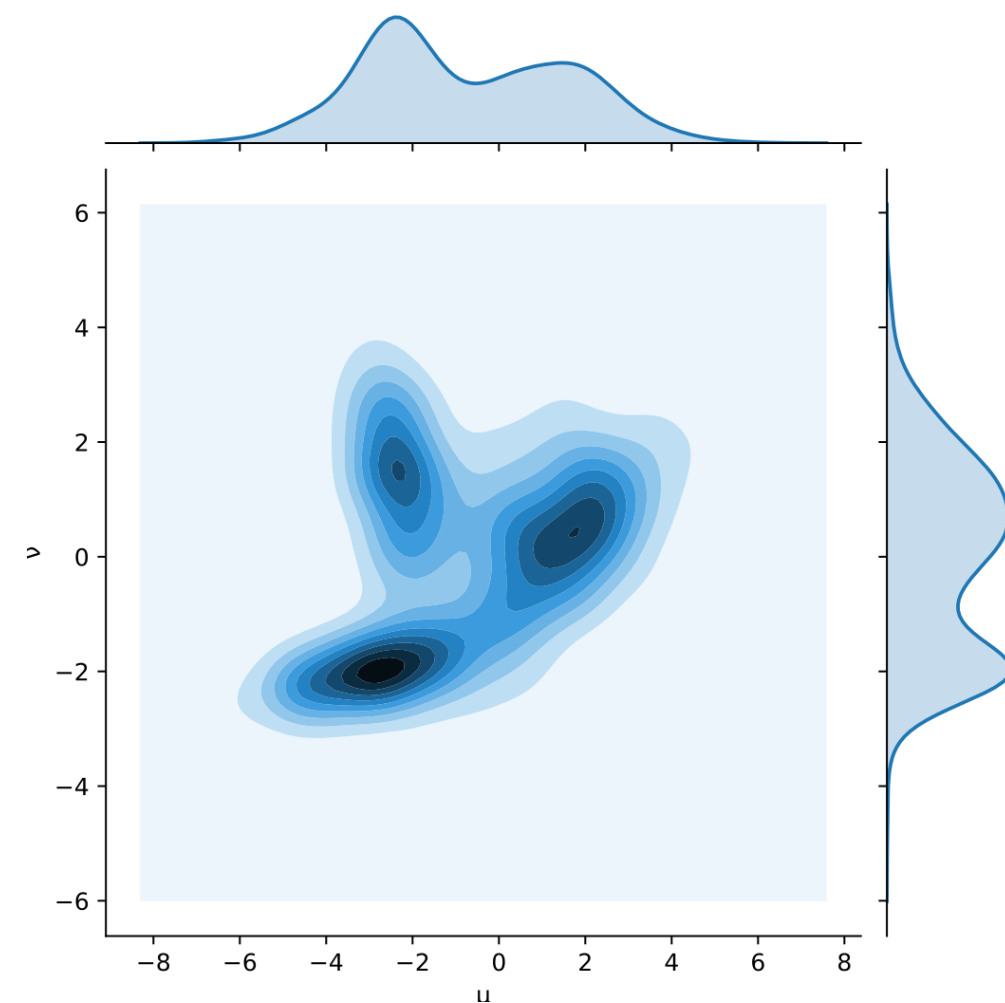
Metric on the feature-label space Ξ

Distance between $\xi = (x, y) \in \Xi$ and $\xi' = (x', y') \in \Xi$ is

$$d(\xi, \xi') = \|x - x'\|_q + \kappa \cdot \mathbf{1}[y \neq y']$$

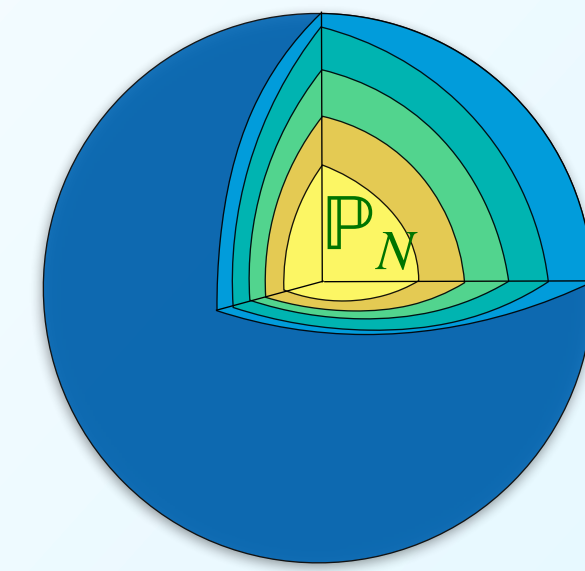
Wasserstein distance

$$W(Q, Q') = \inf_{\Pi \in \mathcal{C}(Q, Q')} \mathbb{E}_{\Pi}[d(\xi, \xi')]$$



Wasserstein ball

$$\mathfrak{B}_{\varepsilon}(\mathbb{P}_N) = \{Q : W(\mathbb{P}_N, Q) \leq \varepsilon\}$$



$$\mathbb{P}_N = \frac{1}{N} \sum_{i \in [N]} \delta_{\xi^i}$$

Distributionally Robust Optimization

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N}[\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0}[\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}}[\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0}[\ell_{\beta}(\mathbf{x}, y)]$

Distributionally Robust Optimization

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N}[\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0}[\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}}[\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0}[\ell_{\beta}(\mathbf{x}, y)]$

1 Access Training Set

$$\{\xi^i = (x^i, y^i)\}_{i \in [N]}$$

2 Optimize Expected ℓ_β

$$\text{minimize}_\beta \mathbb{E}_{(x,y) \sim \mathbb{P}_N}[\ell_\beta(x, y)]$$

β^\star

3 Deploy/Test Solution

$$\mathbb{E}_{(x,y) \sim \mathbb{P}^0}[\ell_{\beta^\star}(x, y)]$$

1 Access Training Set

$$\{\xi^i = (x^i, y^i)\}_{i \in [N]}$$

2 Optimize Expected ℓ_β

$$\text{minimize}_\beta \mathbb{E}_{(x,y) \sim \mathbb{P}_N} [\ell_\beta(x, y)]$$

β^\star



3 Deploy/Test Solution

$$\mathbb{E}_{(x,y) \sim \mathbb{P}^0} [\ell_{\beta^\star}(x, y)]$$

1 Access Training Set

$$\{\xi^i = (x^i, y^i)\}_{i \in [N]}$$

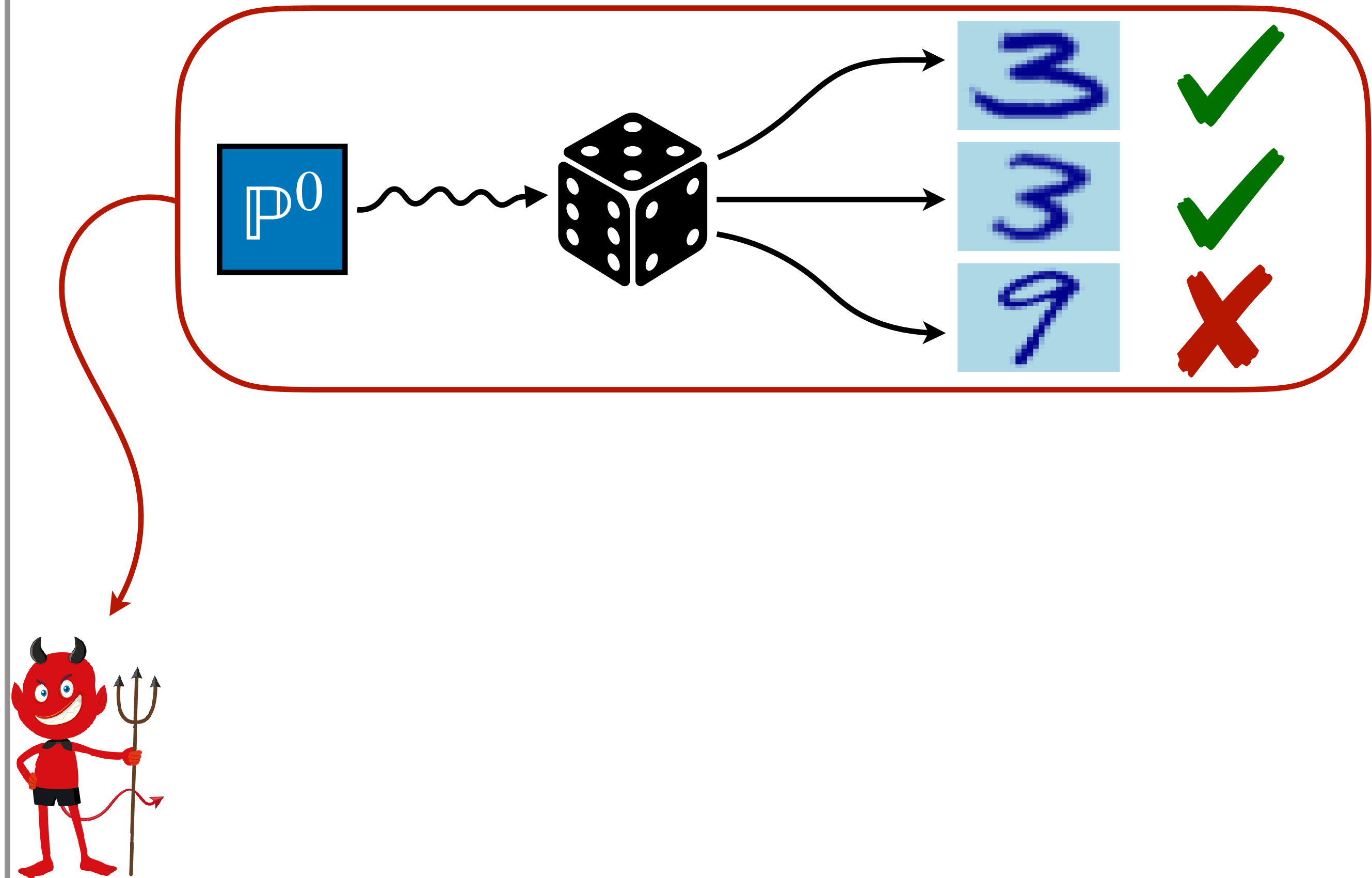
2 Optimize Expected ℓ_β

$$\text{minimize}_\beta \mathbb{E}_{(x,y) \sim \mathbb{P}_N} [\ell_\beta(x, y)]$$

β^\star

3 Deploy/Test Solution

$$\mathbb{E}_{(x,y) \sim \mathbb{P}^0} [\ell_{\beta^\star}(x, y)]$$



1 Access Training Set

$$\{\xi^i = (x^i, y^i)\}_{i \in [N]}$$

2 Optimize Expected ℓ_β

$$\text{minimize}_\beta \mathbb{E}_{(x,y) \sim \mathbb{P}_N} [\ell_\beta(x, y)]$$

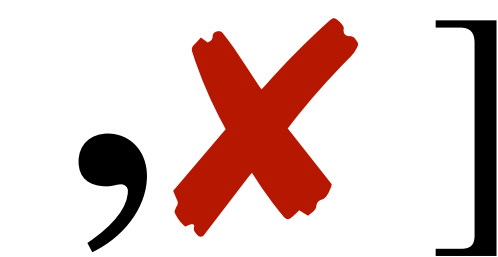
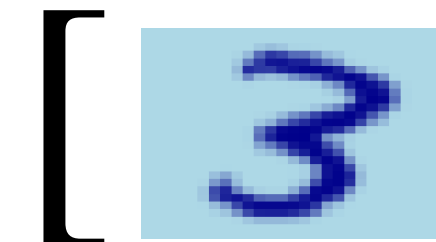
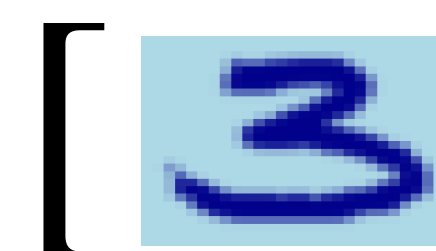
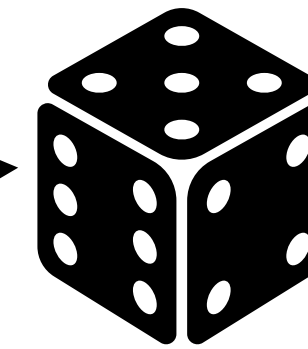
β^\star

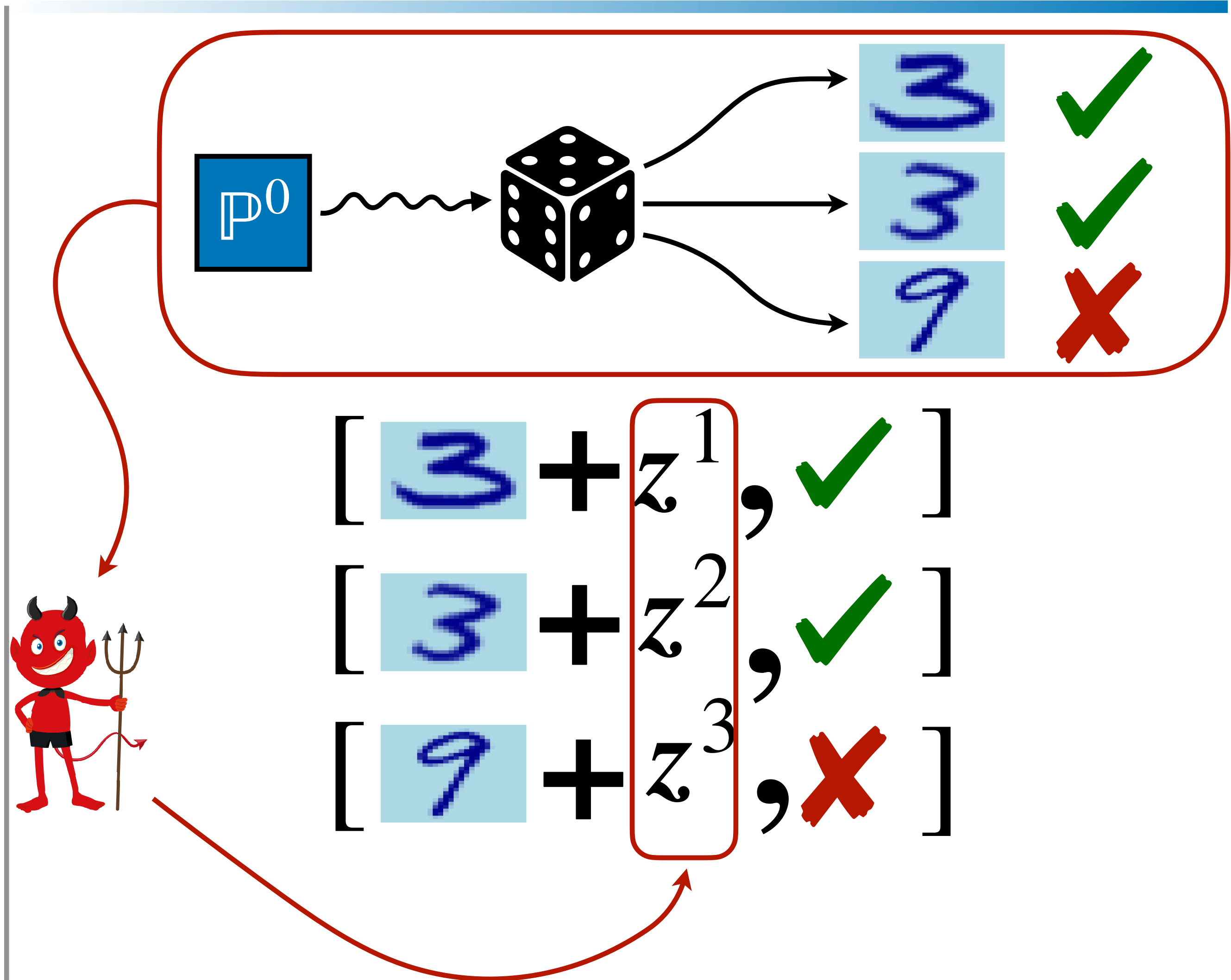
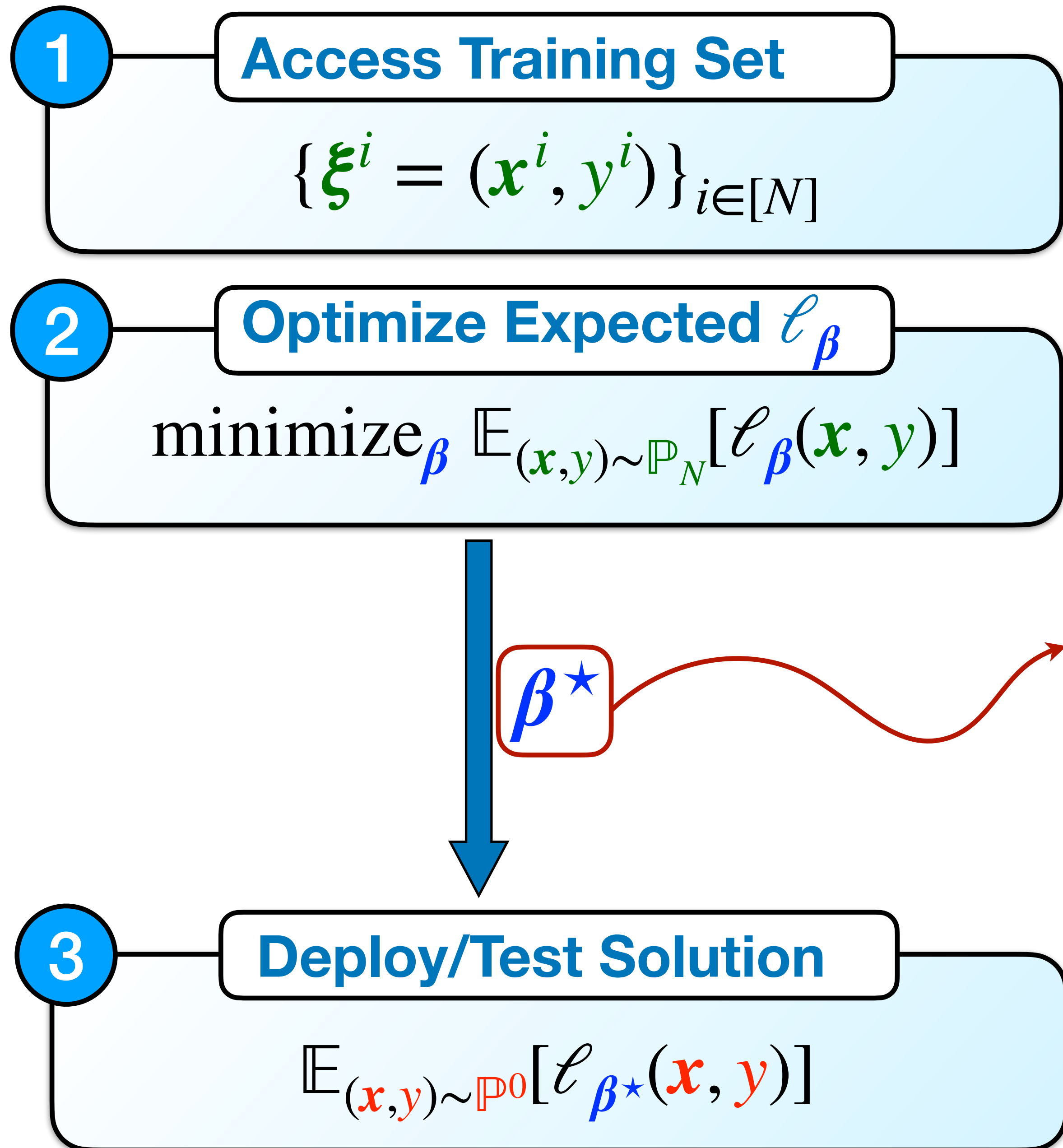
3 Deploy/Test Solution

$$\mathbb{E}_{(x,y) \sim \mathbb{P}^0} [\ell_{\beta^\star}(x, y)]$$



\mathbb{P}^0





1 Access Training Set

$$\{\xi^i = (x^i, y^i)\}_{i \in [N]}$$

2 Optimize Expected ℓ_β

$$\text{minimize}_\beta \mathbb{E}_{(x,y) \sim \mathbb{P}_N} [\ell_\beta(x, y)]$$

β^\star

3 Deploy/Test Solution

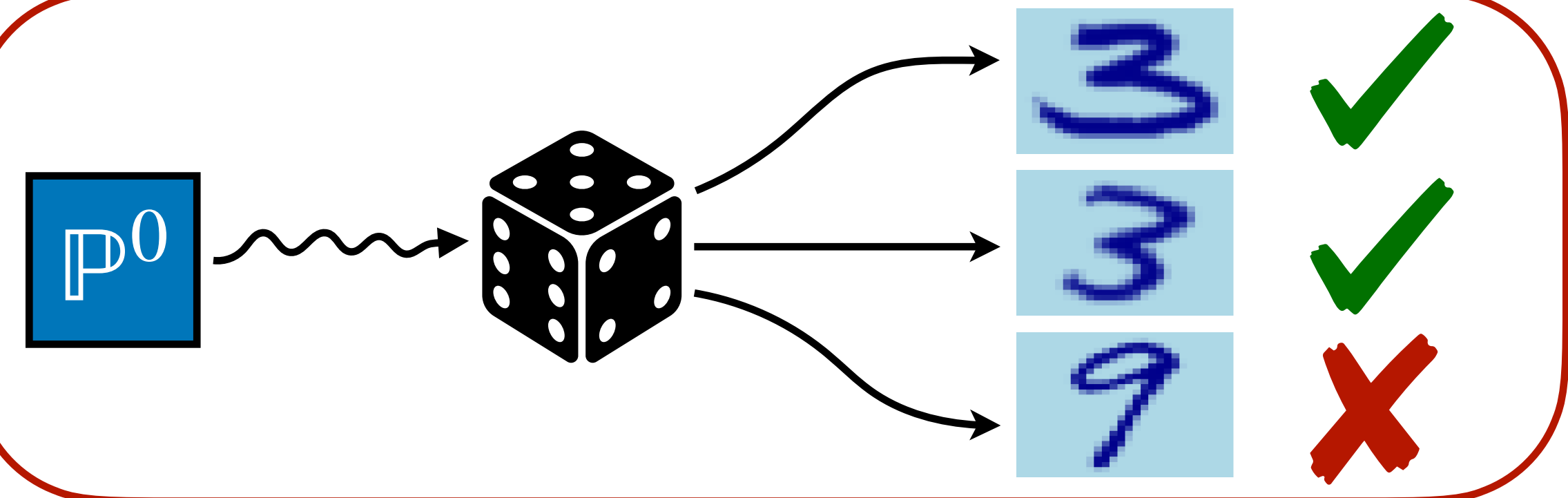
$$\mathbb{E}_{(x,y) \sim \mathbb{P}^0} [\ell_{\beta^\star}(x, y)]$$



3

Adversarial Attacks

$$\mathbb{E}_{(x,y) \sim \mathbb{P}^0} \left[\sup_{\|z\|_p \leq \alpha} \ell_{\beta^\star}(x + z, y) \right]$$



ℓ_2 -attack for β^{ERM}



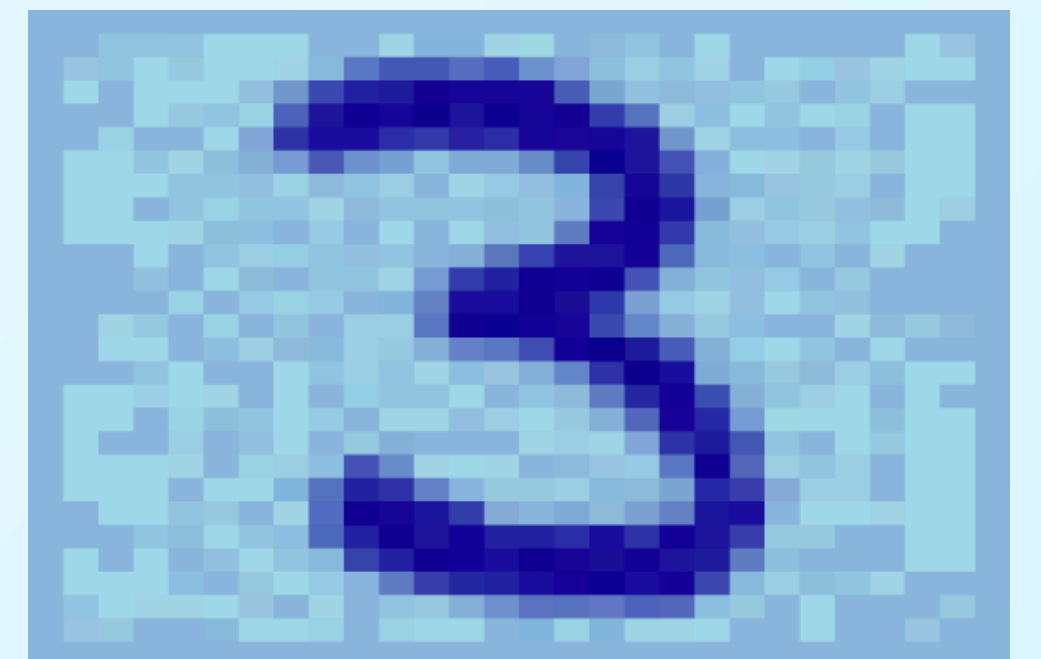
$$\beta^{\text{ERM}\top} \mathbf{x} = 34$$



$$\beta^{\text{ERM}\top} (\mathbf{x} + \mathbf{z}) = 1$$



$$\beta^{\text{ERM}\top} (\mathbf{x} + \mathbf{z}) = -76$$



$$\beta^{\text{ERM}\top} (\mathbf{x} + \mathbf{z}) = -408$$

Stronger attack radius α

Distributionally Robust Optimization

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Adversarially RO		$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} \left[\sup_{\ \mathbf{z}\ _p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y) \right]$

Distributionally Robust Optimization

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Adversarially RO	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\sup_{\ \mathbf{z}\ _p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\sup_{\ \mathbf{z}\ _p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y)]$

ℓ_2 -attack for β^{ARO}



$$\beta^{\text{ARO}\top} \mathbf{x} = 1.90$$



$$\beta^{\text{ARO}\top} (\mathbf{x} + \mathbf{z}) = 1.80$$



$$\beta^{\text{ARO}\top} (\mathbf{x} + \mathbf{z}) = 1.56$$



$$\beta^{\text{ARO}\top} (\mathbf{x} + \mathbf{z}) = 0.53$$

Stronger attack radius α

Distributionally Robust Optimization

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Adversarially RO	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} \left[\sup_{\ \mathbf{z}\ _p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y) \right]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} \left[\sup_{\ \mathbf{z}\ _p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y) \right]$

Paradigm	Robust Overfitting		Empirical Risk
Empirical Risk Min	<ul style="list-style-type: none">• ARO models overfit despite being “robust”• Even more severe than overfitting of ERM• We want DRO and ARO simultaneously		$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N}[\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO			$\sup_{Q \in \mathcal{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim Q}[\ell_{\beta}(\mathbf{x}, y)]$
Adversarially RO			$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_0}[\ell_{\beta}(\mathbf{x}, y)]$
	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N}[\sup_{\ \mathbf{z}\ _p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y)]$		$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_0}[\sup_{\ \mathbf{z}\ _p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y)]$

$$\sup_{\|\mathbf{z}\|_p \leq \alpha} \ell_{\boldsymbol{\beta}}(\mathbf{x} + \mathbf{z}, y) = \sup_{\|\mathbf{z}\|_p \leq \alpha} \log(1 + \exp(-y \cdot \boldsymbol{\beta}^\top (\mathbf{x} + \mathbf{z})))$$

$$\begin{aligned} \sup_{\|\mathbf{z}\|_p \leq \alpha} \ell_{\boldsymbol{\beta}}(\mathbf{x} + \mathbf{z}, y) &= \sup_{\|\mathbf{z}\|_p \leq \alpha} \log(1 + \exp(-y \cdot \boldsymbol{\beta}^\top (\mathbf{x} + \mathbf{z}))) \\ &\quad \downarrow \\ &= \log(1 + \exp(\sup_{\|\mathbf{z}\|_p \leq \alpha} \{-y \cdot \boldsymbol{\beta}^\top (\mathbf{x} + \mathbf{z})\})) \end{aligned}$$

Monotonicity: supremum goes inside

$$\begin{aligned} \sup_{\|\mathbf{z}\|_p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y) &= \sup_{\|\mathbf{z}\|_p \leq \alpha} \log(1 + \exp(-y \cdot \beta^{\top}(\mathbf{x} + \mathbf{z}))) \\ &\downarrow \\ &= \log(1 + \exp(\sup_{\|\mathbf{z}\|_p \leq \alpha} \{-y \cdot \beta^{\top}(\mathbf{x} + \mathbf{z})\})) \\ &\downarrow \\ &= \log(1 + \exp(-y \cdot \beta^{\top} \mathbf{x} + \sup_{\|\mathbf{z}\|_p \leq \alpha} \{-y \cdot \beta^{\top} \mathbf{z}\})) \end{aligned}$$

Rewrite: take constants out

$$\begin{aligned}
 \sup_{\|\mathbf{z}\|_p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y) &= \sup_{\|\mathbf{z}\|_p \leq \alpha} \log(1 + \exp(-y \cdot \beta^{\top}(\mathbf{x} + \mathbf{z}))) \\
 &\downarrow \\
 &= \log(1 + \exp(\sup_{\|\mathbf{z}\|_p \leq \alpha} \{-y \cdot \beta^{\top}(\mathbf{x} + \mathbf{z})\})) \\
 &\downarrow \\
 &= \log(1 + \exp(-y \cdot \beta^{\top} \mathbf{x} + \sup_{\|\mathbf{z}\|_p \leq \alpha} \{-y \cdot \beta^{\top} \mathbf{z}\})) \\
 &\downarrow \\
 &= \log(1 + \exp(-y \cdot \beta^{\top} \mathbf{x} + \alpha \cdot \|-y \cdot \beta\|_{p^*}))
 \end{aligned}$$

Dual norm: use definition

$$\begin{aligned}
 \sup_{\|\mathbf{z}\|_p \leq \alpha} \ell_{\boldsymbol{\beta}}(\mathbf{x} + \mathbf{z}, y) &= \sup_{\|\mathbf{z}\|_p \leq \alpha} \log(1 + \exp(-y \cdot \boldsymbol{\beta}^\top (\mathbf{x} + \mathbf{z}))) \\
 &\downarrow \\
 &= \log(1 + \exp(\sup_{\|\mathbf{z}\|_p \leq \alpha} \{-y \cdot \boldsymbol{\beta}^\top (\mathbf{x} + \mathbf{z})\})) \\
 &\downarrow \\
 &= \log(1 + \exp(-y \cdot \boldsymbol{\beta}^\top \mathbf{x} + \sup_{\|\mathbf{z}\|_p \leq \alpha} \{-y \cdot \boldsymbol{\beta}^\top \mathbf{z}\})) \\
 &\downarrow \\
 &= \log(1 + \exp(-y \cdot \boldsymbol{\beta}^\top \mathbf{x} + \alpha \cdot \|\boldsymbol{\beta}\|_{p^*}))
 \end{aligned}$$

$$\begin{aligned}
 \sup_{\|\mathbf{z}\|_p \leq \alpha} \ell_{\boldsymbol{\beta}}(\mathbf{x} + \mathbf{z}, y) &= \sup_{\|\mathbf{z}\|_p \leq \alpha} \log(1 + \exp(-y \cdot \boldsymbol{\beta}^\top (\mathbf{x} + \mathbf{z}))) \\
 &\downarrow \\
 &= \log(1 + \exp(\sup_{\|\mathbf{z}\|_p \leq \alpha} \{-y \cdot \boldsymbol{\beta}^\top (\mathbf{x} + \mathbf{z})\})) \\
 &\downarrow \\
 &= \log(1 + \exp(-y \cdot \boldsymbol{\beta}^\top \mathbf{x} + \sup_{\|\mathbf{z}\|_p \leq \alpha} \{-y \cdot \boldsymbol{\beta}^\top \mathbf{z}\})) \\
 &\downarrow \\
 &= \log(1 + \exp(-y \cdot \boldsymbol{\beta}^\top \mathbf{x} + \alpha \cdot \|\boldsymbol{\beta}\|_{p^\star})) =: \ell_{\boldsymbol{\beta}}^\alpha(\mathbf{x}, y)
 \end{aligned}$$

$$\sup_{\|\mathbf{z}\|_p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y) = \sup_{\|\mathbf{z}\|_p \leq \alpha} \log(1 + \exp(-y \cdot \beta^{\top}(\mathbf{x} + \mathbf{z})))$$

↓

$$= \log(1 + \exp(-y \cdot \beta^{\top} \mathbf{x} + \alpha \cdot \|\beta\|_{p^*}))$$

Adversarial loss

- Can be interpreted as a new loss function
- Convex and Lipschitz
- Existing Lipschitz Wasserstein DRO theory is applicable

↓

$$= \log(1 + \exp(-y \cdot \beta^{\top} \mathbf{x} + \alpha \cdot \|\beta\|_{p^*})) =: \ell_{\beta}^{\alpha}(\mathbf{x}, y)$$

Adversarially Robust Optimization

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Adversarially RO	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} \left[\sup_{\ \mathbf{z}\ _p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y) \right]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} \left[\sup_{\ \mathbf{z}\ _p \leq \alpha} \ell_{\beta}(\mathbf{x} + \mathbf{z}, y) \right]$

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N}[\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0}[\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}}[\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0}[\ell_{\beta}(\mathbf{x}, y)]$
Adversarially RO	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N}[\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0}[\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$

Adversarially Robust Optimization

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Adversarially RO	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$

**ARO calibrates
the loss**

Adversarially Robust Optimization

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{Q \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim Q} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Adversarially RO	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$

Statistical error
 $W(\mathbb{P}_N, \mathbb{P}^0)$ stays

ARO calibrates
the loss

Distributionally & Adversarially Robust Optimization

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Adversarially RO	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$
Distributionally & Adversarially RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$

Exact Convex Reformulation

Distributionally and Adversarially Robust Optimization Problem **(DR-ARO)**

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \quad \sup_{Q \in \mathfrak{B}_\varepsilon(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim Q} \left[\sup_{\|\mathbf{z}\|_p \leq \alpha} \{ \ell_\beta(\mathbf{x} + \mathbf{z}, \mathbf{y}) \} \right]$$

Exact Convex Reformulation

Distributionally and Adversarially Robust Optimization Problem **(DR-ARO)**

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \quad \sup_{\mathbb{Q} \in \mathfrak{B}_\varepsilon(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbb{Q}} \left[\sup_{\|\mathbf{z}\|_p \leq \alpha} \{ \ell_\beta(\mathbf{x} + \mathbf{z}, \mathbf{y}) \} \right]$$

$$\begin{aligned} &\underset{\beta, \lambda, \mathbf{s}}{\text{minimize}} && \varepsilon \lambda + \frac{1}{N} \sum_{i=1}^N s_i \\ &\text{subject to} && \ell_\beta^\alpha(\mathbf{x}^i, \mathbf{y}^i) \leq s_i && \forall i \in [N] \\ &&& \ell_\beta^\alpha(\mathbf{x}^i, -\mathbf{y}^i) - \lambda \kappa \leq s_i && \forall i \in [N] \\ &&& \|\beta\|_{q^\star} \leq \lambda \\ &&& \beta \in \mathbb{R}^n, \lambda \geq 0, \mathbf{s} \in \mathbb{R}_+^N. \end{aligned}$$

Exact Convex Reformulation

Distributionally and Adversarially Robust Optimization Problem **(DR-ARO)**

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \quad \sup_{\mathbb{Q} \in \mathfrak{B}_\varepsilon(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbb{Q}} \left[\sup_{\|\mathbf{z}\|_p \leq \alpha} \{ \ell_\beta(\mathbf{x} + \mathbf{z}, \mathbf{y}) \} \right]$$

minimize
 β, λ, s

$$\varepsilon \lambda + \frac{1}{N} \sum_{i=1}^N s_i \longrightarrow \text{Linear obj.}$$

subject to

$$\ell_\beta^\alpha(\mathbf{x}^i, \mathbf{y}^i) \leq s_i \quad \forall i \in [N]$$

$$\ell_\beta^\alpha(\mathbf{x}^i, -\mathbf{y}^i) - \lambda \kappa \leq s_i \quad \forall i \in [N]$$

$$\|\beta\|_{q^\star} \leq \lambda$$

$$\beta \in \mathbb{R}^n, \lambda \geq 0, s \in \mathbb{R}_+^N.$$

Exact Convex Reformulation

Distributionally and Adversarially Robust Optimization Problem **(DR-ARO)**

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \quad \sup_{\mathbb{Q} \in \mathfrak{B}_\varepsilon(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbb{Q}} \left[\sup_{\|\mathbf{z}\|_p \leq \alpha} \{ \ell_\beta(\mathbf{x} + \mathbf{z}, \mathbf{y}) \} \right]$$

$$\begin{aligned} &\underset{\beta, \lambda, s}{\text{minimize}} && \varepsilon \lambda + \frac{1}{N} \sum_{i=1}^N s_i \\ &\text{subject to} && \ell_\beta^\alpha(\mathbf{x}^i, \mathbf{y}^i) \leq s_i \\ &&& \ell_\beta^\alpha(\mathbf{x}^i, -\mathbf{y}^i) - \lambda \kappa \leq s_i \\ &&& \|\beta\|_{q^\star} \leq \lambda \end{aligned}$$

$$\forall i \in [N]$$

$$\forall i \in [N]$$

$$\beta \in \mathbb{R}^n, \lambda \geq 0, s \in \mathbb{R}_+^N \longrightarrow \mathcal{O}(N) \text{ many}$$

Exact Convex Reformulation

Distributionally and Adversarially Robust Optimization Problem **(DR-ARO)**

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \quad \sup_{\mathbb{Q} \in \mathfrak{B}_\varepsilon(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbb{Q}} \left[\sup_{\|\mathbf{z}\|_p \leq \alpha} \{ \ell_\beta(\mathbf{x} + \mathbf{z}, \mathbf{y}) \} \right]$$

$$\begin{aligned} &\underset{\beta, \lambda, \mathbf{s}}{\text{minimize}} && \varepsilon \lambda + \frac{1}{N} \sum_{i=1}^N s_i \\ &\text{subject to} && \ell_\beta^\alpha(\mathbf{x}^i, \mathbf{y}^i) \leq s_i && \forall i \in [N] \\ &&& \ell_\beta^\alpha(\mathbf{x}^i, -\mathbf{y}^i) - \lambda \kappa \leq s_i && \forall i \in [N] \\ &\text{Convex for} && \|\beta\|_{q^\star} \leq \lambda \\ &q \in \{1, 2, \infty\} && \beta \in \mathbb{R}^n, \lambda \geq 0, \mathbf{s} \in \mathbb{R}_+^N. \end{aligned}$$

Exact Convex Reformulation

Distributionally and Adversarially Robust Optimization Problem **(DR-ARO)**

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \quad \sup_{\mathbb{Q} \in \mathfrak{B}_\varepsilon(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbb{Q}} \left[\sup_{\|\mathbf{z}\|_p \leq \alpha} \{ \ell_\beta(\mathbf{x} + \mathbf{z}, \mathbf{y}) \} \right]$$

$$\begin{aligned} &\underset{\beta, \lambda, \mathbf{s}}{\text{minimize}} && \varepsilon \lambda + \frac{1}{N} \sum_{i=1}^N s_i \\ &\text{subject to} && \ell_\beta^\alpha(\mathbf{x}^i, \mathbf{y}^i) \leq s_i && \forall i \in [N] \\ &\text{Exponential} && \ell_\beta^\alpha(\mathbf{x}^i, -\mathbf{y}^i) - \lambda \kappa \leq s_i && \forall i \in [N] \\ &\text{cone repr.} && \|\beta\|_{q^\star} \leq \lambda \\ &&& \beta \in \mathbb{R}^n, \lambda \geq 0, \mathbf{s} \in \mathbb{R}_+^N. \end{aligned}$$

Exact Convex Reformulation

Distributionally and Adversarially Robust Optimization Problem **(DR-ARO)**

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \quad \sup_{\mathbb{Q} \in \mathfrak{B}_\varepsilon(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathbb{Q}} \left[\sup_{\|\mathbf{z}\|_p \leq \alpha} \{ \ell_\beta(\mathbf{x} + \mathbf{z}, \mathbf{y}) \} \right]$$

$$\begin{aligned} &\underset{\beta, \lambda, \mathbf{s}}{\text{minimize}} && \varepsilon \lambda + \frac{1}{N} \sum_{i=1}^N s_i \\ &\text{subject to} && \ell_\beta^\alpha(\mathbf{x}^i, \mathbf{y}^i) \leq s_i \\ &&& \ell_\beta^\alpha(\mathbf{x}^i, -\mathbf{y}^i) - \lambda \kappa \leq s_i \\ &&& \|\beta\|_{q^\star} \leq \lambda \end{aligned}$$

$$\beta \in \mathbb{R}^n, \lambda \geq 0, \mathbf{s} \in \mathbb{R}_+^N.$$

- Adversarial loss being ℓ_β^α
- ℓ_β^α being convex & Lipschitz
- Shafieezadeh-Abadeh (2019)

Distributionally & Adversarially Robust Optimization

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Adversarially RO	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$
Distributionally & Adversarially RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$

How we address Robust Overfitting

Distributionally & Adversarially Robust Optimization

Paradigm	Training Risk	True Risk
Empirical Risk Min.	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Distributionally RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\ell_{\beta}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}(\mathbf{x}, y)]$
Adversarially RO	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_N} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$
Distributionally & Adversarially RO	$\sup_{\mathbb{Q} \in \mathfrak{B}_{\varepsilon}(\mathbb{P}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{Q}} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$	$\mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}^0} [\ell_{\beta}^{\alpha}(\mathbf{x}, y)]$

What other approaches?

Data

$$\{\xi^i = (x^i, y^i)\}_{i \in [N]}$$

$\xrightarrow{\sim \text{iid } \mathbb{P}^0}$

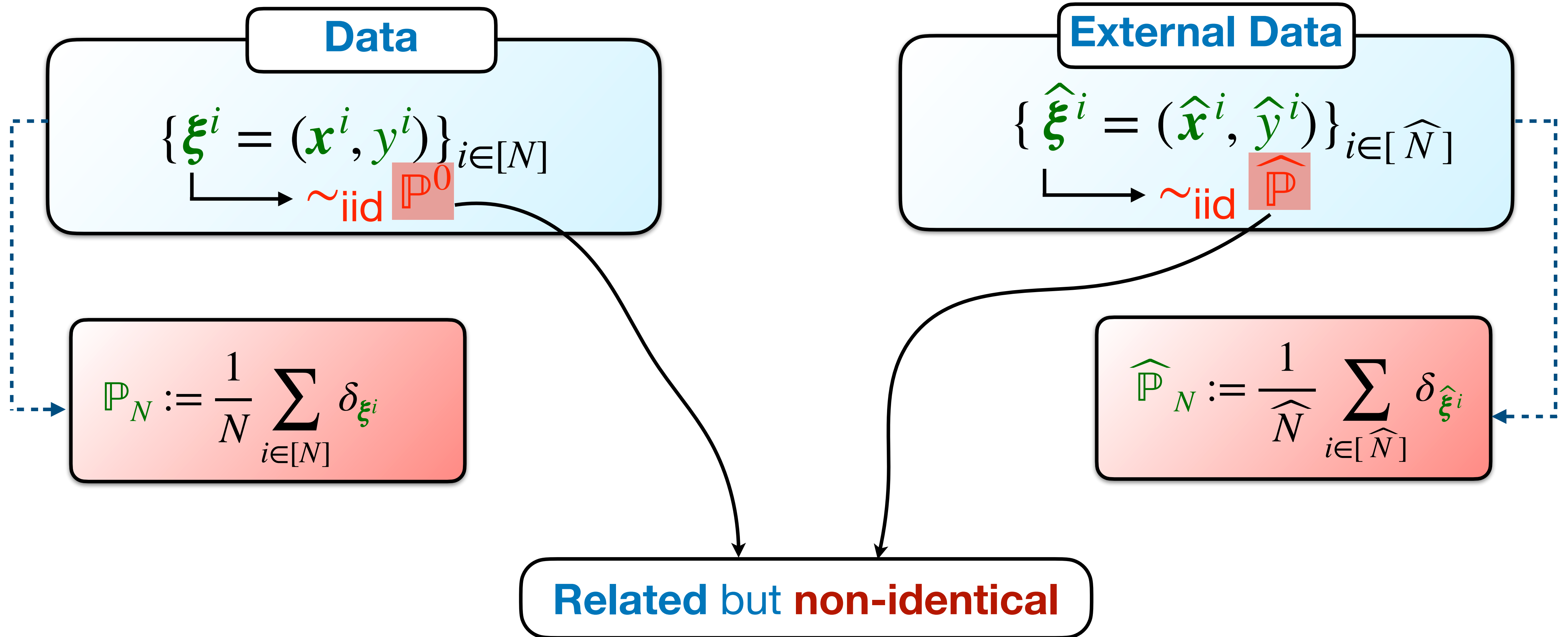
$$\mathbb{P}_N := \frac{1}{N} \sum_{i \in [N]} \delta_{\xi^i}$$

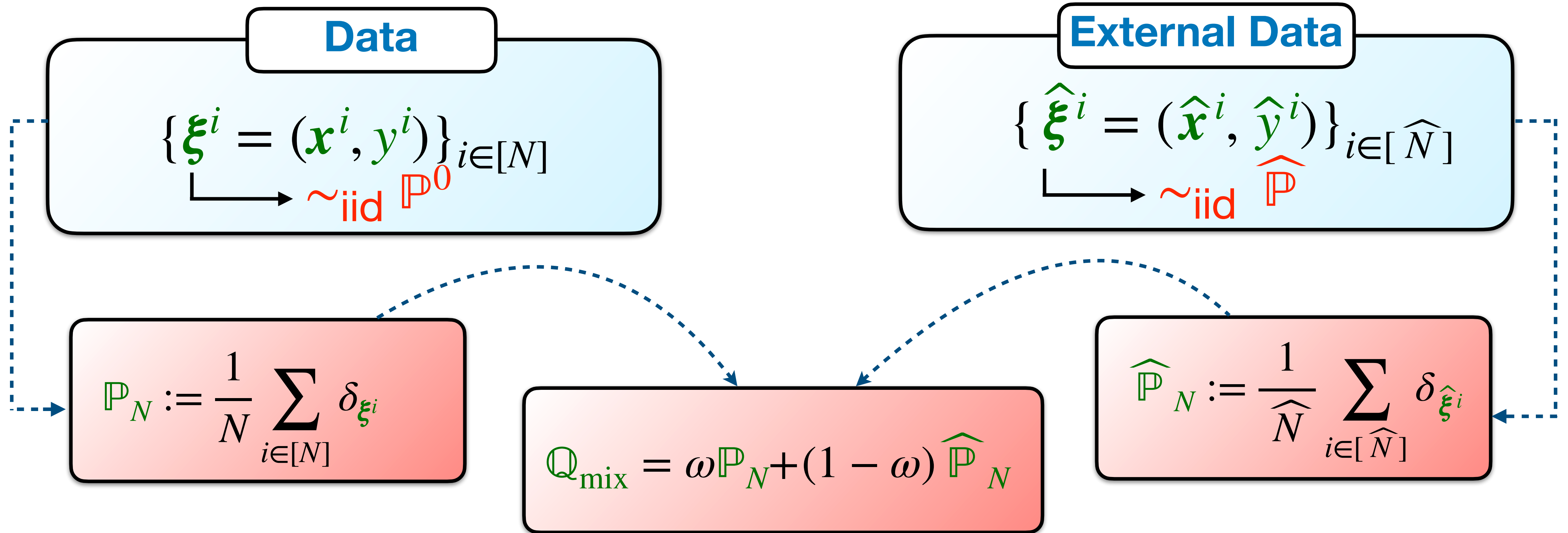
External Data

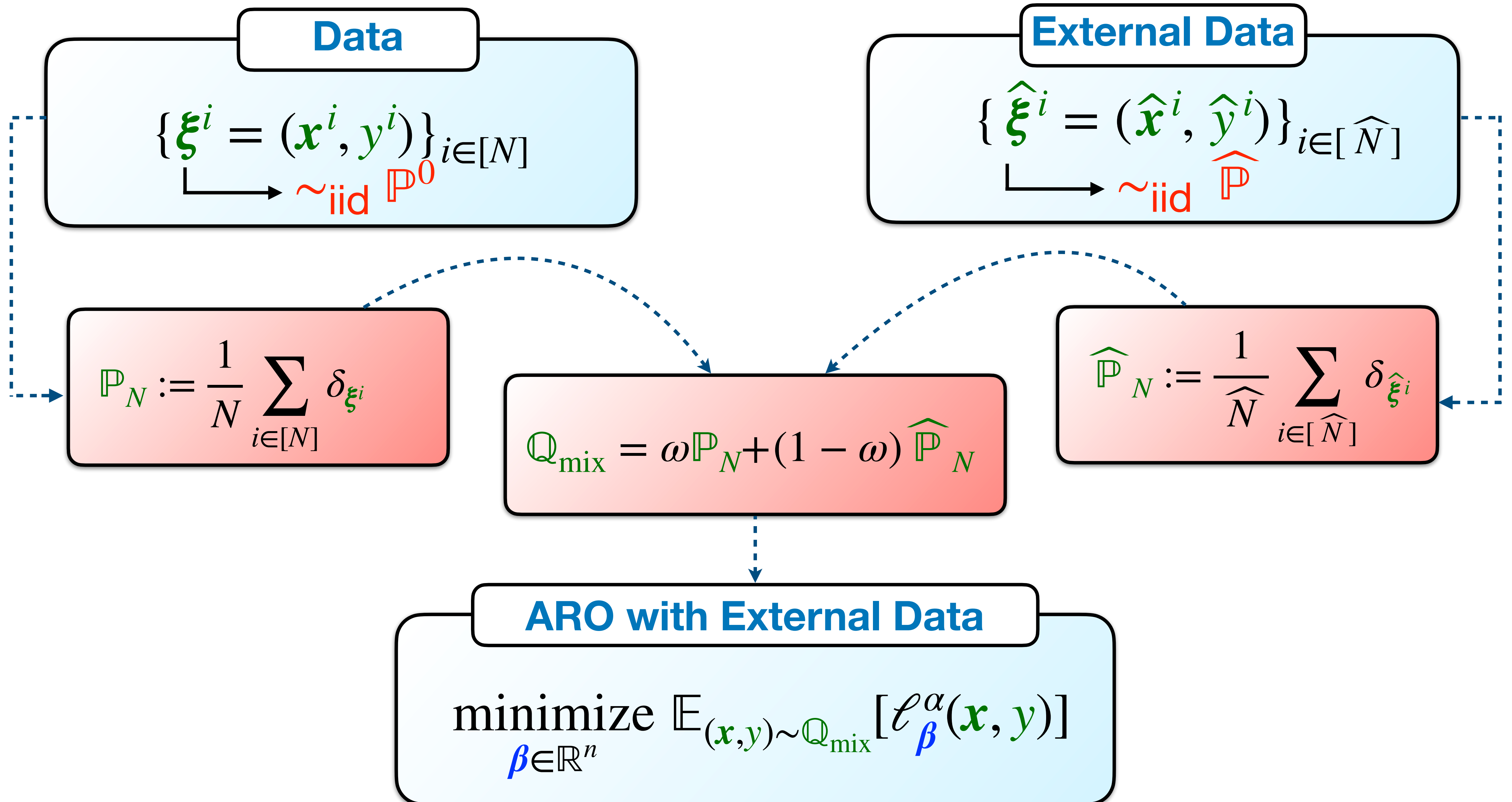
$$\{\hat{\xi}^i = (\hat{x}^i, \hat{y}^i)\}_{i \in [\hat{N}]}$$

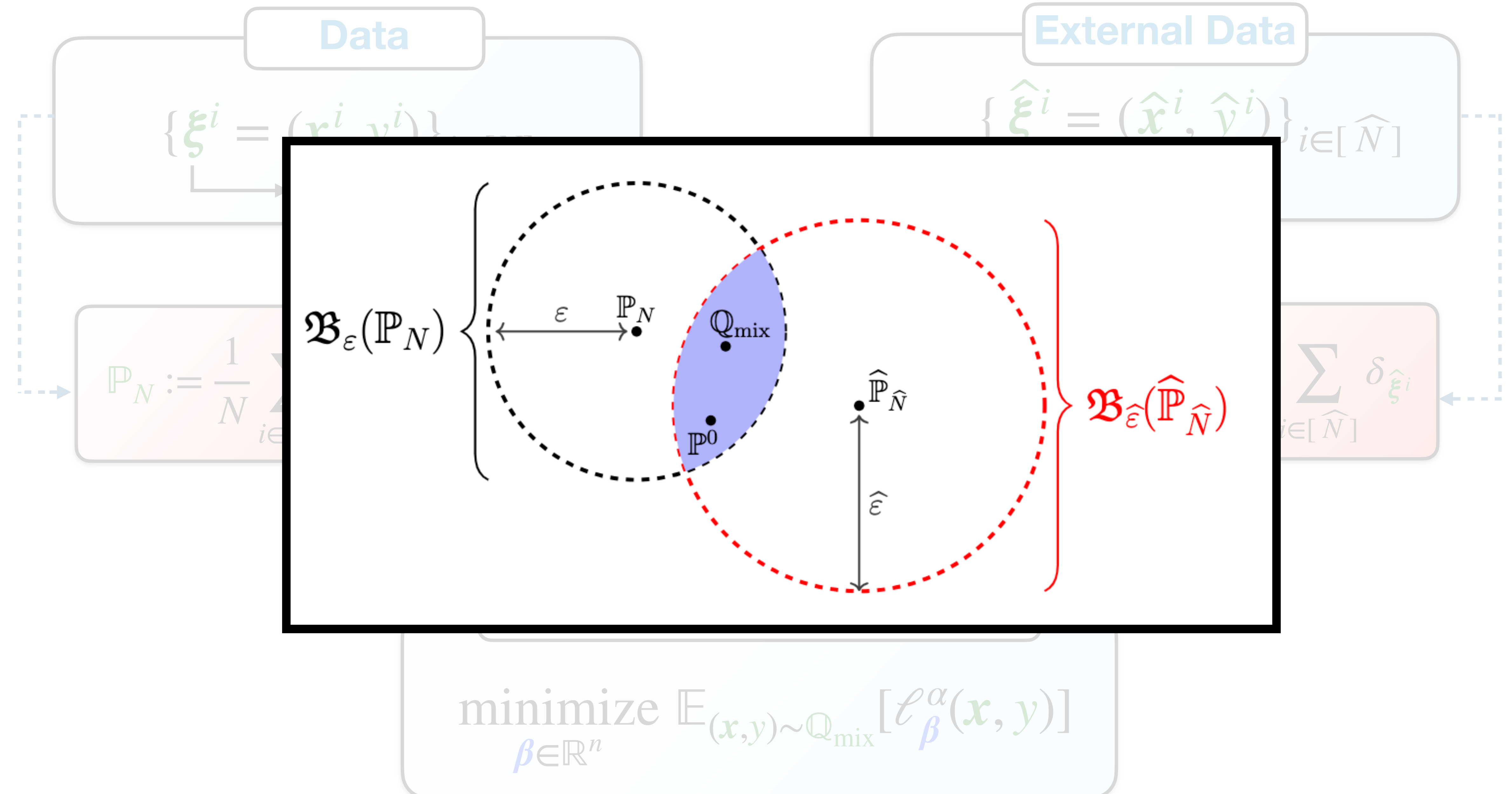
$\xrightarrow{\sim \text{iid } \hat{\mathbb{P}}}$

$$\hat{\mathbb{P}}_N := \frac{1}{\hat{N}} \sum_{i \in [\hat{N}]} \delta_{\hat{\xi}^i}$$









Exact Reformulation

ARO over intersection of Wasserstein balls **(Inter-ARO)**:

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \quad \sup_{Q \in \mathfrak{B}_\varepsilon(\mathbb{P}_N) \cap \mathfrak{B}_{\hat{\varepsilon}}(\hat{\mathbb{P}}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim Q} \left[\sup_{\|\mathbf{z}\|_p \leq \alpha} \{\ell_\beta(\mathbf{x} + \mathbf{z}, y)\} \right]$$

Exact Reformulation

ARO over intersection of Wasserstein balls (**Inter-ARO**):

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \quad \sup_{Q \in \mathfrak{B}_\varepsilon(\mathbb{P}_N) \cap \mathfrak{B}_{\hat{\varepsilon}}(\hat{\mathbb{P}}_N)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim Q} \left[\sup_{\|\mathbf{z}\|_p \leq \alpha} \{\ell_\beta(\mathbf{x} + \mathbf{z}, \mathbf{y})\} \right]$$

1

NP-hard even if $N = 1$ and $\hat{N} = 1$.

Exact Reformulation

ARO over intersection of Wasserstein balls (**Inter-ARO**):

$$\underset{\beta \in \mathbb{R}^n}{\text{minimize}} \quad \sup_{Q \in \mathfrak{B}_\varepsilon(\mathbb{P}_N) \cap \mathfrak{B}_{\hat{\varepsilon}}(\hat{\mathbb{P}}_N)} \mathbb{E}_{(\mathbf{x}, y) \sim Q} \left[\sup_{\|\mathbf{z}\|_p \leq \alpha} \{\ell_\beta(\mathbf{x} + \mathbf{z}, y)\} \right]$$

1

NP-hard even if $N = 1$ and $\hat{N} = 1$.

2

Would admit an exact **tractable** reformulation if

- Squared-loss function (regression)
- Wasserstein ball around first and second moments
- No attack ($\alpha = 0$)

Taskesen et al. (2021)

Static Relaxation Technique

We consult to the **adjustable RO** literature for a **relaxation** that is:

- Convex with $\mathcal{O}(N \cdot \widehat{N})$ exponential conic constraints

Static Relaxation Technique

We consult to the **adjustable RO** literature for a **relaxation** that is:

- Convex with $\mathcal{O}(N \cdot \hat{N})$ exponential conic constraints
- Coincides with the **exact** formulation when $\hat{\varepsilon} \rightarrow \infty$

Static Relaxation Technique

We consult to the **adjustable RO** literature for a **relaxation** that is:

- Convex with $\mathcal{O}(N \cdot \hat{N})$ exponential conic constraints
- Coincides with the **exact** formulation when $\hat{\varepsilon} \rightarrow \infty$
- “**Not learning** from the external distribution” is **always feasible**

Static Relaxation Technique

We consult to the **adjustable RO** literature for a **relaxation** that is:

- Convex with $\mathcal{O}(N \cdot \hat{N})$ exponential conic constraints
- Coincides with the **exact** formulation when $\hat{\varepsilon} \rightarrow \infty$
- “**Not learning** from the external distribution” is **always feasible**
- Only uses the **external data** if it improves the objective

Static Relaxation Technique

We consult to the **adjustable RO** literature for a **relaxation** that is:

- Convex with $\mathcal{O}(N \cdot \hat{N})$ exponential conic constraints
- Coincides with the **exact** formulation when $\hat{\varepsilon} \rightarrow \infty$
- “**Not learning** from the external distribution” is **always feasible**
- Only uses the **external data** if it improves the objective
- Recovers the presented model from the literature as a **special case**

3 Sets of Numerical Experiments

Artificial experiments

- External data is **artificially** generated
- Direct control over distributions

UCI experiments

- Most popular UCI classification datasets
- External data via **synthetic** data generation

MNIST experiments

- Digit recognition (e.g., 3 vs 9)
- External data is digits of high school students

Attack	<u>ERM</u>	<u>ARO</u>	<u>ARO+Aux</u>	<u>DRO+ARO</u>	<u>DRO+ARO+Aux</u>
No attack ($\alpha = 0$)	1.55%	1.55%	1.19%	0.64%	0.53%
ℓ_1 ($\alpha = 68/255$)	2.17%	1.84%	1.33%	0.66%	0.57%
ℓ_2 ($\alpha = 128/255$)	99.93%	3.36%	2.54%	2.40%	2.12%
ℓ_∞ ($\alpha = 8/255$)	100.00%	2.60%	2.38%	2.20%	1.95%

Our DRO models

Attack	<u>ERM</u>	<u>ARO</u>	<u>ARO+Aux</u>	<u>DRO+ARO</u>	<u>DRO+ARO+Aux</u>
No attack ($\alpha = 0$)	1.55%	1.55%	1.19%	0.64%	0.53%
ℓ_1 ($\alpha = 68/255$)	2.17%	1.84%	1.33%	0.66%	0.57%
ℓ_2 ($\alpha = 128/255$)	99.93%	3.36%	2.54%	2.40%	2.12%
ℓ_∞ ($\alpha = 8/255$)	100.00%	2.60%	2.38%	2.20%	1.95%

Adversarially robust optimization

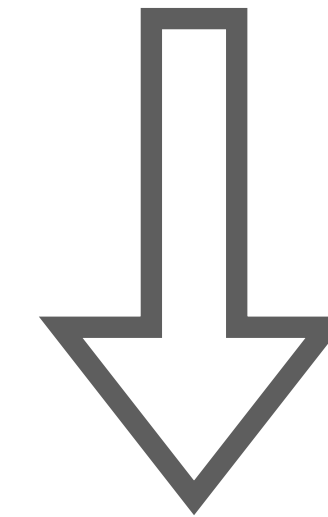
Attack	<u>ERM</u>	<u>ARO</u>	<u>ARO+Aux</u>	<u>DRO+ARO</u>	<u>DRO+ARO+Aux</u>
No attack ($\alpha = 0$)	1.55%	1.55%	1.19%	0.64%	0.53%
ℓ_1 ($\alpha = 68/255$)	2.17%	1.84%	1.33%	0.66%	0.57%
ℓ_2 ($\alpha = 128/255$)	99.93%	3.36%	2.54%	2.40%	2.12%
ℓ_∞ ($\alpha = 8/255$)	100.00%	2.60%	2.38%	2.20%	1.95%

Adversarially robust optimization (over its mixture with external data)

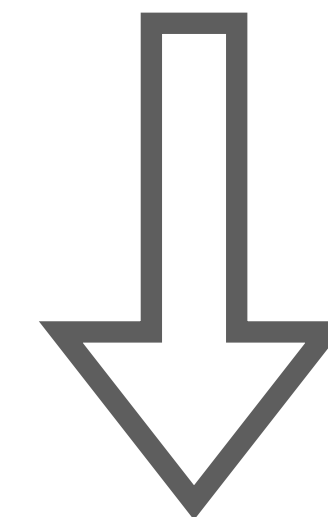
Attack	<u>ERM</u>	<u>ARO</u>	<u>ARO+Aux</u>	<u>DRO+ARO</u>	<u>DRO+ARO+Aux</u>
No attack ($\alpha = 0$)	1.55%	1.55%	1.19%	0.64%	0.53%
ℓ_1 ($\alpha = 68/255$)	2.17%	1.84%	1.33%	0.66%	0.57%
ℓ_2 ($\alpha = 128/255$)	99.93%	3.36%	2.54%	2.40%	2.12%
ℓ_∞ ($\alpha = 8/255$)	100.00%	2.60%	2.38%	2.20%	1.95%

Thank you for listening!

IMPERIAL



**PRINCETON
UNIVERSITY**



Personal Webpage



Future Work

- Different loss functions
- Intersection of more balls
- Comparison of different relaxation techniques
- Specialised algorithms for $\ell_1, \ell_2, \ell_\infty$ norms in the feature-label metric
- Derive your own algorithm instead of using MOSEK
- Ball around \mathbb{Q}_{mix} directly

Key Lemma for Tractability

$L(z)$ is convex with $\text{lip}(L) = 1$, $\omega, a \in \mathbb{R}^n$ and $\lambda > 0$. Then:

$$\sup_{x \in \mathbb{R}^n} \{L(\omega^\top x) - \lambda \|a - x\|_q\}$$

Key Lemma for Tractability

$L(z)$ is convex with $\text{lip}(L) = 1$, $\omega, a \in \mathbb{R}^n$ and $\lambda > 0$. Then:

$$\sup_{x \in \mathbb{R}^n} \{L(\omega^\top x) - \lambda \|a - x\|_q\} \\ = \begin{cases} L(a^\top \omega) & \text{if } \|\omega\|_{q^*} \leq \lambda \\ +\infty & \text{otherwise.} \end{cases}$$

Key Lemma for Tractability

$L(z)$ is convex with $\text{lip}(L) = 1$, $\omega, a \in \mathbb{R}^n$ and $\lambda > 0$. Then:

$$\sup_{x \in \mathbb{R}^n} \{L(\omega^\top x) - \lambda \|a - x\|_q\}$$

**DC
Maximization**

$$= \begin{cases} L(a^\top \omega) & \text{if } \|\omega\|_{q^*} \leq \lambda \\ +\infty & \text{otherwise.} \end{cases}$$

Key Lemma for Tractability

$L(z)$ is convex with $\text{lip}(L) = 1$, $\omega, a \in \mathbb{R}^n$ and $\lambda > 0$. Then:

$$\sup_{x \in \mathbb{R}^n} \{L(\omega^\top x) - \lambda \|a - x\|_q\}$$

Convex fn.
in ω

$$= \begin{cases} L(a^\top \omega) & \text{if } \|\omega\|_{q^*} \leq \lambda \\ +\infty & \text{otherwise.} \end{cases}$$

Convex
constraint on ω

Key Reason for non-tractability

$L(z)$ is convex with $\text{lip}(L) = 1$, $\omega, a, \hat{a} \in \mathbb{R}^n$ and $\lambda, \hat{\lambda} > 0$. Then:

$$\sup_{x \in \mathbb{R}^n} \{L(\omega^\top x) - \lambda \|a - x\|_q - \hat{\lambda} \|\hat{a} - x\|_q\}$$

Key Reason for non-tractability

$L(z)$ is convex with $\text{lip}(L) = 1$, $\omega, a, \hat{a} \in \mathbb{R}^n$ and $\lambda, \hat{\lambda} > 0$. Then:

$$\sup_{x \in \mathbb{R}^n} \{L(\omega^\top x) - \lambda \|a - x\|_q - \hat{\lambda} \|\hat{a} - x\|_q\}$$

$$=: g(\omega)$$

Key Reason for non-tractability

$L(z)$ is convex with $\text{lip}(L) = 1$, $\omega, a, \hat{a} \in \mathbb{R}^n$ and $\lambda, \hat{\lambda} > 0$. Then:

$$\sup_{x \in \mathbb{R}^n} \{L(\omega^\top x) - \lambda \|a - x\|_q - \hat{\lambda} \|\hat{a} - x\|_q\}$$

$$=: g(\omega)$$

We have $\mathcal{O}(N \cdot \hat{N})$ constraints of type $g(\omega) \leq \text{constant}$

Key Reason for non-tractability

$L(z)$ is convex with $\text{lip}(L) = 1$, $\omega, a, \hat{a} \in \mathbb{R}^n$ and $\lambda, \hat{\lambda} > 0$. Then:

$$\sup_{x \in \mathbb{R}^n} \{ L(\omega^\top x) - \lambda \|a - x\|_q - \hat{\lambda} \|\hat{a} - x\|_q \}$$

**DC
Maximization**

Key Reason for non-tractability

$L(z)$ is convex with $\text{lip}(L) = 1$, $\omega, a, \hat{a} \in \mathbb{R}^n$ and $\lambda, \hat{\lambda} > 0$. Then:

$$\begin{aligned} & \sup_{x \in \mathbb{R}^n} \{ L(\omega^\top x) - \lambda \|a - x\|_q - \hat{\lambda} \|\hat{a} - x\|_q \} \\ &= \sup_{\theta \in \text{dom}(L^*)} -L^*(\theta) + \theta \cdot \omega^\top a + \\ & \quad \inf_{z \in \mathbb{R}^n} \{ \theta \cdot z^\top (\hat{a} - a) : |\theta| \cdot \|\omega - z\|_{q^*} \leq \lambda, |\theta| \cdot \|z\|_{q^*} \leq \hat{\lambda} \} \end{aligned}$$

Key Reason for non-tractability

$L(z)$ is convex with $\text{lip}(L) = 1$, $\omega, a, \hat{a} \in \mathbb{R}^n$ and $\lambda, \hat{\lambda} > 0$. Then:

$$\sup_{x \in \mathbb{R}^n} \{L(\omega^\top x) - \lambda \|a - x\|_q - \hat{\lambda} \|\hat{a} - x\|_q\}$$

$$= \sup_{\theta \in \text{dom}(L^*)} -L^*(\theta) + \theta \cdot \omega^\top a +$$

$$\inf_{z \in \mathbb{R}^n} \{ \theta \cdot z^\top (\hat{a} - a) : |\theta| \cdot \|\omega - z\|_{q^*} \leq \lambda, |\theta| \cdot \|z\|_{q^*} \leq \hat{\lambda} \}$$

**Minimax theorem
not applicable**

Key Reason for non-tractability

$L(z)$ is convex with $\text{lip}(L) = 1$, $\omega, a, \hat{a} \in \mathbb{R}^n$ and $\lambda, \hat{\lambda} > 0$. Then:

$$\sup_{x \in \mathbb{R}^n} \{ L(\omega^\top x) - \lambda \|a - x\|_q - \hat{\lambda} \|\hat{a} - x\|_q \}$$

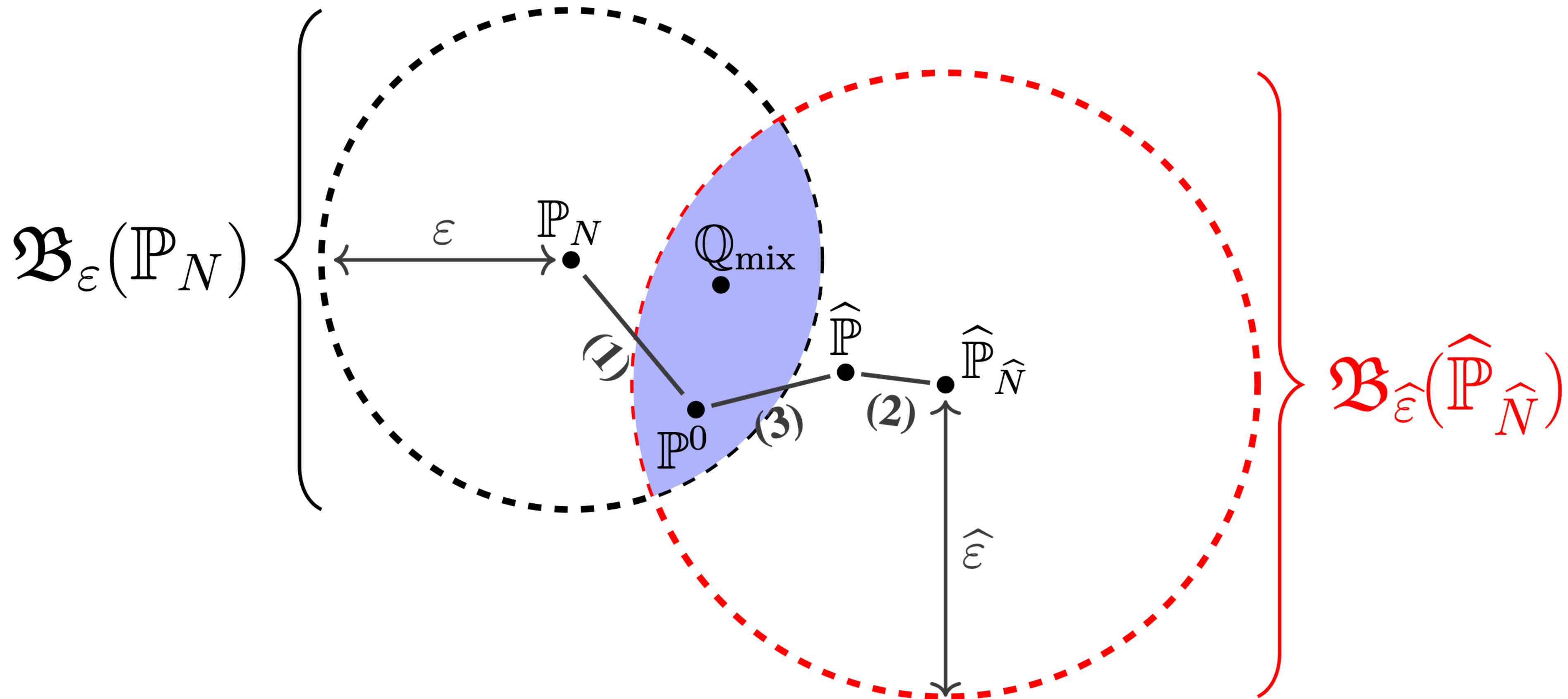
$$= \sup_{\theta \in \text{dom}(L^*)} -L^*(\theta) + \theta \cdot \omega^\top a +$$

$$\inf_{z \in \mathbb{R}^n} \{ \theta \cdot z^\top (\hat{a} - a) : |\theta| \cdot \|\omega - z\|_{q^*} \leq \lambda, |\theta| \cdot \|z\|_{q^*} \leq \hat{\lambda} \}$$

$\forall \theta$

$\exists z$

Can be viewed as an adjustable RO constraint



Case 1: Fix $\hat{\varepsilon} \rightarrow \infty$

Theorem 6.1 (abridged). *For light-tailed \mathbb{P}^0 , if $\varepsilon \geq \mathcal{O}(\frac{\log(\eta)}{N})^{1/n}$ for $\eta \in (0, 1)$, then:*

- $\mathbb{P}^0 \in \mathfrak{B}_\varepsilon(\mathbb{P}_N)$ with $1 - \eta$ confidence;
- *AdvDRO* overestimates true loss with $1 - \eta$ confidence;
- *AdvDRO* is asymptotically consistent \mathbb{P}^0 -a.s.;
- Worst case distributions for optimal solutions of *AdvDRO* are supported on at most $N + 1$ outcomes.

Case 2: Simultaneous

Theorem 6.2 (abridged). *For light-tailed \mathbb{P}^0 and $\hat{\mathbb{P}}$, if $\varepsilon \geq \mathcal{O}(\frac{\log(\eta_1)}{N})^{1/n}$ and $\hat{\varepsilon} \geq W(\mathbb{P}^0, \hat{\mathbb{P}}) + \mathcal{O}(\frac{\log(\eta_2)}{\hat{N}})^{1/n}$ for $\eta_1, \eta_2 \in (0, 1)$ with $\eta := \eta_1 + \eta_2 < 1$, then:*

- $\mathbb{P}^0 \in \mathfrak{B}_\varepsilon(\mathbb{P}_N) \cap \mathfrak{B}_{\hat{\varepsilon}}(\hat{\mathbb{P}}_{\hat{N}})$ with $1 - \eta$ confidence;
- *Synth* overestimates true loss with $1 - \eta$ confidence.

Case 2: Simultaneous

Big assumption!

Theorem 6.2 (abridged). *For light-tailed \mathbb{P}^0 and $\hat{\mathbb{P}}$, if $\varepsilon \geq \mathcal{O}(\frac{\log(\eta_1)}{N})^{1/n}$ and $\hat{\varepsilon} \geq \boxed{W(\mathbb{P}^0, \hat{\mathbb{P}})} + \mathcal{O}(\frac{\log(\eta_2)}{\hat{N}})^{1/n}$ for $\eta_1, \eta_2 \in (0, 1)$ with $\eta := \eta_1 + \eta_2 < 1$, then:*

- $\mathbb{P}^0 \in \mathfrak{B}_\varepsilon(\mathbb{P}_N) \cap \mathfrak{B}_{\hat{\varepsilon}}(\hat{\mathbb{P}}_{\hat{N}})$ with $1 - \eta$ confidence;
- *Synth* overestimates true loss with $1 - \eta$ confidence.

Case 2: Simultaneous

Big assumption!

Theorem 6.2 (abridged). *For light-tailed \mathbb{P}^0 and $\hat{\mathbb{P}}$, if $\varepsilon \geq \mathcal{O}(\frac{\log(\eta_1)}{N})^{1/n}$ and $\hat{\varepsilon} \geq \boxed{W(\mathbb{P}^0, \hat{\mathbb{P}})} + \mathcal{O}(\frac{\log(\eta_2)}{\hat{N}})^{1/n}$ for $\eta_1, \eta_2 \in (0, 1)$ with $\eta := \eta_1 + \eta_2 < 1$, then:*

- $\mathbb{P}^0 \in \mathfrak{B}_\varepsilon(\mathbb{P}_N) \cap \mathfrak{B}_{\hat{\varepsilon}}(\hat{\mathbb{P}}_{\hat{N}})$ with $1 - \eta$ confidence;
- *Synth* overestimates true loss with $1 - \eta$ confidence.

What can be done beyond cross-validation?

1. Uber vs Lyft ([Taskesen et al, 2021](#))
2. Opt-out data with differential privacy ([Ullman and Vadhan, 2020](#))
3. Wasserstein GANs comes with guarantees on $W(\mathbb{P}_N, \hat{\mathbb{P}})$