

Sequential Domain Adaptation by Synthesizing Distributionally Robust Experts

Bahar Taşkesen

Risk Analytics and Optimization Chair
École Polytechnique Fédérale de Lausanne
rao.epfl.ch

Joint work with

Man-Chung Yue



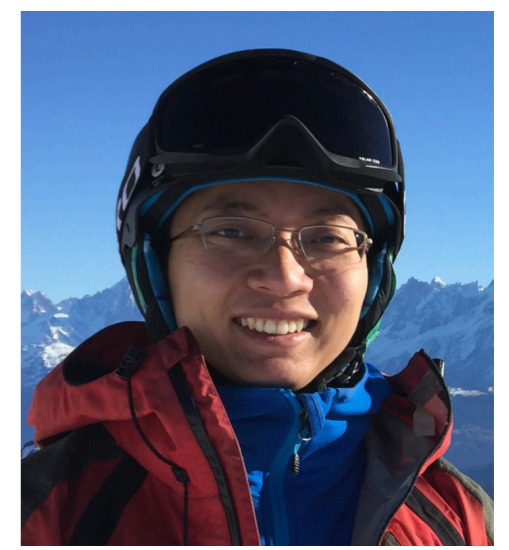
Jose Blanchet



Daniel Kuhn



Viet Anh Nguyen



Domain Adaptation



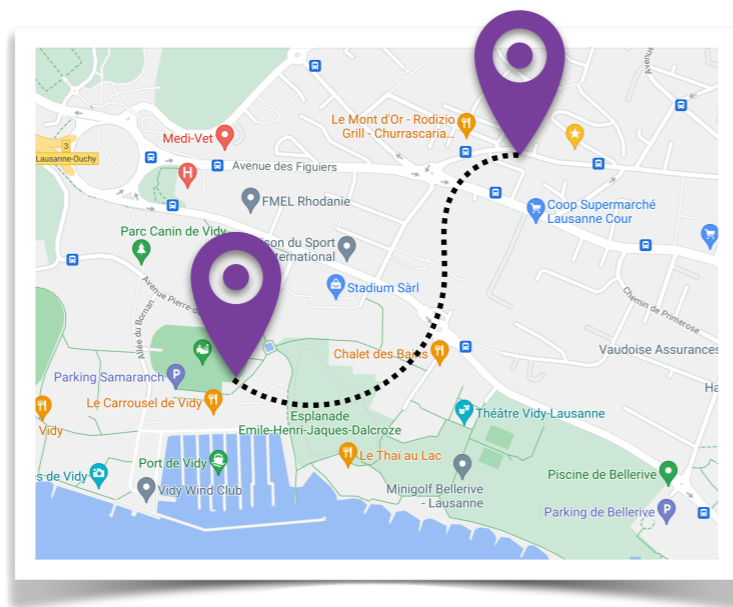
Supervised Domain Adaptation



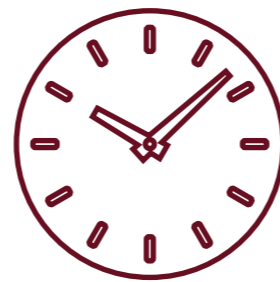
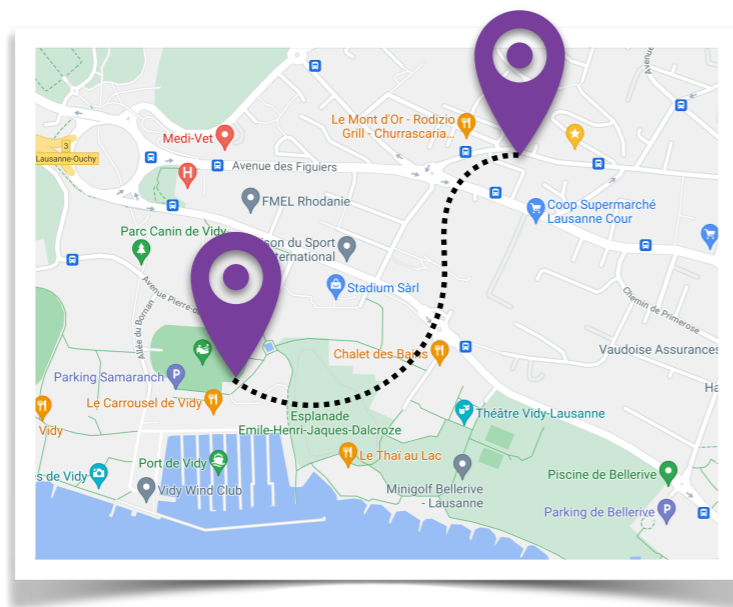
Supervised Domain Adaptation



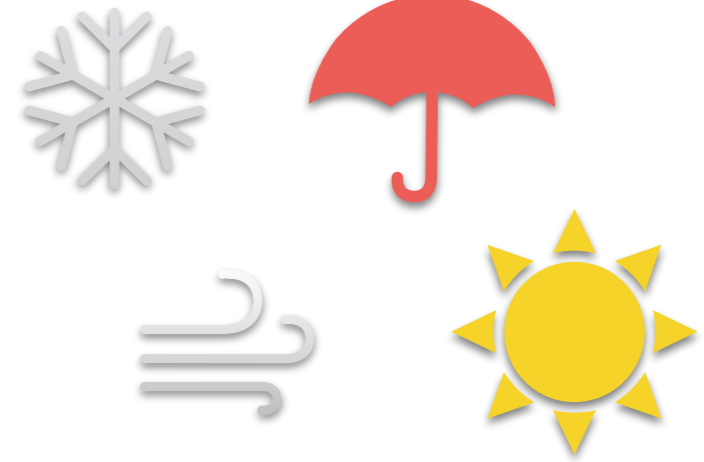
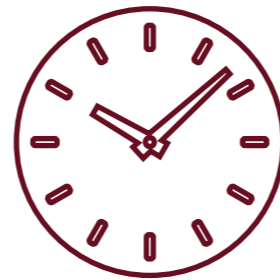
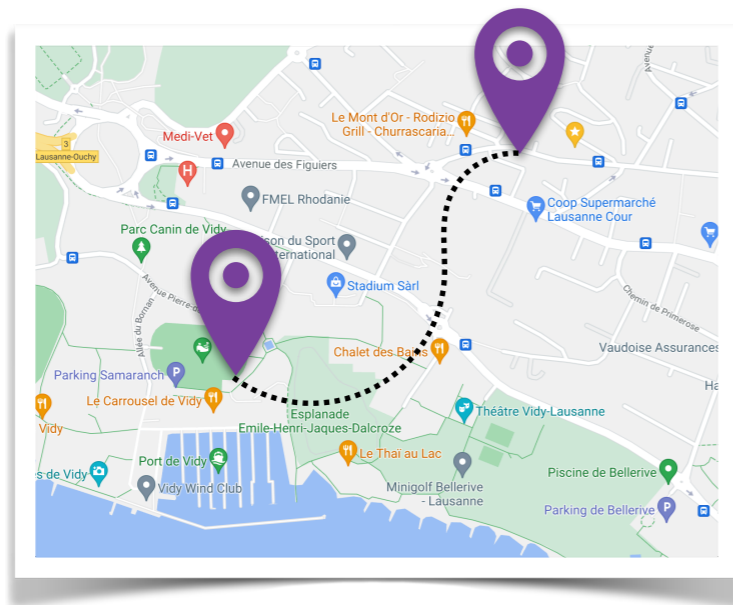
Supervised Domain Adaptation



Supervised Domain Adaptation



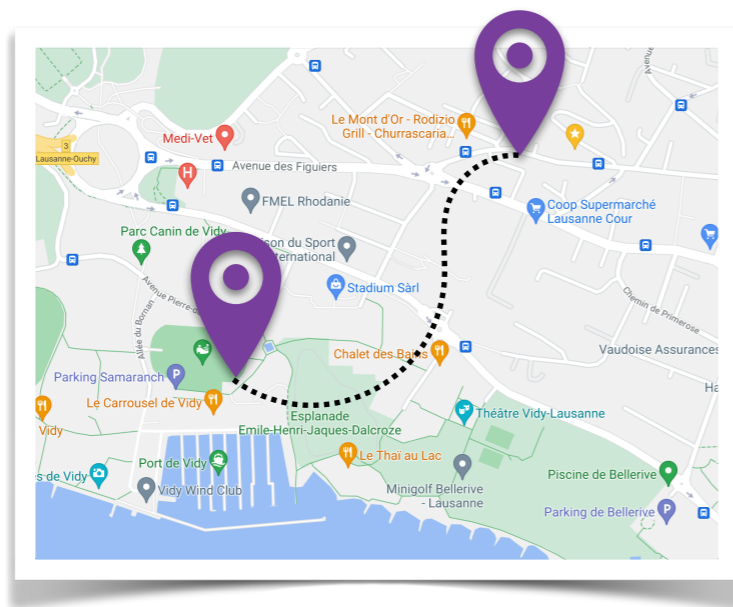
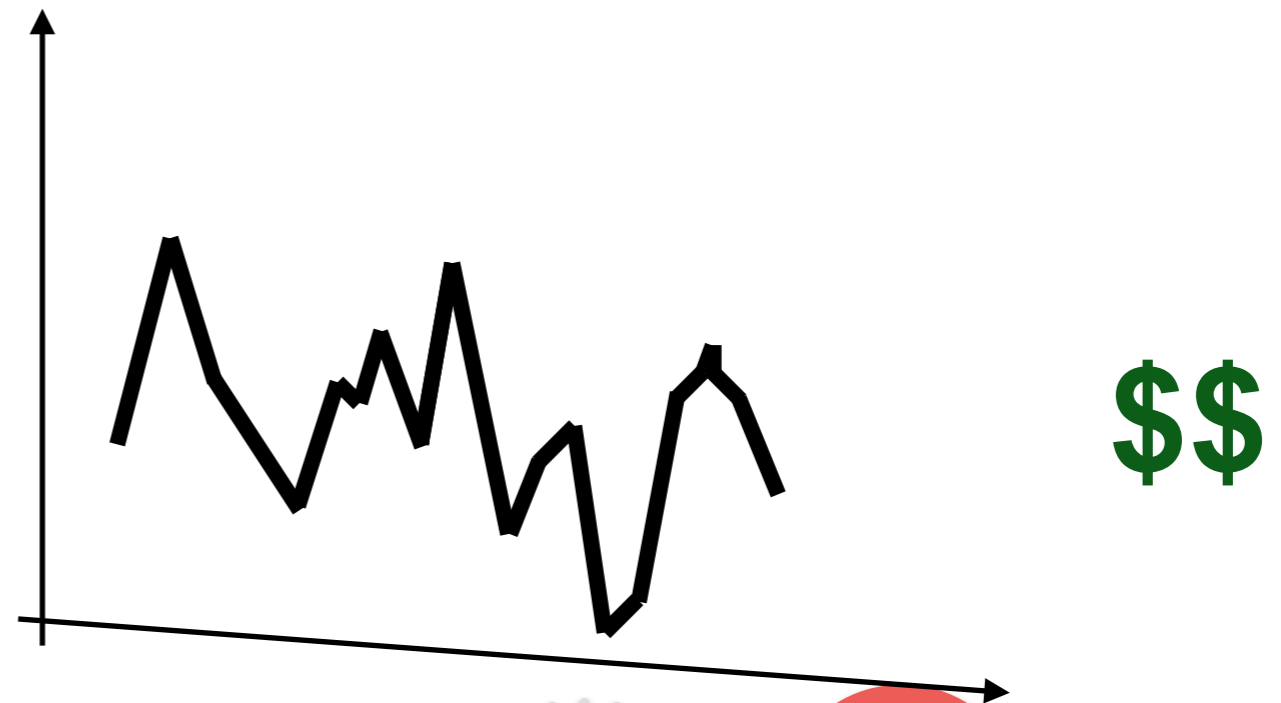
Supervised Domain Adaptation



Supervised Domain Adaptation

Uber

lyft

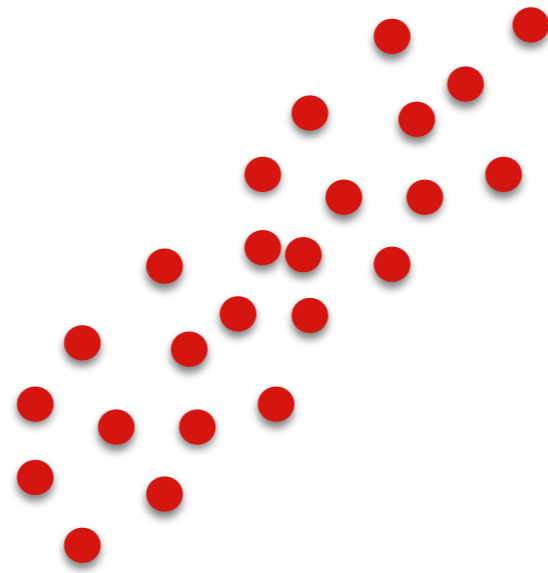


Ridge Regression

X : d -dimensional covariate

Y : univariate response

N : number of samples



$$(\hat{x}_j, \hat{y}_j)_{j=1}^N$$

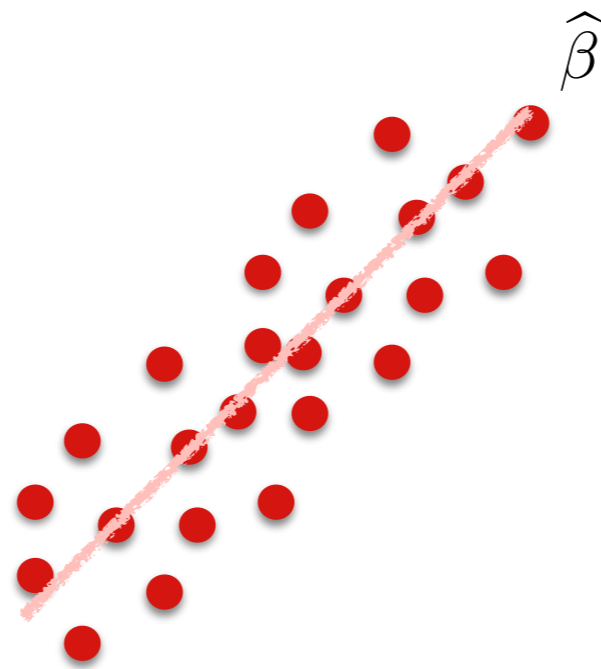
Ridge Regression

X : d -dimensional covariate

Y : univariate response

N : number of samples

$$\min_{\beta \in \mathbb{R}^d} \frac{1}{N} \sum_{j=1}^N (\beta^\top \hat{x}_j - \hat{y}_j)^2 + \eta \|\beta\|_2^2 \quad \longrightarrow \quad \hat{\beta} = \left(\frac{1}{N} \sum_{j=1}^N \hat{x}_j \hat{x}_j^\top + \eta I_d \right)^{-1} \left(\frac{1}{N} \sum_{j=1}^N \hat{x}_j \hat{y}_j \right)$$



$(\hat{x}_j, \hat{y}_j)_{j=1}^N$

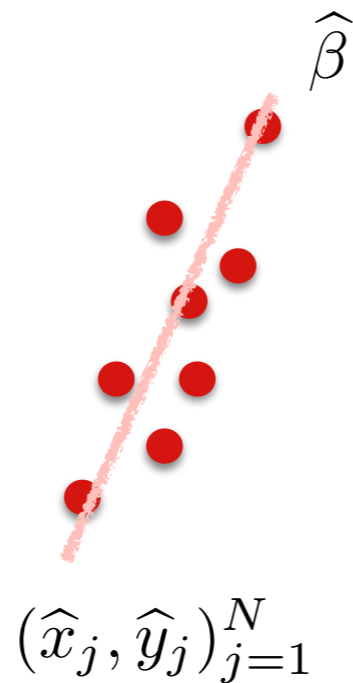
Ridge Regression - Scarce data

X : d -dimensional covariate

Y : univariate response

N : number of samples

$$\min_{\beta \in \mathbb{R}^d} \frac{1}{N} \sum_{j=1}^N (\beta^\top \hat{x}_j - \hat{y}_j)^2 + \eta \|\beta\|_2^2 \quad \rightarrow \quad \hat{\beta} = \left(\frac{1}{N} \sum_{j=1}^N \hat{x}_j \hat{x}_j^\top + \eta I_d \right)^{-1} \left(\frac{1}{N} \sum_{j=1}^N \hat{x}_j \hat{y}_j \right)$$

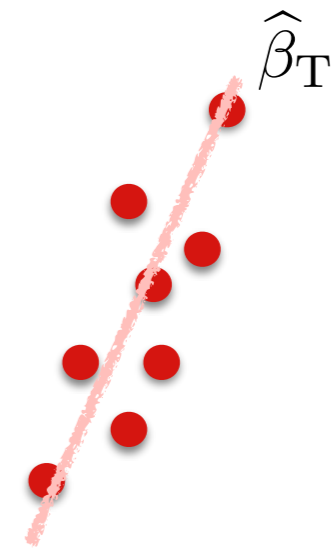
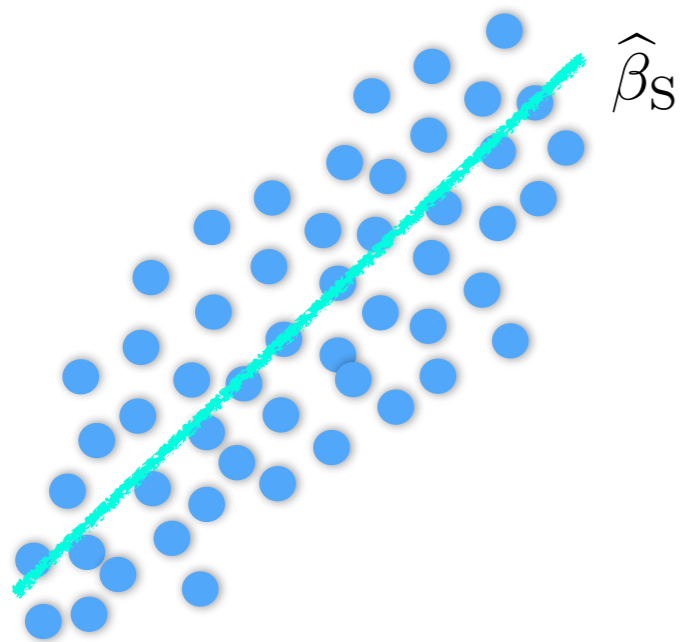


Supervised Domain Adaptation

$$\min_{\beta \in \mathbb{R}^d} \frac{1}{N} \sum_{j=1}^N (\beta^\top \hat{x}_j - \hat{y}_j)^2 + \eta \|\beta\|_2^2$$

Source

Target



Abundant **labelled** data

Scarce **labelled** data

$$(\hat{x}_i, \hat{y}_i)_{i=1}^{N_S}$$

$$(\hat{x}_j, \hat{y}_j)_{j=1}^{N_T}$$

Supervised Domain Adaptation

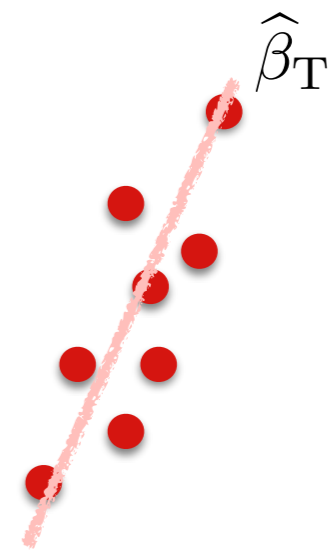
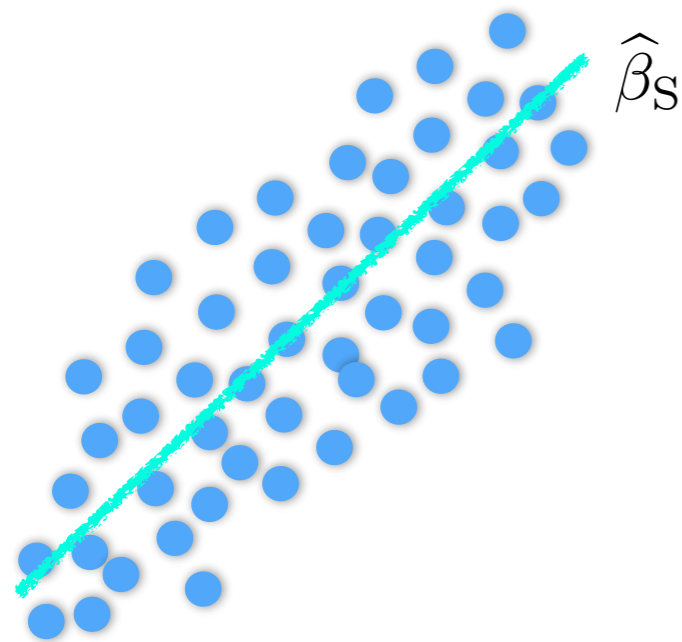
$$\min_{\beta \in \mathbb{R}^d} \frac{1}{N} \sum_{j=1}^N (\beta^\top \hat{x}_j - \hat{y}_j)^2 + \eta \|\beta\|_2^2$$

Source

Target

Aim

- Model works well on target domain



Abundant **labelled** data

Scarce **labelled** data

$$(\hat{x}_i, \hat{y}_i)_{i=1}^{N_S}$$

$$(\hat{x}_j, \hat{y}_j)_{j=1}^{N_T}$$

Supervised Domain Adaptation

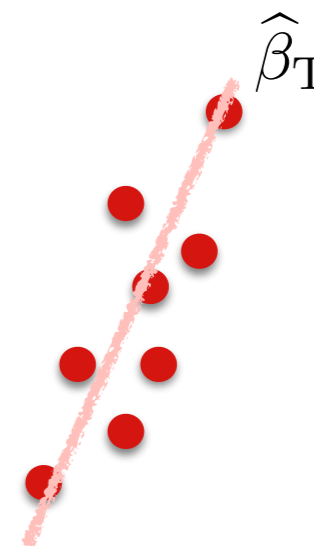
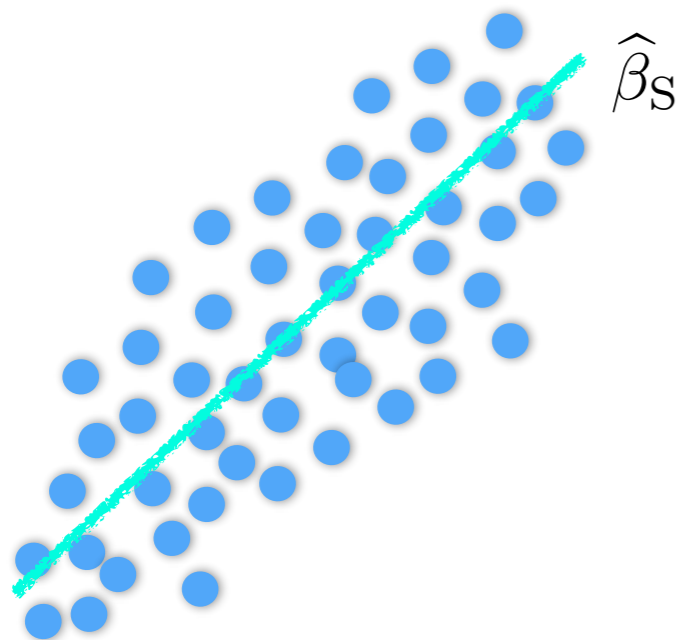
$$\min_{\beta \in \mathbb{R}^d} \frac{1}{N} \sum_{j=1}^N (\beta^\top \hat{x}_j - \hat{y}_j)^2 + \eta \|\beta\|_2^2$$

Source

Target

Assumptions

- Source and target are close
- Unseen target data arrive sequentially



Abundant **labelled** data

Scarce **labelled** data

$$(\hat{x}_i, \hat{y}_i)_{i=1}^{N_S}$$

$$(\hat{x}_j, \hat{y}_j)_{j=1}^{N_T}$$

Supervised Domain Adaptation

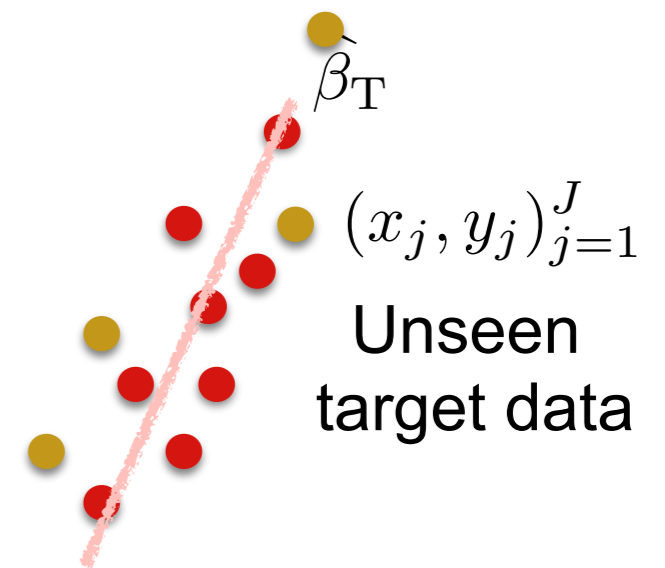
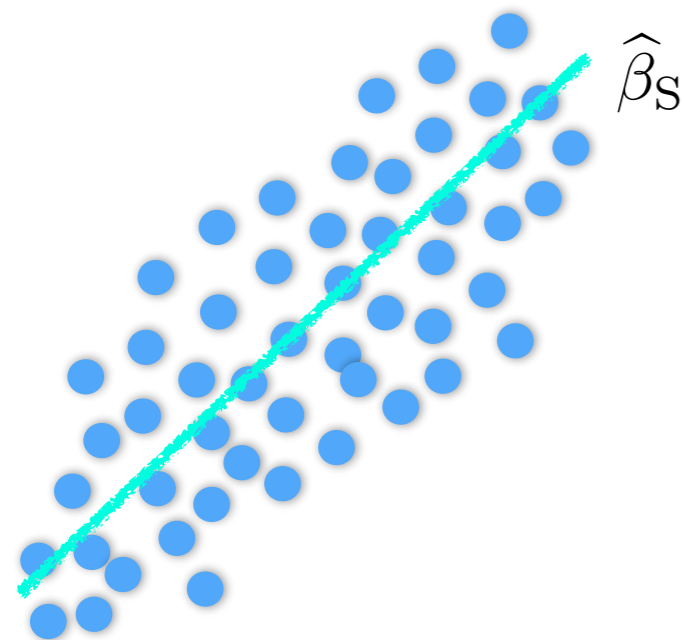
$$\min_{\beta \in \mathbb{R}^d} \frac{1}{N} \sum_{j=1}^N (\beta^\top \hat{x}_j - \hat{y}_j)^2 + \eta \|\beta\|_2^2$$

Source

Target

Assumptions

- Source and target are close
- Unseen target data arrive sequentially



Abundant **labelled** data

$$(\hat{x}_i, \hat{y}_i)_{i=1}^{N_S}$$

Scarce **labelled** data

$$(\hat{x}_j, \hat{y}_j)_{j=1}^{N_T}$$

Supervised Domain Adaptation

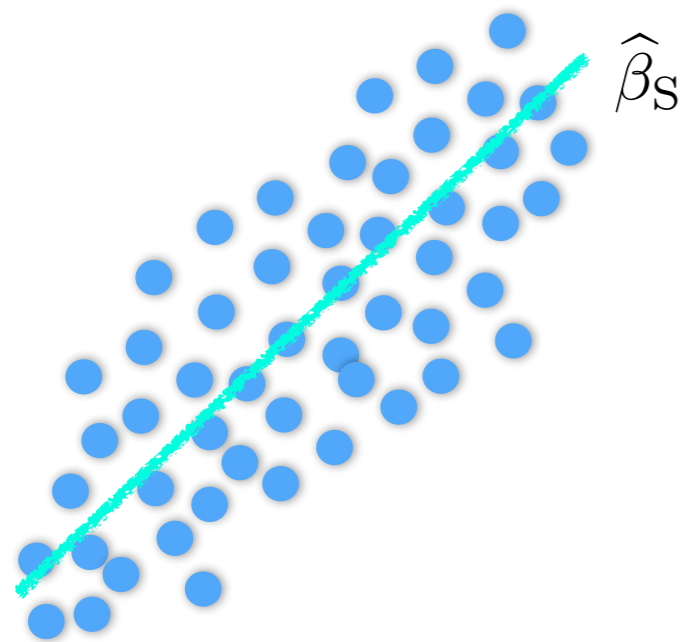
$$\min_{\beta \in \mathbb{R}^d} \frac{1}{N} \sum_{j=1}^N (\beta^\top \hat{x}_j - \hat{y}_j)^2 + \eta \|\beta\|_2^2$$

Source

Target

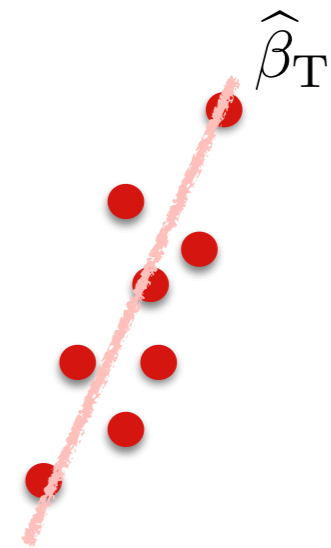
Challenges

1. Exploit source data
2. Tune hyperparameters



Abundant **labelled** data

$$(\hat{x}_i, \hat{y}_i)_{i=1}^{N_S}$$



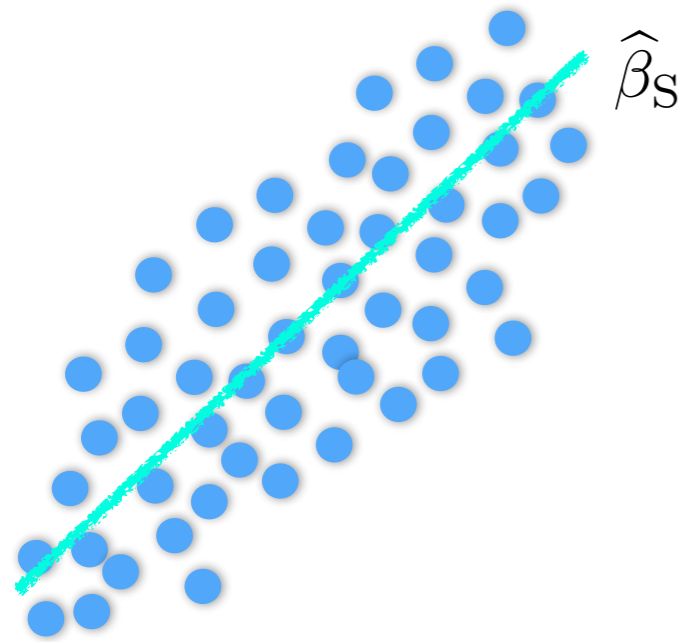
Scarce **labelled** data

$$(\hat{x}_j, \hat{y}_j)_{j=1}^{N_T}$$

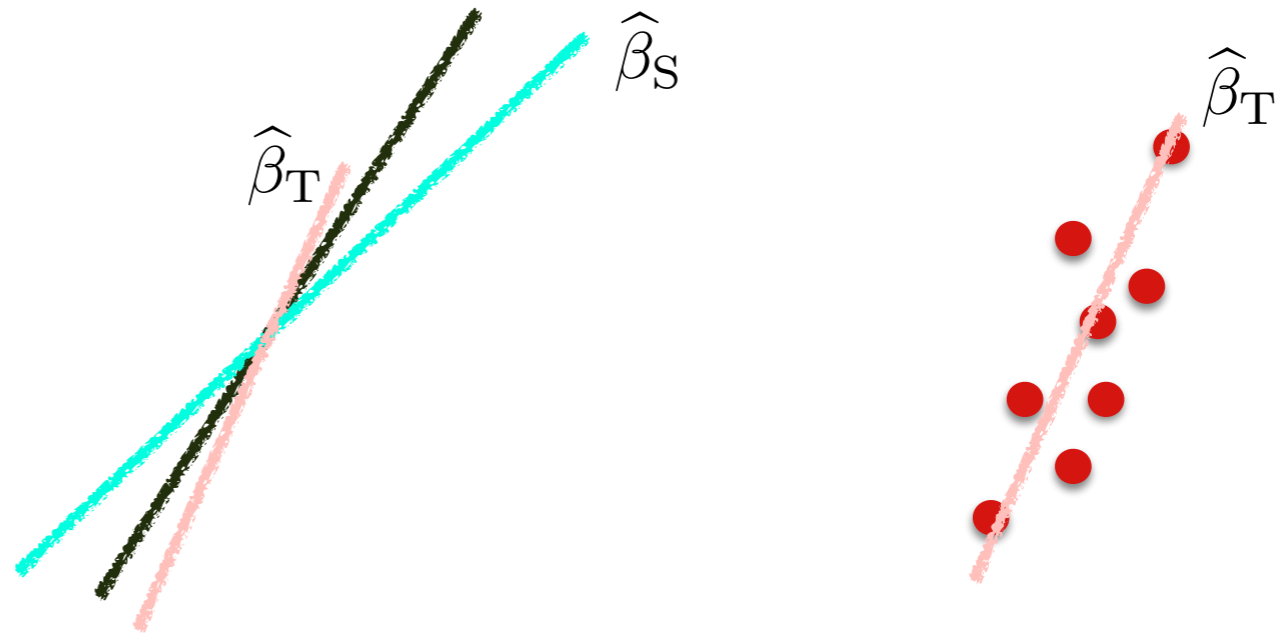
Exploit Source Data

1) Convex Combination Strategy (Baseline)

Source



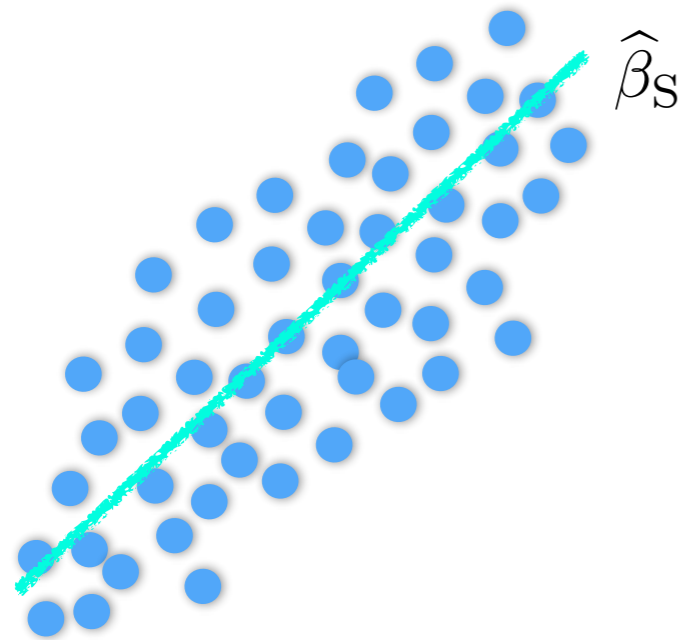
Target



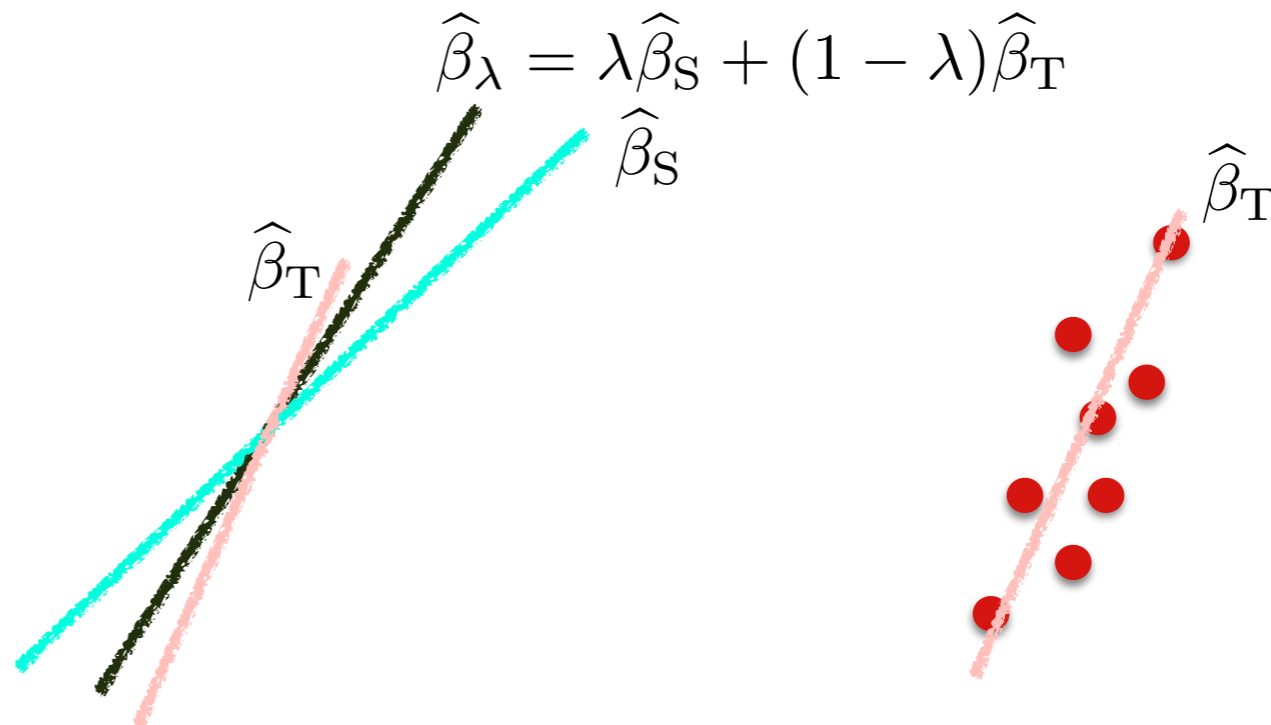
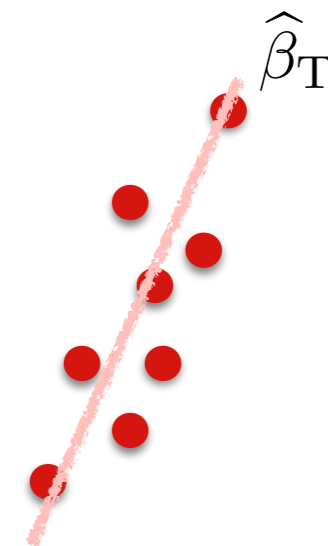
Exploit Source Data

1) Convex Combination Strategy (Baseline)

Source

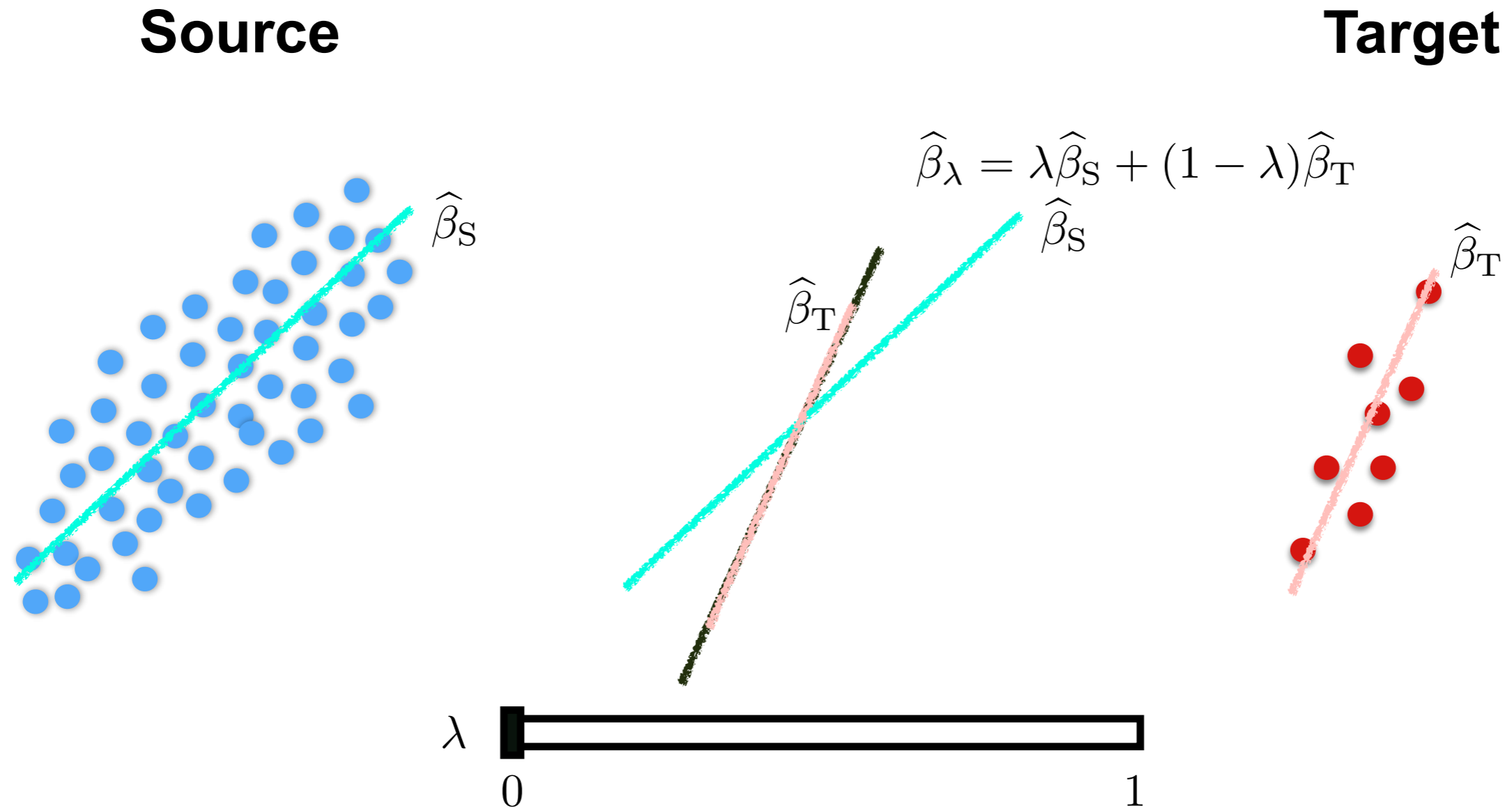


Target



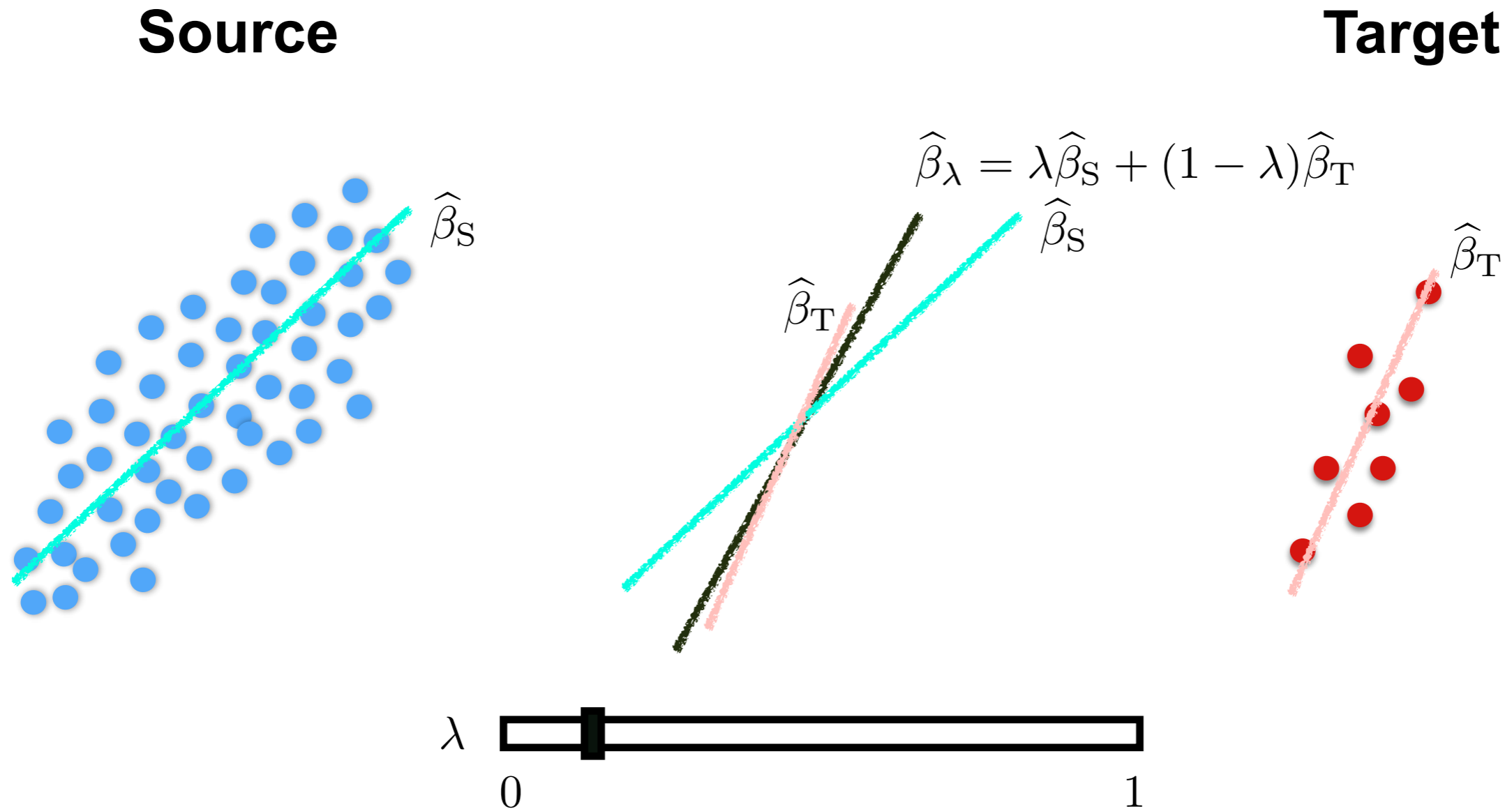
Exploit Source Data

1) Convex Combination Strategy (Baseline)



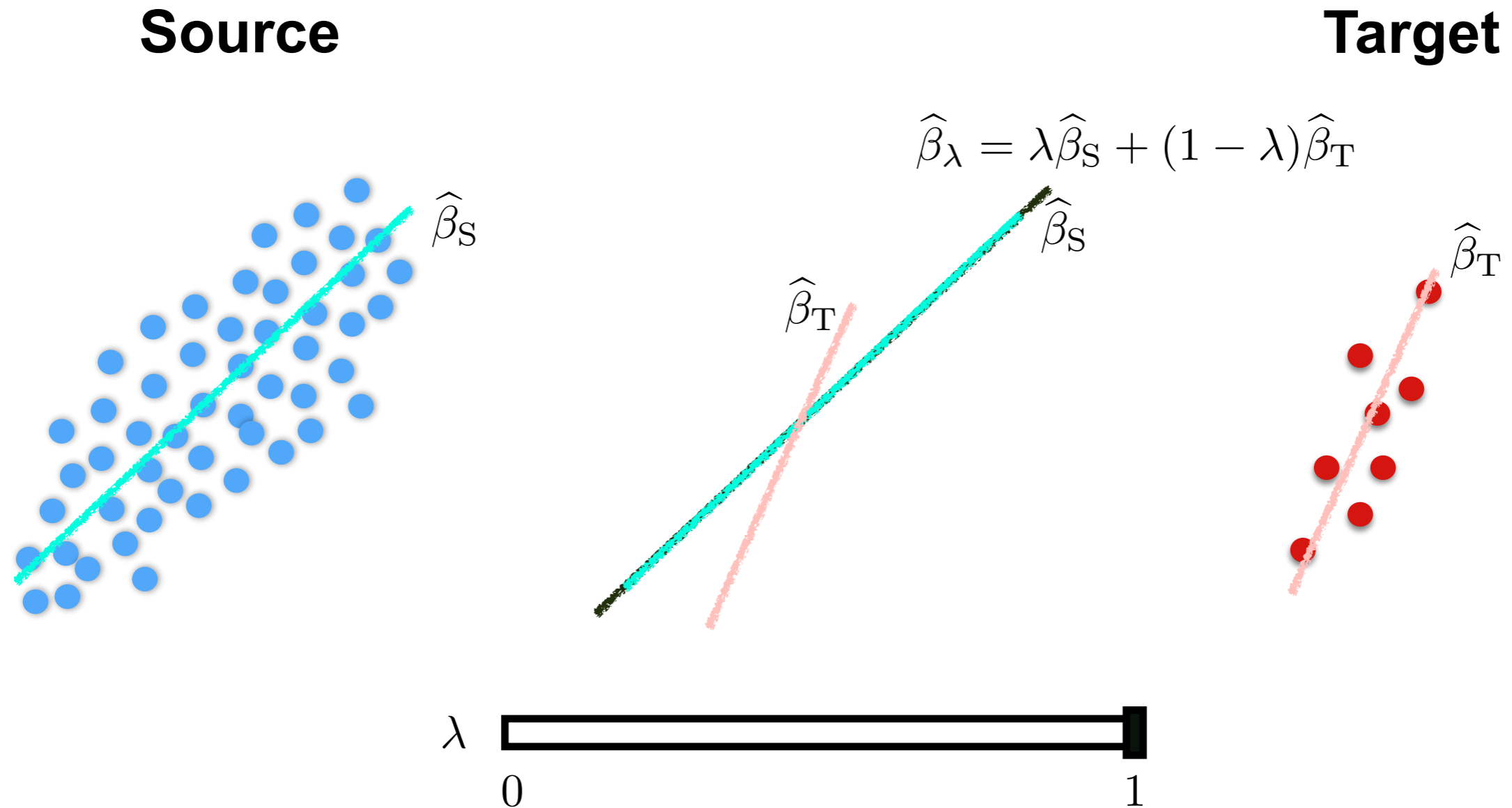
Exploit Source Data

1) Convex Combination Strategy (Baseline)



Exploit Source Data

1) Convex Combination Strategy (Baseline)



How to pick λ ?

Exploit Source Data

2) Reweighting Strategy (RWS)¹⁾

¹⁾ Garcke & Vanck, ECML PKDD, 2014

Exploit Source Data

2) Reweighting Strategy (RWS)¹⁾

$$\min_{\beta \in \mathbb{R}^d} \sum_{i=1}^{N_S} w_{h,i} (\beta^\top \hat{x}_i - \hat{y}_i)^2 + \sum_{j=1}^{N_T} (\beta^\top \hat{x}_j - \hat{y}_j)^2 + \eta \|\beta\|_2^2$$

¹⁾ Garcke & Vanck, ECML PKDD, 2014

Exploit Source Data

2) Reweighting Strategy (RWS)¹⁾

$$\min_{\beta \in \mathbb{R}^d} \sum_{i=1}^{N_S} w_{h,i} (\beta^\top \hat{x}_i - \hat{y}_i)^2 + \sum_{j=1}^{N_T} (\beta^\top \hat{x}_j - \hat{y}_j)^2 + \eta \|\beta\|_2^2$$

$$w_{h,i} = \sum_{l=1}^{N_S} \alpha_l \exp \left(-\frac{\|\hat{x}_i - \hat{x}_l\|_2^2 + (\hat{y}_i - \hat{y}_l)^2}{h^2} \right)$$

¹⁾ Garcke & Vanck, ECML PKDD, 2014

Exploit Source Data

2) Reweighting Strategy (RWS)¹⁾

$$\min_{\beta \in \mathbb{R}^d} \sum_{i=1}^{N_S} w_{h,i} (\beta^\top \hat{x}_i - \hat{y}_i)^2 + \sum_{j=1}^{N_T} (\beta^\top \hat{x}_j - \hat{y}_j)^2 + \eta \|\beta\|_2^2$$

$$w_{h,i} = \sum_{l=1}^{N_S} \alpha_l \exp \left(-\frac{\|\hat{x}_i - \hat{x}_l\|_2^2 + (\hat{y}_i - \hat{y}_l)^2}{h^2} \right)$$

How to pick h ?

¹⁾ Garcke & Vanck, ECML PKDD, 2014

Tune Hyperparameters - Synthesizing Experts

Expert - 1: β_1



Expert - 2: β_2



▪

▪

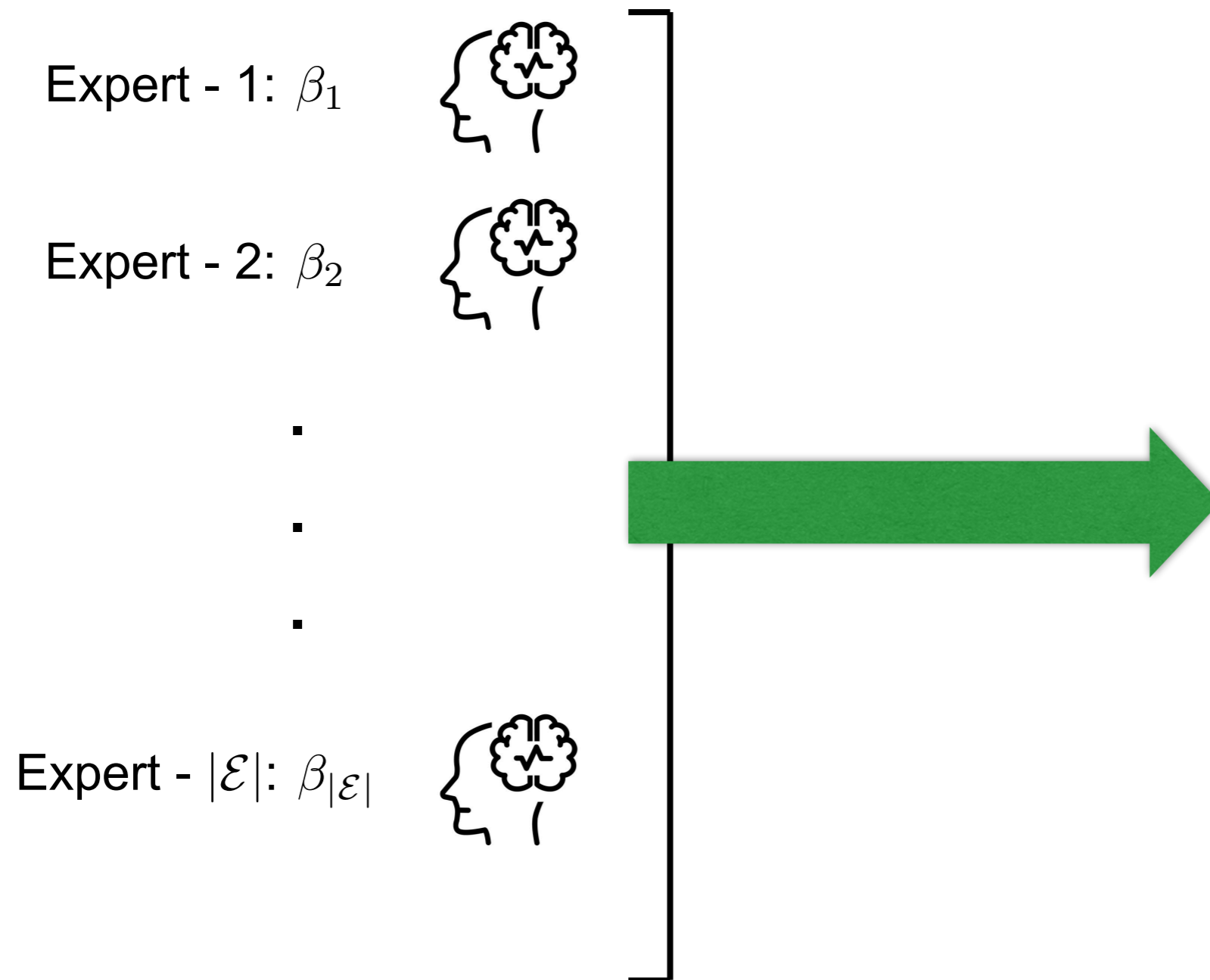
▪

Expert - $|\mathcal{E}|$: $\beta_{|\mathcal{E}|}$



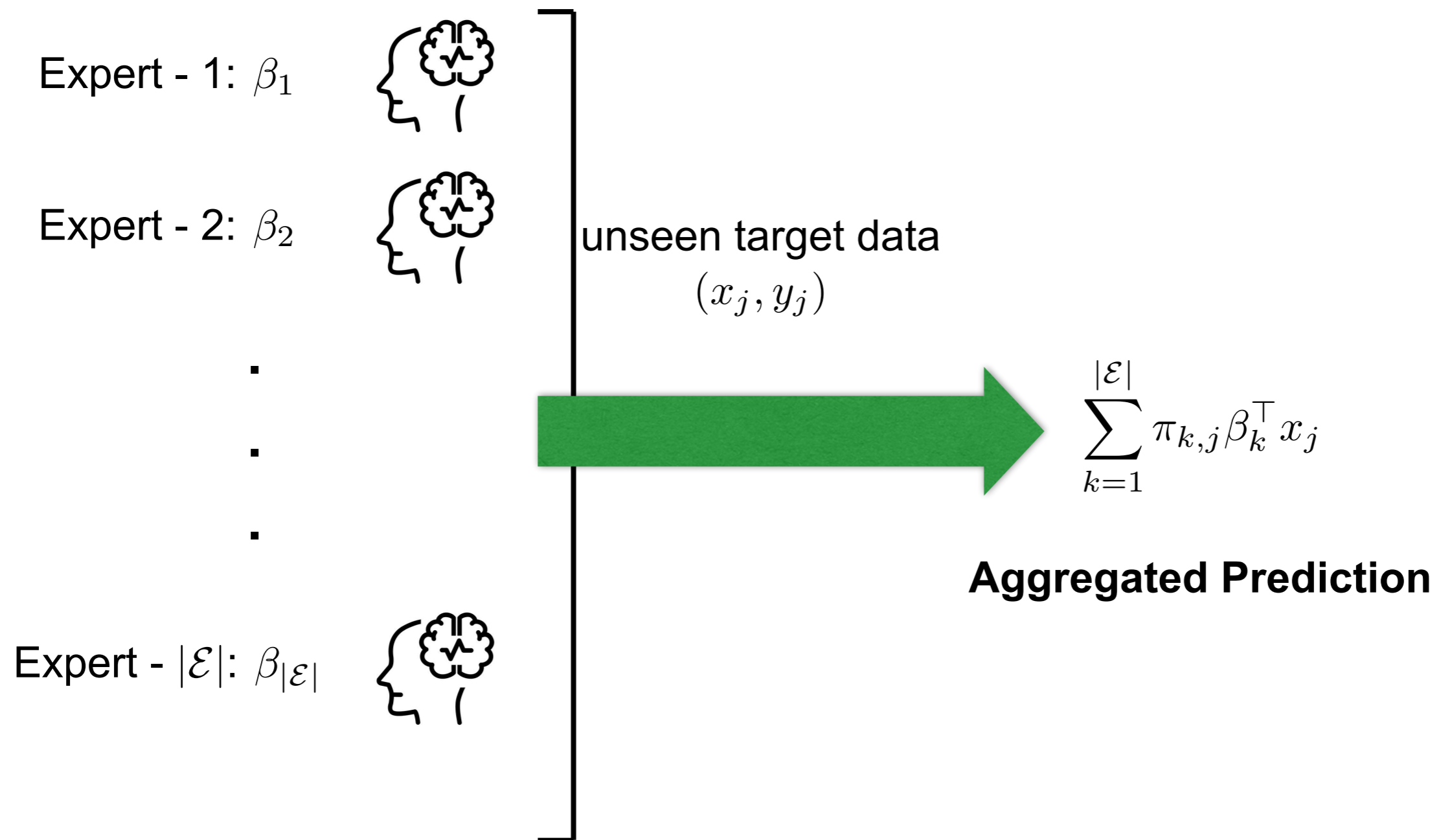
Tune Hyperparameters - Synthesizing Experts

Bernstein Online Aggregation (BOA)¹⁾



Tune Hyperparameters - Synthesizing Experts

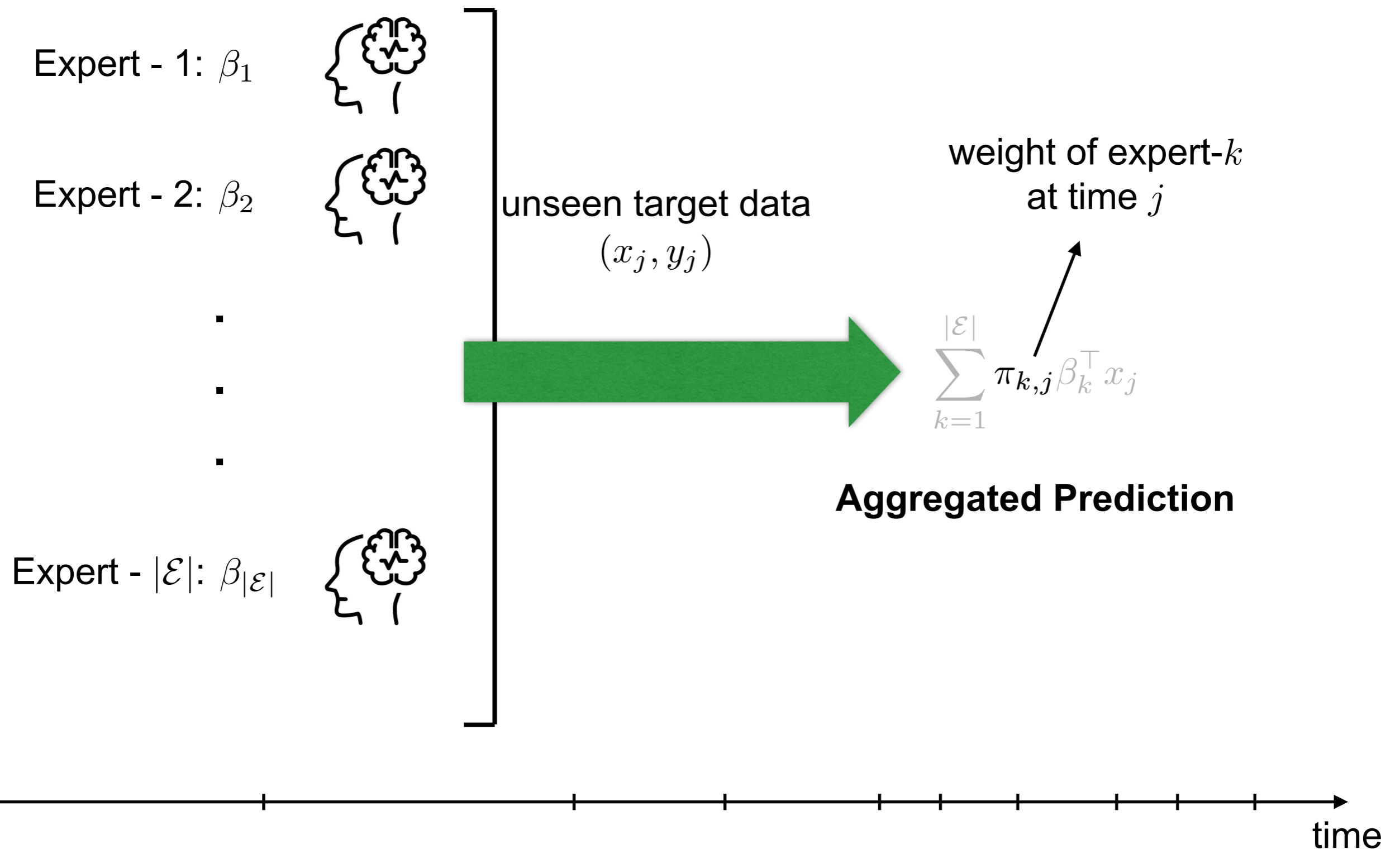
Bernstein Online Aggregation (BOA)¹⁾



¹⁾ Wintenberger, Machine Learning, 2017

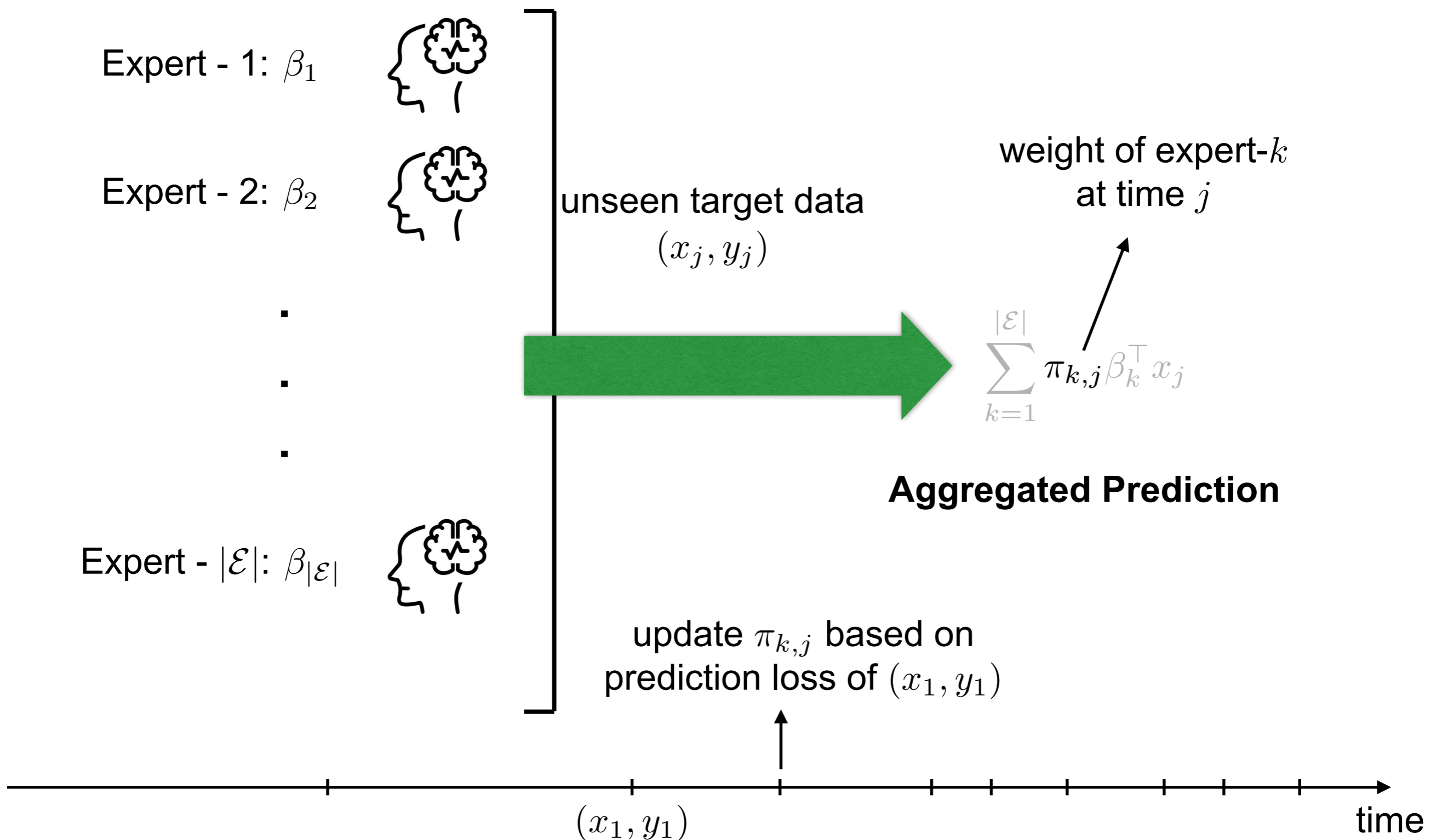
Synthesizing Experts

Bernstein Online Aggregation (BOA)



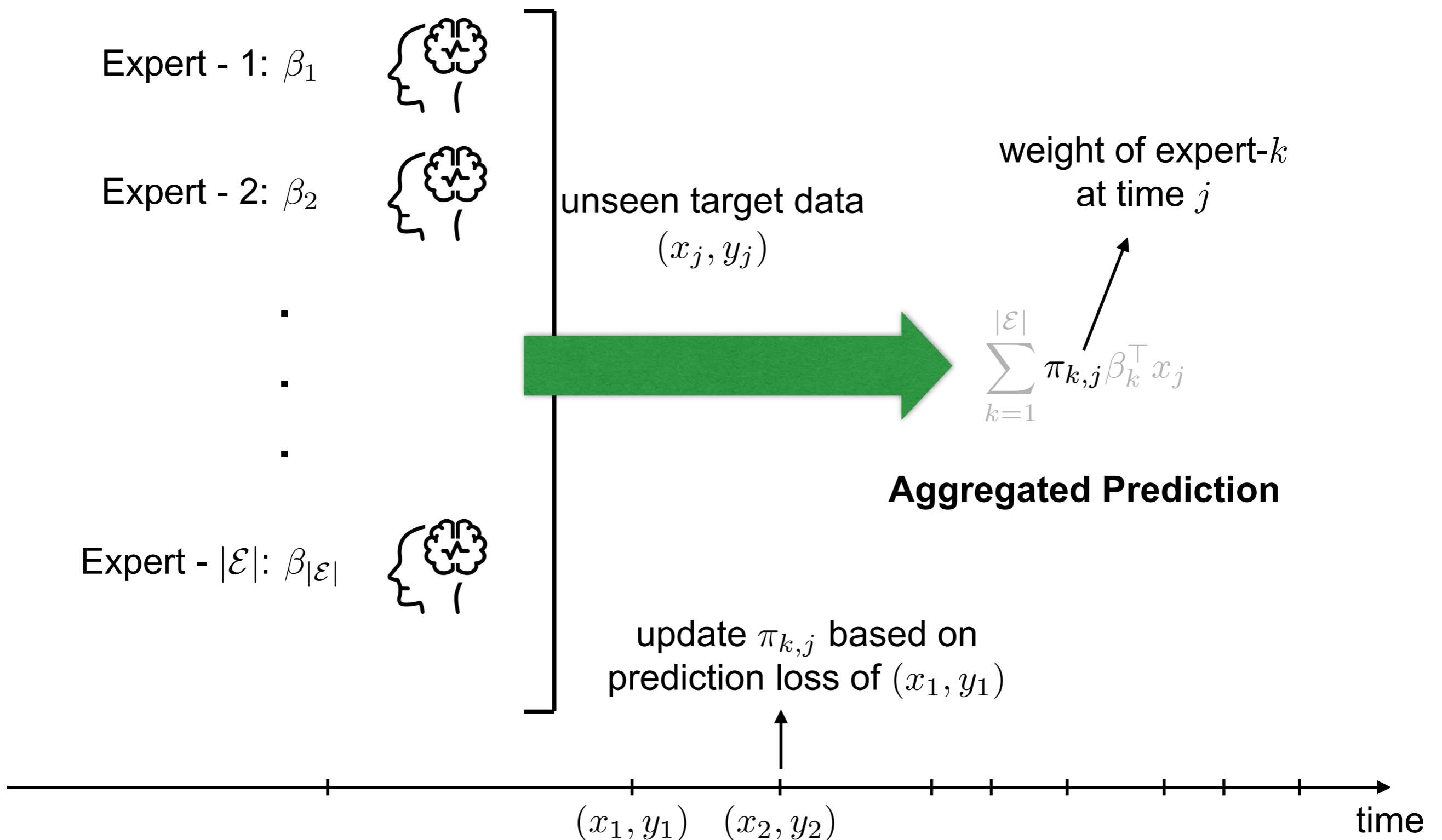
Synthesizing Experts

Bernstein Online Aggregation (BOA)



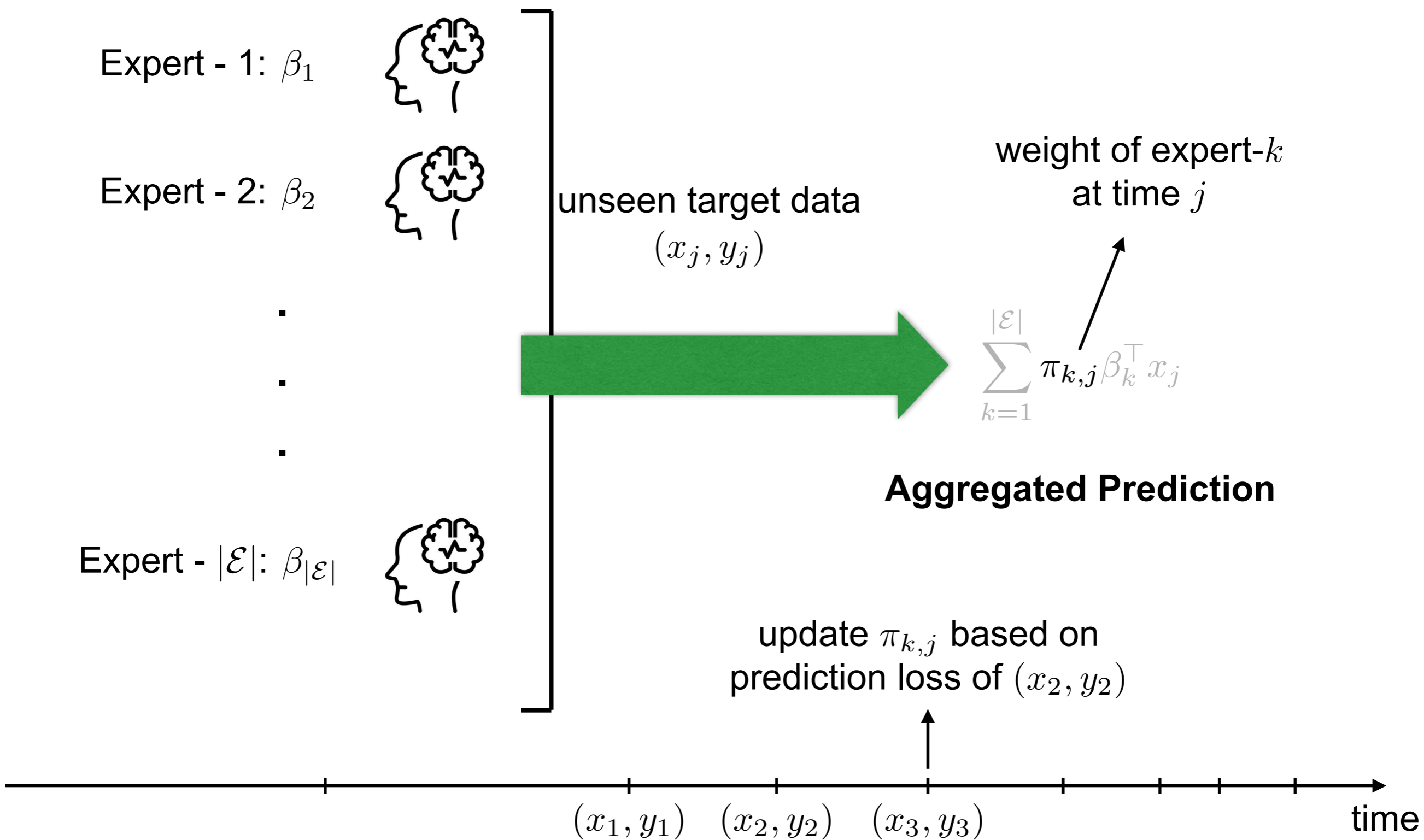
Synthesizing Experts

Bernstein Online Aggregation (BOA)



Synthesizing Experts

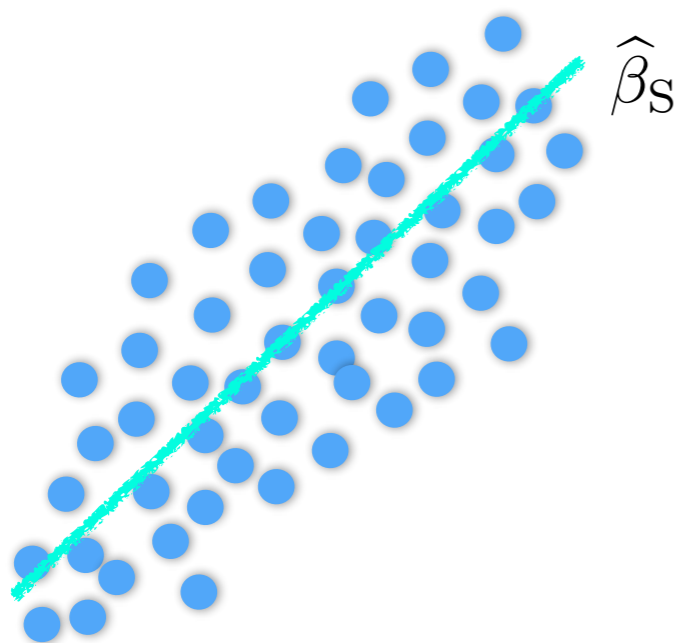
Bernstein Online Aggregation (BOA)



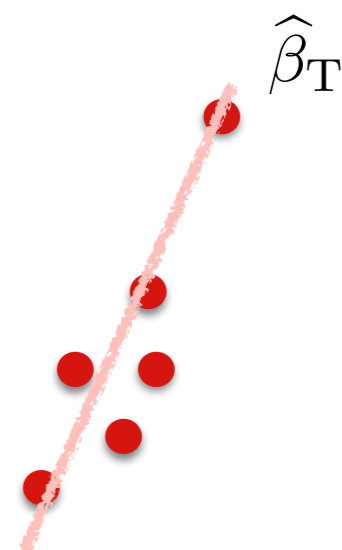
Distributionally Robust Expert Generation

Distributionally Robust Expert Generation

Source



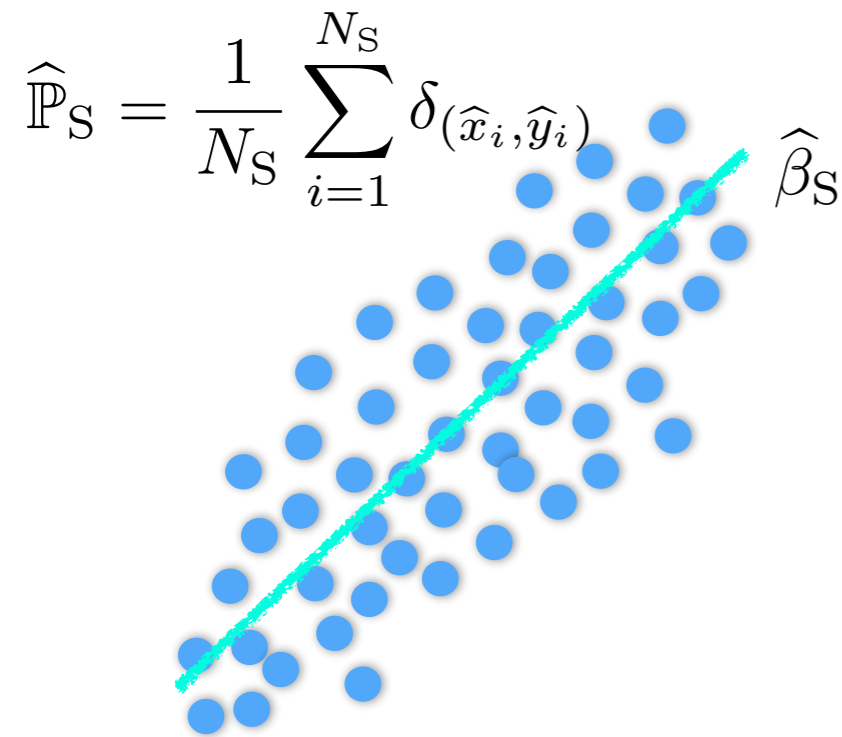
Target



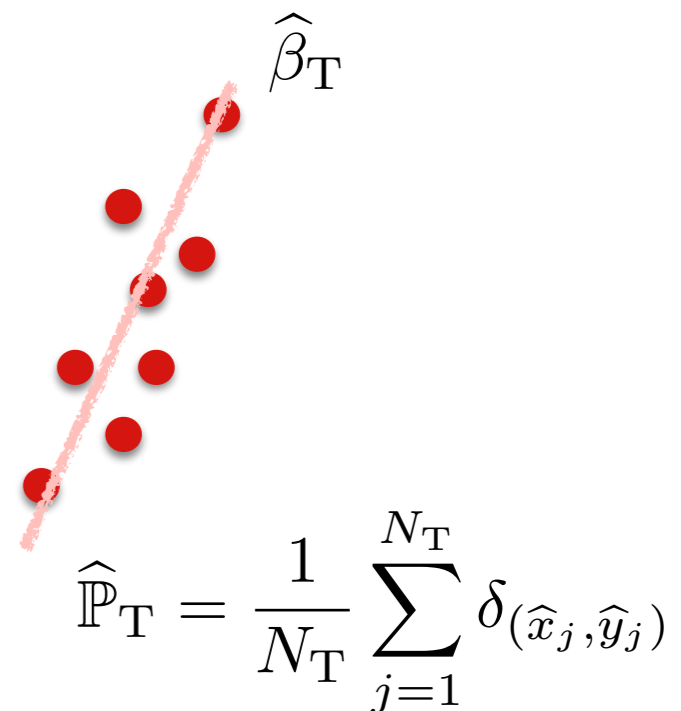
Distributionally Robust Expert Generation

$$\inf_{\beta \in \mathbb{R}^d} \sup_{Q \in \mathbb{B}} \mathbb{E}_Q[(\beta^\top X - Y)^2]$$

Source



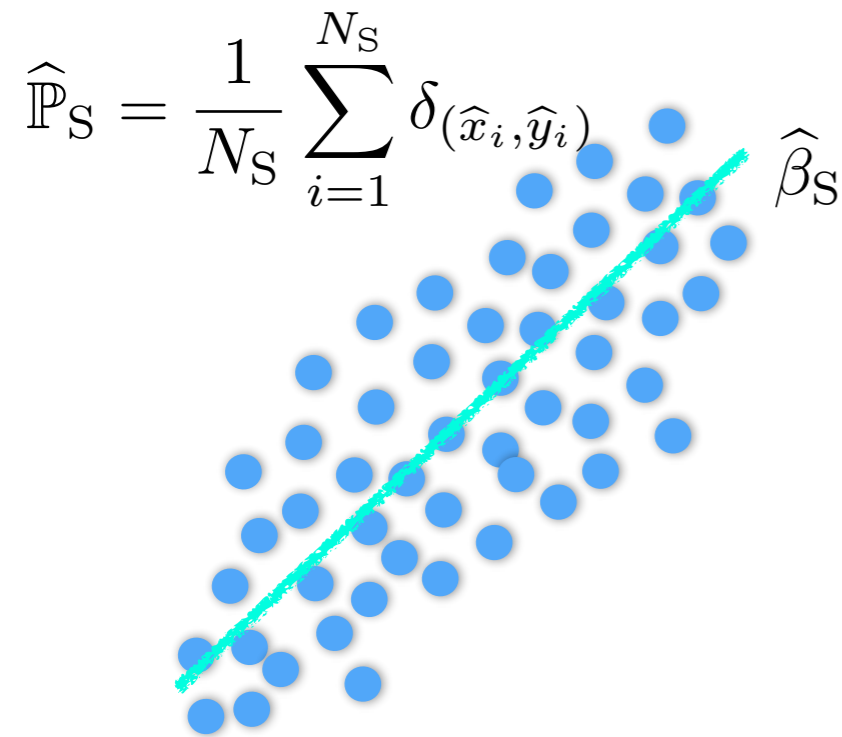
Target



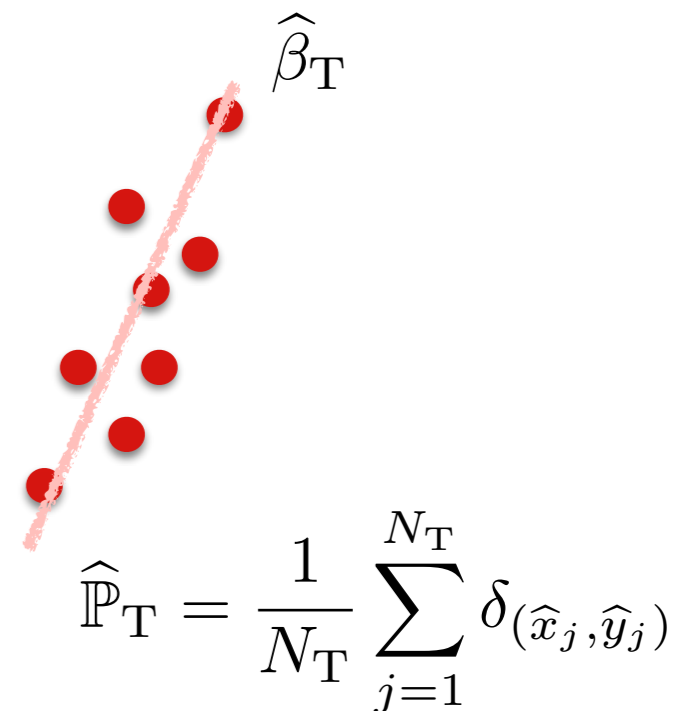
Distributionally Robust Expert Generation

$$\inf_{\beta \in \mathbb{R}^d} \sup_{Q \in \mathbb{B}} \mathbb{E}_Q[(\beta^\top X - Y)^2]$$

Source



Target



1) Distributional probing

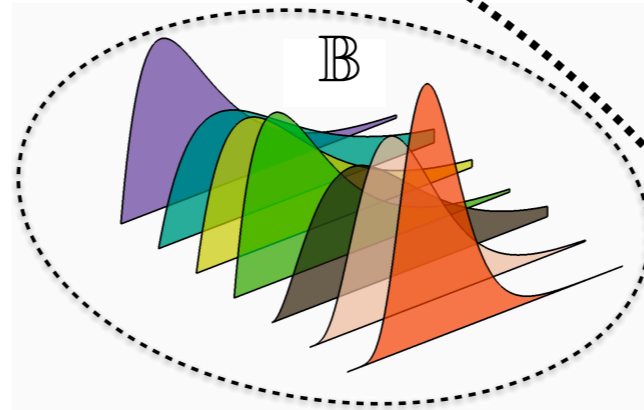
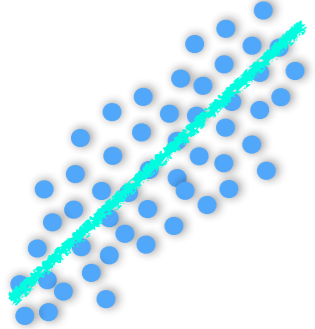
2) Robust Estimation

Distributionally Robust Expert Generation

$$\inf_{\beta \in \mathbb{R}^d} \sup_{Q \in \mathbb{B}} \mathbb{E}_Q [(\beta^\top X - Y)^2]$$

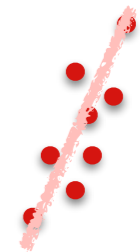
Source

$$\hat{\mathbb{P}}_S = \frac{1}{N_S} \sum_{i=1}^{N_S} \delta_{(\hat{x}_i, \hat{y}_i)}$$



Target

$$\hat{\mathbb{P}}_T = \frac{1}{N_T} \sum_{j=1}^{N_T} \delta_{(\hat{x}_j, \hat{y}_j)}$$



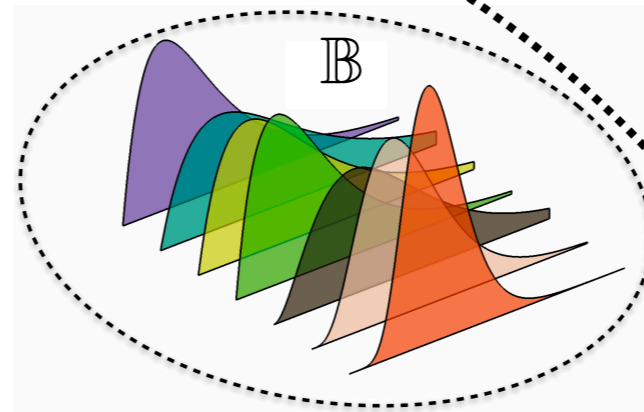
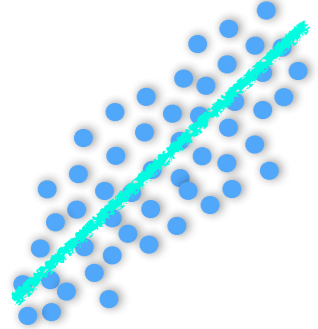
1) **Distributional probing**

Distributionally Robust Expert Generation

$$\inf_{\beta \in \mathbb{R}^d} \sup_{Q \in \mathbb{B}} \mathbb{E}_Q[(\beta^\top X - Y)^2]$$

Source

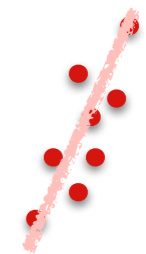
$$\hat{\mathbb{P}}_S = \frac{1}{N_S} \sum_{i=1}^{N_S} \delta_{(\hat{x}_i, \hat{y}_i)}$$



$$\mathbb{B} = \{Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}\}$$

Target

$$\hat{\mathbb{P}}_T = \frac{1}{N_T} \sum_{j=1}^{N_T} \delta_{(\hat{x}_j, \hat{y}_j)}$$



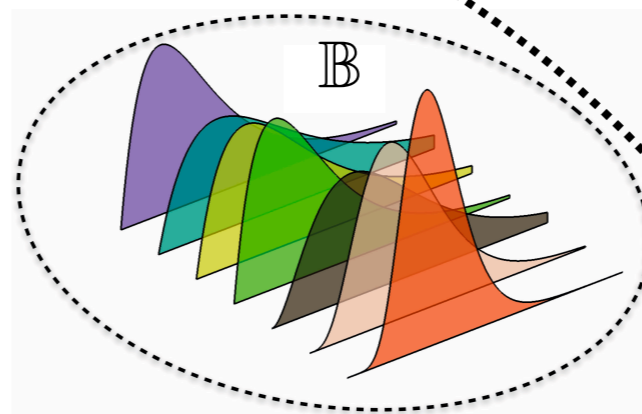
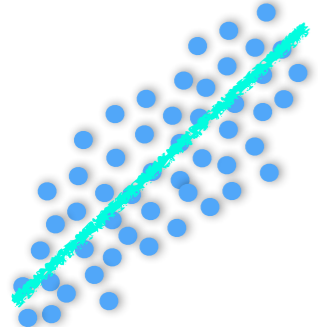
1) Distributional probing

Distributionally Robust Expert Generation

$$\inf_{\beta \in \mathbb{R}^d} \sup_{Q \in \mathbb{B}} \mathbb{E}_Q [(\beta^\top X - Y)^2]$$

Source

$$\hat{\mathbb{P}}_S = \frac{1}{N_S} \sum_{i=1}^{N_S} \delta_{(\hat{x}_i, \hat{y}_i)}$$



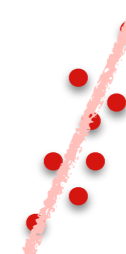
$$p = d + 1$$

$$\mathbb{B} = \{Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}\}$$

moment
information set

Target

$$\hat{\mathbb{P}}_T = \frac{1}{N_T} \sum_{j=1}^{N_T} \delta_{(\hat{x}_j, \hat{y}_j)}$$

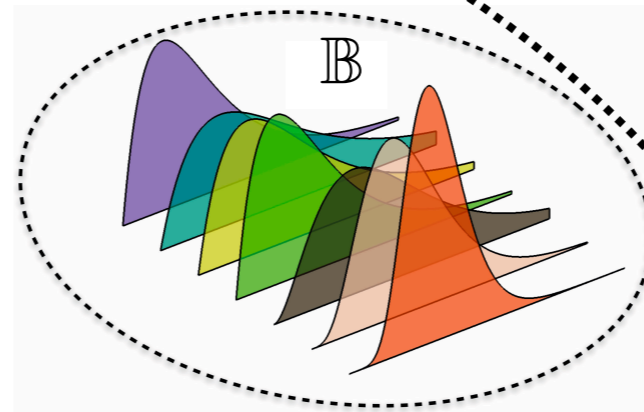
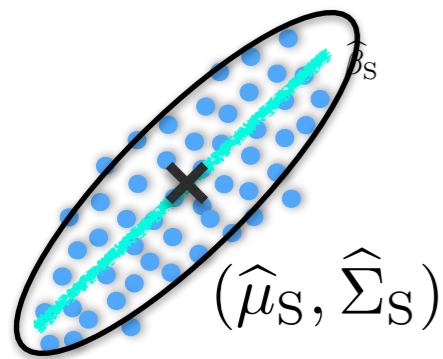


1) Distributional probing

Distributionally Robust Expert Generation

$$\inf_{\beta \in \mathbb{R}^d} \sup_{Q \in \mathbb{B}} \mathbb{E}_Q[(\beta^\top X - Y)^2]$$

Source

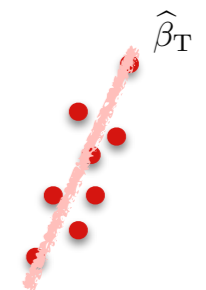


$$p = d + 1$$

$$\mathbb{B} = \{Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}\}$$

moment
information set

Target

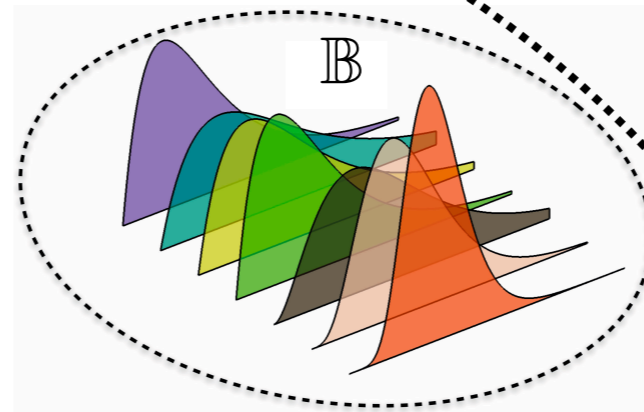
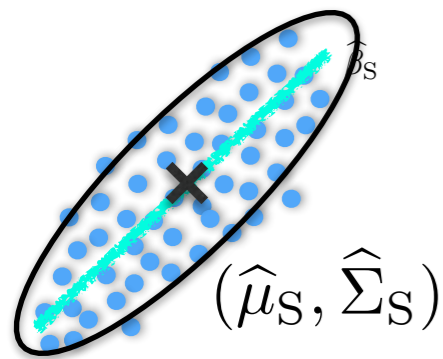


1) Distributional probing

Distributionally Robust Expert Generation

$$\inf_{\beta \in \mathbb{R}^d} \sup_{Q \in \mathbb{B}} \mathbb{E}_Q[(\beta^\top X - Y)^2]$$

Source

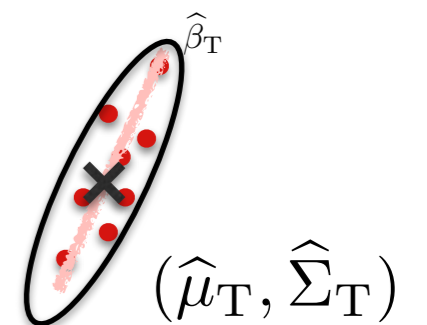


$$p = d + 1$$

$$\mathbb{B} = \{Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}\}$$

moment
information set

Target

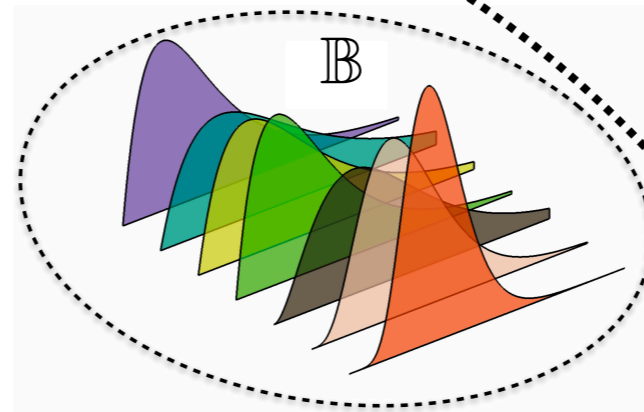
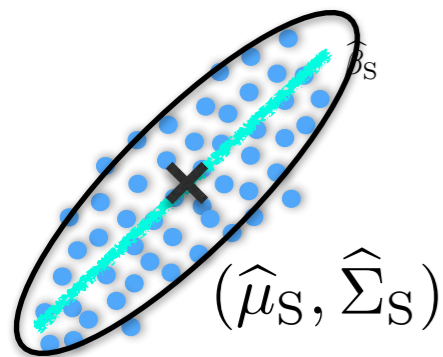


1) Distributional probing

Distributionally Robust Expert Generation

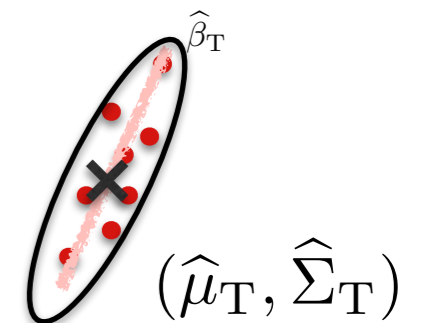
$$\inf_{\beta \in \mathbb{R}^d} \sup_{Q \in \mathbb{B}} \mathbb{E}_Q[(\beta^\top X - Y)^2]$$

Source



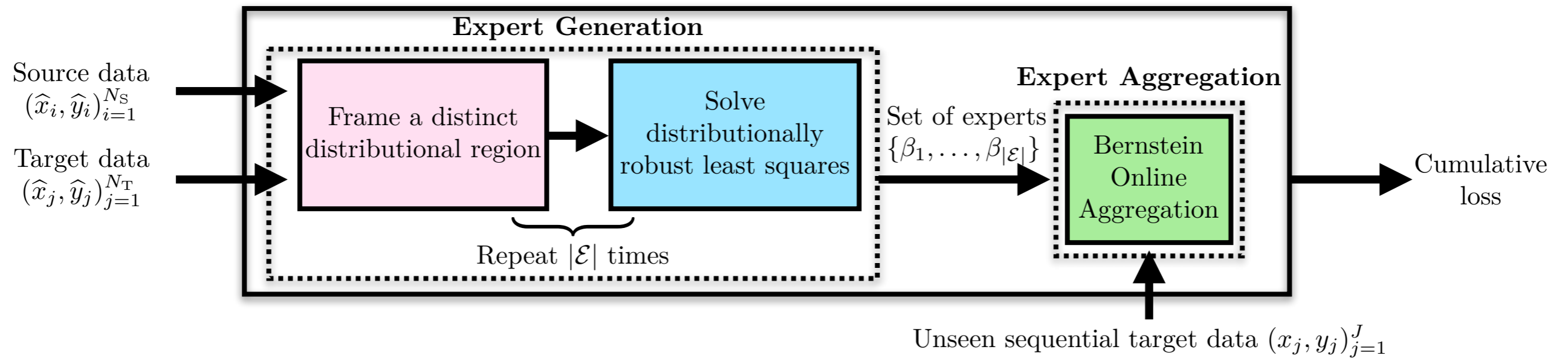
Target

$$\mathbb{B} = \{Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}\}$$

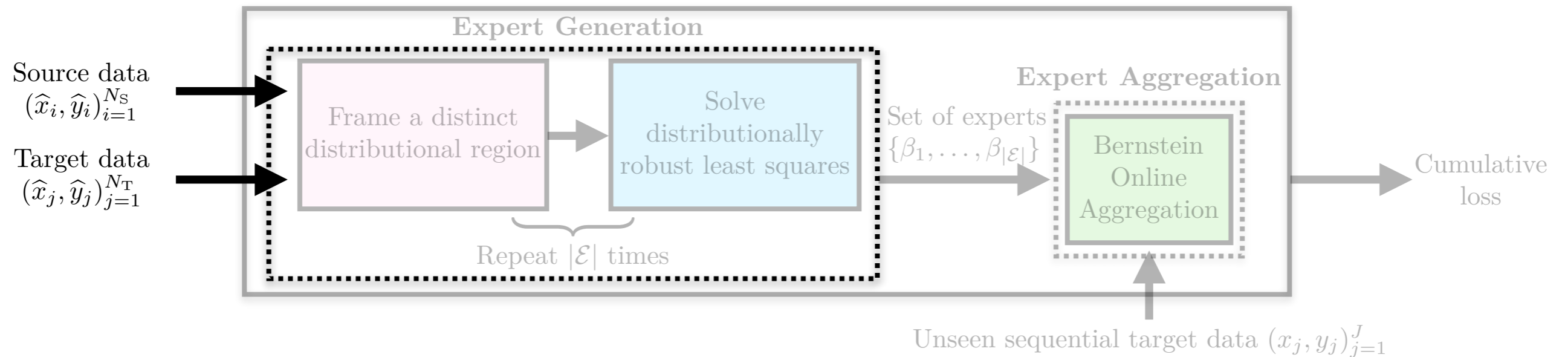


- 1) **Distributional probing**
- 2) **Robust Estimation**

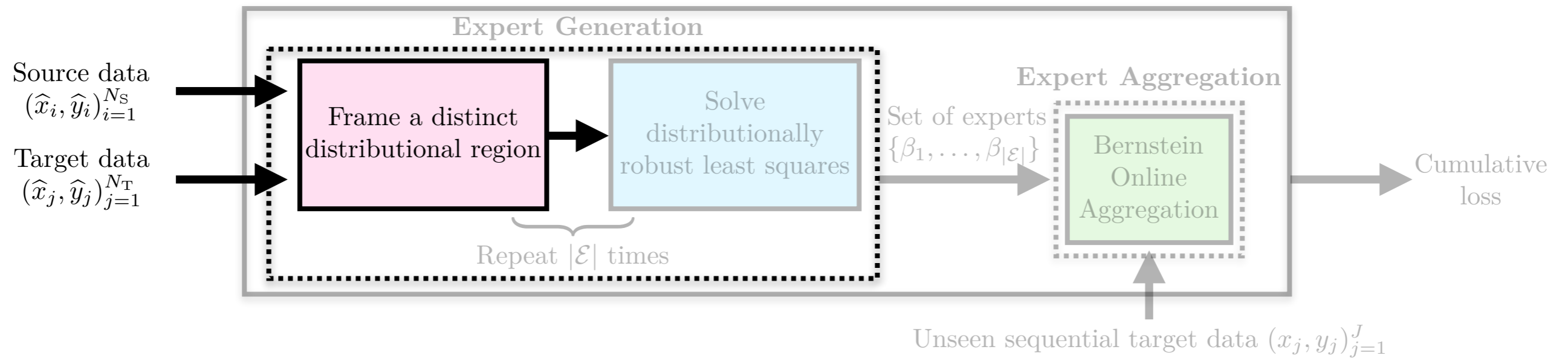
Summary of Framework



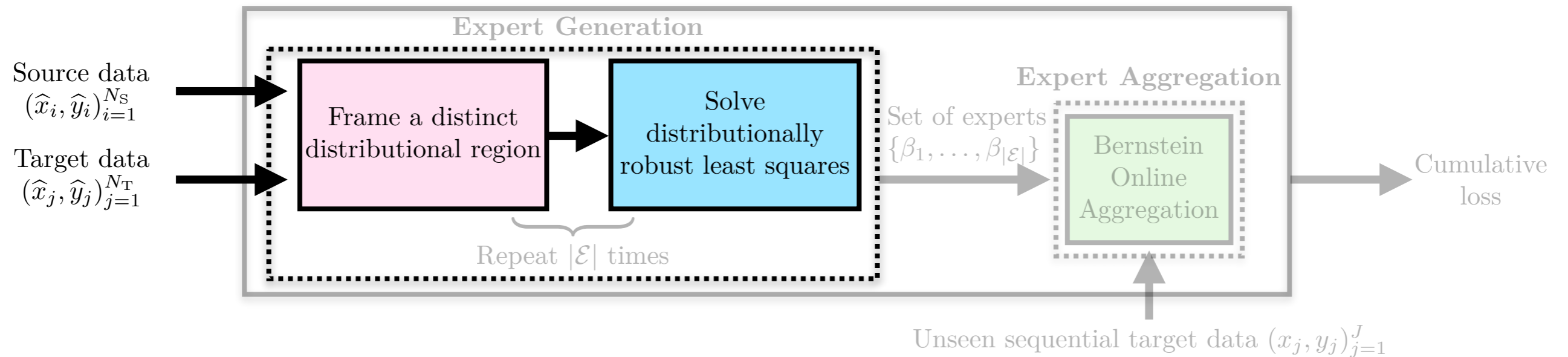
Summary of Framework



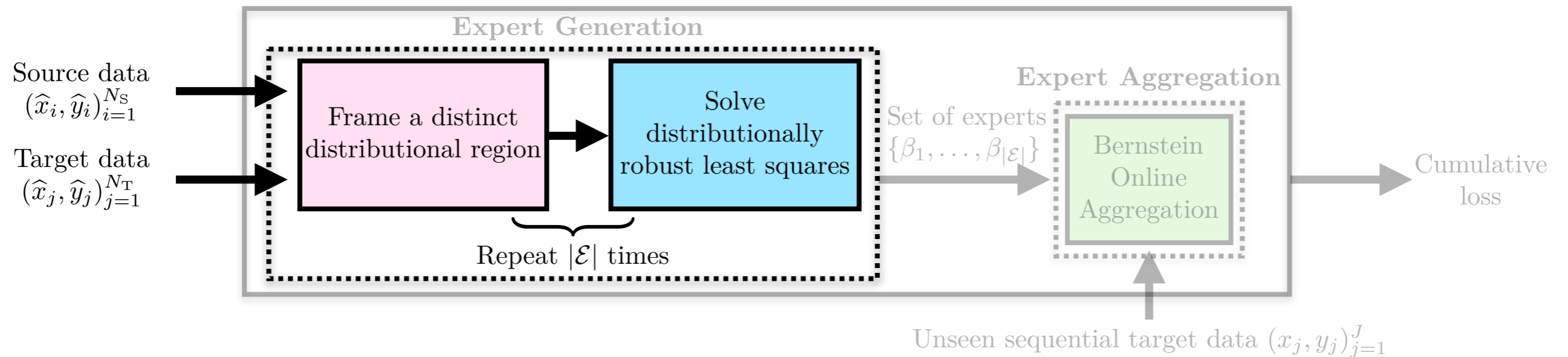
Summary of Framework



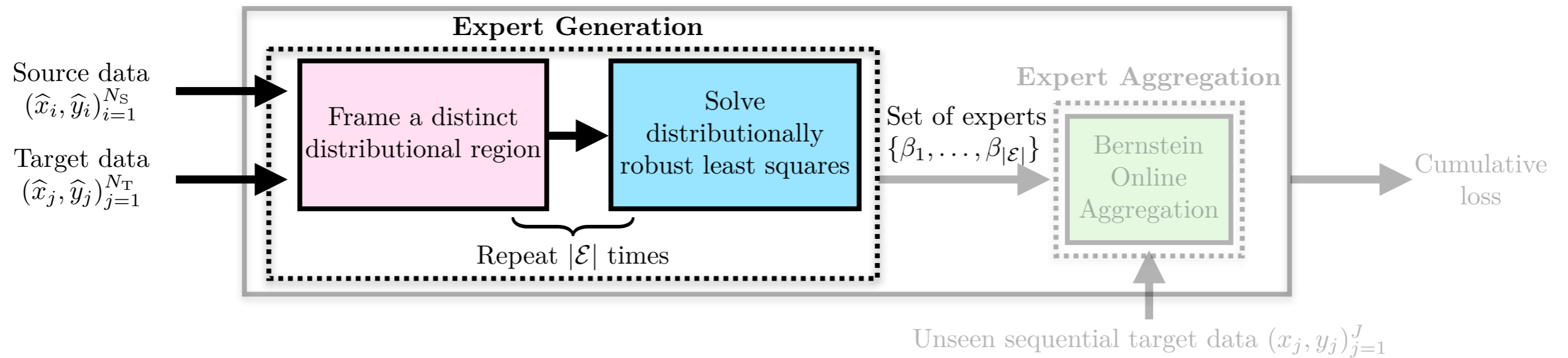
Summary of Framework



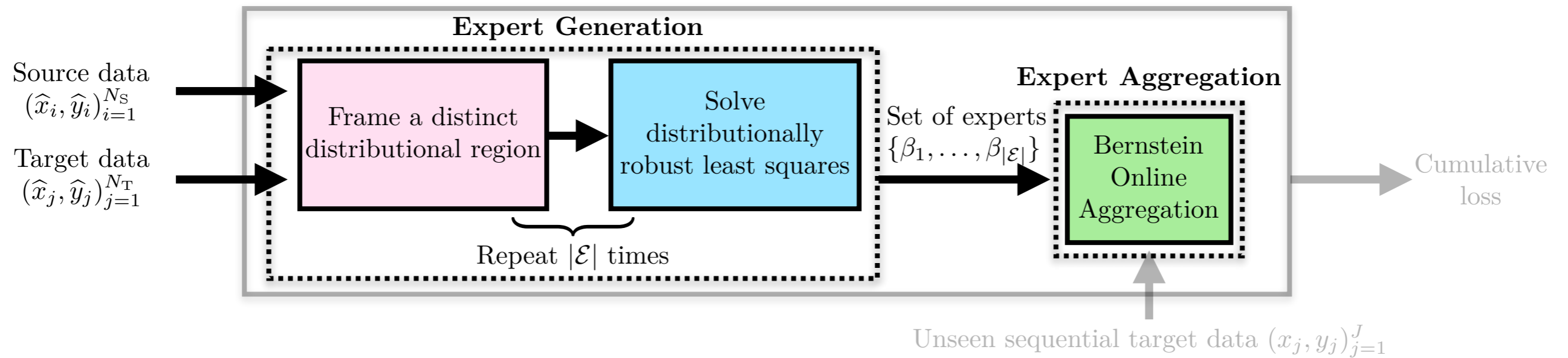
Summary of Framework



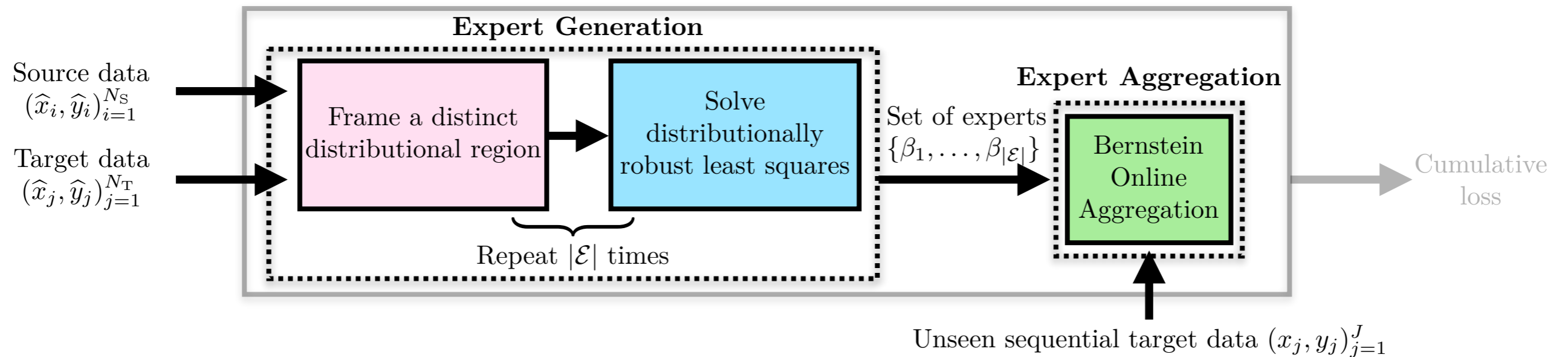
Summary of Framework



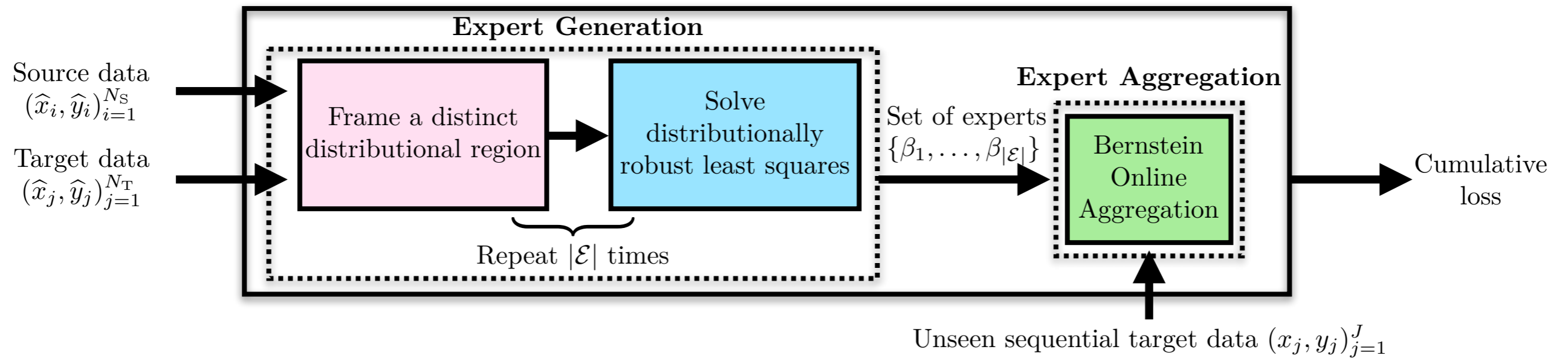
Summary of Framework



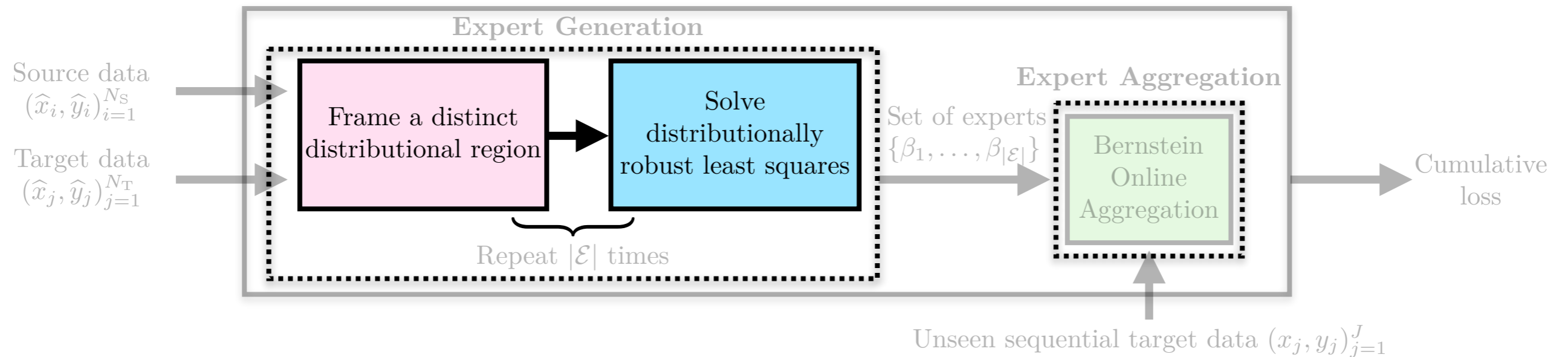
Summary of Framework



Summary of Framework



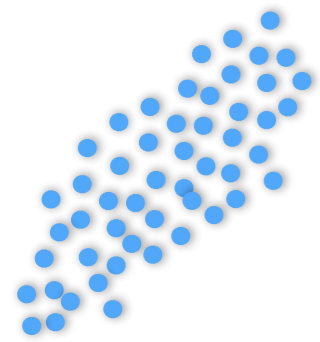
Summary of Framework



- 1) “Interpolate, then Robustify” (IR)
- 2) “Surround, then Intersect” (SI)

1) “Interpolate, then Robustify” (IR) Strategy

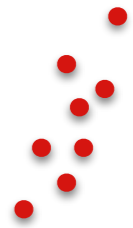
Source



$$(\hat{x}_i, \hat{y}_i)_{i=1}^{N_S}$$

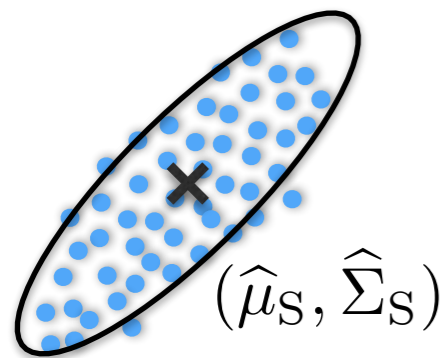
Target

$$(\hat{x}_j, \hat{y}_j)_{j=1}^{N_T}$$



1) “Interpolate, then Robustify” (IR) Strategy

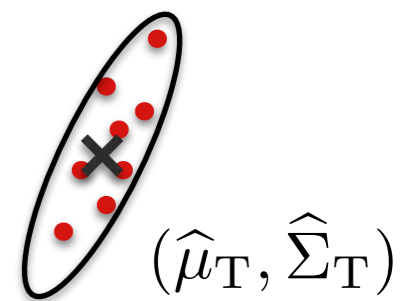
Source



$$\hat{\mu}_S = \frac{1}{N_S} \sum_{i=1}^{N_S} \begin{pmatrix} \hat{x}_i \\ \hat{y}_i \end{pmatrix}$$

$$\hat{\Sigma}_S = \frac{1}{N_S} \sum_{i=1}^{N_S} \begin{pmatrix} \hat{x}_i \\ \hat{y}_i \end{pmatrix} \begin{pmatrix} \hat{x}_i \\ \hat{y}_i \end{pmatrix}^\top - \hat{\mu}_S \hat{\mu}_S^\top$$

Target



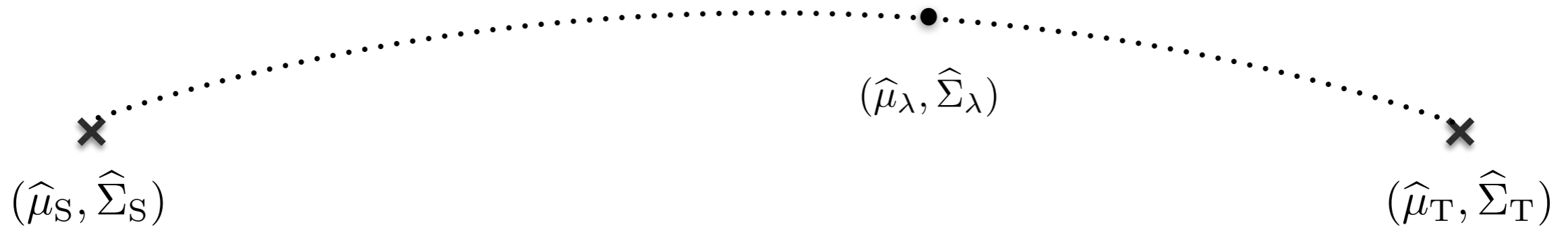
$$\hat{\mu}_T = \frac{1}{N_T} \sum_{j=1}^{N_T} \begin{pmatrix} \hat{x}_j \\ \hat{y}_j \end{pmatrix}$$

$$\hat{\Sigma}_T = \frac{1}{N_T} \sum_{j=1}^{N_T} \begin{pmatrix} \hat{x}_j \\ \hat{y}_j \end{pmatrix} \begin{pmatrix} \hat{x}_j \\ \hat{y}_j \end{pmatrix}^\top - \hat{\mu}_T \hat{\mu}_T^\top$$

1) “Interpolate, then Robustify” (IR) Strategy

Source

Target



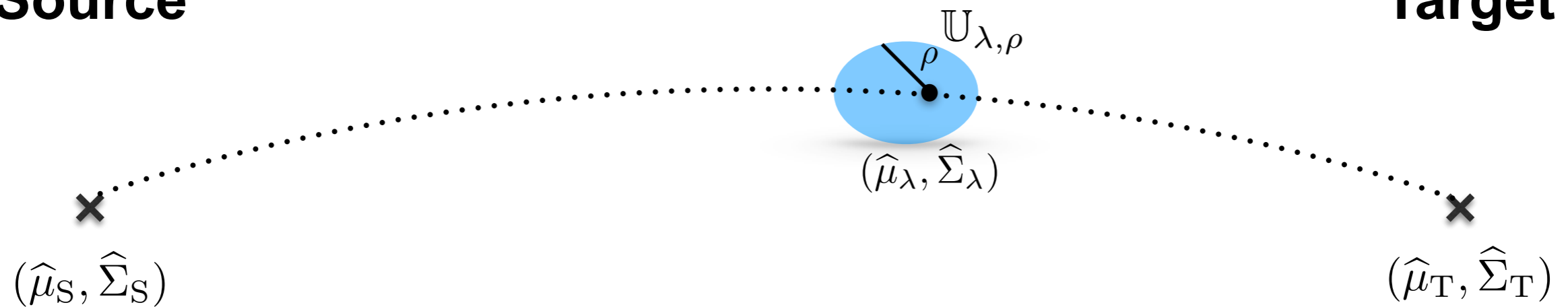
ψ -barycenter

$$(\hat{\mu}_\lambda, \hat{\Sigma}_\lambda) = \operatorname{argmin}_{\mu \in \mathbb{R}^p, \Sigma \in \mathbb{S}_+^p} \lambda \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) + (1 - \lambda) \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T))$$

1) “Interpolate, then Robustify” (IR) Strategy

Source

Target



ψ -barycenter

$$(\hat{\mu}_\lambda, \hat{\Sigma}_\lambda) = \operatorname{argmin}_{\mu \in \mathbb{R}^p, \Sigma \in \mathbb{S}_+^p} \lambda \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) + (1 - \lambda) \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T))$$

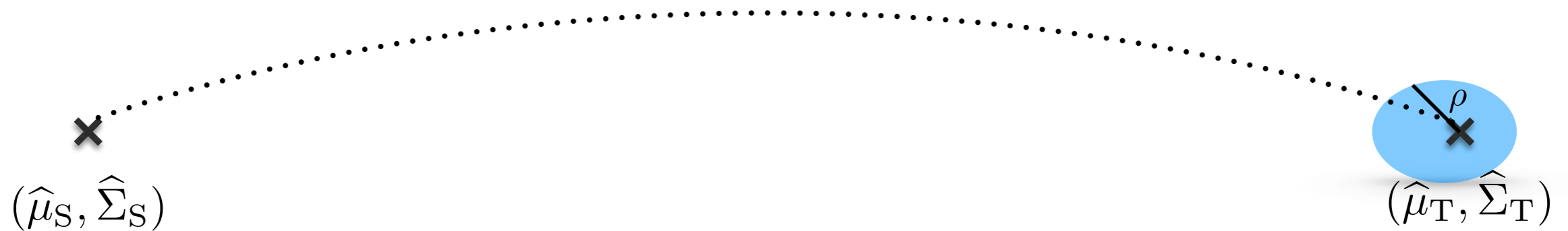
Moment information set

$$\mathbb{U}_{\lambda, \rho} = \{(\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p : \psi((\mu, \Sigma) \| (\hat{\mu}_\lambda, \hat{\Sigma}_\lambda)) \leq \rho\}$$

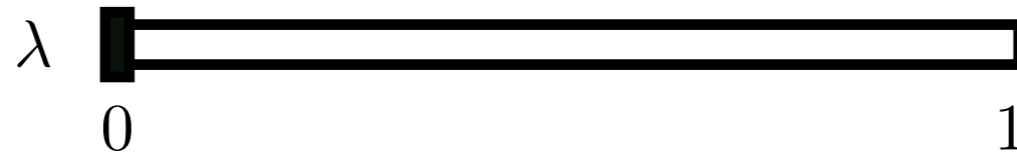
1) “Interpolate, then Robustify” (IR) Strategy

Source

Target



ψ -barycenter



$$(\hat{\mu}_\lambda, \hat{\Sigma}_\lambda) = \operatorname{argmin}_{\mu \in \mathbb{R}^p, \Sigma \in \mathbb{S}_+^p} \lambda \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) + (1 - \lambda) \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T))$$

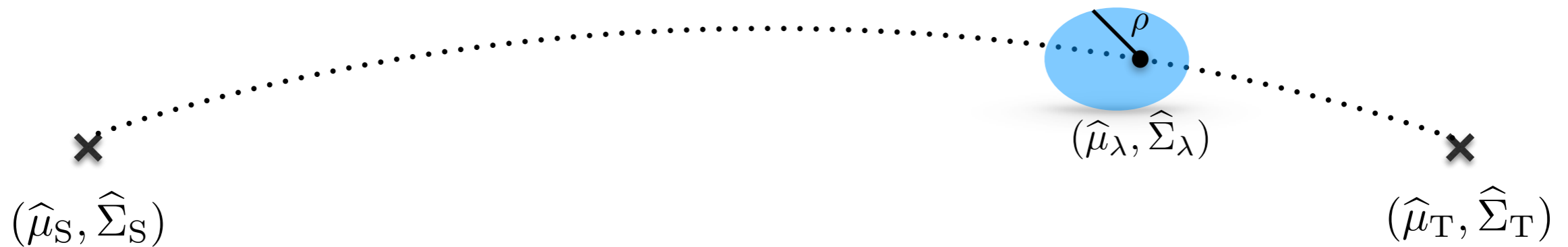
Moment information set

$$\mathbb{U}_{\lambda, \rho} = \{(\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p : \psi((\mu, \Sigma) \| (\hat{\mu}_\lambda, \hat{\Sigma}_\lambda)) \leq \rho\}$$

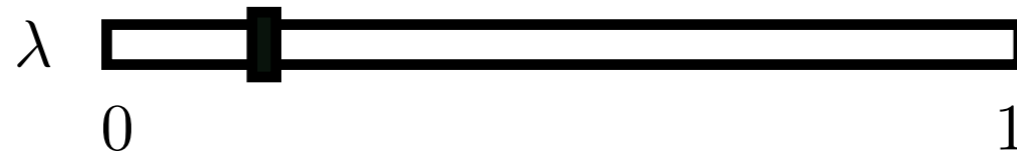
1) “Interpolate, then Robustify” (IR) Strategy

Source

Target



ψ -barycenter



$$(\hat{\mu}_\lambda, \hat{\Sigma}_\lambda) = \operatorname{argmin}_{\mu \in \mathbb{R}^p, \Sigma \in \mathbb{S}_+^p} \lambda \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) + (1 - \lambda) \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T))$$

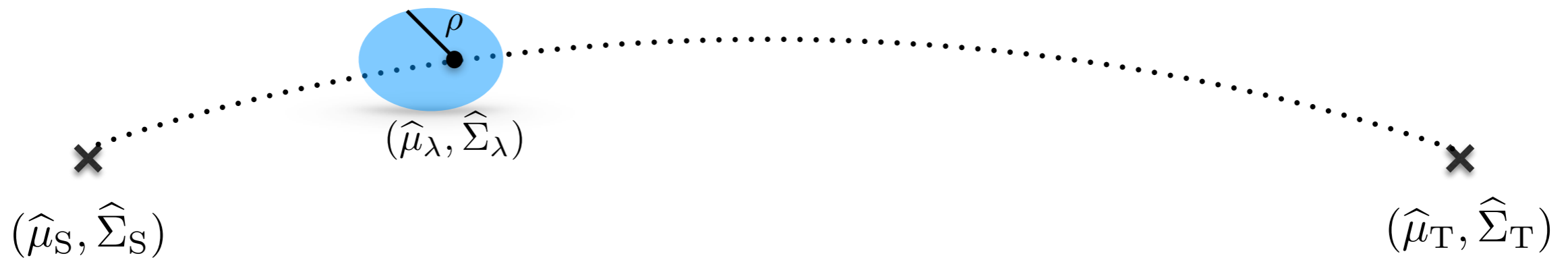
Moment information set

$$\mathbb{U}_{\lambda, \rho} = \{(\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p : \psi((\mu, \Sigma) \| (\hat{\mu}_\lambda, \hat{\Sigma}_\lambda)) \leq \rho\}$$

1) “Interpolate, then Robustify” (IR) Strategy

Source

Target



ψ -barycenter



$$(\hat{\mu}_\lambda, \hat{\Sigma}_\lambda) = \operatorname{argmin}_{\mu \in \mathbb{R}^p, \Sigma \in \mathbb{S}_+^p} \lambda \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) + (1 - \lambda) \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T))$$

Moment information set

$$\mathbb{U}_{\lambda, \rho} = \{(\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p : \psi((\mu, \Sigma) \| (\hat{\mu}_\lambda, \hat{\Sigma}_\lambda)) \leq \rho\}$$

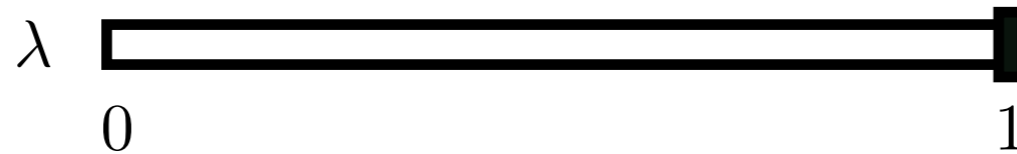
1) “Interpolate, then Robustify” (IR) Strategy

Source

Target



ψ -barycenter



$$(\hat{\mu}_\lambda, \hat{\Sigma}_\lambda) = \operatorname{argmin}_{\mu \in \mathbb{R}^p, \Sigma \in \mathbb{S}_+^p} \lambda \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) + (1 - \lambda) \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T))$$

Moment information set

$$\mathbb{U}_{\lambda, \rho} = \{(\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p : \psi((\mu, \Sigma) \| (\hat{\mu}_\lambda, \hat{\Sigma}_\lambda)) \leq \rho\}$$

1) “Interpolate, then Robustify” (IR) Strategy

i) Kullback-Leibler (KL)-type Divergence

$$\mathbb{D}((\mu, \Sigma) \| (\hat{\mu}, \hat{\Sigma})) = (\hat{\mu} - \mu)^\top \hat{\Sigma}^{-1} (\hat{\mu} - \mu) + \text{Tr}[\Sigma \hat{\Sigma}^{-1}] - \log \det(\Sigma \hat{\Sigma}^{-1}) - p$$

KL barycenter:

$$\hat{\Sigma}_\lambda = (\lambda \hat{\Sigma}_S^{-1} + (1 - \lambda) \hat{\Sigma}_T^{-1})^{-1} \quad \hat{\mu}_\lambda = \hat{\Sigma}_\lambda (\lambda \hat{\Sigma}_S^{-1} \hat{\mu}_S + (1 - \lambda) \hat{\Sigma}_T^{-1} \hat{\mu}_T)$$

Moment information set:

$$\mathbb{U}_{\lambda, \rho} = \{(\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p : \mathbb{D}((\mu, \Sigma) \| (\hat{\mu}_\lambda, \hat{\Sigma}_\lambda)) \leq \rho\}$$

Distribution set:

$$\mathbb{B}_{\lambda, \rho} = \{Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}_{\lambda, \rho}\}$$

Robust Estimation:

$$\min_{\beta \in \mathbb{R}^d} \left\{ \sup_{Q \in \mathbb{B}_{\lambda, \rho}} \mathbb{E}_Q [(\beta^\top X - Y)^2] \right\}$$

convex, continuously differentiable, and locally smooth



Global optimality via adaptive gradient descent algorithm

1) “Interpolate, then Robustify” (IR) Strategy

ii) Wasserstein-type Divergence

$$\mathbb{W}((\mu, \Sigma) \| (\hat{\mu}, \hat{\Sigma})) = \|\mu - \hat{\mu}\|_2^2 + \text{Tr}[\Sigma + \hat{\Sigma} - 2(\hat{\Sigma}^{\frac{1}{2}} \Sigma \hat{\Sigma}^{\frac{1}{2}})^{\frac{1}{2}}]$$

Wasserstein interpolation:

$$\hat{\mu}_\lambda = \lambda \hat{\mu}_S + (1 - \lambda) \hat{\mu}_T$$

$$\hat{\Sigma}_\lambda = (\lambda I_p + (1 - \lambda)L) \hat{\Sigma}_S (\lambda I_p + (1 - \lambda)L)$$

$$L = \hat{\Sigma}_T^{\frac{1}{2}} (\hat{\Sigma}_T^{\frac{1}{2}} \hat{\Sigma}_S \hat{\Sigma}_T^{\frac{1}{2}})^{-\frac{1}{2}} \hat{\Sigma}_T^{\frac{1}{2}}$$

Moment information set:

$$\mathbb{U}_{\lambda, \rho} = \{(\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p : \mathbb{W}((\mu, \Sigma) \| (\hat{\mu}_\lambda, \hat{\Sigma}_\lambda)) \leq \rho\}$$

Distribution set:

$$\mathbb{B}_{\lambda, \rho} = \{Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}_{\lambda, \rho}\}$$

Robust Estimation:

$$\min_{\beta \in \mathbb{R}^d} \sup_{Q \in \mathbb{B}_{\lambda, \rho}} \mathbb{E}_Q[(\beta^\top X - Y)^2] \equiv \min_{\beta \in \mathbb{R}^d} \left\| (\hat{\Sigma}_\lambda + \hat{\mu}_\lambda \hat{\mu}_\lambda^\top)^{\frac{1}{2}} \begin{bmatrix} \beta \\ -1 \end{bmatrix} \right\|_2 + \sqrt{\rho} \left\| \begin{bmatrix} \beta \\ -1 \end{bmatrix} \right\|_2$$



2) “Surround, then Intersect” (SI) Strategy

Source

Target



$(\hat{\mu}_S, \hat{\Sigma}_S)$



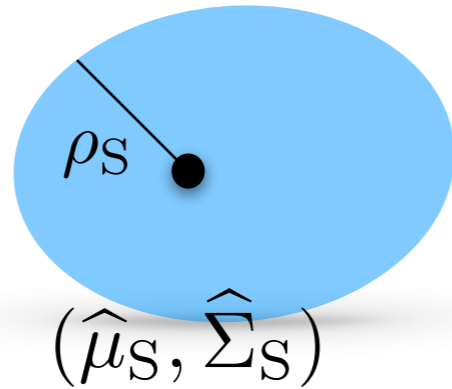
$(\hat{\mu}_T, \hat{\Sigma}_T)$

Moment information set:

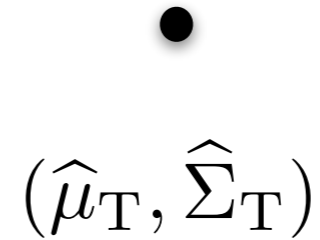
$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

2) “Surround, then Intersect” (SI) Strategy

Source



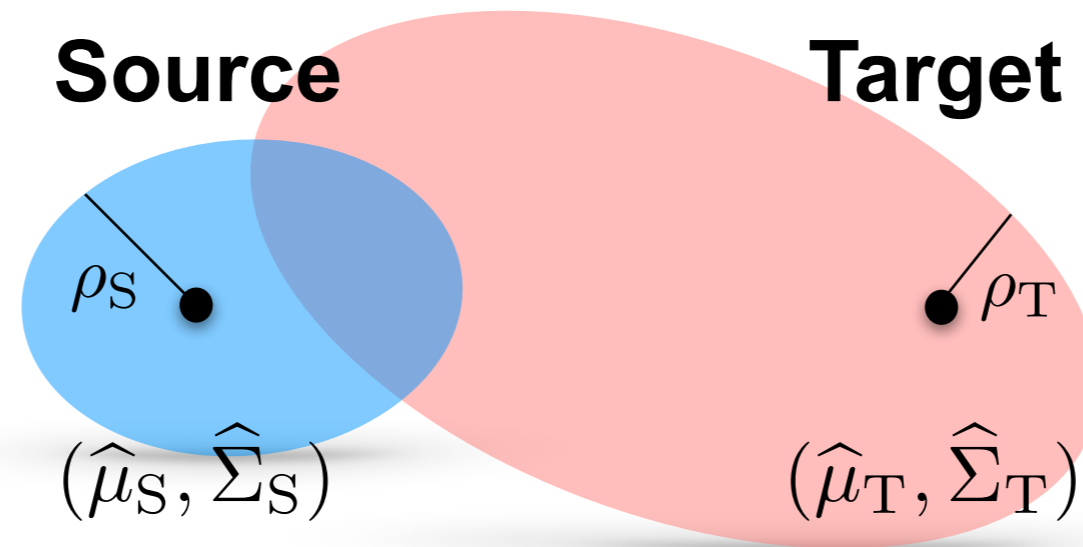
Target



Moment information set:

$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

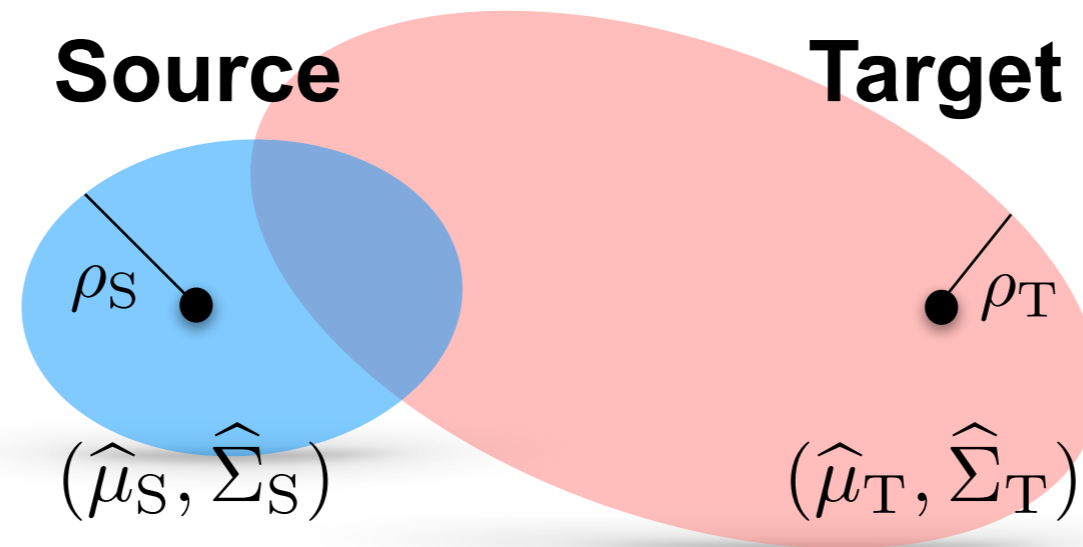
2) “Surround, then Intersect” (SI) Strategy



Moment information set:

$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

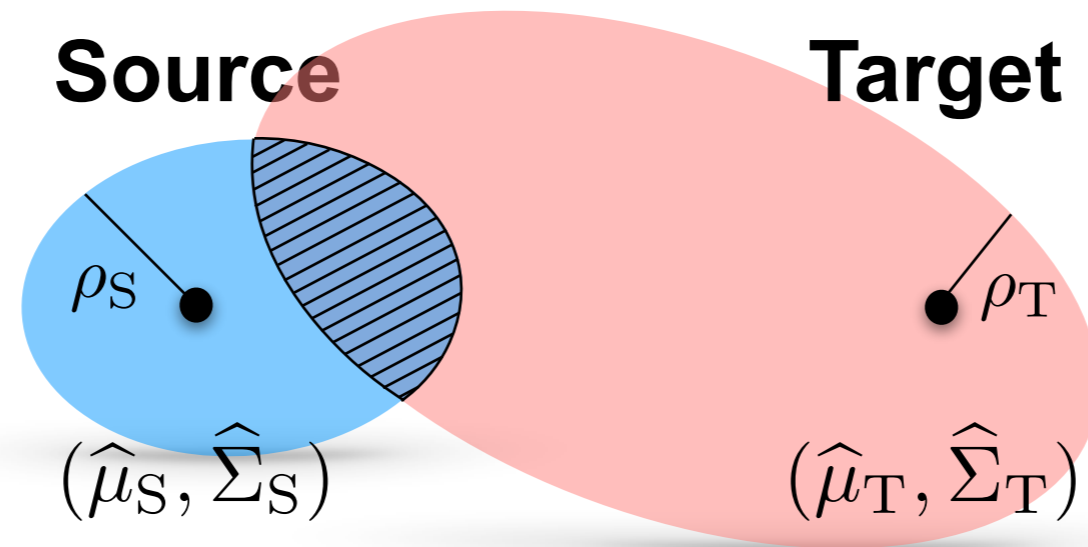
2) “Surround, then Intersect” (SI) Strategy



Moment information set:

$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

2) “Surround, then Intersect” (SI) Strategy



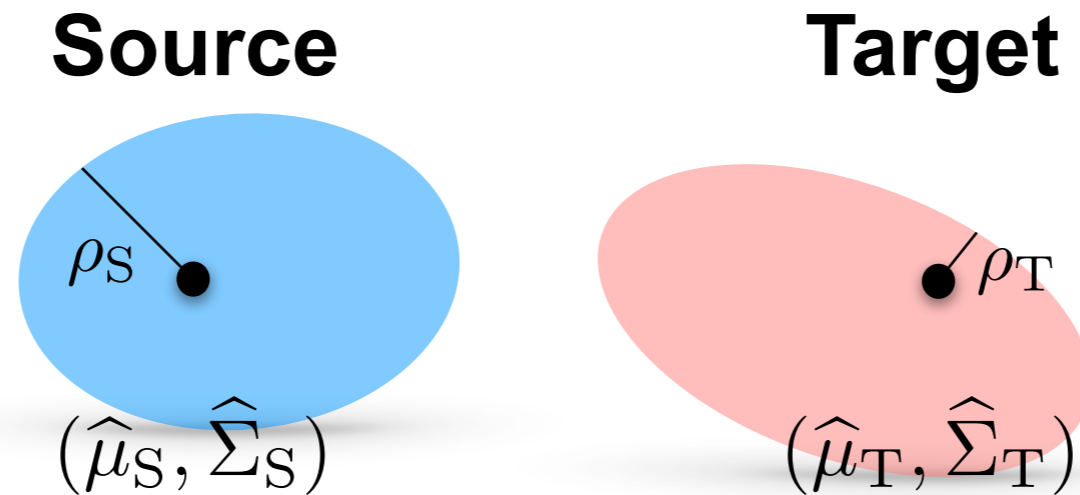
Moment information set:

$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

Distribution set:

$$\mathbb{B}_{\rho_S, \rho_T} = \{ Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}_{\rho_S, \rho_T} \}$$

2) “Surround, then Intersect” (SI) Strategy



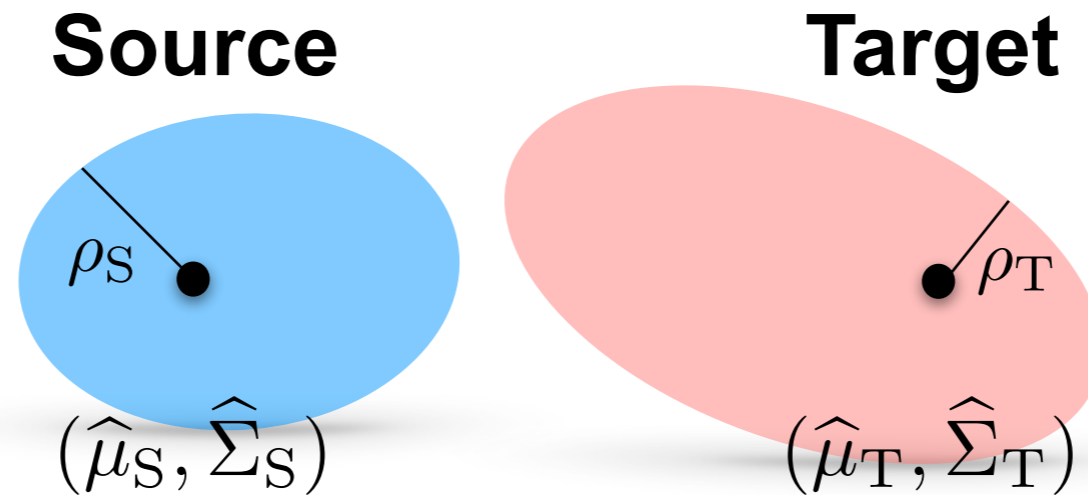
Moment information set:

$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

Distribution set:

$$\mathbb{B}_{\rho_S, \rho_T} = \{ Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}_{\rho_S, \rho_T} \}$$

2) “Surround, then Intersect” (SI) Strategy



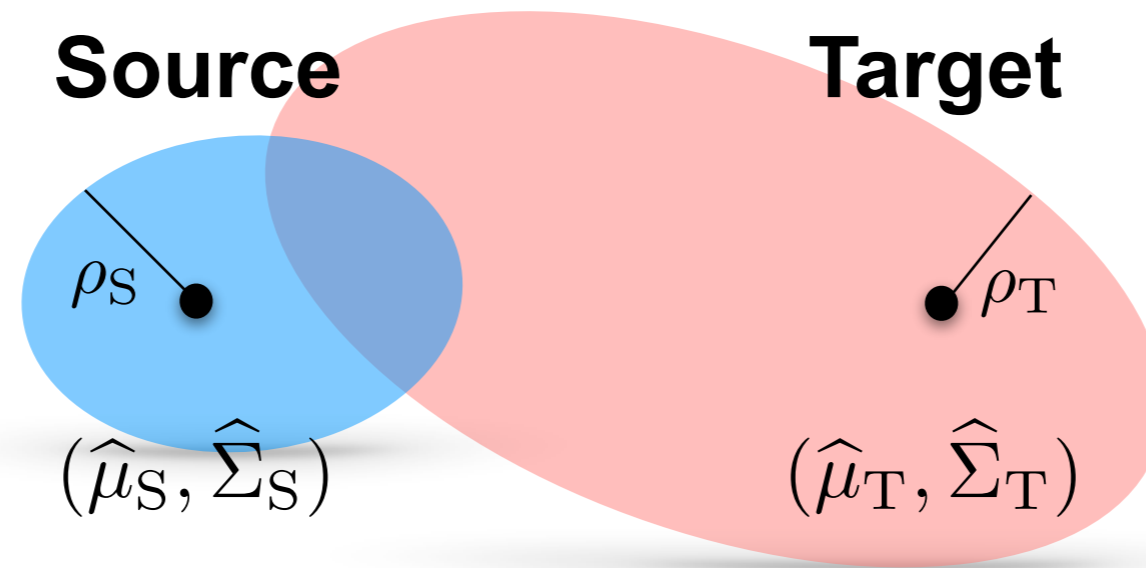
Moment information set:

$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

Distribution set:

$$\mathbb{B}_{\rho_S, \rho_T} = \{ Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}_{\rho_S, \rho_T} \}$$

2) “Surround, then Intersect” (SI) Strategy



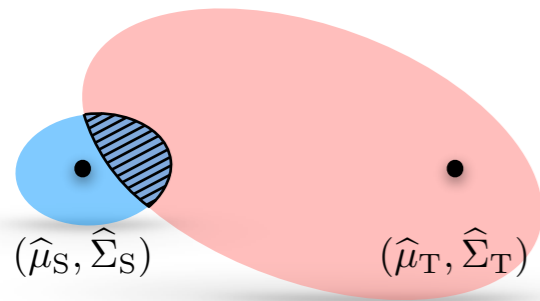
Moment information set:

$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

Distribution set:

$$\mathbb{B}_{\rho_S, \rho_T} = \{ Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}_{\rho_S, \rho_T} \}$$

2) “Surround, then Intersect” (SI) Strategy



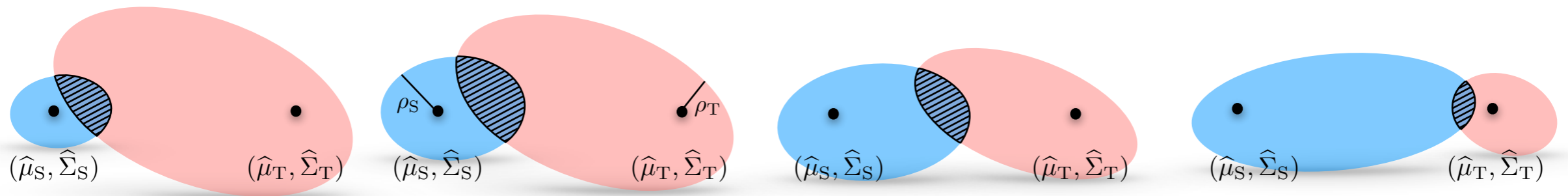
Moment information set:

$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

Distribution set:

$$\mathbb{B}_{\rho_S, \rho_T} = \{ Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}_{\rho_S, \rho_T} \}$$

2) “Surround, then Intersect” (SI) Strategy



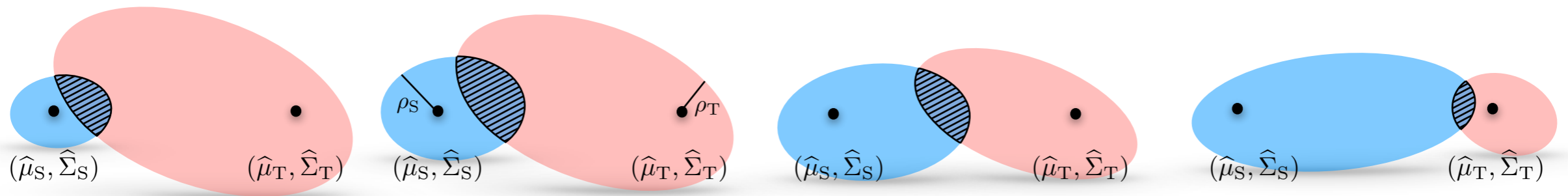
Moment information set:

$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \psi((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \psi((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

Distribution set:

$$\mathbb{B}_{\rho_S, \rho_T} = \{ Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}_{\rho_S, \rho_T} \}$$

2) “Surround, then Intersect” (SI) Strategy



i) Kullback-Leibler (KL)-type Divergence

Moment information set:

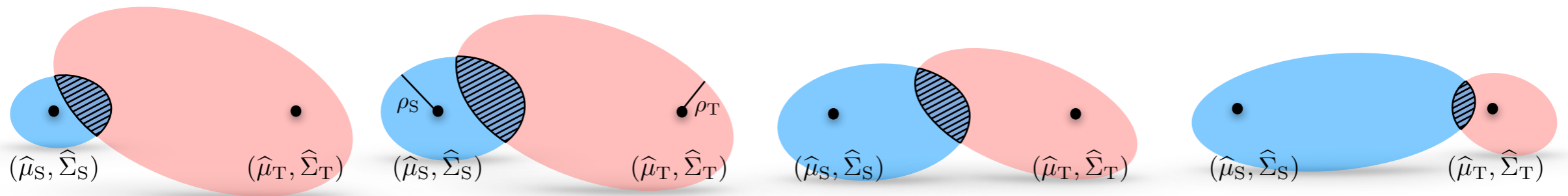
$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \mathbb{D}((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \mathbb{D}((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

Distribution set:

$$\mathbb{B}_{\rho_S, \rho_T} = \{ Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}_{\rho_S, \rho_T} \}$$

Robust Estimation: $\beta^* = (M_{XX}^*)^{-1} M_{XY}^*$ solves $\inf_{\beta \in \mathbb{R}^d} \sup_{Q \in \mathbb{B}_{\rho_S, \rho_T}} \mathbb{E}_Q[(\beta^\top X - Y)^2]$,
 where (M_{XX}^*, M_{XY}^*) is a solution of a convex semidefinite program.

2) “Surround, then Intersect” (SI) Strategy



ii) Wasserstein-type Divergence

Moment information set:

$$\mathbb{U}_{\rho_S, \rho_T} \triangleq \left\{ (\mu, \Sigma) \in \mathbb{R}^p \times \mathbb{S}_+^p \text{ such that: } \mathbb{W}((\mu, \Sigma) \| (\hat{\mu}_S, \hat{\Sigma}_S)) \leq \rho_S, \mathbb{W}((\mu, \Sigma) \| (\hat{\mu}_T, \hat{\Sigma}_T)) \leq \rho_T, \Sigma + \mu\mu^\top \succeq \varepsilon I_p \right\}$$

Distribution set:

$$\mathbb{B}_{\rho_S, \rho_T} = \{ Q \in \mathcal{M}(\mathbb{R}^p) : Q \sim (\mu, \Sigma), (\mu, \Sigma) \in \mathbb{U}_{\rho_S, \rho_T} \}$$

Robust Estimation: $\beta^* = (M_{XX}^*)^{-1} M_{XY}^*$ solves $\inf_{\beta \in \mathbb{R}^d} \sup_{Q \in \mathbb{B}_{\rho_S, \rho_T}} \mathbb{E}_Q[(\beta^\top X - Y)^2]$,
 where (M_{XX}^*, M_{XY}^*) is a solution of a linear semidefinite program.

Numerical Experiments

Data Set	Time	IR-KL	IR-WASS	SI-KL	SI-WASS	CC-L	CC-TL	CC-SL	CC-TE	CC-SE	RWS	LSE-T	LSE-T&S
Uber&Lyft	5	17.65	1.00	199.28	1.01	34.04	98.43	12.03	155.71	1.74	1.45	119.65	11.08
	10	13.67	1.00	111.52	1.01	30.85	99.22	11.40	161.72	1.58	1.34	137.15	6.32
	50	13.39	1.00	60.29	1.01	25.87	85.06	9.72	147.45	1.42	1.16	57.85	2.12
	100	15.24	1.00	59.06	1.01	26.01	85.77	9.91	148.49	1.41	1.12	31.25	1.57
US Births (2018)	5	79.83	1.02	44.71	1.00	64.99	257.60	25.13	432.09	2.07	4.50	727.88	39.17
	10	115.47	1.02	39.35	1.00	45.59	195.14	18.33	339.11	1.60	3.29	524.39	19.28
	50	107.40	1.01	40.04	1.00	42.74	192.46	13.12	361.51	1.31	2.00	191.27	5.20
	100	117.03	1.01	53.13	1.00	45.35	208.65	12.94	397.33	1.22	1.75	104.75	3.19
Life Expectancy	5	33.18	1.00	6.24	1.03	17.24	77.06	7.38	125.71	1.46	1.15	255.08	20.72
	10	25.59	1.00	5.45	1.02	12.49	60.19	5.50	104.00	1.40	1.15	167.15	10.73
	50	19.81	1.00	8.70	1.01	7.57	44.00	3.10	84.98	1.38	1.10	39.83	3.15
	100	19.02	1.00	8.25	1.005	6.82	41.40	2.68	83.60	1.38	1.08	20.42	2.10
House Prices in KC	5	1.58	1.00	1.21	1.01	3.98	8.87	2.12	13.31	1.29	1.23	11.75	3.70
	10	1.52	1.00	1.20	1.01	3.58	7.77	2.02	11.70	1.27	1.23	6.93	2.25
	50	1.34	1.00	1.31	1.01	2.79	6.52	1.86	10.37	1.27	1.20	3.91	1.30
	100	1.34	1.00	1.30	1.01	2.65	6.54	1.91	10.74	1.27	1.18	2.72	1.12
California Housing	5	63.33	1.05	3.31	1.00	27.63	102.82	9.60	181.52	1.35	1.17	96.43	54.34
	10	68.08	1.04	2.42	1.00	20.57	91.86	6.23	169.87	1.19	1.17	45.64	24.76
	50	70.08	1.01	1.97	1.00	11.79	81.72	2.49	170.18	1.05	1.13	10.17	5.63
	100	72.80	1.003	1.90	1.00	9.71	79.19	1.83	173.96	1.04	1.14	5.81	3.39

Numerical Experiments

Data Set	Time	IR-KL	IR-WASS	SI-KL	SI-WASS	CC-L	CC-TL	CC-SL	CC-TE	CC-SE	RWS	LSE-T	LSE-T&S
Uber&Lyft	5	17.65	1.00	199.28	1.01	34.04	98.43	12.03	155.71	1.74	1.45	119.65	11.08
	10	13.67	1.00	111.52	1.01	30.85	99.22	11.40	161.72	1.58	1.34	137.15	6.32
	50	13.39	1.00	60.29	1.01	25.87	85.06	9.72	147.45	1.42	1.16	57.85	2.12
	100	15.24	1.00	59.06	1.01	26.01	85.77	9.91	148.49	1.41	1.12	31.25	1.57
US Births (2018)	5	79.83	1.02	44.71	1.00	64.99	257.60	25.13	432.09	2.07	4.50	727.88	39.17
	10	115.47	1.02	39.35	1.00	45.59	195.14	18.33	339.11	1.60	3.29	524.39	19.28
	50	107.40	1.01	40.04	1.00	42.74	192.46	13.12	361.51	1.31	2.00	191.27	5.20
	100	117.03	1.01	53.13	1.00	45.35	208.65	12.94	397.33	1.22	1.75	104.75	3.19
Life Expectancy	5	33.18	1.00	6.24	1.03	17.24	77.06	7.38	125.71	1.46	1.15	255.08	20.72
	10	25.59	1.00	5.45	1.02	12.49	60.19	5.50	104.00	1.40	1.15	167.15	10.73
	50	19.81	1.00	8.70	1.01	7.57	44.00	3.10	84.98	1.38	1.10	39.83	3.15
	100	19.02	1.00	8.25	1.005	6.82	41.40	2.68	83.60	1.38	1.08	20.42	2.10
House Prices in KC	5	1.58	1.00	1.21	1.01	3.98	8.87	2.12	13.31	1.29	1.23	11.75	3.70
	10	1.52	1.00	1.20	1.01	3.58	7.77	2.02	11.70	1.27	1.23	6.93	2.25
	50	1.34	1.00	1.31	1.01	2.79	6.52	1.86	10.37	1.27	1.20	3.91	1.30
	100	1.34	1.00	1.30	1.01	2.65	6.54	1.91	10.74	1.27	1.18	2.72	1.12
California Housing	5	63.33	1.05	3.31	1.00	27.63	102.82	9.60	181.52	1.35	1.17	96.43	54.34
	10	68.08	1.04	2.42	1.00	20.57	91.86	6.23	169.87	1.19	1.17	45.64	24.76
	50	70.08	1.01	1.97	1.00	11.79	81.72	2.49	170.18	1.05	1.13	10.17	5.63
	100	72.80	1.003	1.90	1.00	9.71	79.19	1.83	173.96	1.04	1.14	5.81	3.39

Numerical Experiments

Data Set	Time	IR-KL	IR-WASS	SI-KL	SI-WASS	CC-L	CC-TL	CC-SL	CC-TE	CC-SE	RWS	LSE-T	LSE-T&S
Uber&Lyft	5	17.65	1.00	199.28	1.01	34.04	98.43	12.03	155.71	1.74	1.45	119.65	11.08
	10	13.67	1.00	111.52	1.01	30.85	99.22	11.40	161.72	1.58	1.34	137.15	6.32
	50	13.39	1.00	60.29	1.01	25.87	85.06	9.72	147.45	1.42	1.16	57.85	2.12
	100	15.24	1.00	59.06	1.01	26.01	85.77	9.91	148.49	1.41	1.12	31.25	1.57
US Births (2018)	5	79.83	1.02	44.71	1.00	64.99	257.60	25.13	432.09	2.07	4.50	727.88	39.17
	10	115.47	1.02	39.35	1.00	45.59	195.14	18.33	339.11	1.60	3.29	524.39	19.28
	50	107.40	1.01	40.04	1.00	42.74	192.46	13.12	361.51	1.31	2.00	191.27	5.20
	100	117.03	1.01	53.13	1.00	45.35	208.65	12.94	397.33	1.22	1.75	104.75	3.19
Life Expectancy	5	33.18	1.00	6.24	1.03	17.24	77.06	7.38	125.71	1.46	1.15	255.08	20.72
	10	25.59	1.00	5.45	1.02	12.49	60.19	5.50	104.00	1.40	1.15	167.15	10.73
	50	19.81	1.00	8.70	1.01	7.57	44.00	3.10	84.98	1.38	1.10	39.83	3.15
	100	19.02	1.00	8.25	1.005	6.82	41.40	2.68	83.60	1.38	1.08	20.42	2.10
House Prices in KC	5	1.58	1.00	1.21	1.01	3.98	8.87	2.12	13.31	1.29	1.23	11.75	3.70
	10	1.52	1.00	1.20	1.01	3.58	7.77	2.02	11.70	1.27	1.23	6.93	2.25
	50	1.34	1.00	1.31	1.01	2.79	6.52	1.86	10.37	1.27	1.20	3.91	1.30
	100	1.34	1.00	1.30	1.01	2.65	6.54	1.91	10.74	1.27	1.18	2.72	1.12
California Housing	5	63.33	1.05	3.31	1.00	27.63	102.82	9.60	181.52	1.35	1.17	96.43	54.34
	10	68.08	1.04	2.42	1.00	20.57	91.86	6.23	169.87	1.19	1.17	45.64	24.76
	50	70.08	1.01	1.97	1.00	11.79	81.72	2.49	170.18	1.05	1.13	10.17	5.63
	100	72.80	1.003	1.90	1.00	9.71	79.19	1.83	173.96	1.04	1.14	5.81	3.39

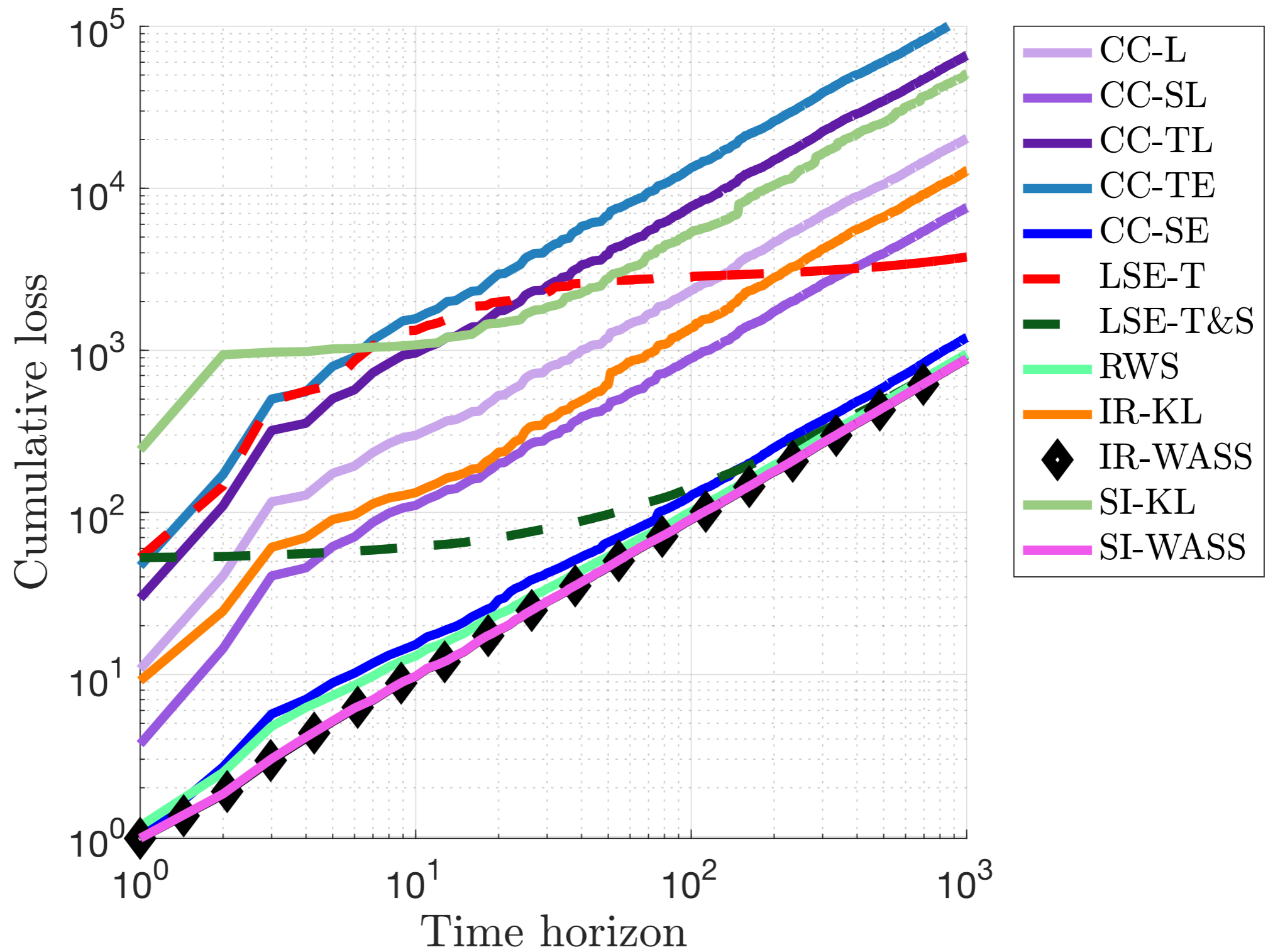
Numerical Experiments

Data Set	Time	IR-KL	IR-WASS	SI-KL	SI-WASS	CC-L	CC-TL	CC-SL	CC-TE	CC-SE	RWS	LSE-T	LSE-T&S
Uber&Lyft	5	17.65	1.00	199.28	1.01	34.04	98.43	12.03	155.71	1.74	1.45	119.65	11.08
	10	13.67	1.00	111.52	1.01	30.85	99.22	11.40	161.72	1.58	1.34	137.15	6.32
	50	13.39	1.00	60.29	1.01	25.87	85.06	9.72	147.45	1.42	1.16	57.85	2.12
	100	15.24	1.00	59.06	1.01	26.01	85.77	9.91	148.49	1.41	1.12	31.25	1.57
US Births (2018)	5	79.83	1.02	44.71	1.00	64.99	257.60	25.13	432.09	2.07	4.50	727.88	39.17
	10	115.47	1.02	39.35	1.00	45.59	195.14	18.33	339.11	1.60	3.29	524.39	19.28
	50	107.40	1.01	40.04	1.00	42.74	192.46	13.12	361.51	1.31	2.00	191.27	5.20
	100	117.03	1.01	53.13	1.00	45.35	208.65	12.94	397.33	1.22	1.75	104.75	3.19
Life Expectancy	5	33.18	1.00	6.24	1.03	17.24	77.06	7.38	125.71	1.46	1.15	255.08	20.72
	10	25.59	1.00	5.45	1.02	12.49	60.19	5.50	104.00	1.40	1.15	167.15	10.73
	50	19.81	1.00	8.70	1.01	7.57	44.00	3.10	84.98	1.38	1.10	39.83	3.15
	100	19.02	1.00	8.25	1.005	6.82	41.40	2.68	83.60	1.38	1.08	20.42	2.10
House Prices in KC	5	1.58	1.00	1.21	1.01	3.98	8.87	2.12	13.31	1.29	1.23	11.75	3.70
	10	1.52	1.00	1.20	1.01	3.58	7.77	2.02	11.70	1.27	1.23	6.93	2.25
	50	1.34	1.00	1.31	1.01	2.79	6.52	1.86	10.37	1.27	1.20	3.91	1.30
	100	1.34	1.00	1.30	1.01	2.65	6.54	1.91	10.74	1.27	1.18	2.72	1.12
California Housing	5	63.33	1.05	3.31	1.00	27.63	102.82	9.60	181.52	1.35	1.17	96.43	54.34
	10	68.08	1.04	2.42	1.00	20.57	91.86	6.23	169.87	1.19	1.17	45.64	24.76
	50	70.08	1.01	1.97	1.00	11.79	81.72	2.49	170.18	1.05	1.13	10.17	5.63
	100	72.80	1.003	1.90	1.00	9.71	79.19	1.83	173.96	1.04	1.14	5.81	3.39

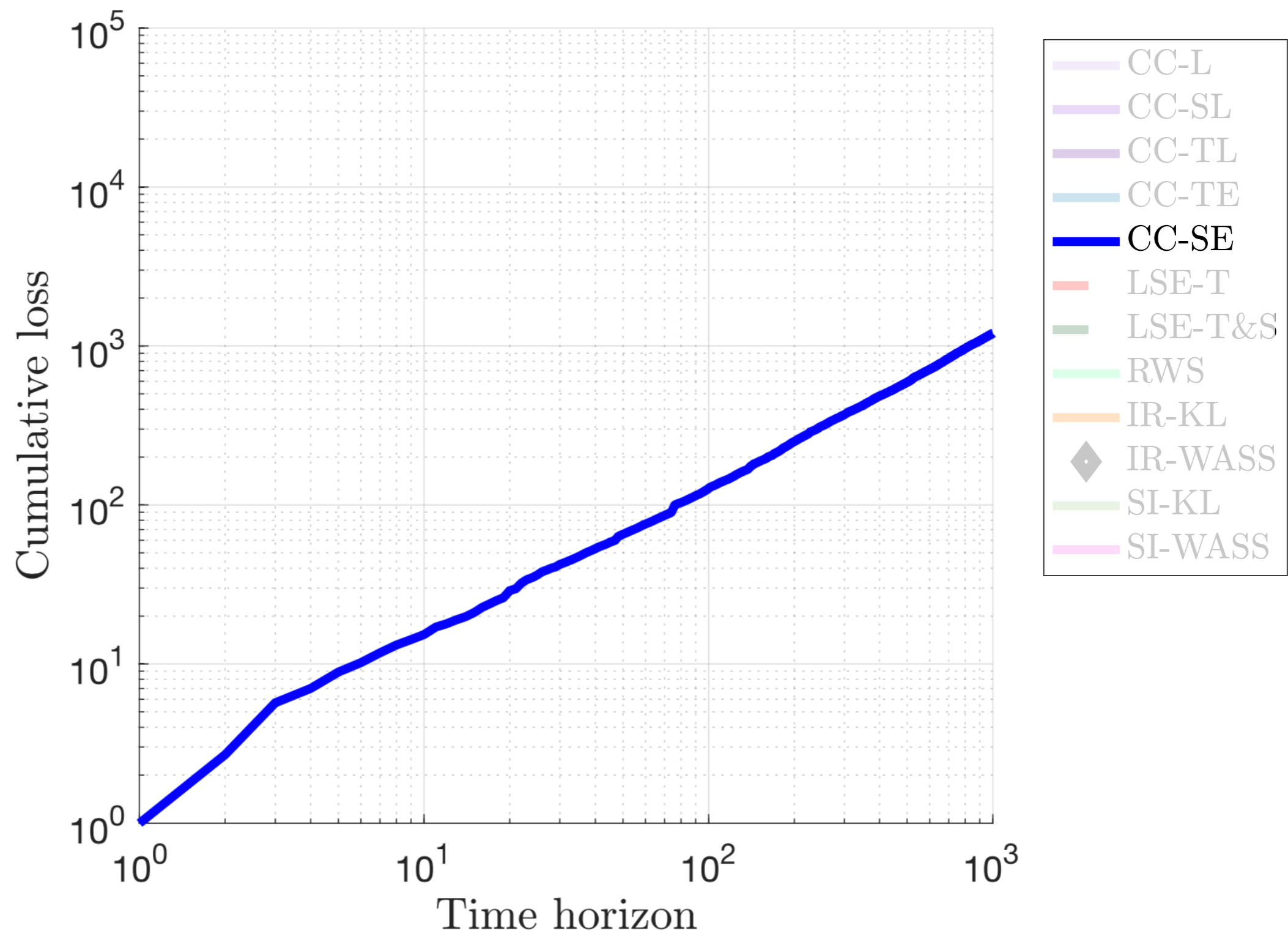
Numerical Experiments

Data Set	Time	IR-KL	IR-WASS	SI-KL	SI-WASS	CC-L	CC-TL	CC-SL	CC-TE	CC-SE	RWS	LSE-T	LSE-T&S
Uber&Lyft	5	17.65	1.00	199.28	1.01	34.04	98.43	12.03	155.71	1.74	1.45	119.65	11.08
	10	13.67	1.00	111.52	1.01	30.85	99.22	11.40	161.72	1.58	1.34	137.15	6.32
	50	13.39	1.00	60.29	1.01	25.87	85.06	9.72	147.45	1.42	1.16	57.85	2.12
	100	15.24	1.00	59.06	1.01	26.01	85.77	9.91	148.49	1.41	1.12	31.25	1.57
US Births (2018)	5	79.83	1.02	44.71	1.00	64.99	257.60	25.13	432.09	2.07	4.50	727.88	39.17
	10	115.47	1.02	39.35	1.00	45.59	195.14	18.33	339.11	1.60	3.29	524.39	19.28
	50	107.40	1.01	40.04	1.00	42.74	192.46	13.12	361.51	1.31	2.00	191.27	5.20
	100	117.03	1.01	53.13	1.00	45.35	208.65	12.94	397.33	1.22	1.75	104.75	3.19
Life Expectancy	5	33.18	1.00	6.24	1.03	17.24	77.06	7.38	125.71	1.46	1.15	255.08	20.72
	10	25.59	1.00	5.45	1.02	12.49	60.19	5.50	104.00	1.40	1.15	167.15	10.73
	50	19.81	1.00	8.70	1.01	7.57	44.00	3.10	84.98	1.38	1.10	39.83	3.15
	100	19.02	1.00	8.25	1.005	6.82	41.40	2.68	83.60	1.38	1.08	20.42	2.10
House Prices in KC	5	1.58	1.00	1.21	1.01	3.98	8.87	2.12	13.31	1.29	1.23	11.75	3.70
	10	1.52	1.00	1.20	1.01	3.58	7.77	2.02	11.70	1.27	1.23	6.93	2.25
	50	1.34	1.00	1.31	1.01	2.79	6.52	1.86	10.37	1.27	1.20	3.91	1.30
	100	1.34	1.00	1.30	1.01	2.65	6.54	1.91	10.74	1.27	1.18	2.72	1.12
California Housing	5	63.33	1.05	3.31	1.00	27.63	102.82	9.60	181.52	1.35	1.17	96.43	54.34
	10	68.08	1.04	2.42	1.00	20.57	91.86	6.23	169.87	1.19	1.17	45.64	24.76
	50	70.08	1.01	1.97	1.00	11.79	81.72	2.49	170.18	1.05	1.13	10.17	5.63
	100	72.80	1.003	1.90	1.00	9.71	79.19	1.83	173.96	1.04	1.14	5.81	3.39

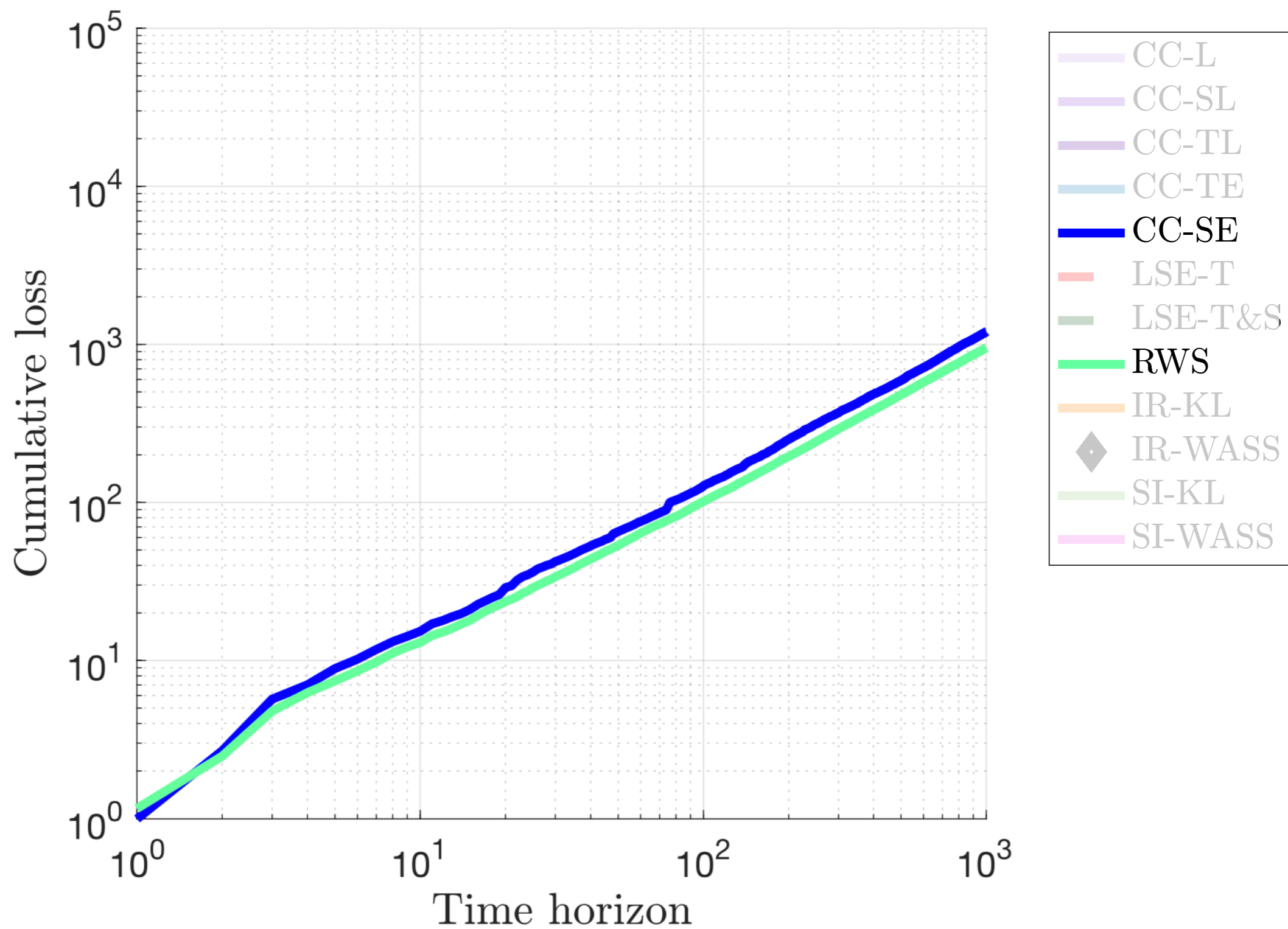
Numerical Experiments - Uber&Lyft



Numerical Experiments



Numerical Experiments



Concluding Remarks

- KL-type divergence based methods are **less** numerically stable
- Asymmetry of the KL-type divergence
- Extrapolating schemes
- Favourable properties of different divergences
- Theoretical guarantees for the out-of-sample performance

References

- [1] Taskesen B., Yue M.C., Blanchet J., Kuhn D, Nguyen V.A.. **Sequential Domain Adaptation by Synthesizing Distributionally Robust Experts**, *ICML 2021*
- [2] Garcke, J. and Vanck, T. Importance weighted inductive transfer learning for regression. *In Joint European conference on machine learning and knowledge discovery in databases*, pp. 466–481, 2014
- [3] Wintenberger, O. Optimal learning with Bernstein online aggregation. *Machine Learning*, 106(1):119–141, 2017.



Thanks a lot for your attention!