# Toward a theory of multi-objective learning

Geelon So, UC San Diego

October 12 - 16, 2025
Symposium on Mathematical Foundations of Trustworthy Learning



Tobias Wegel



Junhyung Park



Fanny Yang

# Motivating example: self-driving car



Goal: train a model for a self-driving car.

$$\min_{h \in \mathcal{H}} R(h)$$

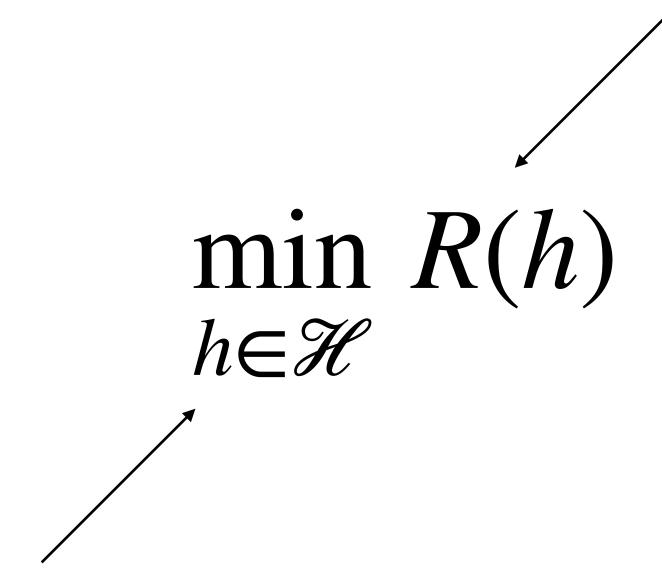
$$\min_{h \in \mathcal{H}} R(h)$$

 $h: X \to Y$  is a **model** from a model class  $\mathcal{H}$ 

"Under situation x, the car should do y."

R(h) measures the **population risk** of the model h

"This is the expected fuel efficiency of h on highways."

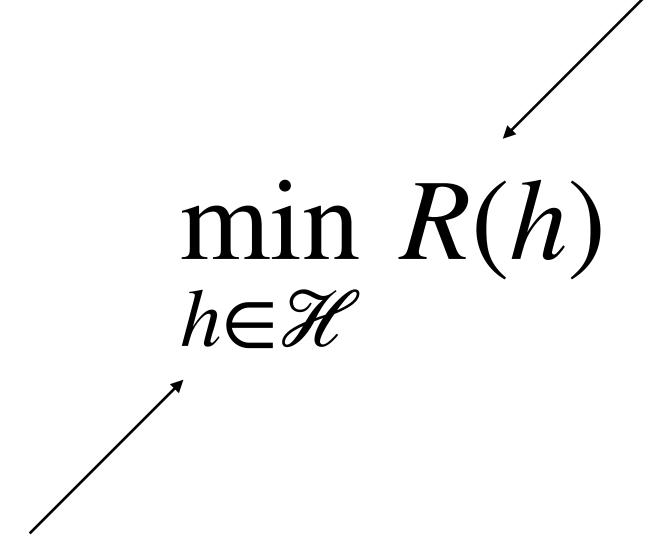


 $h: X \to Y$  is a **model** from a model class  $\mathcal{H}$ 

"Under situation x, the car should do y."

R(h) measures the **population risk** of the model h

"This is the expected fuel efficiency of h on highways."



The learning problem

Directly optimizing R is not possible since we only have sample access to it.

 $h: X \to Y$  is a **model** from a model class  $\mathcal{H}$ 

"Under situation x, the car should do y."

## Motivating example: self-driving car



Goal: train a model for a self-driving car.

#### What we might care about:

- Safety
- Efficiency
- Comfort
- [many more]

#### Multi-objective risk minimization

$$\min_{h \in \mathcal{H}} \mathbf{R}(h) \equiv \left( R_1(h), \dots, R_K(h) \right)$$

Now, we care about many types of risks  $\mathbf{R}(h)$ .

#### Multi-objective risk minimization

$$\min_{h \in \mathcal{H}} \mathbf{R}(h) \equiv \left( R_1(h), \dots, R_K(h) \right)$$

Now, we care about many types of risks  $\mathbf{R}(h)$ .

• It is not usually possible to minimize all risks simultaneously.

Let  $X \times Y$  be a data space.

Let  $X \times Y$  be a data space. For each objective  $k \in [K]$ :

• Let  $R_k$  measure the risk of a standard supervised learning task:

Let  $X \times Y$  be a data space. For each objective  $k \in [K]$ :

- Let  $R_k$  measure the risk of a standard supervised learning task:
  - $P_k$  a data distribution

Let  $X \times Y$  be a data space. For each objective  $k \in [K]$ :

- Let  $R_k$  measure the risk of a standard supervised learning task:
  - $P_k$  a data distribution
  - $\ell_k(y,\hat{y})$  measures the loss of predicting  $\hat{y}$  when correct answer is y

$$R_k(h) = \mathbb{E}_{P_k} \left[ \mathcal{C}_k (y, h(x)) \right]$$

Let  $X \times Y$  be a data space. For each objective  $k \in [K]$ :

- Let  $R_k$  measure the risk of a standard supervised learning task:
  - $P_k$  a data distribution
  - $\ell_k(y,\hat{y})$  measures the loss of predicting  $\hat{y}$  when correct answer is y

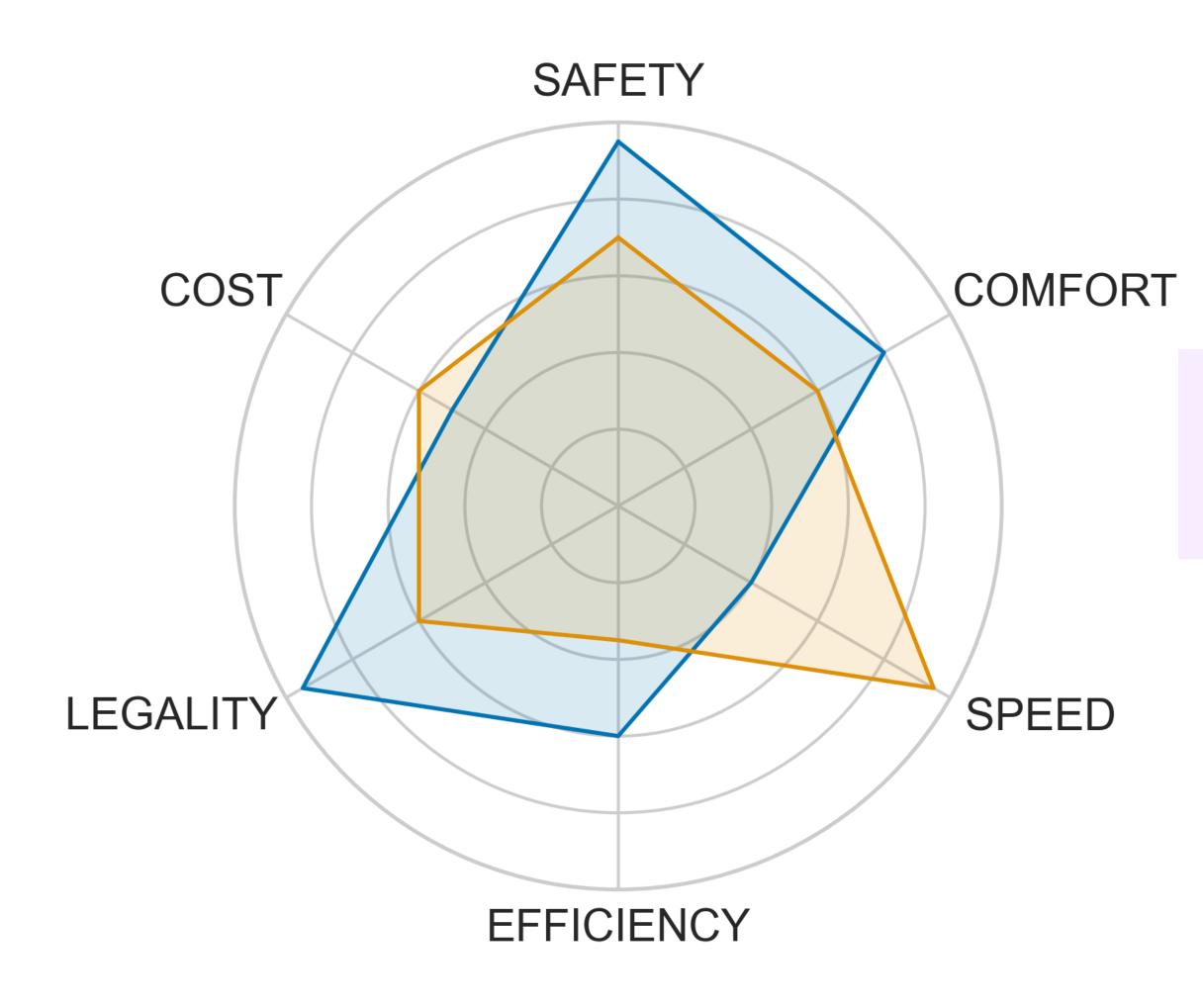
$$R_k(h) = \mathbb{E}_{P_k} \left[ \mathscr{C}_k \big( y, h(x) \big) \right]$$

•  $f_k^{\star}$  is the Bayes-optimal model minimizing  $R_k$ 

# I. Background

Question: what does it even mean to minimize  $\mathbf{R}(h)$ ?

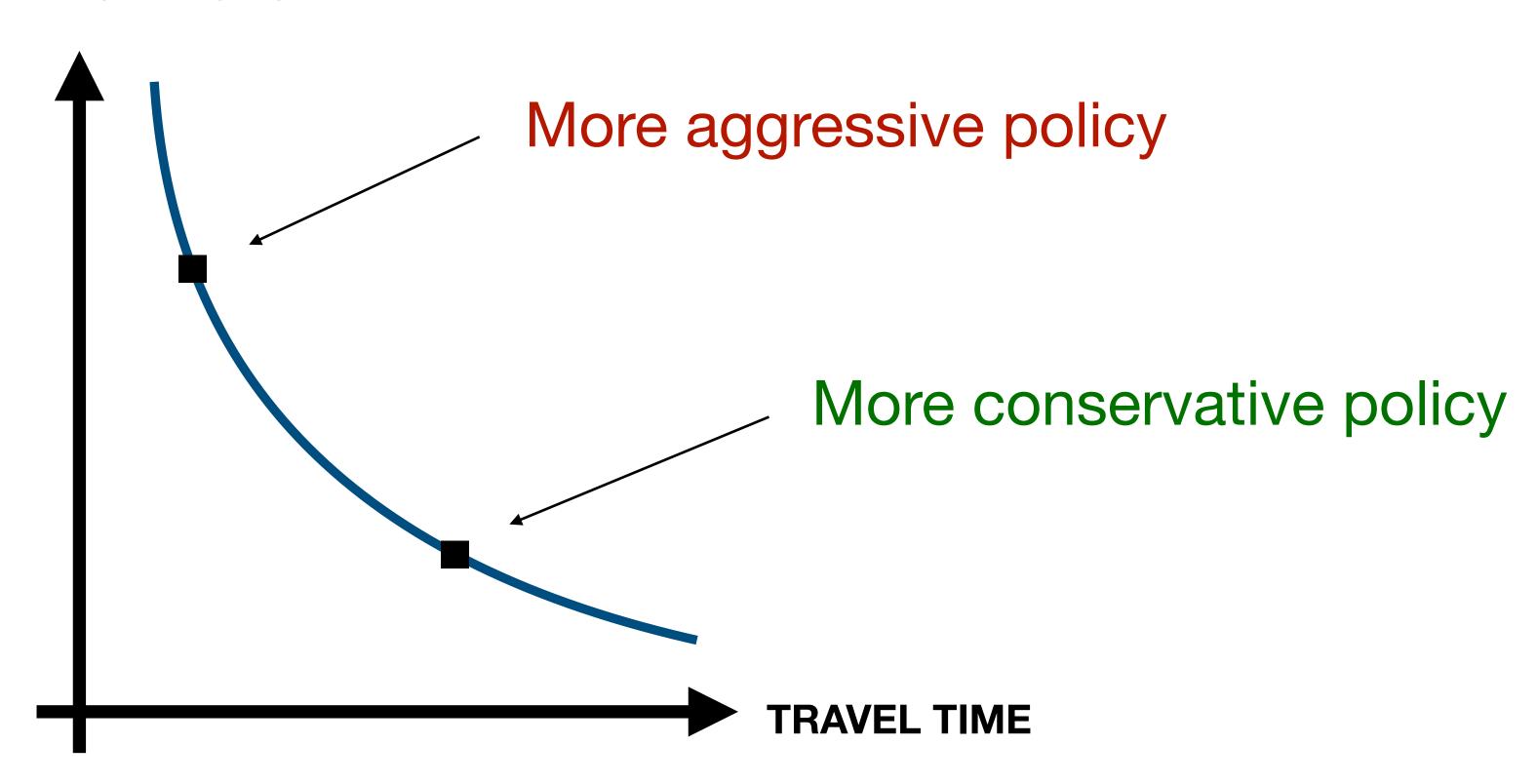
## Solution concept: Pareto optimality



A model is **Pareto optimal** if improving one objective must come at a cost of another.

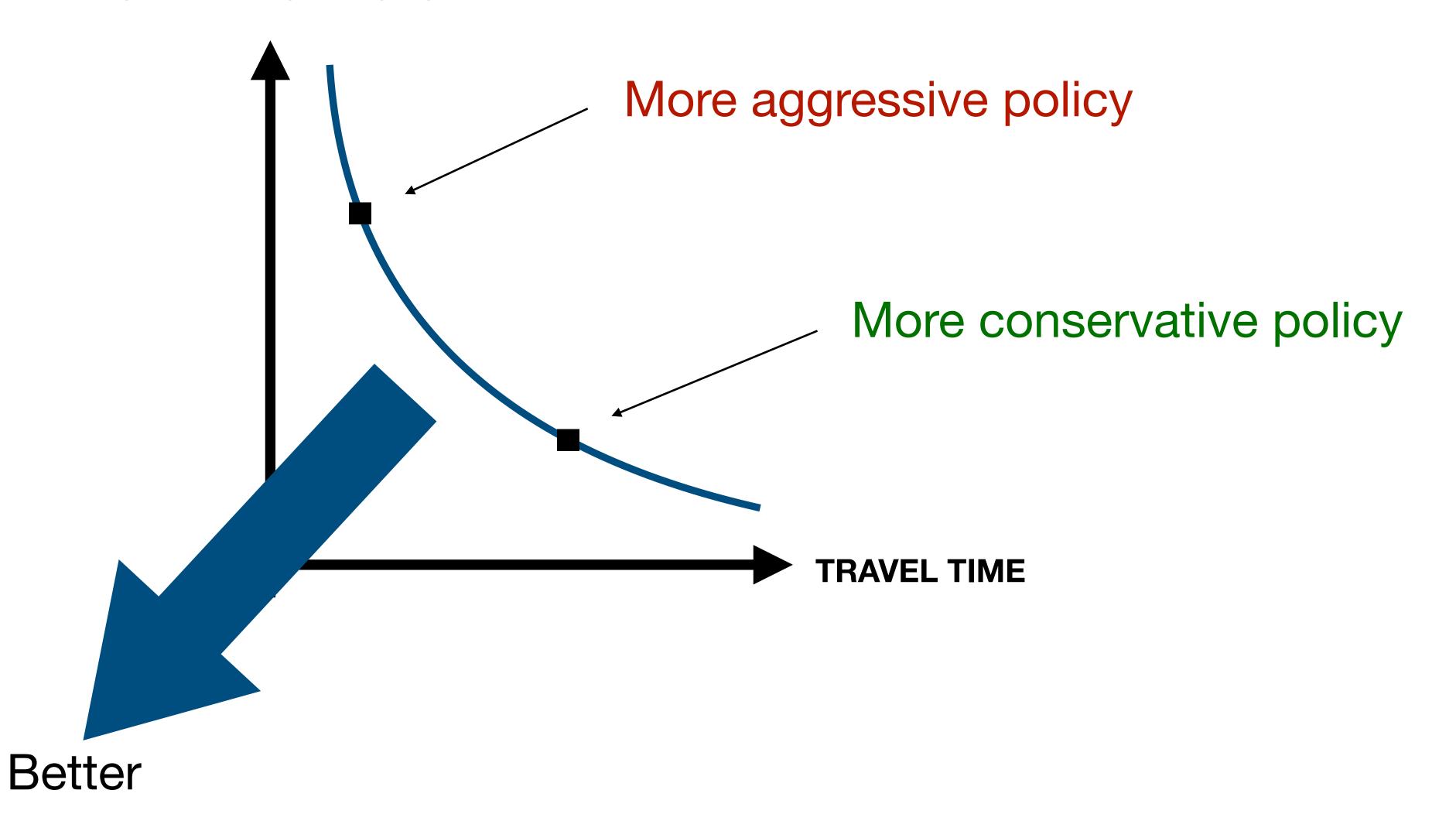
#### Visualization: Pareto front

#### **SAFETY VIOLATIONS**



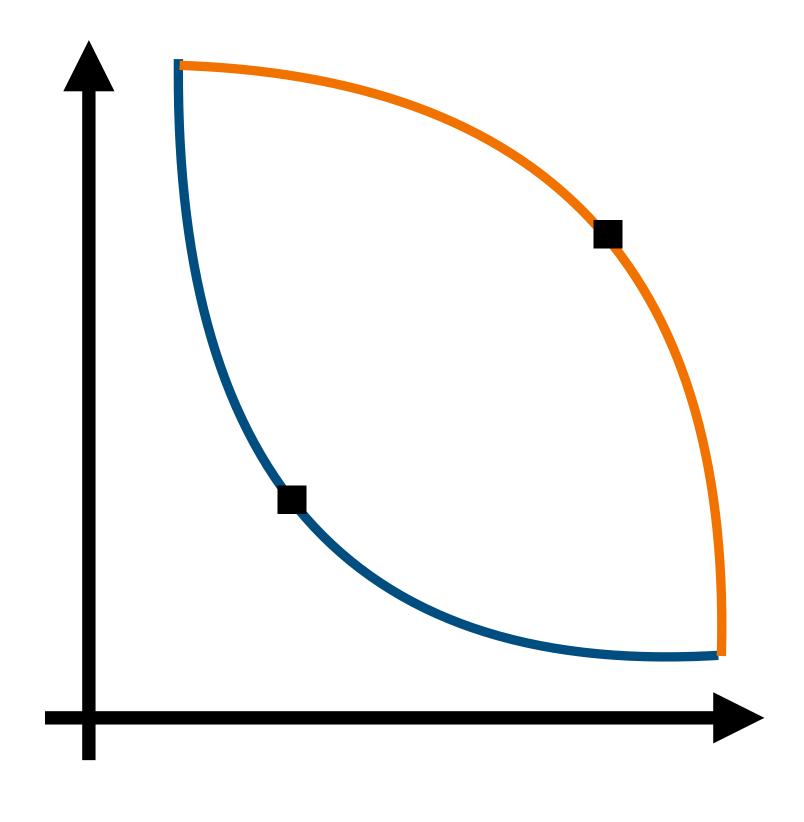
#### Visualization: Pareto front

#### SAFETY VIOLATIONS



#### Which trade-off to pick?

#### **SAFETY VIOLATIONS**



The specific type of *trade-off* that is preferred is not usually known beforehand.

Preference may depend on the Pareto set itself.

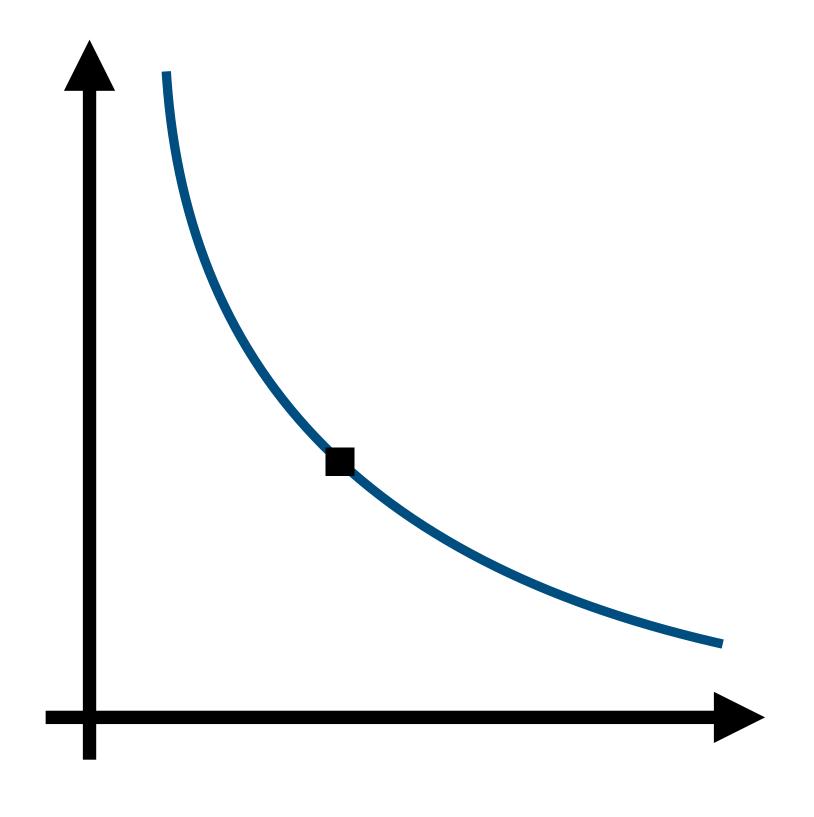
TRAVEL TIME

#### Main goal: learn all Pareto-optimal models

(Allows for downstream user preferences to dictate which one to deploy).

Question: how do you find Pareto-optimal models?

#### Scalarization



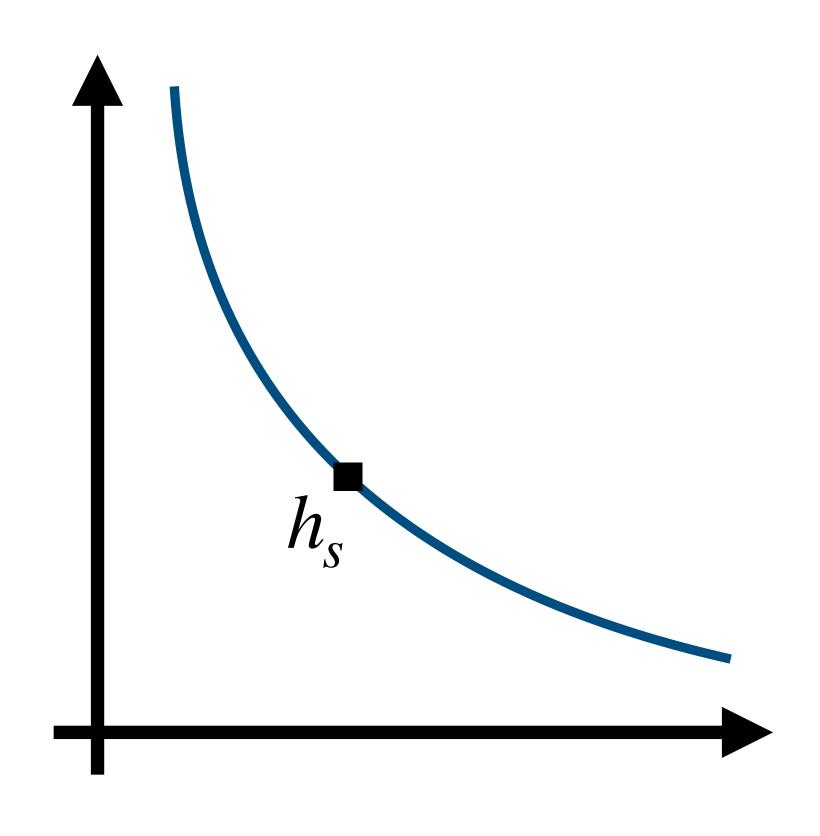
Pareto-optimal models often corresponds to solutions of *scalarized* objectives.

#### **Examples:**

$$\min_{h \in \mathcal{H}} \frac{1}{K} \sum_{k \in [K]} R_k(h) \quad \text{or} \quad \min_{h \in \mathcal{H}} \max_{k \in [K]} R_k(h) \,.$$
 Minimize task-averaged risk

Minimize worst risk

#### The s-trade-off solutions



- Let  $s: \mathbb{R}^K \to \mathbb{R}$  be a scalarizer.
- The s-trade-off solution minimizes:

$$h_s = \arg\min_{h \in \mathcal{H}} s(\mathbf{R}(h)).$$

• The Pareto set can be parameterized by a family  $\{h_{\scriptscriptstyle S}:s\in\mathcal{S}\}$ 

## What learning the Pareto set means

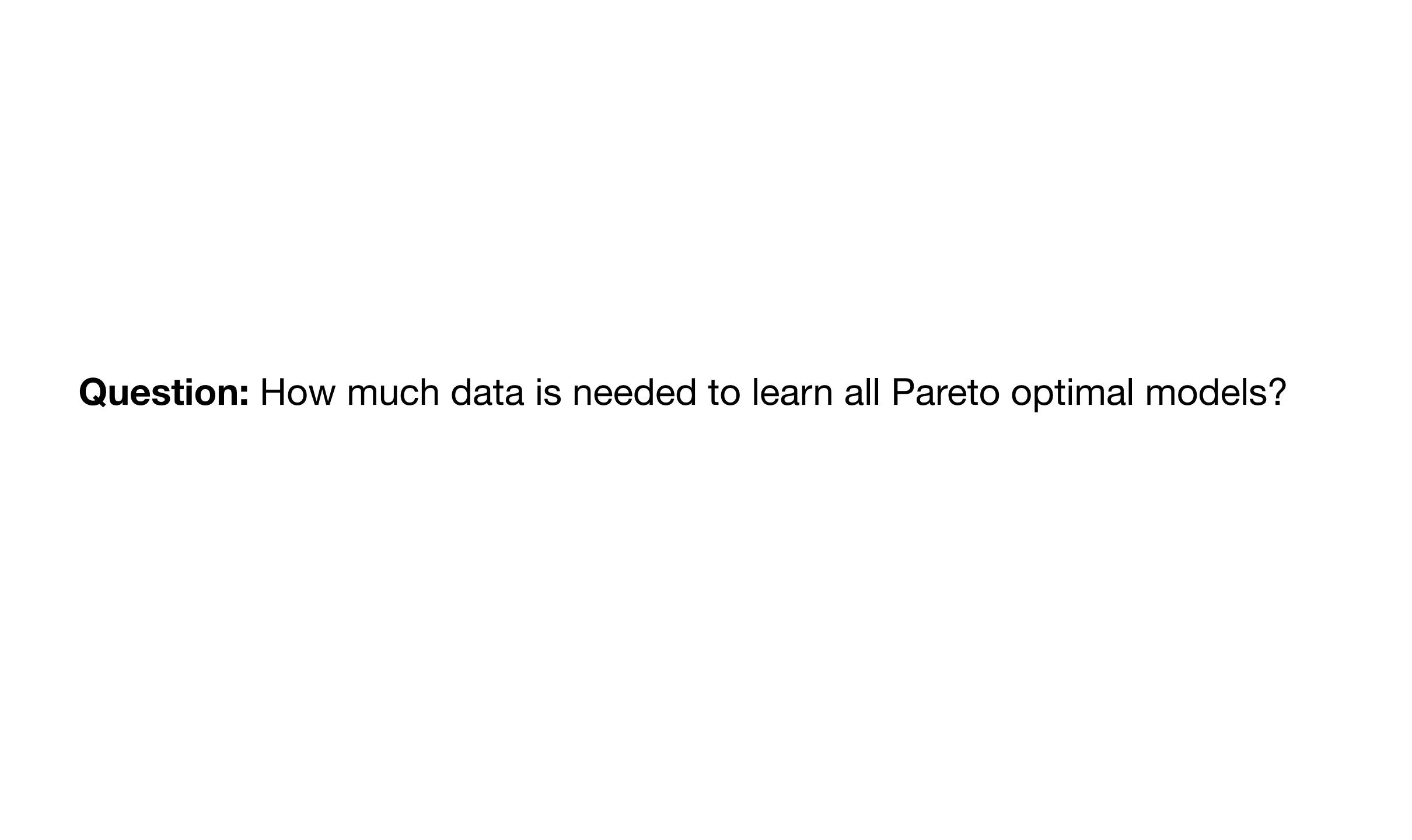
$$s(\mathbf{R}(\hat{h}_s)) < \min_{h \in \mathcal{H}} s(\mathbf{R}(h)) + \varepsilon.$$

• Return  $\varepsilon$ -optimal model for each s-trade-off where  $s \in \mathcal{S}$ .

#### What learning the Pareto set means

$$\Pr\bigg(\forall s \in \mathcal{S}, \quad s\big(\mathbf{R}(\hat{h}_s)\big) < \min_{h \in \mathcal{H}} s\big(\mathbf{R}(h)\big) + \varepsilon\bigg) > 1 - \delta.$$

- Return  $\varepsilon$ -optimal model for each s-trade-off where  $s \in \mathcal{S}$ .
- Succeed on the whole Pareto set with high probability.



#### **Supervised multi-objective learning**

C. Cortes, M. Mohri, J. Gonzalvo, and D. Storcheus. *Agnostic learning with multiple objectives*. NeurIPS 2020.

P. Súkeník and C. Lampert. *Generalization in multi-objective machine learning*. Neural Computing and Applications 2024.

Informal Theorem (Súkeník & Lampert 2024). Let  ${\mathscr H}$  be a model class.

To  $\varepsilon$ -learn all Pareto optimal models in  $\mathcal{H}$ , we need:

$$\tilde{O}\left(\frac{\mathrm{VC}(\mathcal{H})\cdot K}{\varepsilon^2}\right)$$
 samples,

where K is the number of tasks.

#### Supervised multi-objective learning

C. Cortes, M. Mohri, J. Gonzalvo, and D. Storcheus. *Agnostic learning with multiple objectives*. NeurIPS 2020.

P. Súkeník and C. Lampert. *Generalization in multi-objective machine learning*. Neural Computing and Applications 2024.

Informal Theorem (Súkeník & Lampert 2024). Let  ${\mathscr H}$  be a model class.

To  $\varepsilon$ -learn all Pareto optimal models in  $\mathcal{H}$ , we need:

$$\tilde{\Theta}\left(\frac{\mathrm{VC}(\mathcal{H})\cdot K}{\varepsilon^2}\right) \text{ samples,}$$

where K is the number of tasks.

Tight because solving each task alone costs: 
$$\Omega\left(\frac{\mathrm{VC}(\mathcal{H})}{\varepsilon^2}\right)$$
 samples.

Informal Theorem (Súkeník & Lampert 2024). Let  ${\mathscr H}$  be a model class.

To  $\varepsilon$ -learn all Pareto optimal models in  $\mathcal{H}$ , we need:

$$\tilde{\Theta}\left(\frac{\mathrm{VC}(\mathcal{H})\cdot K}{\varepsilon^2}\right) \text{ samples,}$$

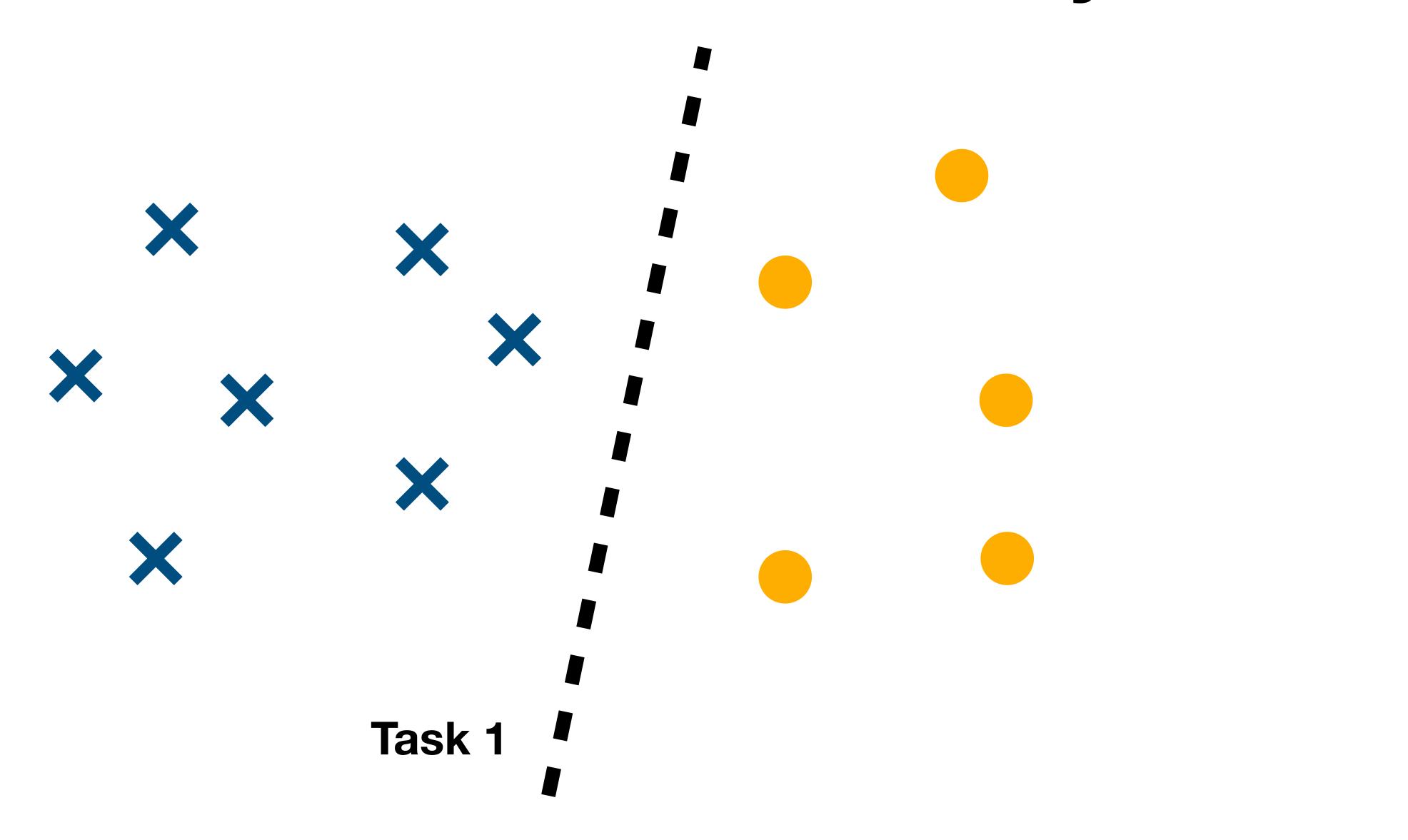
where K is the number of tasks.

Tight because **solving each task alone** costs: 
$$\Omega\left(\frac{\mathrm{VC}(\mathcal{H})}{\varepsilon^2}\right)$$
 samples.

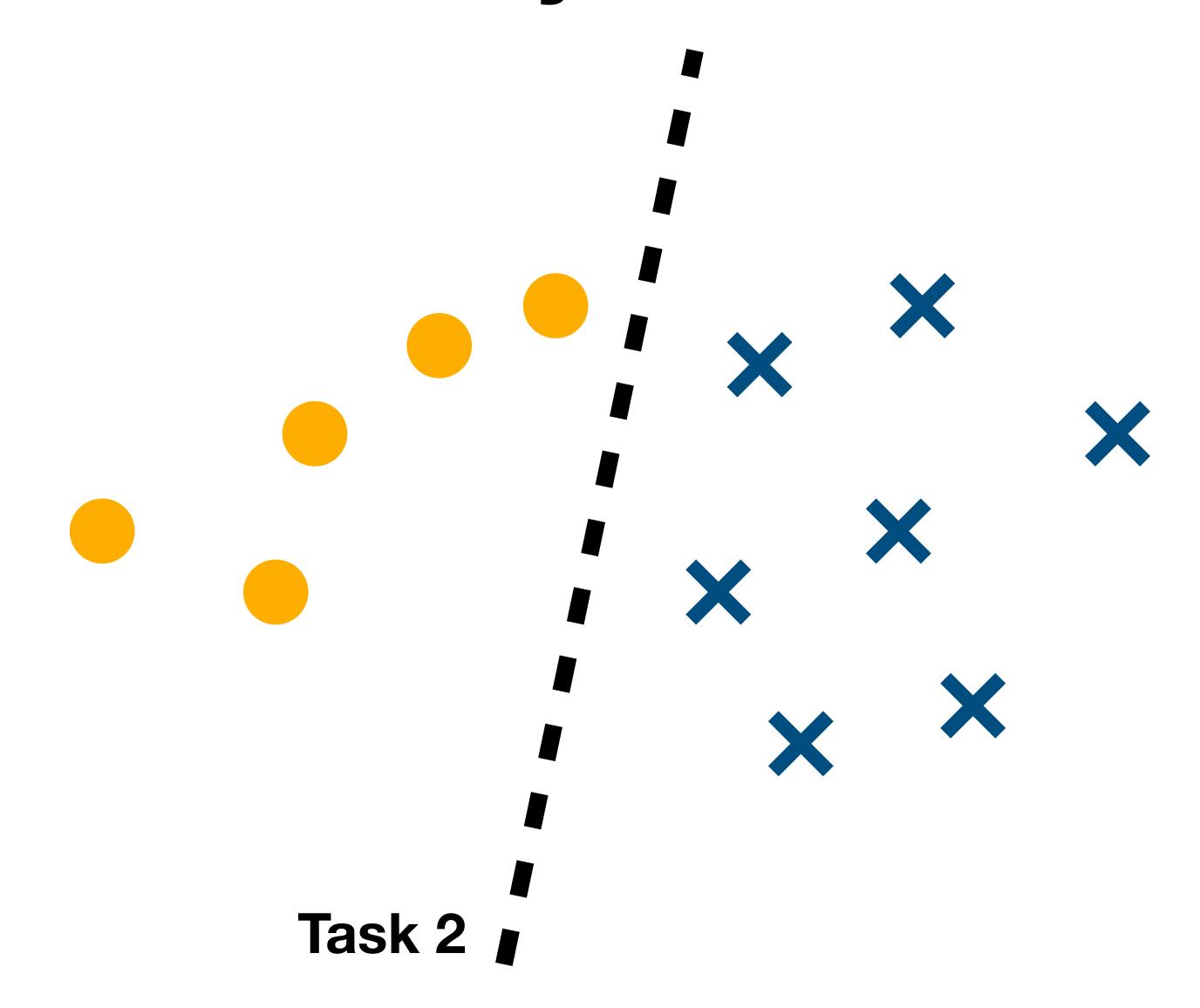
Is this really where the statistical hardness is coming from?

## II. Multi-objective learning for simple tasks?

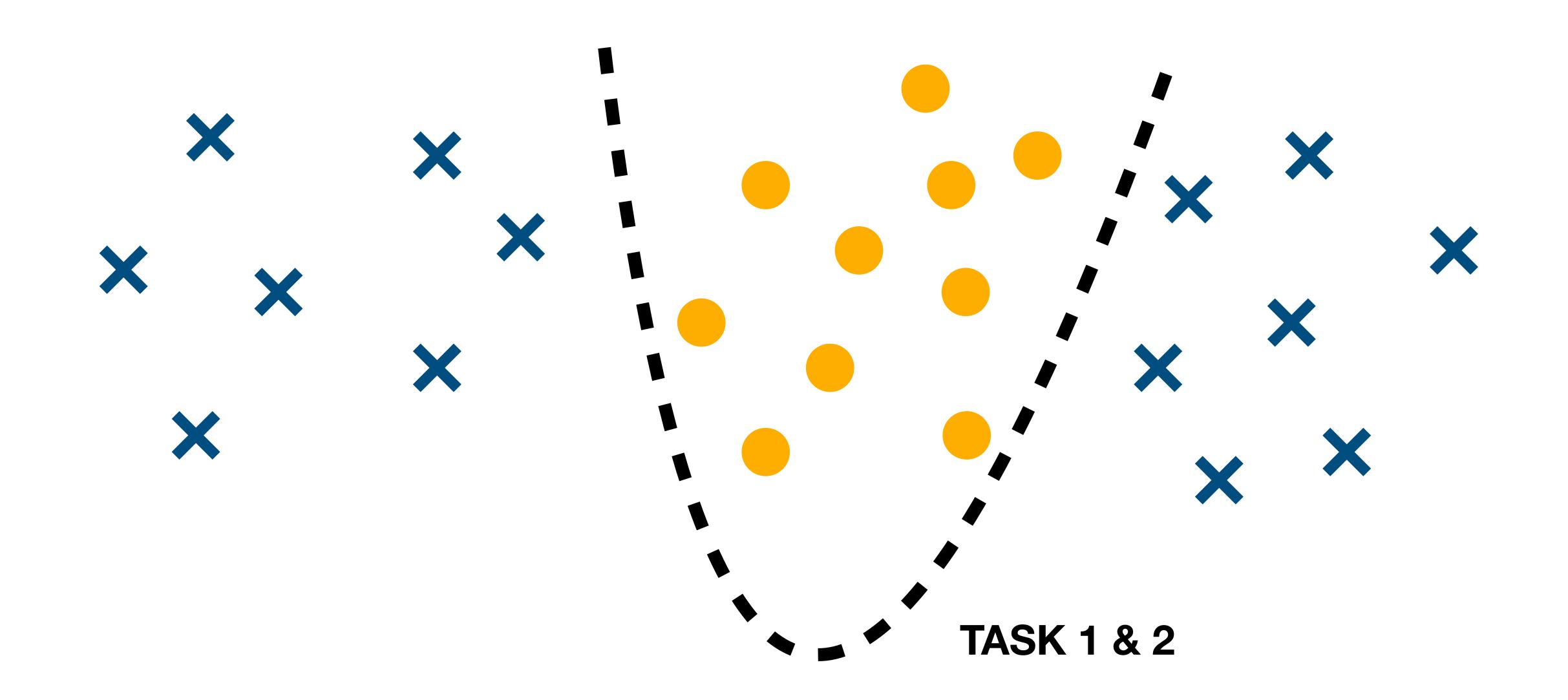
#### Even if individual tasks are easy to learn



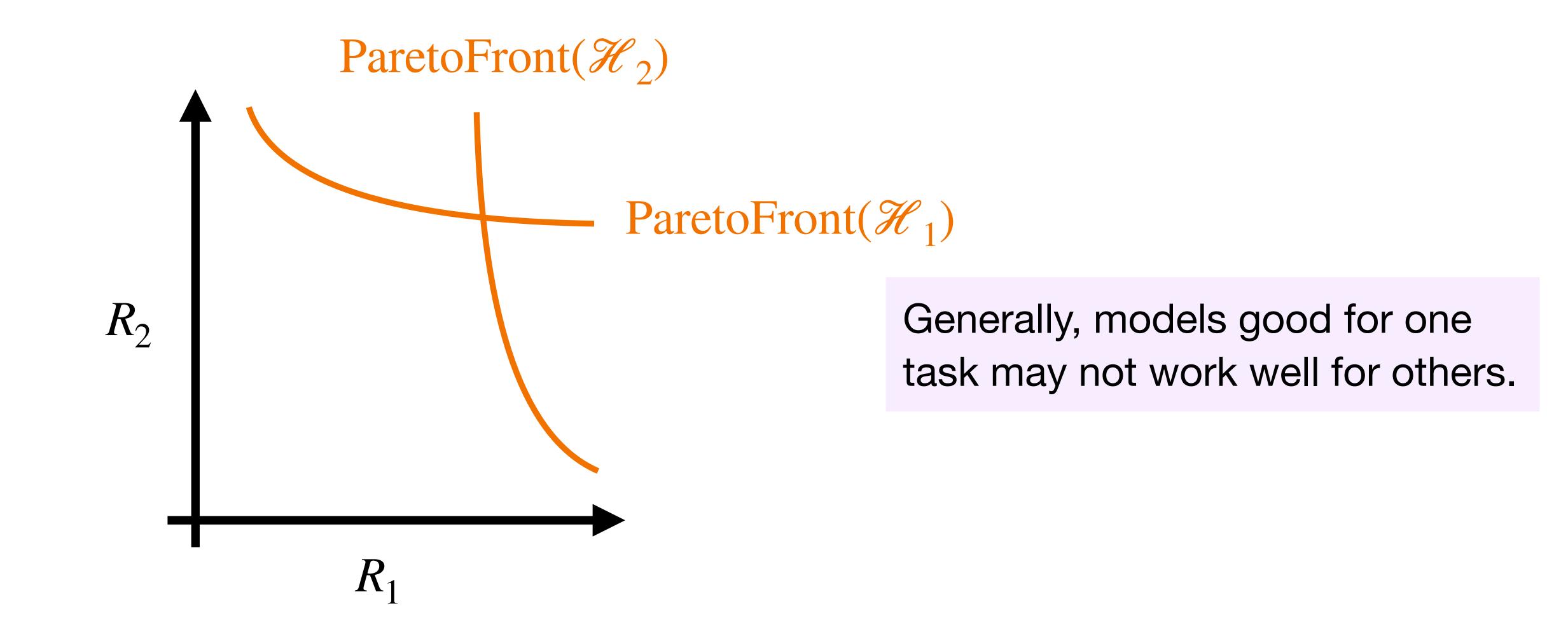
#### Even if individual tasks are easy to learn



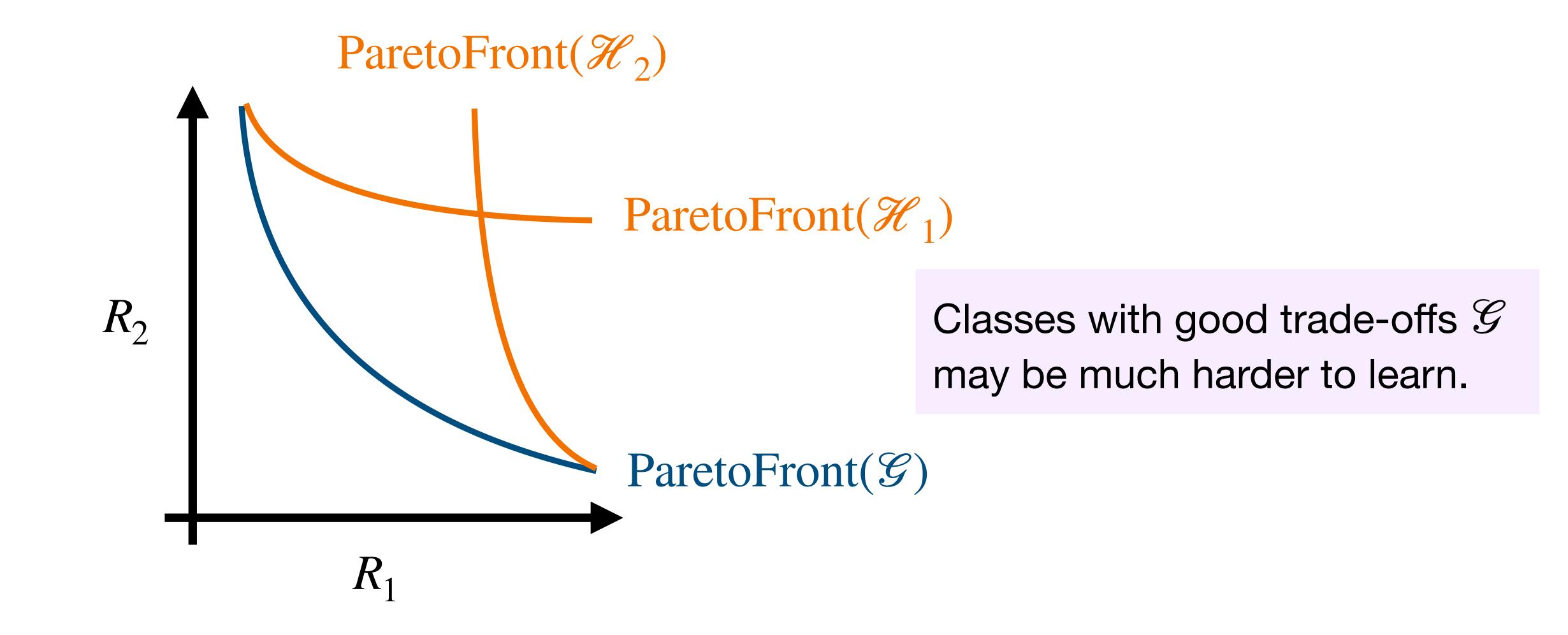
#### Joint tasks may require more complex classes

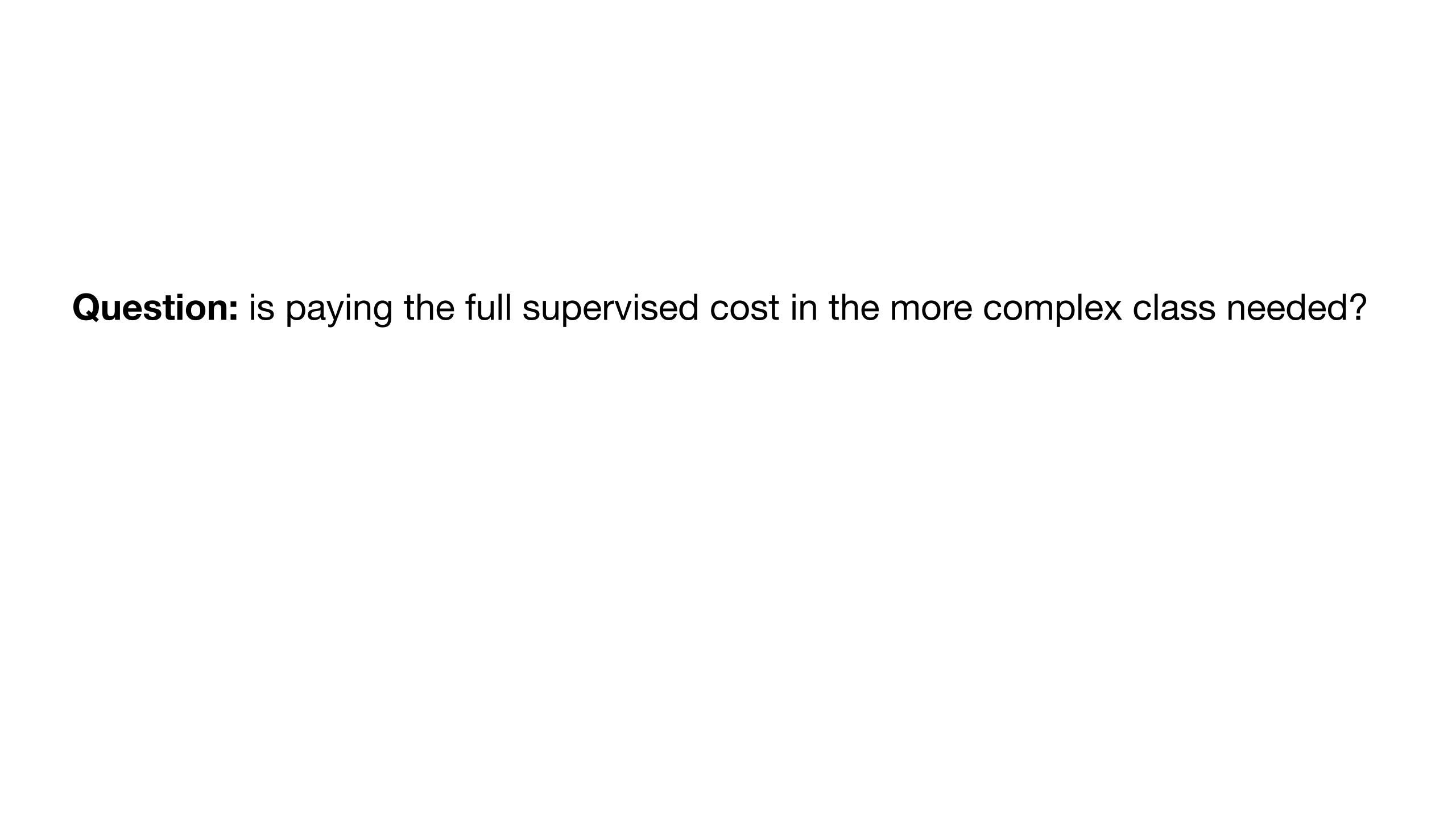


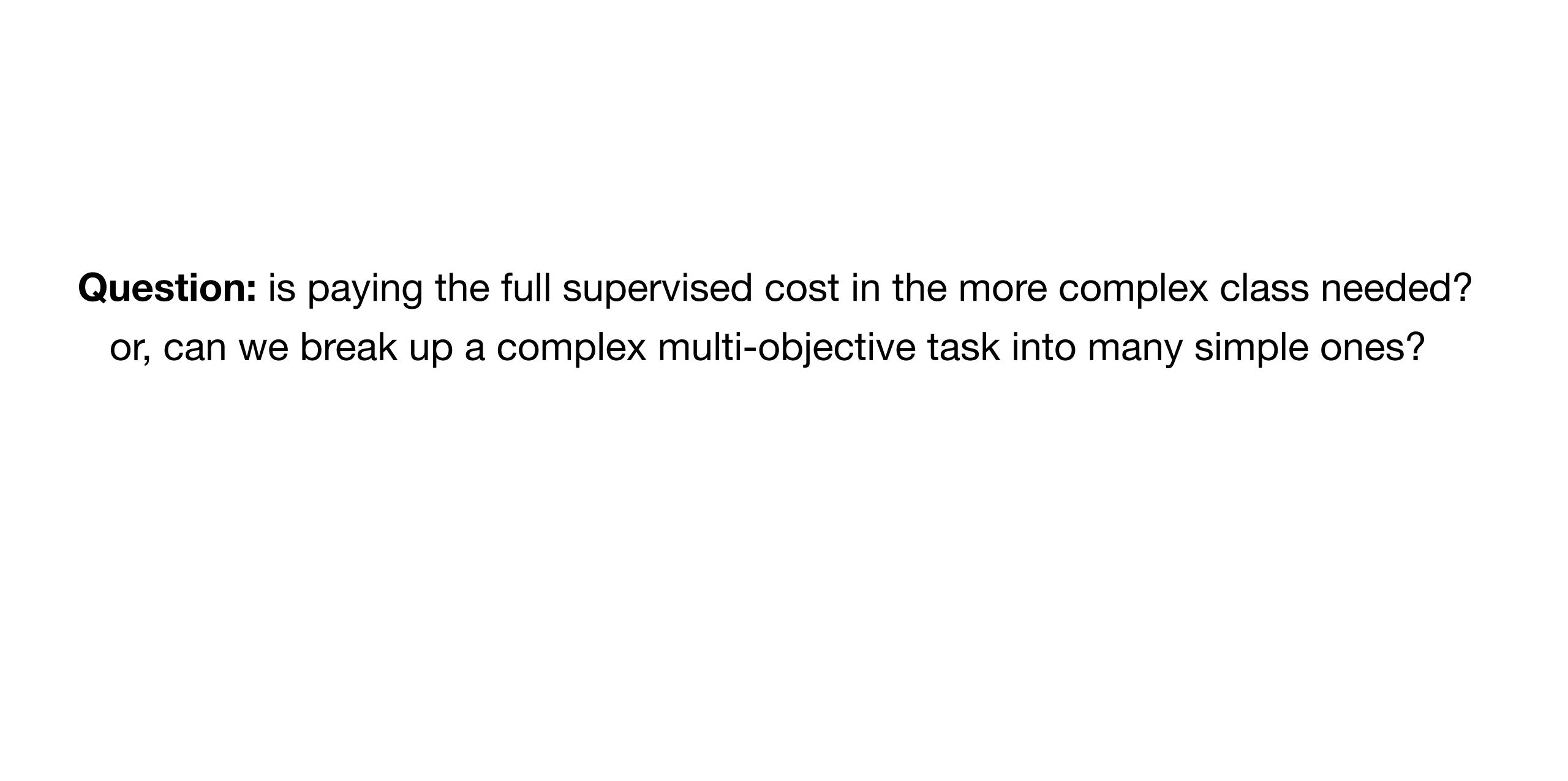
### From individual to joint model classes



### From individual to joint model classes







**Question:** is paying the full supervised cost in the more complex class needed? or, can we break up a complex multi-objective task into many simple ones?

We want the best of both worlds:

- low statistical cost to learn
- efficient trade-offs across risks

**Assumptions:** 

#### **Assumptions:**

• We have access to simple model classes  $\mathcal{H}_k$  for each task  $k \in [K]$ .

#### **Assumptions:**

- We have access to simple model classes  $\mathcal{H}_k$  for each task  $k \in [K]$ .
- These classes are guaranteed to perform well individually. In fact,

$$f_k^{\star} \in \mathcal{H}_k$$

#### **Assumptions:**

- We have access to simple model classes  $\mathcal{H}_k$  for each task  $k \in [K]$ .
- These classes are guaranteed to perform well individually. In fact,

$$f_k^{\star} \in \mathcal{H}_k$$

• We also have access to unlabeled data from  $P_X^k$ .

**Theorem.** Let  ${\mathcal G}$  be a model class. To  ${\mathcal E}$ -learn all Pareto optimal models, we need:

$$\tilde{\Theta}\left(\frac{\mathrm{VC}(\mathscr{G})\cdot K}{\varepsilon^2}\right) \text{ samples,}$$

**Theorem.** Let  $\mathscr{G}$  be a model class. To  $\varepsilon$ -learn all Pareto optimal models, we need:

$$\tilde{\Theta}\left(\frac{\mathrm{VC}(\mathscr{G})\cdot K}{\varepsilon^2}\right) \text{ samples,}$$

even if we know  $f_k^*$  and have infinite unlabeled data from  $P_k$  for each  $k \in [K]$ .

**Theorem.** Let  $\mathscr{G}$  be a model class. To  $\varepsilon$ -learn all Pareto optimal models, we need:

$$\tilde{\Theta}\left(\frac{\mathrm{VC}(\mathcal{G})\cdot K}{\varepsilon^2}\right) \text{ samples,}$$

even if we know  $f_k^*$  and have infinite unlabeled data from  $P_k$  for each  $k \in [K]$ .

**Intuition:** ability to drive fast + ability to be safe  $\Rightarrow$  ability to drive fast safely.

**Theorem.** Let  $\mathscr G$  be a model class. To  $\varepsilon$ -learn all Pareto optimal models, we need:

$$\tilde{\Theta}\left(\frac{\mathrm{VC}(\mathscr{G})\cdot K}{\varepsilon^2}\right) \text{ samples,}$$

even if we know  $f_k^*$  and have infinite unlabeled data from  $P_k$  for each  $k \in [K]$ .

**Intuition:** ability to drive fast + ability to be safe  $\Rightarrow$  ability to drive fast safely.

Implication: learning trade-offs can be arbitrarily harder than solving individual tasks.

Task 1: predict outcome of flip of Coin 1

Task 2: predict outcome of flip of Coin 2

Multi-objective task: predict the majority vote

Task 1: predict outcome of flip of Coin 1

Task 2: predict outcome of flip of Coin 2

Multi-objective task: predict the majority vote

#### **Promise:**

$$p_1 \in \left\{ \frac{1}{2}, \frac{1}{2} + \varepsilon \right\}$$
 and  $p_2 \in \left\{ \frac{1}{2} - \varepsilon, \frac{1}{2} \right\}$ .

**Task 1:** predict outcome of flip of Coin 1 (optimal prediction = H)

Task 2: predict outcome of flip of Coin 2

Multi-objective task: predict the majority vote

#### **Promise:**

$$p_1 \in \left\{ \frac{1}{2}, \frac{1}{2} + \varepsilon \right\}$$
 and  $p_2 \in \left\{ \frac{1}{2} - \varepsilon, \frac{1}{2} \right\}$ .

**Task 1:** predict outcome of flip of Coin 1 (optimal prediction = H)

**Task 2:** predict outcome of flip of Coin 2 (optimal prediction = T)

Multi-objective task: predict the majority vote

#### **Promise:**

$$p_1 \in \left\{ \frac{1}{2}, \frac{1}{2} + \varepsilon \right\}$$
 and  $p_2 \in \left\{ \frac{1}{2} - \varepsilon, \frac{1}{2} \right\}$ .

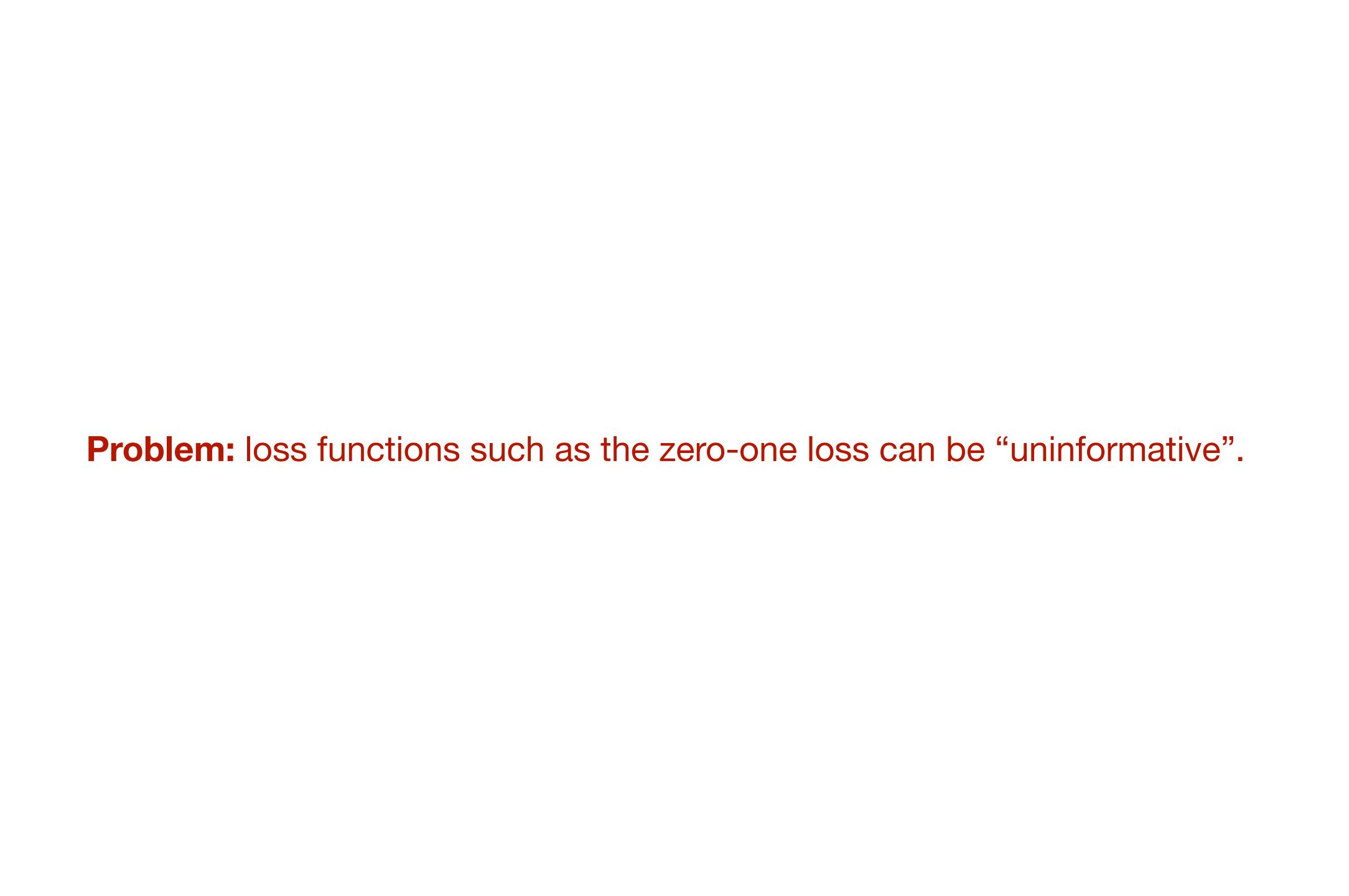
**Task 1:** predict outcome of flip of Coin 1 (optimal prediction = H)

**Task 2:** predict outcome of flip of Coin 2 (optimal prediction = T)

Multi-objective task: predict the majority vote (requires learning bias of coins)

#### **Promise:**

$$p_1 \in \left\{ \frac{1}{2}, \frac{1}{2} + \varepsilon \right\}$$
 and  $p_2 \in \left\{ \frac{1}{2} - \varepsilon, \frac{1}{2} \right\}$ .



### II. Multi-objective learning with nice losses

#### **Assumptions:**

- We have access to simple model classes  $\mathcal{H}_k$  for each task  $k \in [K]$ .
- These classes are guaranteed to perform well individually. In fact,

$$f_k^{\star} \in \mathcal{H}_k$$

- We also have access to unlabeled data from  $P_X^k$ .
- We have nice losses  $\mathcal{C}_k$ .

### Bregman losses

There are standard losses that are much nicer than the zero-one loss:

- Square loss
- Logistic loss/cross-entropy loss
- etc.

### Bregman losses

There are standard losses that are much nicer than the zero-one loss:

- Square loss
- Logistic loss/cross-entropy loss
- etc.

#### Risk Decomposition for Bregman Losses

$$R(h) = \mathbb{E}\left[\ell(Y, f^{\star}(X))\right] + \mathbb{E}\left[\ell(f^{\star}(X), h(X))\right]$$

### Bregman losses

There are standard losses that are much nicer than the zero-one loss:

- Square loss
- Logistic loss/cross-entropy loss
- etc.

#### Risk Decomposition for Bregman Losses

$$R(h) = \mathbb{E}\left[\mathcal{E}(Y, f^{\star}(X))\right] + \mathbb{E}\left[\mathcal{E}(f^{\star}(X), h(X))\right]$$

This expectation is only over unlabeled data  ${\cal P}_{\cal X}$ 

### Risk functional estimation

#### Risk Decomposition for Bregman Losses

$$R(h) = \mathbb{E}\left[\mathcal{E}(Y, f^{\star}(X))\right] + \mathbb{E}\left[\mathcal{E}(f^{\star}(X), h(X))\right]$$

 $\Longrightarrow$  knowledge of the risk minimizer  $f^*$  tells us a lot about the risk R.

### Risk functional estimation

#### Risk Decomposition for Bregman Losses

$$R(h) = \mathbb{E}\left[\ell(Y, f^{\star}(X))\right] + \mathbb{E}\left[\ell(f^{\star}(X), h(X))\right]$$

#### A plug-in estimator

$$\hat{f} \leftarrow \widehat{\arg\min} \, \hat{\mathbb{E}} \left[ \ell(Y, h(X)) \right]$$

$$h \in \mathcal{H}$$

### Risk functional estimation

#### Risk Decomposition for Bregman Losses

$$R(h) = \mathbb{E}\left[\mathcal{E}(Y, f^{\star}(X))\right] + \mathbb{E}\left[\mathcal{E}(f^{\star}(X), h(X))\right]$$

#### A plug-in estimator

$$\hat{f} \leftarrow \widehat{\arg\min} \, \hat{\mathbb{E}} \left[ \ell(Y, h(X)) \right]$$

$$h \in \mathcal{H}$$

• 
$$\hat{R}(h) = \hat{\mathbb{E}}\left[\mathcal{E}(Y,\hat{f}(X))\right] + \hat{\mathbb{E}}\left[\mathcal{E}(\hat{f}(X),h(X))\right]$$

### For each task $k \in [K]$ :

• Use labeled data to approximately recover  $f_k^\star$  from  $\mathcal{H}_k$ .

### For each task $k \in [K]$ :

- Use labeled data to approximately recover  $f_k^{\star}$  from  $\mathcal{H}_k$ .
- Generate a lot of pseudo-label data  $(X,\hat{f}_k(X))$  to estimate  $\hat{R}_k$ .

### For each task $k \in [K]$ :

- Use labeled data to approximately recover  $f_k^{\star}$  from  $\mathcal{H}_k$ .
- Generate a lot of pseudo-label data  $(X,\hat{f}_k(X))$  to estimate  $\hat{R}_k$ .

This only requires additional unlabeled data from  $P_X^k$ .

### For each task $k \in [K]$ :

- Use labeled data to approximately recover  $f_k^{\star}$  from  $\mathcal{H}_k$ .
- Generate a lot of pseudo-label data  $(X,\hat{f}_k(X))$  to estimate  $\hat{R}_k$ .

For each  $s \in \mathcal{S}$  (parametrizing the Pareto set):

### For each task $k \in [K]$ :

- Use labeled data to approximately recover  $f_k^{\star}$  from  $\mathcal{H}_k$ .
- Generate a lot of pseudo-label data  $(X,\hat{f}_k(X))$  to estimate  $\hat{R}_k$ .

### For each $s \in \mathcal{S}$ (parametrizing the Pareto set):

Solve the optimization problem:

$$\hat{g}_s \leftarrow \arg\min_{g \in \mathcal{G}} s(\hat{\mathbf{R}}(g)).$$

**Theorem.** Let  $\mathcal G$  be a joint model class,

**Theorem.** Let  $\mathcal{G}$  be a joint model class, and let  $\mathcal{H}_k \ni f_k^*$  be model class that contains the Bayes-optimal model.

**Theorem.** Let  $\mathcal{G}$  be a joint model class, and let  $\mathcal{H}_k \ni f_k^*$  be model class that contains the Bayes-optimal model. If the losses  $\ell_k$  are Bregman losses:

**Theorem.** Let  $\mathscr{G}$  be a joint model class, and let  $\mathscr{H}_k \ni f_k^*$  be model class that contains the Bayes-optimal model. If the losses  $\mathscr{C}_k$  are Bregman losses:

$$O\left(rac{\sum_k \mathrm{VC}(\mathscr{H}_k)}{arepsilon^4}
ight)$$
 labeled samples,  $O\left(rac{\mathrm{VC}(\mathscr{G})}{arepsilon^2}
ight)$  unlabeled samples

are enough to  $\varepsilon$ -learn all Pareto-optimal models.

**Theorem.** Let  $\mathscr{G}$  be a joint model class, and let  $\mathscr{H}_k \ni f_k^*$  be model class that contains the Bayes-optimal model. If the losses  $\mathscr{C}_k$  are Bregman losses:

$$O\left(\frac{\sum_k \text{VC}(\mathcal{H}_k)}{\varepsilon^4}\right) \text{ labeled samples, } O\left(\frac{\text{VC}(\mathcal{G})}{\varepsilon^2}\right) \text{ unlabeled samples}$$

are enough to  $\varepsilon$ -learn all Pareto-optimal models.

**Importantly,** the label sample complexity does not the complexity of the joint class  $\mathcal{G}$  in which the good trade-offs are possible.

**Theorem.** Let  $\mathscr{G}$  be a joint model class, and let  $\mathscr{H}_k \ni f_k^*$  be model class that contains the Bayes-optimal model. If the losses  $\mathscr{C}_k$  are Bregman losses:

$$O\left(\frac{\sum_k \text{VC}(\mathcal{H}_k)}{\varepsilon^4}\right) \text{ labeled samples, } O\left(\frac{\text{VC}(\mathcal{G})}{\varepsilon^2}\right) \text{ unlabeled samples}$$

are enough to  $\varepsilon$ -learn all Pareto-optimal models.

**Importantly,** the label sample complexity does not the complexity of the joint class  $\mathcal{G}$  in which the good trade-offs are possible.

Tighter, data-dependent bounds are also possible. See paper/poster.

### Takeaways

- Multi-objective learning (MOL) problems are ubiquitous in practice, but subtle.
- Learning good trade-offs can be much harder than solving the individual tasks.
- Structure in loss/feedback important for efficient multi-objective generalization.

#### Thanks!

On the sample complexity of semi-supervised multi-objective learning (NeurIPS 2025)

https://arxiv.org/abs/2508.17152