

# Robust Optimization and Machine Learning under Distribution Shift

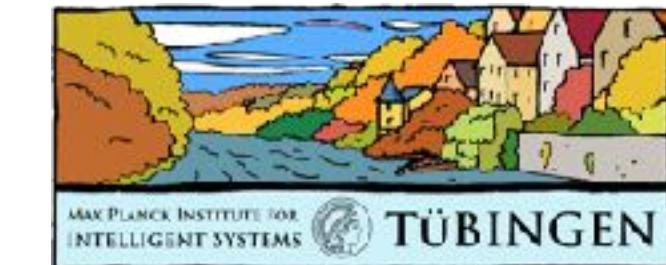
J.J. (Jia-Jie) Zhu

[jj-zhu.github.io](https://jj-zhu.github.io)

WIAS, Berlin & MPI-IS, Tübingen



Weierstraß-Institut für  
Angewandte Analysis und Stochastik



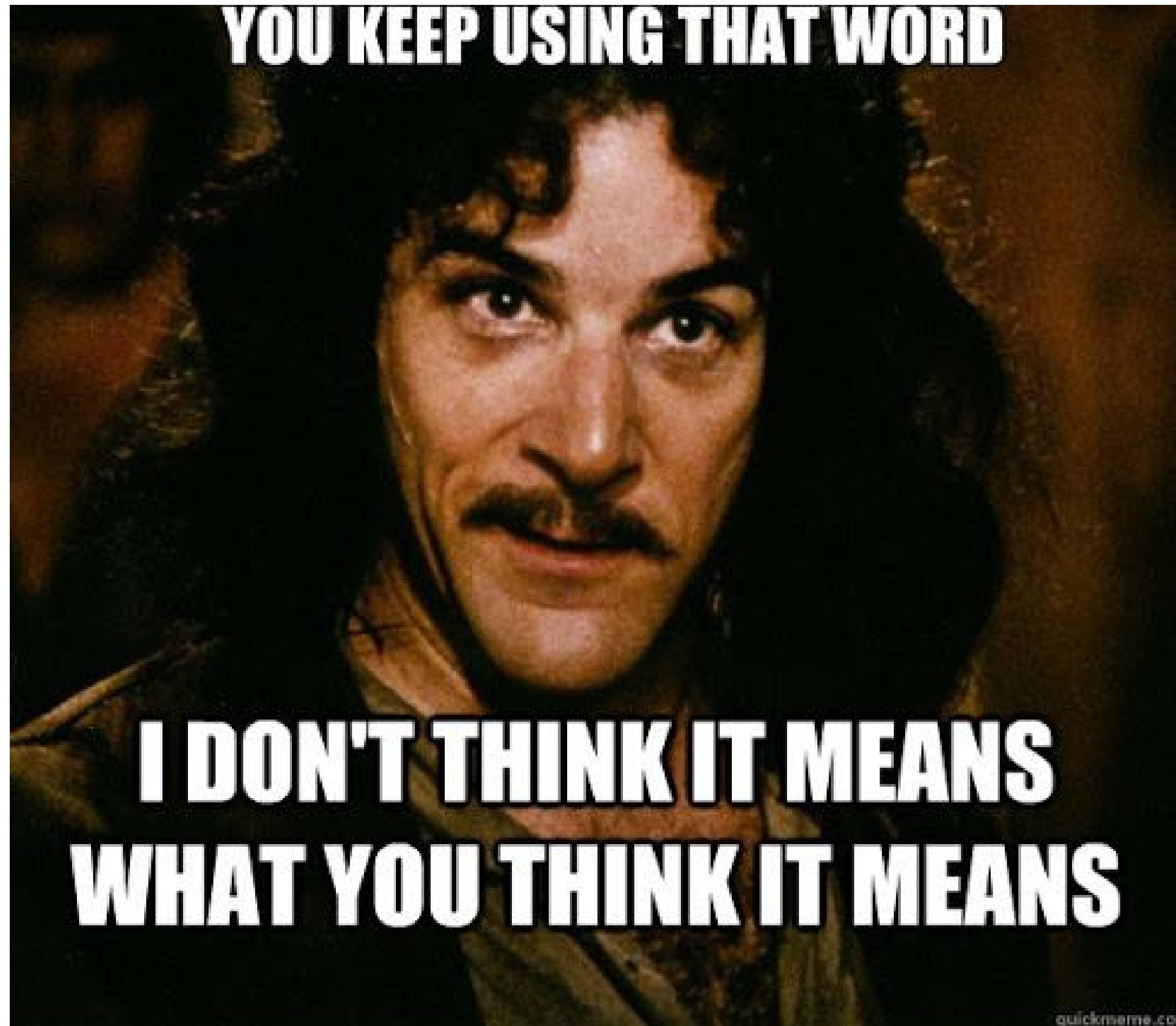
Leibniz MMS Summer School, Dagstuhl  
August 26, 2021

# Distributional Robustness

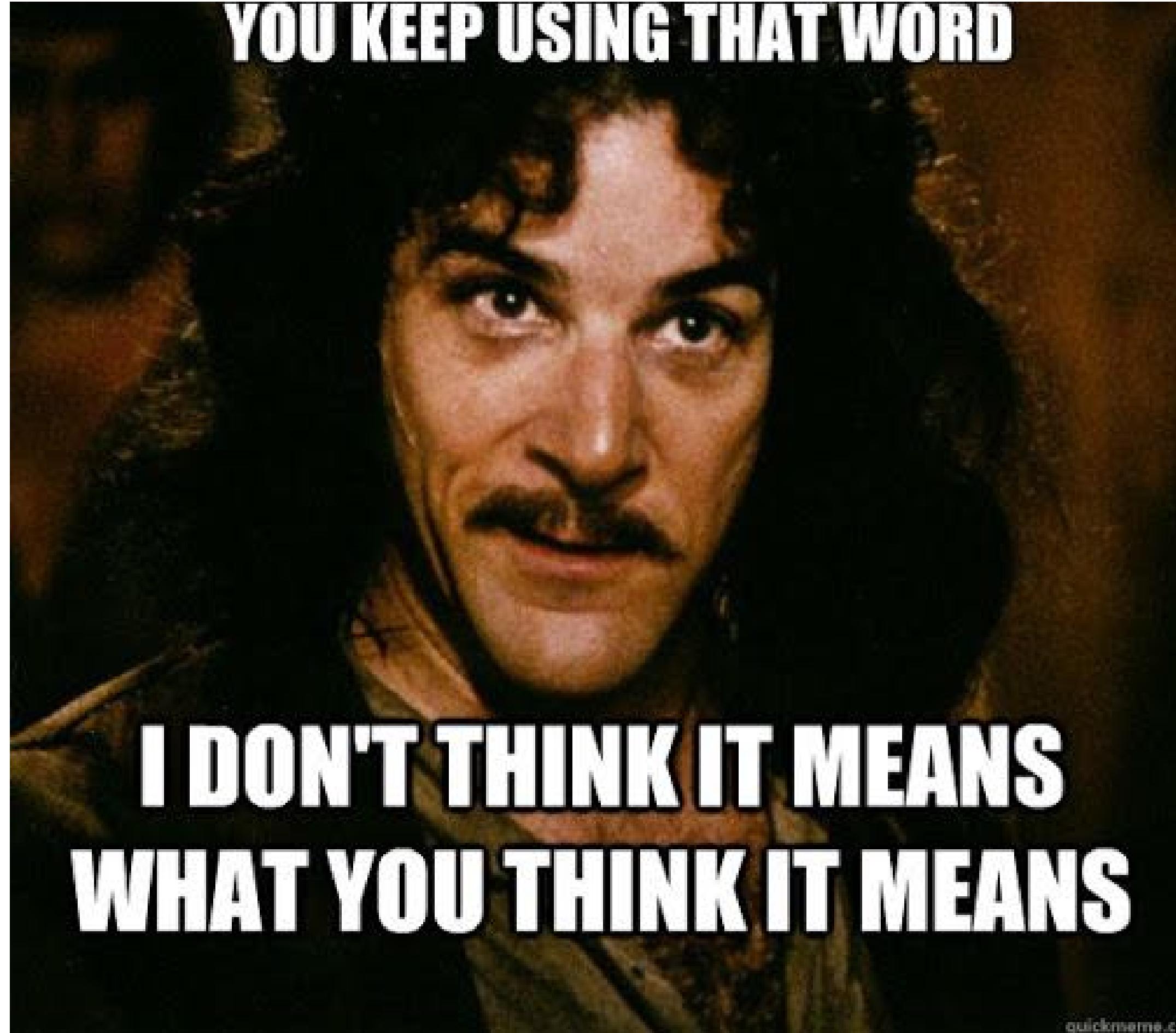
# Distributional Robustness

# **What is robustness?**

# What is robustness?

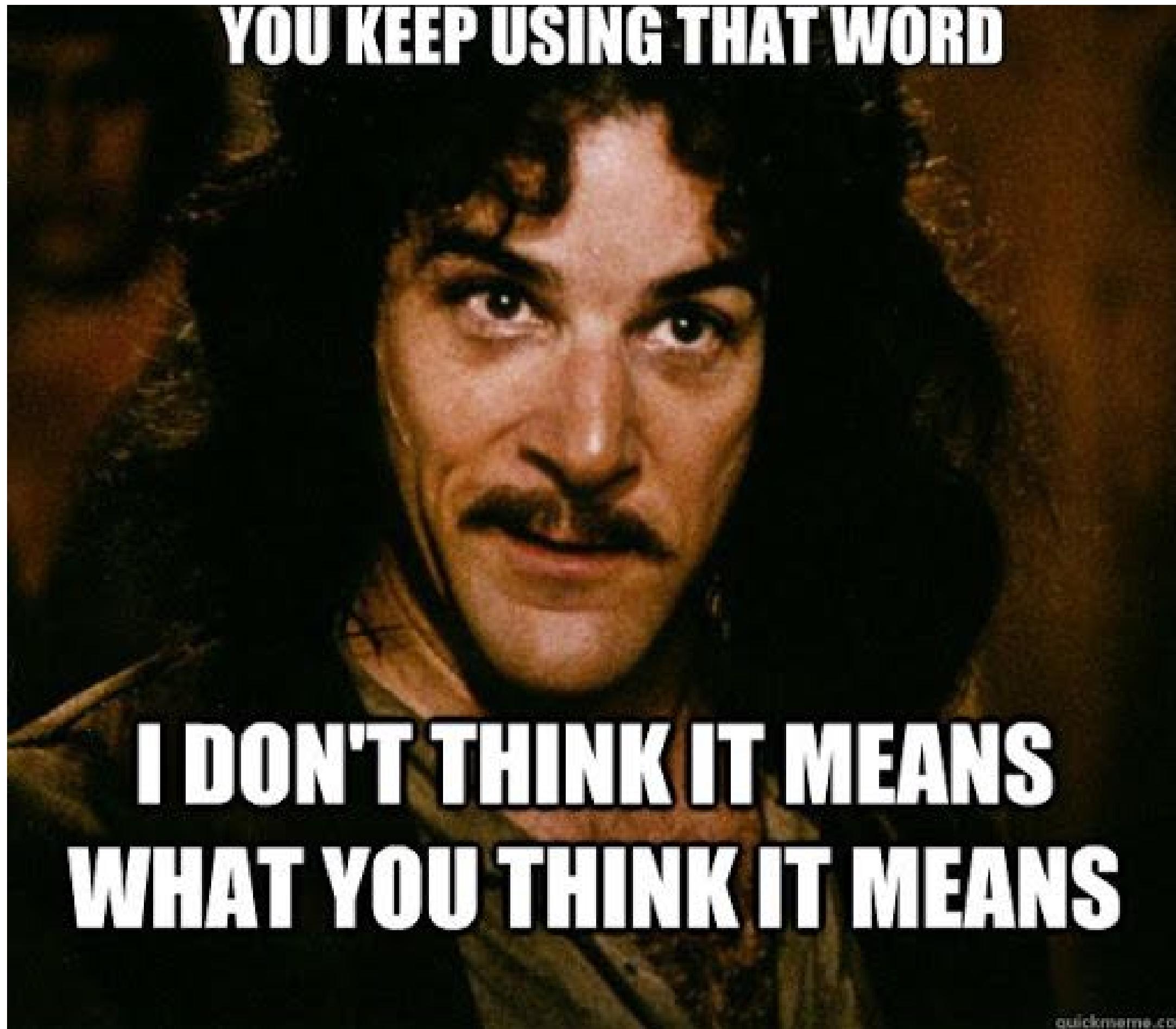


# What is robustness?



- Many fields: ...robust statistics, robust control, robust optimization, adversarial robustness, robust learning...

# What is robustness?

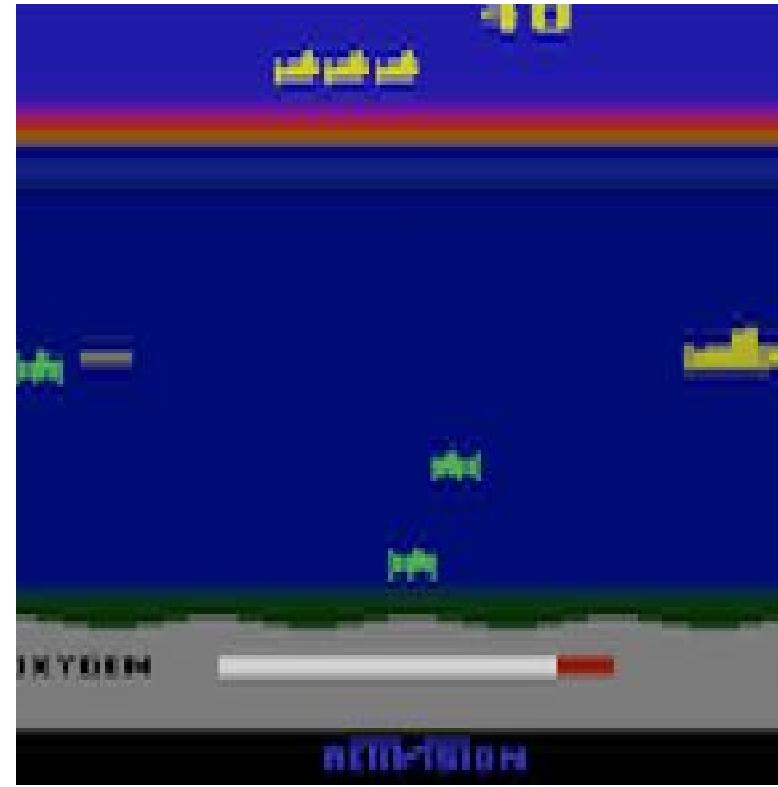


- Many fields: ...robust statistics, robust control, robust optimization, adversarial robustness, robust learning...
- Robustness is a principle for decision-making

# Optimization under uncertainty

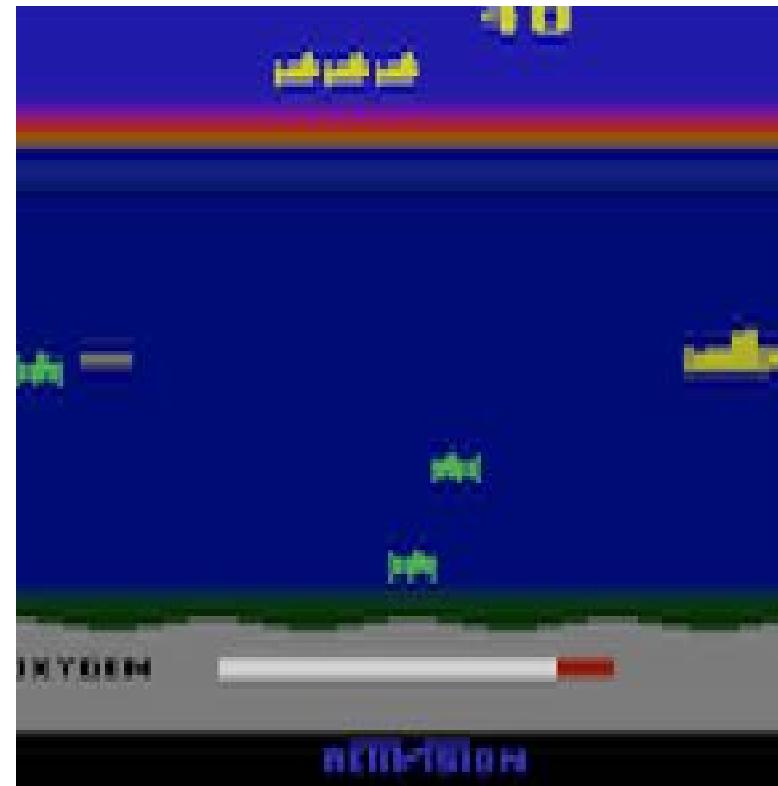
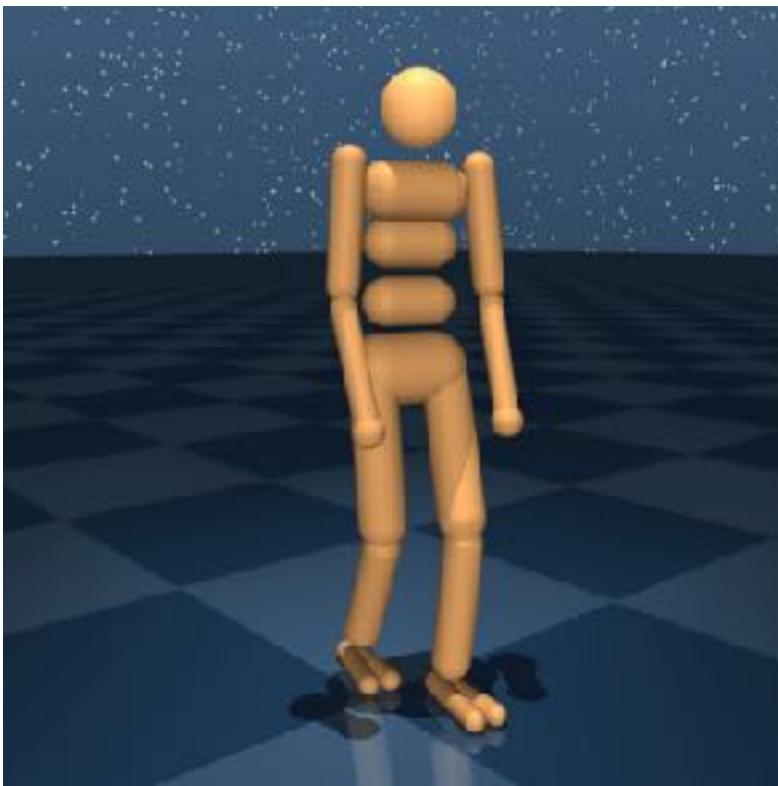
# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))



# Optimization under uncertainty

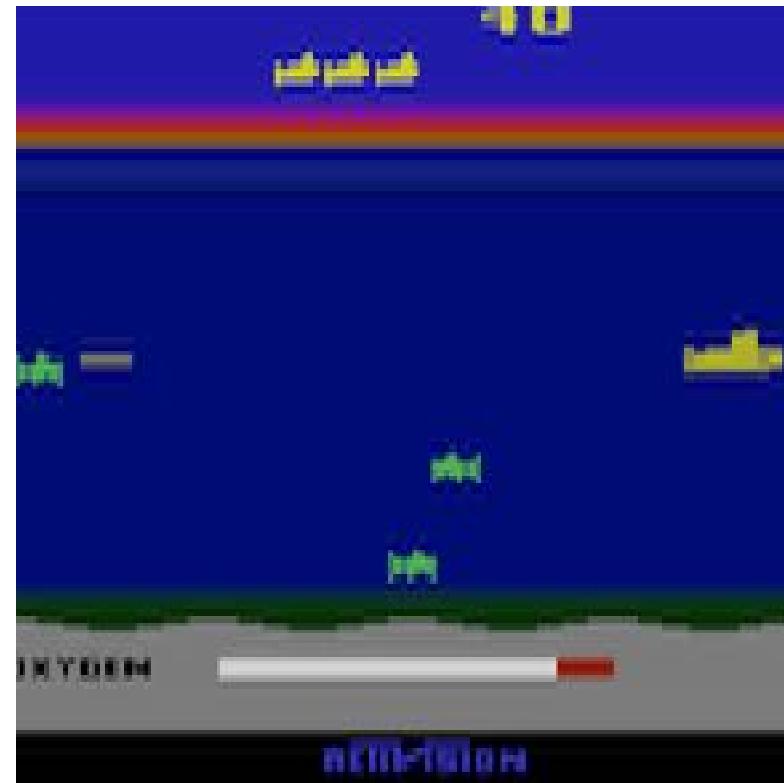
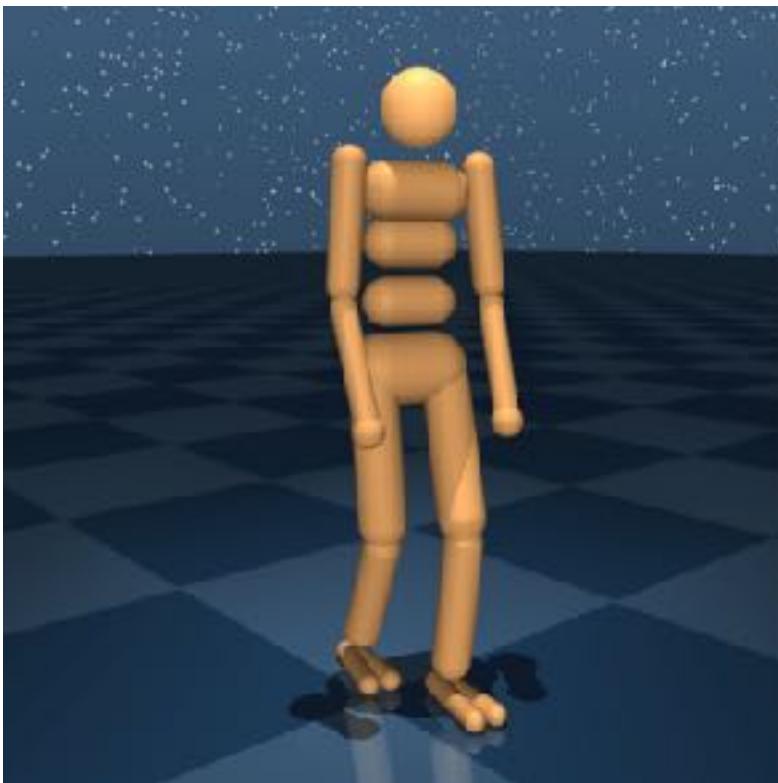
Empirical risk minimization (ERM)  
(sample average approximation (SAA))



$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))



$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

loss/cost

# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))

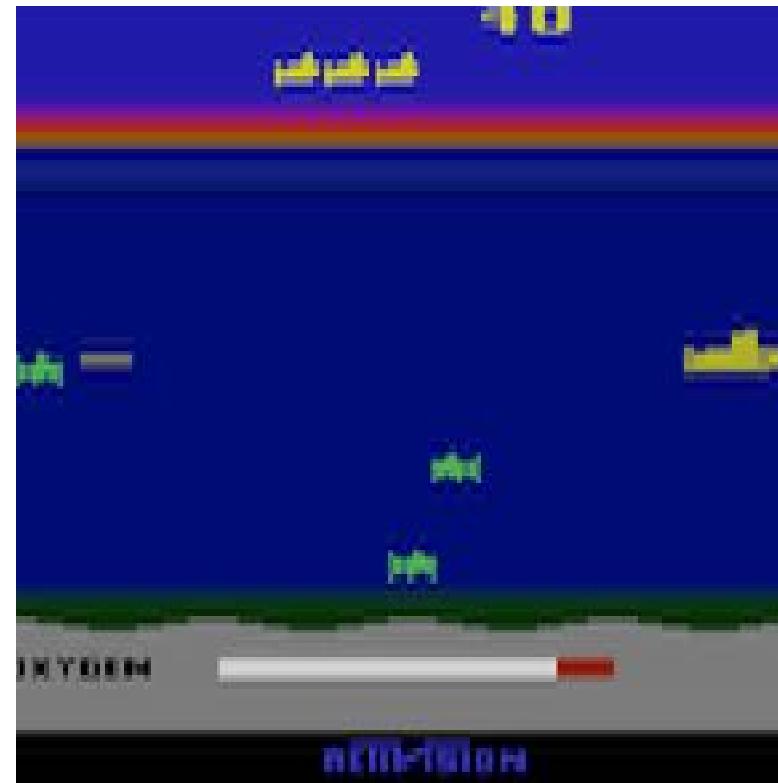
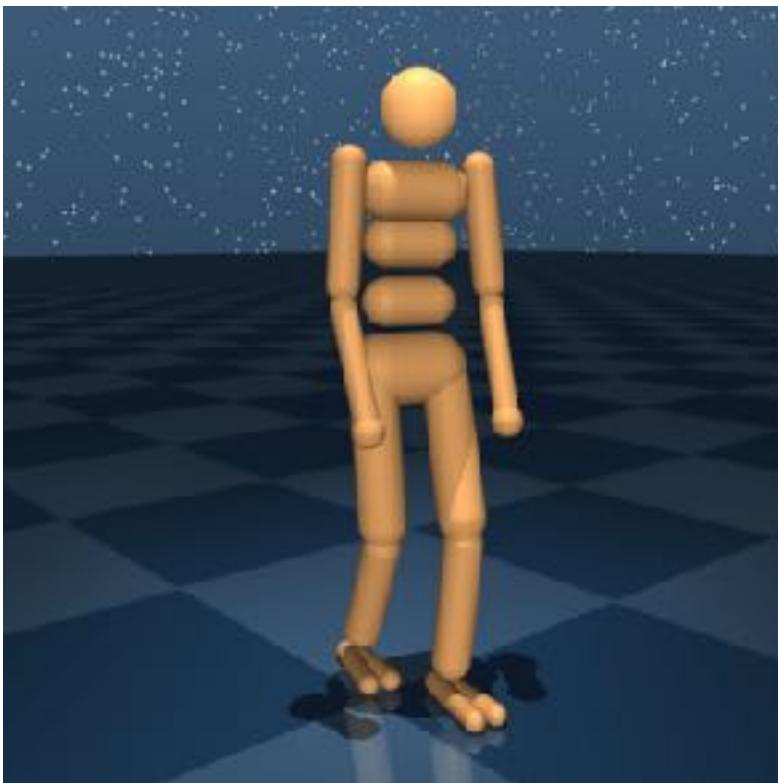


$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

loss/cost →  $\mathbb{E}$  ← uncertain

# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))

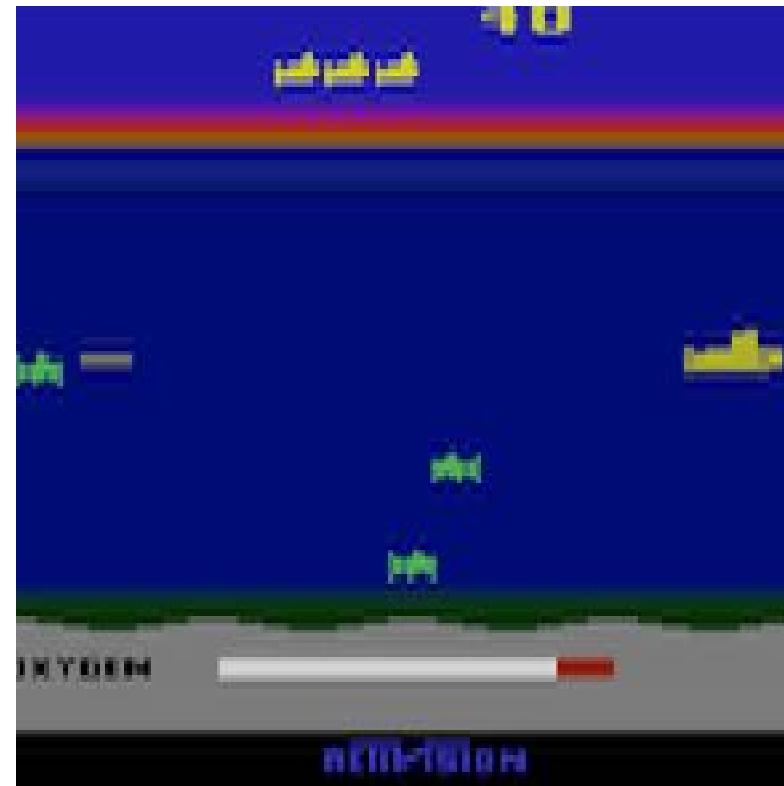


$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

loss/cost →  $\hat{P}$  ← uncertain  
Empirical dist.  $\hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$

# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))



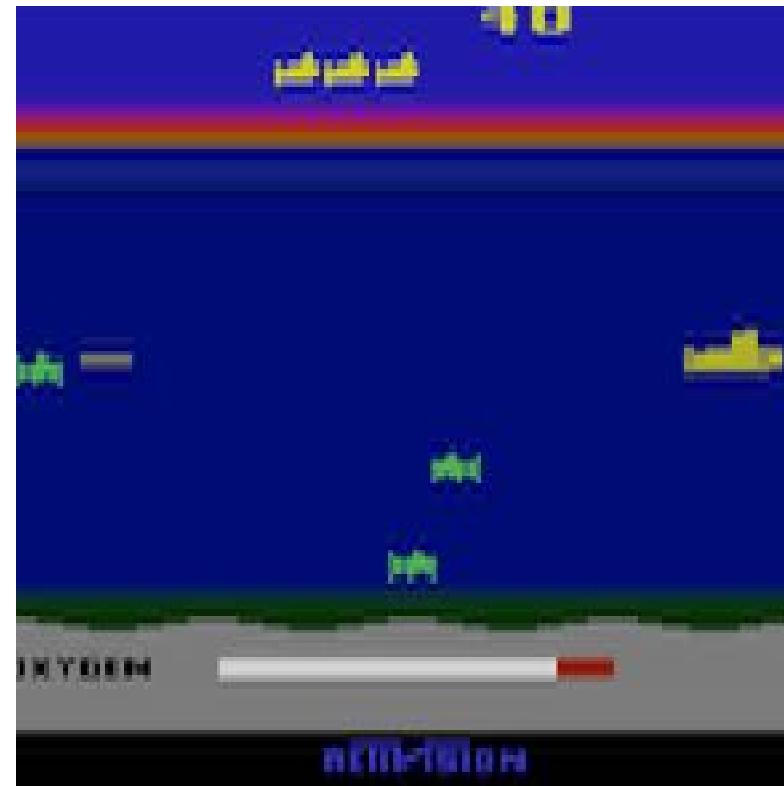
$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

loss/cost →  $\mathbb{E}$  → uncertain  
→  $\xi \sim \hat{P}$  → Empirical dist.  $\hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$

- Do well on **average**

# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))



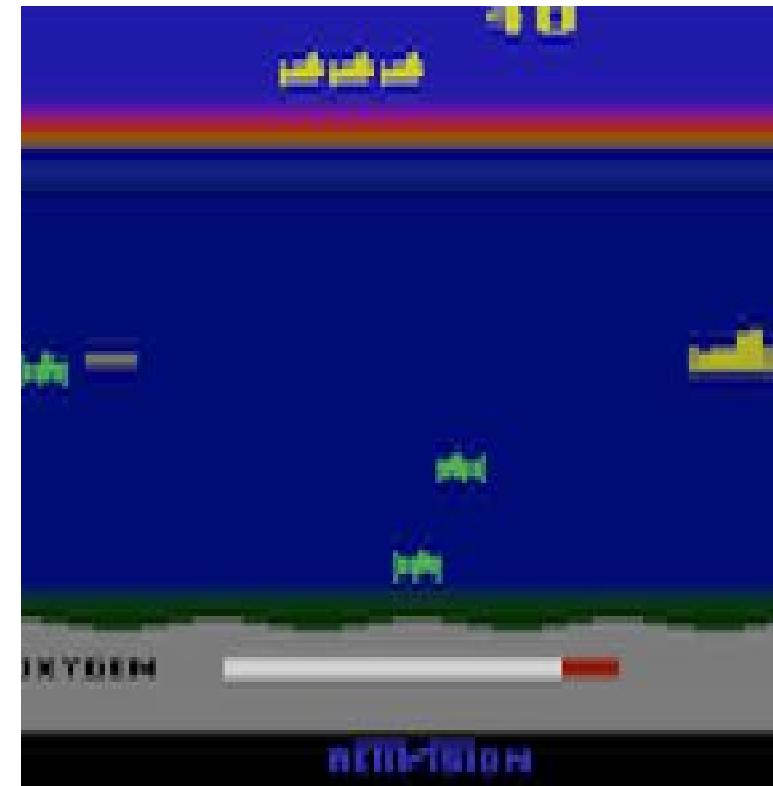
$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

loss/cost                          uncertain  
    Empirical dist.  $\hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$

- Do well on **average**
- Strength: high-performance (**optimal**)

# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))



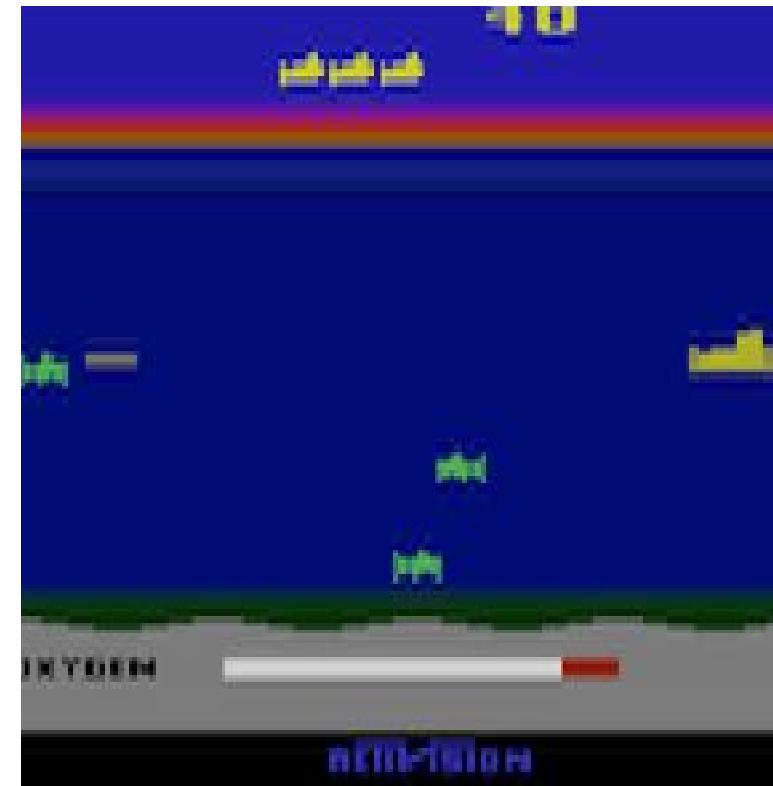
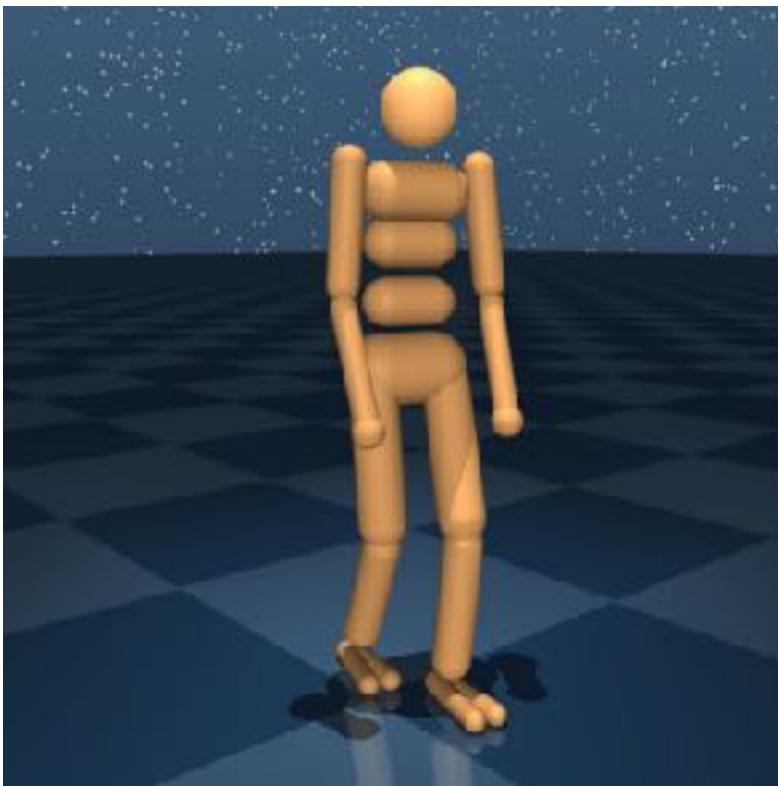
$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

loss/cost →  $\mathbb{E}$  → uncertain  
→  $\xi \sim \hat{P}$  → Empirical dist.  $\hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$

- Do well on **average**
- Strength: high-performance (**optimal**)
- Weakness: **fragile** – adversarial attacks, sim2real, off-policy/offline RL

# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))

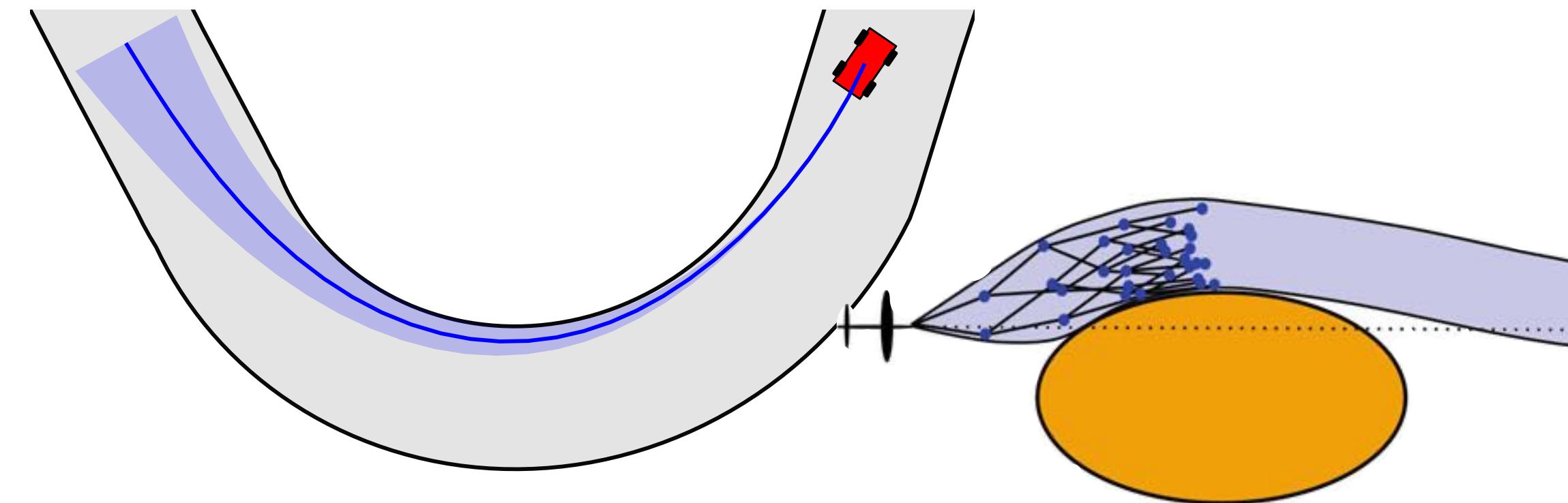


$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

$$\text{Empirical dist. } \hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$$

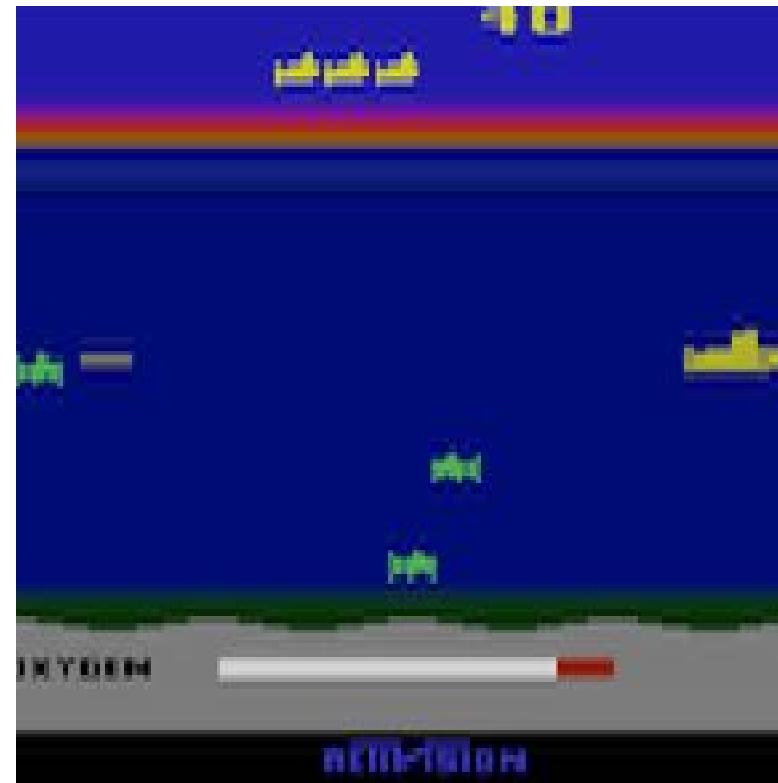
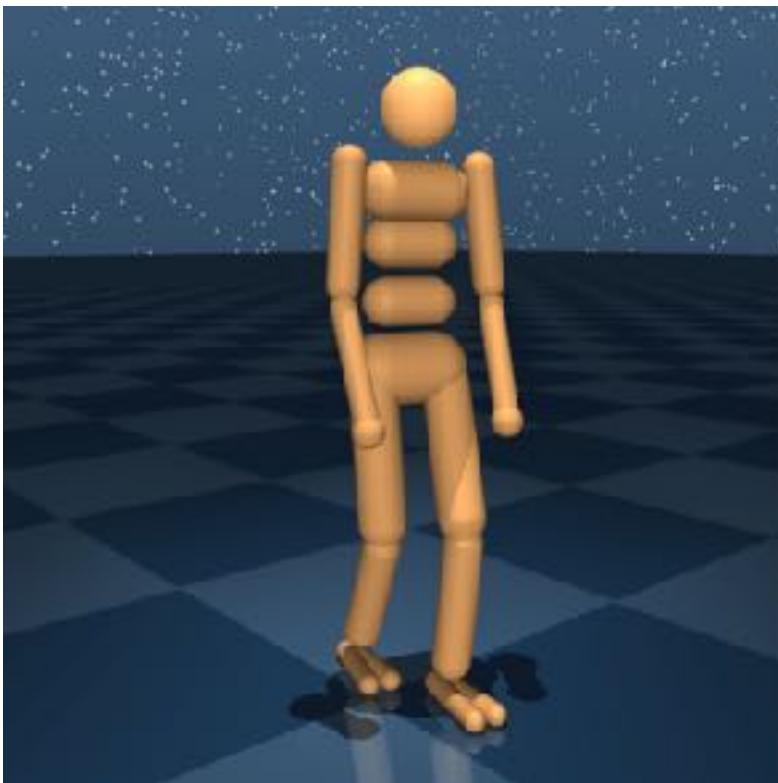
- Do well on **average**
- Strength: high-performance (**optimal**)
- Weakness: **fragile** – adversarial attacks, sim2real, off-policy/offline RL

Robust optimization (RO)  
(robust control, games)



# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))

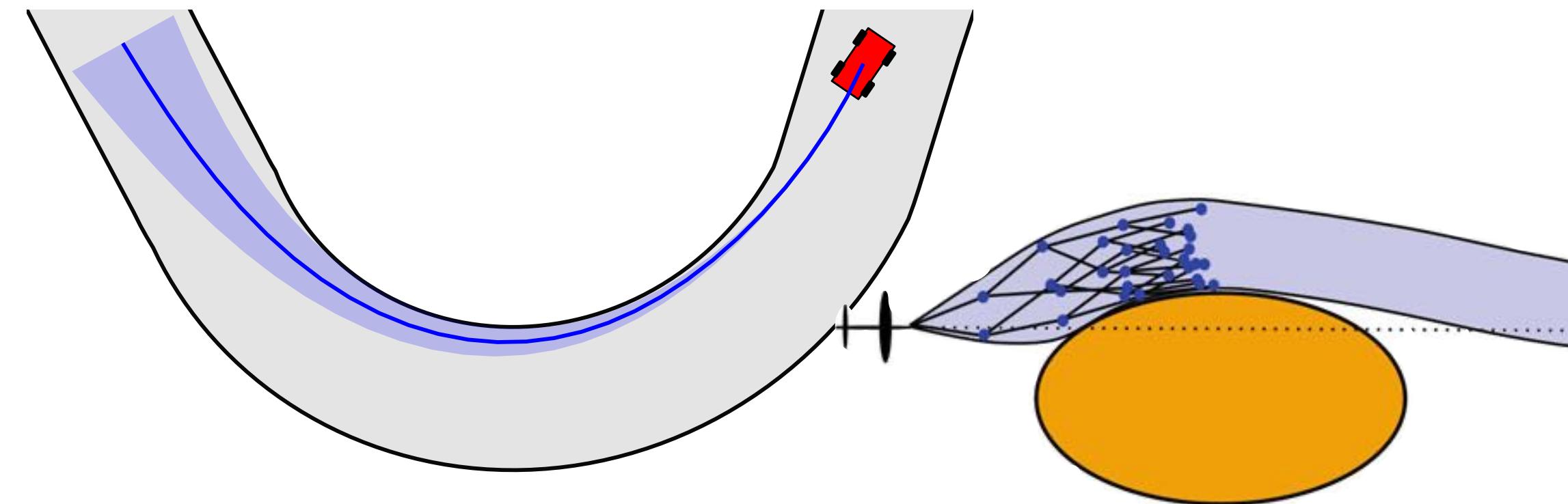


$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

$$\text{Empirical dist. } \hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$$

- Do well on **average**
- Strength: high-performance (**optimal**)
- Weakness: **fragile** – adversarial attacks, sim2real, off-policy/offline RL

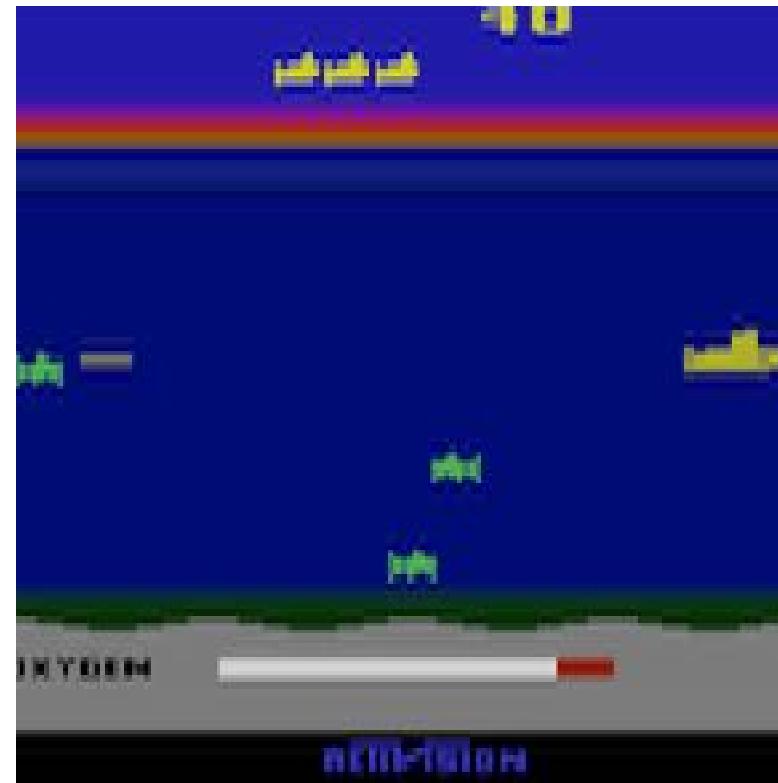
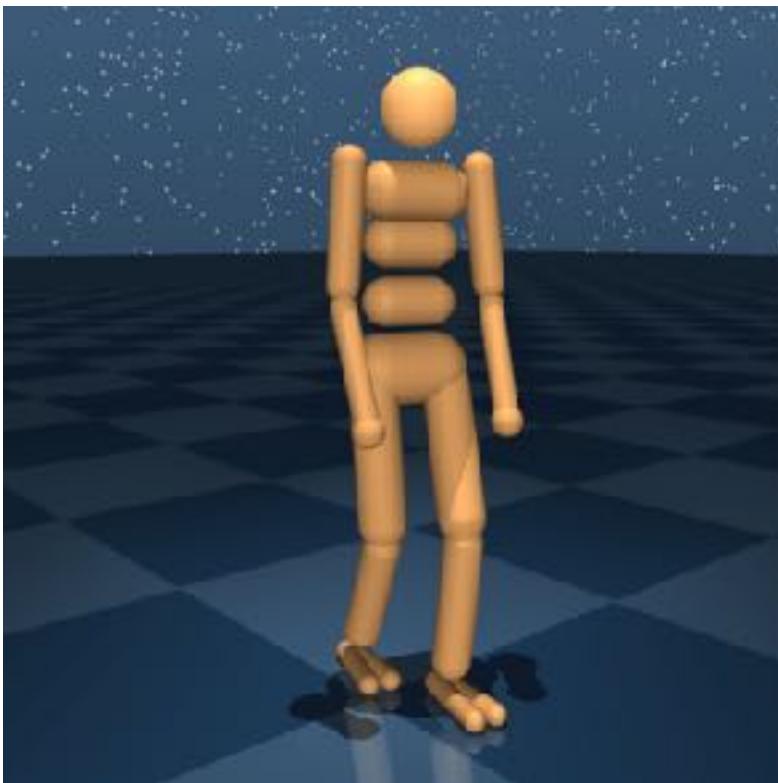
Robust optimization (RO)  
(robust control, games)



$$\min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$

# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))

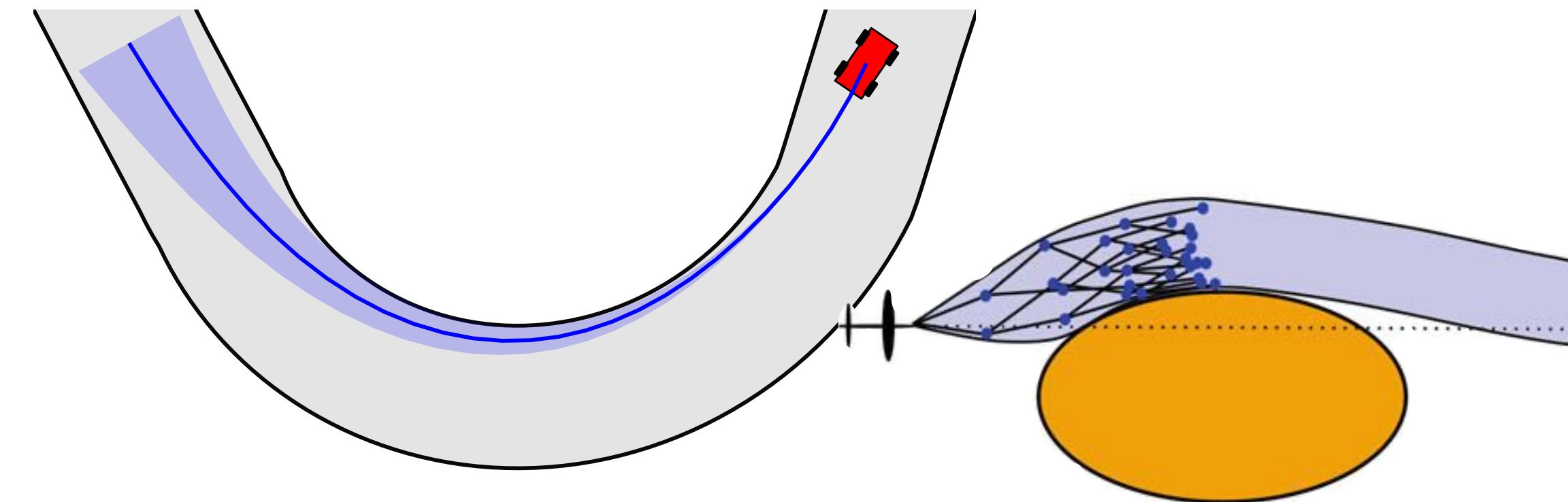


$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

← Empirical dist.  $\hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$

- Do well on **average**
- Strength: high-performance (**optimal**)
- Weakness: **fragile** – adversarial attacks, sim2real, off-policy/offline RL

Robust optimization (RO)  
(robust control, games)

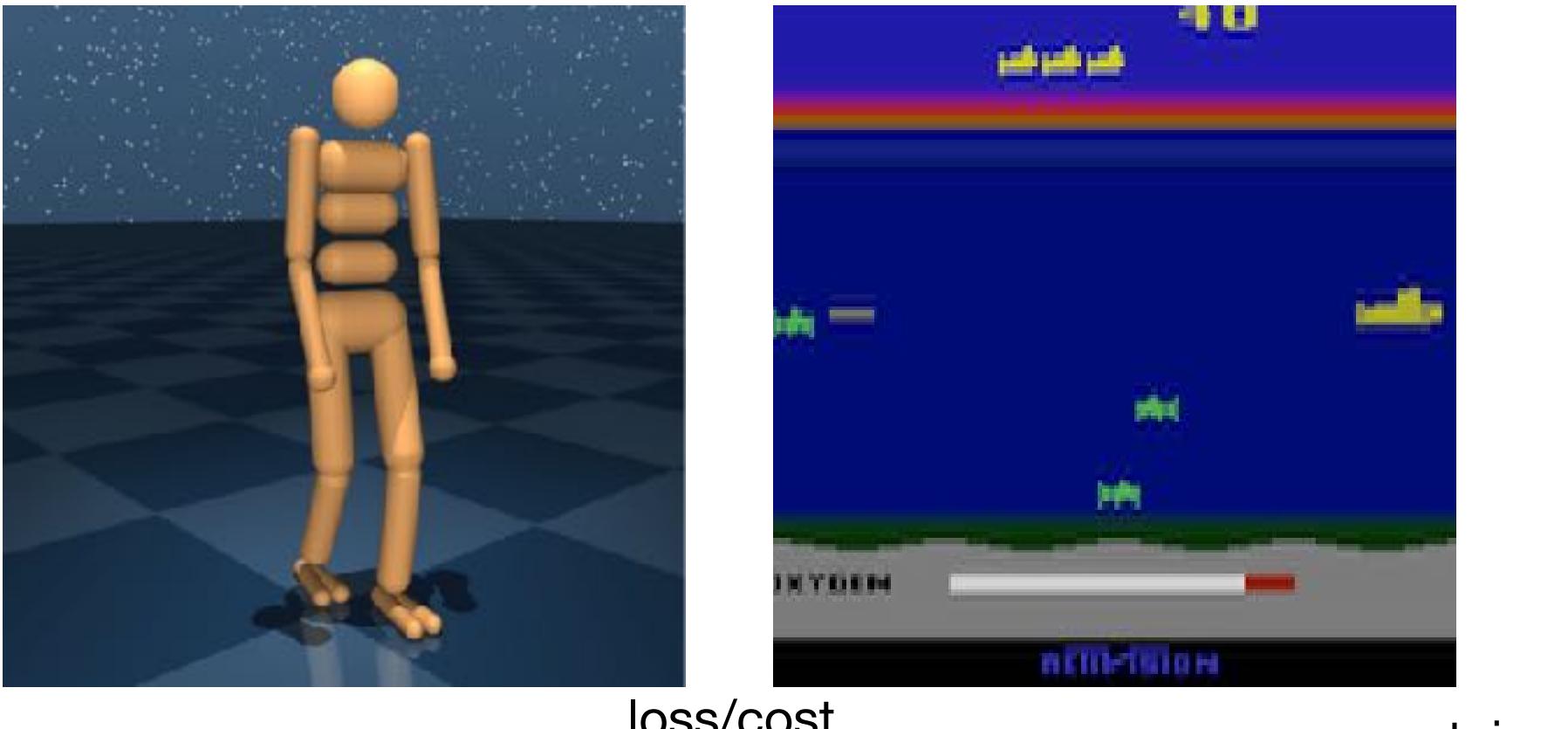


$$\min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$

- Do well in the **worst case**

# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))

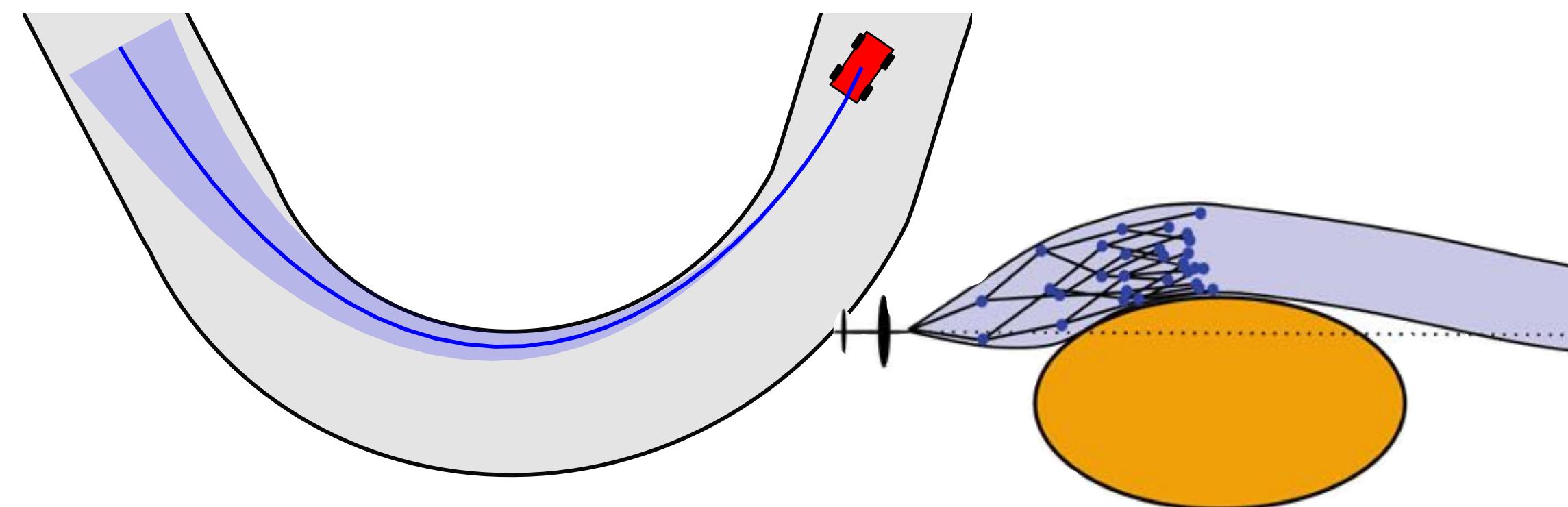


$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

$$\text{Empirical dist. } \hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$$

- Do well on **average**
- Strength: high-performance (**optimal**)
- Weakness: **fragile** – adversarial attacks, sim2real, off-policy/offline RL

Robust optimization (RO)  
(robust control, games)

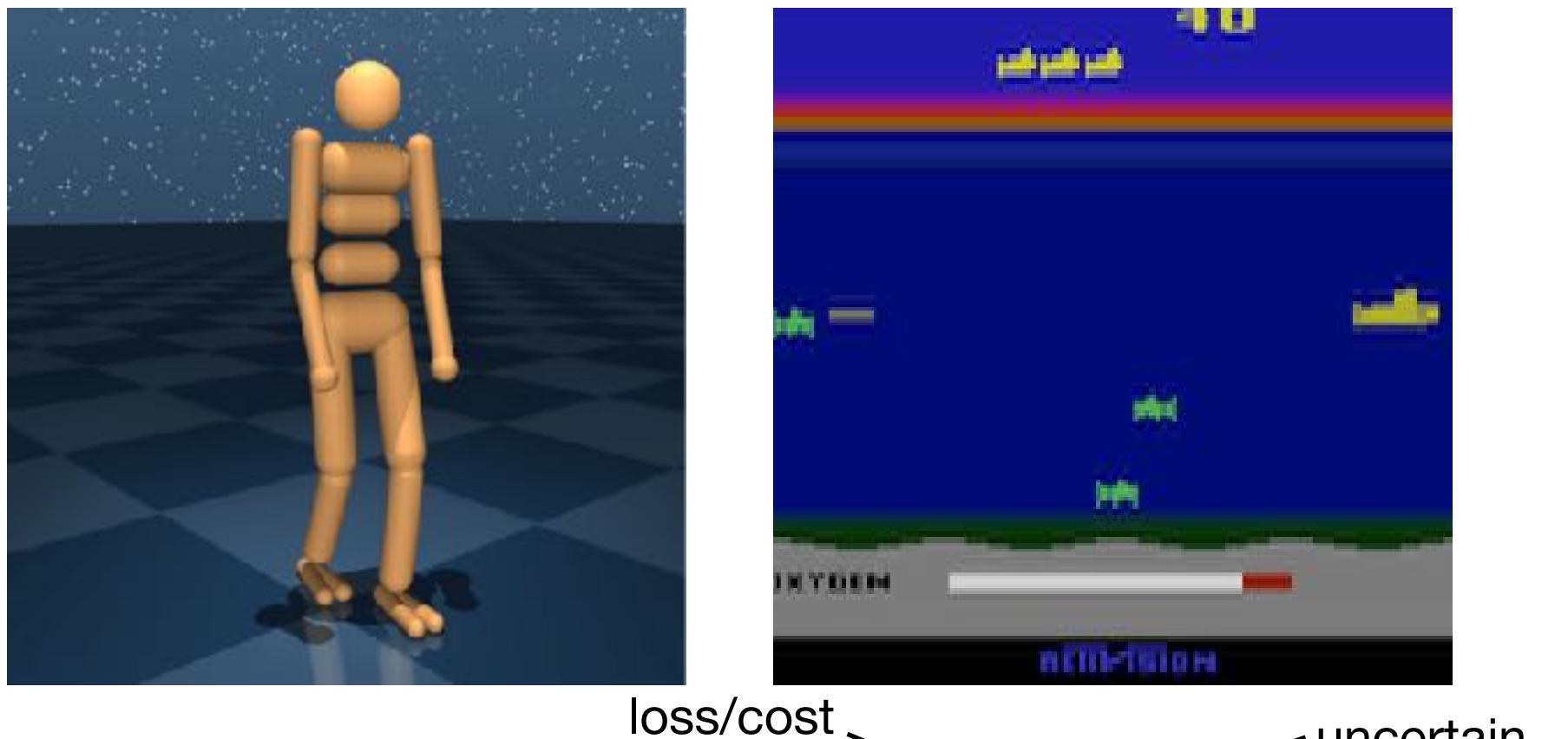


$$\min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$

- Do well in the **worst case**
- Strength: **robustness**

# Optimization under uncertainty

Empirical risk minimization (ERM)  
(sample average approximation (SAA))

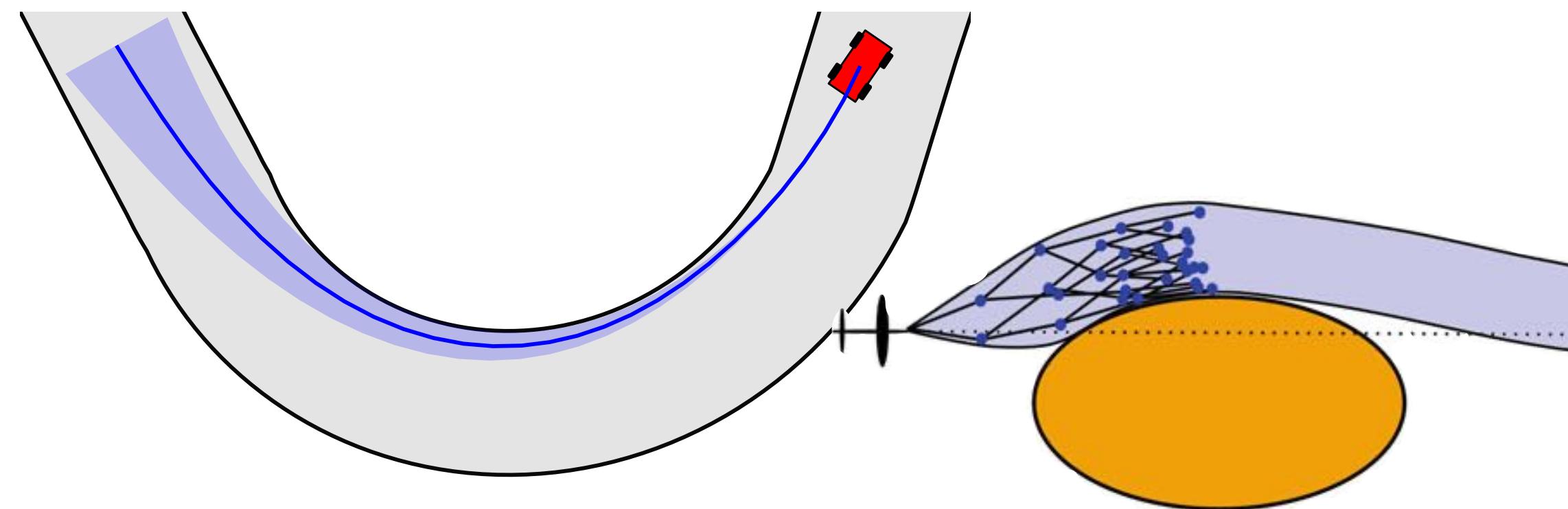


$$\min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

$$\text{Empirical dist. } \hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$$

- Do well on **average**
- Strength: **high-performance (optimal)**
- Weakness: **fragile** – adversarial attacks, sim2real, off-policy/offline RL

Robust optimization (RO)  
(robust control, games)



- Do well in the **worst case**
- Strength: **robustness**
- Weakness: **conservative** – worst case doesn't often happen

# Various robustness formulation in optimization

# Various robustness formulation in optimization

- Nominal constraint:  $c(z) \leq 0$

# Various robustness formulation in optimization

- Nominal constraint:  $c(z) \leq 0$
- Robust constraint:  $c(z, e) \leq 0, \forall e \in \mathcal{E}$

# Various robustness formulation in optimization

- Nominal constraint:  $c(z) \leq 0$
- Robust constraint:  $c(z, e) \leq 0, \forall e \in \mathcal{E}$
- Chance constraint:  $P(c(z, e) \leq 0) \geq 0.95$

# Various robustness formulation in optimization

- Nominal constraint:  $c(z) \leq 0$
- Robust constraint:  $c(z, e) \leq 0, \forall e \in \mathcal{E}$
- Chance constraint:  $P(c(z, e) \leq 0) \geq 0.95$
- Scenario optimization:  $c(z, e_i) \leq 0$ , for  $i = 1, 2, \dots, N$

# **Worst-case robust optimization (RO)**

# *Worst-case robust optimization (RO)*

$$\min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$

# *Worst-case robust optimization (RO)*

$$\min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$

- That looks robust, but how do I solve it?

# *Worst-case robust optimization (RO)*

$$\min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$

- That looks robust, but how do I solve it?

# Worst-case robust optimization (RO)

$$\min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$

- That looks robust, but how do I solve it?
- A simple case

$$(P) \min_{\theta} \sup_{\|p - \hat{p}\|_a \leq \epsilon} p^T \theta$$

# Worst-case robust optimization (RO)

$$\min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$

- That looks robust, but how do I solve it?
- A simple case

$$(P) \min_{\theta} \sup_{\|p - \hat{p}\|_a \leq \epsilon} p^T \theta$$

- But this formulation actually contains infinitely many optimization problems!

# Worst-case robust optimization (RO)

$$\min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$

- That looks robust, but how do I solve it?
- A simple case

$$(P) \min_{\theta} \sup_{\|p - \hat{p}\|_a \leq \epsilon} p^T \theta$$

- But this formulation actually contains infinitely many optimization problems!
- We solve this RO using the (strong) *duality of convex optimization*:  
*The primal (worst-case) RO problem (P) is equivalent to the dual problem*

# Worst-case robust optimization (RO)

$$\min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$

- That looks robust, but how do I solve it?
- A simple case

$$(P) \min_{\theta} \sup_{\|p - \hat{p}\|_a \leq \epsilon} p^T \theta$$

- But this formulation actually contains infinitely many optimization problems!
- We solve this RO using the (strong) *duality of convex optimization*:  
*The primal (worst-case) RO problem (P) is equivalent to the dual problem*

$$(D) \min_{\theta} \hat{p}^T \theta + \epsilon \|\theta\|_2$$

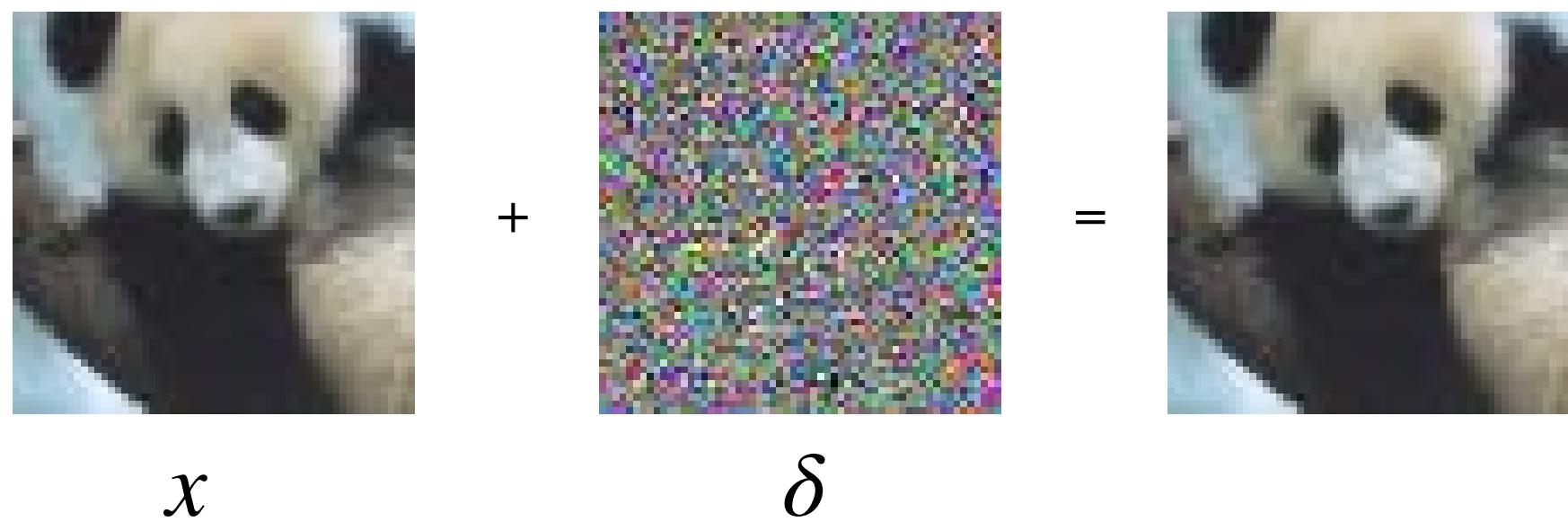
i.e.,  $(P) = (D)$ .

# **Application of worst-case RO: adversarial robustness of DNN**

# Application of worst-case RO: adversarial robustness of DNN

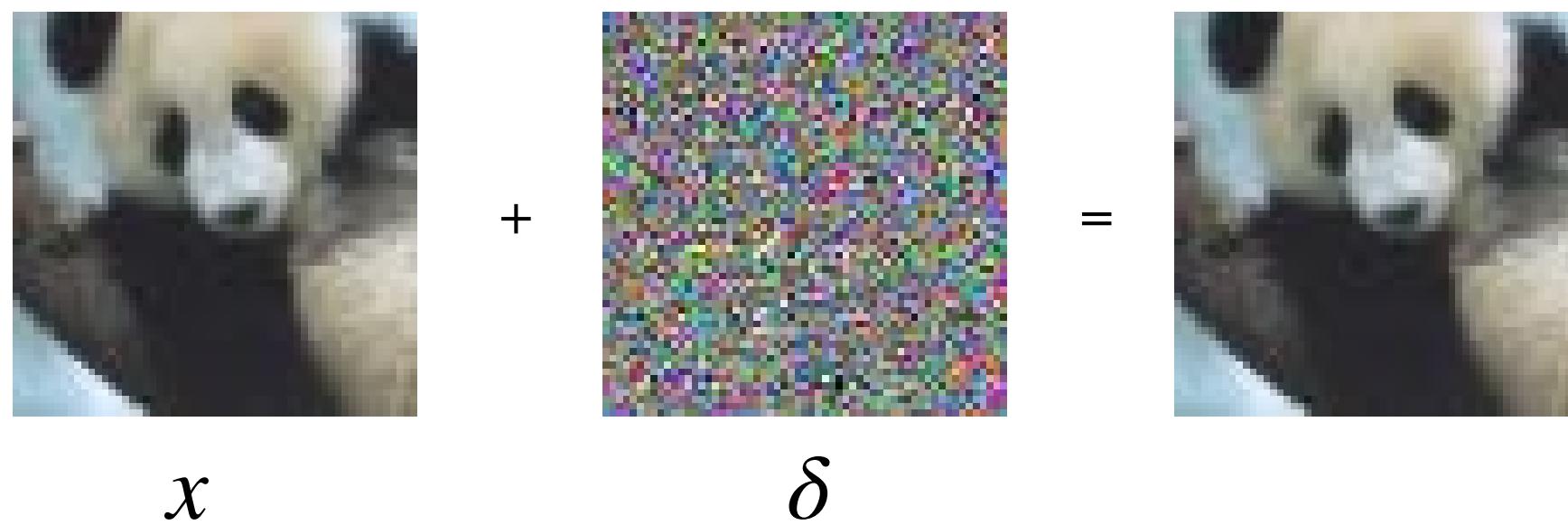
$$x + \delta = \text{adversarial image}$$

# Application of worst-case RO: adversarial robustness of DNN



$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \sup_{\|\delta\| \leq \epsilon} l(f_{\theta}(x_i + \delta), y_i)$$

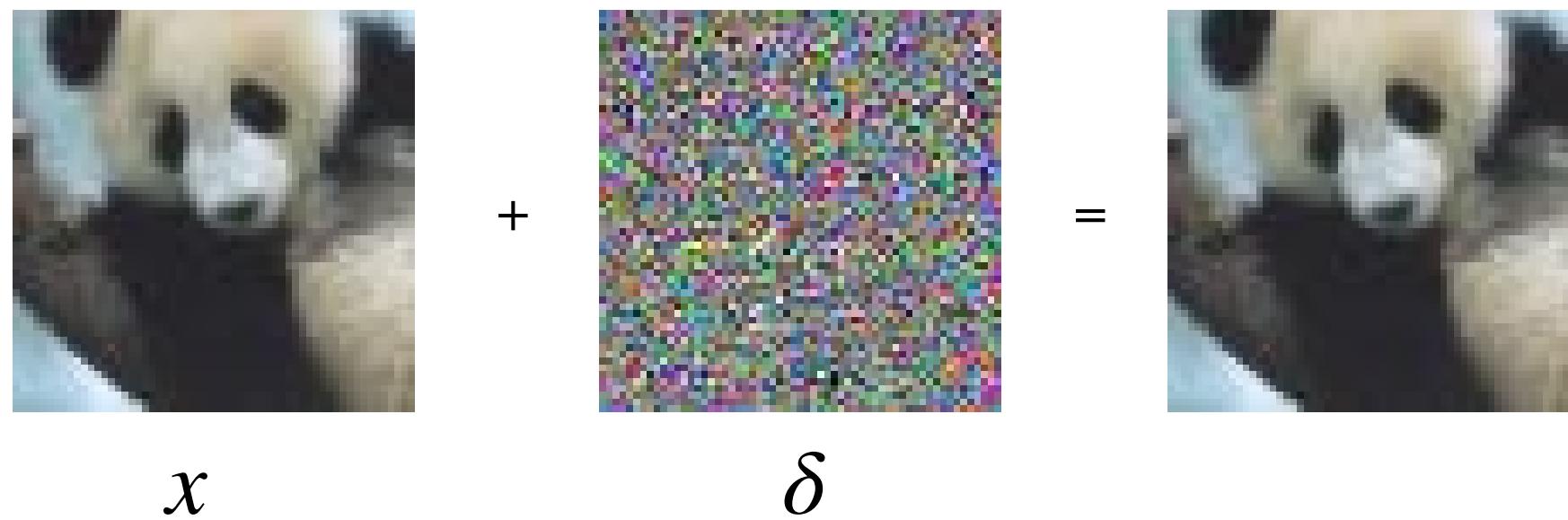
# Application of worst-case RO: adversarial robustness of DNN



$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \sup_{\|\delta\| \leq \epsilon} l(f_{\theta}(x_i + \delta), y_i)$$

- Since deep model typically don't have the nice structure we have seen in the last slide, we are often left to performing SGD and hope for the alignment of the stars.

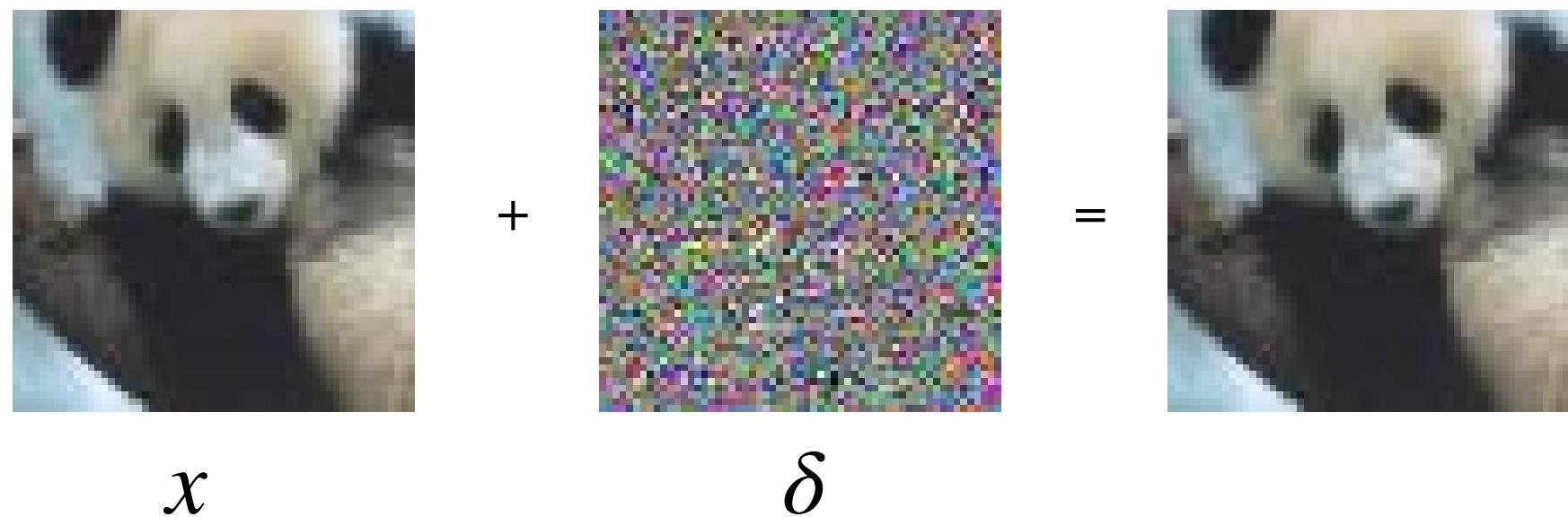
# Application of worst-case RO: adversarial robustness of DNN



$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \sup_{\|\delta\| \leq \epsilon} l(f_{\theta}(x_i + \delta), y_i)$$

- Since deep model typically don't have the nice structure we have seen in the last slide, we are often left to performing SGD and hope for the alignment of the stars.
- Projected gradient descent:  $\delta \leftarrow \text{proj}(\delta + \alpha \nabla_{\delta} l)$

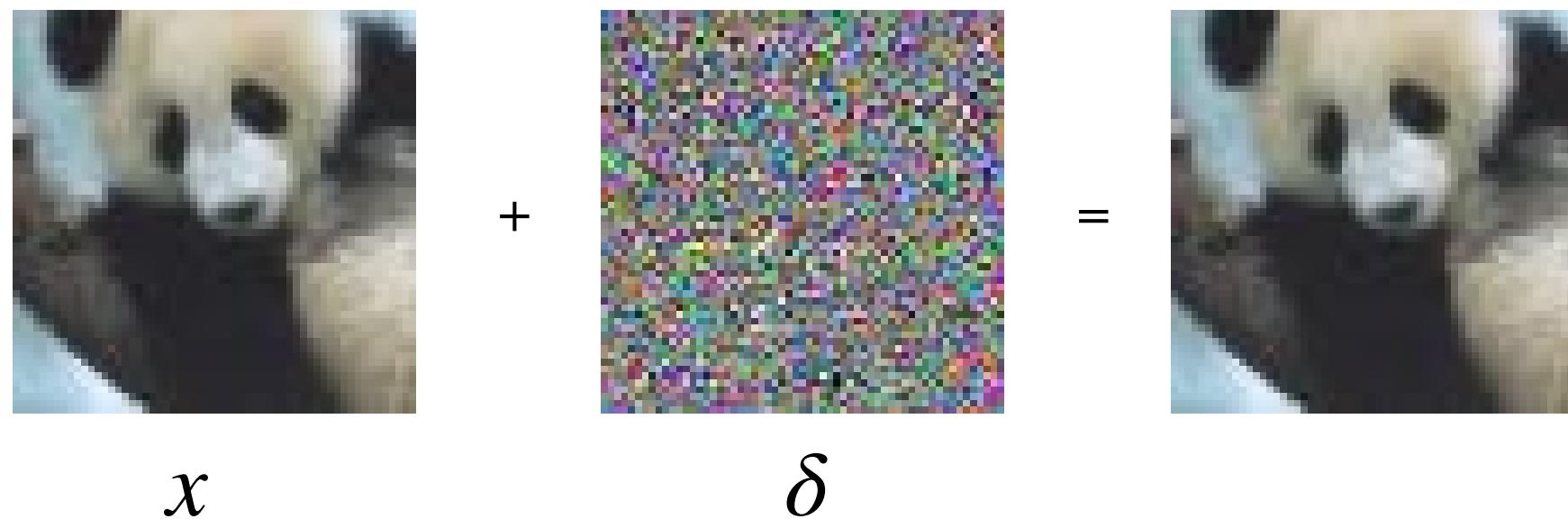
# Application of worst-case RO: adversarial robustness of DNN



$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \sup_{\|\delta\| \leq \epsilon} l(f_{\theta}(x_i + \delta), y_i)$$

- Since deep model typically don't have the nice structure we have seen in the last slide, we are often left to performing SGD and hope for the alignment of the stars.
- Projected gradient descent:  $\delta \leftarrow \text{proj}(\delta + \alpha \nabla_{\delta} l)$ 
  - sample  $(x_i, y_i)$

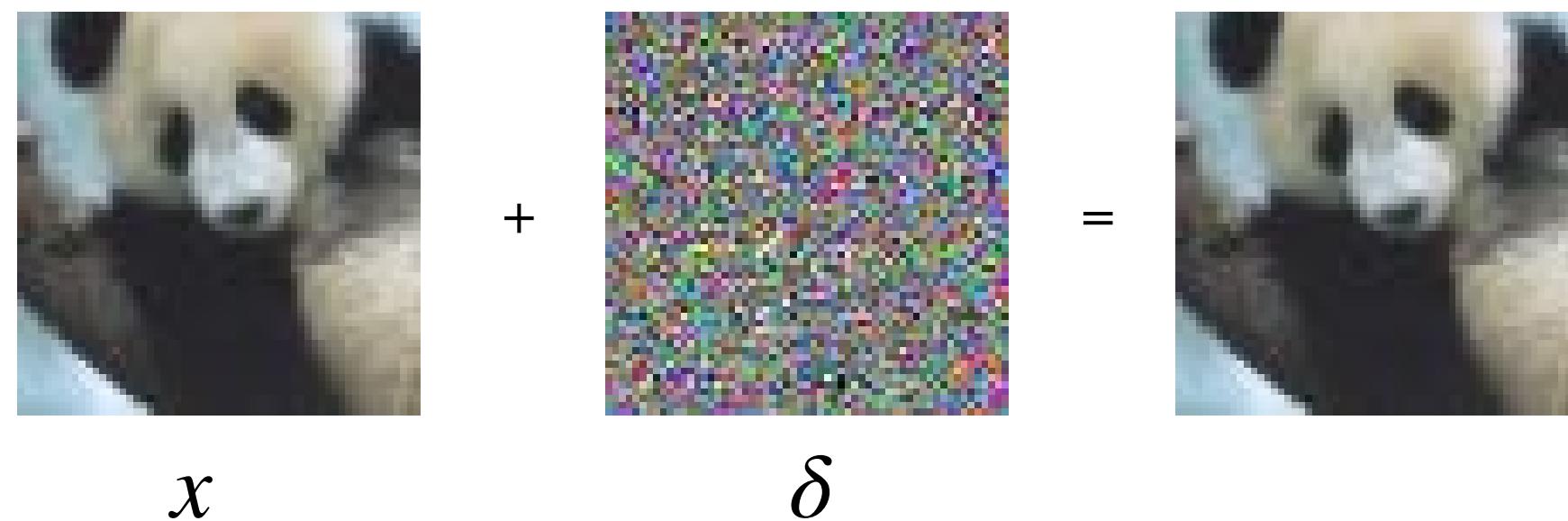
# Application of worst-case RO: adversarial robustness of DNN



$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \sup_{\|\delta\| \leq \epsilon} l(f_{\theta}(x_i + \delta), y_i)$$

- Since deep model typically don't have the nice structure we have seen in the last slide, we are often left to performing SGD and hope for the alignment of the stars.
- Projected gradient descent:  $\delta \leftarrow \text{proj}(\delta + \alpha \nabla_{\delta} l)$ 
  - sample  $(x_i, y_i)$
  - perform gradient ascent for the inner maximization problem over  $\delta$

# Application of worst-case RO: adversarial robustness of DNN

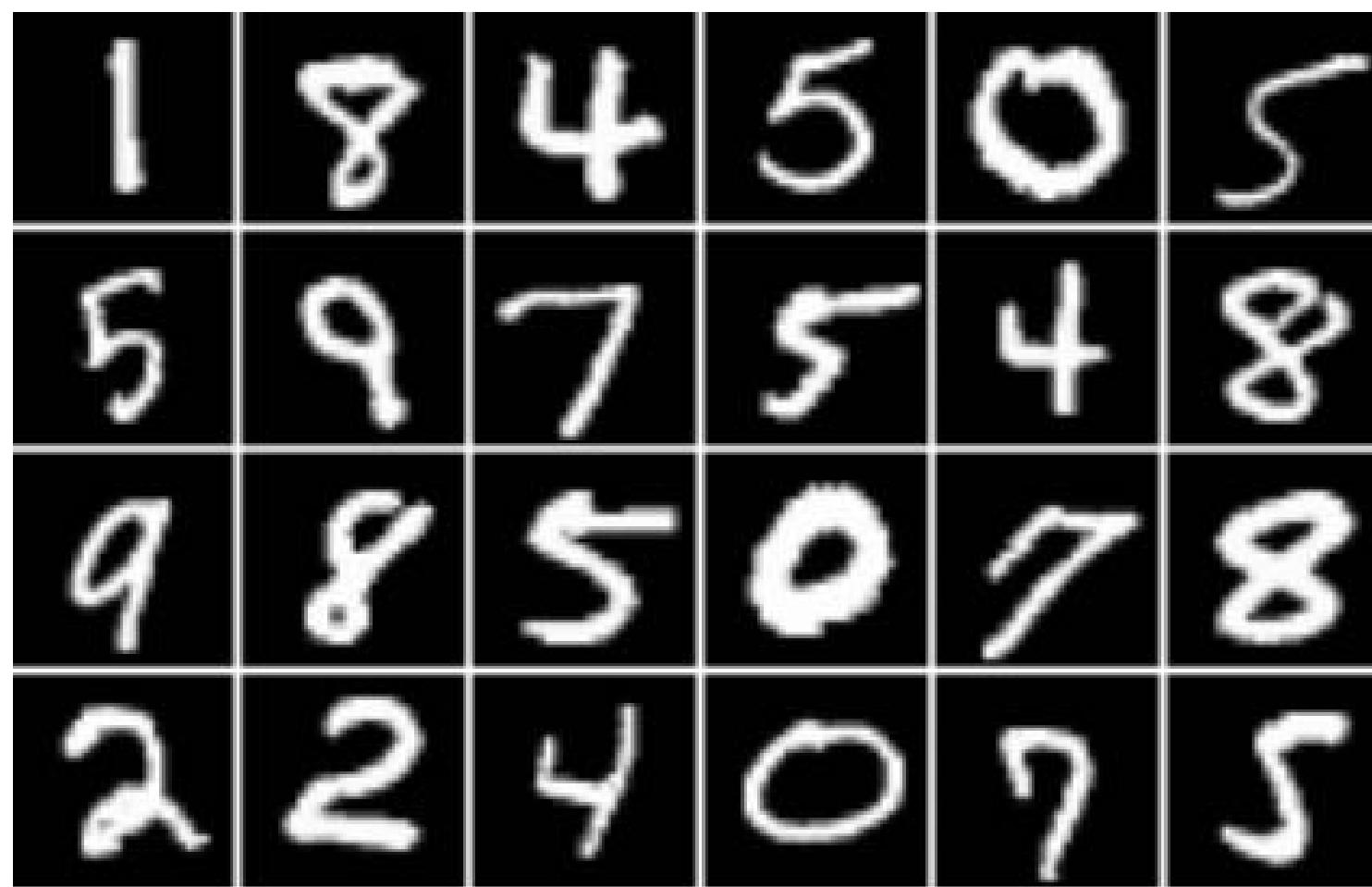


$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \sup_{\|\delta\| \leq \epsilon} l(f_{\theta}(x_i + \delta), y_i)$$

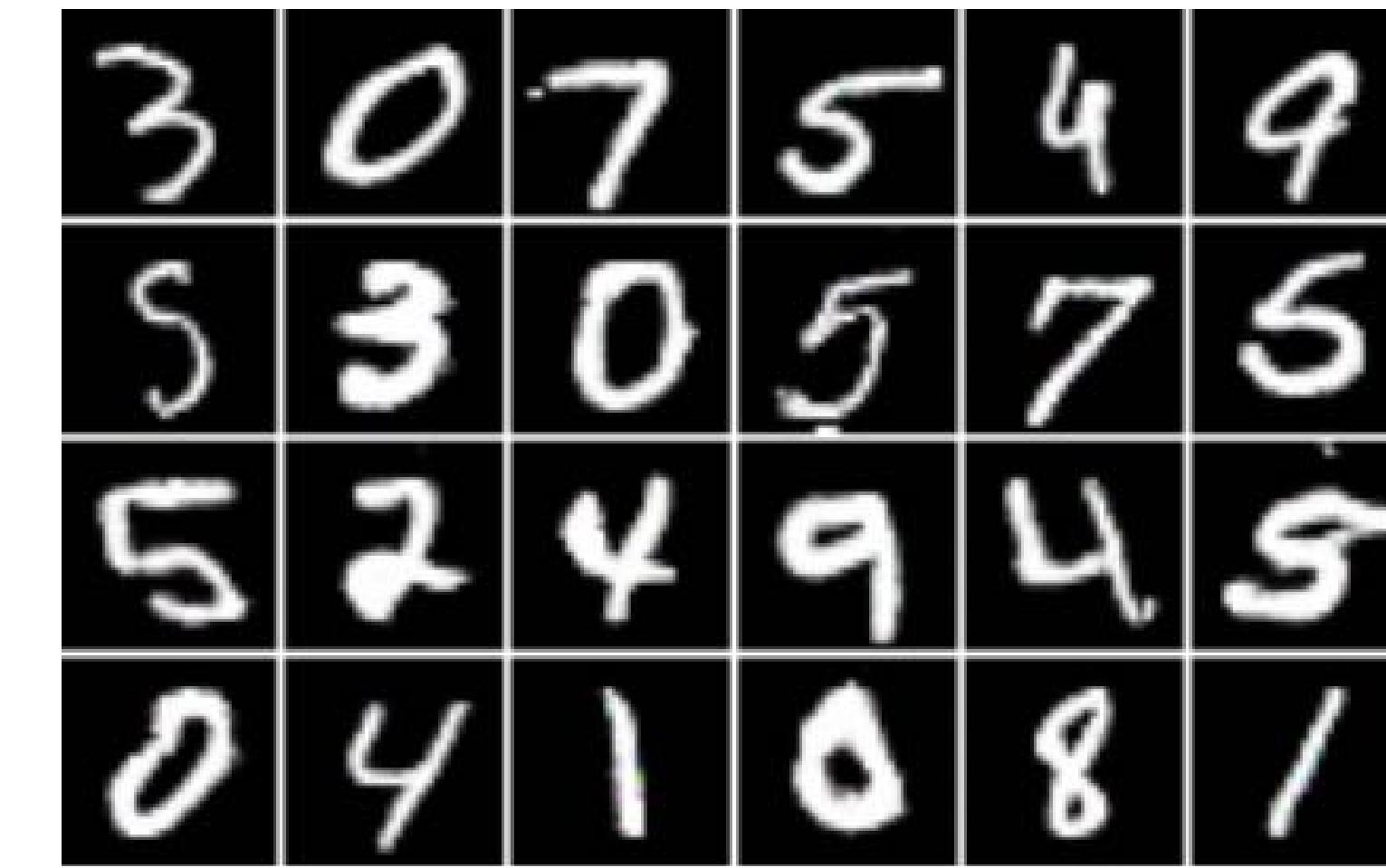
- Since deep model typically don't have the nice structure we have seen in the last slide, we are often left to performing SGD and hope for the alignment of the stars.
- Projected gradient descent:  $\delta \leftarrow \text{proj}(\delta + \alpha \nabla_{\delta} l)$ 
  - sample  $(x_i, y_i)$
  - perform gradient ascent for the inner maximization problem over  $\delta$
  - project back to the norm-ball

# Distributional Robustness

# How do we compare two sets of data?



Observed MNIST handwritten digits.  $\mathbf{X}$ .

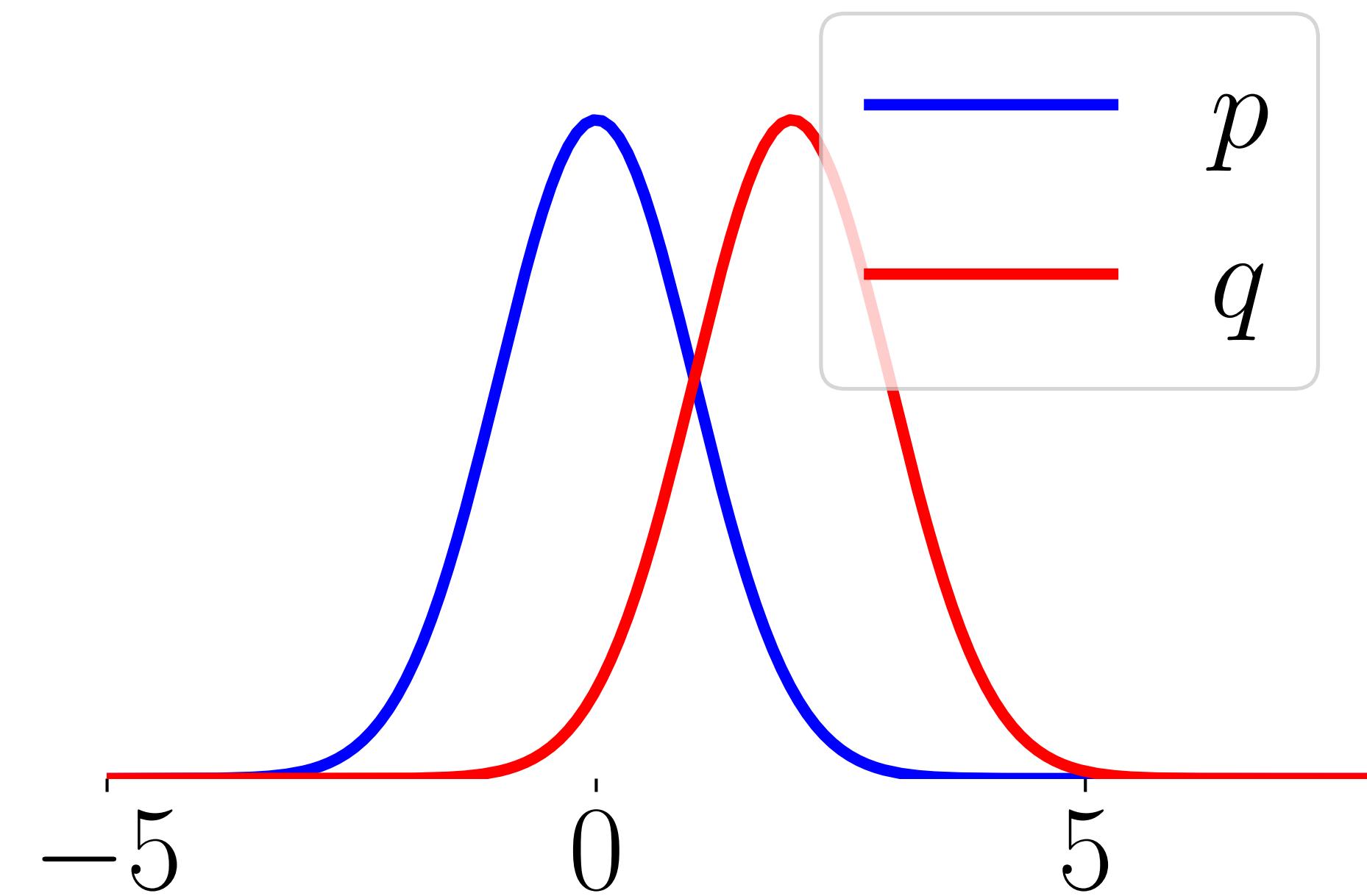


Generated images from a model.  $\mathbf{Y}$ .

Is  $\mathbf{Y}$  similar to  $\mathbf{X}$ ?

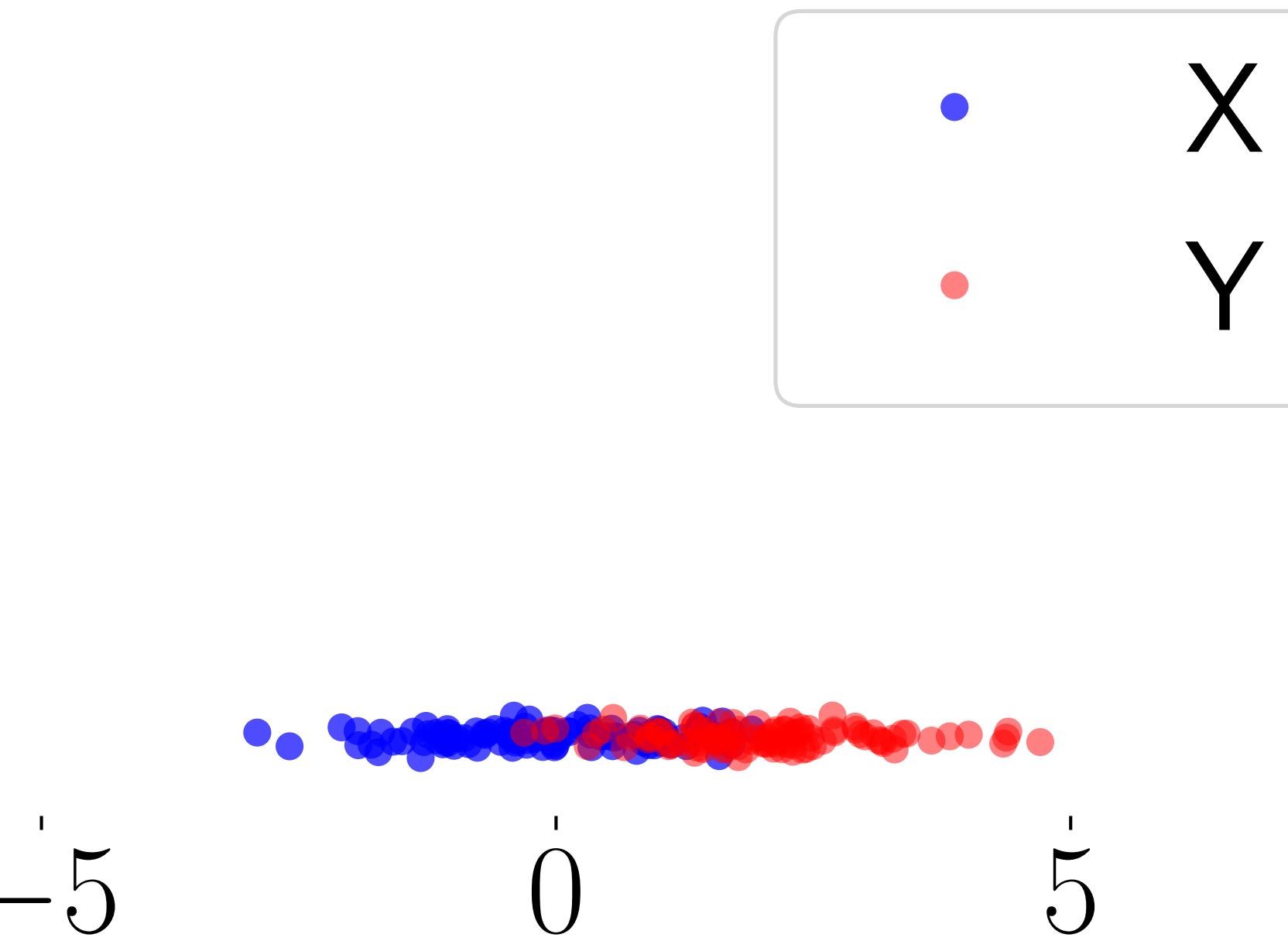
- Distance between distributions can be used to train generative models.

## Case 1: Simple Mean Shift in 1D



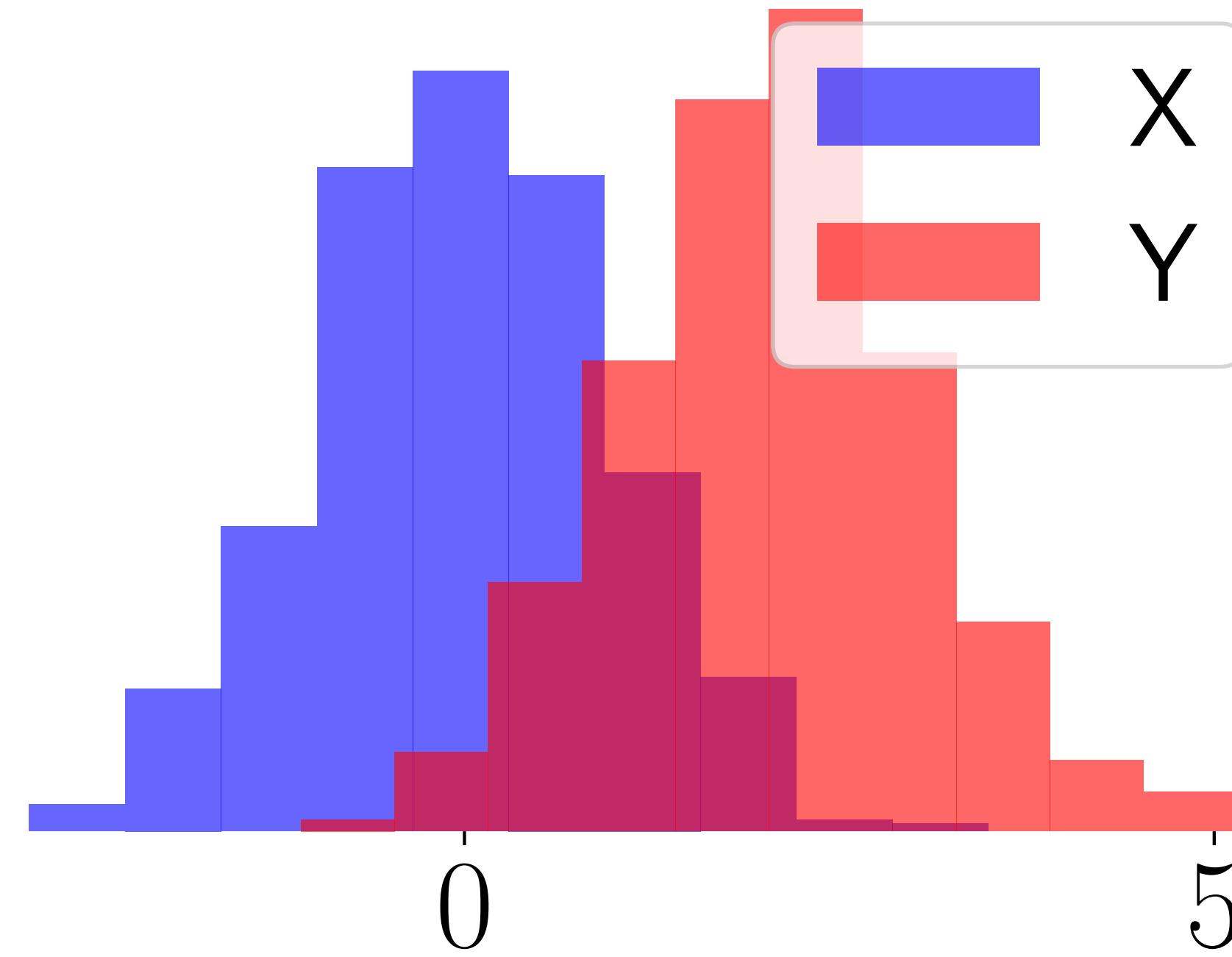
- Two Gaussian distributions.

## Case 1: Simple Mean Shift in 1D



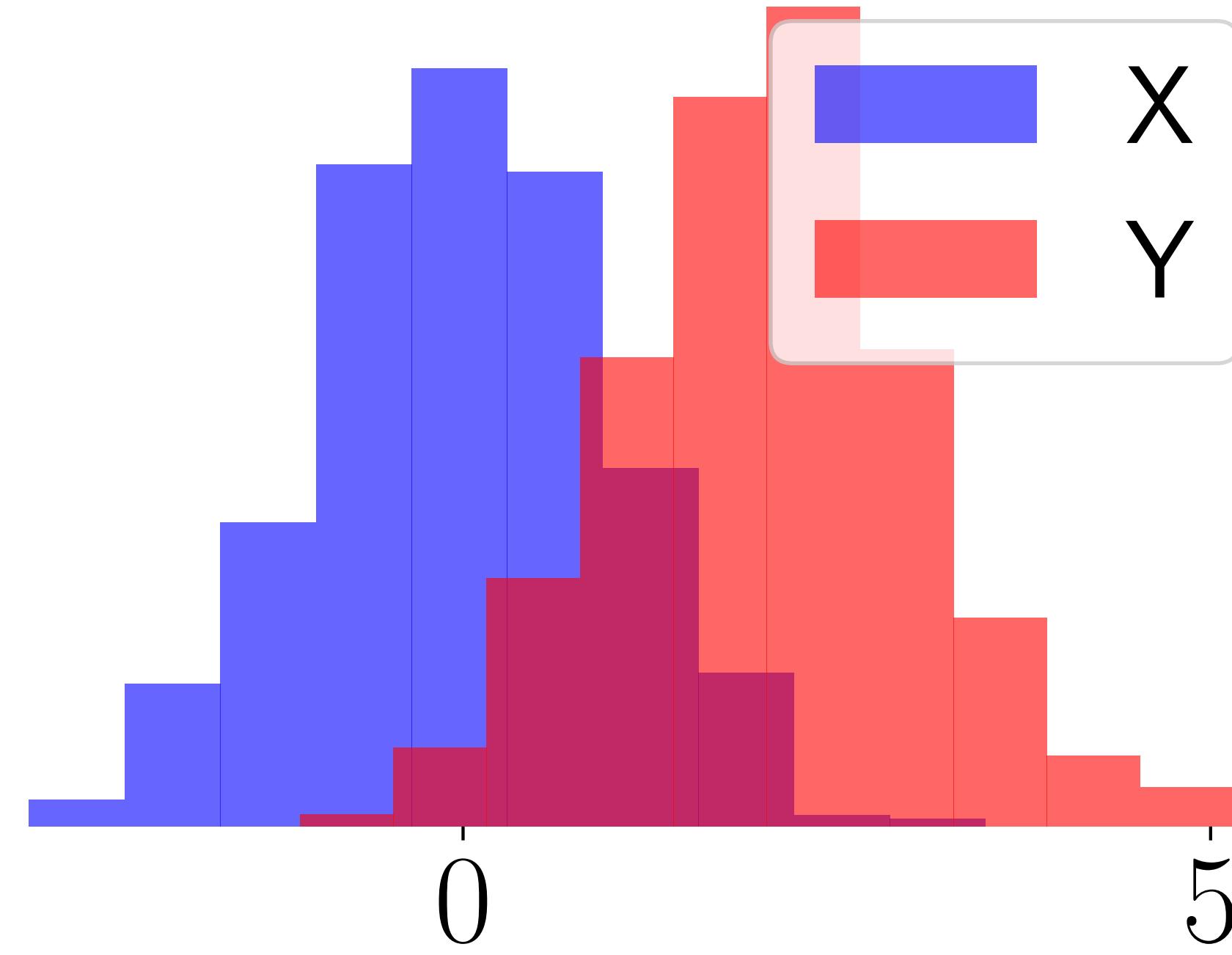
- We have only samples  $X \sim p$  and  $Y \sim q$ .
- $X = \{x_1, \dots, x_n\}$  and  $Y = \{y_1, \dots, y_n\}$ . Sets of numbers.

## Case 1: Simple Mean Shift in 1D



- We have only samples  $X \sim p$  and  $Y \sim q$ .
- $X = \{x_1, \dots, x_n\}$  and  $Y = \{y_1, \dots, y_n\}$ . Sets of numbers.

## Case 1: Simple Mean Shift in 1D

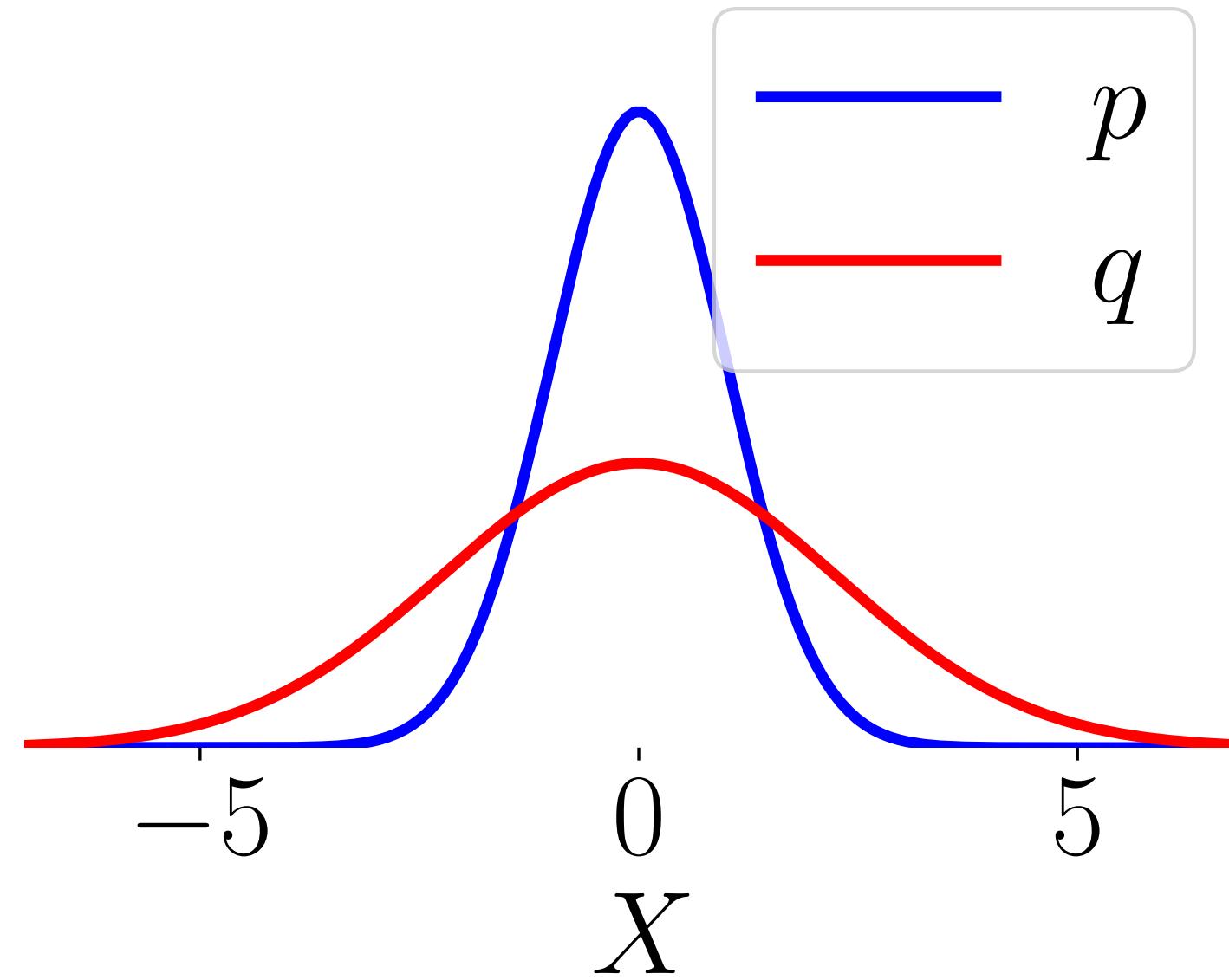


- Assume no difference in high-order moments.
- “Distance” = difference in the means. T-test.

$$\text{(population)} \ D_1(p, q) := |\mathbb{E}_{X \sim p}[X] - \mathbb{E}_{Y \sim q}[Y]|$$

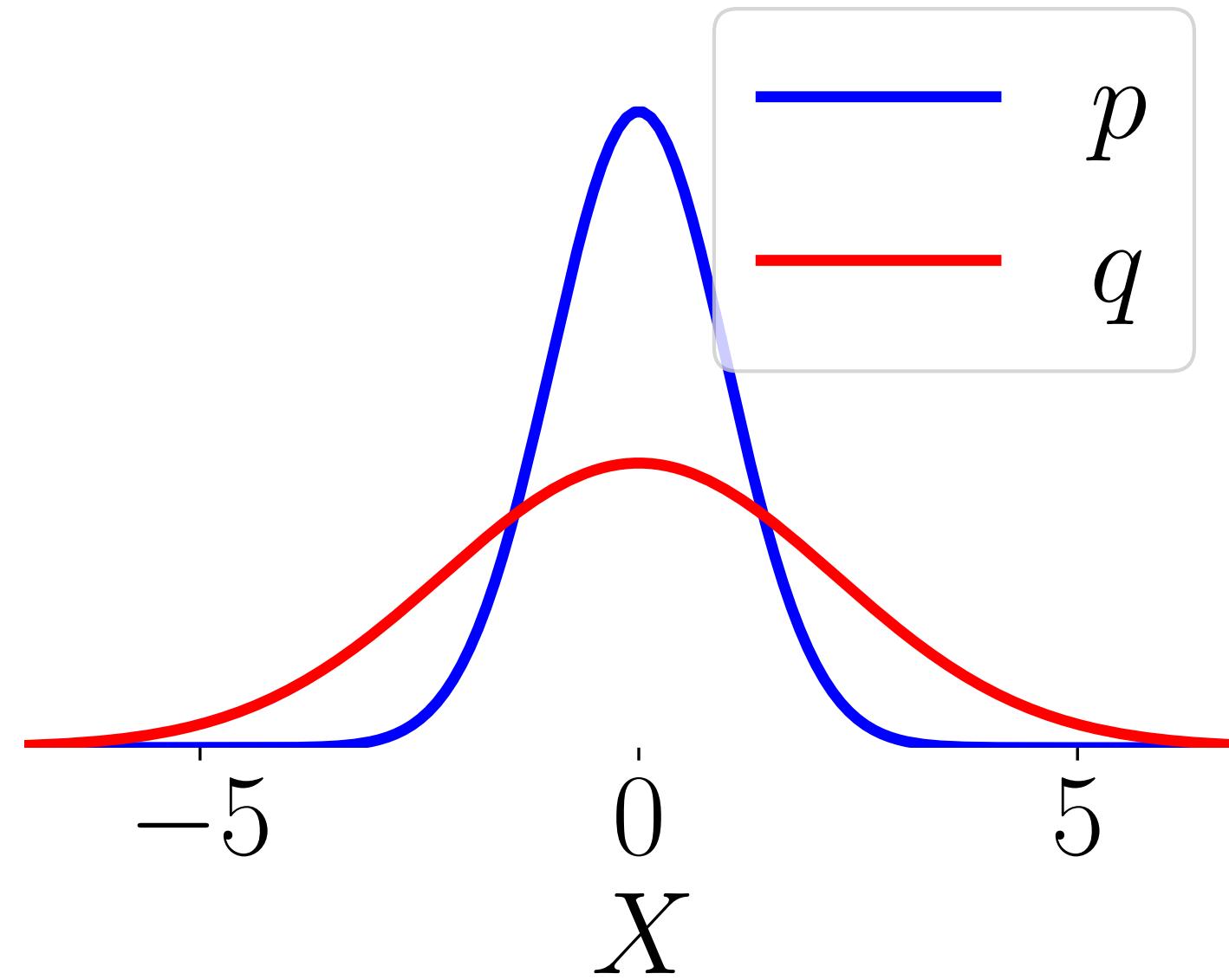
$$\text{(empirical)} \ \hat{D}_1(\mathbf{X}, \mathbf{Y}) = \left| \frac{1}{n} \sum_{i=1}^n x_i - \frac{1}{n} \sum_{j=1}^n y_j \right|$$

## Case 2: Same Mean, Different Variances



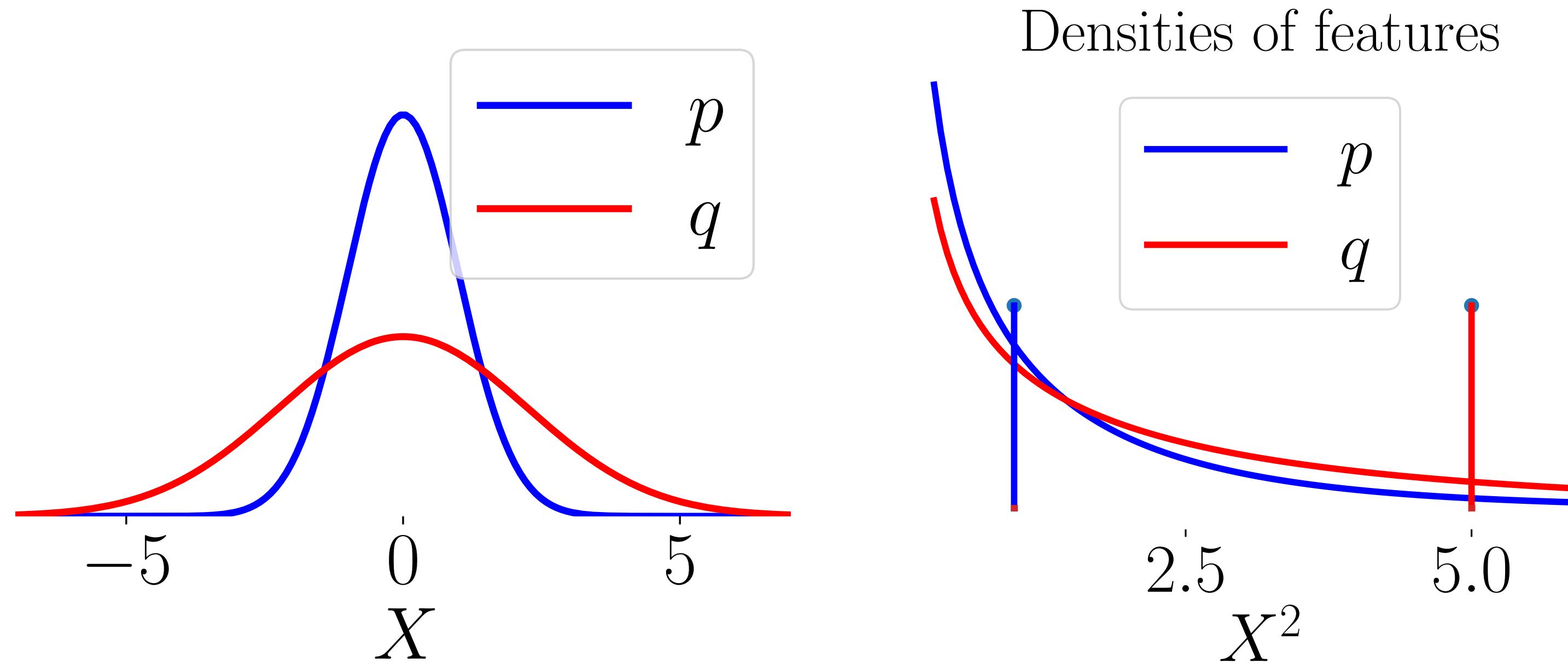
- $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] - \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference.  
Why?
- Idea: look at difference in means of features  $\phi(\cdot)$  of  $X$  and  $Y$ .

## Case 2: Same Mean, Different Variances



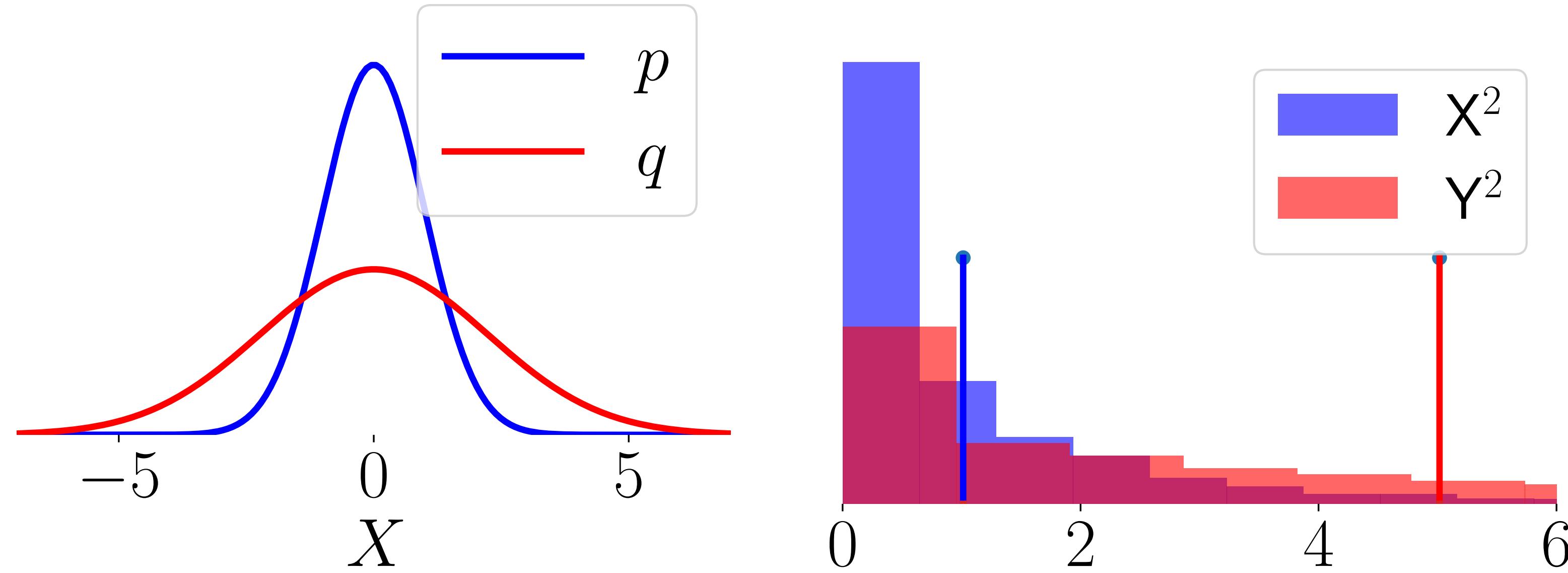
- $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] - \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference.  
Why?
- Idea: look at difference in means of features  $\phi(\cdot)$  of  $X$  and  $Y$ .

## Case 2: Same Mean, Different Variances



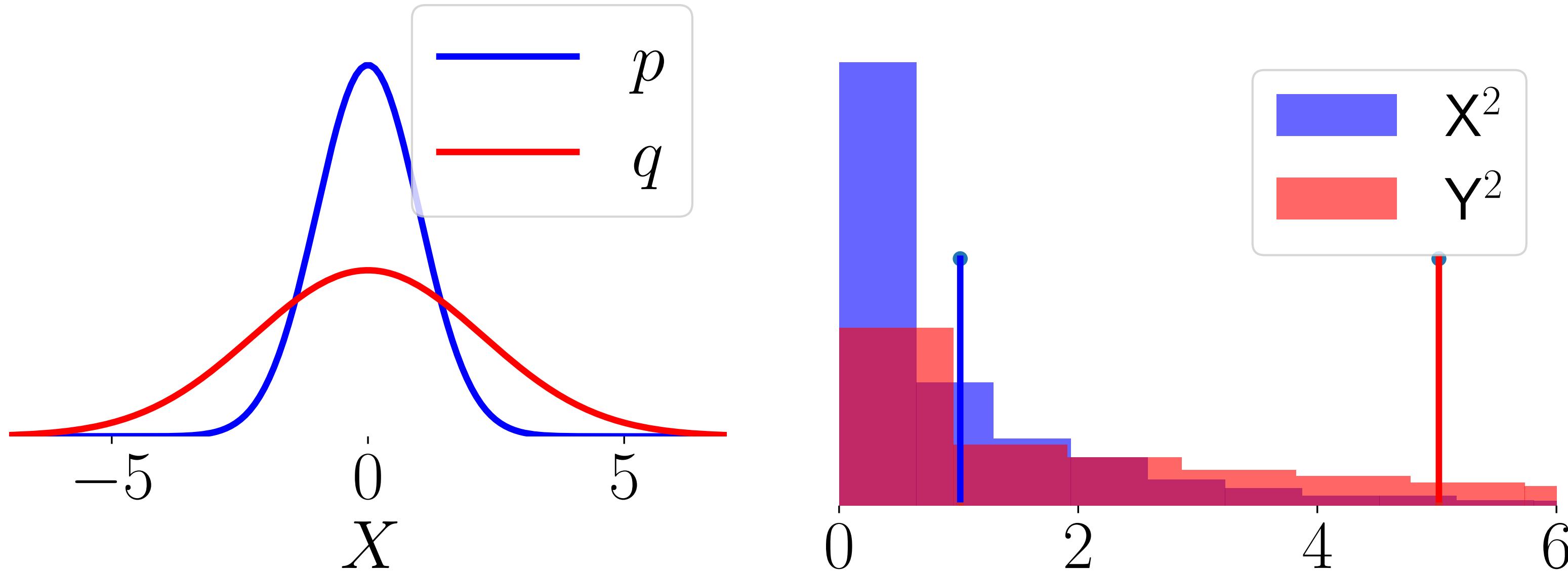
- $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] - \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference.  
Why?
- Idea: look at difference in means of features  $\phi(\cdot)$  of  $X$  and  $Y$ .

## Case 2: Same Mean, Different Variances



- $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] - \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference.  
Why?
- Idea: look at difference in means of features  $\phi(\cdot)$  of  $X$  and  $Y$ .

## Case 2: Same Mean, Different Variances

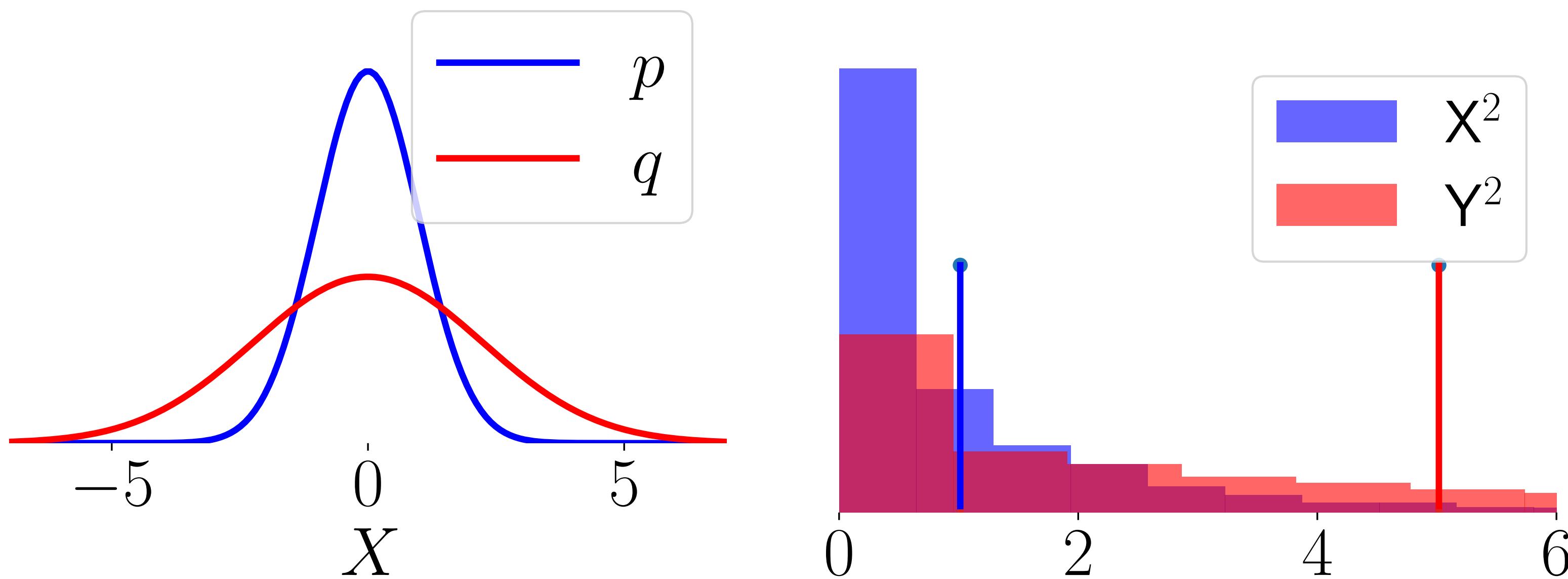


- $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] - \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference.  
Why?
- Idea: look at difference in means of features  $\phi(\cdot)$  of  $X$  and  $Y$ .
- New “distance”:

$$D_2(p, q) = \|\mathbb{E}_{X \sim p}[\phi(X)] - \mathbb{E}_{Y \sim q}[\phi(X)]\|,$$

where  $\phi(x) = (x, x^2)^\top$ .

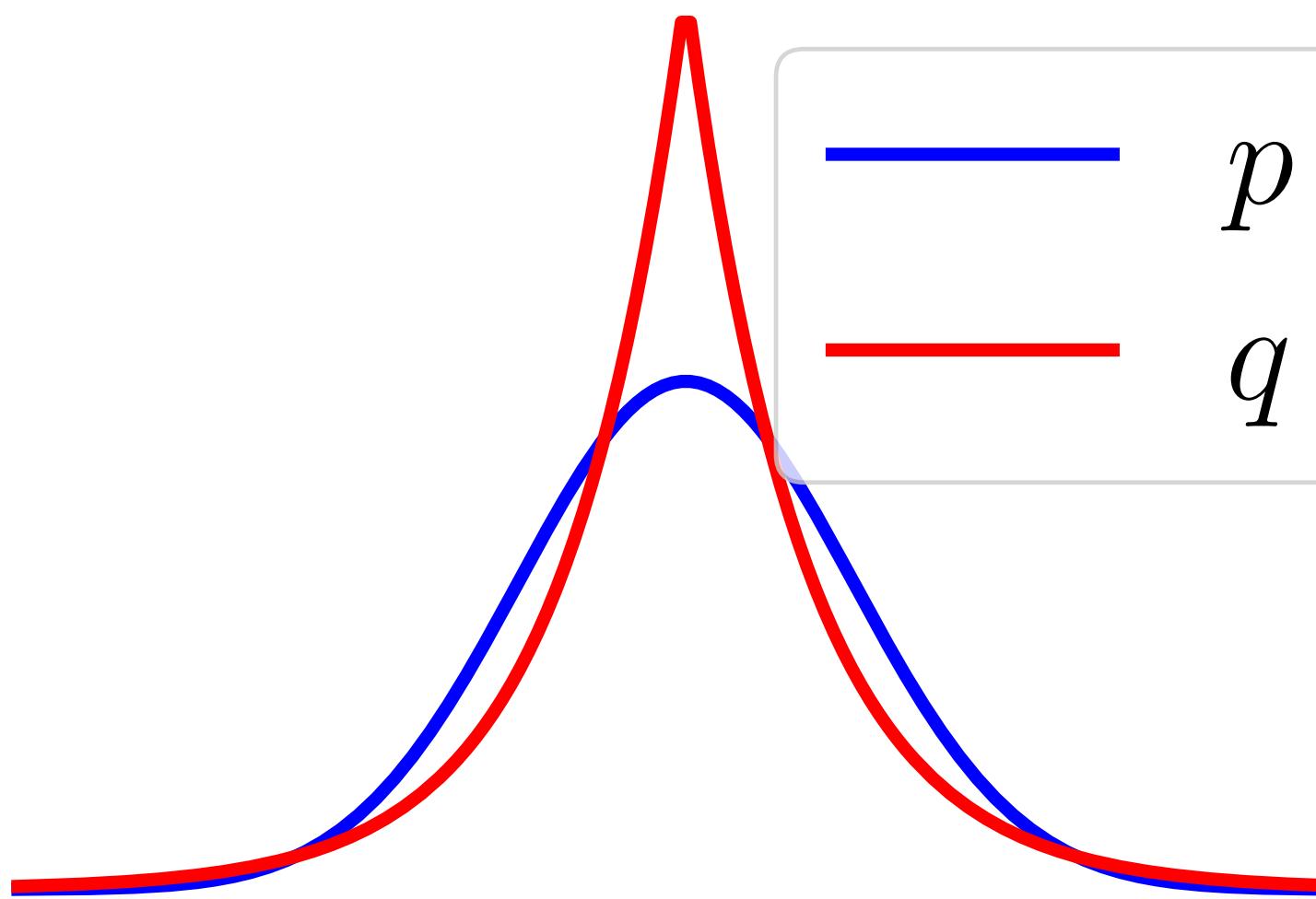
## Case 2: Same Mean, Different Variances



- $D_1(p, q) := |\mathbb{E}_{X \sim p}[X] - \mathbb{E}_{Y \sim q}[Y]|$  cannot detect the difference.  
Why?
- Idea: look at difference in means of features  $\phi(\cdot)$  of  $X$  and  $Y$ .

$$D_2(p, q) = \left\| \begin{pmatrix} \mathbb{E}_{X \sim p}[X] \\ \mathbb{E}_{X \sim p}[X^2] \end{pmatrix} - \begin{pmatrix} \mathbb{E}_{X \sim q}[X] \\ \mathbb{E}_{X \sim q}[X^2] \end{pmatrix} \right\|.$$

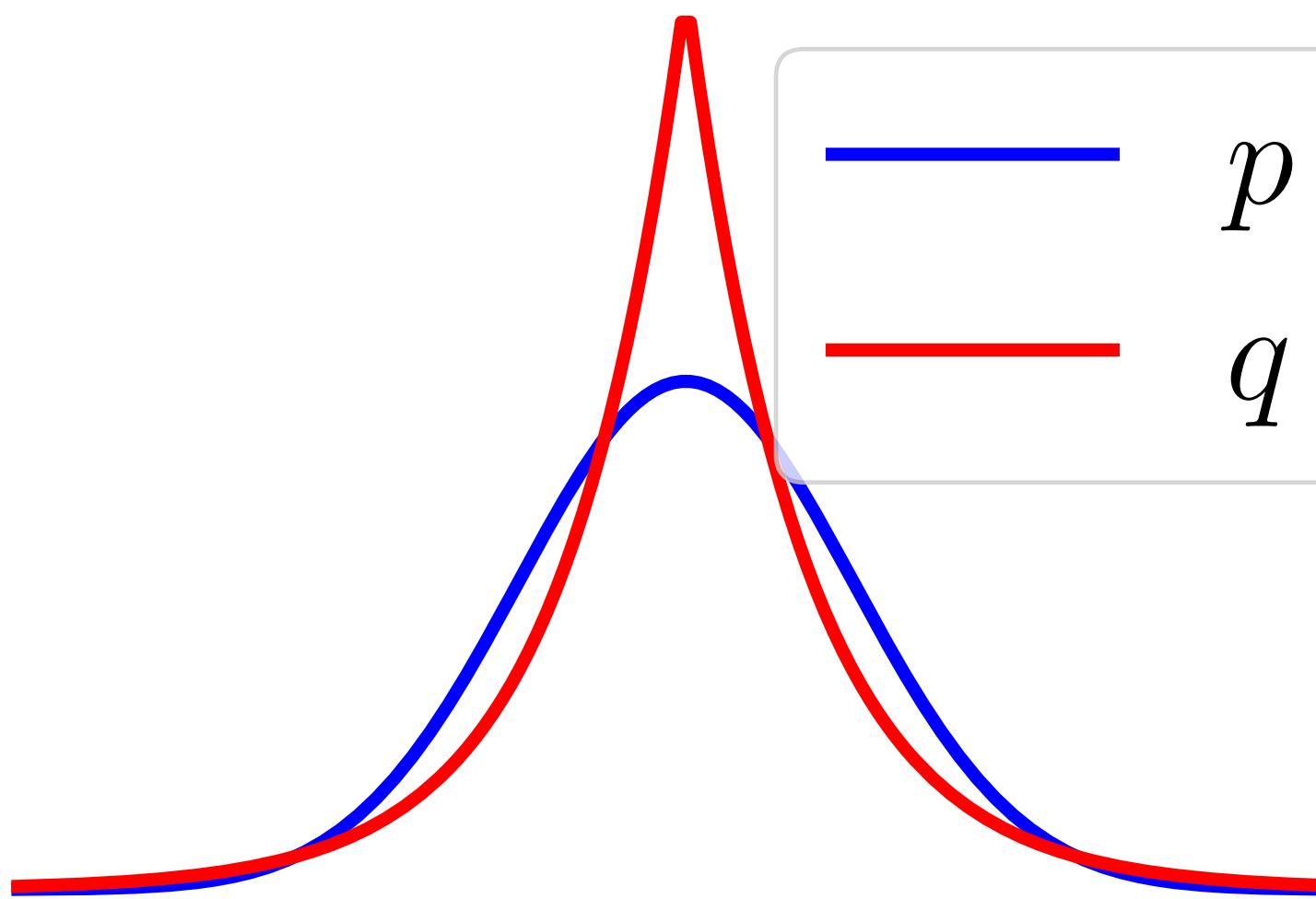
## Case 3: Difference in High-Order Moments



- $p = \text{Gaussian distribution},$   
 $q = \text{Laplace distribution}.$
- Same mean and variance.
- $D_2$  fails.

- $\phi(x) = (x, x^2, x^4)^\top$  works. Difference is in kurtosis ( $4^{\text{th}}$  moment).
- $\phi(x) = (x, x^2, x^4, \cos x, e^x, \dots)^\top$ . But, when to stop?
- **Solution:** Use an infinite-dimensional feature map  $\phi(\cdot)$  with the kernel trick.

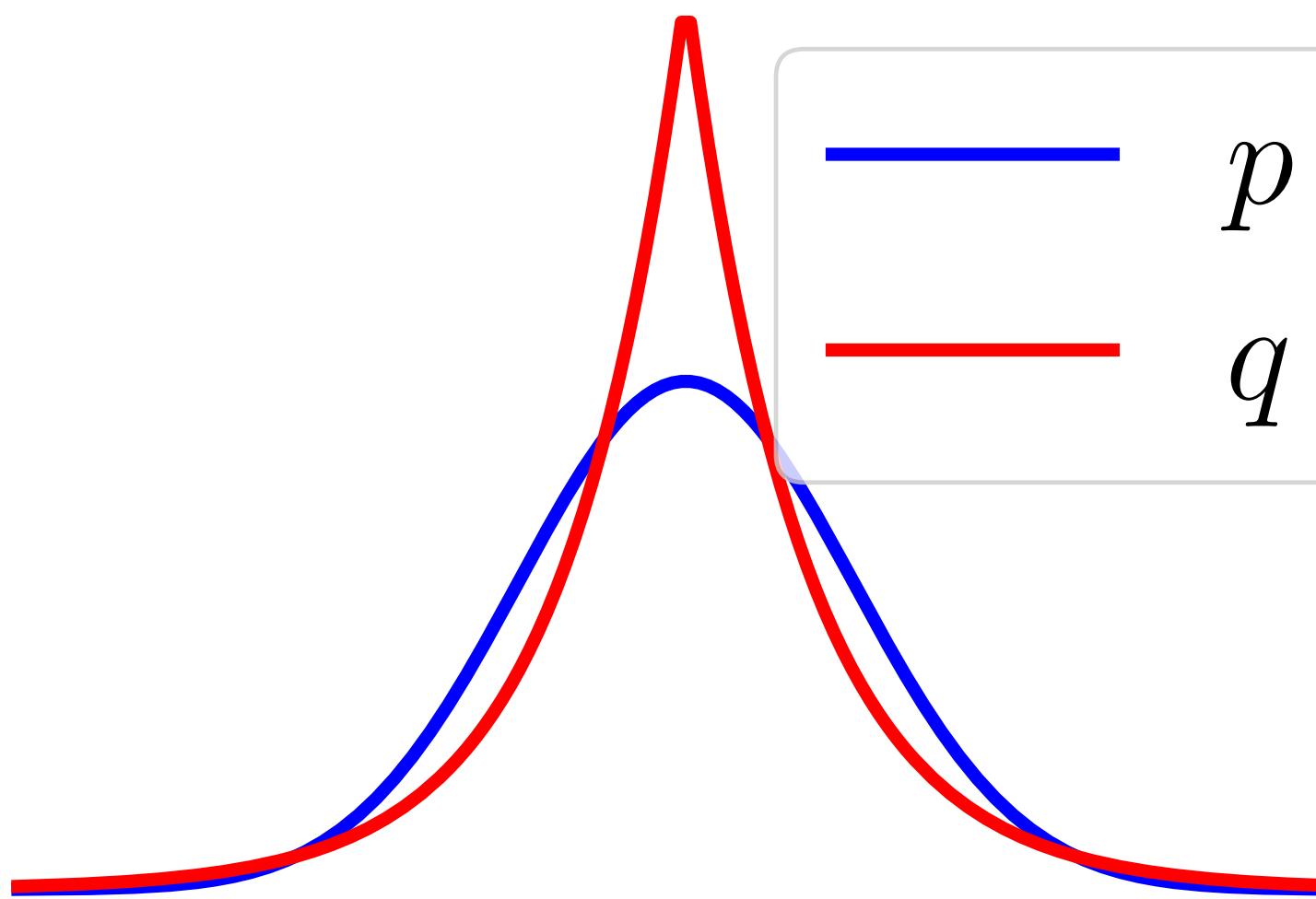
## Case 3: Difference in High-Order Moments



- $p = \text{Gaussian distribution}$ ,  
 $q = \text{Laplace distribution}$ .
- Same mean and variance.
- $D_2$  fails.

- $\phi(x) = (x, x^2, x^4)^\top$  works. Difference is in kurtosis ( $4^{\text{th}}$  moment).
- $\phi(x) = (x, x^2, x^4, \cos x, e^x, \dots)^\top$ . But, when to stop?
- **Solution:** Use an infinite-dimensional feature map  $\phi(\cdot)$  with the kernel trick.

## Case 3: Difference in High-Order Moments



- $p = \text{Gaussian distribution}$ ,  
 $q = \text{Laplace distribution}$ .
- Same mean and variance.
- $D_2$  fails.

- $\phi(x) = (x, x^2, x^4)^\top$  works. Difference is in kurtosis ( $4^{\text{th}}$  moment).
- $\phi(x) = (x, x^2, x^4, \cos x, e^x, \dots)^\top$ . But, when to stop?
- **Solution:** Use an infinite-dimensional feature map  $\phi(\cdot)$  with the kernel trick.

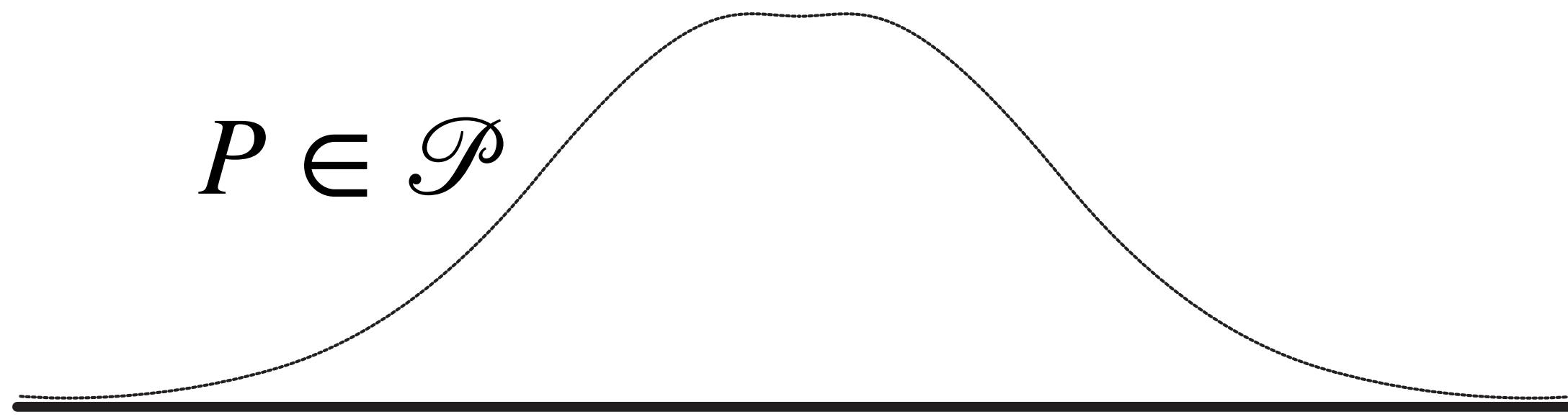
# Working with probability measures

# Working with probability measures

The probability “simplex”  $\mathcal{P}$

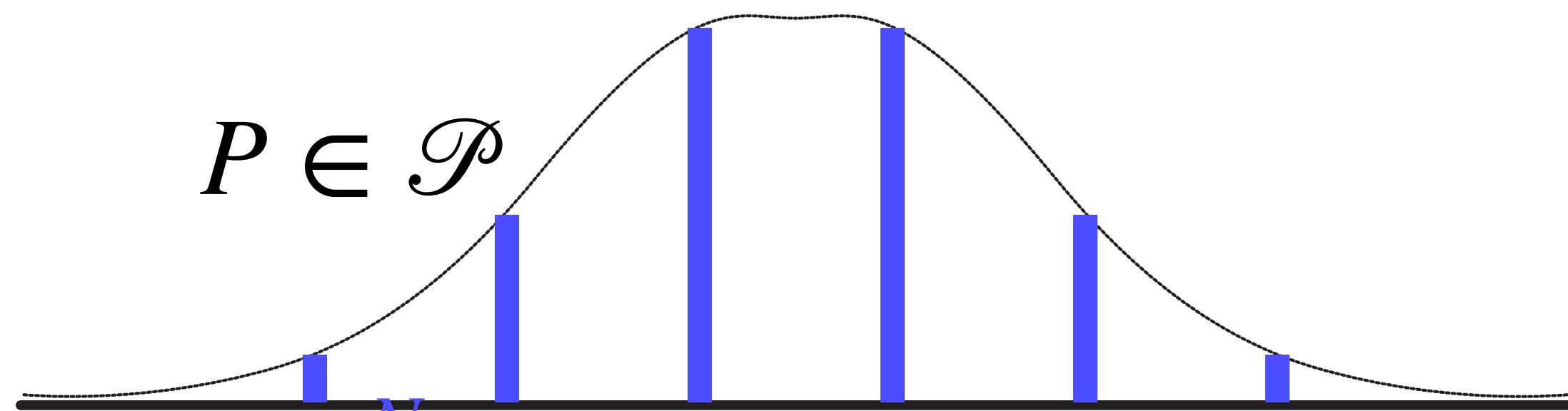
# Working with probability measures

The probability “simplex”  $\mathcal{P}$

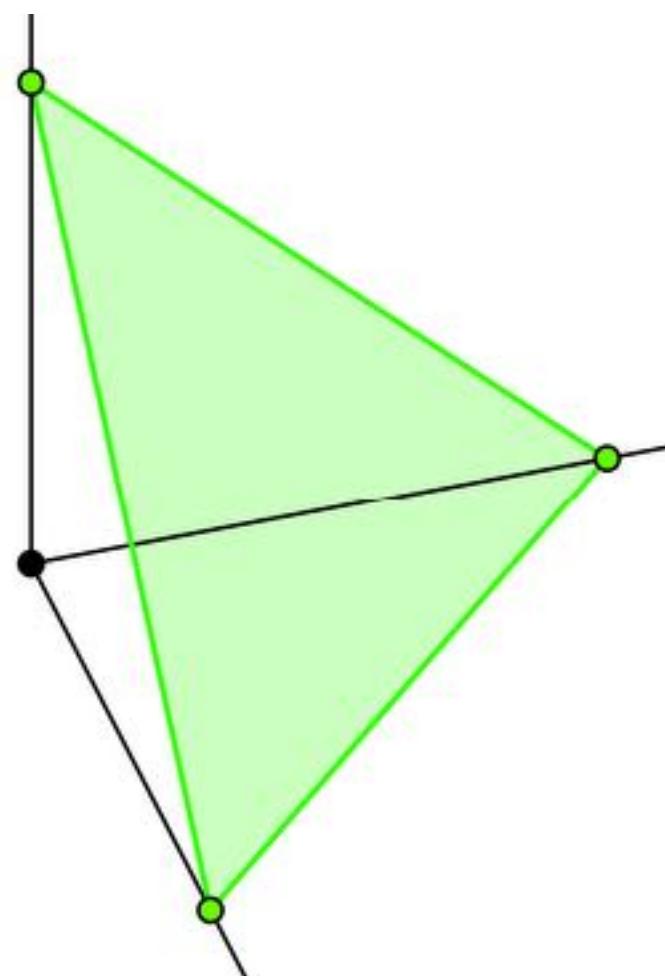


# Working with probability measures

The probability “simplex”  $\mathcal{P}$

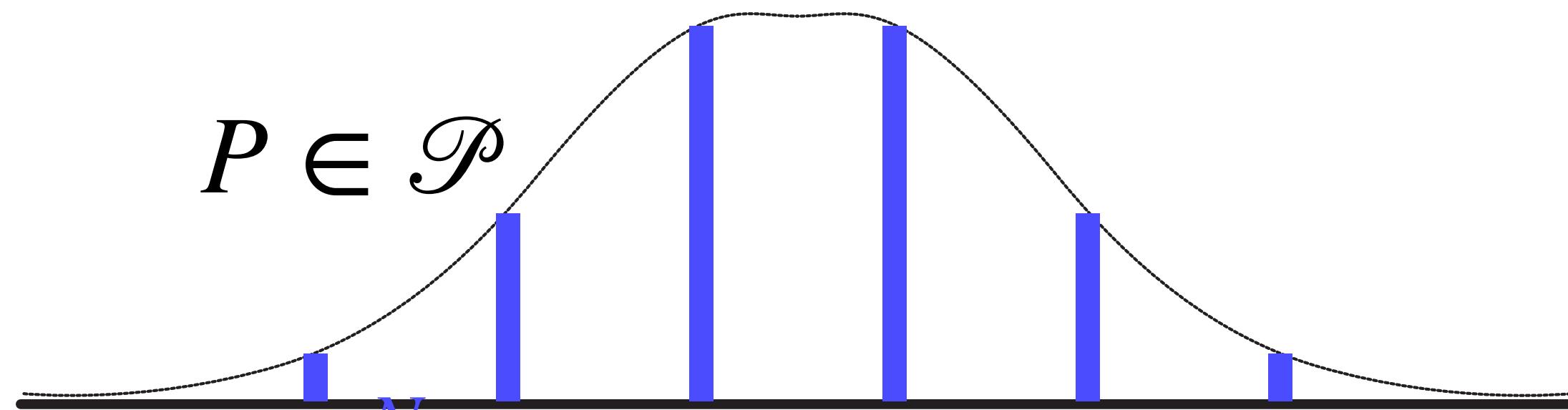


$$\sum_{i=1}^N \gamma_i \delta_{\xi_i} \quad \gamma \in \Delta_d \subset \mathbb{R}^d$$



# Working with probability measures

The probability “simplex”  $\mathcal{P}$

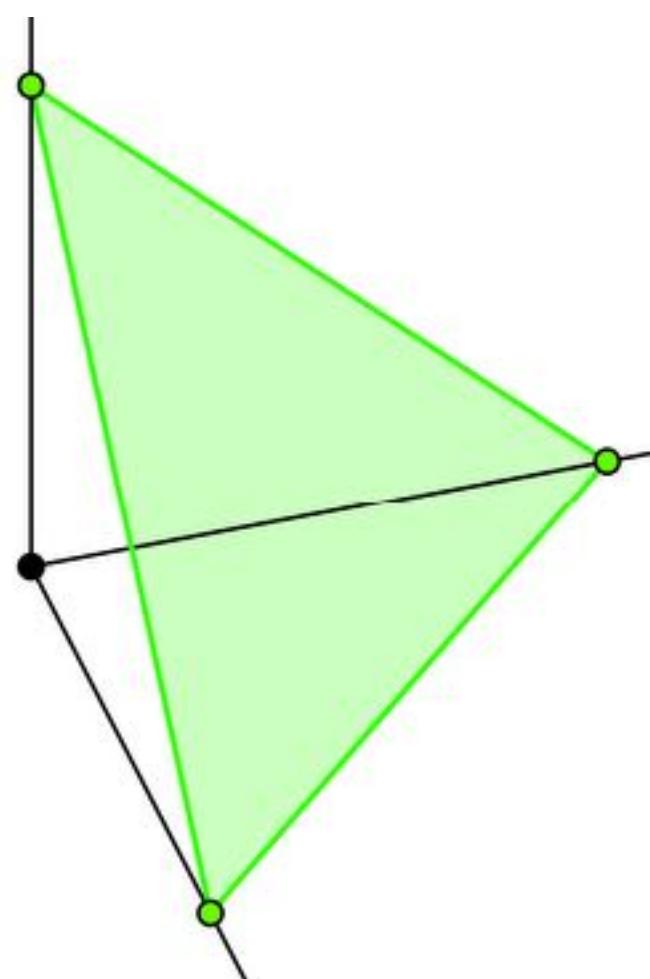


$$\sum_{i=1}^N \gamma_i \delta_{\xi_i} \quad \gamma \in \Delta_d \subset \mathbb{R}^d$$

Empirical data distribution

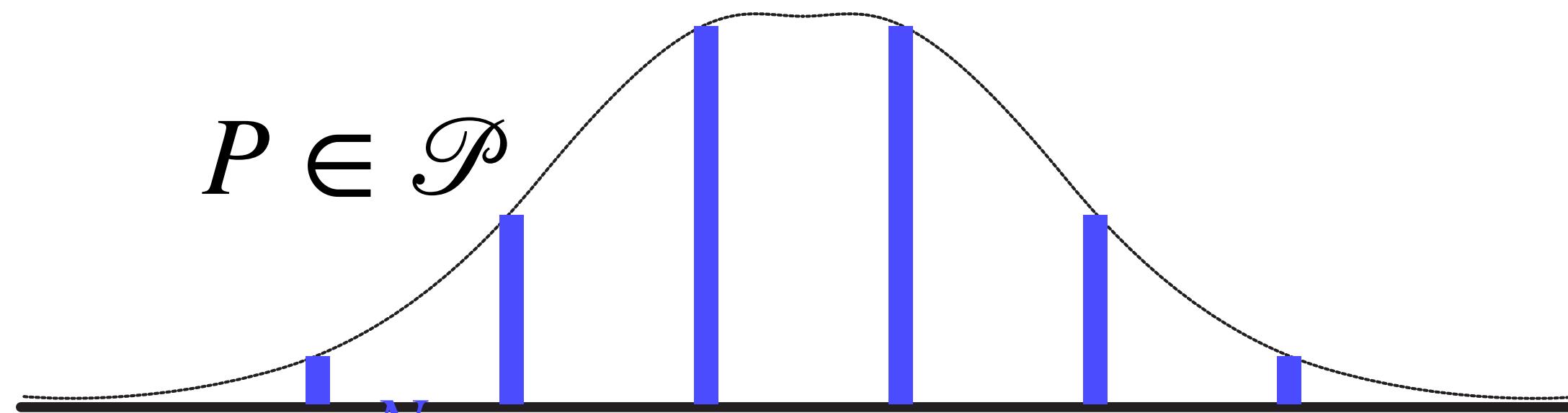
Given data samples:  $\xi_1, \xi_2, \dots, \xi_N$

$$\text{empirical data distribution } \hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$$



# Working with probability measures

The probability “simplex”  $\mathcal{P}$

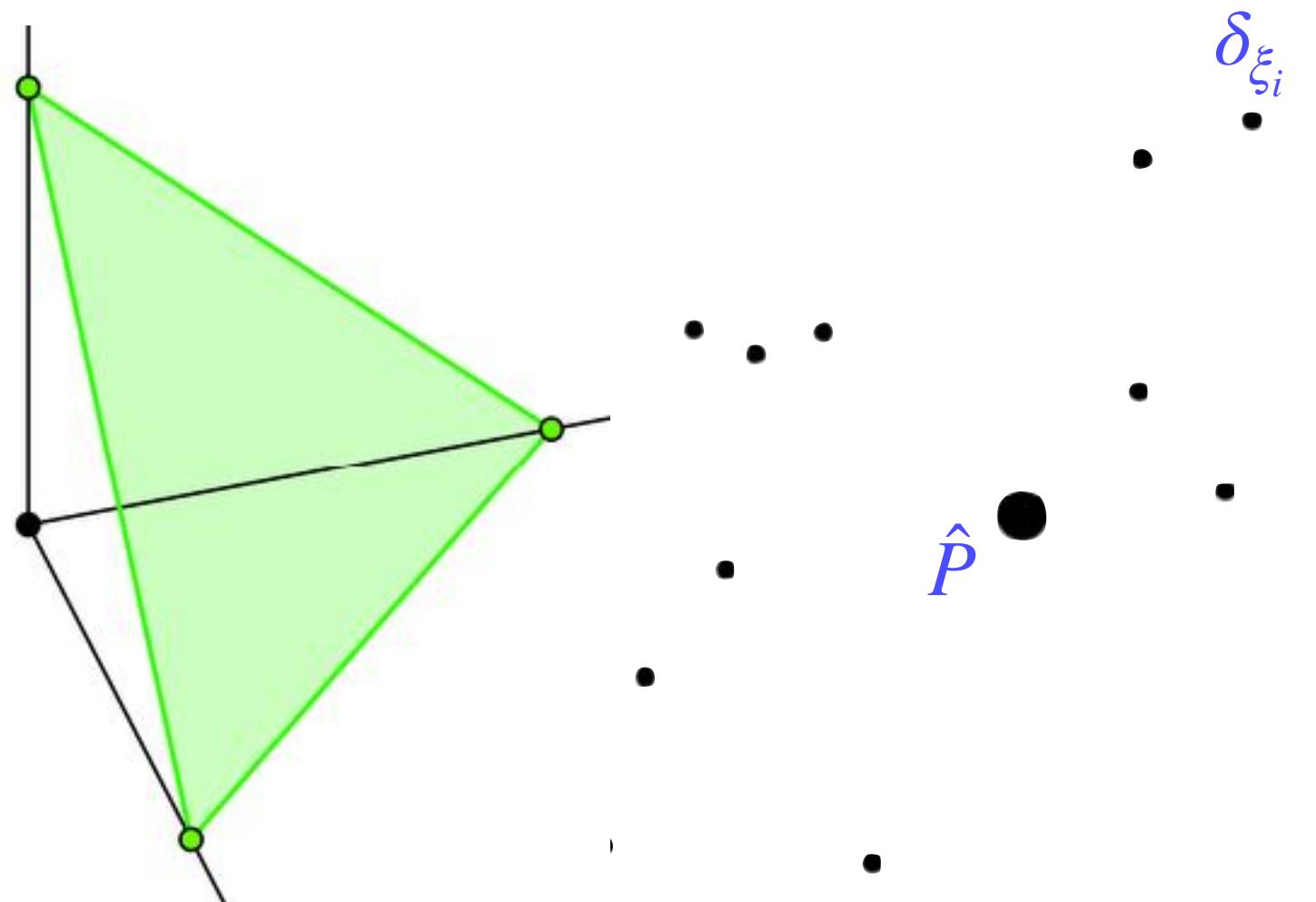


$$\sum_{i=1}^N \gamma_i \delta_{\xi_i} \quad \gamma \in \Delta_d \subset \mathbb{R}^d$$

Empirical data distribution

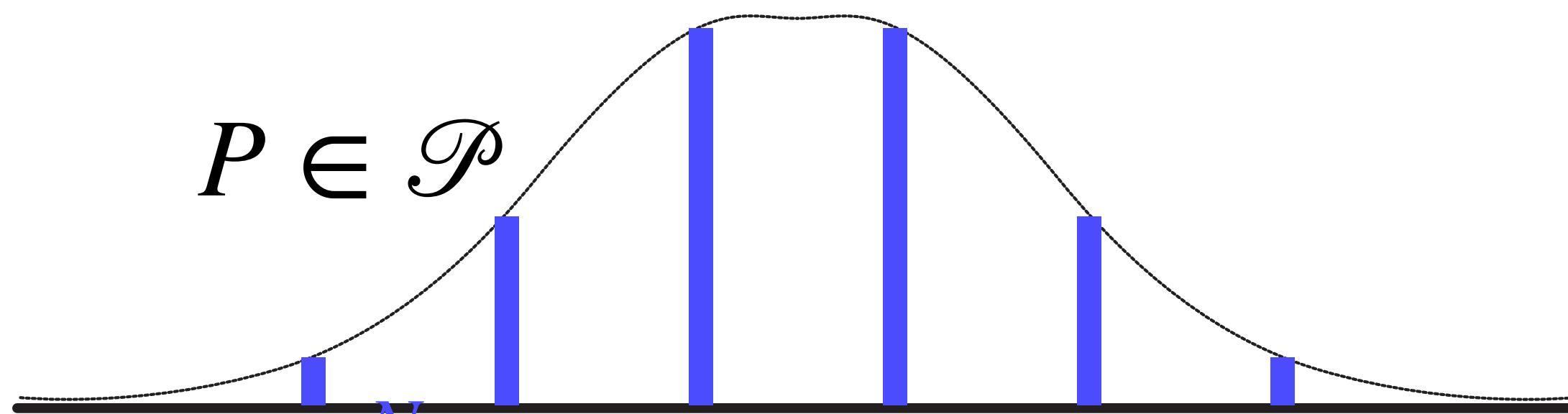
Given data samples:  $\xi_1, \xi_2, \dots, \xi_N$

$$\text{empirical data distribution } \hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$$

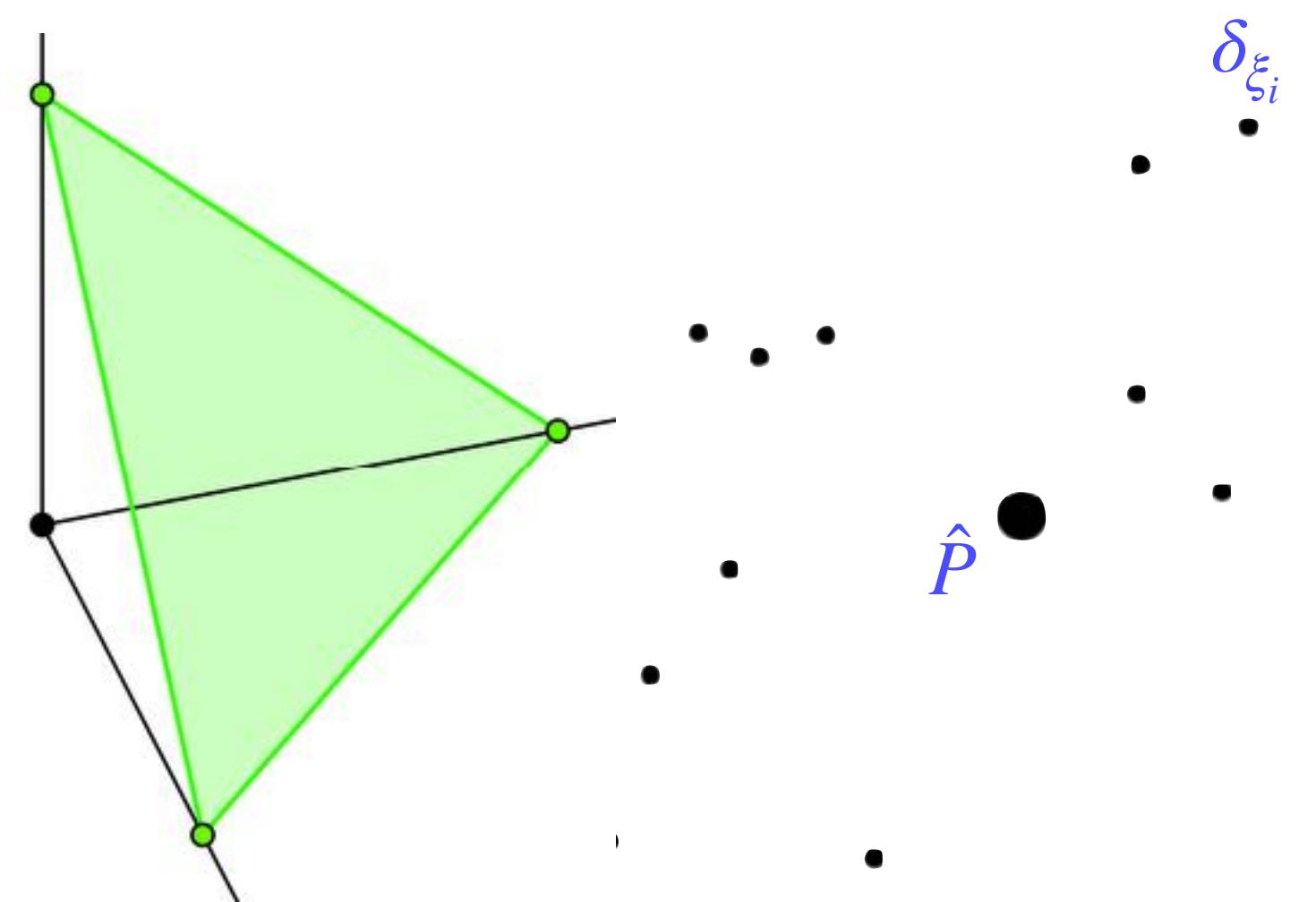


# Working with probability measures

The probability “simplex”  $\mathcal{P}$



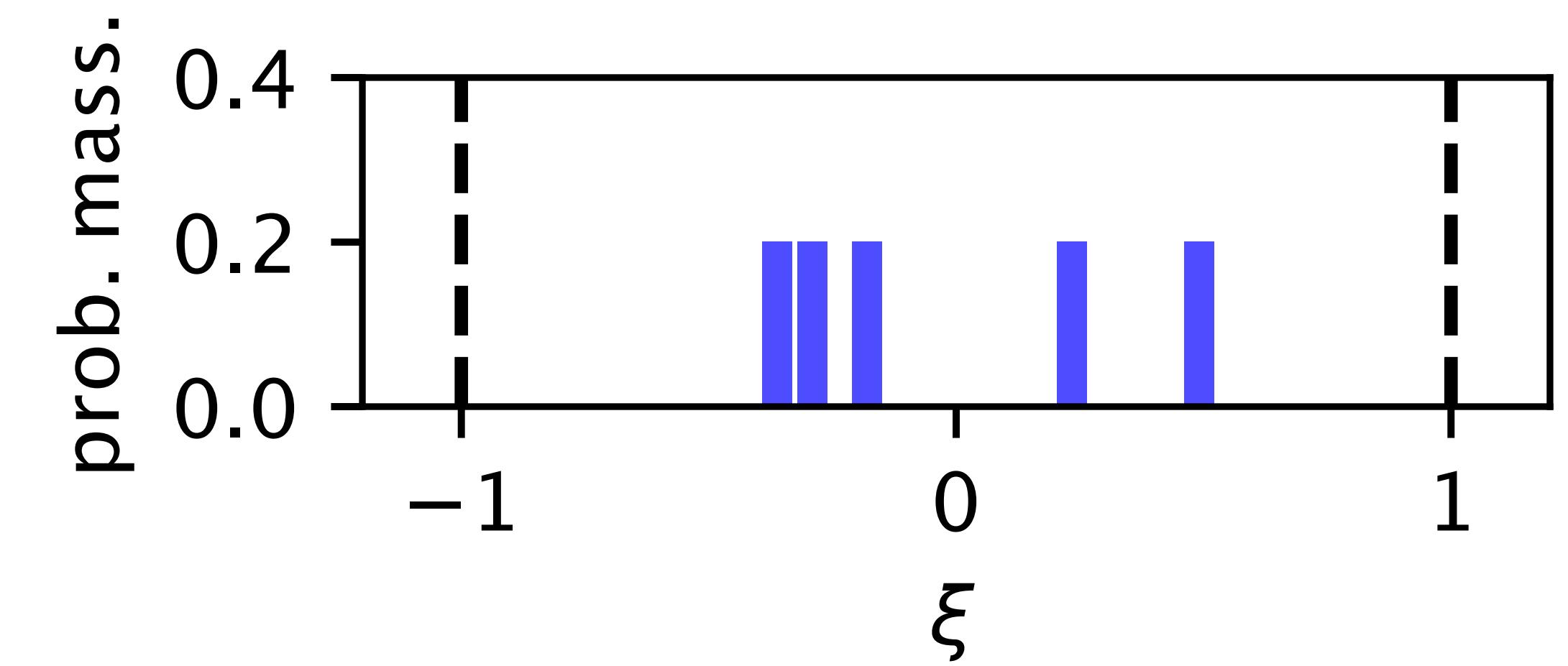
$$\sum_{i=1}^N \gamma_i \delta_{\xi_i} \quad \gamma \in \Delta_d \subset \mathbb{R}^d$$



Empirical data distribution

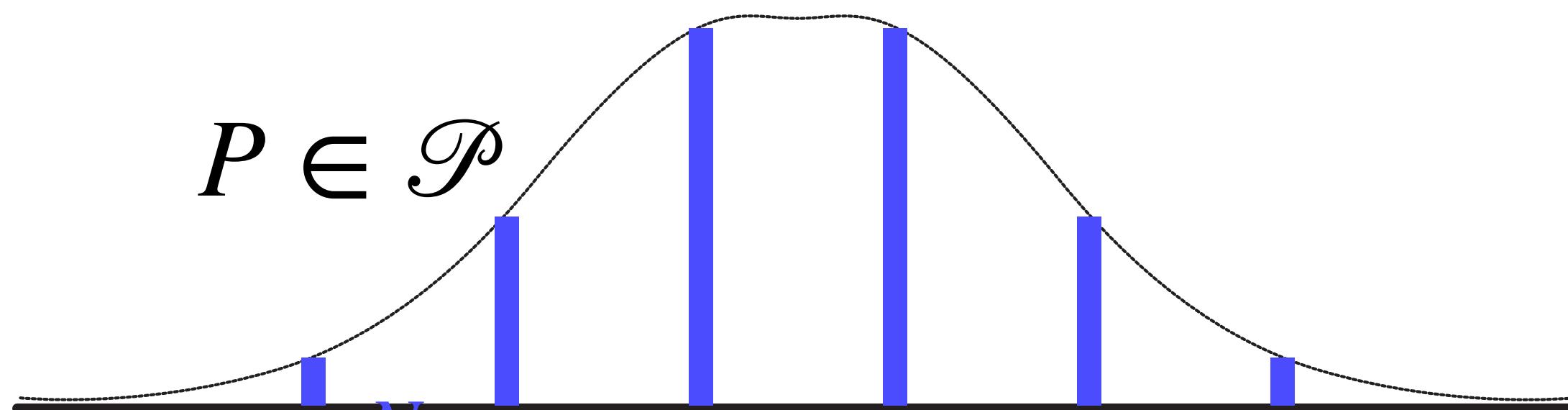
Given data samples:  $\xi_1, \xi_2, \dots, \xi_N$

$$\text{empirical data distribution } \hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$$

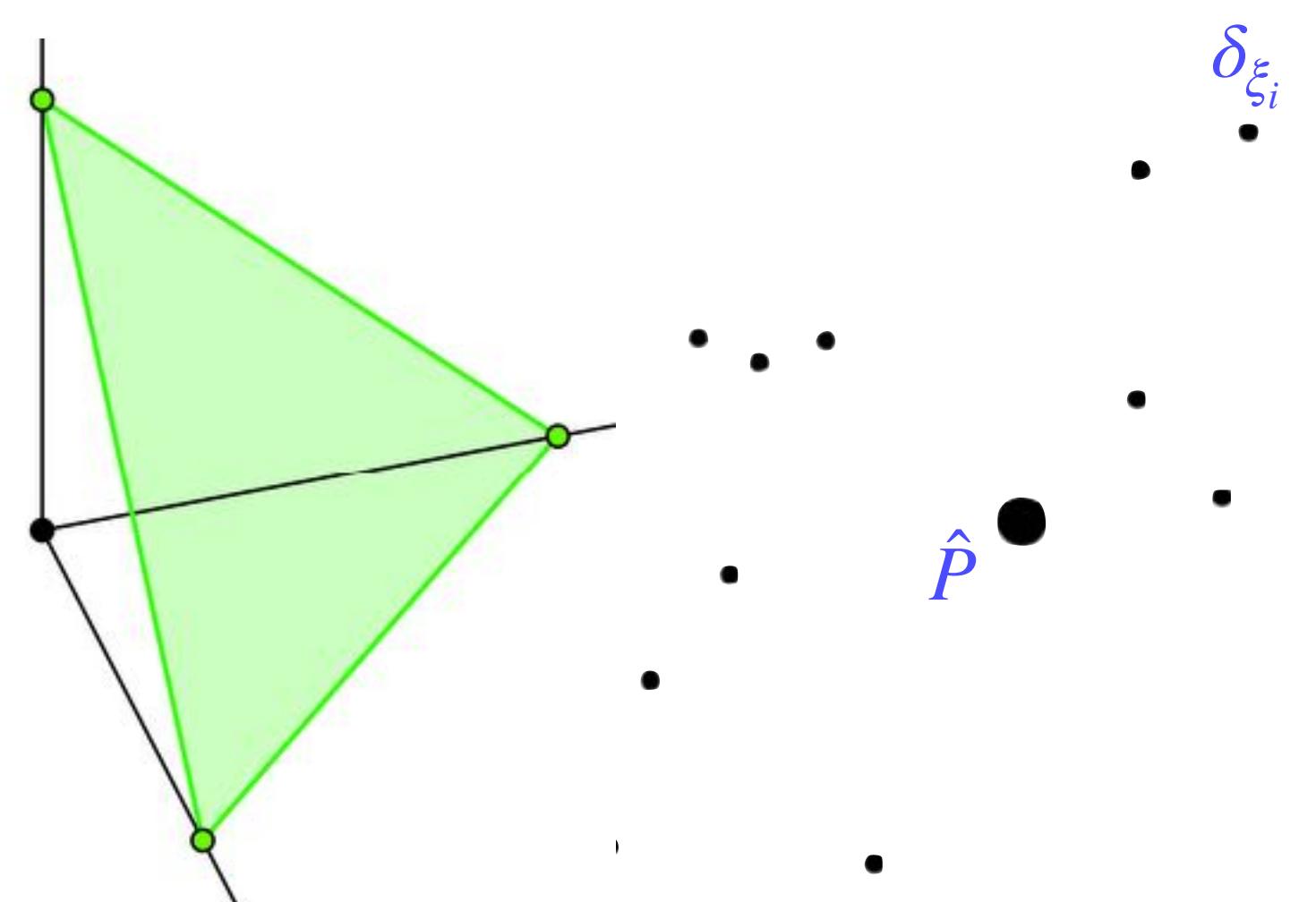


# Working with probability measures

The probability “simplex”  $\mathcal{P}$



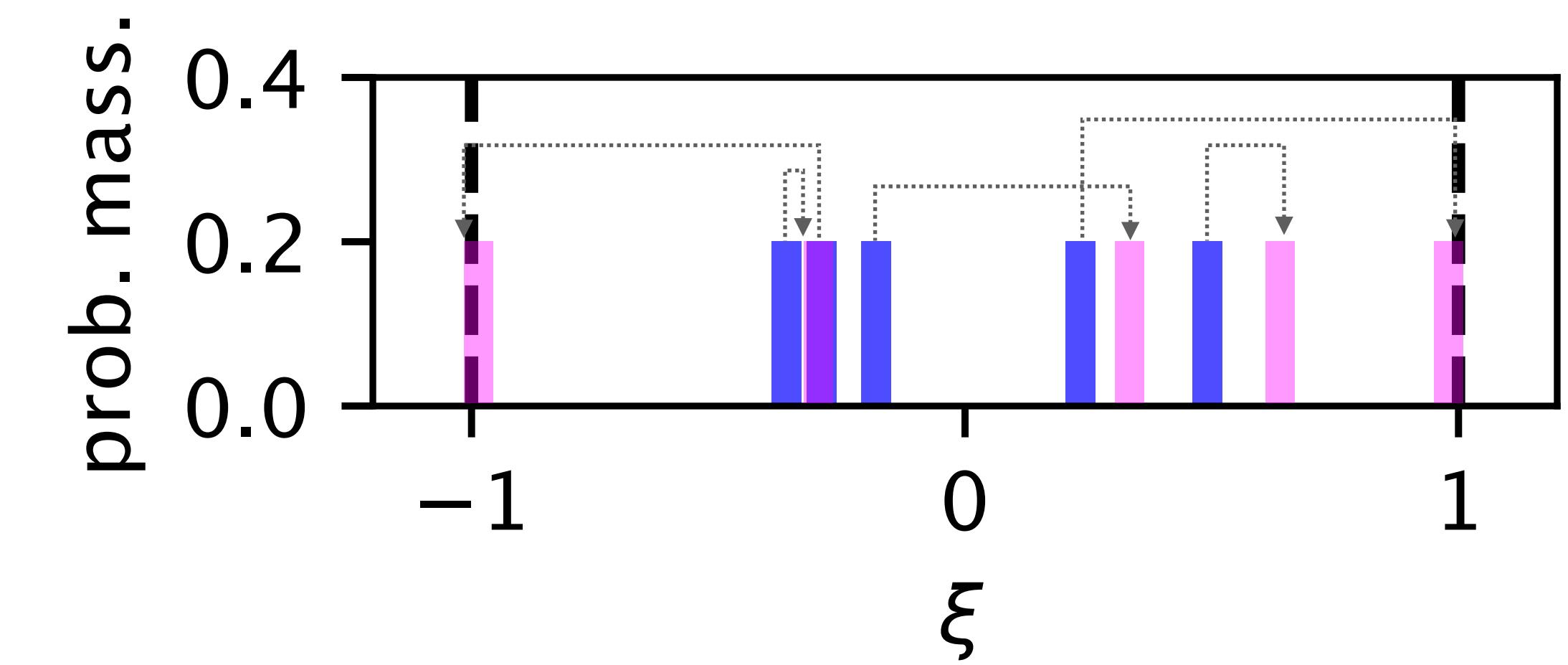
$$\sum_{i=1}^N \gamma_i \delta_{\xi_i} \quad \gamma \in \Delta_d \subset \mathbb{R}^d$$



Empirical data distribution

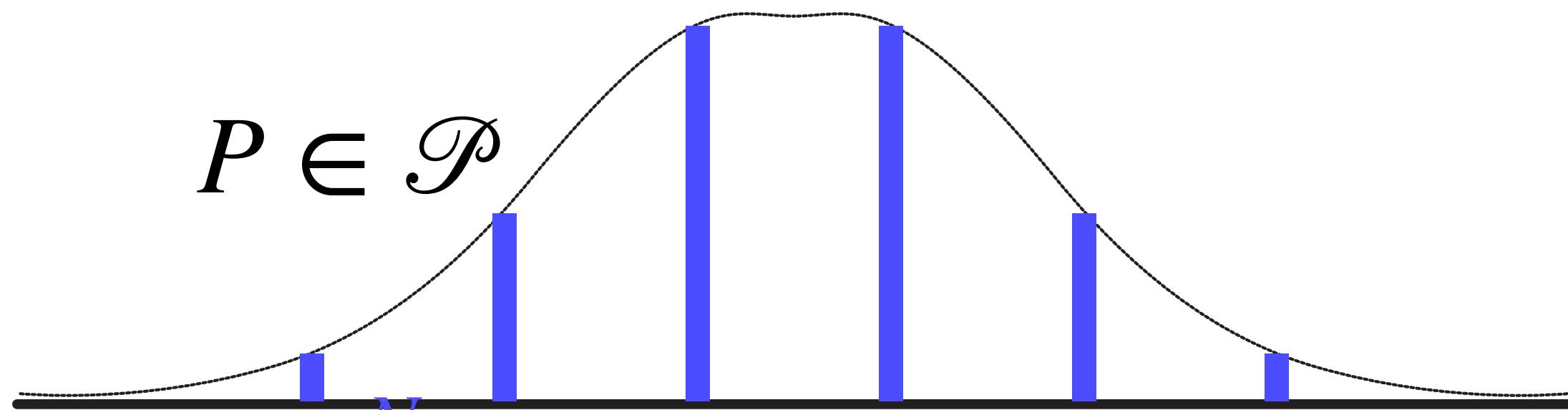
Given data samples:  $\xi_1, \xi_2, \dots, \xi_N$

$$\text{empirical data distribution } \hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$$

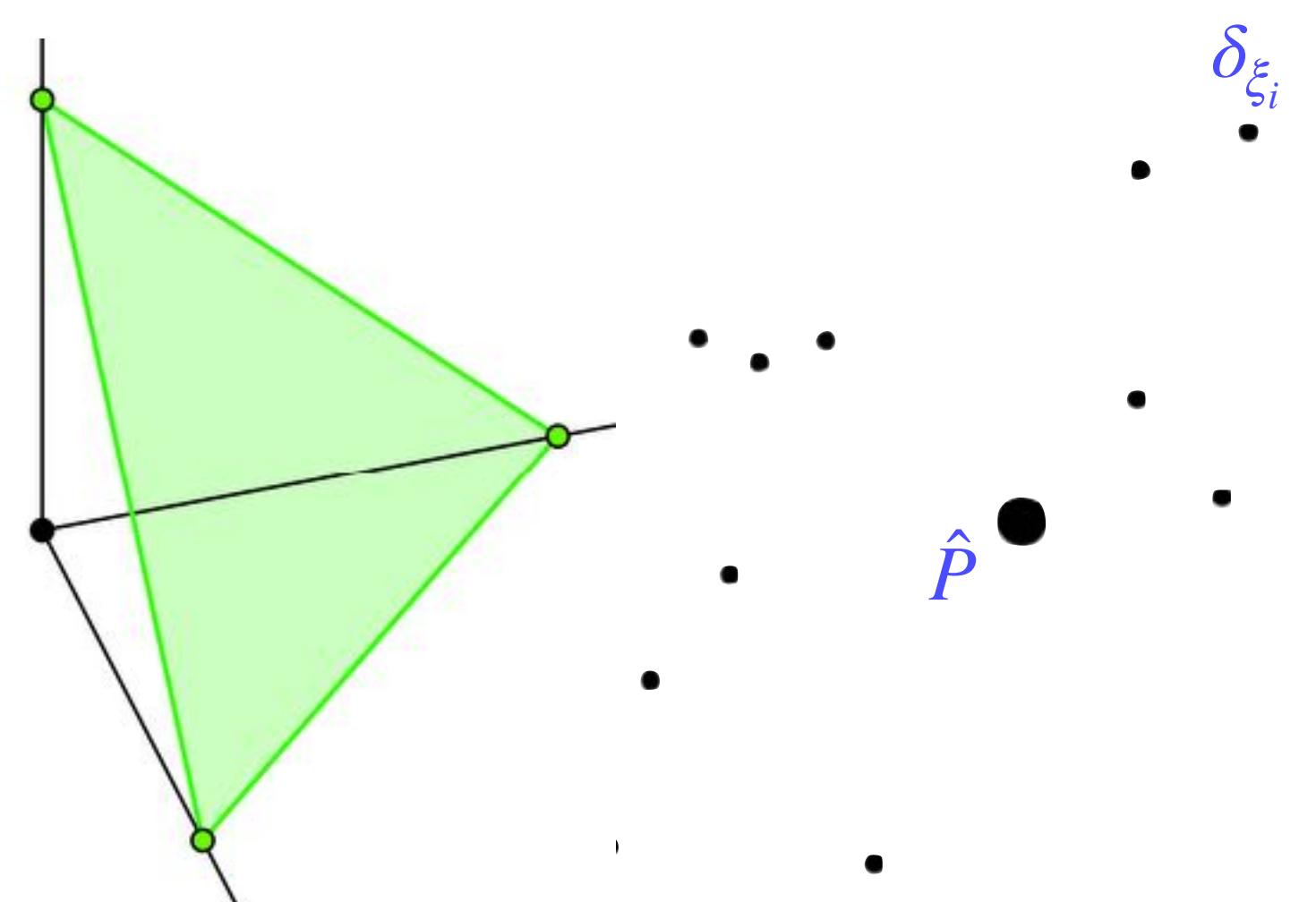


# Working with probability measures

The probability “simplex”  $\mathcal{P}$



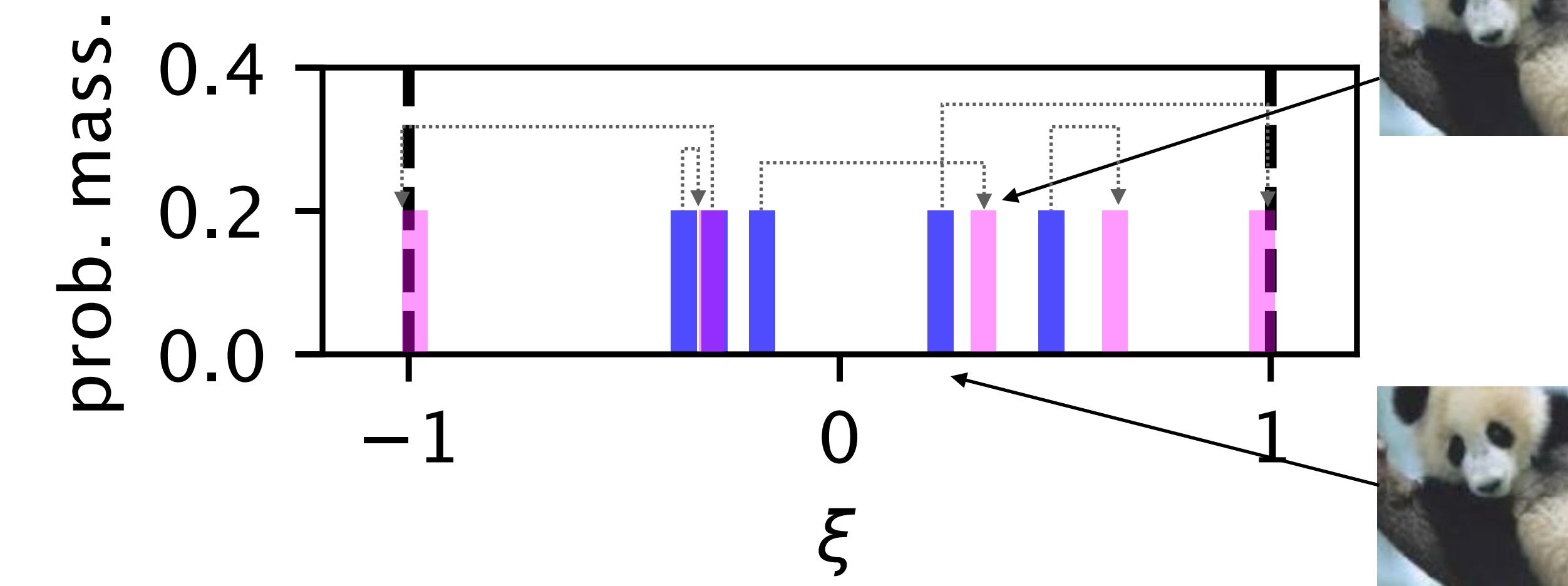
$$\sum_{i=1}^N \gamma_i \delta_{\xi_i} \quad \gamma \in \Delta_d \subset \mathbb{R}^d$$



Empirical data distribution

Given data samples:  $\xi_1, \xi_2, \dots, \xi_N$

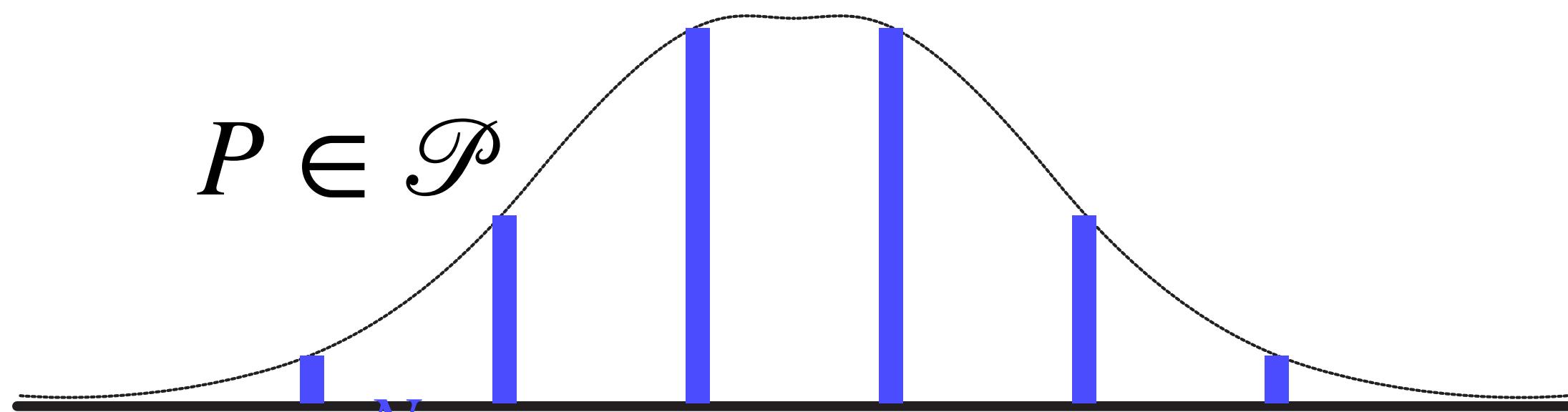
$$\text{empirical data distribution } \hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$$



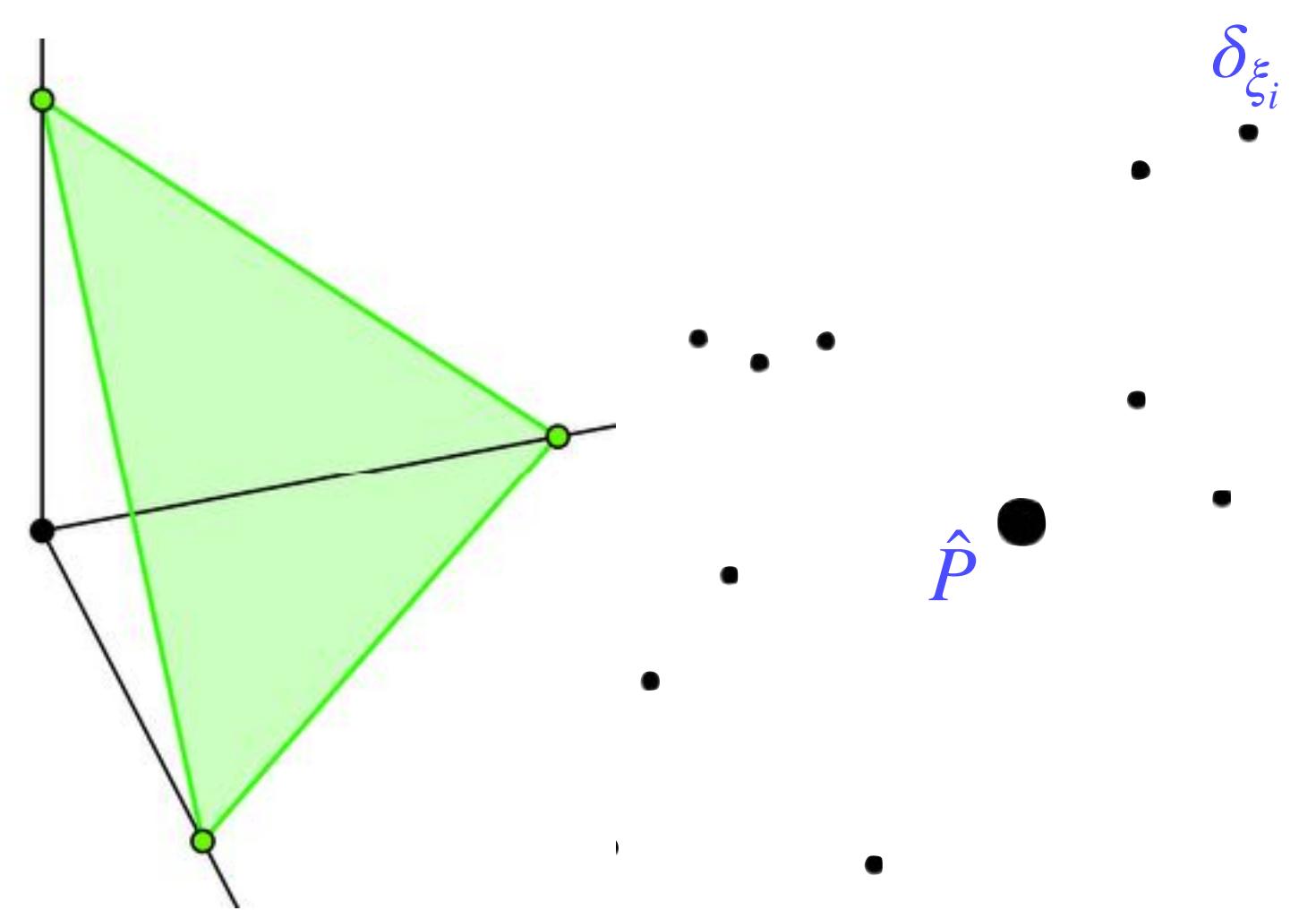
Adversarial attack: Goodfellow et al. 2015

# Working with probability measures

The probability “simplex”  $\mathcal{P}$



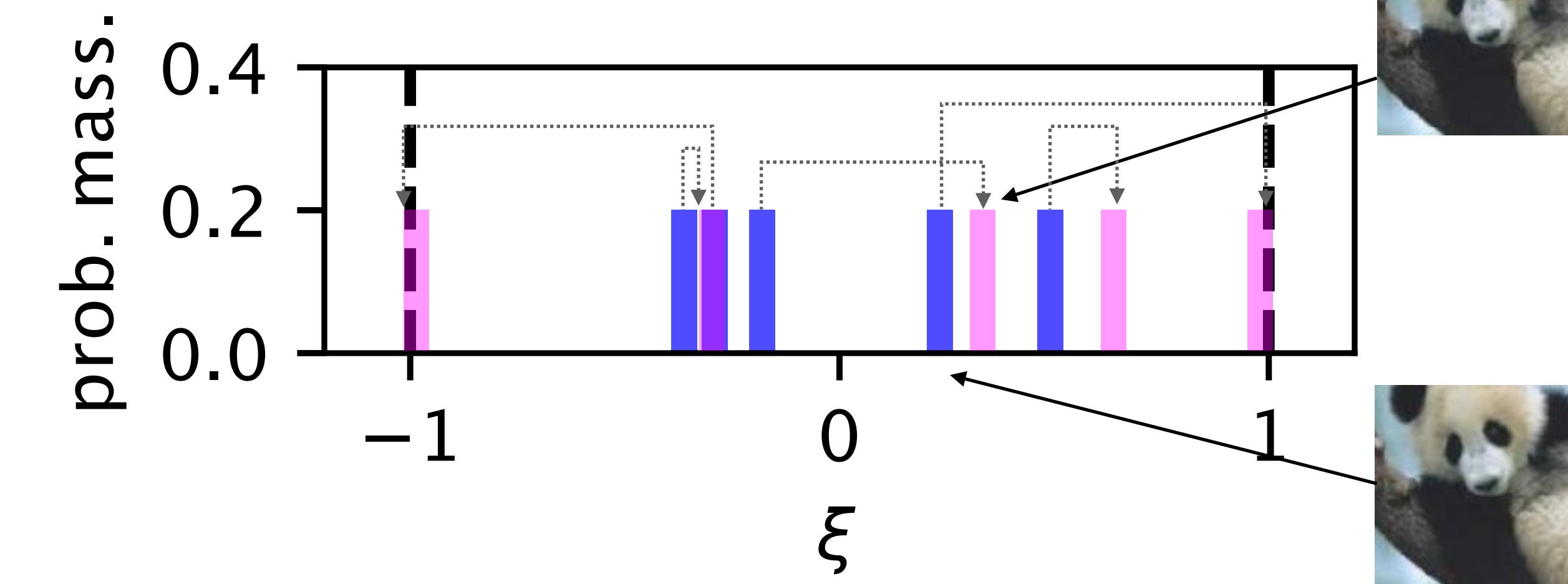
$$\sum_{i=1}^N \gamma_i \delta_{\xi_i} \quad \gamma \in \Delta_d \subset \mathbb{R}^d$$



Empirical data distribution

Given data samples:  $\xi_1, \xi_2, \dots, \xi_N$

$$\text{empirical data distribution } \hat{P} = \sum_{i=1}^N \frac{1}{N} \delta_{\xi_i}$$



Adversarial attack: Goodfellow et al. 2015

$P$  might only differ slightly from  $\hat{P}$ , but can break the system.

# Wasserstein distance

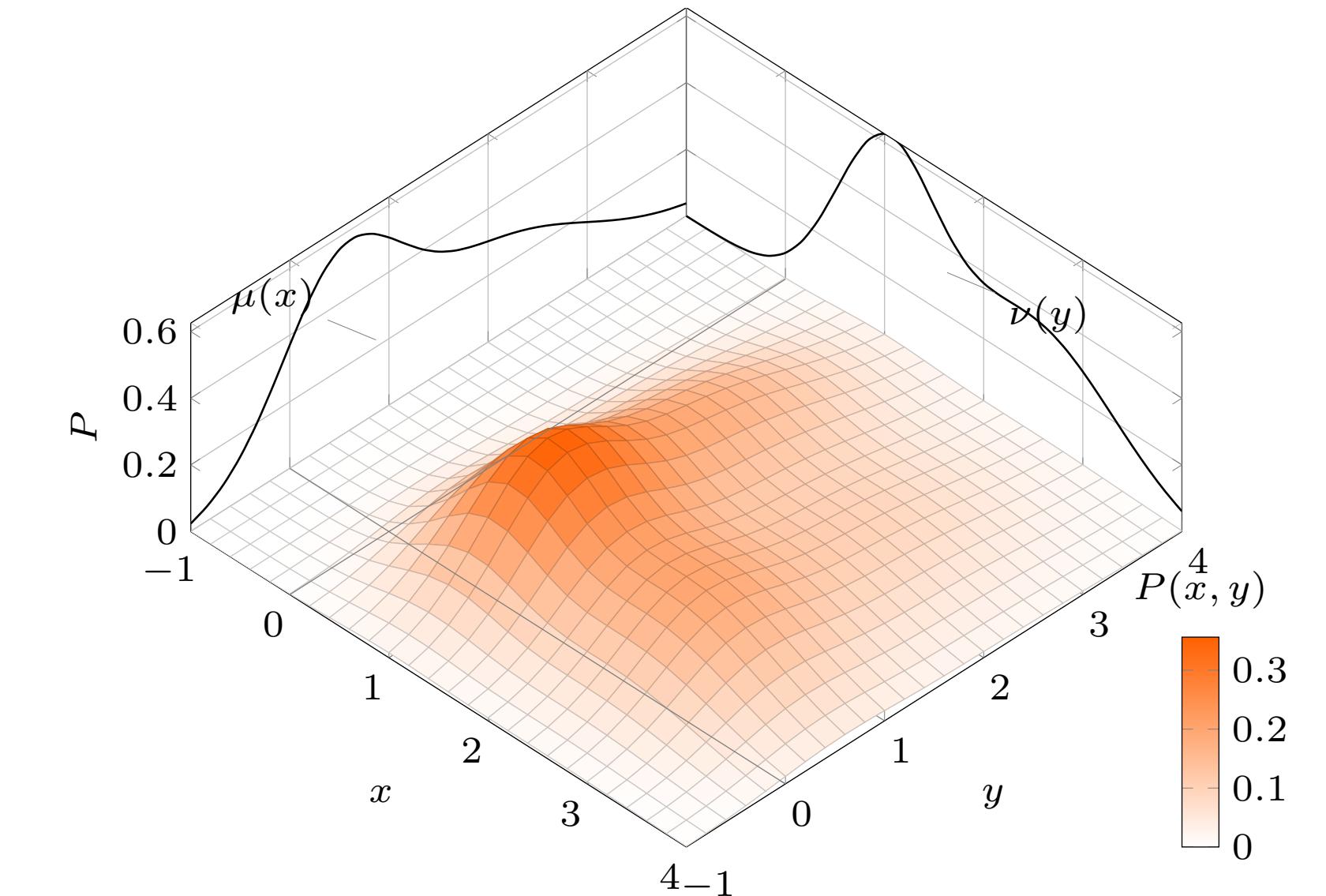
The optimal value of the Kantorovich problem

$$\mathcal{W}_p(P, Q) = \inf_{\pi \in \Pi(P, Q)} \left\{ \int \|u - v\|^p d\pi(u, v) \right\}^{1/p}$$

where  $p \geq 1$ ,  $\pi$  is a transport plan (coupling) that belongs to the set of valid plans

$$\Pi(\mu, \nu) = \{\pi \in P(X \times X) : (\pi_0)_\# \gamma = \mu, (\pi_1)_\# \gamma = \nu, \}.$$

$\pi_0$  and  $\pi_1$  are the two projections onto the marginal distributions.



# Wasserstein distance

The optimal value of the Kantorovich problem

$$\mathcal{W}_p(P, Q) = \inf_{\pi \in \Pi(P, Q)} \left\{ \int \|u - v\|^p d\pi(u, v) \right\}^{1/p}$$

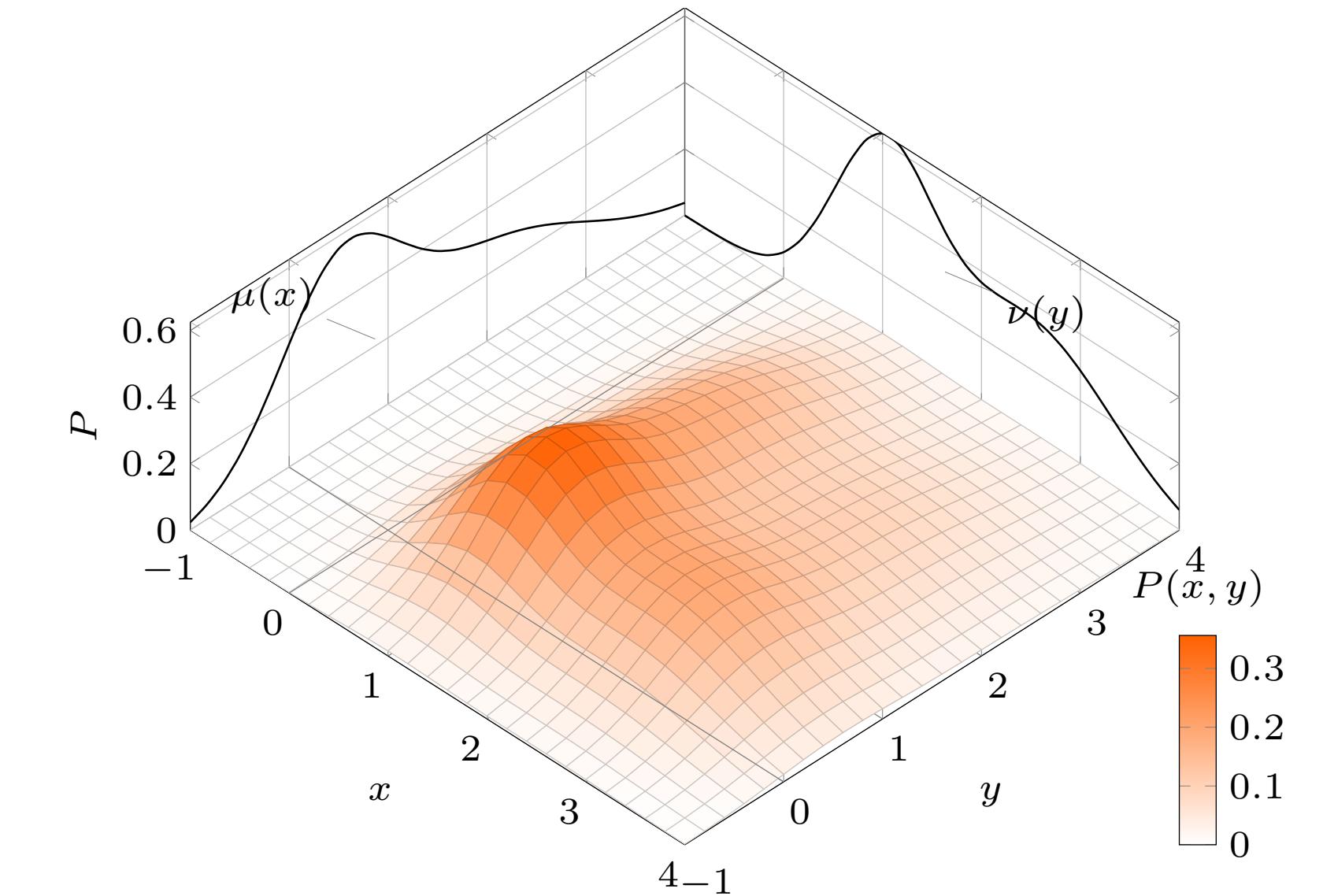
where  $p \geq 1$ ,  $\pi$  is a transport plan (coupling) that belongs to the set of valid plans

$$\Pi(\mu, \nu) = \{\pi \in P(X \times X) : (\pi_0)_\# \gamma = \mu, (\pi_1)_\# \gamma = \nu, \}.$$

$\pi_0$  and  $\pi_1$  are the two projections onto the marginal distributions.

If  $p = 1$ , this is equivalent to IPM with 1-Lipschitz family

$$W_1(P, Q) = \sup_{\text{lip}(f) \leq 1} \int f d(P - Q)$$



# Wasserstein distance

The optimal value of the Kantorovich problem

$$\mathcal{W}_p(P, Q) = \inf_{\pi \in \Pi(P, Q)} \left\{ \int \|u - v\|^p d\pi(u, v) \right\}^{1/p}$$

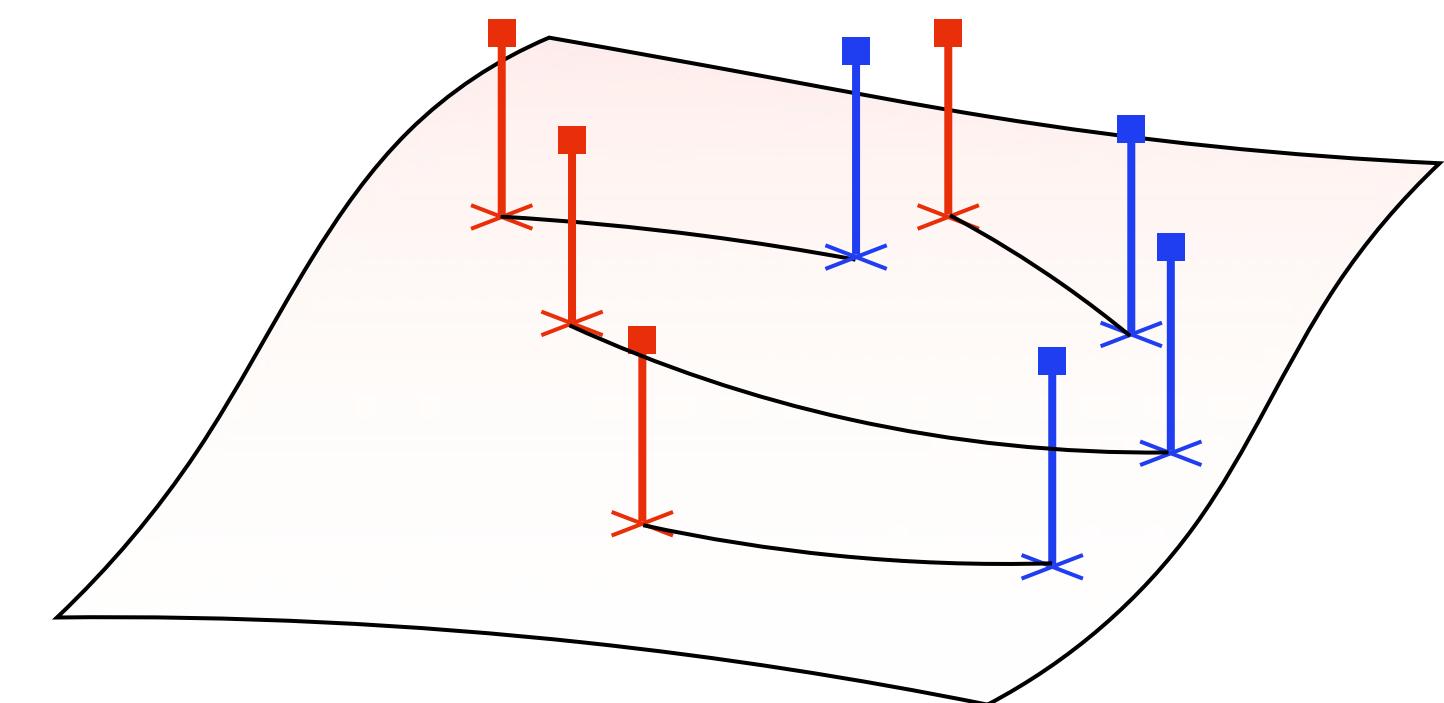
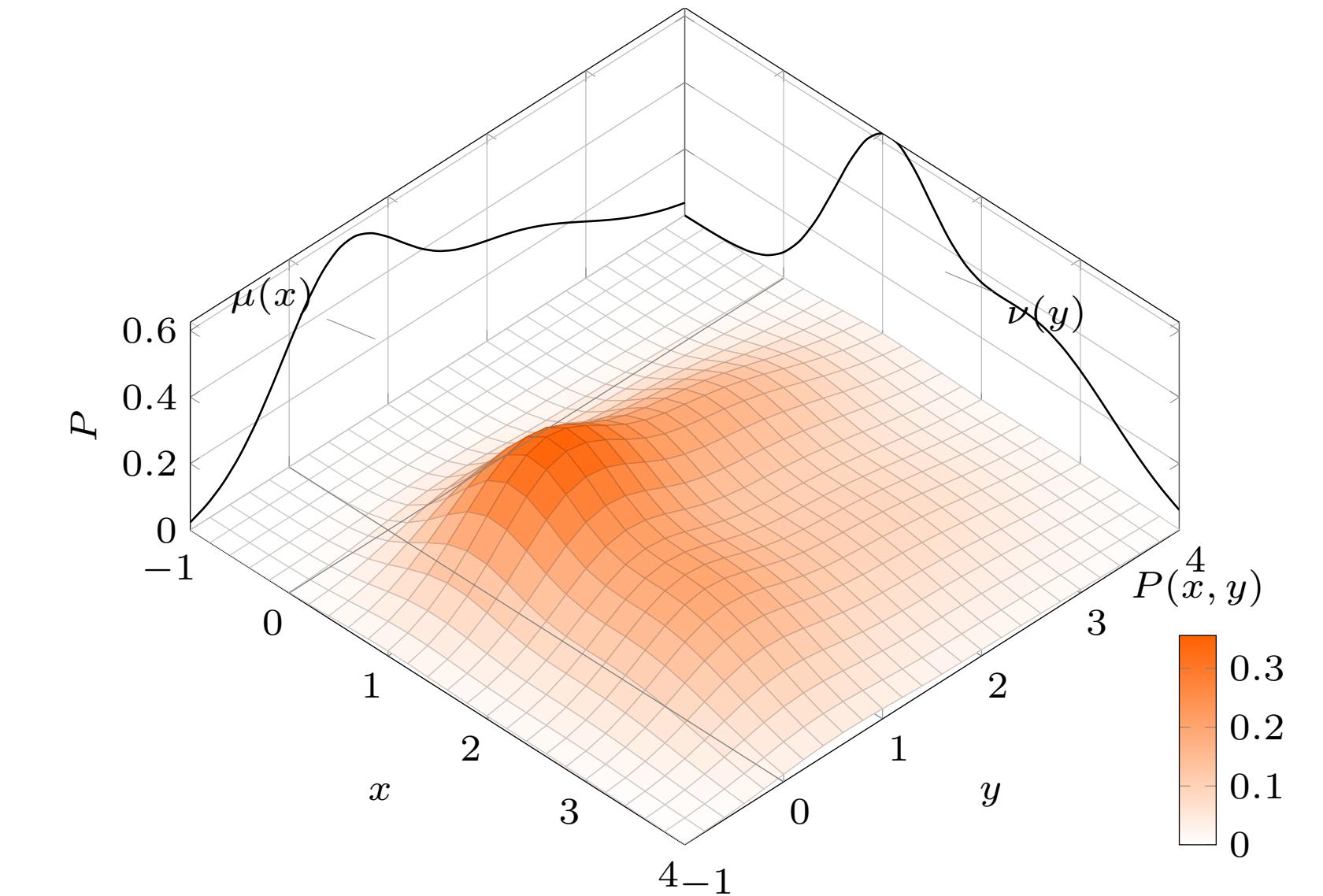
where  $p \geq 1$ ,  $\pi$  is a transport plan (coupling) that belongs to the set of valid plans

$$\Pi(\mu, \nu) = \{\pi \in P(X \times X) : (\pi_0)_\# \gamma = \mu, (\pi_1)_\# \gamma = \nu, \}.$$

$\pi_0$  and  $\pi_1$  are the two projections onto the marginal distributions.

If  $p = 1$ , this is equivalent to IPM with 1-Lipschitz family

$$W_1(P, Q) = \sup_{\text{lip}(f) \leq 1} \int f d(P - Q)$$



# Application to generative modeling: Wasserstein auto-encoder (WAE)

# Application to generative modeling: Wasserstein auto-encoder (WAE)

Minimize the regularized reconstruction loss

$$\inf_{Q(Z|X) \in \mathcal{Q}} \mathbb{E}_{P_X} \mathbb{E}_{Q(Z|X)} [c(X, G(Z))] + \lambda \cdot \mathcal{D}_Z(Q_Z, P_Z).$$

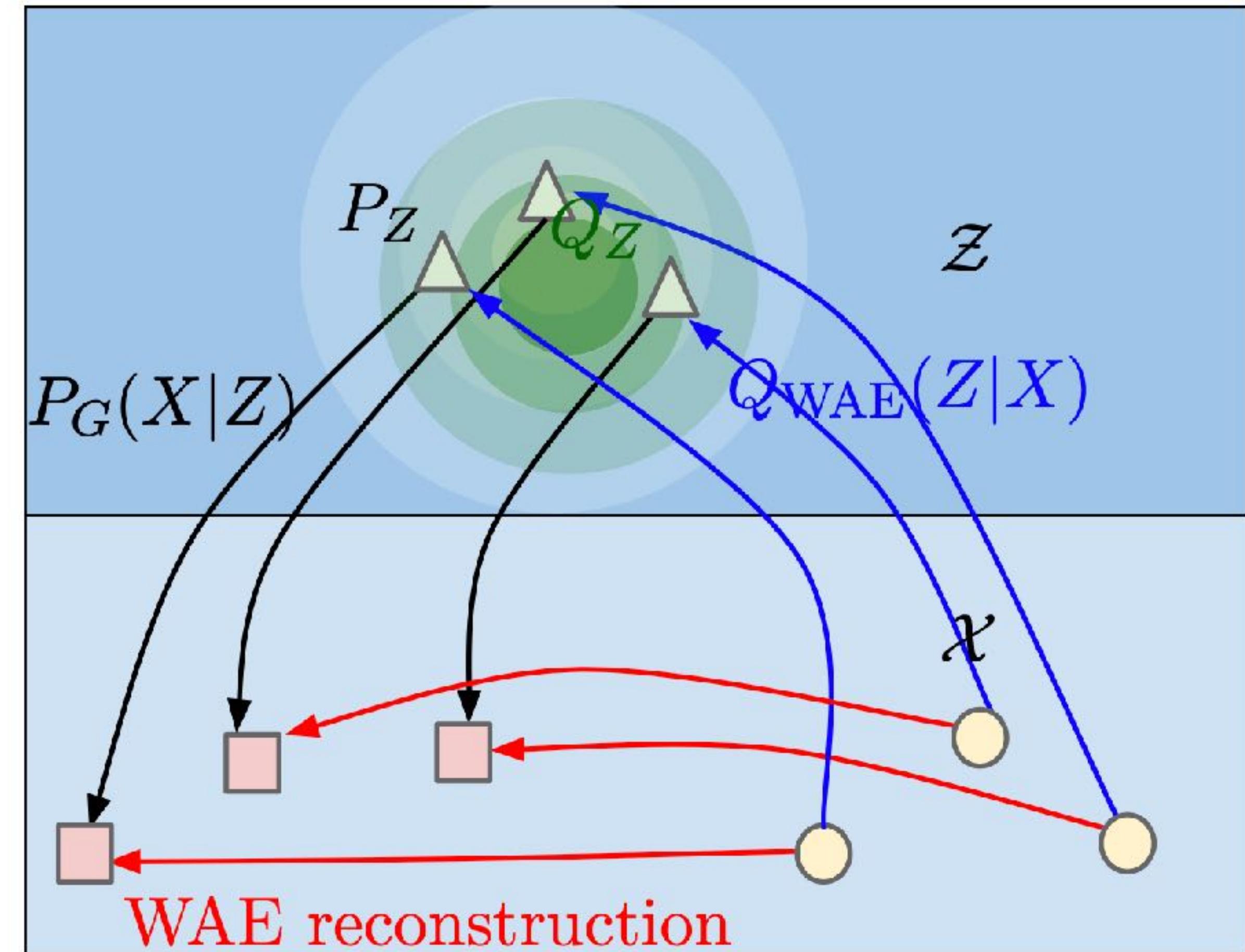
- the posterior distribution  $Q(Z|X)$  is the encoder,  $X$  is data
- $\mathcal{D}_Z$  is some discrepancy measure in the latent space, e.g., MMD
- $P_Z$  is the fixed target latent distribution, a.k.a. prior, such as a Gaussian or Uniform.
- In practice, we take the marginal distribution  $P_X$  to be the empirical data distribution  $\frac{1}{N} \sum_{i=1}^N \delta_{x_i}$ .

# Application to generative modeling: Wasserstein auto-encoder (WAE)

Minimize the regularized reconstruction loss

$$\inf_{Q(Z|X) \in \mathcal{Q}} \mathbb{E}_{P_X} \mathbb{E}_{Q(Z|X)} [c(X, G(Z))] + \lambda \cdot \mathcal{D}_Z(Q_Z, P_Z).$$

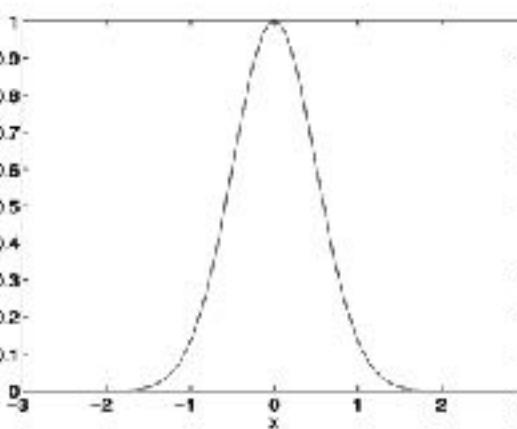
- the posterior distribution  $Q(Z|X)$  is the encoder,  $X$  is data
- $\mathcal{D}_Z$  is some discrepancy measure in the latent space, e.g., MMD
- $P_Z$  is the fixed target latent distribution, a.k.a. prior, such as a Gaussian or Uniform.
- In practice, we take the marginal distribution  $P_X$  to be the empirical data distribution  $\frac{1}{N} \sum_{i=1}^N \delta_{x_i}$ .



# Learning with kernels (the modern way)

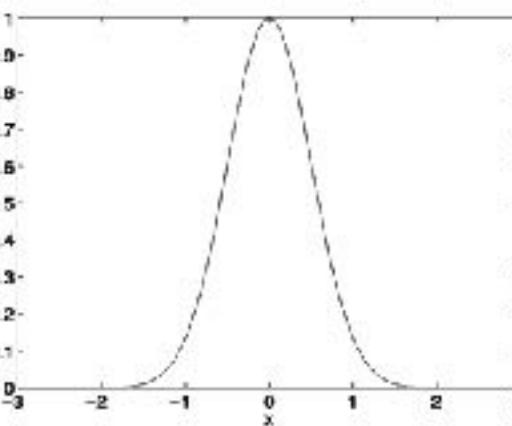
# Learning with kernels (the modern way)

- A kernel is a symmetric function  
 $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , e.g., Gaussian kernel  
$$k(x, x') = \exp\left(-\|x - x'\|_2^2 / 2\sigma^2\right).$$



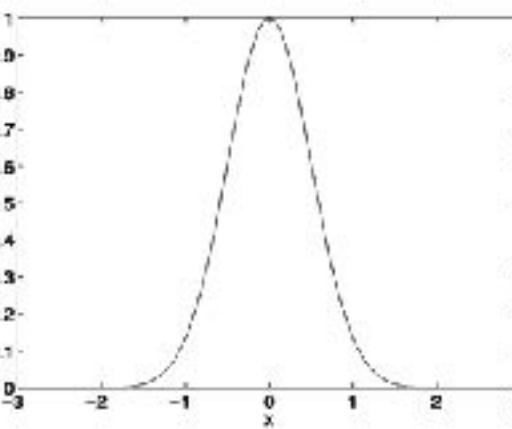
# Learning with kernels (the modern way)

- A kernel is a symmetric function  
 $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , e.g., Gaussian kernel  
$$k(x, x') = \exp\left(-\|x - x'\|_2^2 / 2\sigma^2\right).$$
- A p.d.  $k$  corresponds to a Hilbert space  $\mathcal{H}$  (RKHS), which satisfies the **reproducing property**  
$$f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \forall f \in \mathcal{H}, x \in \mathcal{X},$$
  
$$\phi(x) := k(x, \cdot)$$
 is the **canonical feature** of  $\mathcal{H}$ .



# Learning with kernels (the modern way)

- A kernel is a symmetric function  
 $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , e.g., Gaussian kernel  
$$k(x, x') = \exp\left(-\|x - x'\|_2^2 / 2\sigma^2\right).$$
- A p.d.  $k$  corresponds to a Hilbert space  $\mathcal{H}$  (RKHS), which satisfies the **reproducing property**  
$$f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \forall f \in \mathcal{H}, x \in \mathcal{X},$$
  
$$\phi(x) := k(x, \cdot)$$
 is the **canonical feature** of  $\mathcal{H}$ .



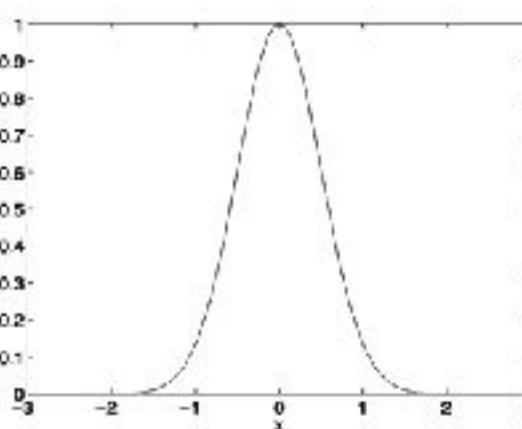
$\sim P$   
 $\sim Q$



$\mathcal{P}$

# Learning with kernels (the modern way)

- A kernel is a symmetric function  
 $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , e.g., Gaussian kernel  
$$k(x, x') = \exp\left(-\|x - x'\|_2^2 / 2\sigma^2\right).$$
- A p.d.  $k$  corresponds to a Hilbert space  $\mathcal{H}$  (RKHS), which satisfies the **reproducing property**  
$$f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \forall f \in \mathcal{H}, x \in \mathcal{X},$$
  
$$\phi(x) := k(x, \cdot)$$
 is the **canonical feature** of  $\mathcal{H}$ .

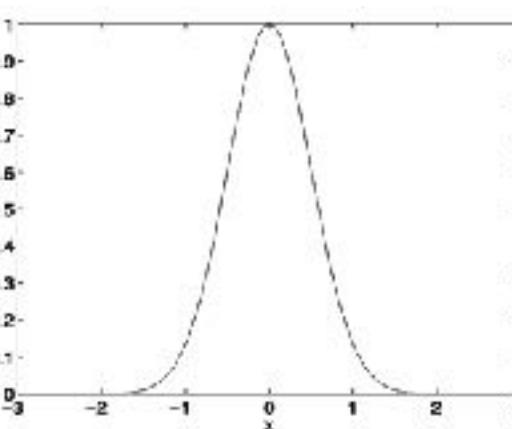


$$\begin{array}{c} \sim P \\ \longrightarrow \\ \sim Q \end{array}$$

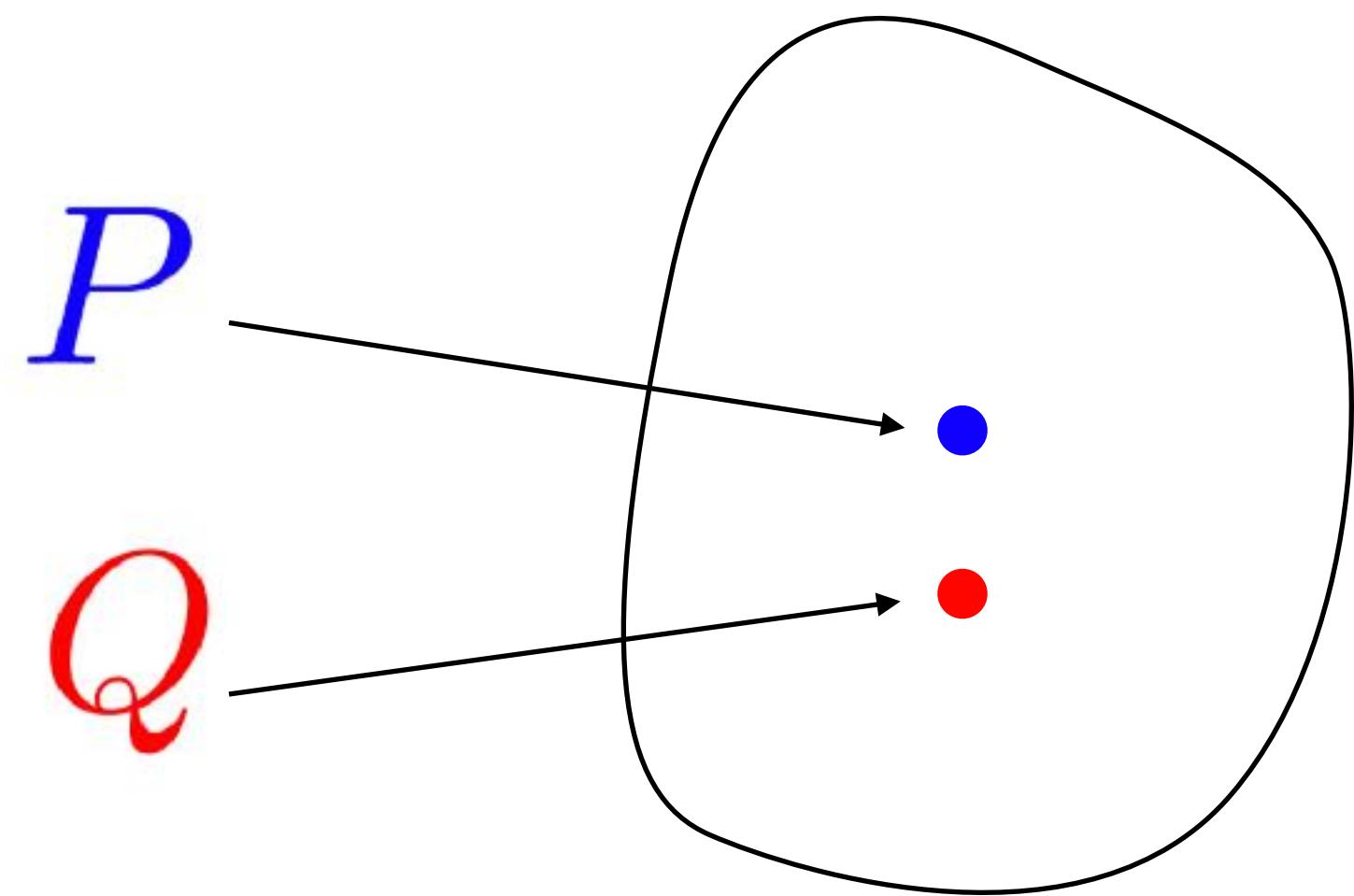
$$\mathcal{P} \quad \mu_P := \int \phi \, dP$$

# Learning with kernels (the modern way)

- A kernel is a symmetric function  
 $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , e.g., Gaussian kernel  
$$k(x, x') = \exp\left(-\|x - x'\|_2^2 / 2\sigma^2\right).$$
- A p.d.  $k$  corresponds to a Hilbert space  $\mathcal{H}$  (RKHS), which satisfies the **reproducing property**  
$$f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \forall f \in \mathcal{H}, x \in \mathcal{X},$$
  
$$\phi(x) := k(x, \cdot)$$
 is the **canonical feature** of  $\mathcal{H}$ .



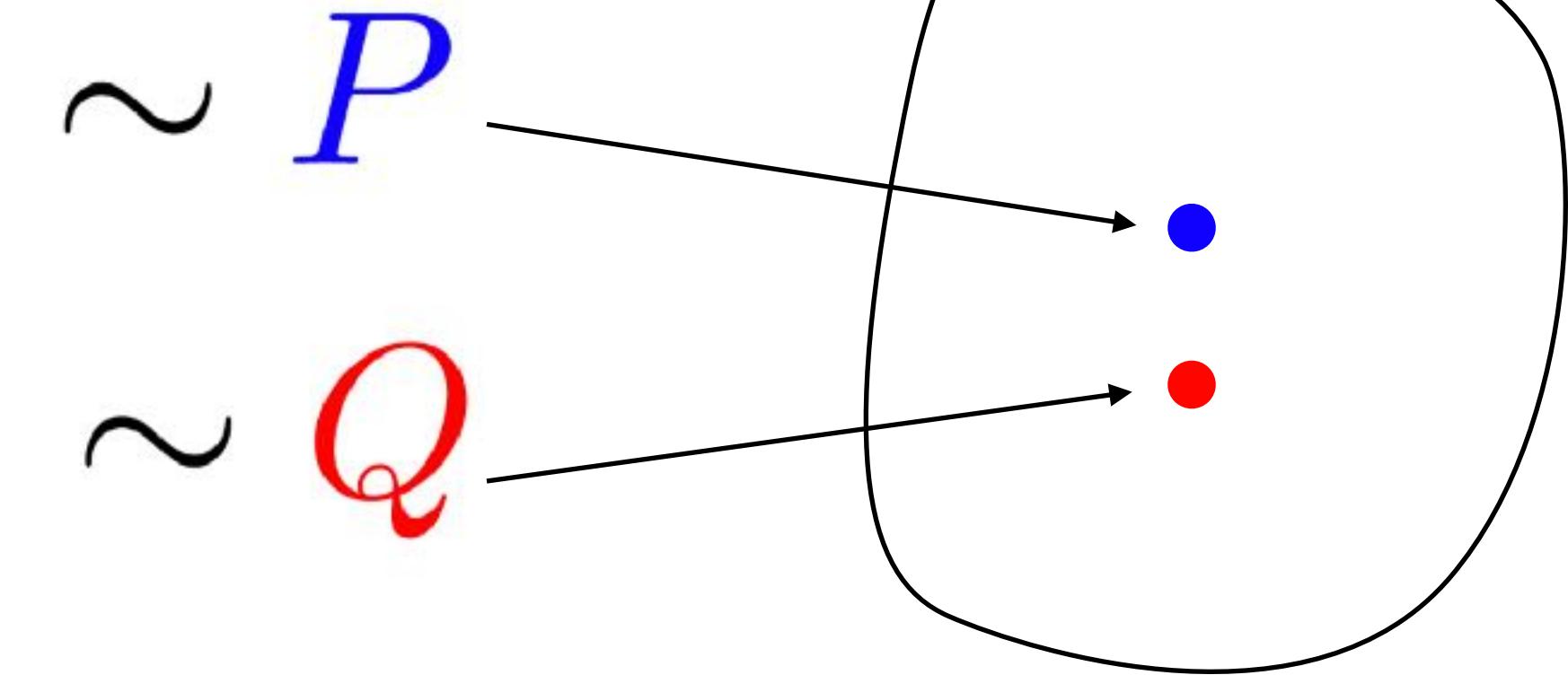
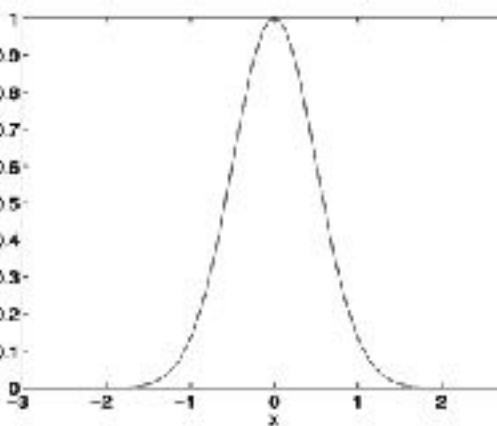
$$\sim P \quad \sim Q$$



$$\mathcal{P} \quad \mu_P := \int \phi \, dP \quad \mathcal{H}$$

# Learning with kernels (the modern way)

- A kernel is a symmetric function  
 $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , e.g., Gaussian kernel  
 $k(x, x') = \exp(-\|x - x'\|_2^2 / 2\sigma^2).$
- A p.d.  $k$  corresponds to a Hilbert space  $\mathcal{H}$  (RKHS), which satisfies the **reproducing property**  
 $f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \forall f \in \mathcal{H}, x \in \mathcal{X},$   
 $\phi(x) := k(x, \cdot)$  is the **canonical feature** of  $\mathcal{H}$ .



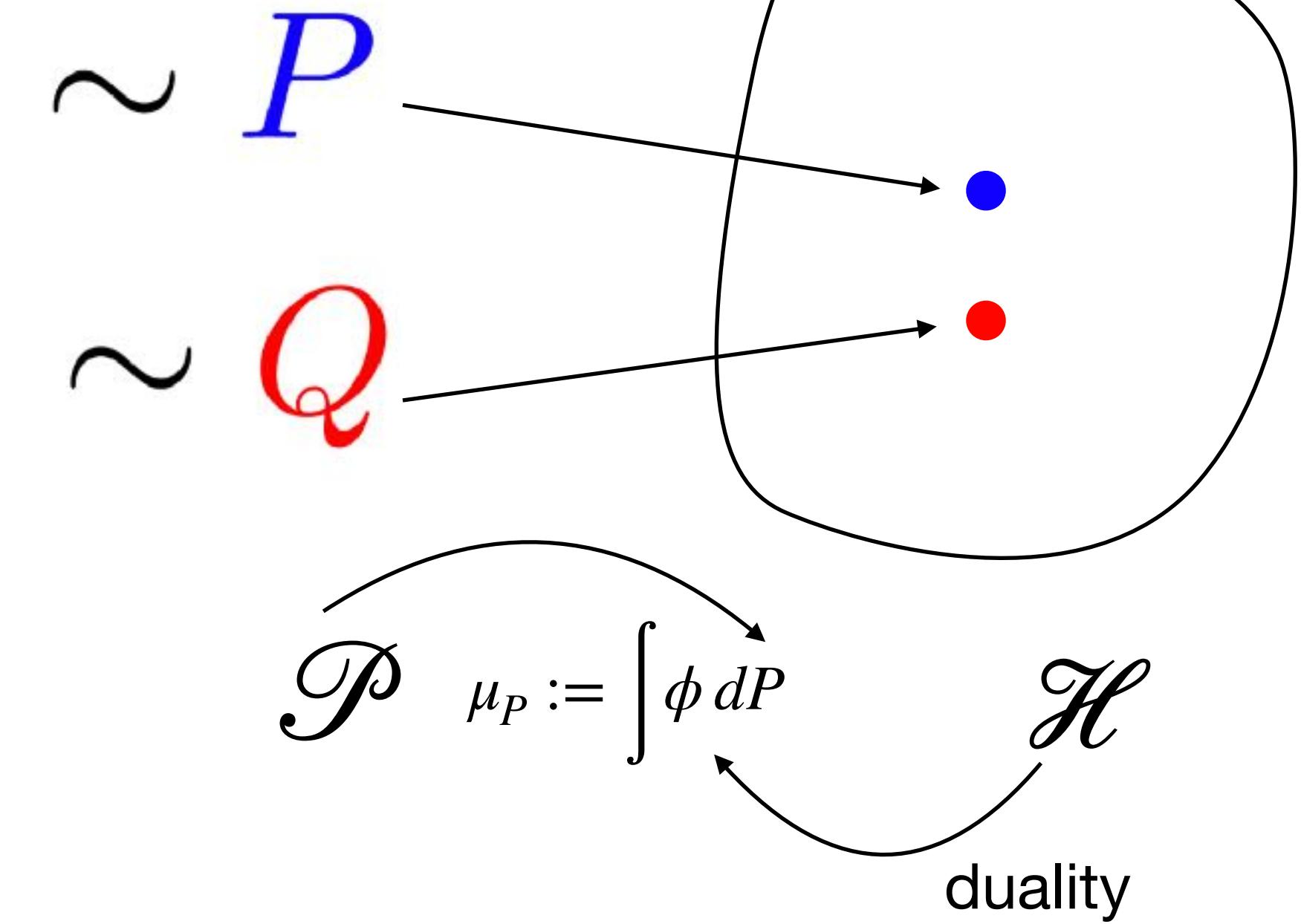
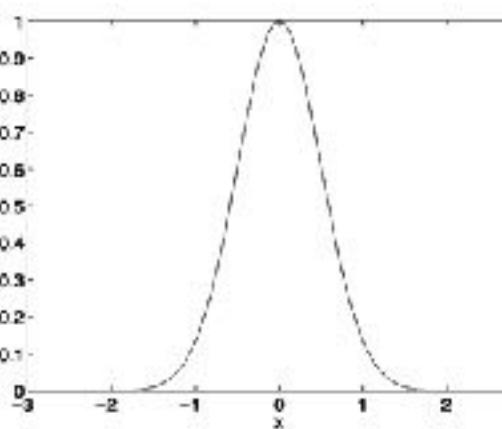
$$\mathcal{P} \quad \mu_P := \int \phi \, dP \qquad \mathcal{H}$$

$\mu := \int \phi \, dP$  is the (*kernel*) **mean embedding** of  $P$  in  $\mathcal{H}$ .

$\mu$  can be viewed as a generalized moment vector  
e.g., let  $\phi(x) = [x, x^2]^\top$

# Learning with kernels (the modern way)

- A kernel is a symmetric function  
 $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , e.g., Gaussian kernel  
 $k(x, x') = \exp(-\|x - x'\|_2^2 / 2\sigma^2).$
- A p.d.  $k$  corresponds to a Hilbert space  $\mathcal{H}$  (RKHS), which satisfies the **reproducing property**  
 $f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \forall f \in \mathcal{H}, x \in \mathcal{X},$   
 $\phi(x) := k(x, \cdot)$  is the **canonical feature** of  $\mathcal{H}$ .

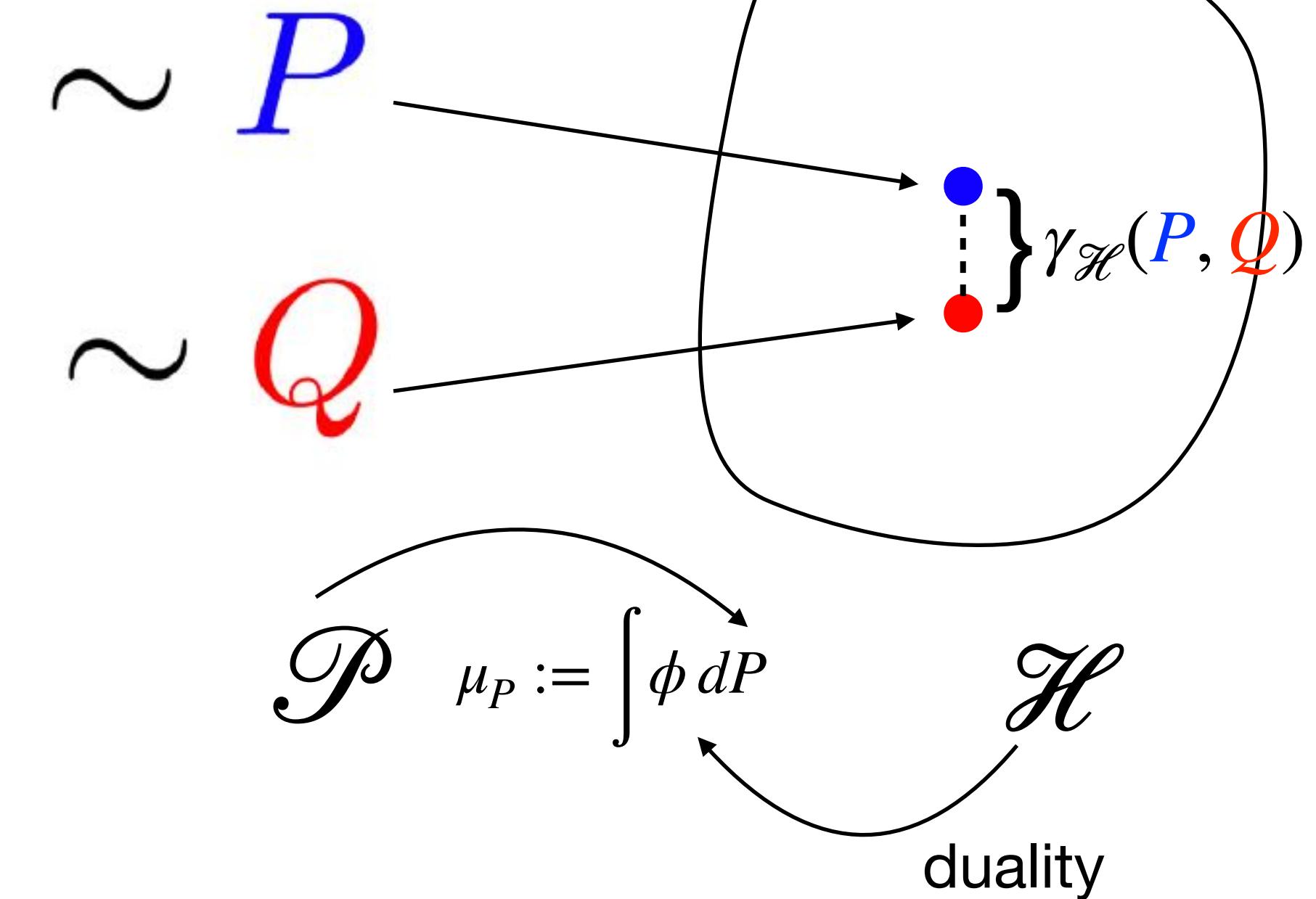
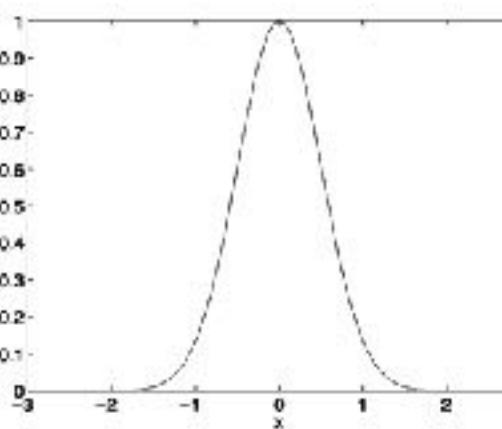


$\mu := \int \phi \, dP$  is the (*kernel*) **mean embedding** of  $P$  in  $\mathcal{H}$ .

$\mu$  can be viewed as a generalized moment vector  
e.g., let  $\phi(x) = [x, x^2]^\top$

# Learning with kernels (the modern way)

- A kernel is a symmetric function  
 $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , e.g., Gaussian kernel  
 $k(x, x') = \exp(-\|x - x'\|_2^2 / 2\sigma^2).$
- A p.d.  $k$  corresponds to a Hilbert space  $\mathcal{H}$  (RKHS), which satisfies the **reproducing property**  
 $f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \forall f \in \mathcal{H}, x \in \mathcal{X},$   
 $\phi(x) := k(x, \cdot)$  is the **canonical feature** of  $\mathcal{H}$ .

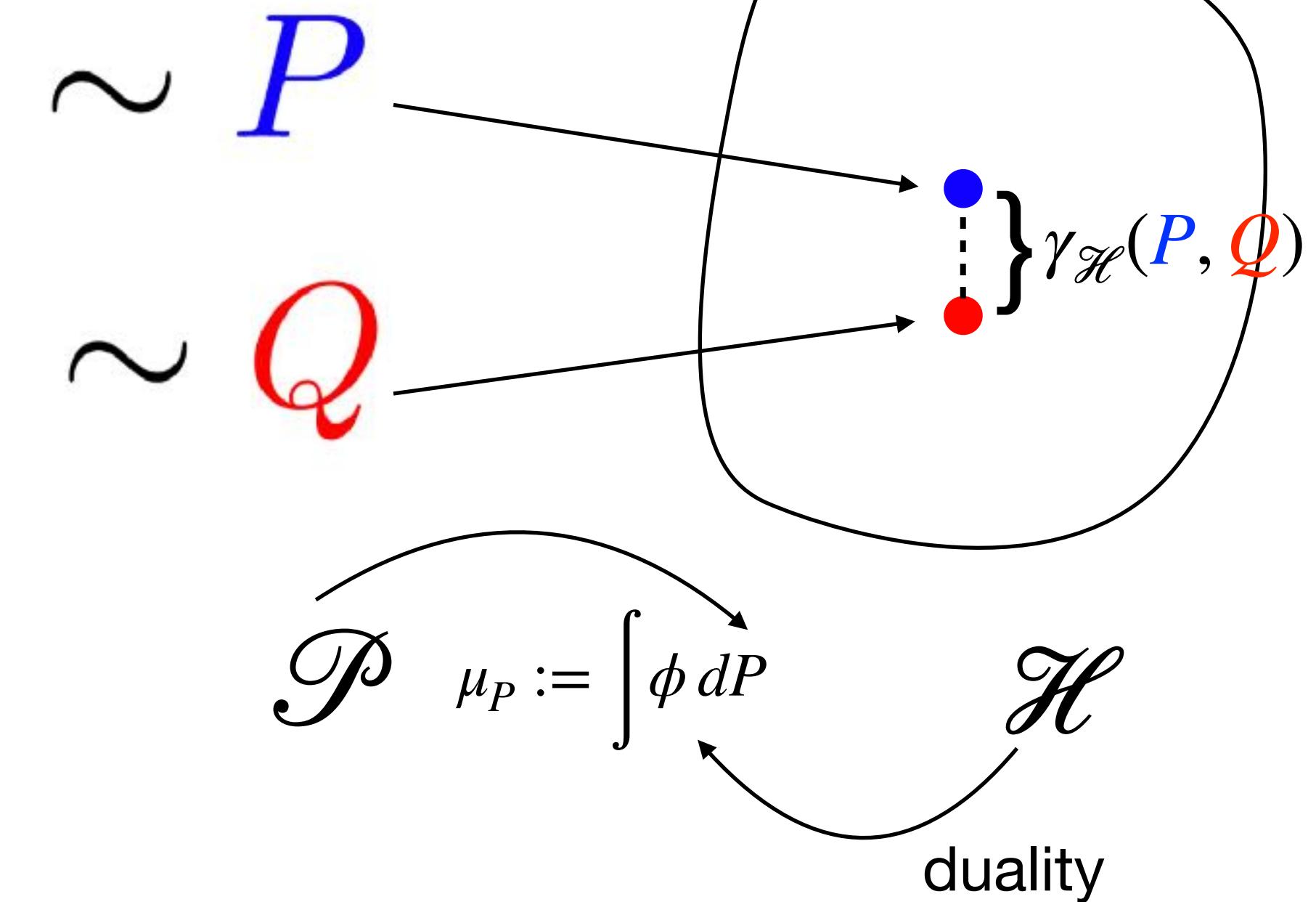
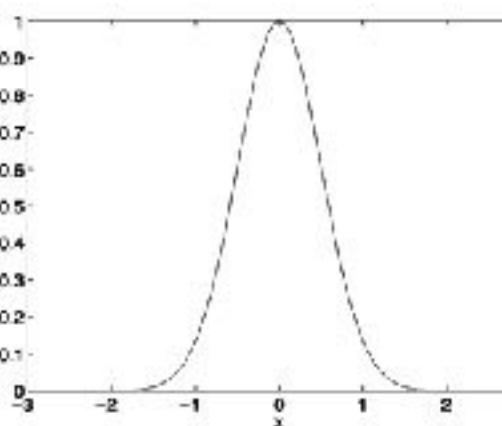


$\mu := \int \phi \, dP$  is the (*kernel*) **mean embedding** of  $P$  in  $\mathcal{H}$ .

$\mu$  can be viewed as a generalized moment vector  
e.g., let  $\phi(x) = [x, x^2]^T$

# Learning with kernels (the modern way)

- A kernel is a symmetric function  
 $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , e.g., Gaussian kernel  
 $k(x, x') = \exp(-\|x - x'\|_2^2 / 2\sigma^2).$
- A p.d.  $k$  corresponds to a Hilbert space  $\mathcal{H}$  (RKHS), which satisfies the **reproducing property**  
 $f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \forall f \in \mathcal{H}, x \in \mathcal{X},$   
 $\phi(x) := k(x, \cdot)$  is the **canonical feature** of  $\mathcal{H}$ .
- If  $\mathcal{H}$  is a large (dense in  $C$ ),  $\gamma_{\mathcal{H}}$  is a metric on  $\mathcal{P}$ .



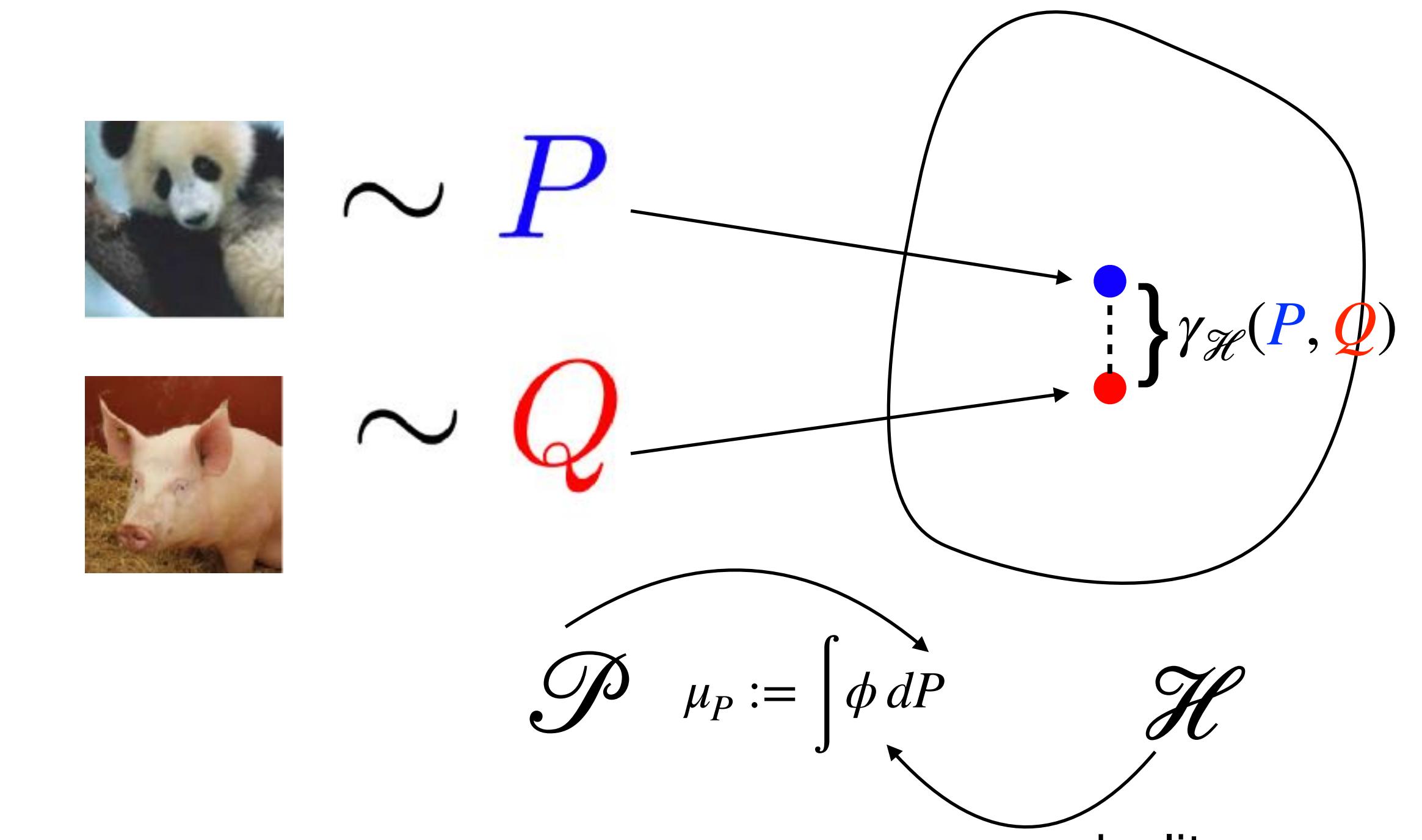
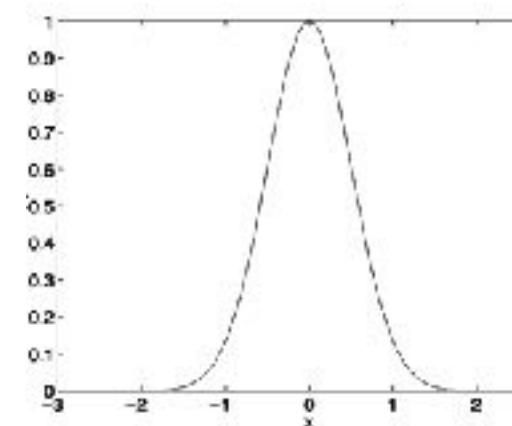
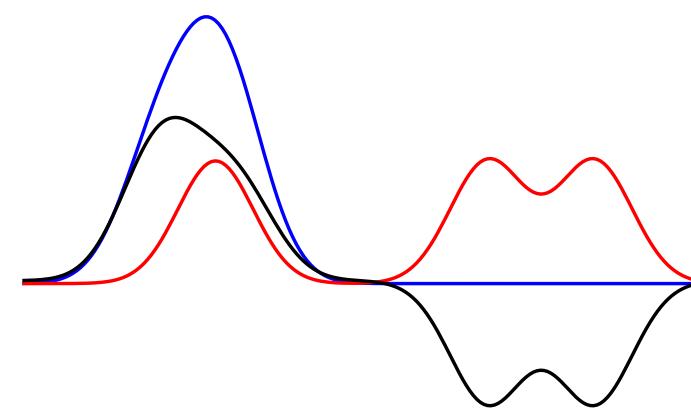
$\mu := \int \phi dP$  is the (*kernel*) **mean embedding** of  $P$  in  $\mathcal{H}$ .

$\mu$  can be viewed as a generalized moment vector  
e.g., let  $\phi(x) = [x, x^2]^T$

# Learning with kernels (the modern way)

- A kernel is a symmetric function  
 $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , e.g., Gaussian kernel  
 $k(x, x') = \exp(-\|x - x'\|_2^2 / 2\sigma^2).$
- A p.d.  $k$  corresponds to a Hilbert space  $\mathcal{H}$  (RKHS), which satisfies the **reproducing property**  
 $f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \forall f \in \mathcal{H}, x \in \mathcal{X},$   
 $\phi(x) := k(x, \cdot)$  is the **canonical feature** of  $\mathcal{H}$ .
- If  $\mathcal{H}$  is a large (dense in  $C$ ),  $\gamma_{\mathcal{H}}$  is a metric on  $\mathcal{P}$ .
- We can generalize to the more general **integral probability metric** (IPM)

$$\text{IPM}(\mathcal{F}; P, Q) := \sup_{f \in \mathcal{F}} \int f d(P - Q).$$



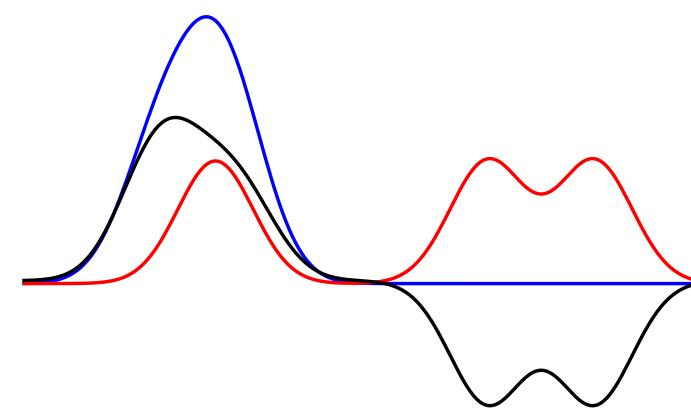
$\mu := \int \phi dP$  is the (*kernel*) **mean embedding** of  $P$  in  $\mathcal{H}$ .

$\mu$  can be viewed as a generalized moment vector  
e.g., let  $\phi(x) = [x, x^2]^T$

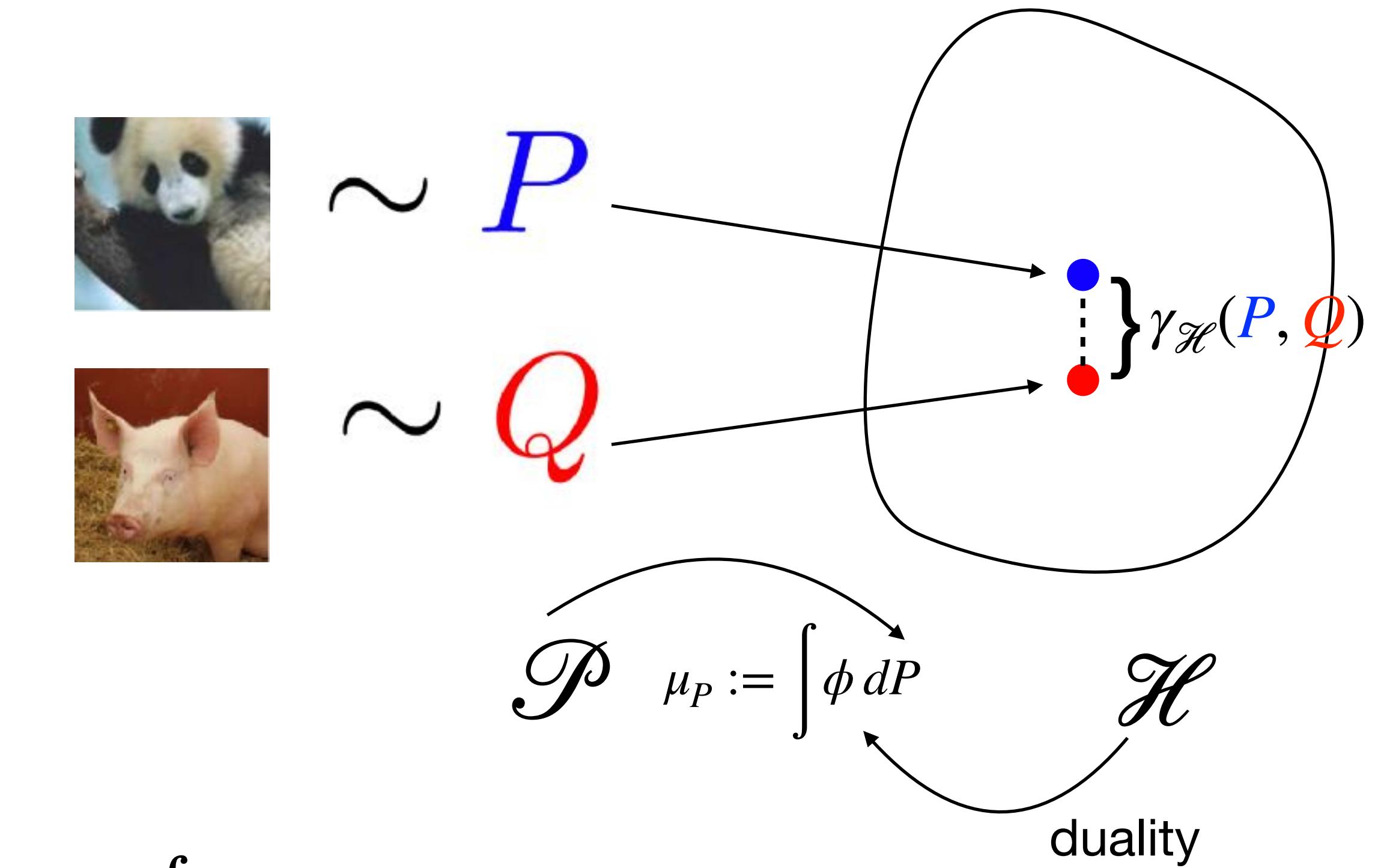
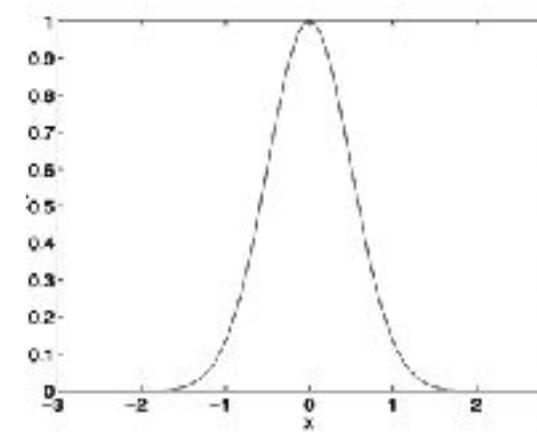
# Learning with kernels (the modern way)

- A kernel is a symmetric function  
 $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , e.g., Gaussian kernel  
 $k(x, x') = \exp(-\|x - x'\|_2^2 / 2\sigma^2).$
- A p.d.  $k$  corresponds to a Hilbert space  $\mathcal{H}$  (RKHS), which satisfies the **reproducing property**  
 $f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}, \forall f \in \mathcal{H}, x \in \mathcal{X},$   
 $\phi(x) := k(x, \cdot)$  is the **canonical feature** of  $\mathcal{H}$ .
- If  $\mathcal{H}$  is a large (dense in  $C$ ),  $\gamma_{\mathcal{H}}$  is a metric on  $\mathcal{P}$ .
- We can generalize to the more general **integral probability metric** (IPM)

$$\text{IPM}(\mathcal{F}; P, Q) := \sup_{f \in \mathcal{F}} \int f d(P - Q).$$



- Special cases:  
 $\mathcal{F} = \{f: \|f\|_{\mathcal{H}} \leq 1\} \rightarrow$  Maximum Mean Discrepancy (MMD)  
 $\mathcal{F} = \{f: \|f\|_{\text{lip}} \leq 1\} \rightarrow$  Wasserstein (type-1)



$\mu := \int \phi dP$  is the (*kernel*) **mean embedding** of  $P$  in  $\mathcal{H}$ .

$\mu$  can be viewed as a generalized moment vector  
e.g., let  $\phi(x) = [x, x^2]^T$

# The MMD



$\sim P$



$\sim Q$

# The MMD



$\sim P$



$\sim Q$

## Maximum Mean Discrepancy (MMD) [Gretton et al., 2012a]

$$\text{MMD}(p, q) := \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}},$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ .

- Recall  $\|a - b\|^2 = \langle a, a \rangle - 2 \langle a, b \rangle + \langle b, b \rangle$ .

$$\begin{aligned}\text{MMD}^2(p, q) &= \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^2 \\ &= \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{x' \sim p}[\phi(x')] \rangle - 2 \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{y \sim q}[\phi(y)] \rangle \\ &\quad + \langle \mathbb{E}_{y \sim q}[\phi(y)], \mathbb{E}_{y' \sim q}[\phi(y')] \rangle \\ &= \mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim p} \langle \phi(x), \phi(x') \rangle - 2 \mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} \langle \phi(x), \phi(y) \rangle \\ &\quad + \mathbb{E}_{y \sim q} \mathbb{E}_{y' \sim q} \langle \phi(y), \phi(y') \rangle\end{aligned}$$

- Depend on only the inner product  $\langle \phi(x), \phi(y) \rangle$ .
- Don't need  $\phi(x)$  explicitly (could be  $\infty$ -dimensional!).

## Maximum Mean Discrepancy (MMD) [Gretton et al., 2012a]

$$\text{MMD}(p, q) := \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}},$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ .

- Recall  $\|a - b\|^2 = \langle a, a \rangle - 2 \langle a, b \rangle + \langle b, b \rangle$ .

$$\begin{aligned}\text{MMD}^2(p, q) &= \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^2 \\ &= \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{x' \sim p}[\phi(x')] \rangle - 2 \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{y \sim q}[\phi(y)] \rangle \\ &\quad + \langle \mathbb{E}_{y \sim q}[\phi(y)], \mathbb{E}_{y' \sim q}[\phi(y')] \rangle \\ &= \mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim p} \langle \phi(x), \phi(x') \rangle - 2 \mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} \langle \phi(x), \phi(y) \rangle \\ &\quad + \mathbb{E}_{y \sim q} \mathbb{E}_{y' \sim q} \langle \phi(y), \phi(y') \rangle\end{aligned}$$

- Depend on only the inner product  $\langle \phi(x), \phi(y) \rangle$ .
- Don't need  $\phi(x)$  explicitly (could be  $\infty$ -dimensional!).

## Maximum Mean Discrepancy (MMD) [Gretton et al., 2012a]

$$\text{MMD}(p, q) := \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}},$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ .

- Recall  $\|a - b\|^2 = \langle a, a \rangle - 2 \langle a, b \rangle + \langle b, b \rangle$ .

$$\begin{aligned}\text{MMD}^2(p, q) &= \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^2 \\ &= \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{x' \sim p}[\phi(x')] \rangle - 2 \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{y \sim q}[\phi(y)] \rangle \\ &\quad + \langle \mathbb{E}_{y \sim q}[\phi(y)], \mathbb{E}_{y' \sim q}[\phi(y')] \rangle \\ &= \mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim p} \langle \phi(x), \phi(x') \rangle - 2 \mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} \langle \phi(x), \phi(y) \rangle \\ &\quad + \mathbb{E}_{y \sim q} \mathbb{E}_{y' \sim q} \langle \phi(y), \phi(y') \rangle\end{aligned}$$

- Depend on only the inner product  $\langle \phi(x), \phi(y) \rangle$ .
- Don't need  $\phi(x)$  explicitly (could be  $\infty$ -dimensional!).

## Maximum Mean Discrepancy (MMD) [Gretton et al., 2012a]

$$\text{MMD}(p, q) := \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}},$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ .

- Recall  $\|a - b\|^2 = \langle a, a \rangle - 2 \langle a, b \rangle + \langle b, b \rangle$ .

$$\begin{aligned}\text{MMD}^2(p, q) &= \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^2 \\ &= \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{x' \sim p}[\phi(x')] \rangle - 2 \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{y \sim q}[\phi(y)] \rangle \\ &\quad + \langle \mathbb{E}_{y \sim q}[\phi(y)], \mathbb{E}_{y' \sim q}[\phi(y')] \rangle \\ &= \mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim p} \langle \phi(x), \phi(x') \rangle - 2 \mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} \langle \phi(x), \phi(y) \rangle \\ &\quad + \mathbb{E}_{y \sim q} \mathbb{E}_{y' \sim q} \langle \phi(y), \phi(y') \rangle\end{aligned}$$

- Depend on only the inner product  $\langle \phi(x), \phi(y) \rangle$ .
- Don't need  $\phi(x)$  explicitly (could be  $\infty$ -dimensional!).

## Maximum Mean Discrepancy (MMD) [Gretton et al., 2012a]

$$\text{MMD}(p, q) := \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}},$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ .

- Recall  $\|a - b\|^2 = \langle a, a \rangle - 2 \langle a, b \rangle + \langle b, b \rangle$ .

$$\begin{aligned}\text{MMD}^2(p, q) &= \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^2 \\ &= \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{x' \sim p}[\phi(x')] \rangle - 2 \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{y \sim q}[\phi(y)] \rangle \\ &\quad + \langle \mathbb{E}_{y \sim q}[\phi(y)], \mathbb{E}_{y' \sim q}[\phi(y')] \rangle \\ &= \mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim p} \langle \phi(x), \phi(x') \rangle - 2 \mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} \langle \phi(x), \phi(y) \rangle \\ &\quad + \mathbb{E}_{y \sim q} \mathbb{E}_{y' \sim q} \langle \phi(y), \phi(y') \rangle\end{aligned}$$

- Depend on only the inner product  $\langle \phi(x), \phi(y) \rangle$ .
- Don't need  $\phi(x)$  explicitly (could be  $\infty$ -dimensional!).

## Maximum Mean Discrepancy (MMD) [Gretton et al., 2012a]

$$\text{MMD}(p, q) := \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}},$$

where  $\phi(x)$  is in the Hilbert space  $\mathcal{H}$ .

- Recall  $\|a - b\|^2 = \langle a, a \rangle - 2 \langle a, b \rangle + \langle b, b \rangle$ .

$$\begin{aligned}\text{MMD}^2(p, q) &= \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^2 \\ &= \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{x' \sim p}[\phi(x')] \rangle - 2 \langle \mathbb{E}_{x \sim p}[\phi(x)], \mathbb{E}_{y \sim q}[\phi(y)] \rangle \\ &\quad + \langle \mathbb{E}_{y \sim q}[\phi(y)], \mathbb{E}_{y' \sim q}[\phi(y')] \rangle \\ &= \mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim p} \langle \phi(x), \phi(x') \rangle - 2 \mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} \langle \phi(x), \phi(y) \rangle \\ &\quad + \mathbb{E}_{y \sim q} \mathbb{E}_{y' \sim q} \langle \phi(y), \phi(y') \rangle\end{aligned}$$

- Depend on only the inner product  $\langle \phi(x), \phi(y) \rangle$ .
- Don't need  $\phi(x)$  explicitly (could be  $\infty$ -dimensional!).

## Maximum Mean Discrepancy (MMD) [Gretton et al., 2012a]

- Define  $k(x, x') := \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$  (kernel).

$$\begin{aligned}\text{MMD}_{\mathbf{k}}^2(p, q) &= \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^2 \\ &= \mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim p} k(x, x') - 2 \mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} k(x, y') \\ &\quad + \mathbb{E}_{y \sim q} \mathbb{E}_{y' \sim q} k(y, y').\end{aligned}$$

- Unbiased estimator:

$$\begin{aligned}\widehat{\text{MMD}}_{\mathbf{k}}^2(X, Y) &= \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i} k(x_i, x_j) - \frac{2}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(x_i, y_j) \\ &\quad + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i} k(y_i, y_j).\end{aligned}$$

- $k(x, x') \approx$  similarity between  $x$  and  $x'$ .

## Maximum Mean Discrepancy (MMD) [Gretton et al., 2012a]

- Define  $k(x, x') := \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$  (kernel).

$$\begin{aligned}\text{MMD}_{\mathbf{k}}^2(p, q) &= \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^2 \\ &= \mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim p} k(x, x') - 2 \mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} k(x, y') \\ &\quad + \mathbb{E}_{y \sim q} \mathbb{E}_{y' \sim q} k(y, y').\end{aligned}$$

- Unbiased estimator:

$$\begin{aligned}\widehat{\text{MMD}}_{\mathbf{k}}^2(X, Y) &= \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i} k(x_i, x_j) - \frac{2}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(x_i, y_j) \\ &\quad + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i} k(y_i, y_j).\end{aligned}$$

- $k(x, x') \approx$  similarity between  $x$  and  $x'$ .

## Maximum Mean Discrepancy (MMD) [Gretton et al., 2012a]

- Define  $k(x, x') := \langle \phi(x), \phi(x') \rangle_{\mathcal{H}}$  (kernel).

$$\begin{aligned}\text{MMD}_{\mathbf{k}}^2(p, q) &= \left\| \mathbb{E}_{x \sim p}[\phi(x)] - \mathbb{E}_{y \sim q}[\phi(y)] \right\|_{\mathcal{H}}^2 \\ &= \mathbb{E}_{x \sim p} \mathbb{E}_{x' \sim p} k(x, x') - 2 \mathbb{E}_{x \sim p} \mathbb{E}_{y \sim q} k(x, y') \\ &\quad + \mathbb{E}_{y \sim q} \mathbb{E}_{y' \sim q} k(y, y').\end{aligned}$$

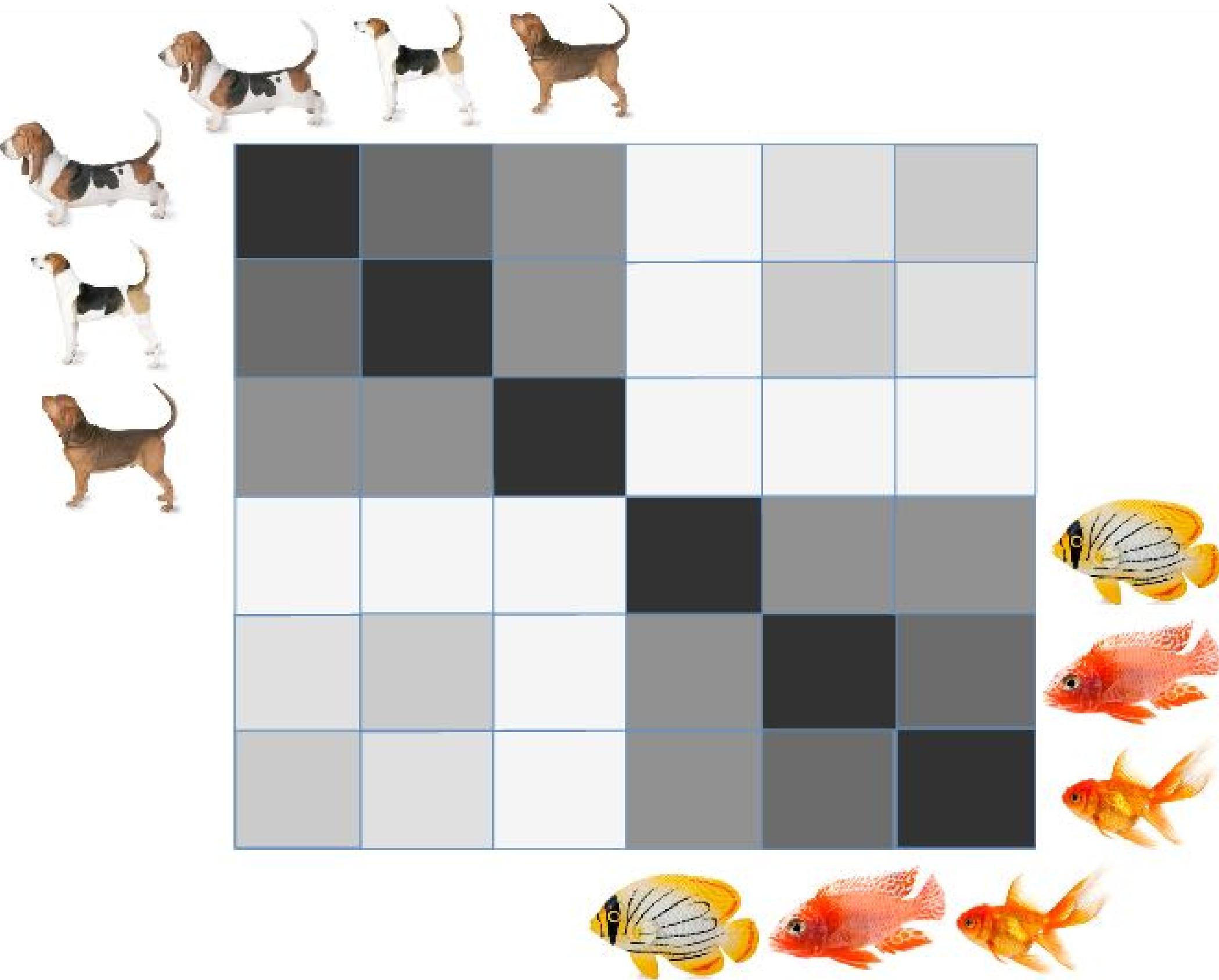
- Unbiased estimator:

$$\begin{aligned}\widehat{\text{MMD}}_{\mathbf{k}}^2(X, Y) &= \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i} k(x_i, x_j) - \frac{2}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(x_i, y_j) \\ &\quad + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i} k(y_i, y_j).\end{aligned}$$

- $k(x, x') \approx$  similarity between  $x$  and  $x'$ .

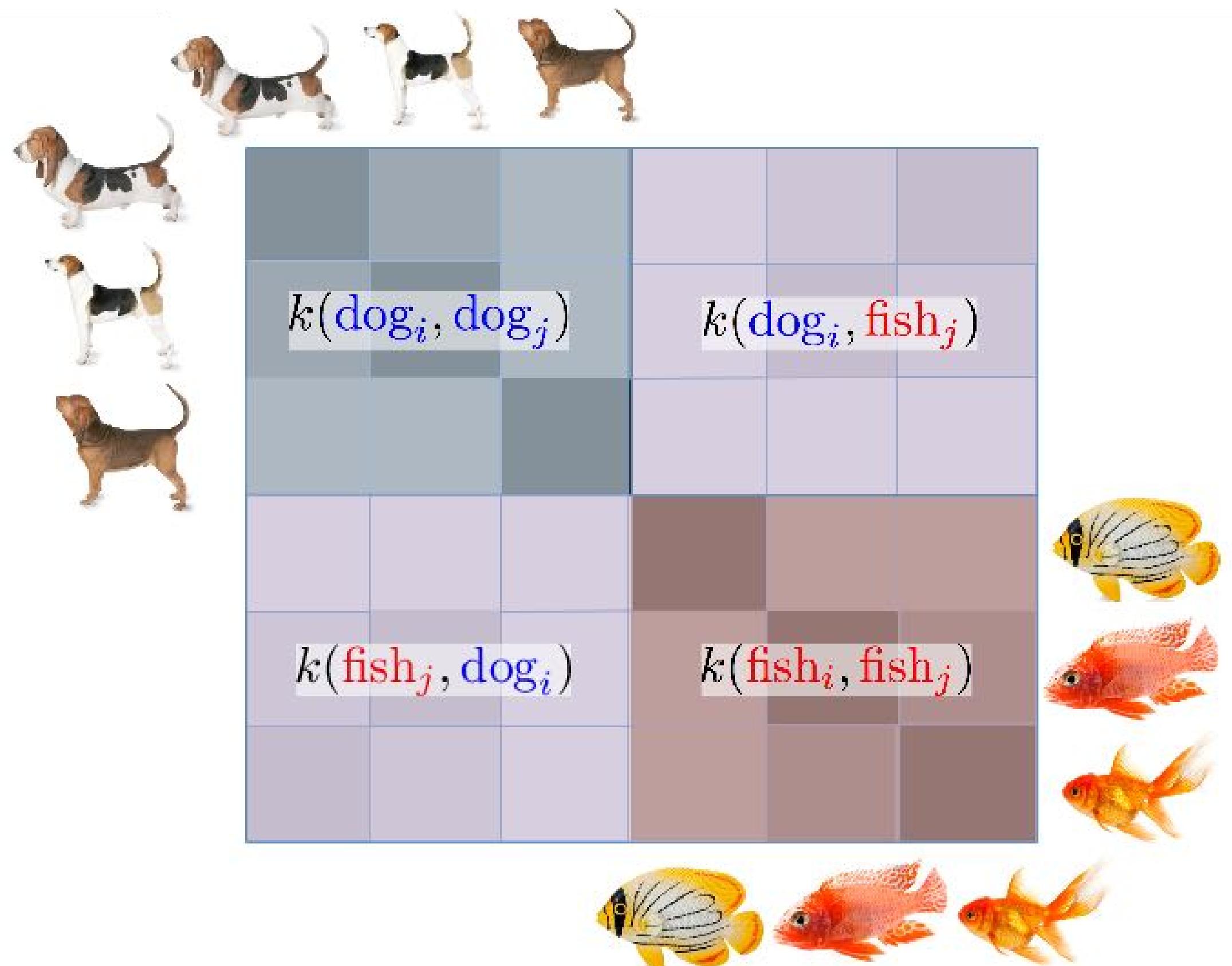
# Intuition for the MMD

- Dogs  $\sim p$  and fish  $\sim q$ .
  - Each entry is one of  $k(\text{dog}_i, \text{dog}_j)$ ,  $k(\text{dog}_i, \text{fish}_j)$ , or  $k(\text{fish}_i, \text{fish}_j)$

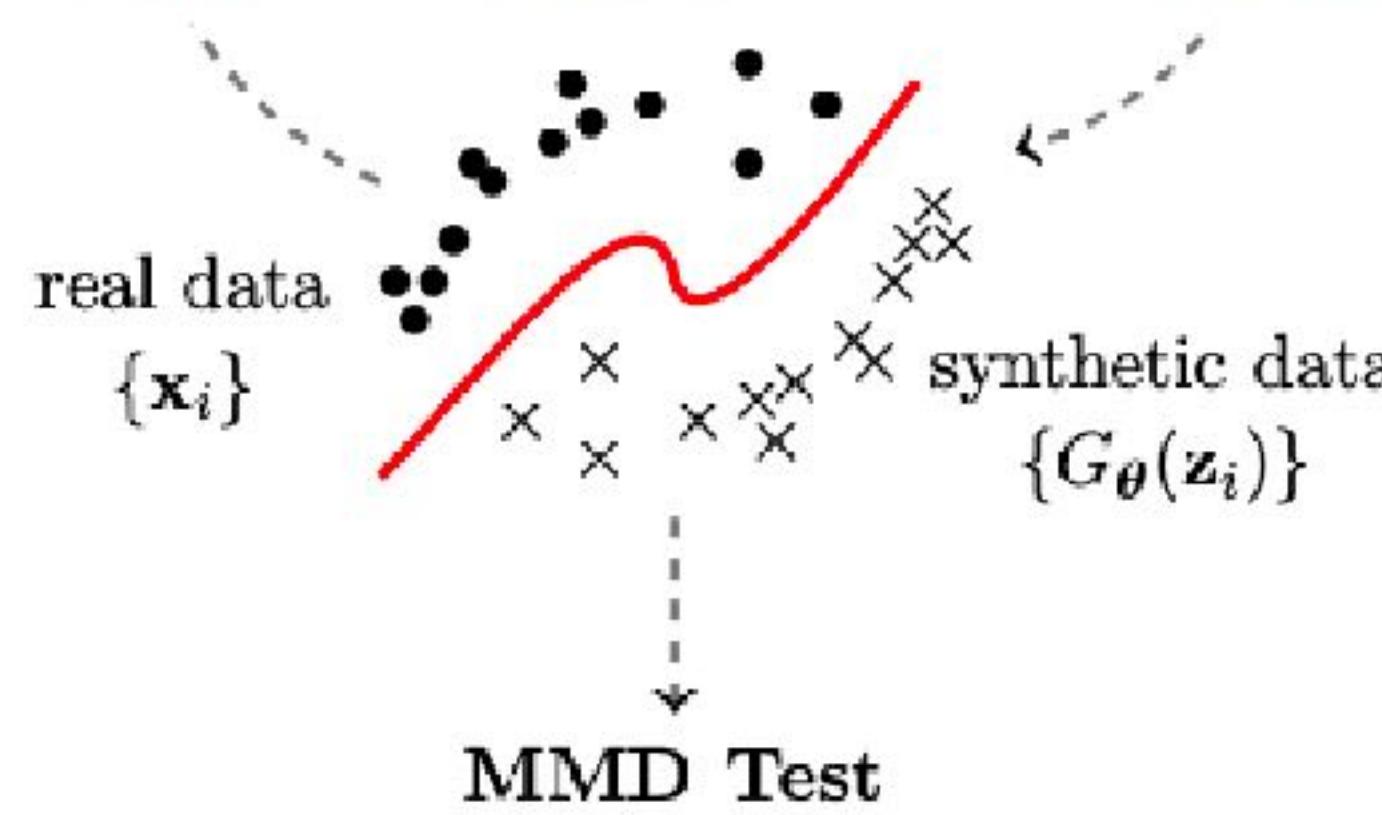
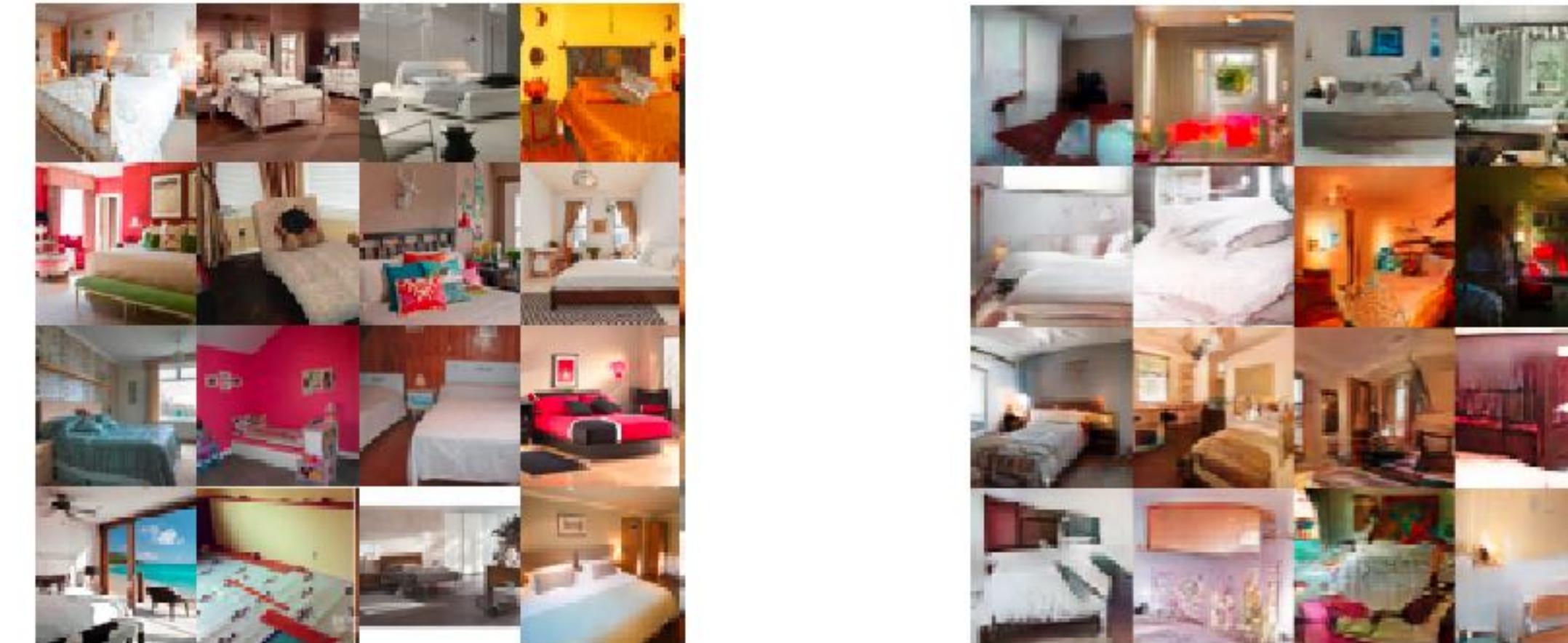


# Intuition for the MMD

$$\widehat{\text{MMD}}_k^2 = \frac{1}{n(n-1)} \sum_{i \neq j} k(\text{dog}_i, \text{dog}_j) - \frac{2}{n^2} \sum_{i,j} k(\text{dog}_i, \text{fish}_j) \\ + \frac{1}{n(n-1)} \sum_{i \neq j} k(\text{fish}_i, \text{fish}_j)$$



# Application to generative modeling: MMD for generative adversarial nets (GAN)



$$\|\hat{\mu}_X - \hat{\mu}_{G_{\theta}(Z)}\|_{\mathcal{H}} \text{ is zero?}$$

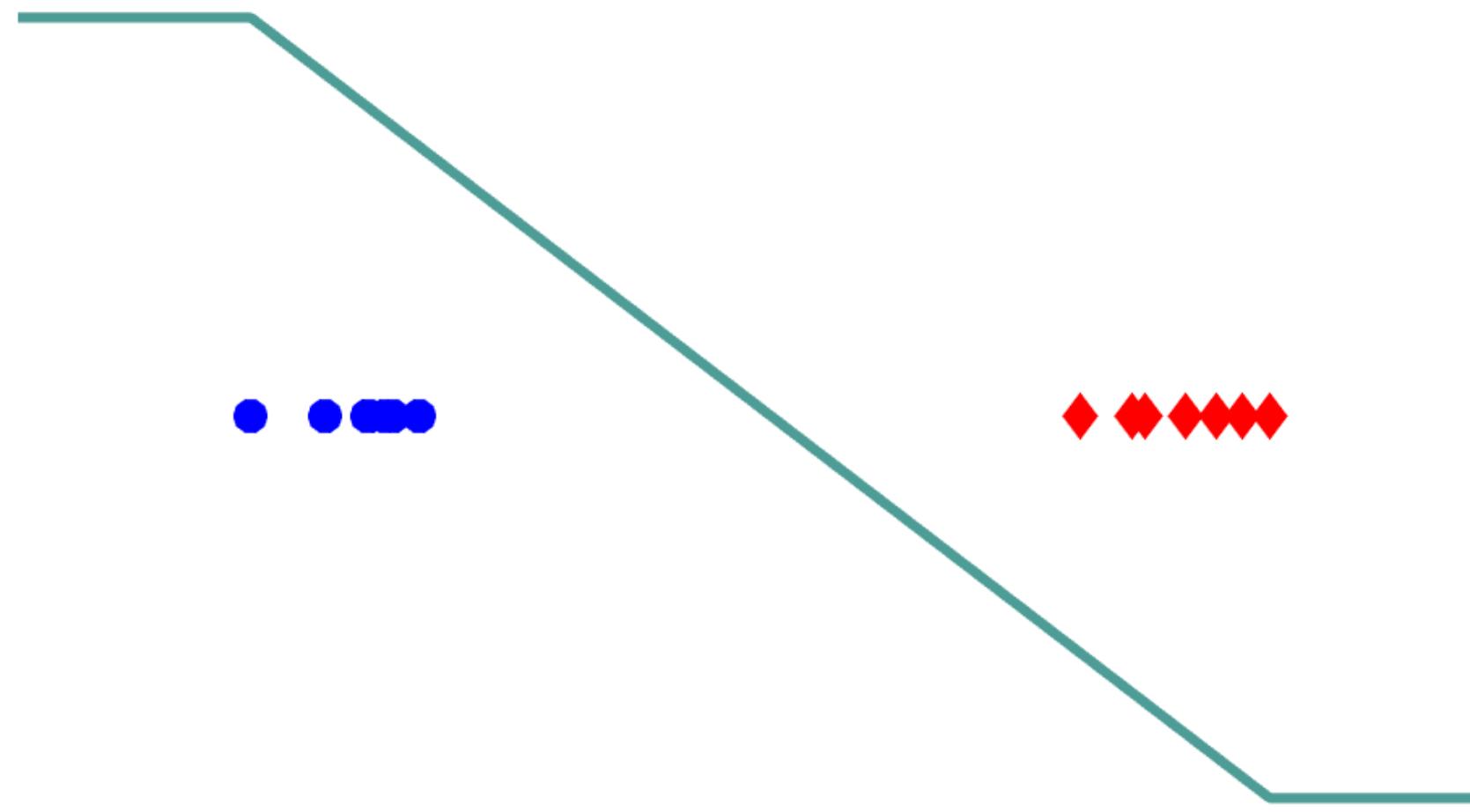
# Comparison: Wasserstein-1 vs. MMD- $k$

# Comparison: Wasserstein-1 vs. MMD- $k$

$$W_1(\mathcal{P}, \mathcal{Q}) = \sup_{\|\mathbf{f}\|_L \leq 1} E_{\mathcal{P}} \mathbf{f}(\mathcal{X}) - E_{\mathcal{Q}} \mathbf{f}(\mathcal{Y}).$$

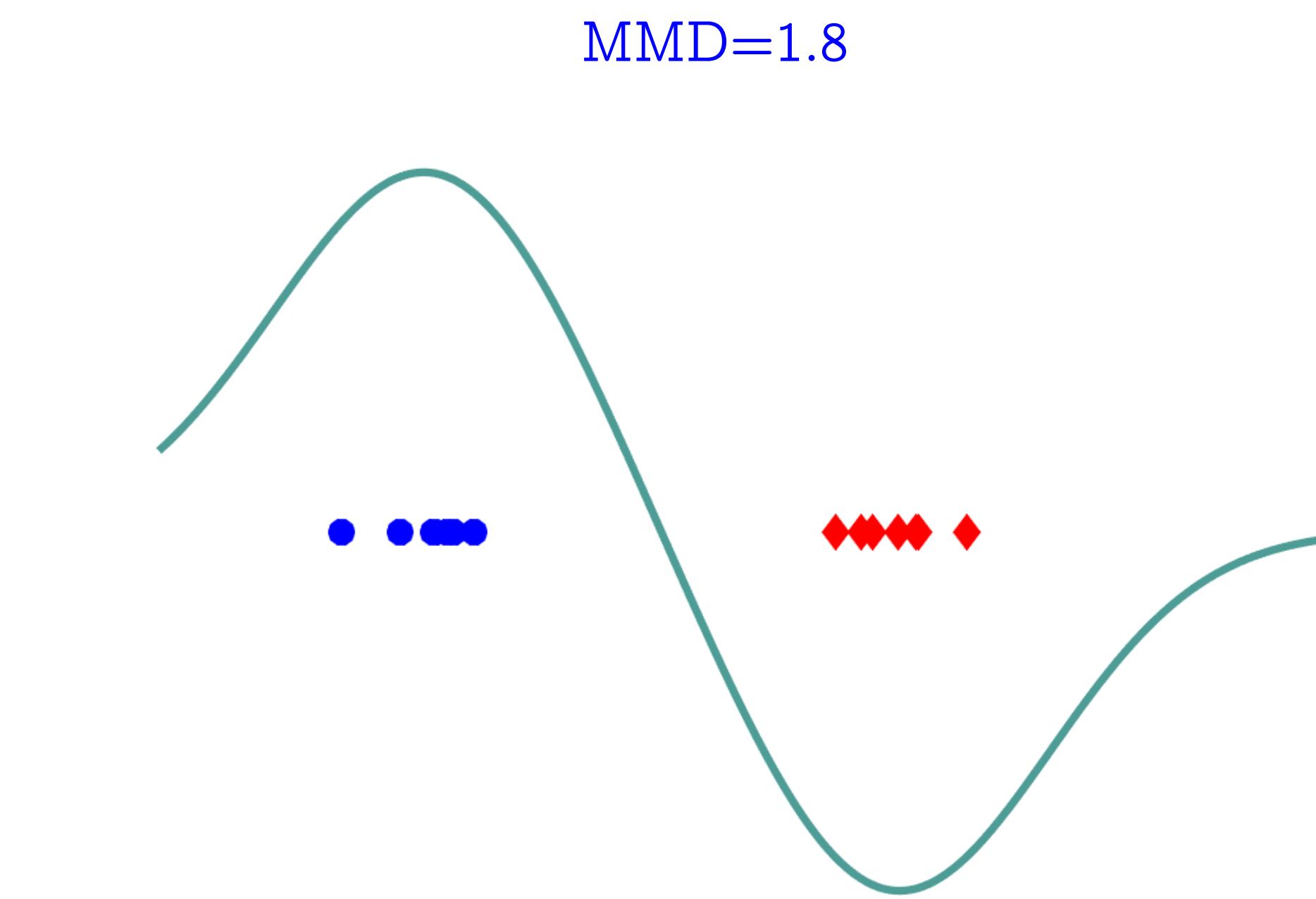
$$\|\mathbf{f}\|_L := \sup_{x \neq y} |f(x) - f(y)| / \|x - y\|$$

$$W_1=0.88$$



# Comparison: Wasserstein-1 vs. MMD- $k$

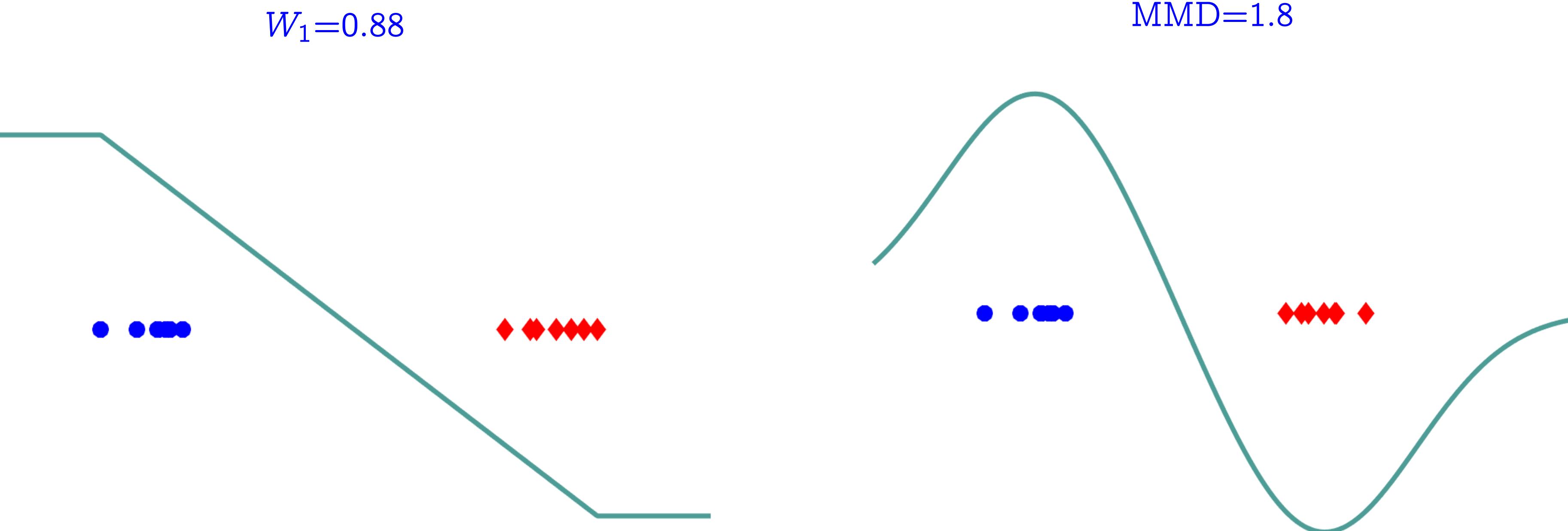
$$MMD(\textcolor{blue}{P}, \textcolor{red}{Q}) = \sup_{\|\textcolor{teal}{f}\|_{\mathcal{F}} \leq 1} E_{\textcolor{blue}{P}} \textcolor{teal}{f}(\textcolor{blue}{X}) - E_{\textcolor{red}{Q}} \textcolor{teal}{f}(\textcolor{red}{Y}).$$



# Comparison: Wasserstein-1 vs. MMD- $k$

$$W_1(\mathcal{P}, \mathcal{Q}) = \sup_{\|\mathbf{f}\|_L \leq 1} E_{\mathcal{P}} \mathbf{f}(\mathbf{X}) - E_{\mathcal{Q}} \mathbf{f}(\mathbf{Y}).$$
$$\|\mathbf{f}\|_L := \sup_{x \neq y} |f(x) - f(y)| / \|x - y\|$$

$$MMD(\mathcal{P}, \mathcal{Q}) = \sup_{\|\mathbf{f}\|_{\mathcal{F}} \leq 1} E_{\mathcal{P}} \mathbf{f}(\mathbf{X}) - E_{\mathcal{Q}} \mathbf{f}(\mathbf{Y}).$$



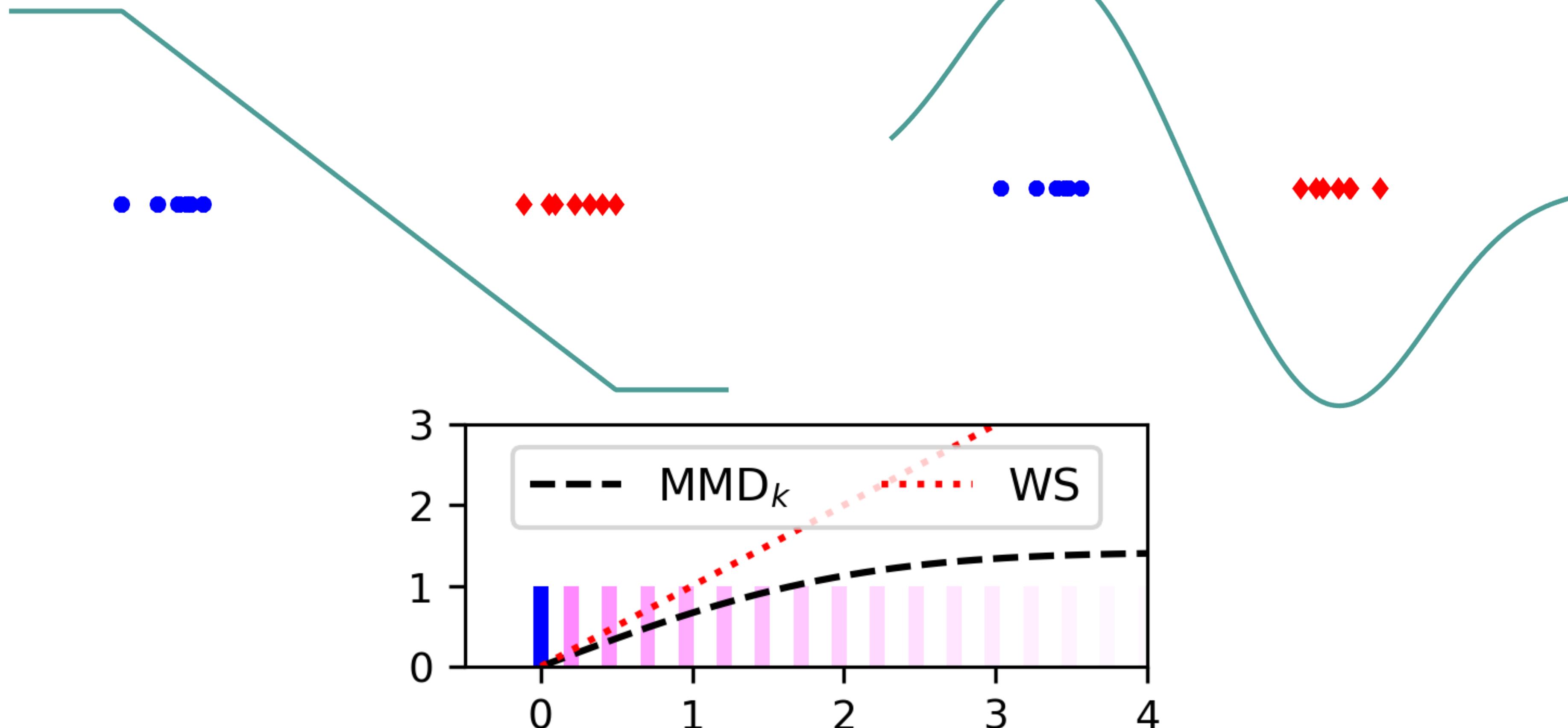
# Comparison: Wasserstein-1 vs. MMD- $k$

$$W_1(\mathcal{P}, \mathcal{Q}) = \sup_{\|\mathbf{f}\|_L \leq 1} E_{\mathcal{P}} \mathbf{f}(\mathbf{X}) - E_{\mathcal{Q}} \mathbf{f}(\mathbf{Y}).$$
$$\|\mathbf{f}\|_L := \sup_{x \neq y} |f(x) - f(y)| / \|x - y\|$$

$$MMD(\mathcal{P}, \mathcal{Q}) = \sup_{\|\mathbf{f}\|_{\mathcal{F}} \leq 1} E_{\mathcal{P}} \mathbf{f}(\mathbf{X}) - E_{\mathcal{Q}} \mathbf{f}(\mathbf{Y}).$$

$W_1=0.88$

MMD=1.8



# Distributional Robustness

**Combine the strengths of ERM and RO:  
distributionally robust optimization (DRO)**

# Combine the strengths of ERM and RO: distributionally robust optimization (DRO)

$$\text{(ERM)} \min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

$$\text{(RO)} \min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$

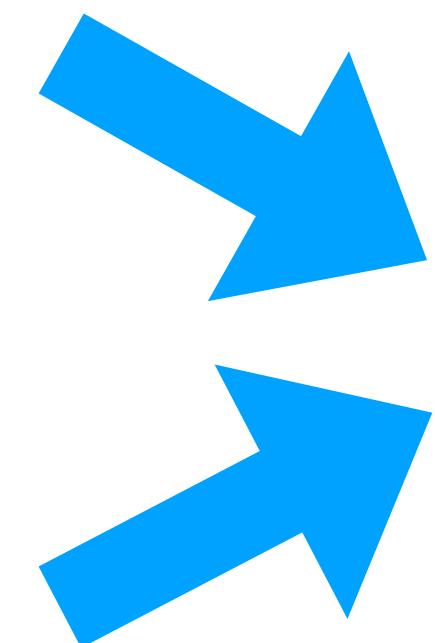
# Combine the strengths of ERM and RO: distributionally robust optimization (DRO)

$$\begin{array}{ll} \text{(ERM)} \min_{\theta} & \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi) \\ & \swarrow \\ \text{(RO)} \min_{\theta} & \sup_{\xi \in \mathcal{U}} l(\theta, \xi) \end{array} \quad \begin{array}{l} \min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad \text{(DRO)} \\ \text{[Delage and Ye 2010, Scarf 1958]} \end{array}$$

# Combine the strengths of ERM and RO: distributionally robust optimization (DRO)

$$\text{(ERM)} \min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

$$\text{(RO)} \min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$



$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad \text{(DRO)}$$

[Delage and Ye 2010, Scarf 1958]

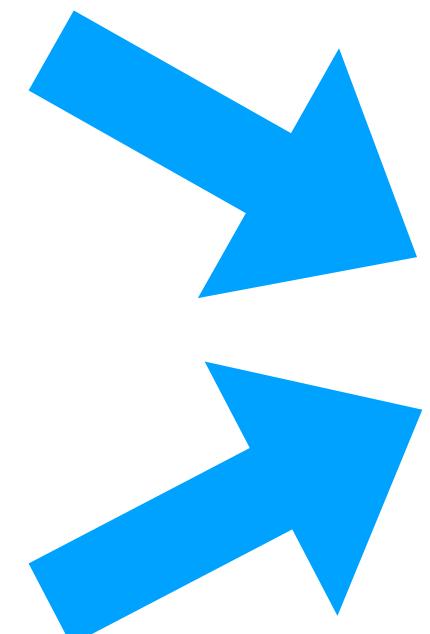
Find the worst-case distribution!

Problem of Moments [Stieltjes, Hausdorff, Hamburger, ...]

# Combine the strengths of ERM and RO: distributionally robust optimization (DRO)

$$(\text{ERM}) \min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

$$(\text{RO}) \min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$



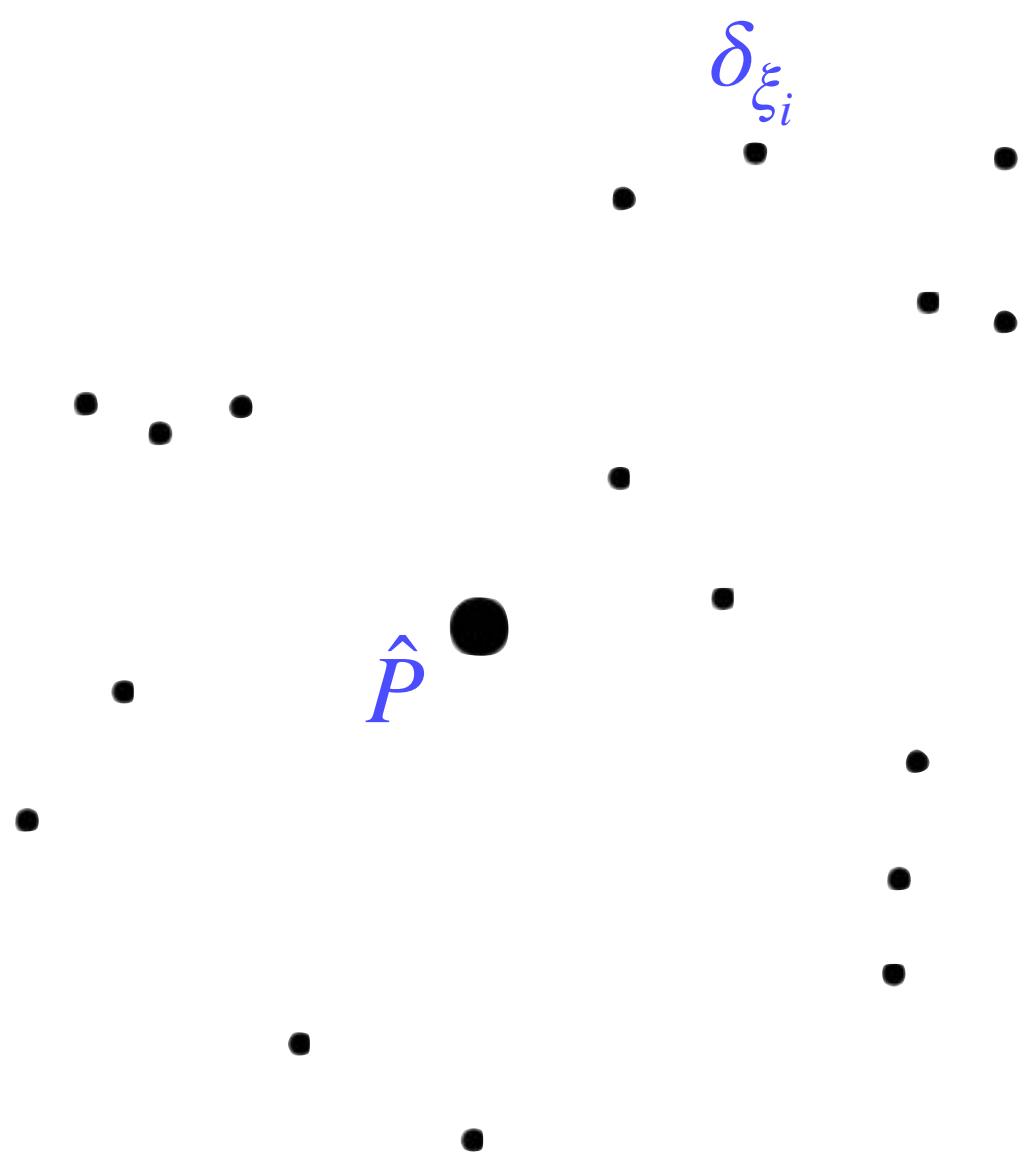
$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad (\text{DRO})$$

[Delage and Ye 2010, Scarf 1958]

Find the worst-case distribution!

Problem of Moments [Stieltjes, Hausdorff, Hamburger, ...]

$$\delta_{\xi_i}$$

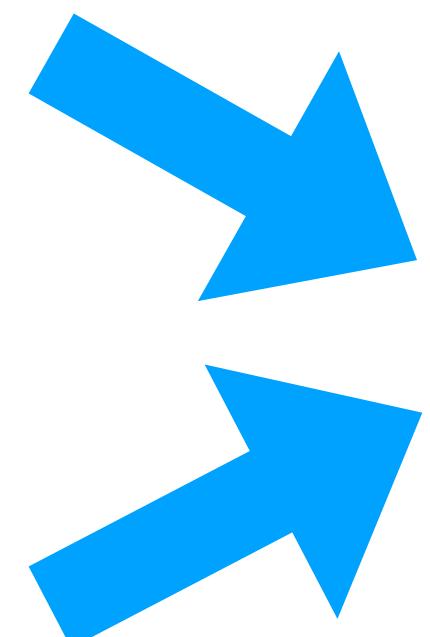


- Robustifies against a set of probability measures  $\mathcal{K}$  (**ambiguity set**), e.g.,

# Combine the strengths of ERM and RO: distributionally robust optimization (DRO)

$$\text{(ERM)} \min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

$$\text{(RO)} \min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$



$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad \text{(DRO)}$$

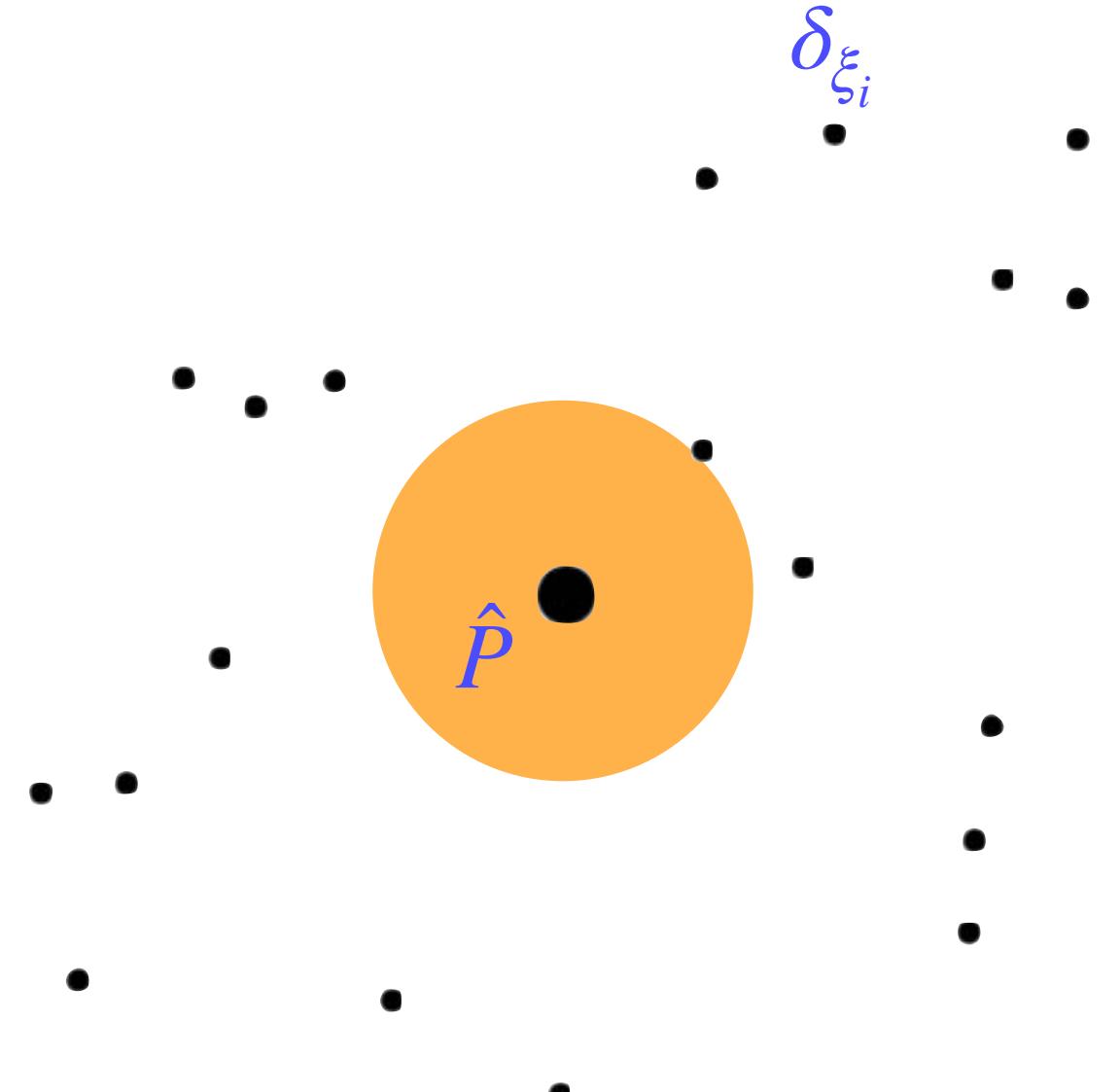
[Delage and Ye 2010, Scarf 1958]

Find the worst-case distribution!

Problem of Moments [Stieltjes, Hausdorff, Hamburger, ...]

$\delta_{\xi_i}$

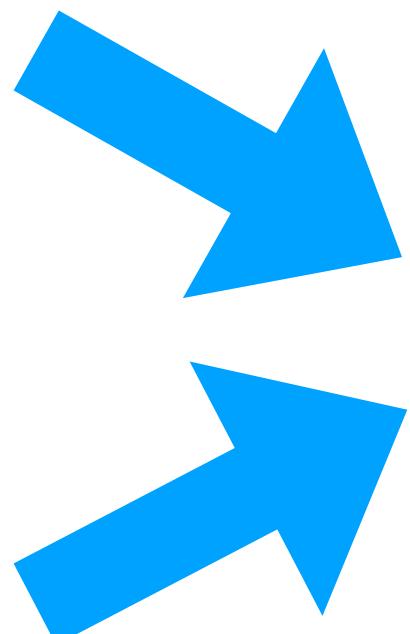
- Robustifies against a set of probability measures  $\mathcal{K}$  (**ambiguity set**), e.g.,
  - $\mathcal{K}$  can be a metric-ball centered at  $\hat{P}$ , e.g., using relative entropy, optimal transport, and kernel methods .



# Combine the strengths of ERM and RO: distributionally robust optimization (DRO)

$$\text{(ERM)} \min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

$$\text{(RO)} \min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$



$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad \text{(DRO)}$$

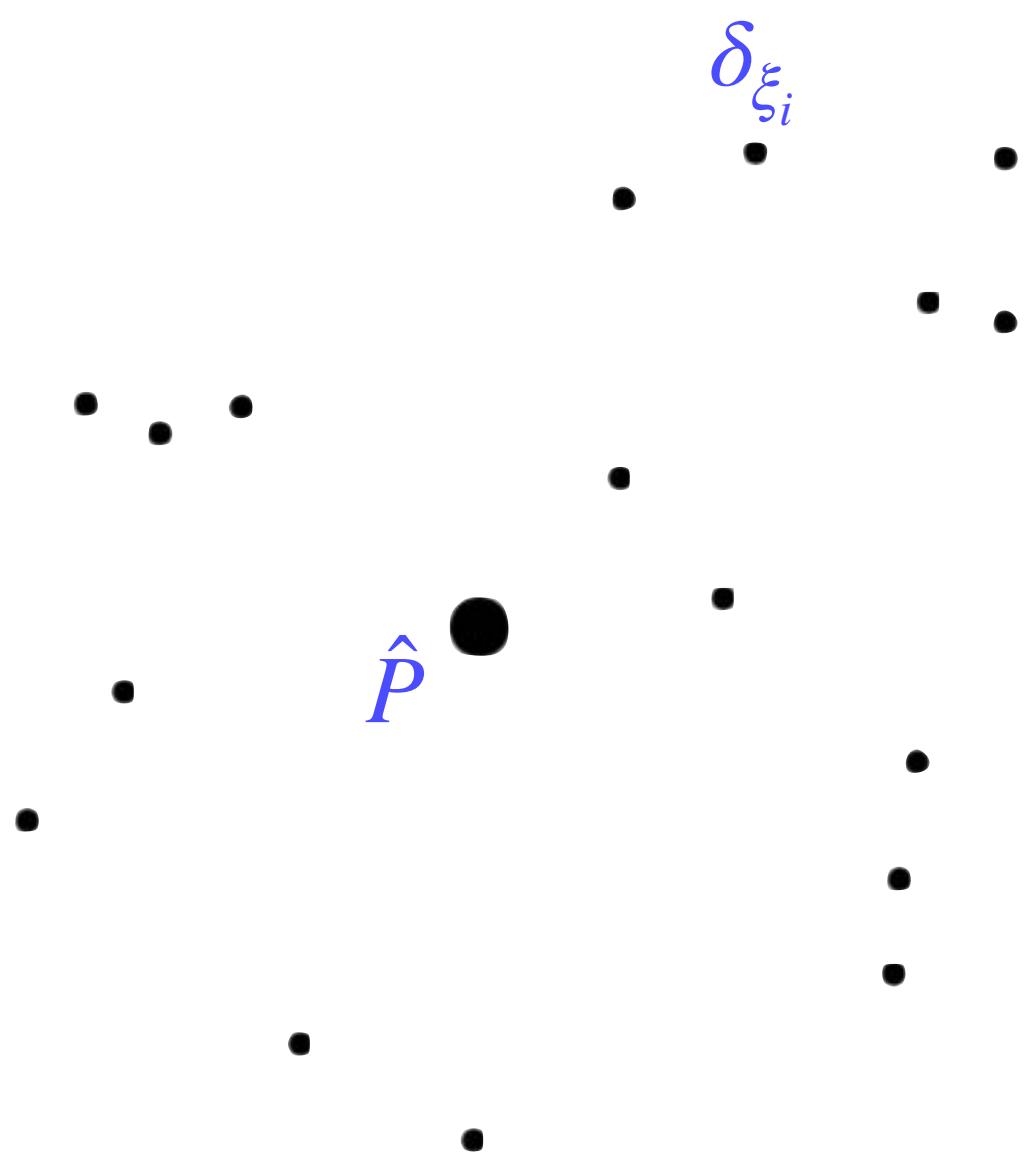
[Delage and Ye 2010, Scarf 1958]

Find the worst-case distribution!

Problem of Moments [Stieltjes, Hausdorff, Hamburger, ...]

$$\delta_{\xi_i}$$

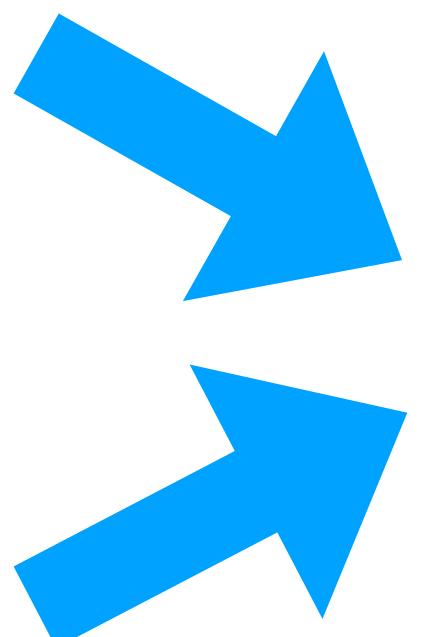
- Robustifies against a set of probability measures  $\mathcal{K}$  (**ambiguity set**), e.g.,
  - $\mathcal{K}$  can be a metric-ball centered at  $\hat{P}$ , e.g., using relative entropy, optimal transport, and kernel methods .



# Combine the strengths of ERM and RO: distributionally robust optimization (DRO)

$$\text{(ERM)} \min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

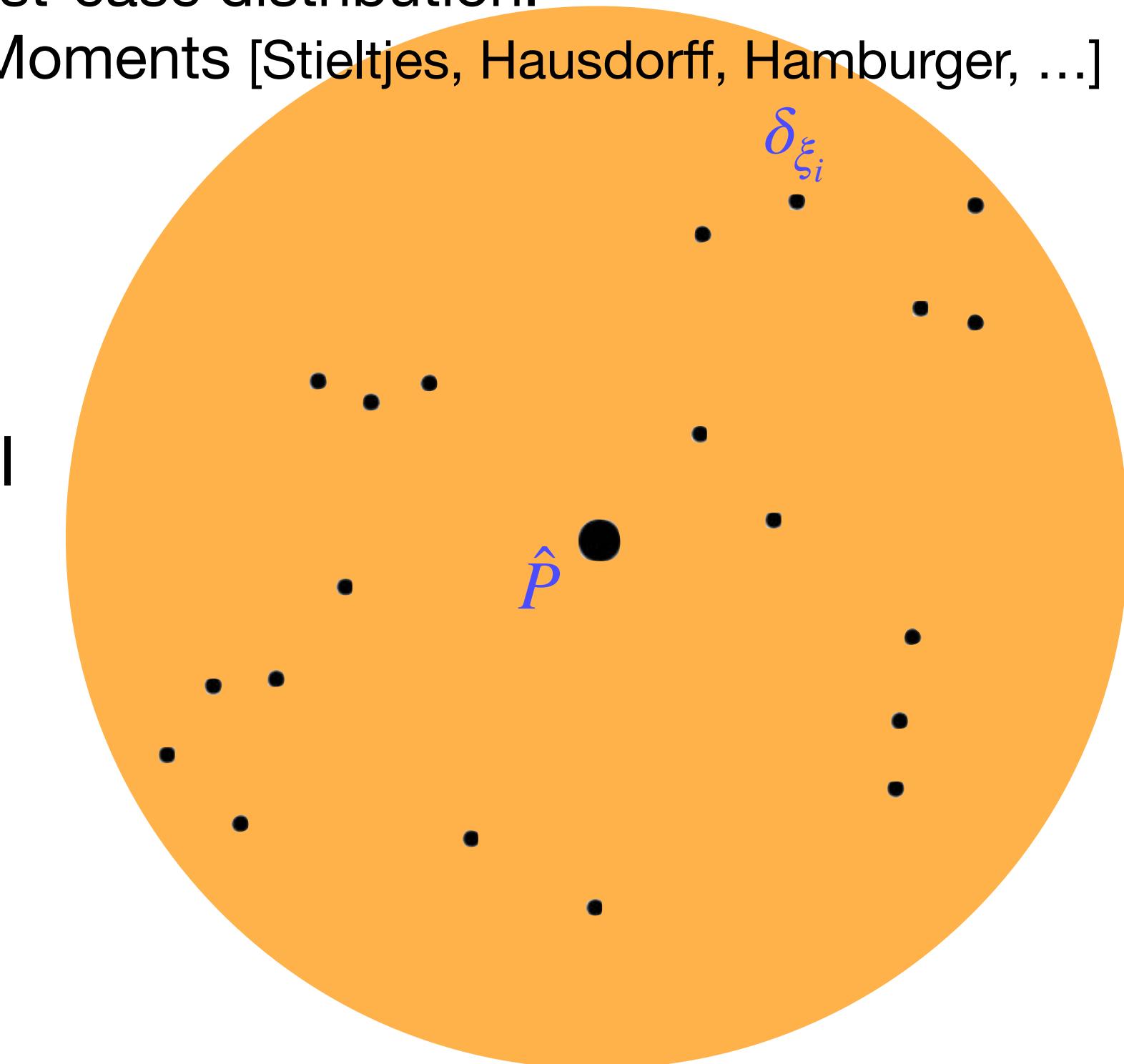
$$\text{(RO)} \min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$



$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad \text{(DRO)}$$

[Delage and Ye 2010, Scarf 1958]

Find the worst-case distribution!  
Problem of Moments [Stieltjes, Hausdorff, Hamburger, ...]

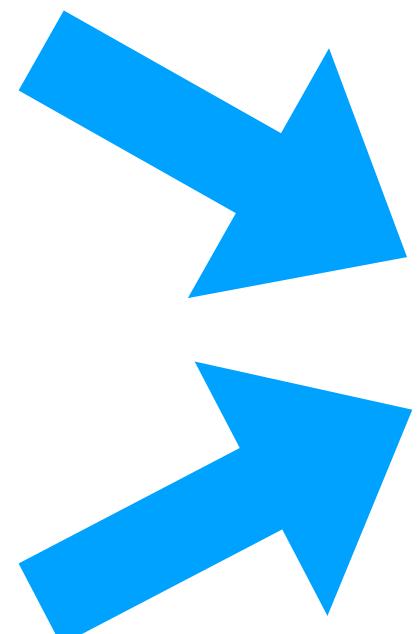


- Robustifies against a set of probability measures  $\mathcal{K}$  (**ambiguity set**), e.g.,
  - $\mathcal{K}$  can be a metric-ball centered at  $\hat{P}$ , e.g., using relative entropy, optimal transport, and kernel methods .

# Combine the strengths of ERM and RO: distributionally robust optimization (DRO)

$$\text{(ERM)} \min_{\theta} \mathbb{E}_{\xi \sim \hat{P}} l(\theta, \xi)$$

$$\text{(RO)} \min_{\theta} \sup_{\xi \in \mathcal{U}} l(\theta, \xi)$$



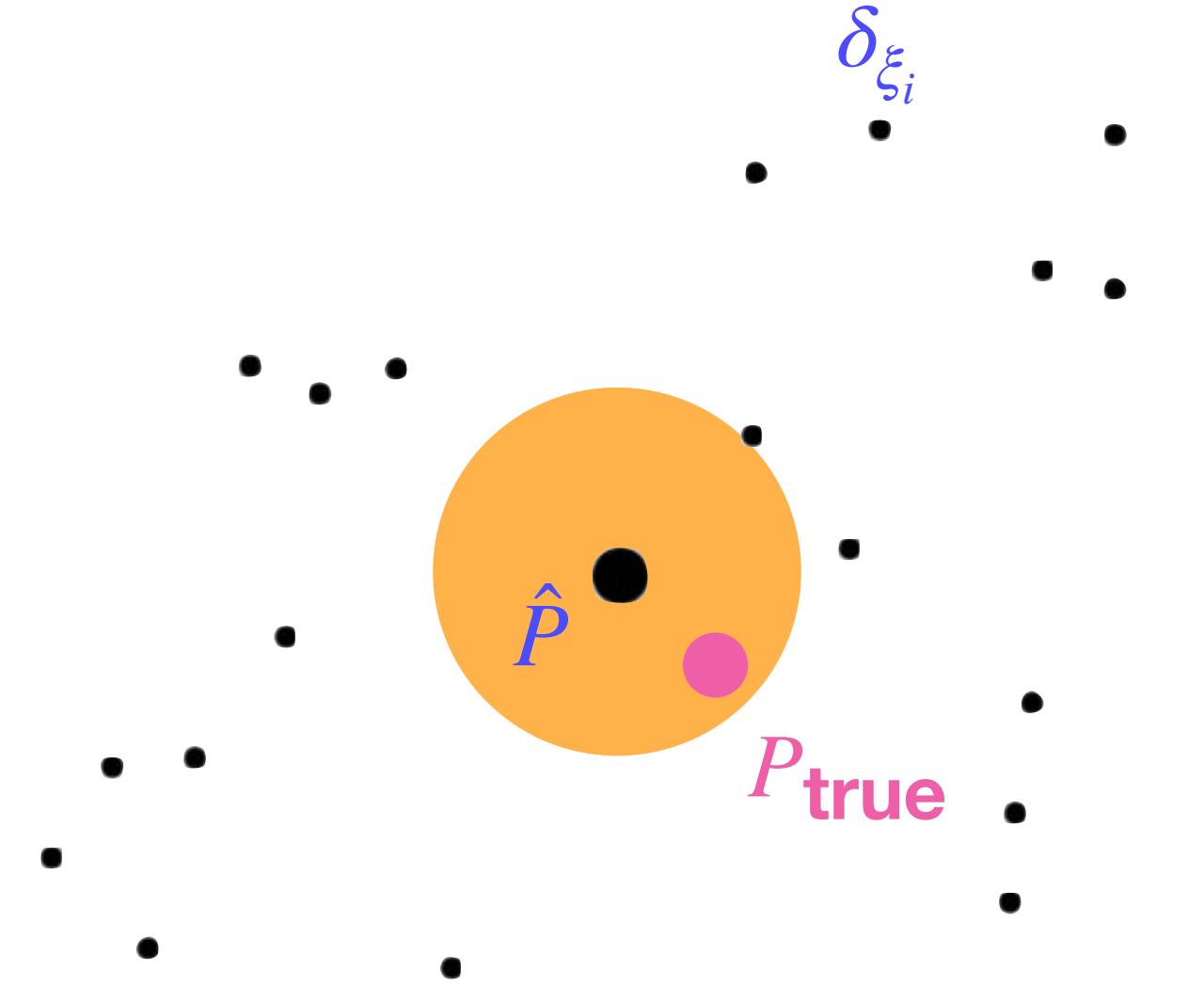
$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad \text{(DRO)}$$

[Delage and Ye 2010, Scarf 1958]

Find the worst-case distribution!  
Problem of Moments [Stieltjes, Hausdorff, Hamburger, ...]

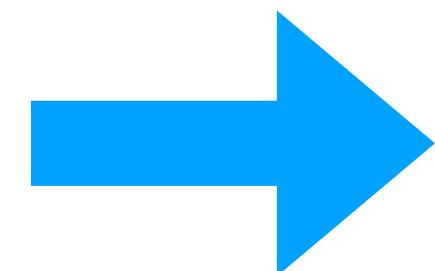
$\delta_{\xi_i}$

- Robustifies against a set of probability measures  $\mathcal{K}$  (**ambiguity set**), e.g.,
  - $\mathcal{K}$  can be a metric-ball centered at  $\hat{P}$ , e.g., using relative entropy, optimal transport, and kernel methods .
  - One way of constructing ambiguity region: one can quantify the empirical mean convergence rate  $\gamma(\hat{P}, P_{\text{true}}) \leq \epsilon$ .



# What distributions does the ambiguity set $\mathcal{K}$ contain?

$$\min_{\theta} \sup_{\xi \in \mathcal{X}} l(\theta, \xi) \quad (\text{RO})$$



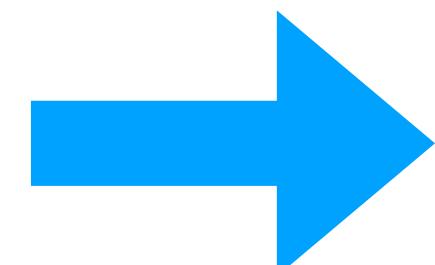
$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad (\text{DRO})$$

Uncertainty set

Ambiguity set

# What distributions does the ambiguity set $\mathcal{K}$ contain?

$$\min_{\theta} \sup_{\xi \in \mathcal{X}} l(\theta, \xi) \quad (\text{RO})$$



$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad (\text{DRO})$$

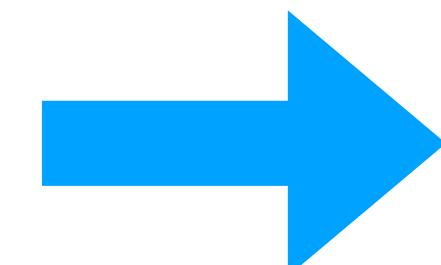
$$\mathcal{K} = \{P : \gamma(P, \hat{P}) \leq \epsilon\}$$

Uncertainty set

Ambiguity set

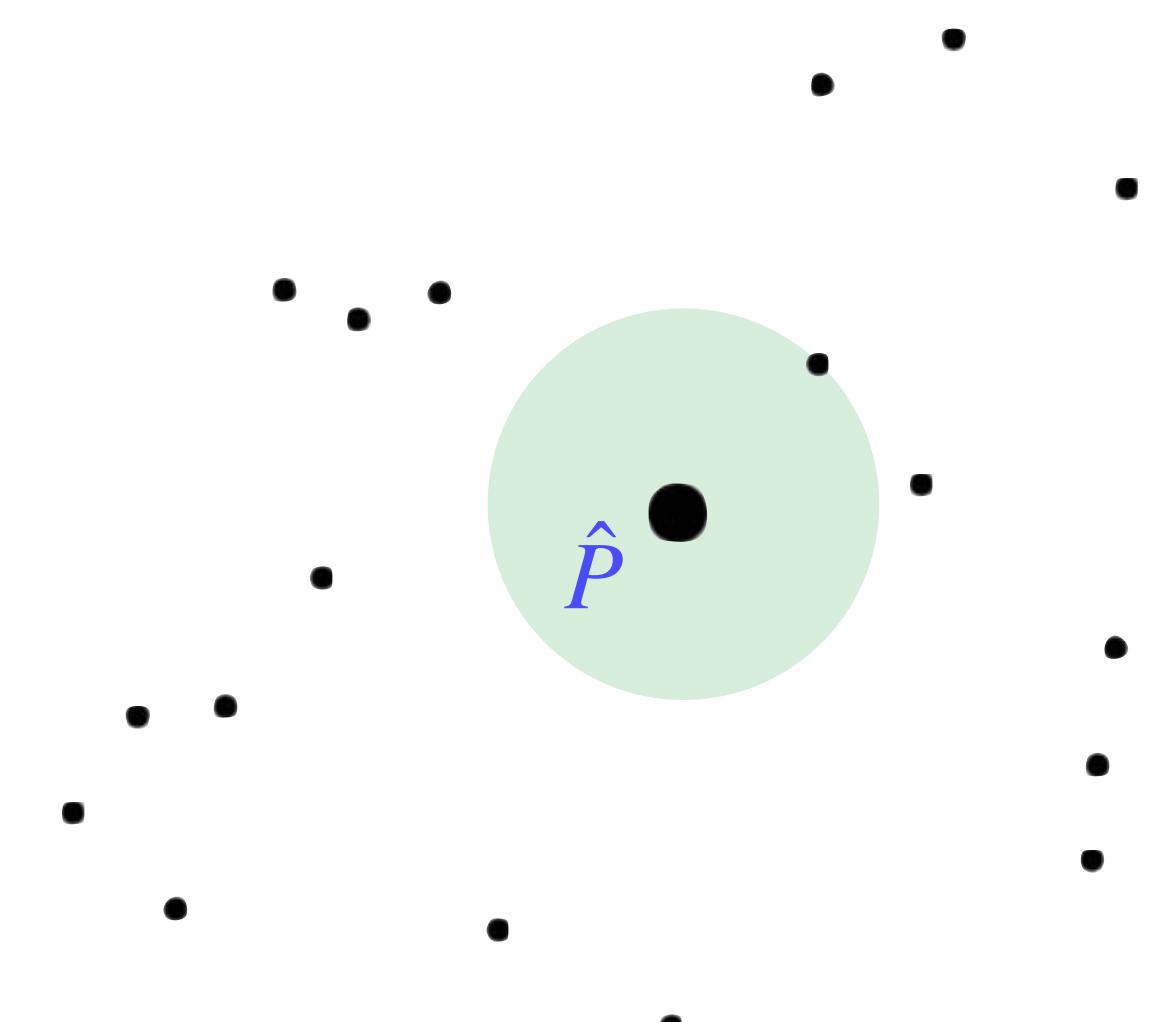
# What distributions does the ambiguity set $\mathcal{K}$ contain?

$$\min_{\theta} \sup_{\xi \in \mathcal{X}} l(\theta, \xi) \quad (\text{RO})$$



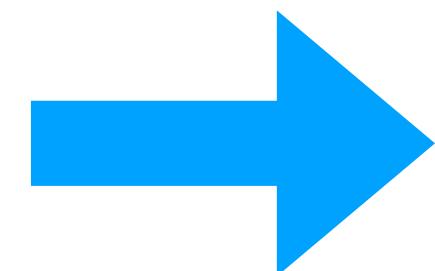
$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad (\text{DRO})$$

$$\mathcal{K} = \{P : \gamma(P, \hat{P}) \leq \epsilon\}$$



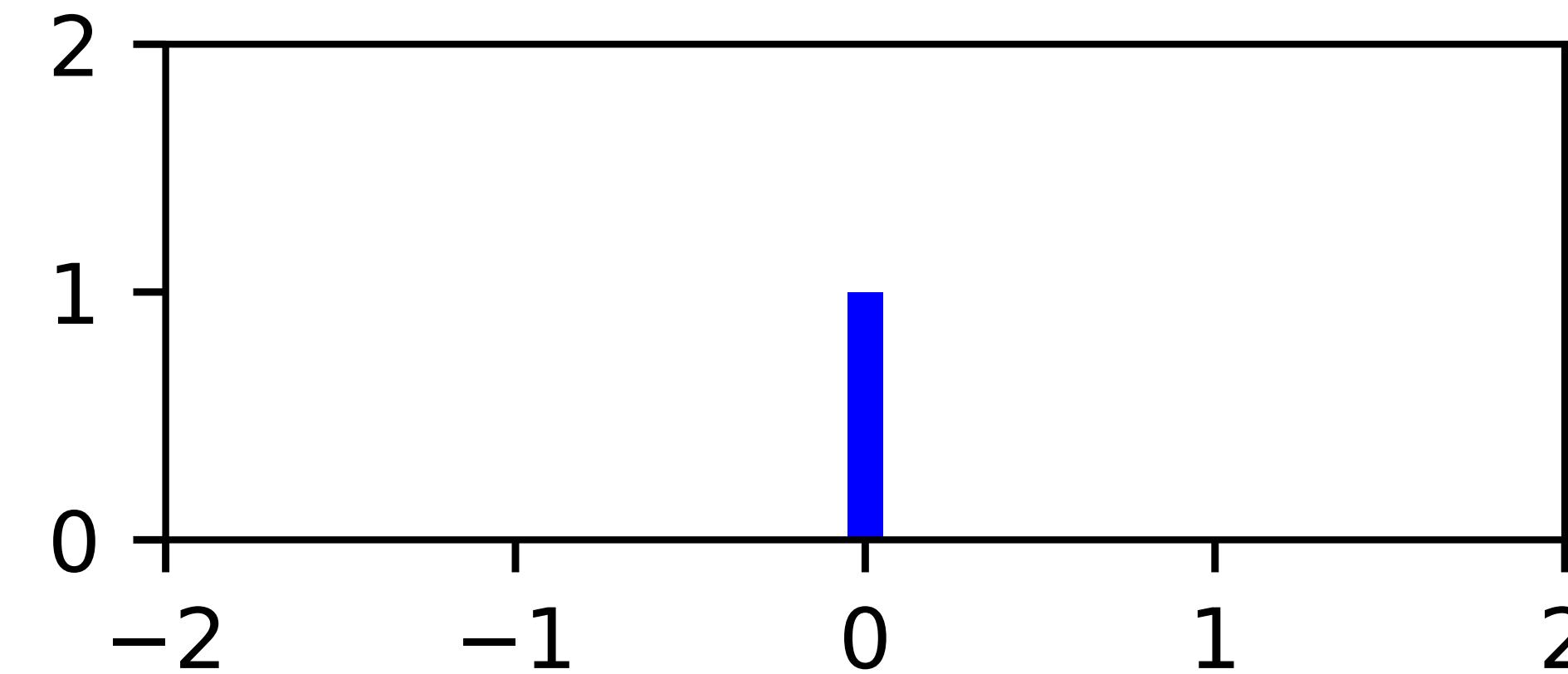
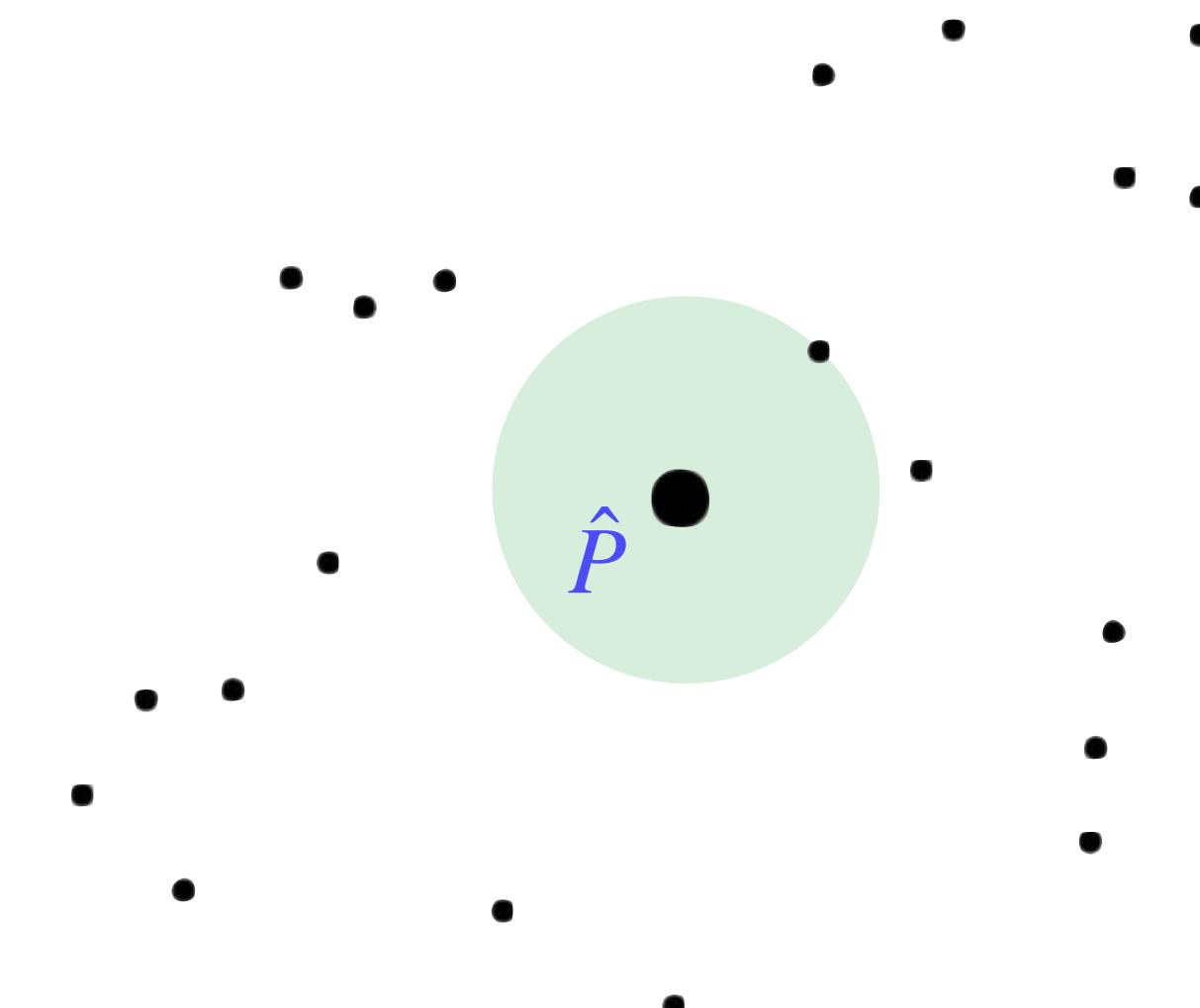
# What distributions does the ambiguity set $\mathcal{K}$ contain?

$$\min_{\theta} \sup_{\xi \in \mathcal{X}} l(\theta, \xi) \quad (\text{RO})$$



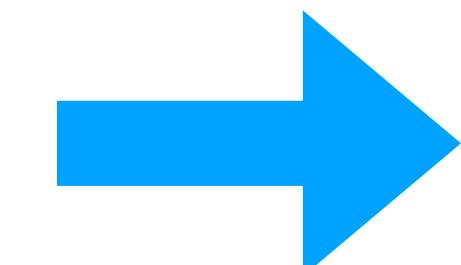
$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad (\text{DRO})$$

$$\mathcal{K} = \{P : \gamma(P, \hat{P}) \leq \epsilon\}$$



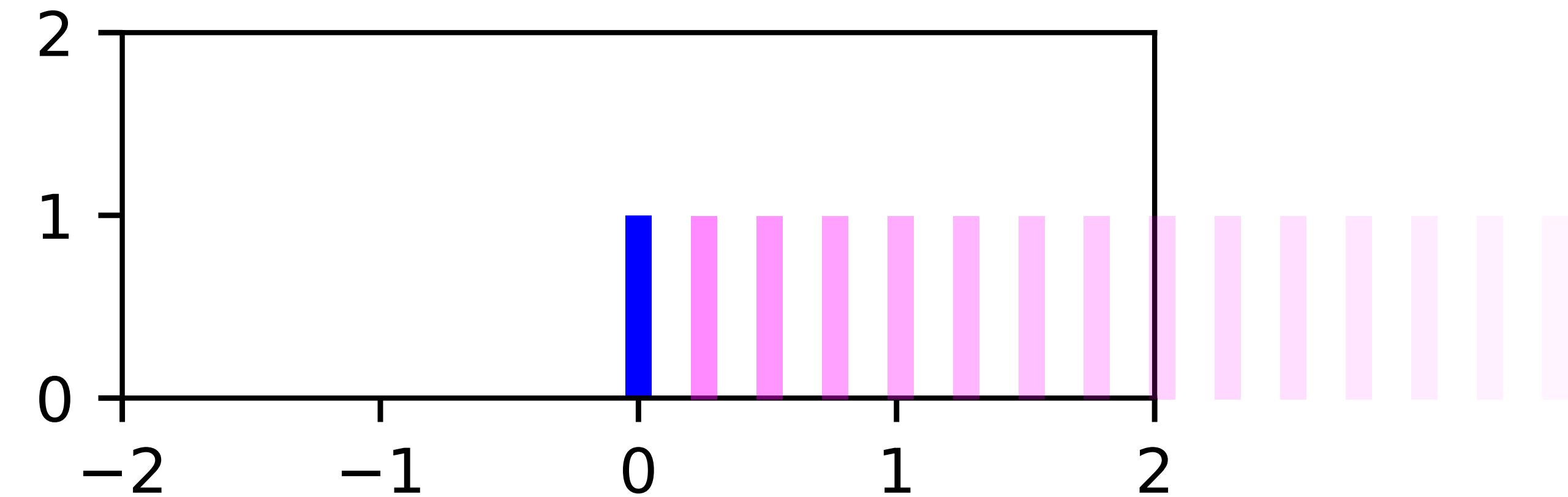
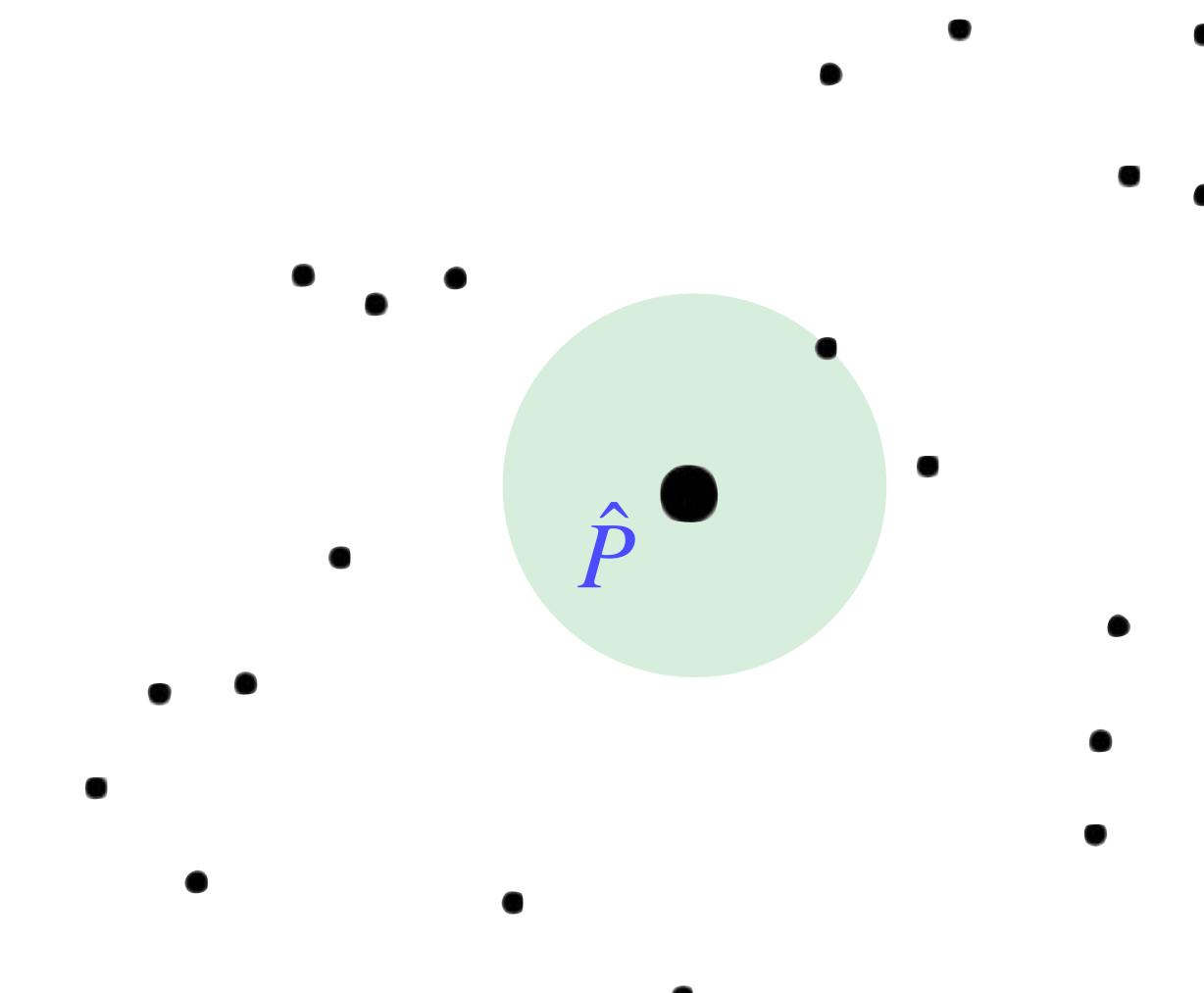
# What distributions does the ambiguity set $\mathcal{K}$ contain?

$$\min_{\theta} \sup_{\xi \in \mathcal{X}} l(\theta, \xi) \quad (\text{RO})$$



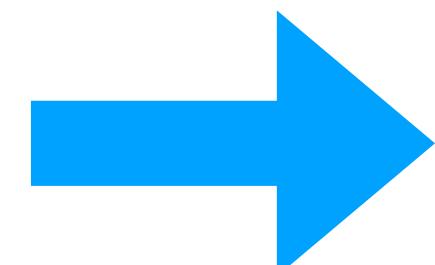
$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad (\text{DRO})$$

$$\mathcal{K} = \{P : \gamma(P, \hat{P}) \leq \epsilon\}$$



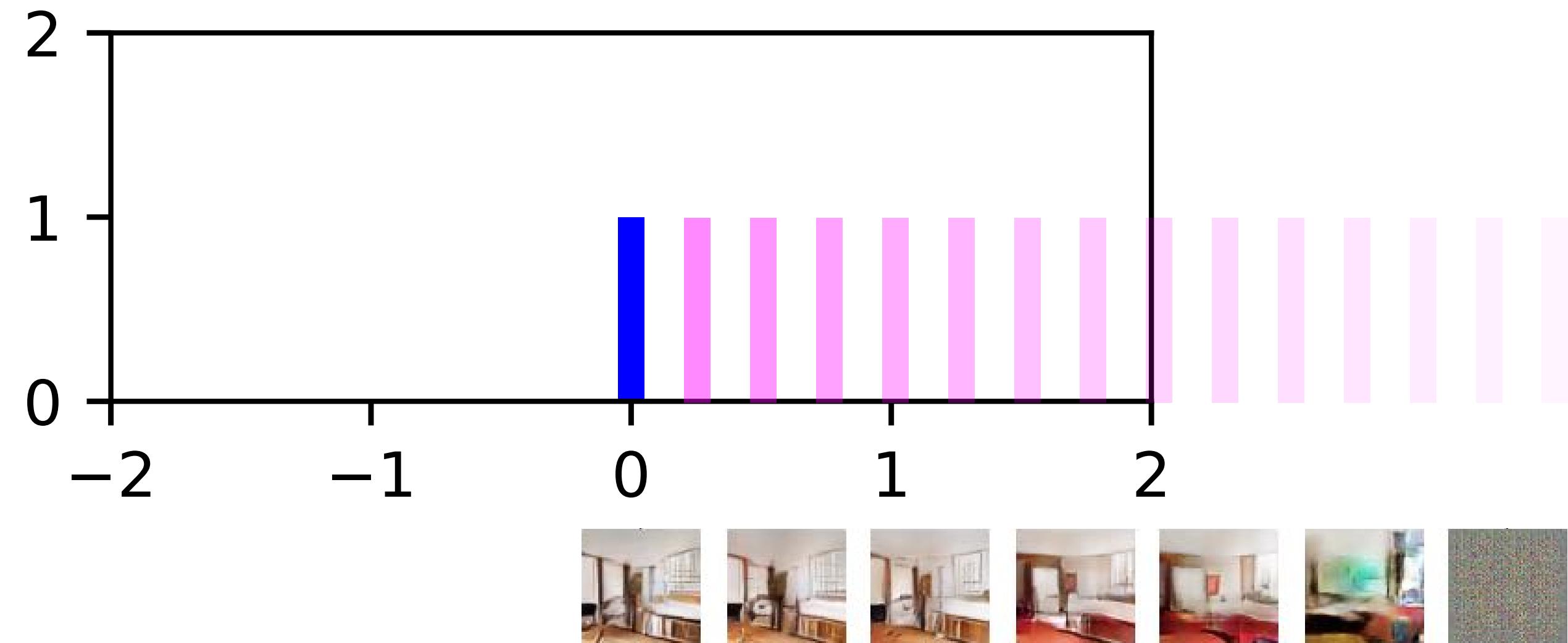
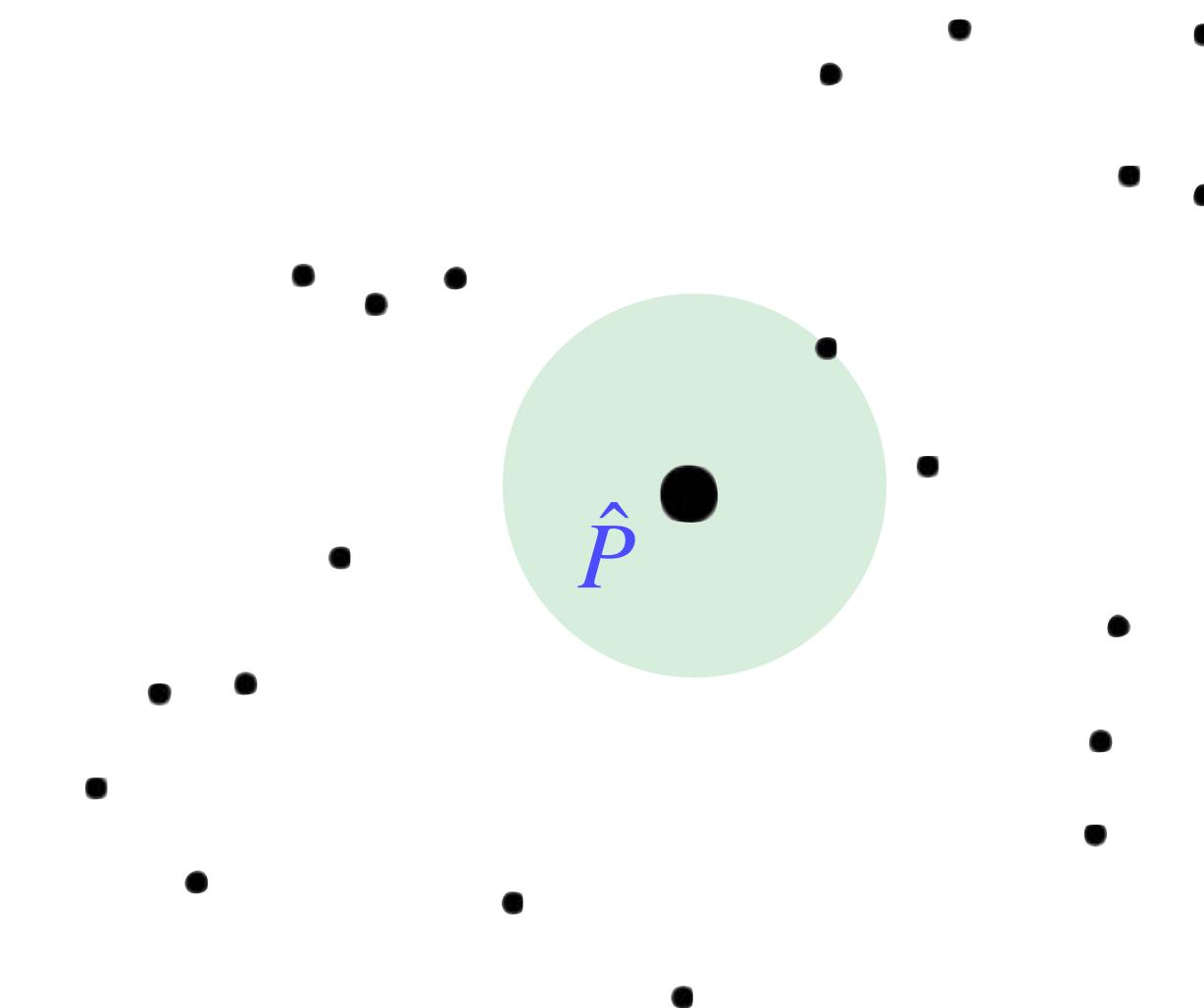
# What distributions does the ambiguity set $\mathcal{K}$ contain?

$$\min_{\theta} \sup_{\xi \in \mathcal{X}} l(\theta, \xi) \quad (\text{RO})$$



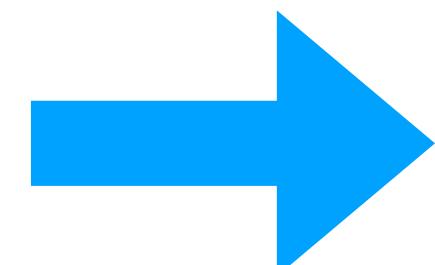
$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad (\text{DRO})$$

$$\mathcal{K} = \{P : \gamma(P, \hat{P}) \leq \epsilon\}$$

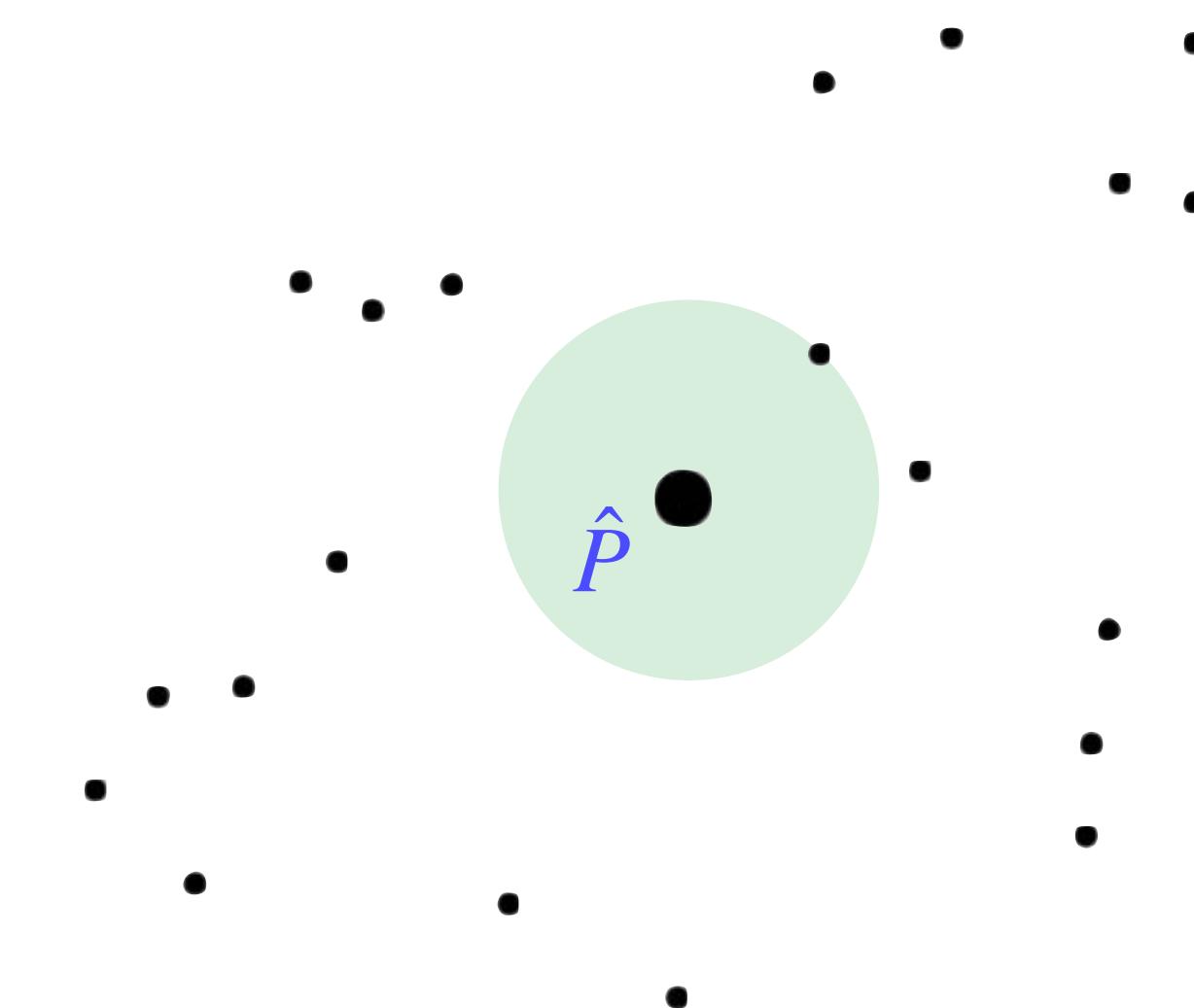


# What distributions does the ambiguity set $\mathcal{K}$ contain?

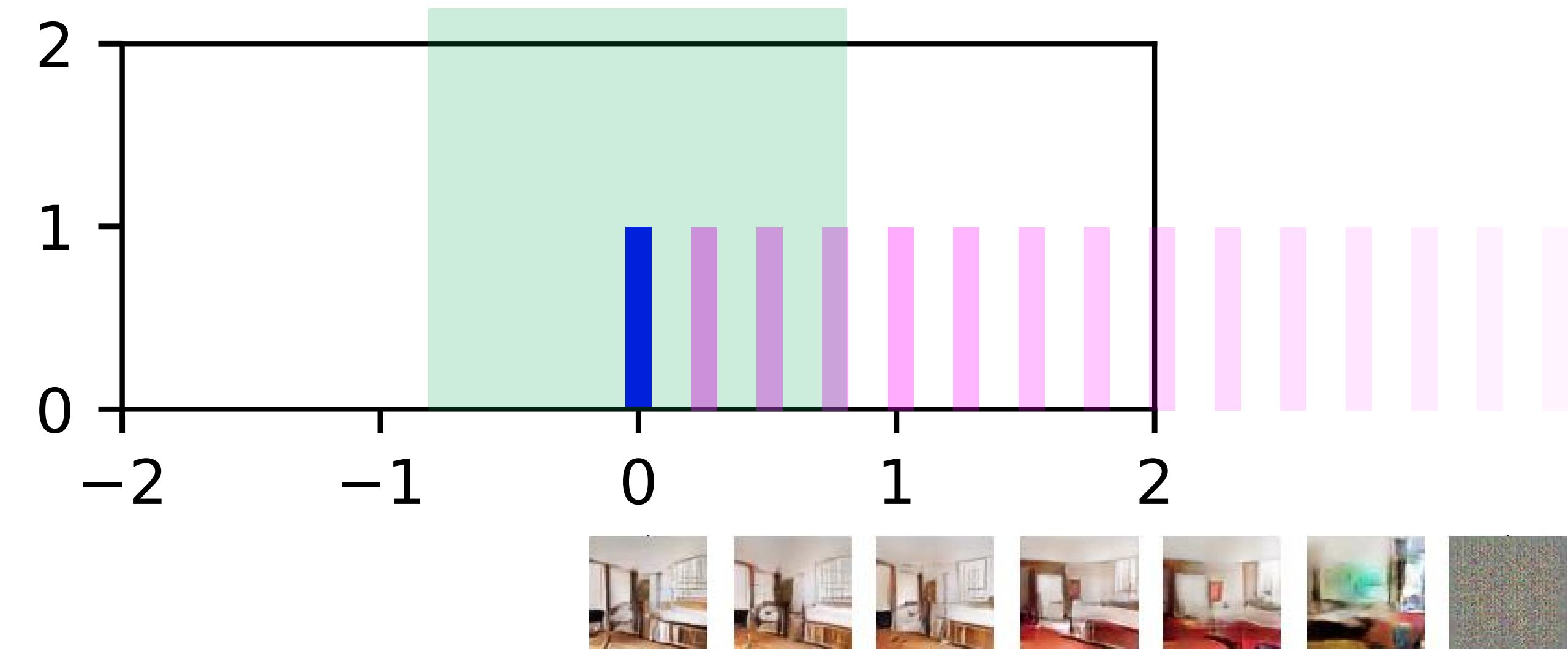
$$\min_{\theta} \sup_{\xi \in \mathcal{X}} l(\theta, \xi) \quad (\text{RO})$$



$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad (\text{DRO})$$

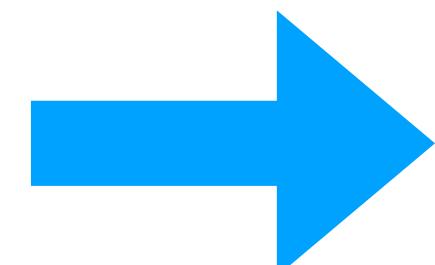


$$\mathcal{K} = \{P : \gamma(P, \hat{P}) \leq \epsilon\}$$

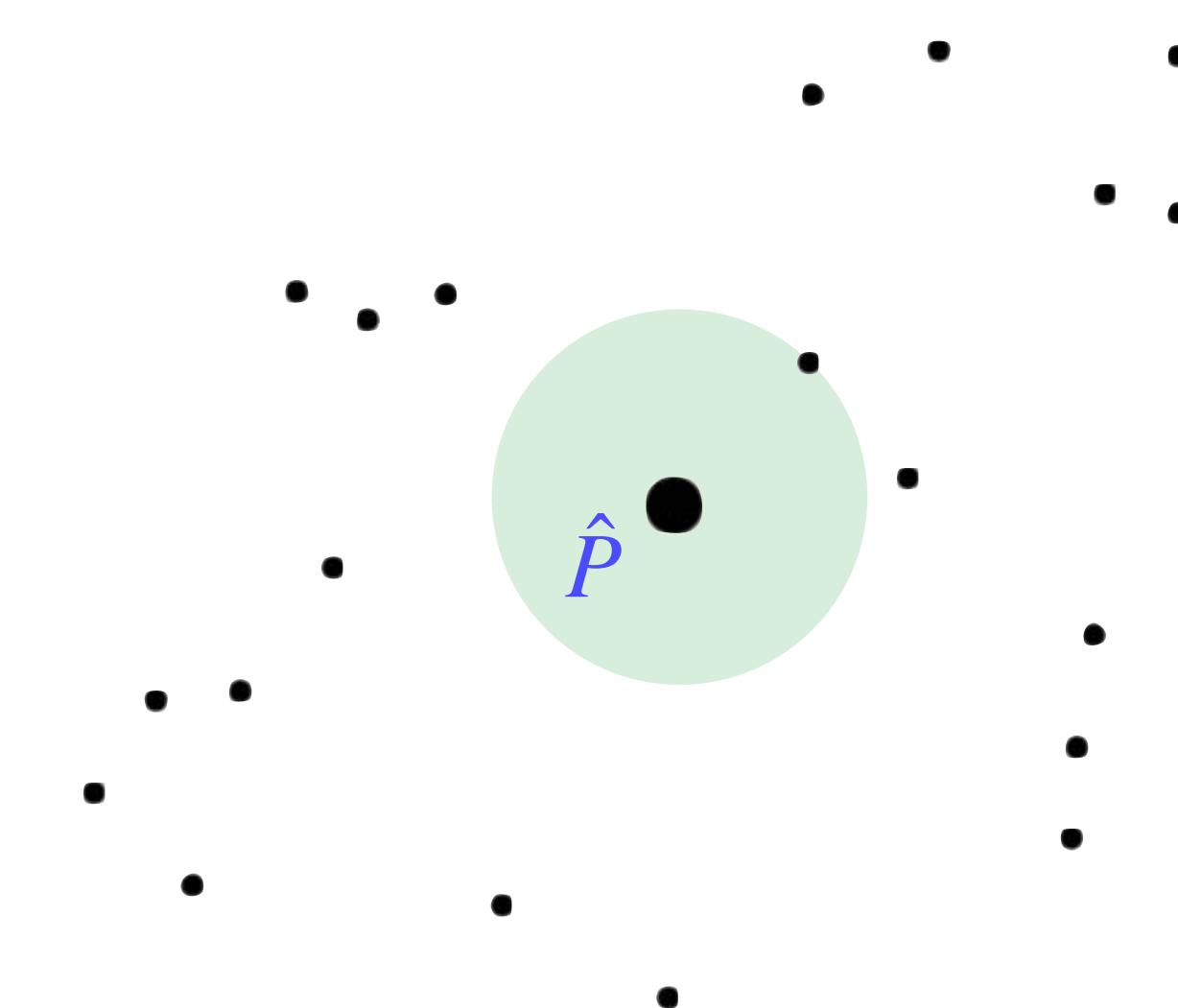


# What distributions does the ambiguity set $\mathcal{K}$ contain?

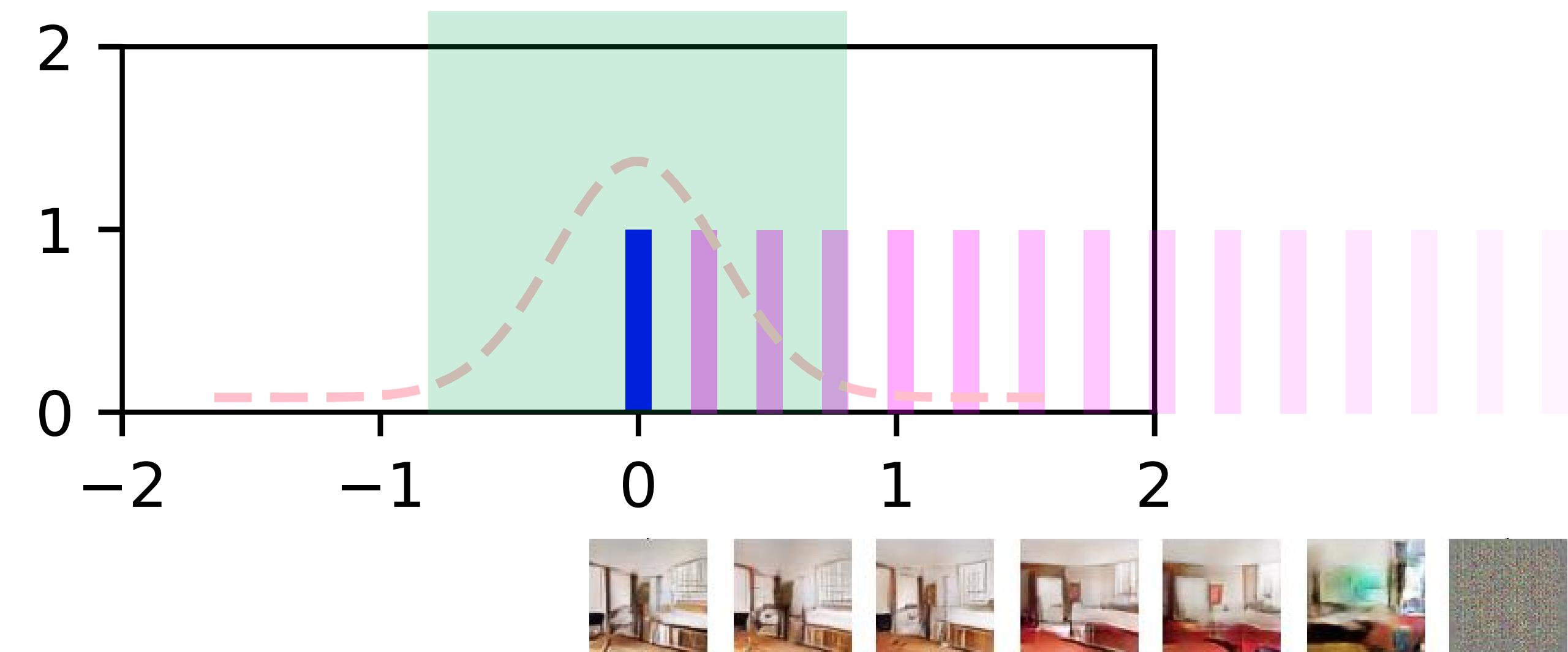
$$\min_{\theta} \sup_{\xi \in \mathcal{X}} l(\theta, \xi) \quad (\text{RO})$$



$$\min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi) \quad (\text{DRO})$$



$$\mathcal{K} = \{P : \gamma(P, \hat{P}) \leq \epsilon\}$$



# Choosing different ambiguity sets

# Choosing different ambiguity sets

$$\text{(DRO)} \quad \min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi)$$

# Choosing different ambiguity sets

$$\text{(DRO)} \quad \min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi)$$

Examples of ambiguity set  $\mathcal{K} \subseteq \mathcal{P}$

# Choosing different ambiguity sets

$$\text{(DRO)} \quad \min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi)$$

Examples of ambiguity set  $\mathcal{K} \subseteq \mathcal{P}$

- singleton  $\mathcal{K} = \left\{ \frac{1}{N} \sum_{i=1}^N \delta_{\xi_i} \right\}$ : ERM

# Choosing different ambiguity sets

$$\text{(DRO)} \quad \min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi)$$

Examples of ambiguity set  $\mathcal{K} \subseteq \mathcal{P}$

- singleton  $\mathcal{K} = \{\frac{1}{N} \sum_{i=1}^N \delta_{\xi_i}\}$ : ERM
- entire space  $\mathcal{K} = \mathcal{P}$ : RO

# Choosing different ambiguity sets

$$\text{(DRO)} \quad \min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi)$$

Examples of ambiguity set  $\mathcal{K} \subseteq \mathcal{P}$

- singleton  $\mathcal{K} = \{\frac{1}{N} \sum_{i=1}^N \delta_{\xi_i}\}$ : ERM
- entire space  $\mathcal{K} = \mathcal{P}$ : RO
- polytope  $\mathcal{K} = \text{conv}\{\delta_{\xi_1}, \dots, \delta_{\xi_N}\}$ : SVMs, scenario opt.

# Choosing different ambiguity sets

$$\text{(DRO)} \quad \min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi)$$

Examples of ambiguity set  $\mathcal{K} \subseteq \mathcal{P}$

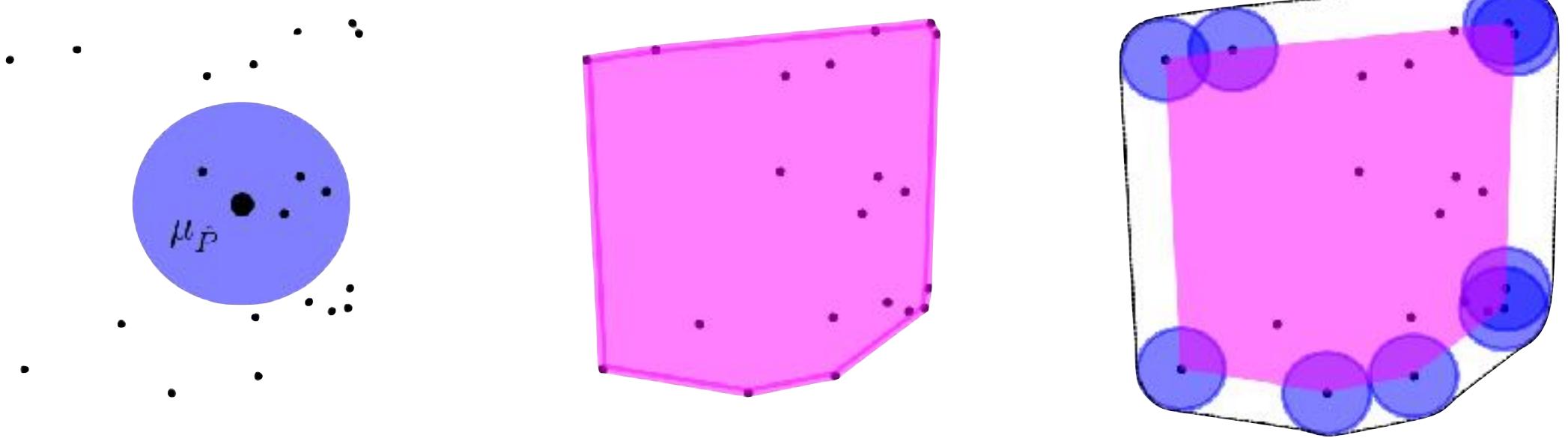
- singleton  $\mathcal{K} = \left\{ \frac{1}{N} \sum_{i=1}^N \delta_{\xi_i} \right\}$ : ERM
- entire space  $\mathcal{K} = \mathcal{P}$ : RO
- polytope  $\mathcal{K} = \text{conv}\{\delta_{\xi_1}, \dots, \delta_{\xi_N}\}$ : SVMs, scenario opt.
- metric-ball  $\mathcal{K} = \{P : \text{MMD}(P, \hat{P}; \mathcal{H}) \leq \epsilon\}$   
 $\mathcal{K} = \{P : W_\alpha(P, \hat{P}) \leq \epsilon\}$

# Choosing different ambiguity sets

$$\text{(DRO)} \quad \min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi)$$

Examples of ambiguity set  $\mathcal{K} \subseteq \mathcal{P}$

- singleton  $\mathcal{K} = \{\frac{1}{N} \sum_{i=1}^N \delta_{\xi_i}\}$ : ERM
- entire space  $\mathcal{K} = \mathcal{P}$ : RO
- polytope  $\mathcal{K} = \text{conv}\{\delta_{\xi_1}, \dots, \delta_{\xi_l}\}$ : SVMs, scenario opt.
- metric-ball  $\mathcal{K} = \{P: \text{MMD}(P, \hat{P}; \mathcal{H}) \leq \epsilon\}$   
 $\mathcal{K} = \{P: W_\alpha(P, \hat{P}) \leq \epsilon\}$

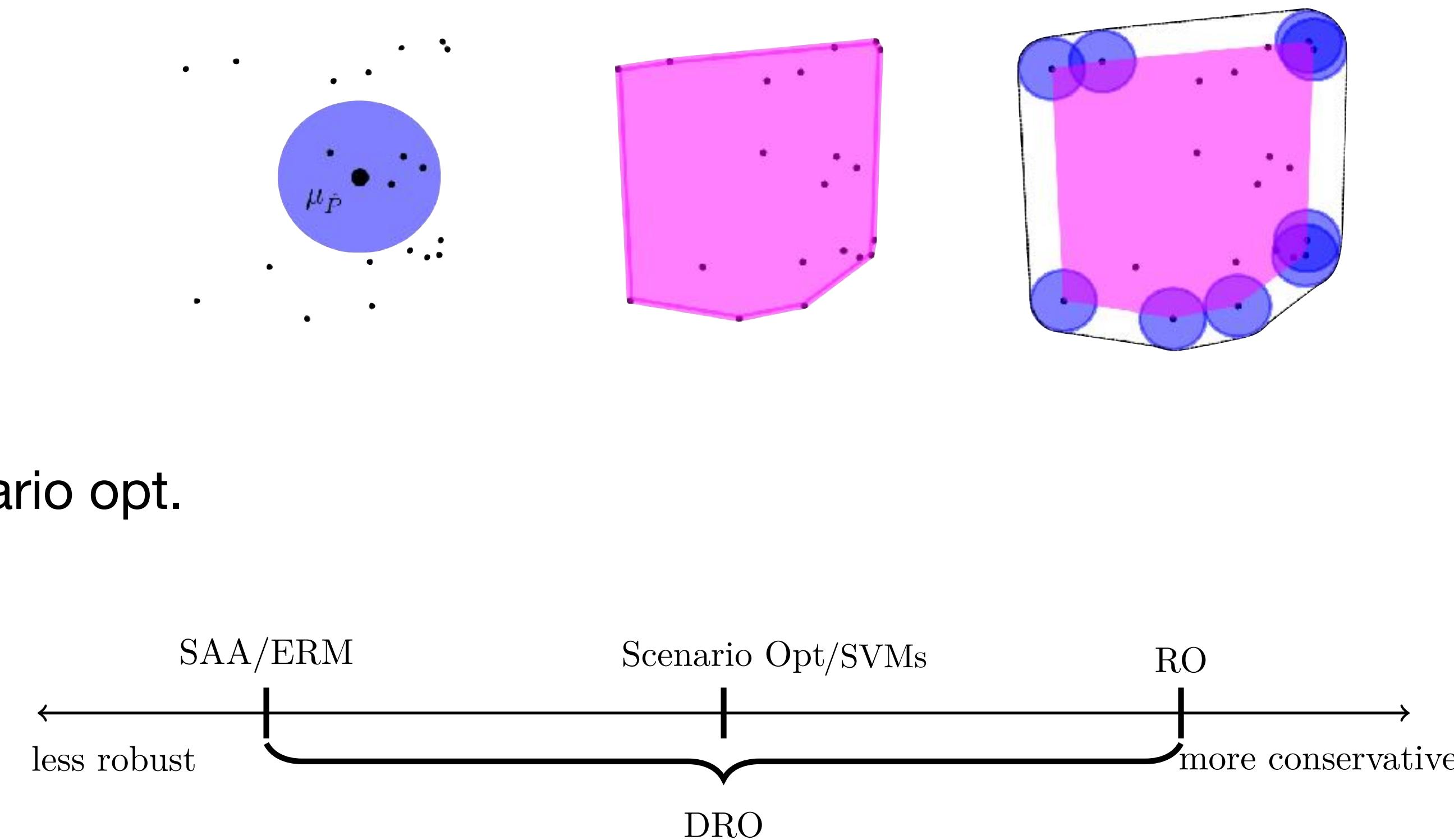


# Choosing different ambiguity sets

$$(DRO) \quad \min_{\theta} \sup_{P \in \mathcal{K}} \mathbb{E}_P l(\theta, \xi)$$

Examples of ambiguity set  $\mathcal{K} \subseteq \mathcal{P}$

- singleton  $\mathcal{K} = \{\frac{1}{N} \sum_{i=1}^N \delta_{\xi_i}\}$ : ERM
- entire space  $\mathcal{K} = \mathcal{P}$ : RO
- polytope  $\mathcal{K} = \text{conv}\{\delta_{\xi_1}, \dots, \delta_{\xi_l}\}$ : SVMs, scenario opt.
- metric-ball  $\mathcal{K} = \{P: \text{MMD}(P, \hat{P}; \mathcal{H}) \leq \epsilon\}$   
 $\mathcal{K} = \{P: W_\alpha(P, \hat{P}) \leq \epsilon\}$



# The lifting trick of DRO

# The lifting trick of DRO

- How do we solve metric-ball constrained DRO?

$$\min_{\theta} \sup_{\gamma(P, \hat{P}) \leq \epsilon} \mathbb{E}_P l(\theta, \xi)$$

# The lifting trick of DRO

- How do we solve metric-ball constrained DRO?

$$\min_{\theta} \sup_{\gamma(P, \hat{P}) \leq \epsilon} \mathbb{E}_P l(\theta, \xi)$$

- Remember the RO we solved earlier?

$$\min_{\theta} \sup_{\|p - \hat{p}\|_a \leq \epsilon} p^T \theta$$

# The lifting trick of DRO

- How do we solve metric-ball constrained DRO?

$$\min_{\theta} \sup_{\gamma(P, \hat{P}) \leq \epsilon} \mathbb{E}_P l(\theta, \xi)$$

- Remember the RO we solved earlier?

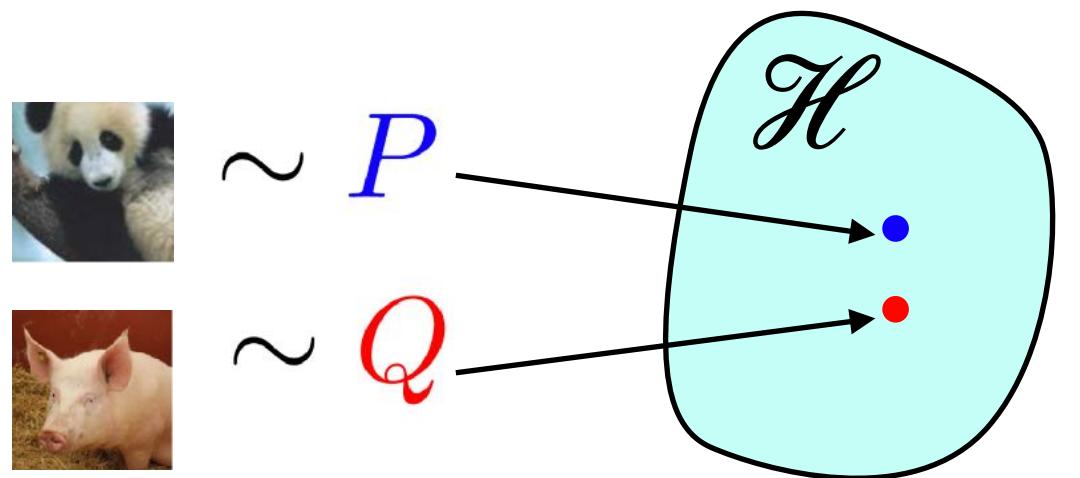
$$\min_{\theta} \sup_{\|p - \hat{p}\|_a \leq \epsilon} p^T \theta$$

- By comparing them, we find that DRO has the same structure as RO! This motivates us to solve DRO with convex duality (but DRO is  $\infty$ -dimensional!).

# Dual program of DRO

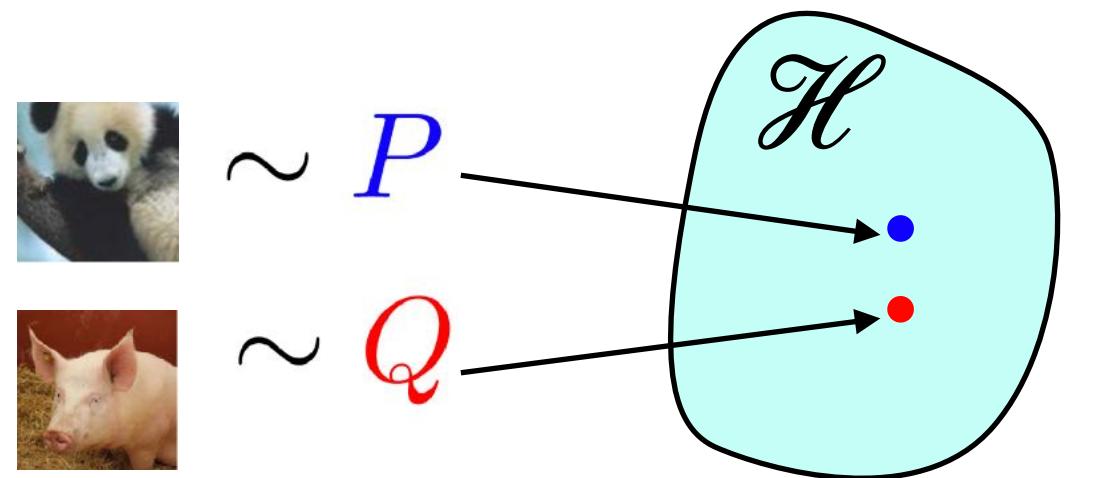
# Dual program of DRO

$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi)$$



# Dual program of DRO

$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi)$$

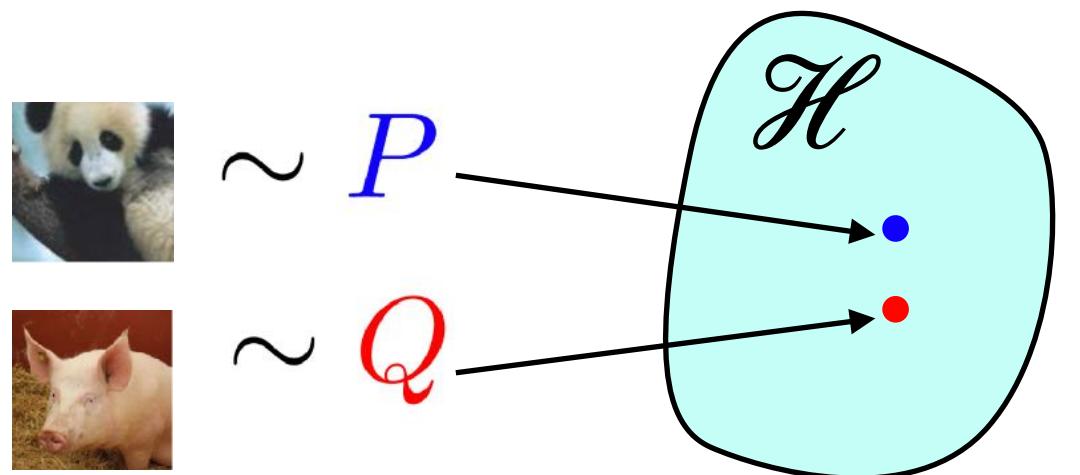


**Theorem** (Simplified: **Generalized DRO duality, Zhu et al. '20**). DRO (P) is equivalent to solving its dual problem

$$(D) \min_{\theta, f \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} f + \epsilon \|f\|_{\mathcal{H}} \quad \text{s.t. } l(\theta, \cdot) \leq f,$$

# Dual program of DRO

$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi)$$

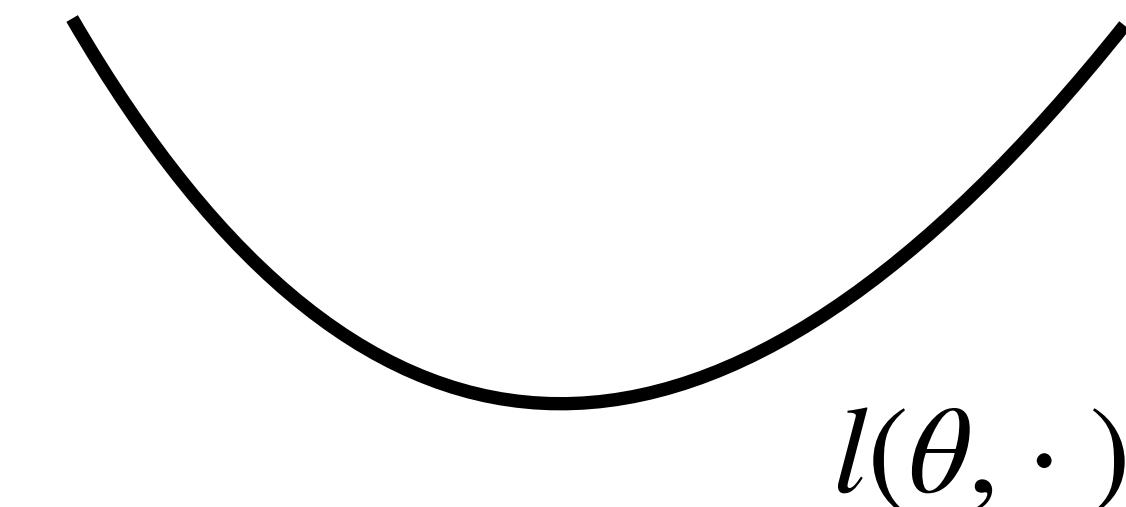


**Theorem** (Simplified: **Generalized DRO duality, Zhu et al.**

'20). DRO (P) is equivalent to solving its dual problem

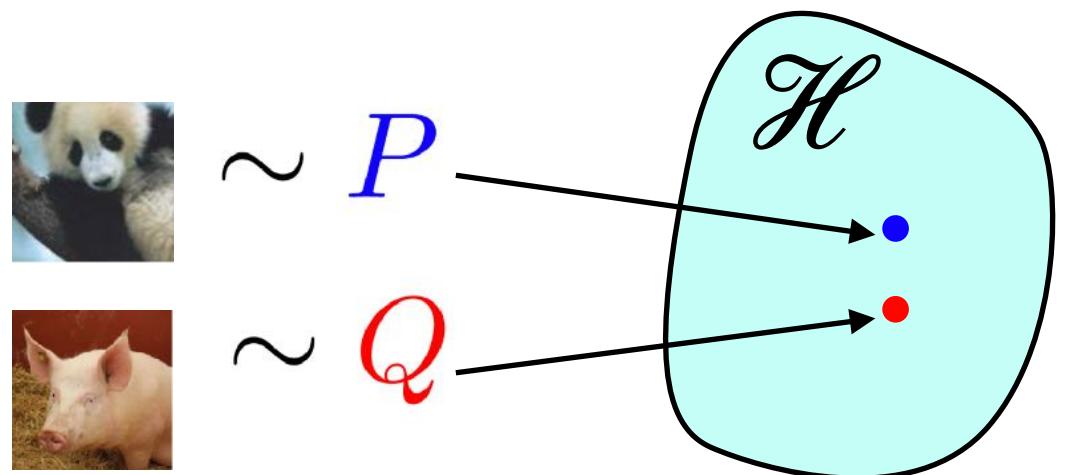
$$(D) \min_{\theta, \mathbf{f} \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} \mathbf{f} + \epsilon \|\mathbf{f}\|_{\mathcal{H}} \quad \text{s.t. } l(\theta, \cdot) \leq \mathbf{f},$$

Geometric intuition



# Dual program of DRO

$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi)$$

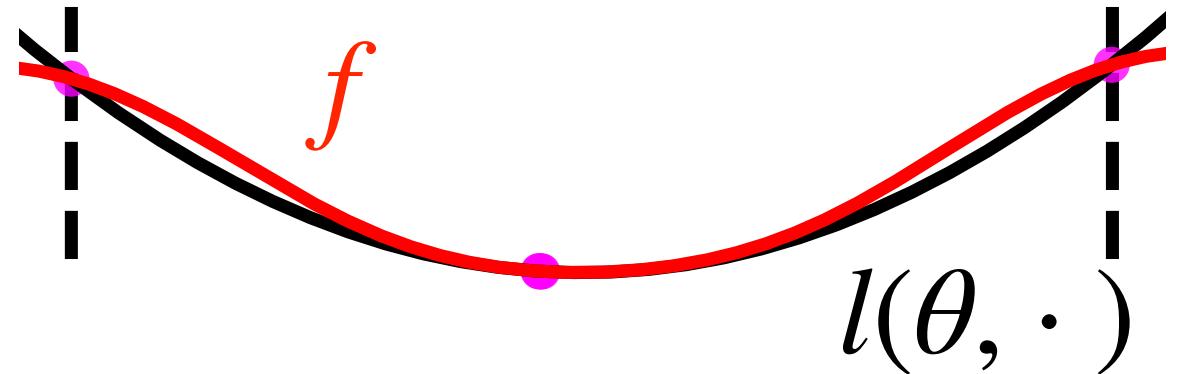


**Theorem** (Simplified: **Generalized DRO duality**, Zhu et al.

'20). DRO (P) is equivalent to solving its dual problem

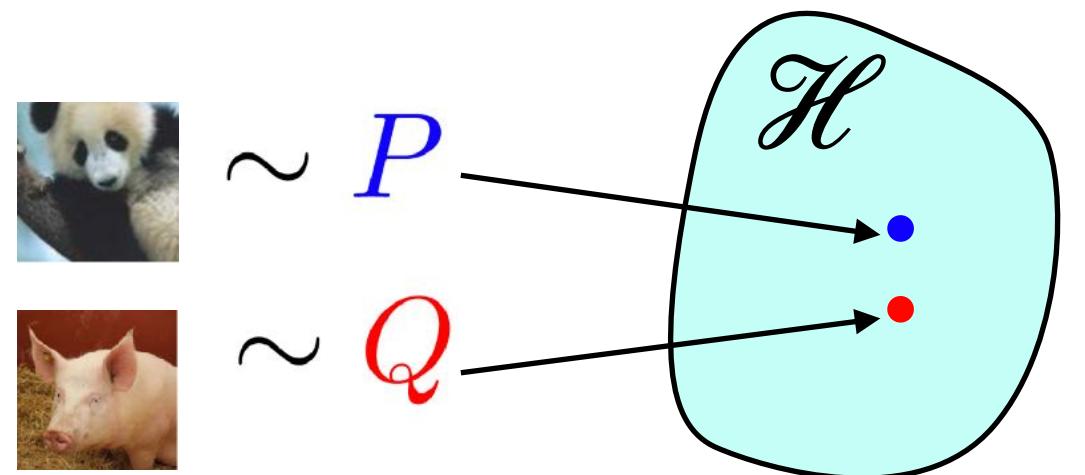
$$(D) \min_{\theta, f \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} f + \epsilon \|f\|_{\mathcal{H}} \quad \text{s.t. } l(\theta, \cdot) \leq f,$$

Geometric intuition



# Dual program of DRO

$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi)$$

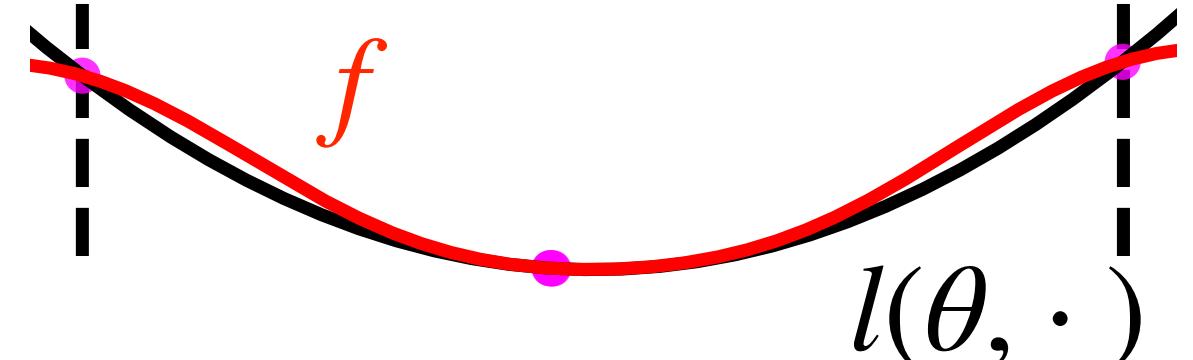


**Theorem** (Simplified: **Generalized DRO duality**, Zhu et al.

'20). DRO (P) is equivalent to solving its dual problem

$$(D) \min_{\theta, f \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} f + \epsilon \|f\|_{\mathcal{H}} \quad \text{s.t. } l(\theta, \cdot) \leq f,$$

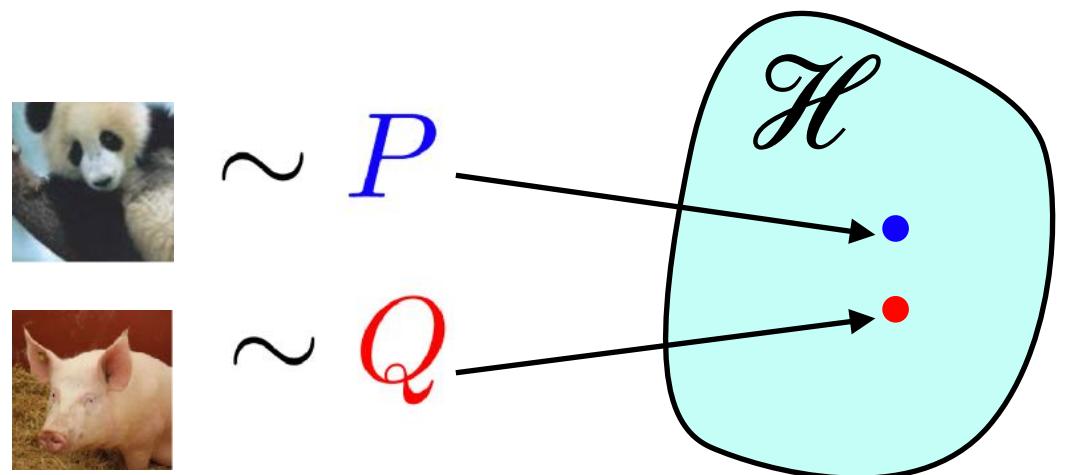
Geometric intuition



Smoothness of  $f \leftrightarrow$  Distributional robustness ( $\leftrightarrow$  Size of  $\mathcal{H}$ )

# Dual program of DRO

$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi)$$

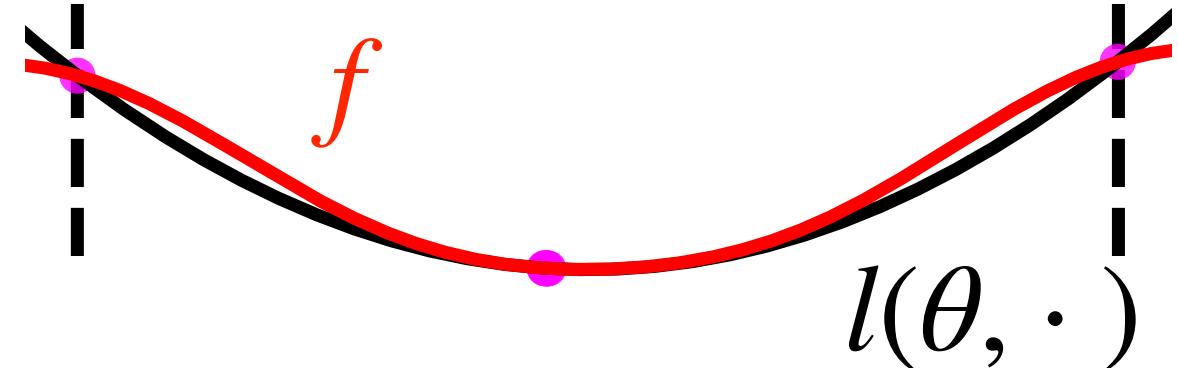


**Theorem** (Simplified: **Generalized DRO duality**, Zhu et al.

'20). DRO (P) is equivalent to solving its dual problem

$$(D) \min_{\theta, f \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} f + \epsilon \|f\|_{\mathcal{H}} \quad \text{s.t. } l(\theta, \cdot) \leq f,$$

Geometric intuition

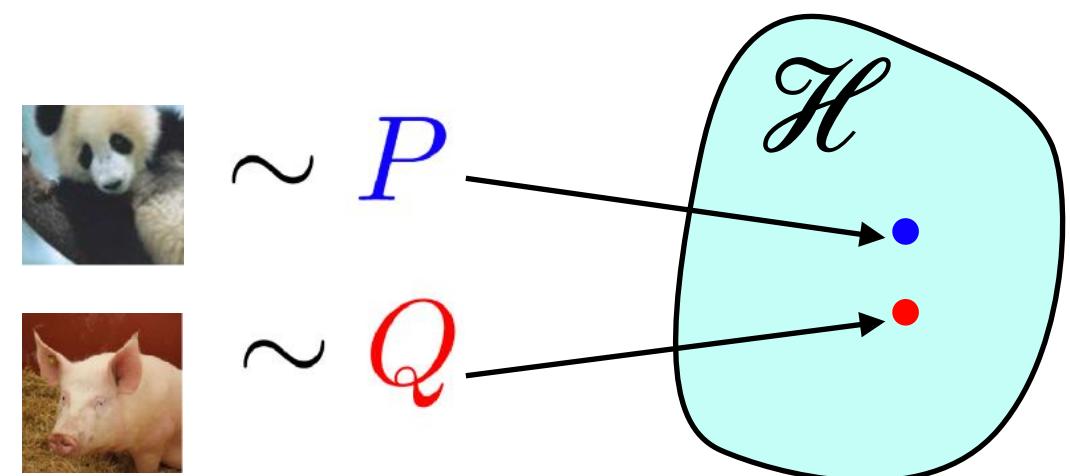


Smoothness of  $f \leftrightarrow$  Distributional robustness ( $\leftrightarrow$  Size of  $\mathcal{H}$ )

Intuition: flatten the curve, smooth is robust

# Dual program of DRO

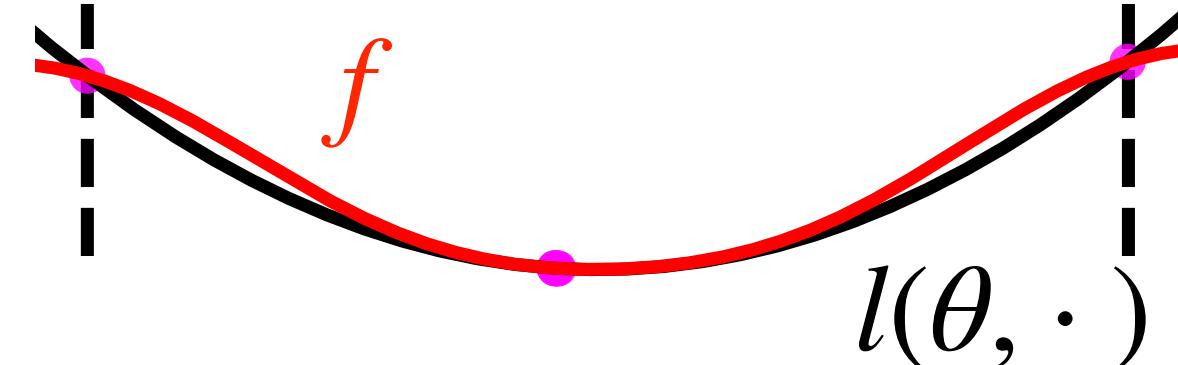
$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi)$$



**Theorem** (Simplified: **Generalized DRO duality**, Zhu et al. '20). DRO (P) is equivalent to solving its dual problem

$$(D) \min_{\theta, f \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} f + \epsilon \|f\|_{\mathcal{H}} \quad \text{s.t. } l(\theta, \cdot) \leq f,$$

Geometric intuition



Smoothness of  $f \leftrightarrow$  Distributional robustness ( $\leftrightarrow$  Size of  $\mathcal{H}$ )

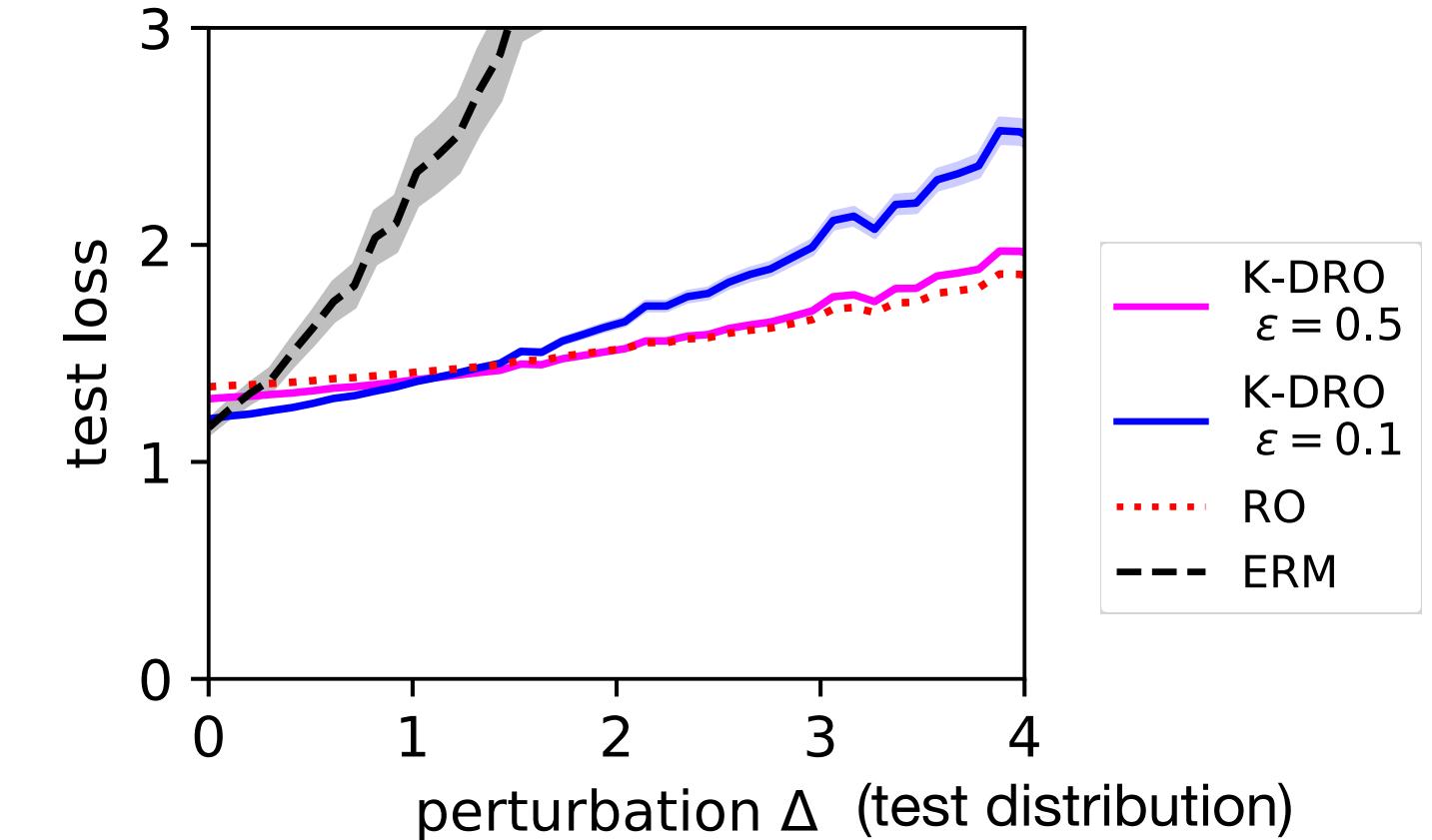
Intuition: flatten the curve, smooth is robust

**Example. Robust least squares**

[El Ghaoui Lebret '97]

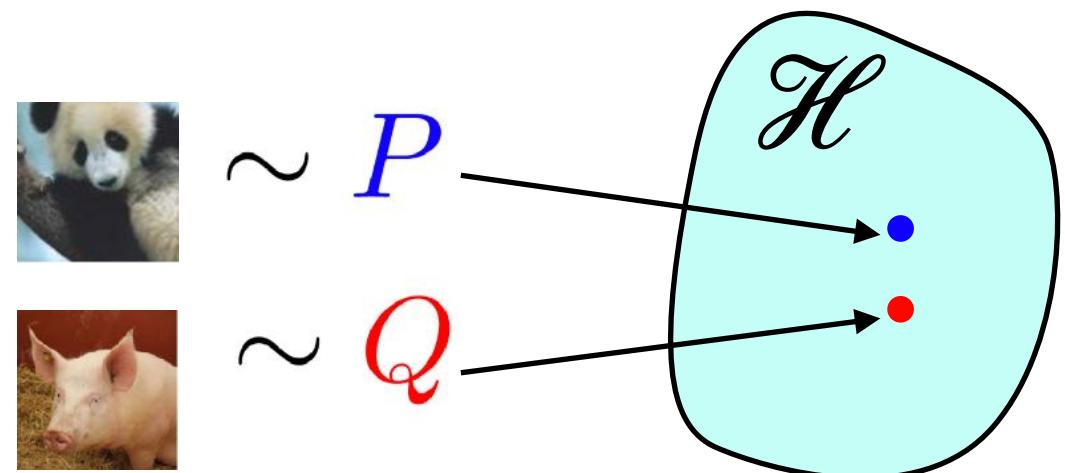
$$\text{minimize } l(\theta, \xi) := \|A(\xi) \cdot \theta - b\|_2^2$$

Given historical samples  $\xi_1, \xi_2, \dots, \xi_N$



# Dual program of DRO

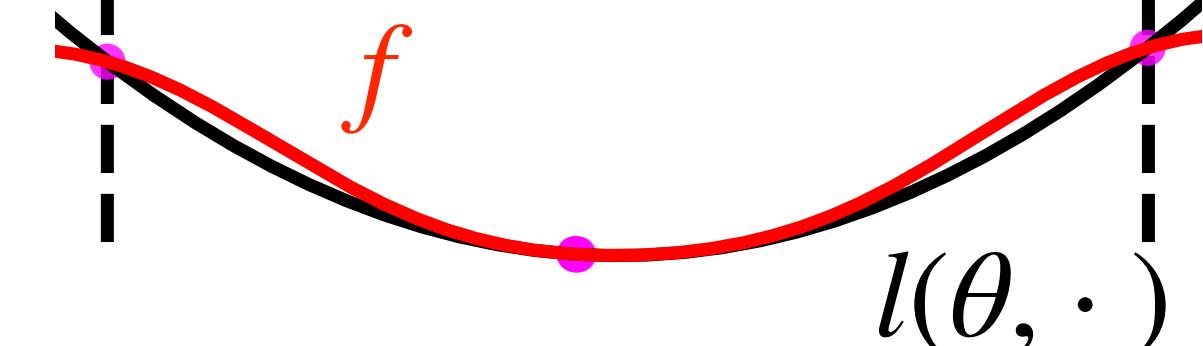
$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi)$$



**Theorem** (Simplified: **Generalized DRO duality**, Zhu et al. '20). DRO (P) is equivalent to solving its dual problem

$$(D) \min_{\theta, f \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} f + \epsilon \|f\|_{\mathcal{H}} \quad \text{s.t. } l(\theta, \cdot) \leq f,$$

Geometric intuition



Smoothness of  $f \leftrightarrow$  Distributional robustness ( $\leftrightarrow$  Size of  $\mathcal{H}$ )

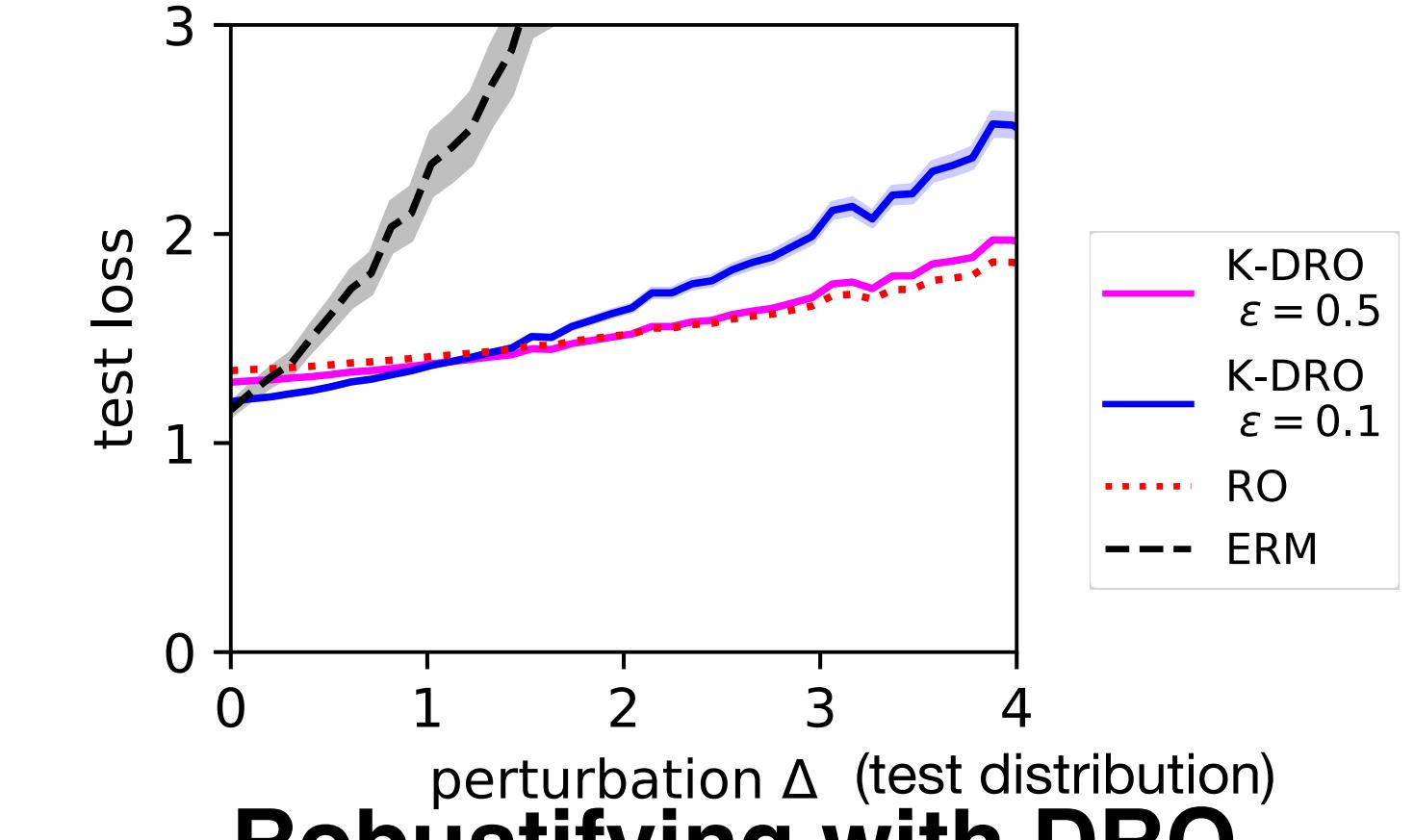
Intuition: flatten the curve, smooth is robust

## Example. Robust least squares

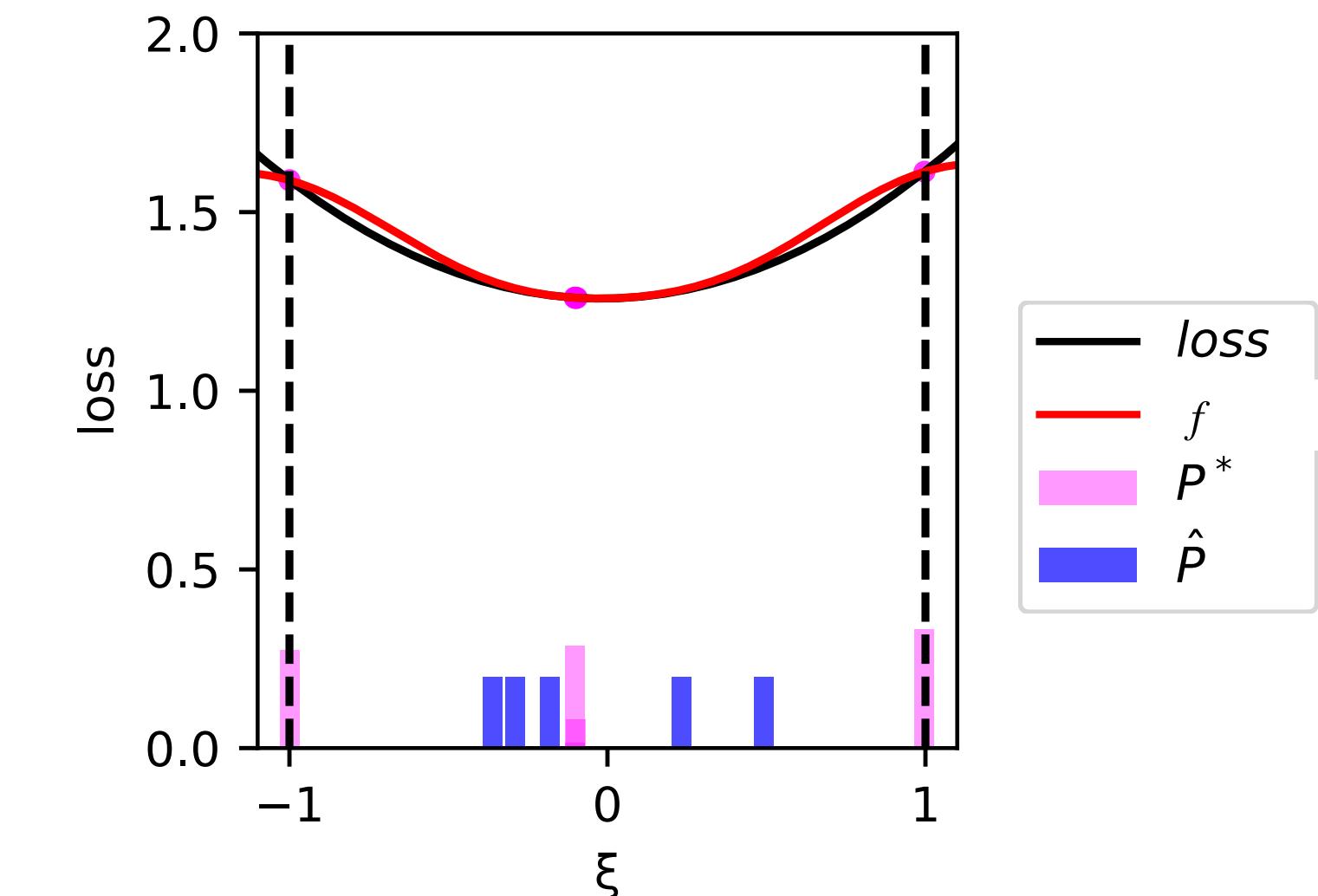
[El Ghaoui Lebret '97]

$$\text{minimize } l(\theta, \xi) := \|A(\xi) \cdot \theta - b\|_2^2$$

Given historical samples  $\xi_1, \xi_2, \dots, \xi_N$

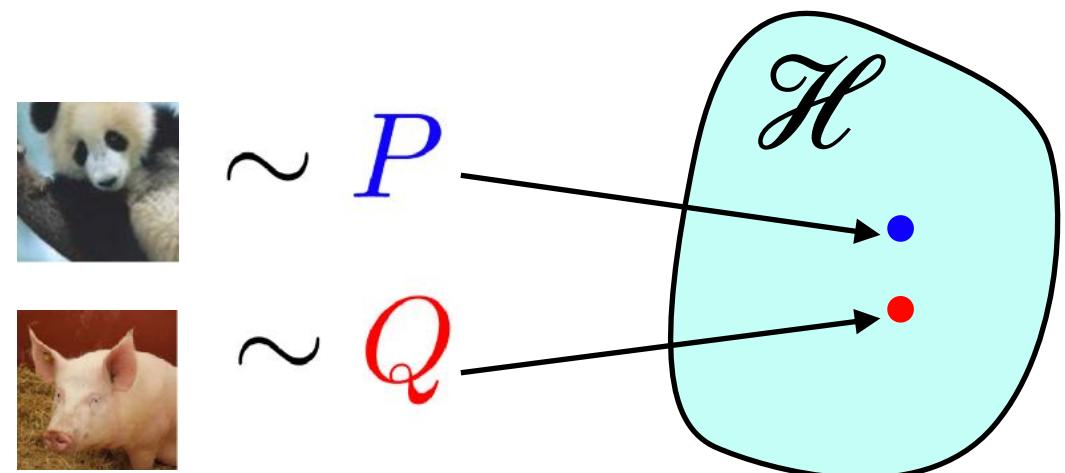


## Robustifying with DRO



# Dual program of DRO

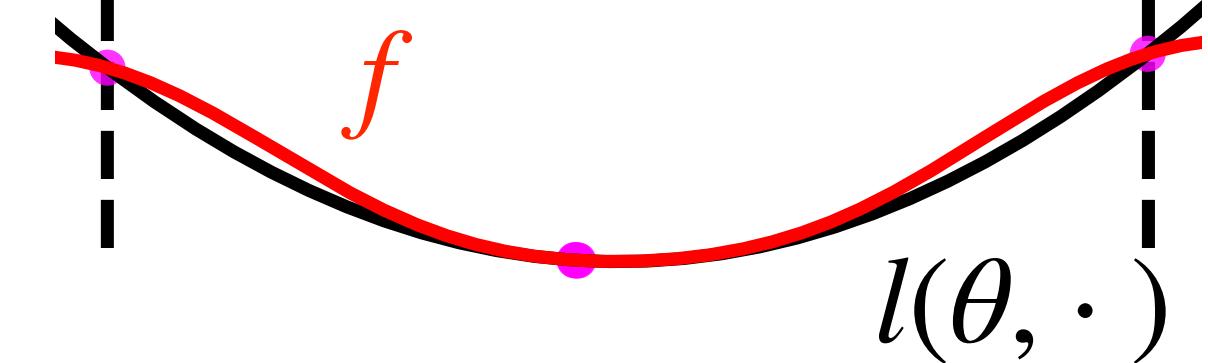
$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi)$$



**Theorem** (Simplified: **Generalized DRO duality**, Zhu et al. '20). DRO (P) is equivalent to solving its dual problem

$$(D) \min_{\theta, f \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} f + \epsilon \|f\|_{\mathcal{H}} \quad \text{s.t. } l(\theta, \cdot) \leq f,$$

Geometric intuition



Smoothness of  $f \leftrightarrow$  Distributional robustness ( $\leftrightarrow$  Size of  $\mathcal{H}$ )

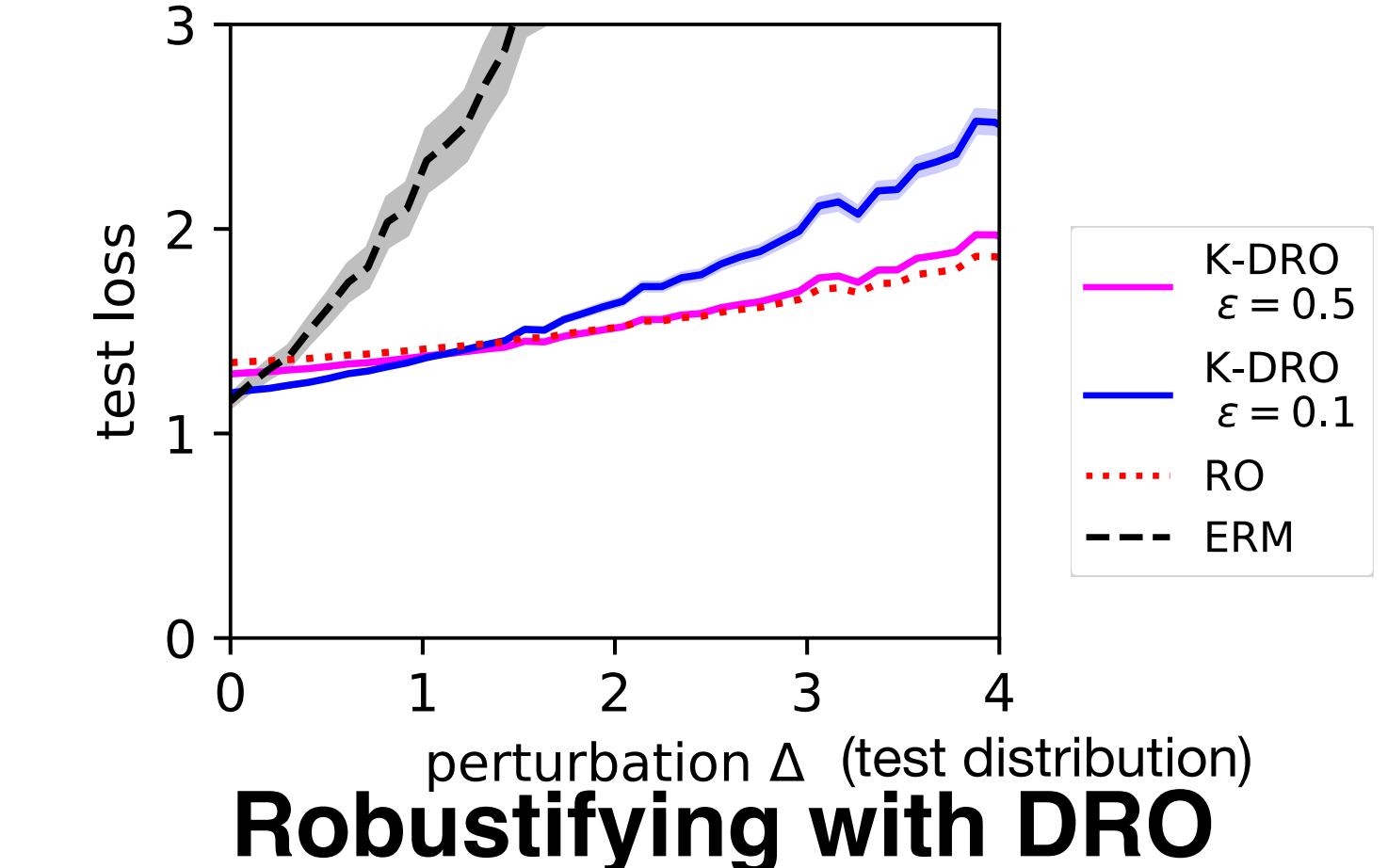
Intuition: flatten the curve, smooth is robust

## Example. Robust least squares

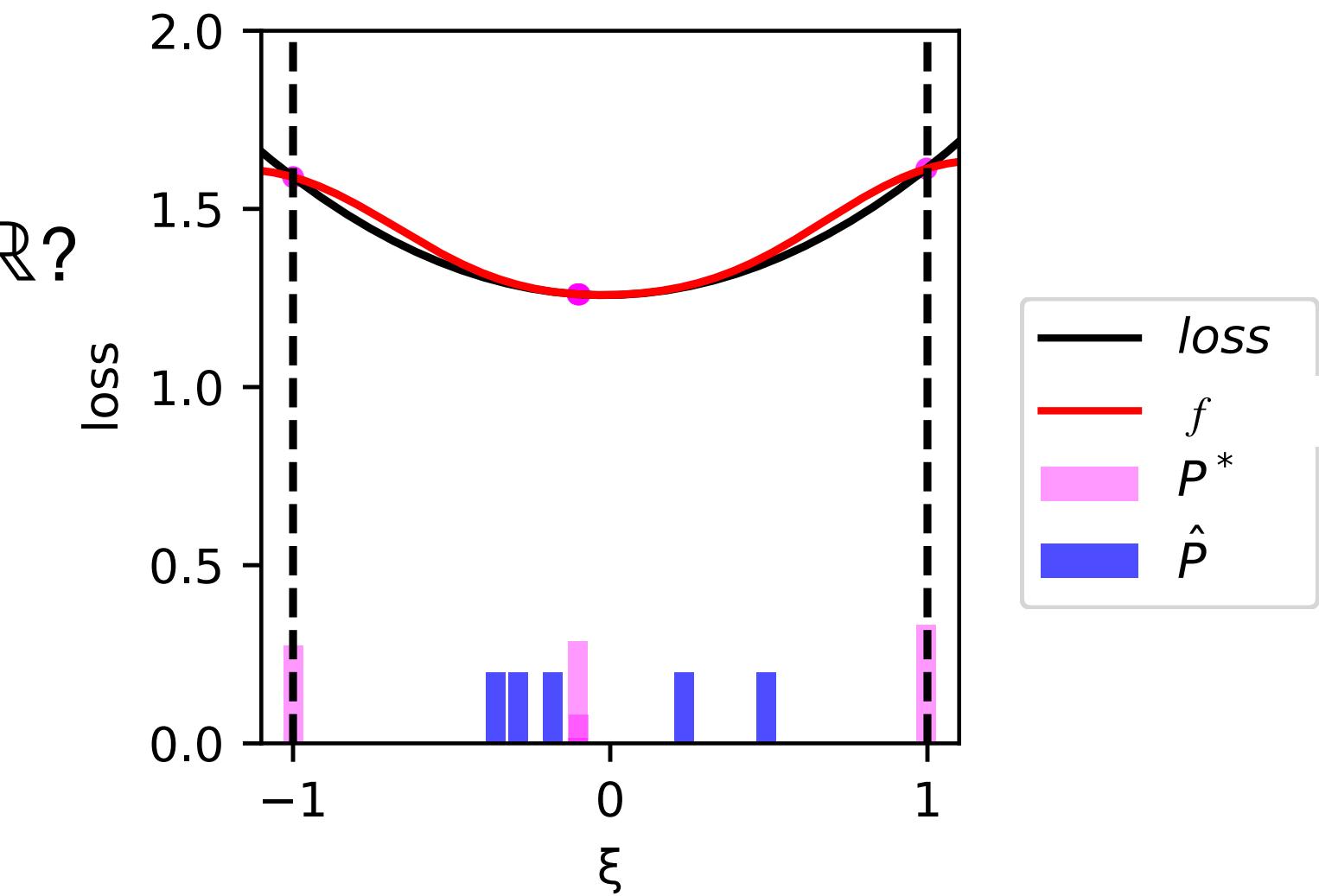
[El Ghaoui Lebret '97]

$$\text{minimize } l(\theta, \xi) := \|A(\xi) \cdot \theta - b\|_2^2$$

Given historical samples  $\xi_1, \xi_2, \dots, \xi_N$

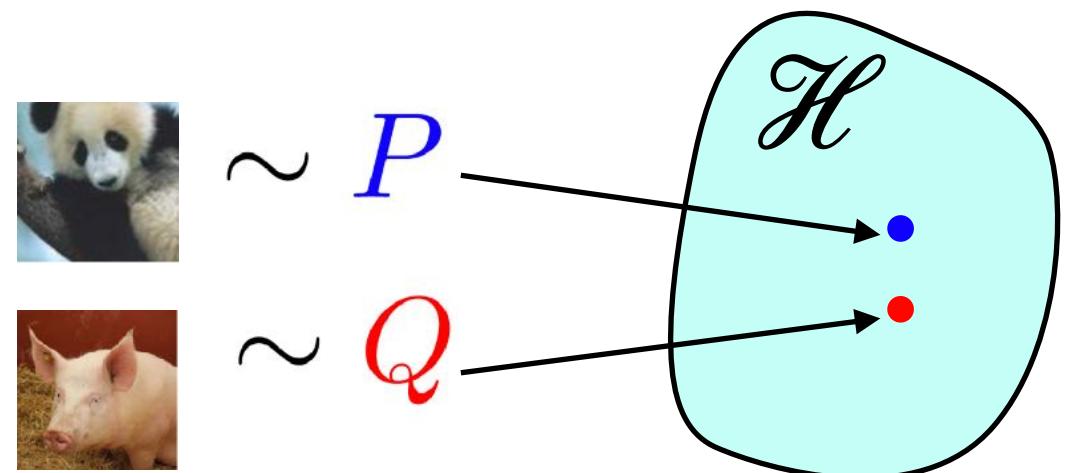


What if  $f \equiv c \in \mathbb{R}$ ?



# Dual program of DRO

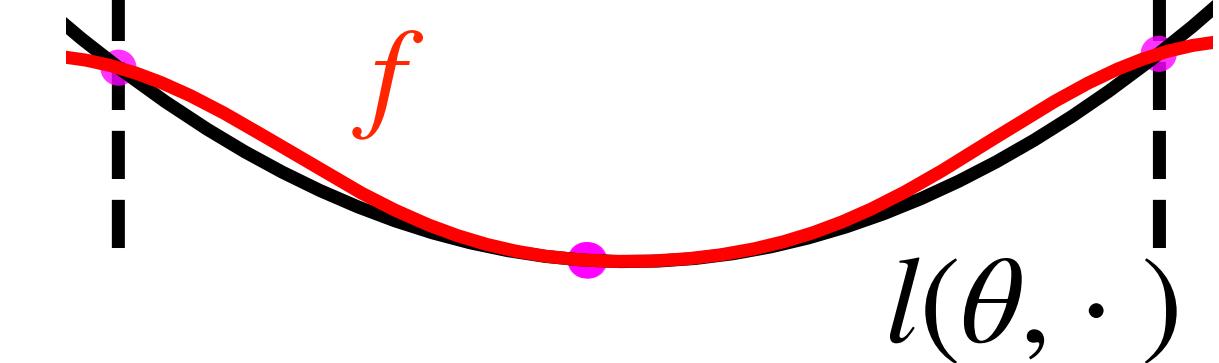
$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi)$$



**Theorem** (Simplified: **Generalized DRO duality**, Zhu et al. '20). DRO (P) is equivalent to solving its dual problem

$$(D) \min_{\theta, f \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} f + \epsilon \|f\|_{\mathcal{H}} \quad \text{s.t. } l(\theta, \cdot) \leq f,$$

Geometric intuition



Smoothness of  $f \leftrightarrow$  Distributional robustness ( $\leftrightarrow$  Size of  $\mathcal{H}$ )

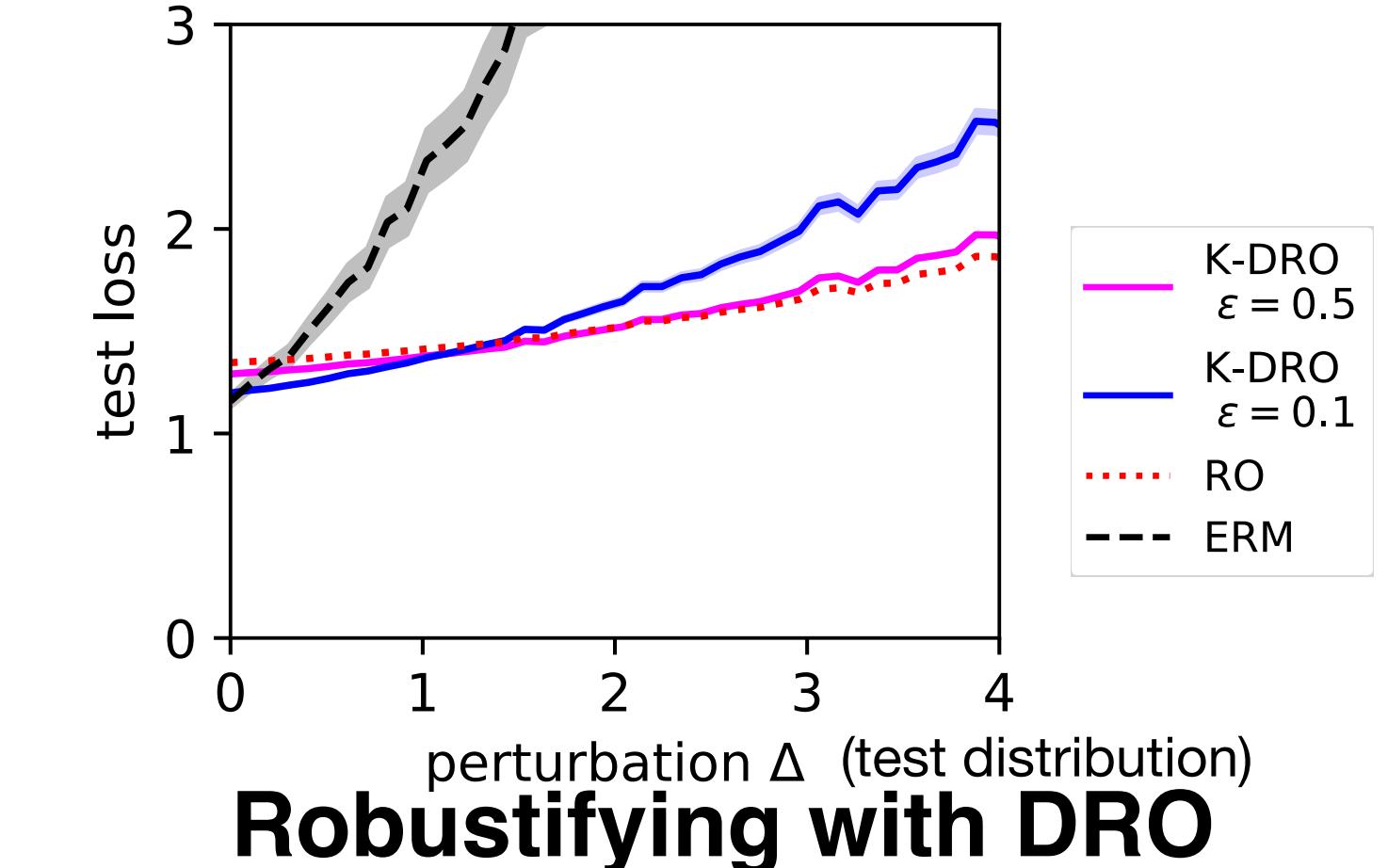
Intuition: flatten the curve, smooth is robust

## Example. Robust least squares

[El Ghaoui Lebret '97]

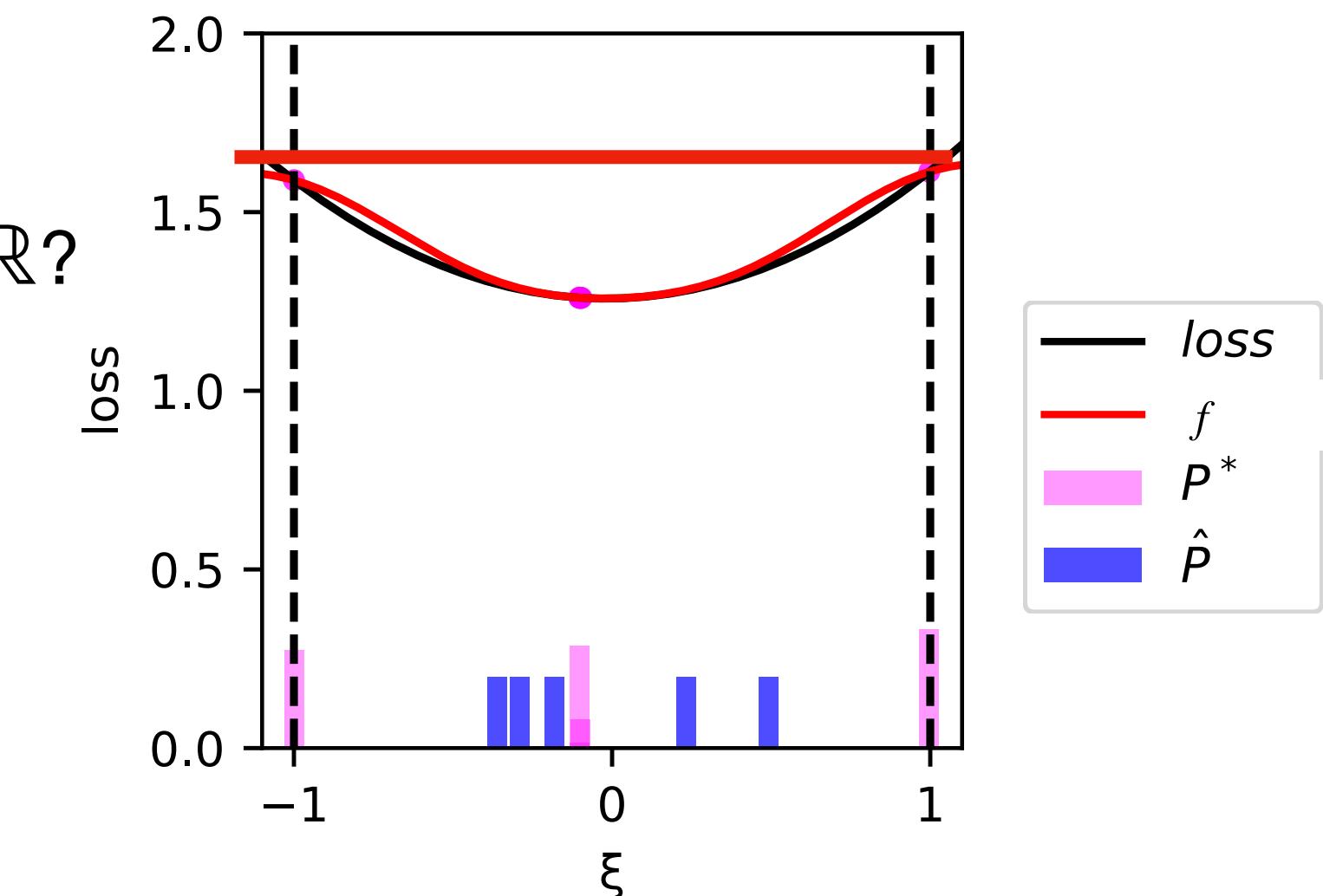
$$\text{minimize } l(\theta, \xi) := \|A(\xi) \cdot \theta - b\|_2^2$$

Given historical samples  $\xi_1, \xi_2, \dots, \xi_N$



## Robustifying with DRO

What if  $f \equiv c \in \mathbb{R}$ ?



# Special cases of DRO dual program

$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi) \iff (D) \min_{\theta, \mathbf{f} \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} \mathbf{f} + \epsilon \|\mathbf{f}\|_{\mathcal{H}} \text{ s.t. } l(\theta, \cdot) \leq \mathbf{f}$$

# Special cases of DRO dual program

$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi) \iff (D) \min_{\theta, \mathbf{f} \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} \mathbf{f} + \epsilon \|\mathbf{f}\|_{\mathcal{H}} \text{ s.t. } l(\theta, \cdot) \leq \mathbf{f}$$

- Type-1 Wasserstein DRO. Suppose the loss function  $l(\theta, \cdot)$  is Lipschitz continuous, and we choose  $\gamma$  to be the W-1 distance. (D) is then equivalent to Lipschitz regularized ERM

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(\theta, \xi_i) + \epsilon \text{Lip}(l).$$

# Special cases of DRO dual program

$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi) \iff (D) \min_{\theta, f \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} f + \epsilon \|f\|_{\mathcal{H}} \text{ s.t. } l(\theta, \cdot) \leq f$$

- Type-1 Wasserstein DRO. Suppose the loss function  $l(\theta, \cdot)$  is Lipschitz continuous, and we choose  $\gamma$  to be the W-1 distance. (D) is then equivalent to Lipschitz regularized ERM

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(\theta, \xi_i) + \epsilon \text{Lip}(l).$$

- Distributionally robust logistic regression (with exact regularization coefficient)

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(\theta^\top x_i, y_i) + \epsilon \cdot \|\theta\|_2.$$

# Special cases of DRO dual program

$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi) \iff (D) \min_{\theta, \mathbf{f} \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} \mathbf{f} + \epsilon \|\mathbf{f}\|_{\mathcal{H}} \text{ s.t. } l(\theta, \cdot) \leq \mathbf{f}$$

- Type-1 Wasserstein DRO. Suppose the loss function  $l(\theta, \cdot)$  is Lipschitz continuous, and we choose  $\gamma$  to be the W-1 distance. (D) is then equivalent to Lipschitz regularized ERM

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(\theta, \xi_i) + \epsilon \text{Lip}(l).$$

- Distributionally robust logistic regression (with exact regularization coefficient)

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(\theta^\top x_i, y_i) + \epsilon \cdot \|\theta\|_2.$$

- Empirical gradient penalty for DNN

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(f_\theta(x_i), y_i) + \lambda \cdot \max_i \|\nabla f_\theta(x_i)\|_2.$$

# Special cases of DRO dual program

$$(P) \min_{\theta} \sup_{\gamma(\mathbf{P}, \hat{\mathbf{P}}) \leq \epsilon} \mathbb{E}_{\mathbf{P}} l(\theta, \xi) \iff (D) \min_{\theta, \mathbf{f} \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} \mathbf{f} + \epsilon \|\mathbf{f}\|_{\mathcal{H}} \text{ s.t. } l(\theta, \cdot) \leq \mathbf{f}$$

- Type-1 Wasserstein DRO. Suppose the loss function  $l(\theta, \cdot)$  is Lipschitz continuous, and we choose  $\gamma$  to be the W-1 distance. (D) is then equivalent to Lipschitz regularized ERM

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(\theta, \xi_i) + \epsilon \text{Lip}(l).$$

- Distributionally robust logistic regression (with exact regularization coefficient)

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(\theta^\top x_i, y_i) + \epsilon \cdot \|\theta\|_2.$$

- Empirical gradient penalty for DNN

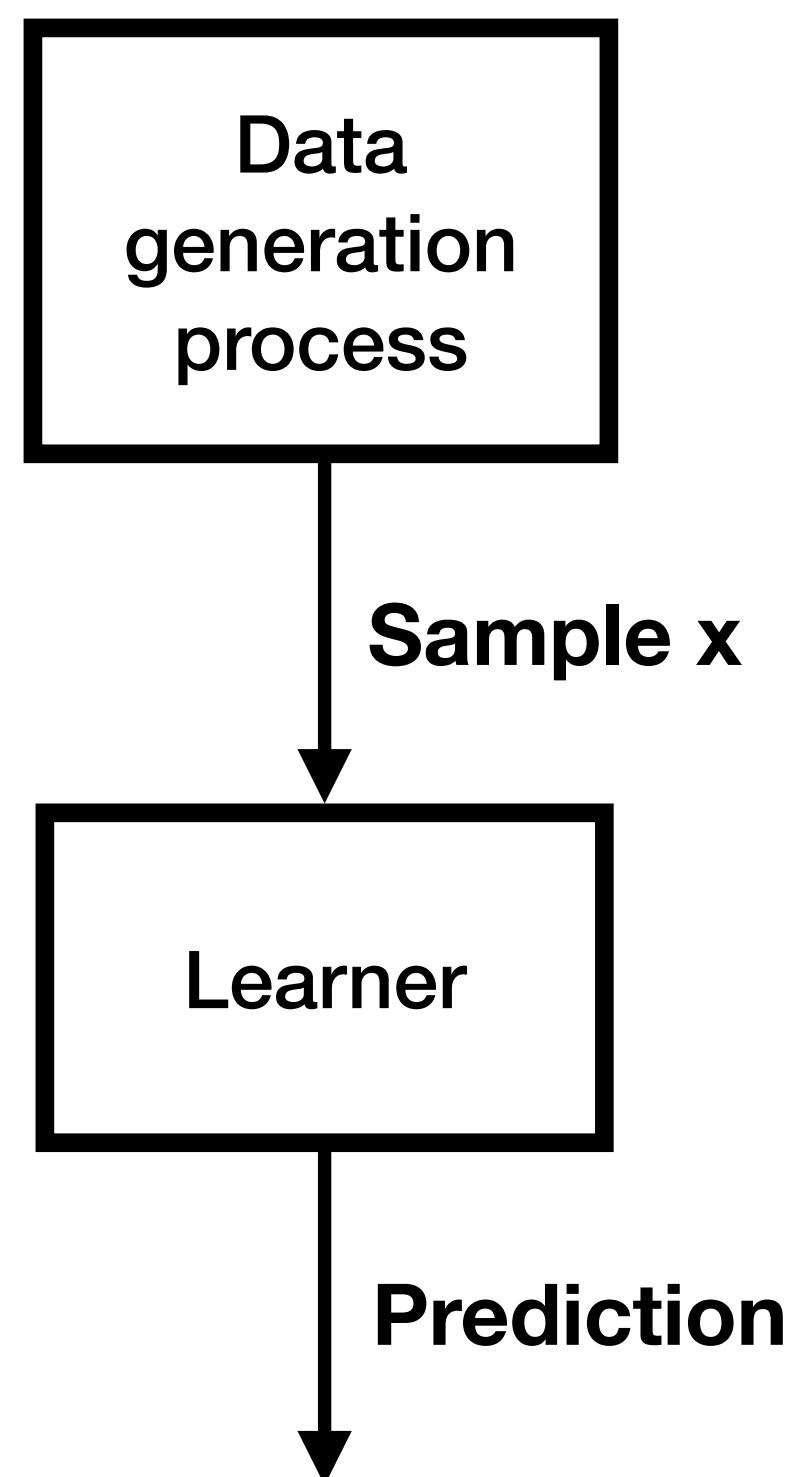
$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l(f_\theta(x_i), y_i) + \lambda \cdot \max_i \|\nabla f_\theta(x_i)\|_2.$$

- Kernel DRO. If we choose to be the MMD, then solving DRO becomes a kernel-based learning

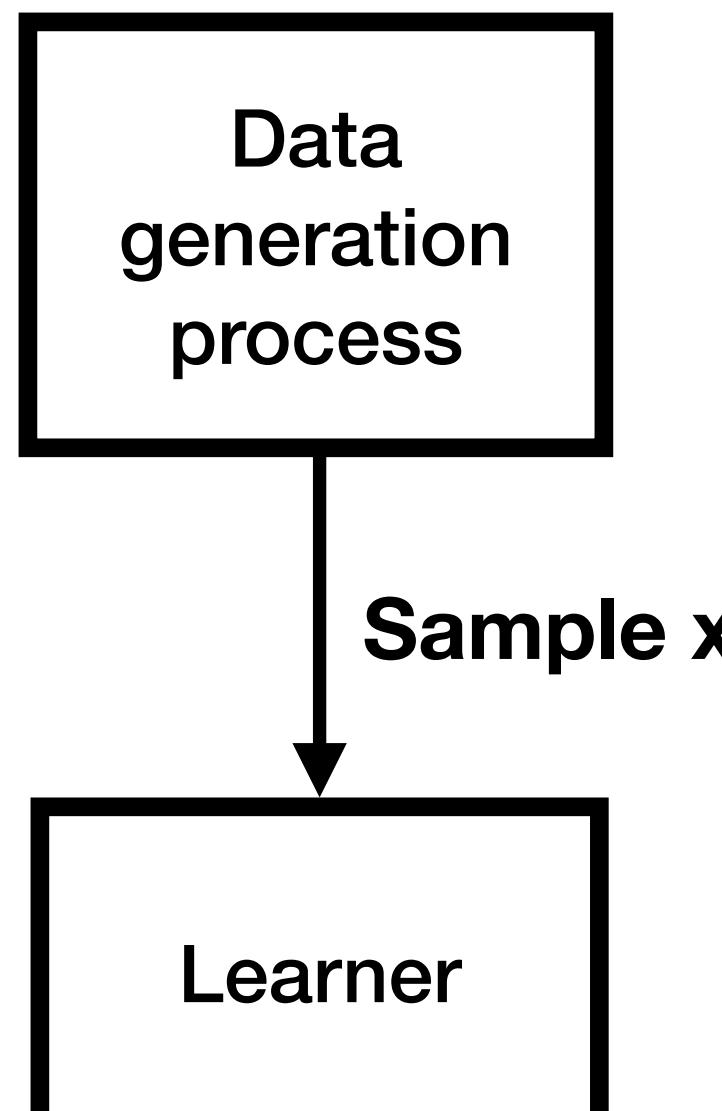
$$(D) \min_{\theta, \mathbf{f} \in \mathcal{H}} \mathbb{E}_{\hat{\mathbf{P}}} \mathbf{f} + \epsilon \|\mathbf{f}\|_{\mathcal{H}} \text{ s.t. } l(\theta, \cdot) \leq \mathbf{f}.$$

# **Multi-stage robust decision-making**

# SL

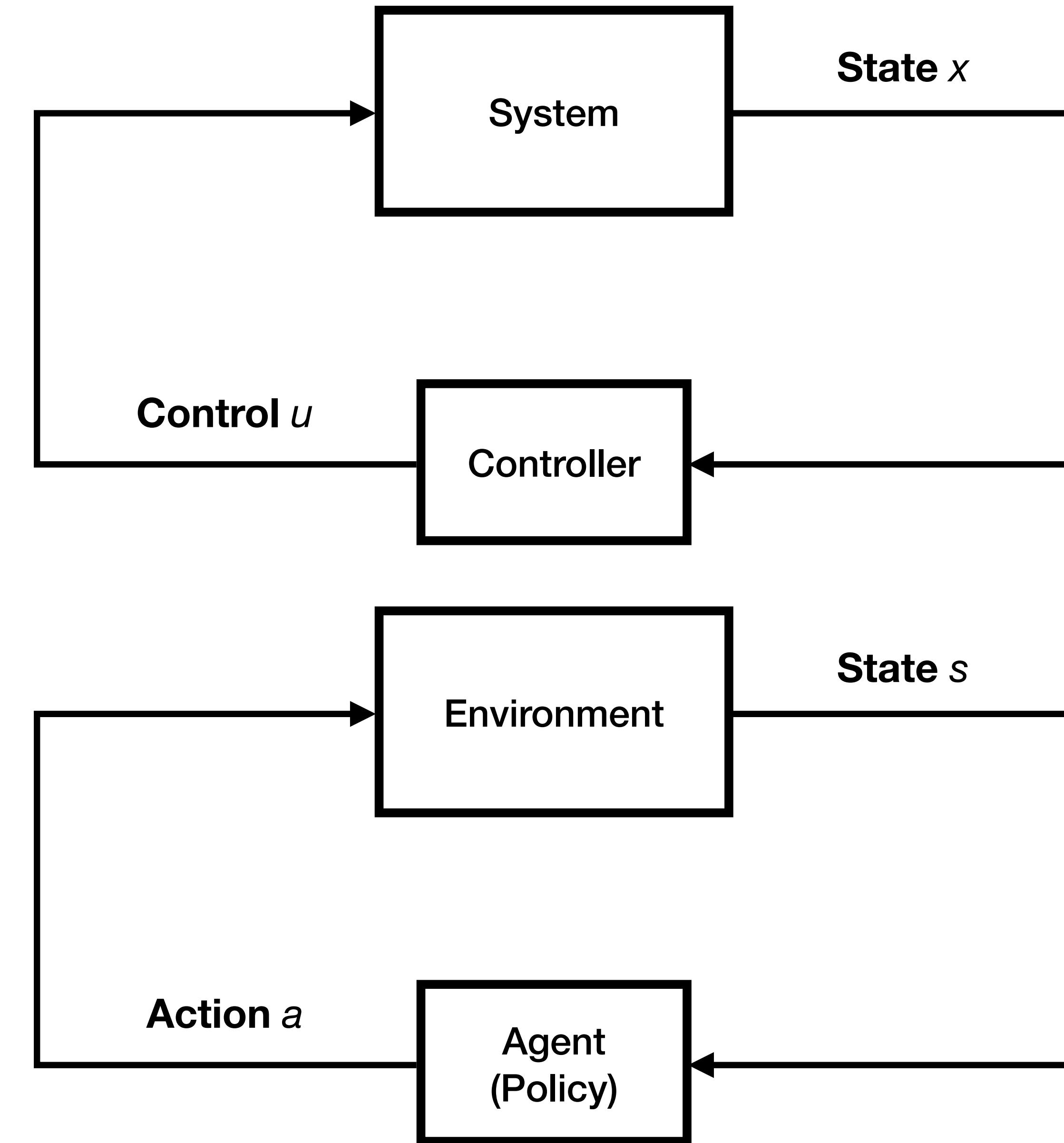


# SL

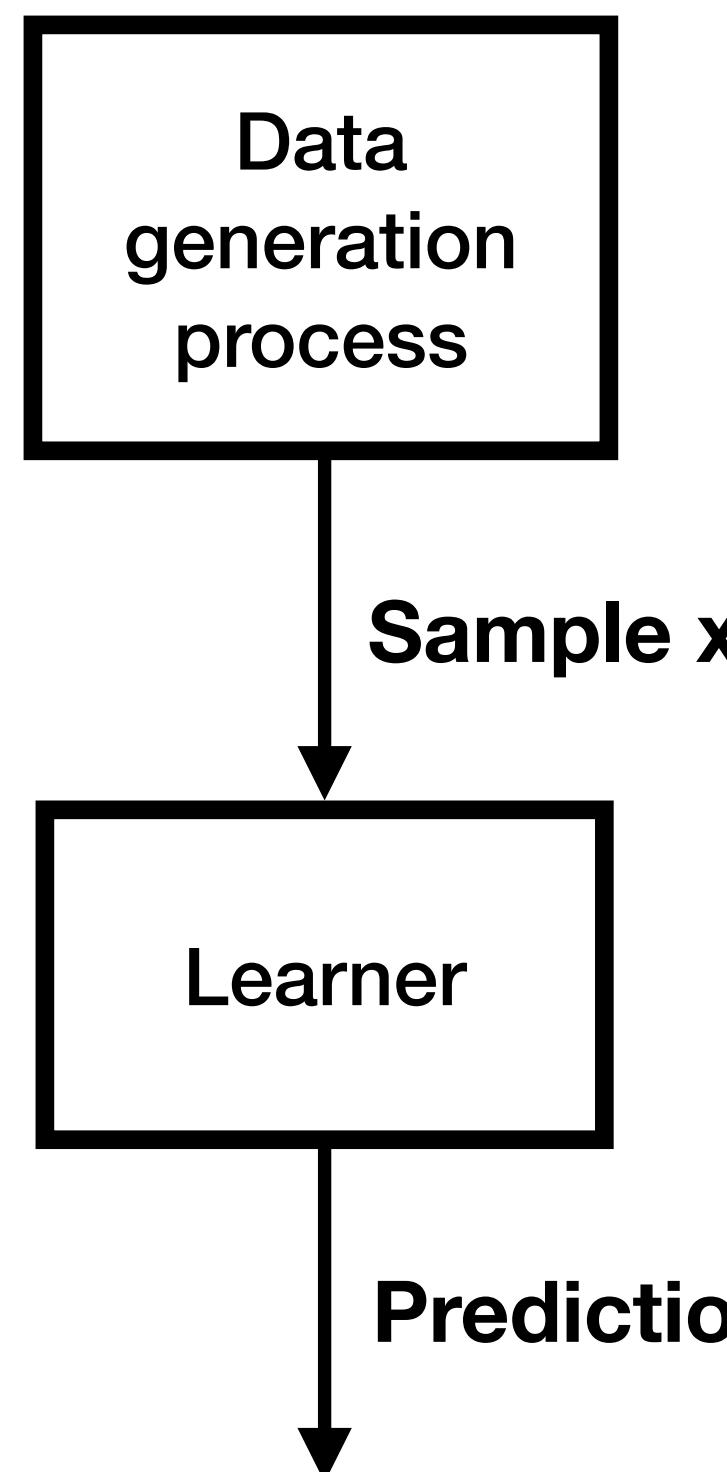


# Control

# RL

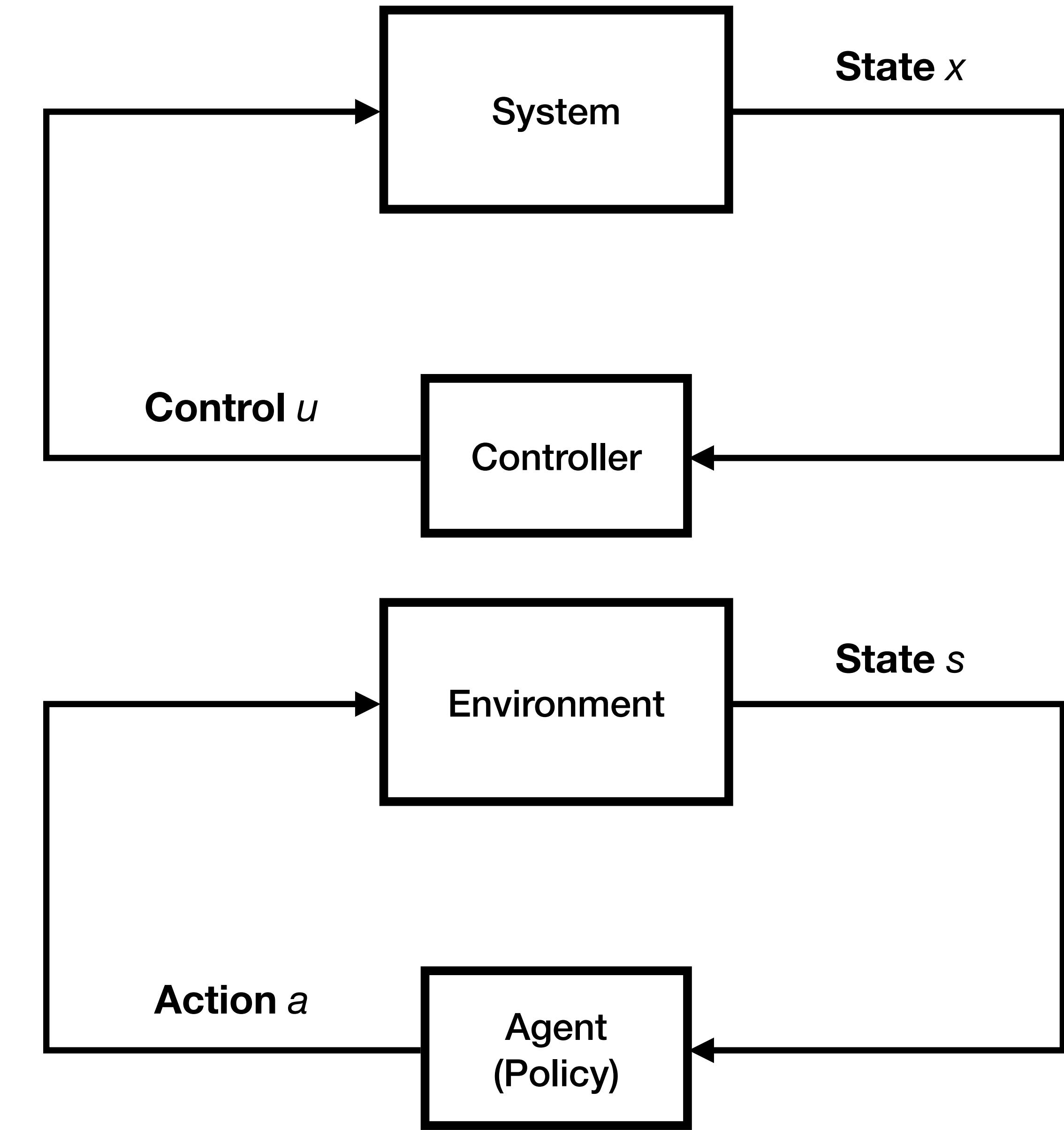


# SL



# Control

# RL



Distinction: feedback in systems causes **distribution shift** in the data-generation processes

# Reinforcement learning problem

$$\max_{\mu} \quad \mathbb{E} \sum_{t=1}^T \mathbf{reward}(x_t, \mu(x_t))$$

# Reinforcement learning problem

$$\max_{\mu} \quad \mathbb{E} \sum_{t=1}^T \mathbf{reward}(x_t, \mu(x_t))$$



$$\min_{u_{1:T}} \quad \sum_{t=1}^T \mathbf{cost}(x_t, u_t)$$

# Optimal control problem

# Reinforcement learning problem

$$\max_{\mu} \quad \mathbb{E} \sum_{t=1}^T \mathbf{reward}(x_t, \mu(x_t))$$

↓

$$\min_{u_{1:T}} \quad \sum_{t=1}^T \mathbf{cost}(x_t, u_t)$$



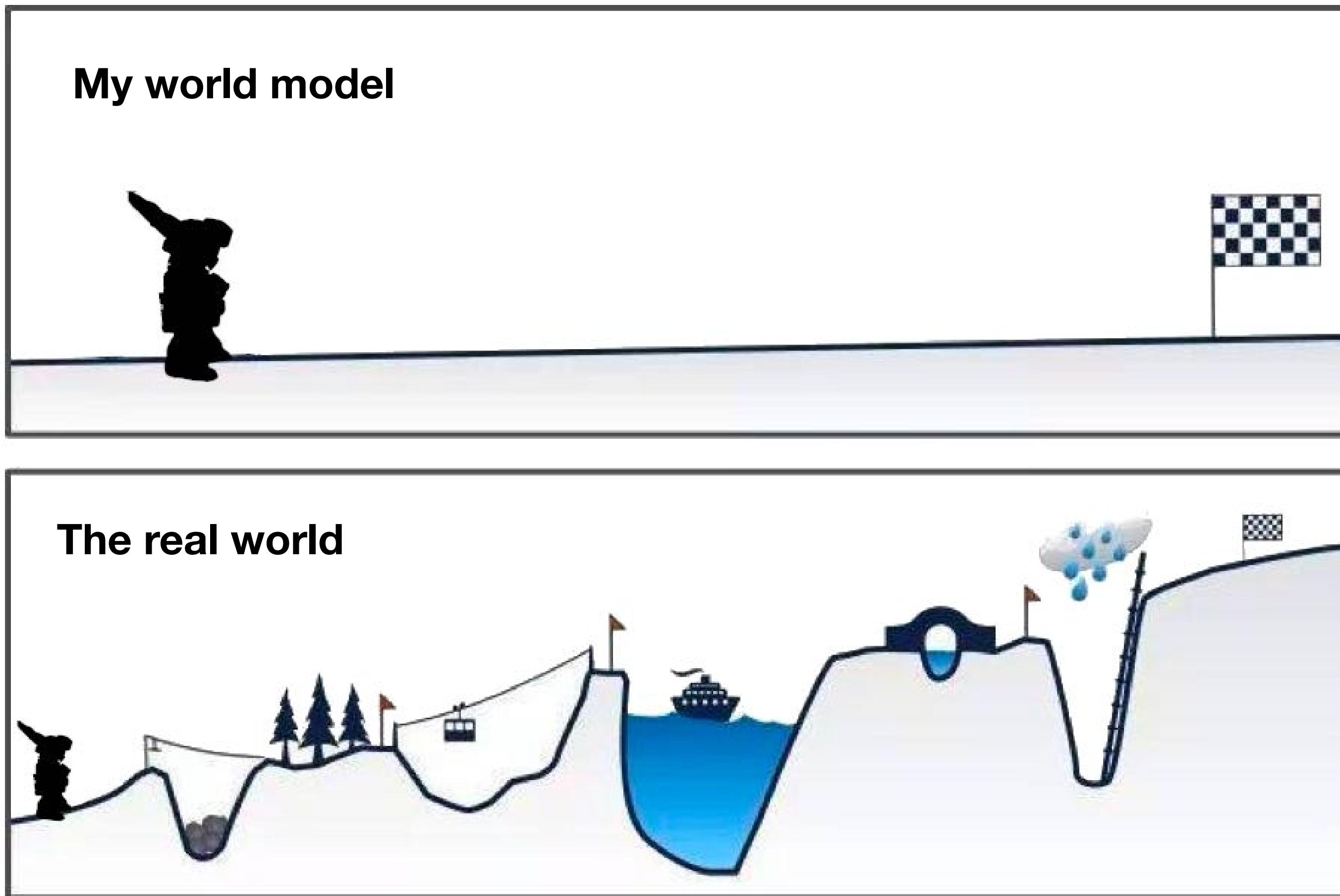
# Optimal control problem

# Why do we need to be robust in RL/control?

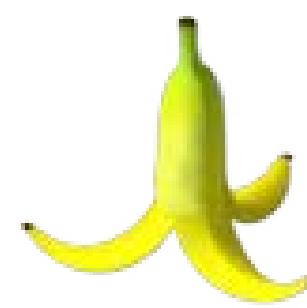
# Why do we need to be robust in RL/control?



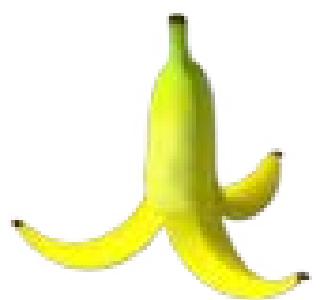
# Why do we need to be robust in RL/control?



# Model Predictive Control



# Model Predictive Control



## Optimal control problem

(discrete-time)

$$\underset{u_0, \dots, u_{N-1}}{\text{minimize}} \quad \sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

subject to

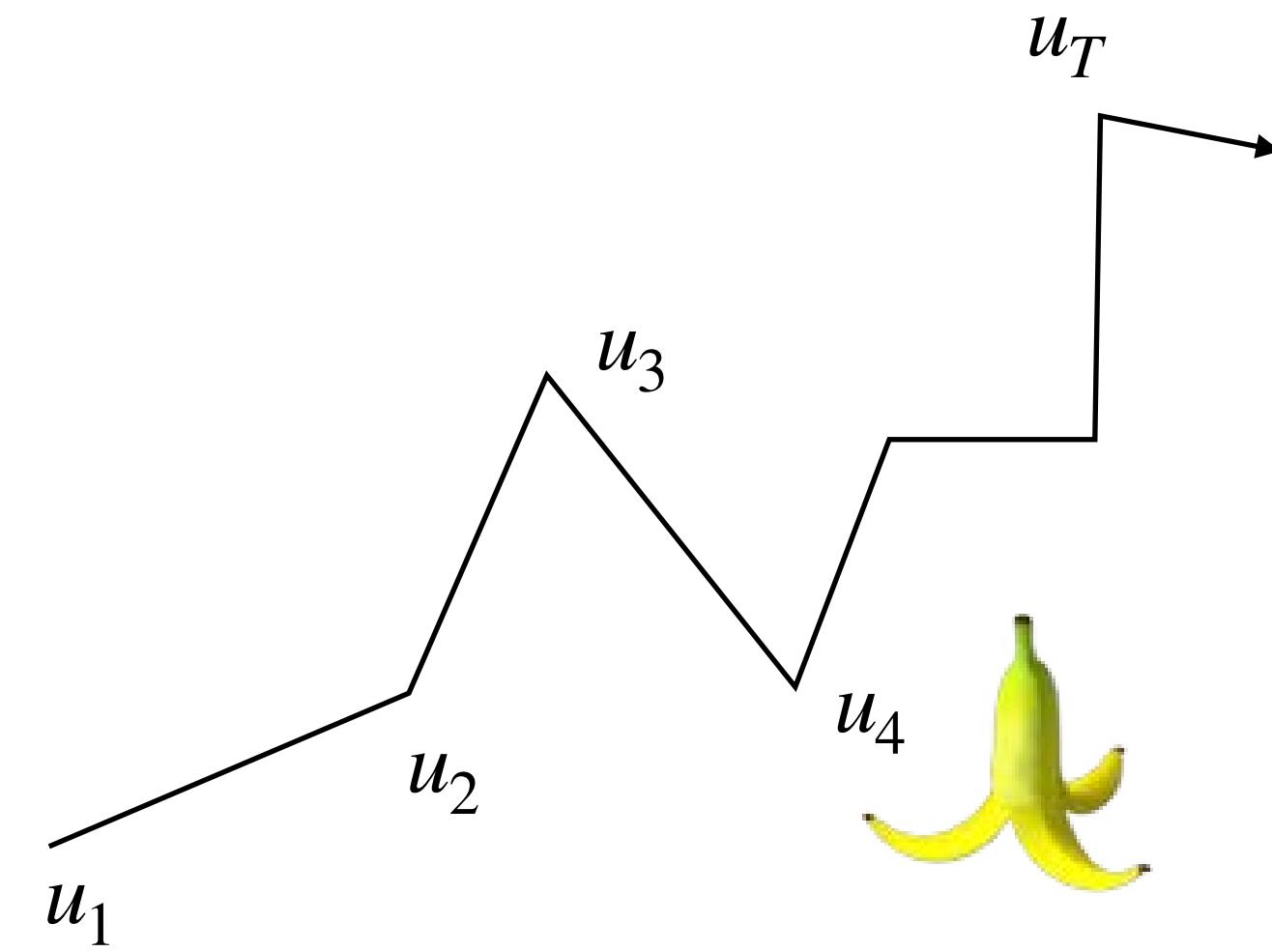
$$x_{t+1} = f(x_t, u_t)$$

$$c(x_t, u_t) \leq 0$$

$$x_0 = s$$

$$t = 0, \dots, N - 1$$

# Model Predictive Control



## Optimal control problem

(discrete-time)

$$\sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

minimize  
 $u_0, \dots, u_{N-1}$

subject to

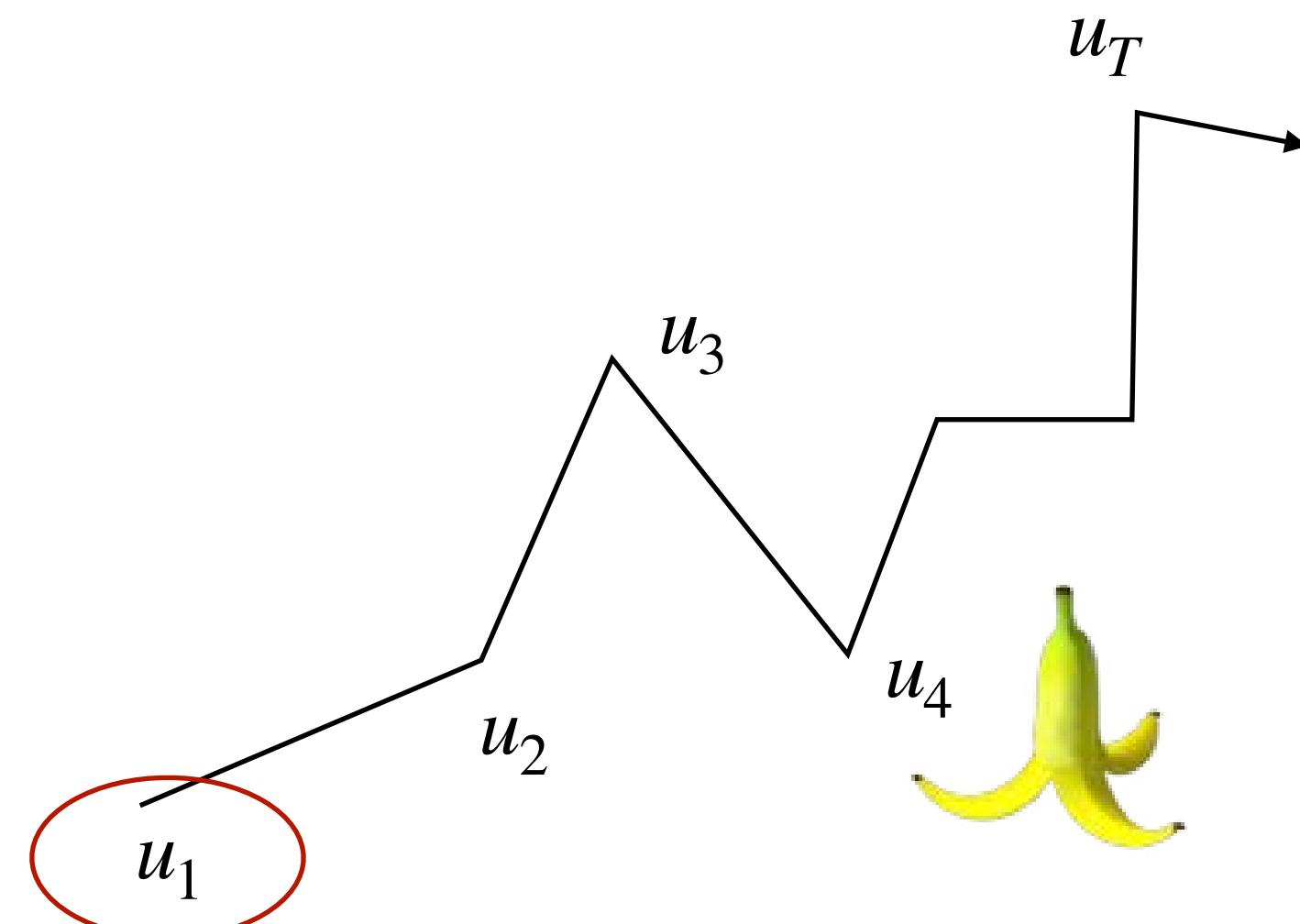
$$x_{t+1} = f(x_t, u_t)$$

$$c(x_t, u_t) \leq 0$$

$$x_0 = s$$

$$t = 0, \dots, N - 1$$

# Model Predictive Control



## Optimal control problem

(discrete-time)

$$\sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

minimize  
 $u_0, \dots, u_{N-1}$

subject to

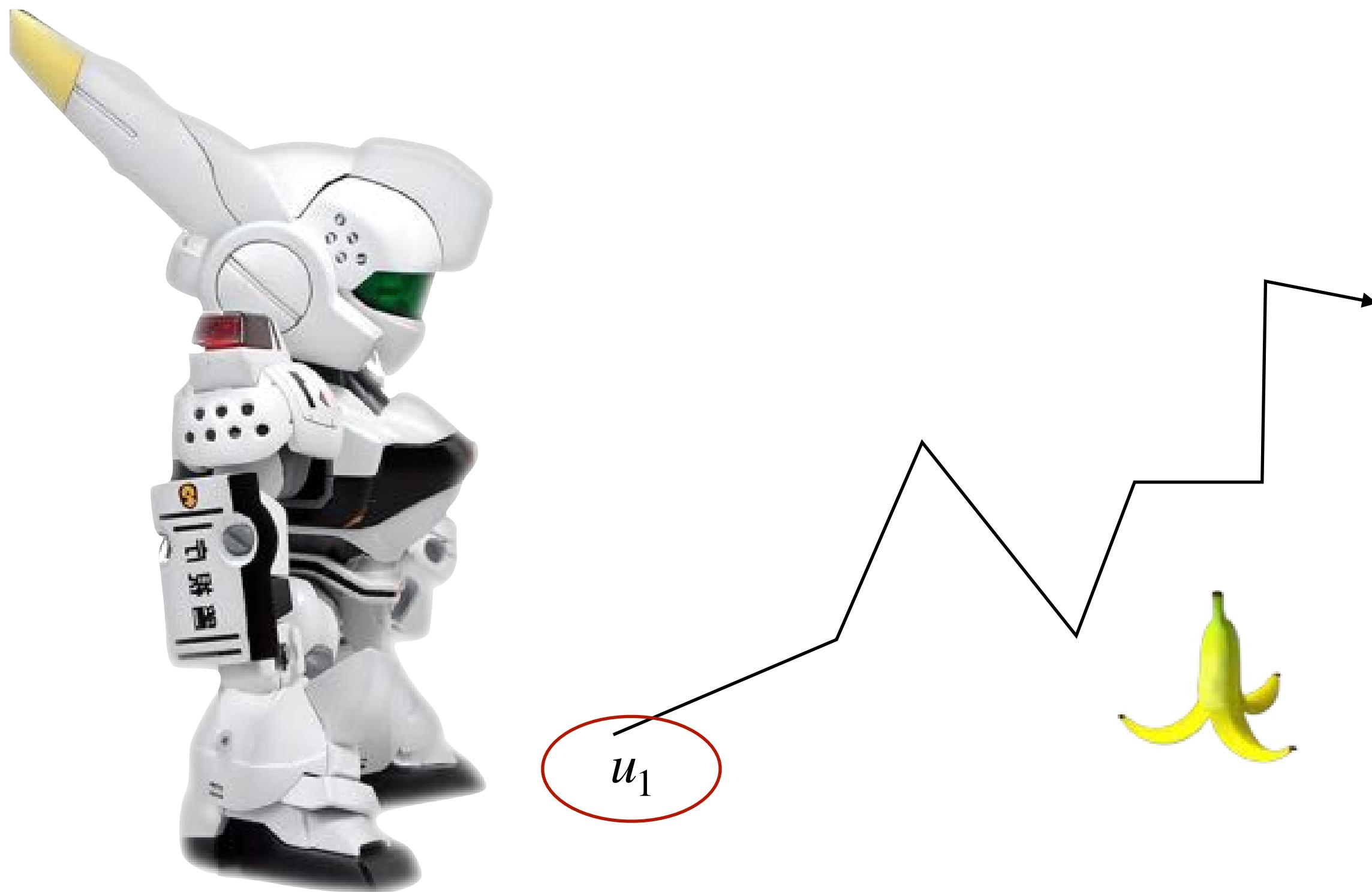
$$x_{t+1} = f(x_t, u_t)$$

$$c(x_t, u_t) \leq 0$$

$$x_0 = s$$

$$t = 0, \dots, N - 1$$

# Model Predictive Control



## Optimal control problem

(discrete-time)

$$\sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

minimize  
 $u_0, \dots, u_{N-1}$

subject to

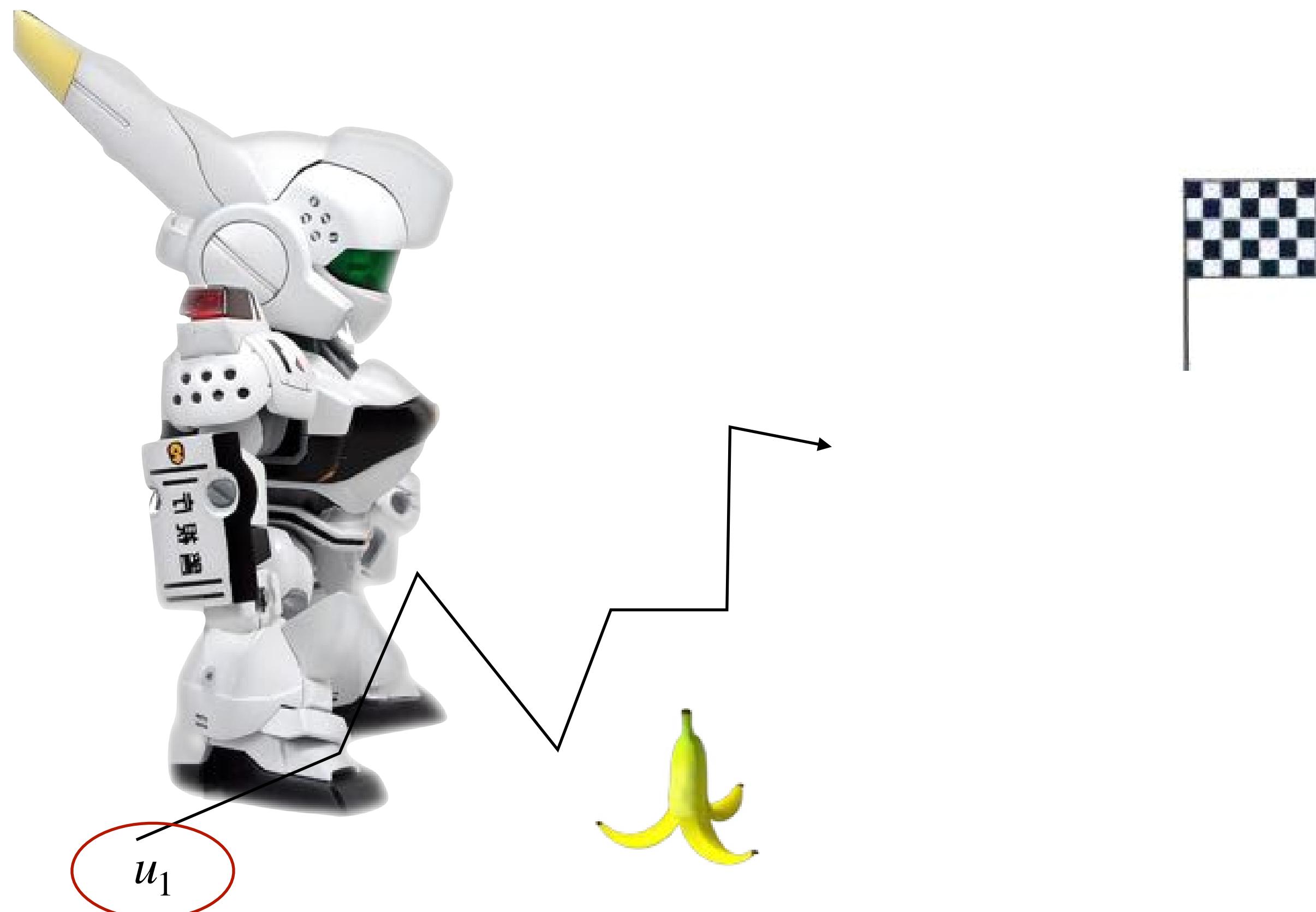
$$x_{t+1} = f(x_t, u_t)$$

$$c(x_t, u_t) \leq 0$$

$$x_0 = s$$

$$t = 0, \dots, N - 1$$

# Model Predictive Control



## Optimal control problem

(discrete-time)

$$\underset{u_0, \dots, u_{N-1}}{\text{minimize}} \quad \sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

subject to

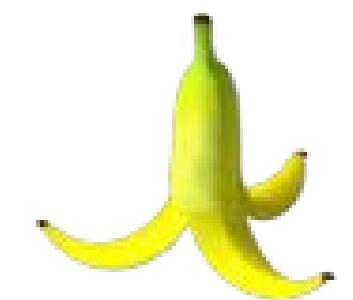
$$x_{t+1} = f(x_t, u_t)$$

$$c(x_t, u_t) \leq 0$$

$$x_0 = s$$

$$t = 0, \dots, N - 1$$

# Model Predictive Control



## Optimal control problem

(discrete-time)

minimize  
 $u_0, \dots, u_{N-1}$

$$\sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

subject to

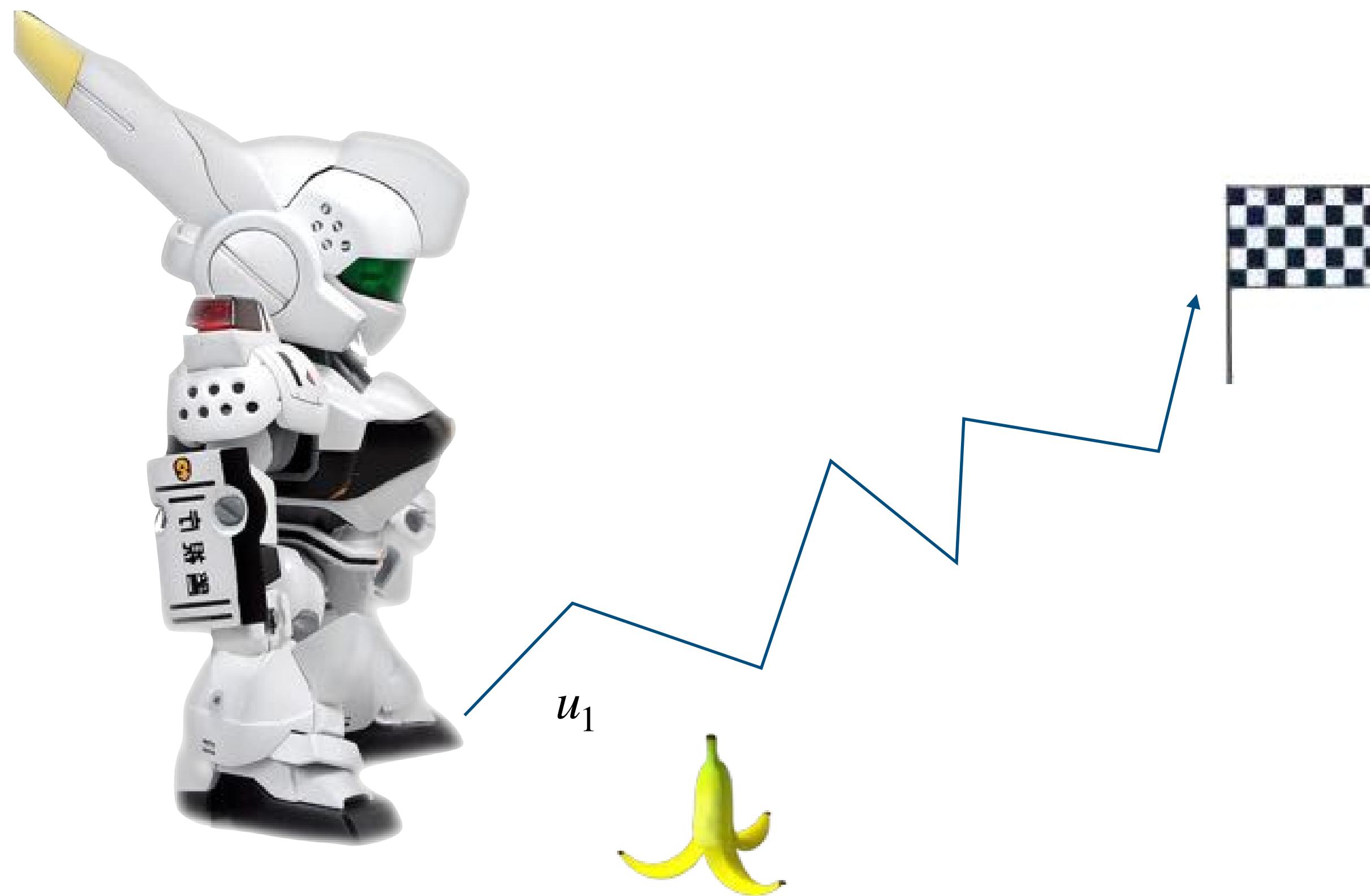
$$x_{t+1} = f(x_t, u_t)$$

$$c(x_t, u_t) \leq 0$$

$$x_0 = s$$

$$t = 0, \dots, N - 1$$

# Model Predictive Control



## Optimal control problem

(discrete-time)

$$\begin{array}{l} \text{minimize} \\ u_0, \dots, u_{N-1} \end{array}$$

subject to

$$\sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

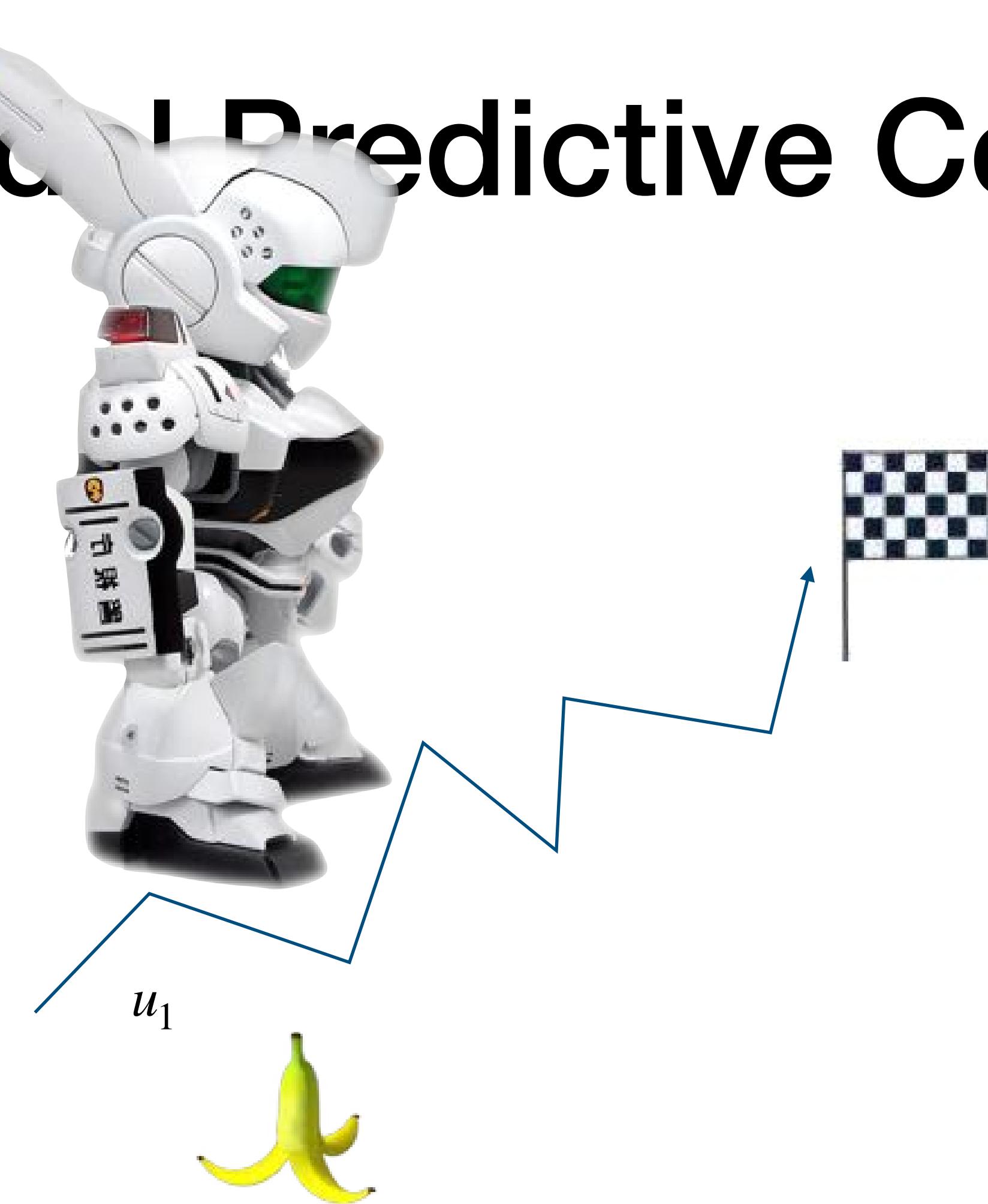
$$x_{t+1} = f(x_t, u_t)$$

$$c(x_t, u_t) \leq 0$$

$$x_0 = s$$

$$t = 0, \dots, N - 1$$

# Model Predictive Control



## Optimal control problem

(discrete-time)

$$\sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

minimize  
 $u_0, \dots, u_{N-1}$

subject to

$$x_{t+1} = f(x_t, u_t)$$

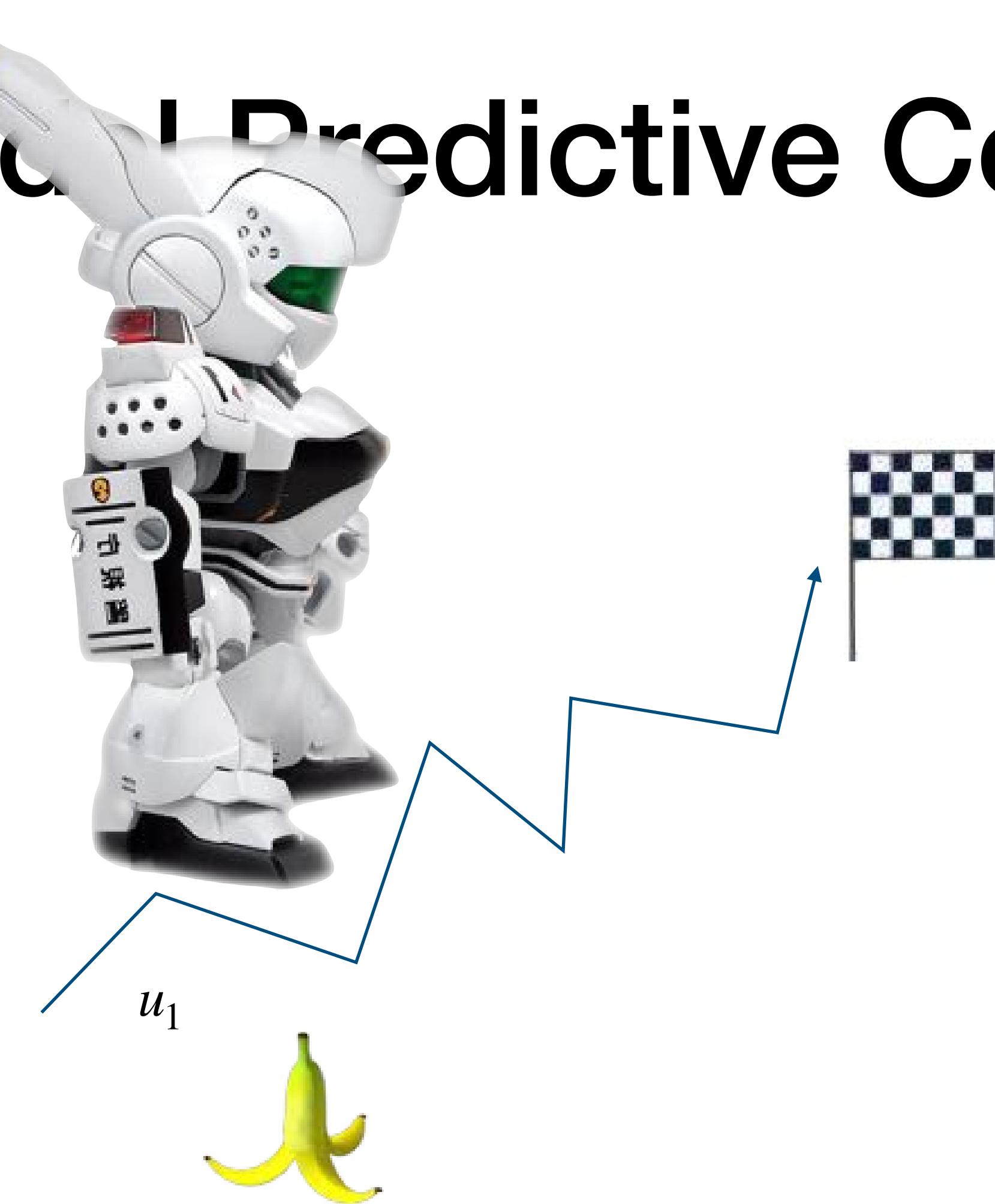
$$c(x_t, u_t) \leq 0$$

$$x_0 = s$$

$$t = 0, \dots, N - 1$$

You are using MPC when you drive or catch a train.

# Model Predictive Control



## Optimal control problem

(discrete-time)

$$\begin{aligned} & \text{minimize}_{u_0, \dots, u_{N-1}} \\ & \quad \sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N) \\ & \text{subject to} \\ & \quad x_{t+1} = f(x_t, u_t) \\ & \quad c(x_t, u_t) \leq 0 \\ & \quad x_0 = s \\ & \quad t = 0, \dots, N-1 \end{aligned}$$

You are using MPC when you drive or catch a train.  
Learning-based MPC is an active research area!

# **Robust and stochastic model predictive control**

# Robust and stochastic model predictive control

## Optimal control problem

(discrete-time)

minimize  
 $u_0, \dots, u_{N-1}$

$$\sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

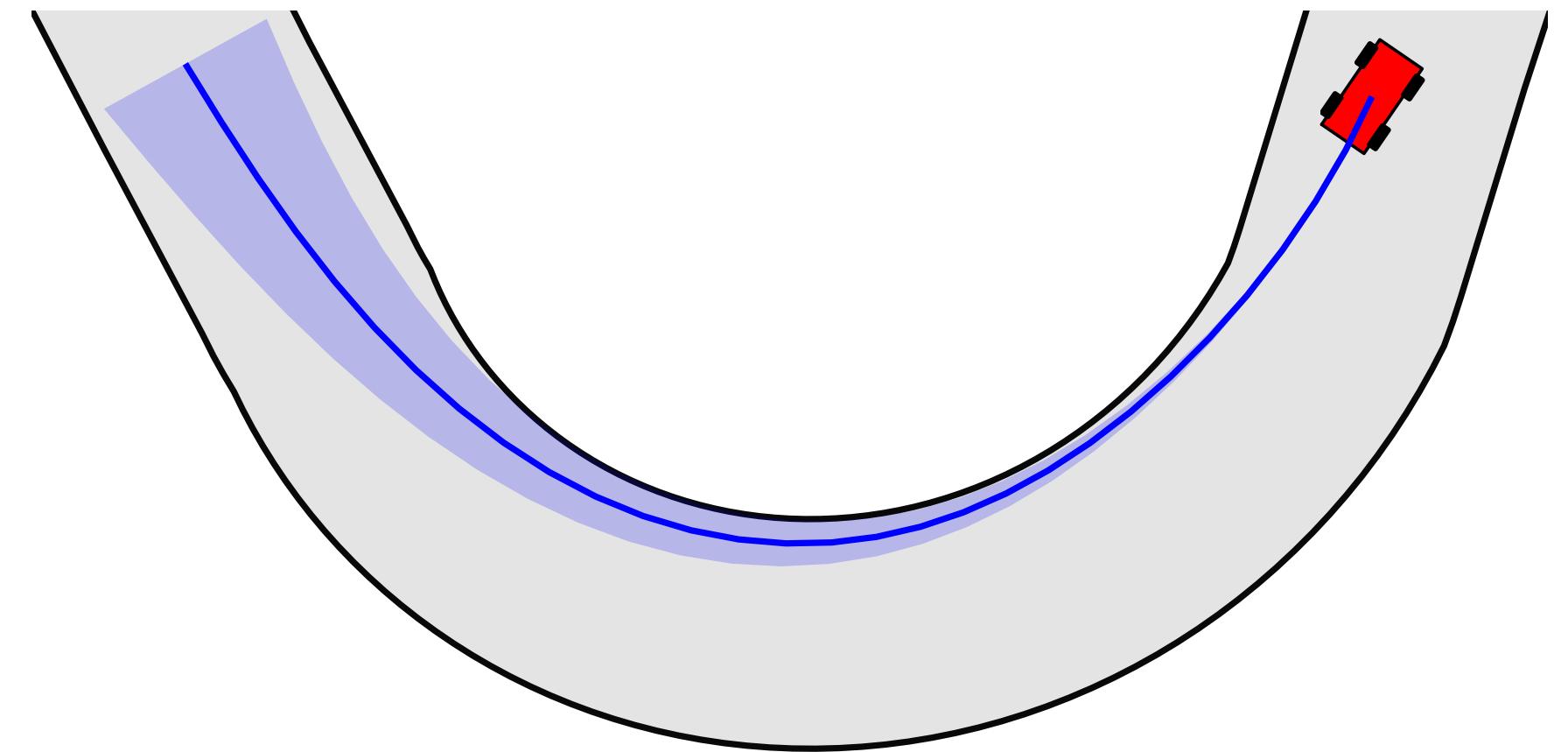
subject to

$$x_{t+1} = f(x_t, u_t)$$

$$c(x_t, u_t) \leq 0$$

$$x_0 = s$$

$$t = 0, \dots, N - 1$$



# Robust and stochastic model predictive control

## Optimal control problem

(discrete-time)

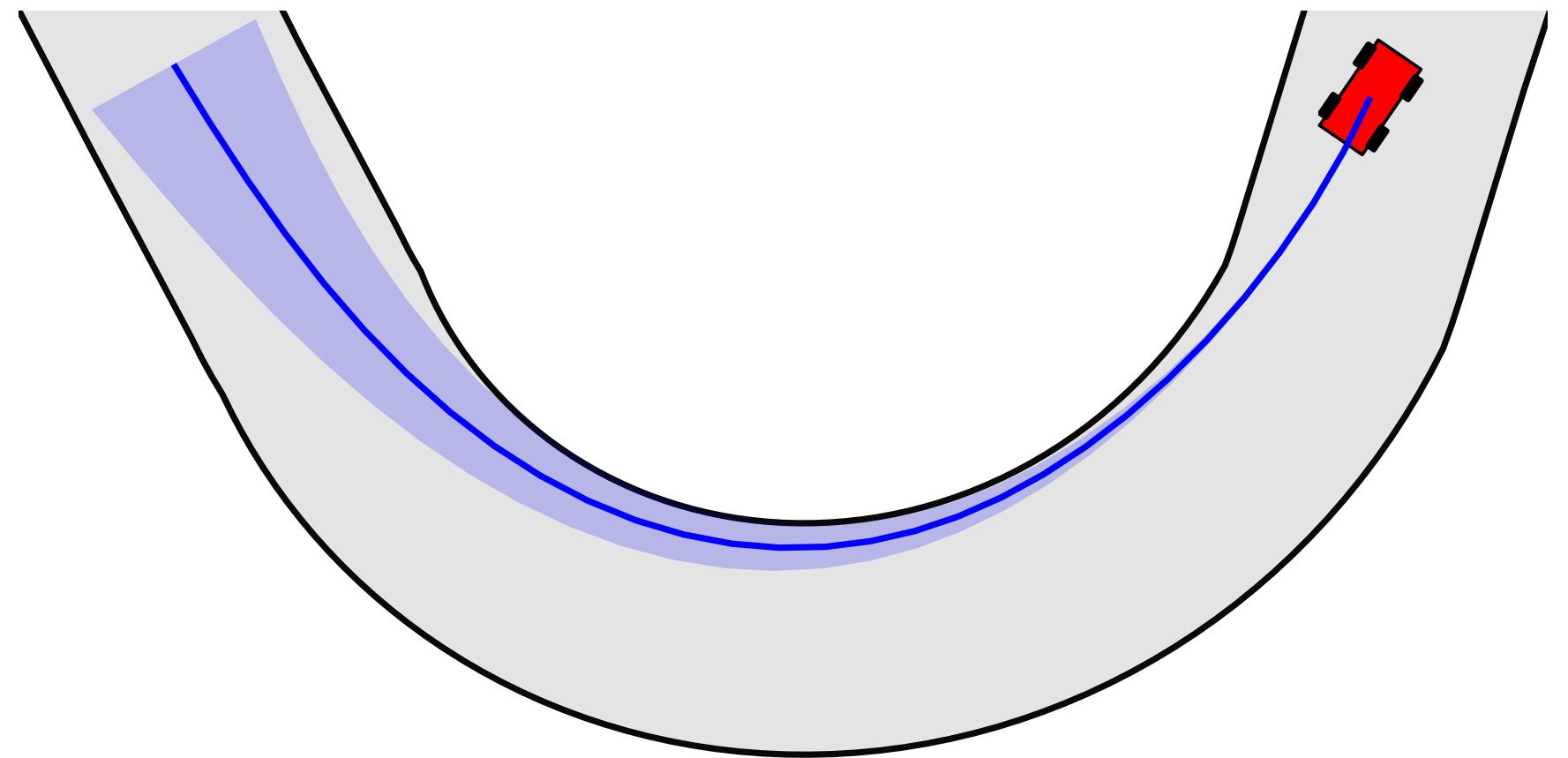
$$\text{minimize}_{u_0, \dots, u_{N-1}} \sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

$$\text{subject to } x_{t+1} = f(x_t, u_t)$$

$$c(x_t, u_t) \leq 0$$

$$x_0 = s$$

$$t = 0, \dots, N - 1$$



## Guaranteed Margins for LQG Regulators

JOHN C. DOYLE

*Abstract*—There are none.

# Robust and stochastic model predictive control

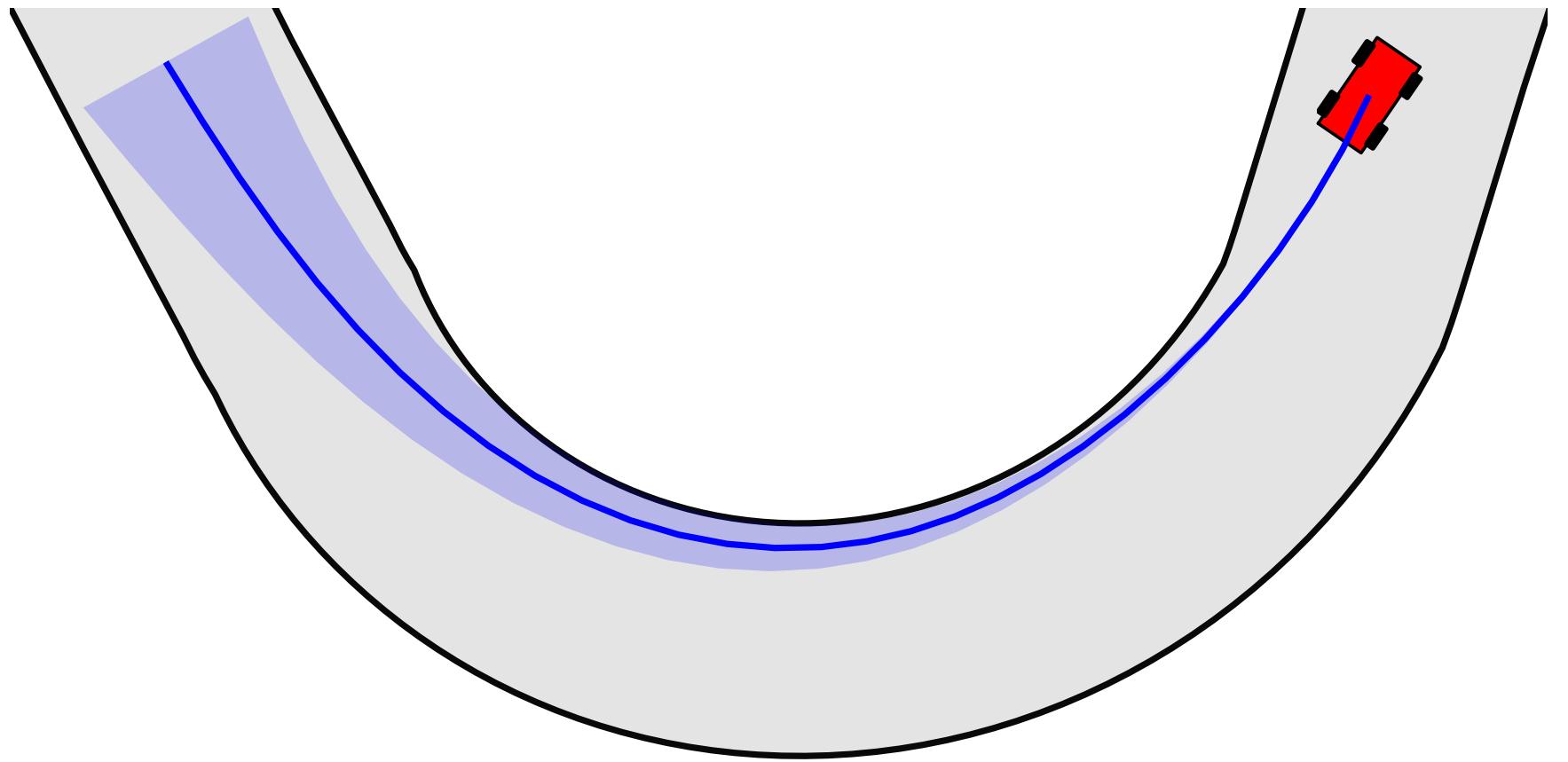
## Optimal control problem

(discrete-time)

minimize  
 $u_0, \dots, u_{N-1}$

$$\sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

subject to  
 $x_{t+1} = f(x_t, u_t)$   
 $c(x_t, u_t) \leq 0$   
 $x_0 = s$   
 $t = 0, \dots, N - 1$



## Guaranteed Margins for LQG Regulators

JOHN C. DOYLE

**Abstract**—There are none.

Optimization (alone) is not robust.



# Robust and stochastic model predictive control

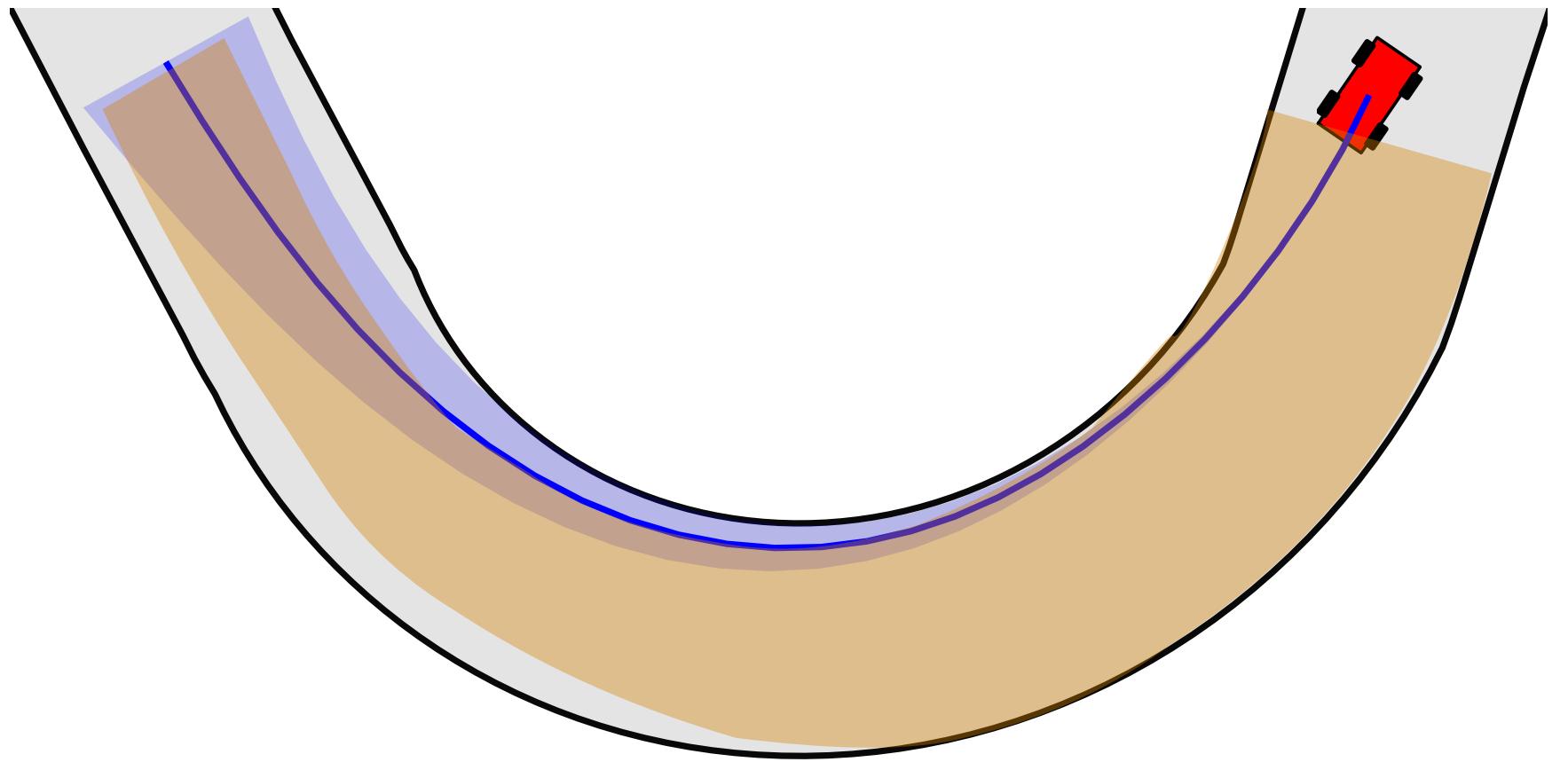
## Optimal control problem

(discrete-time)

minimize  
 $u_0, \dots, u_{N-1}$

$$\sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

subject to  
 $x_{t+1} = f(x_t, u_t)$   
 $c(x_t, u_t) \leq 0$   
 $x_0 = s$   
 $t = 0, \dots, N - 1$



## Guaranteed Margins for LQG Regulators

JOHN C. DOYLE

**Abstract**—There are none.

Optimization (alone) is not robust.



- **Constraint Tightening strategy:** Drive within a narrower “tube” than the original track!

# Robust and stochastic model predictive control

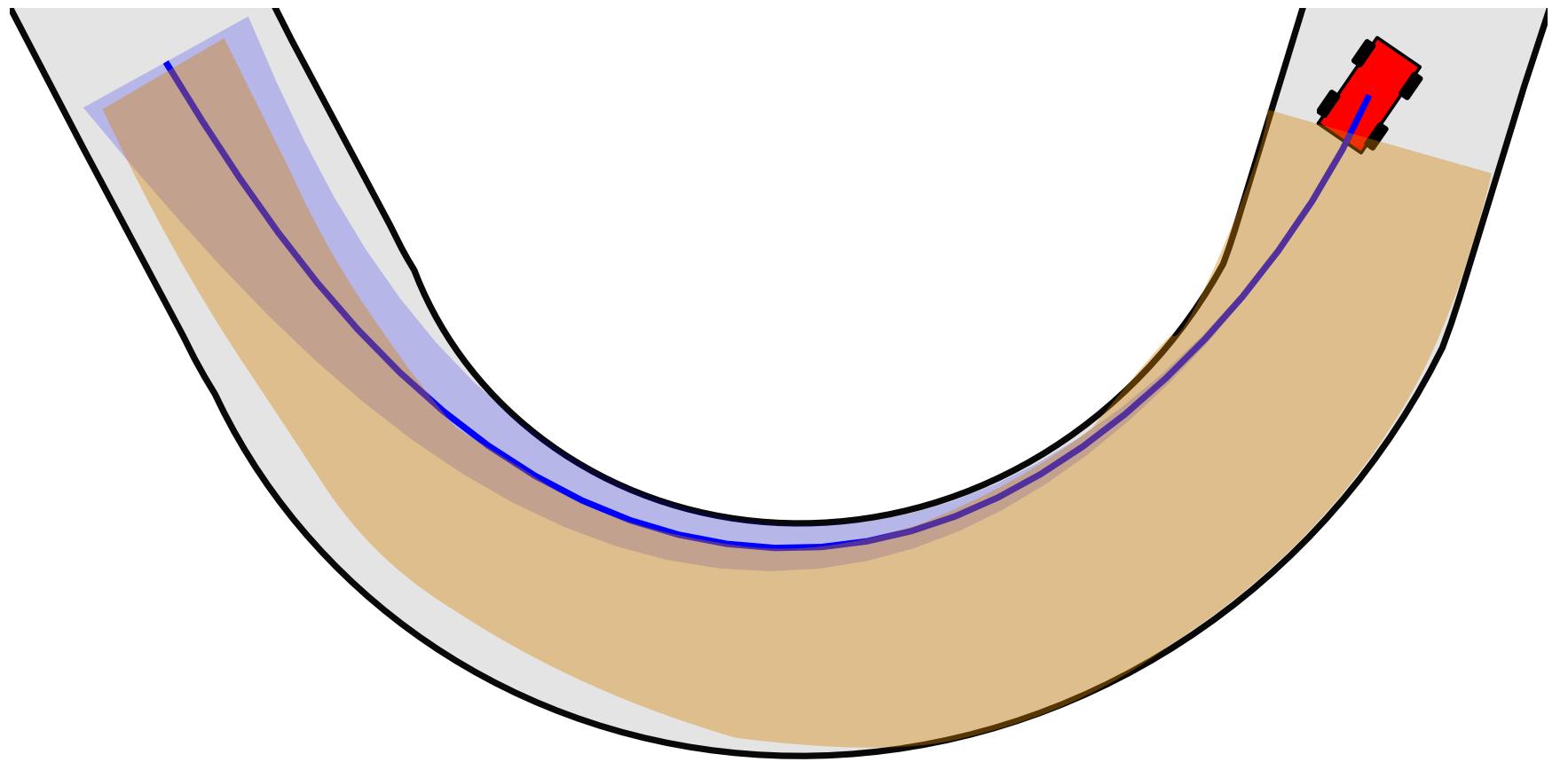
## Optimal control problem

(discrete-time)

minimize  
 $u_0, \dots, u_{N-1}$

$$\sum_{t=1}^{N-1} l_t(x_t, u_t) + l_N(x_N)$$

subject to  
 $x_{t+1} = f(x_t, u_t)$   
 $c(x_t, u_t) \leq 0$   
 $x_0 = s$   
 $t = 0, \dots, N - 1$



## Guaranteed Margins for LQG Regulators

JOHN C. DOYLE

**Abstract**—There are none.

Optimization (alone) is not robust.



- **Constraint Tightening strategy:** Drive within a narrower “tube” than the original track!
- In practice, this often boils down to replacing the constraint  $c(x_t, u_t) \leq 0$  by  $c(x_t, u_t) + b \leq 0$  for some  $b > 0$ .

# Learning for control: Three main classes of methods

# Learning for control: Three main classes of methods

- Learning-based control, a.k.a., model-based reinforcement learning
  - Approximate the system model, then solve OCP

# Learning for control: Three main classes of methods

- Learning-based control, a.k.a., model-based reinforcement learning
  - Approximate the system model, then solve OCP
- Imitation learning
  - Directly learning a control policy from demonstration

# Learning for control: Three main classes of methods

- Learning-based control, a.k.a., model-based reinforcement learning
  - Approximate the system model, then solve OCP
- Imitation learning
  - Directly learning a control policy from demonstration
- Model-free reinforcement learning
  - Policy/value iteration using Bellman, policy search, Q-learning, ...

# Learning for control: Three main classes of methods

- Learning-based control, a.k.a., model-based reinforcement learning
  - Approximate the system model, then solve OCP
- Imitation learning
  - Directly learning a control policy from demonstration
- Model-free reinforcement learning
  - Policy/value iteration using Bellman, policy search, Q-learning, ...

**They all suffer from distribution shift!**

# Covariate shift in imitation learning

# Covariate shift in imitation learning



Train

**Formula1Net**  $\pi_\theta$

# Covariate shift in imitation learning



Train

Formula1Net  $\pi_\theta$

$$(\text{IL}): \min_{\theta \in \Theta} \int \|\pi_\theta(x) - \pi_e(x)\|^2 \mu^{\pi_e}(dx),$$

# Covariate shift in imitation learning



Train

Formula1Net  $\pi_\theta$

$$(\text{IL}): \min_{\theta \in \Theta} \int \|\pi_\theta(x) - \pi_e(x)\|^2 \mu^{\pi_e}(dx),$$

What can go wrong?

# Covariate shift in imitation learning



Train

Formula1Net  $\pi_\theta$

$$(\text{IL}): \min_{\theta \in \Theta} \int \|\pi_\theta(x) - \pi_e(x)\|^2 \mu^{\pi_e}(dx),$$

What can go wrong?

- Shift in the induced state distributions  $\mu^\pi$ .

# Covariate shift in imitation learning



Train

Formula1Net  $\pi_\theta$

$$(\text{IL}): \min_{\theta \in \Theta} \int \|\pi_\theta(x) - \pi_e(x)\|^2 \mu^{\pi_e}(dx),$$

What can go wrong?

- Shift in the induced state distributions  $\mu^\pi$ .
- Since we care about the performance

$$\int l(x, \pi_\theta(x)) \mu^{\pi_\theta}(dx)$$

not  $\int l(x, \pi_\theta(x)) \mu^{\pi_e}(dx)$ .

But  $\pi_\theta, \pi_e$  are only similar under the  $\mu^{\pi_e}$ .

# Covariate shift in imitation learning



Train

Formula1Net  $\pi_\theta$

$$(\text{IL}): \min_{\theta \in \Theta} \int \|\pi_\theta(x) - \pi_e(x)\|^2 \mu^{\pi_e}(dx),$$

What can go wrong?

- Shift in the induced state distributions  $\mu^\pi$ .
  - Since we care about the performance  
$$\int l(x, \pi_\theta(x)) \mu^{\pi_\theta}(dx)$$
  
not 
$$\int l(x, \pi_\theta(x)) \mu^{\pi_e}(dx).$$
- But  $\pi_\theta, \pi_e$  are only similar under the  $\mu^{\pi_e}$ .
- Same goes for learning dynamics for control

# Fix: covariate shift in imitation learning

# Fix: covariate shift in imitation learning

DAgger (Data Aggregation)

1. Train a policy  $\pi_\theta$  from expert data  $D$
2. Run  $\pi_\theta$  (on the robot) to obtain new states
3. Ask the expert to label the states with actions and obtain new data  $D'$
4. Add the labeled data to the training set  $D = D \cup D'$  (aggregation)

# Fix: covariate shift in imitation learning

DAgger (Data Aggregation)

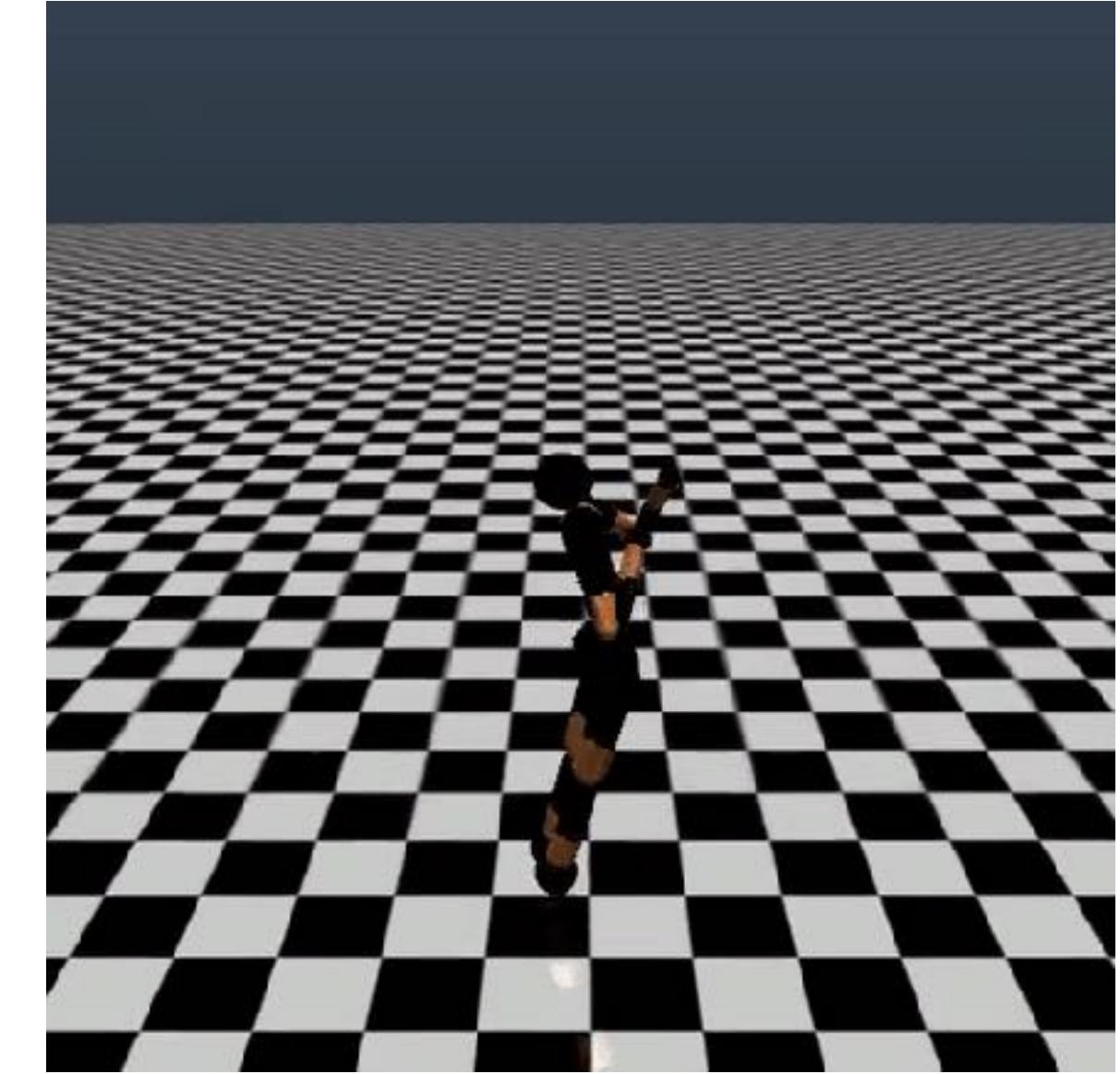
1. Train a policy  $\pi_\theta$  from expert data  $D$
2. Run  $\pi_\theta$  (on the robot) to obtain new states
3. Ask the expert to label the states with actions and obtain new data  $D'$
4. Add the labeled data to the training set  $D = D \cup D'$  (aggregation)

**Want:**  $\min_{\theta \in \Theta} \int \|\pi_\theta(x) - \pi_e(x)\|^2 \mu^{\pi_\theta}(dx)$

# Fix: covariate shift in imitation learning

DAgger (Data Aggregation)

1. Train a policy  $\pi_\theta$  from expert data  $D$
2. Run  $\pi_\theta$  (on the robot) to obtain new states
3. Ask the expert to label the states with actions and obtain new data  $D'$
4. Add the labeled data to the training set  $D = D \cup D'$  (aggregation)



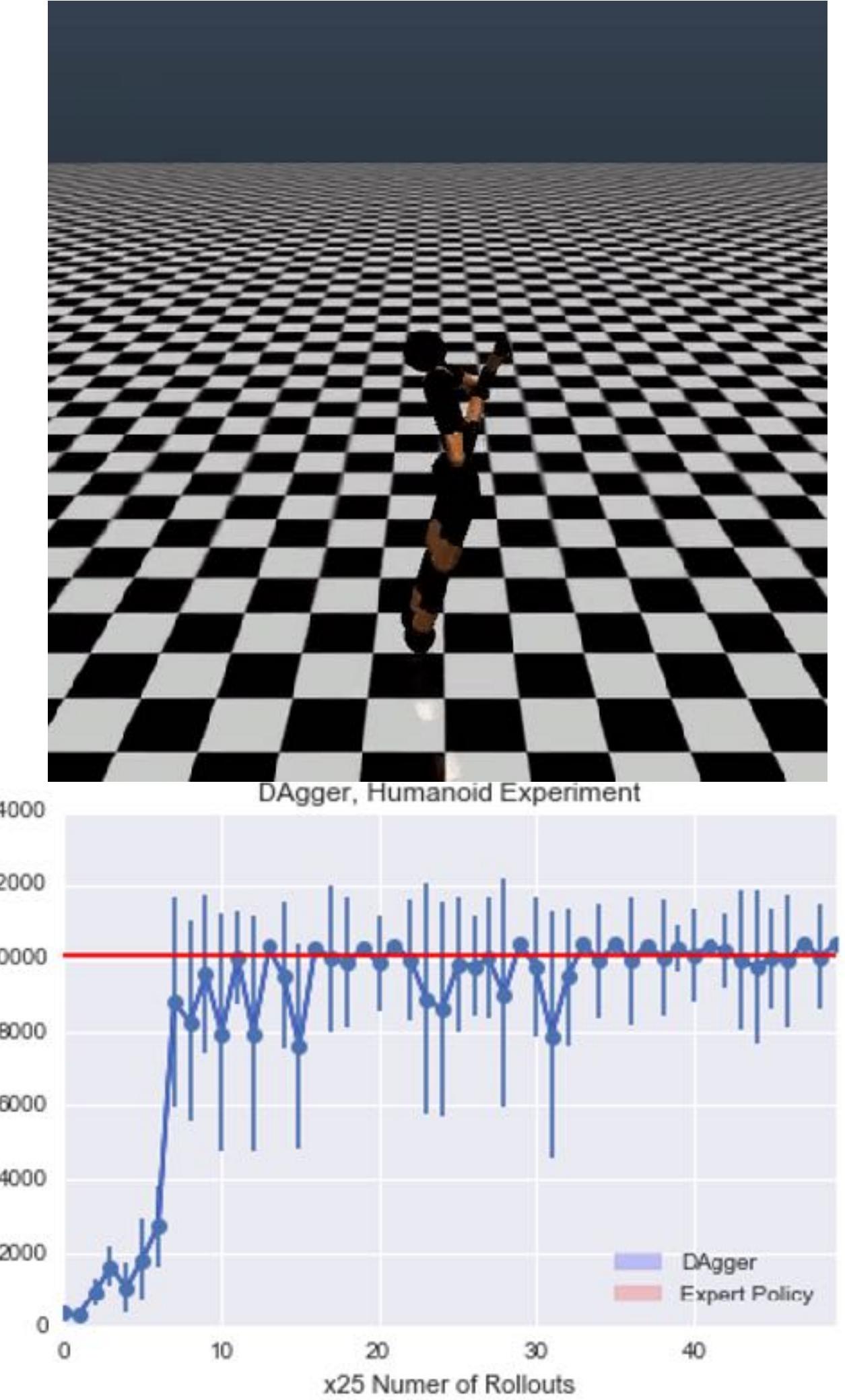
**Want:** 
$$\min_{\theta \in \Theta} \int \|\pi_\theta(x) - \pi_e(x)\|^2 \mu^{\pi_\theta}(dx)$$

# Fix: covariate shift in imitation learning

DAgger (Data Aggregation)

1. Train a policy  $\pi_\theta$  from expert data  $D$
2. Run  $\pi_\theta$  (on the robot) to obtain new states
3. Ask the expert to label the states with actions and obtain new data  $D'$
4. Add the labeled data to the training set  $D = D \cup D'$  (aggregation)

**Want:**  $\min_{\theta \in \Theta} \int \|\pi_\theta(x) - \pi_e(x)\|^2 \mu^{\pi_\theta}(dx)$

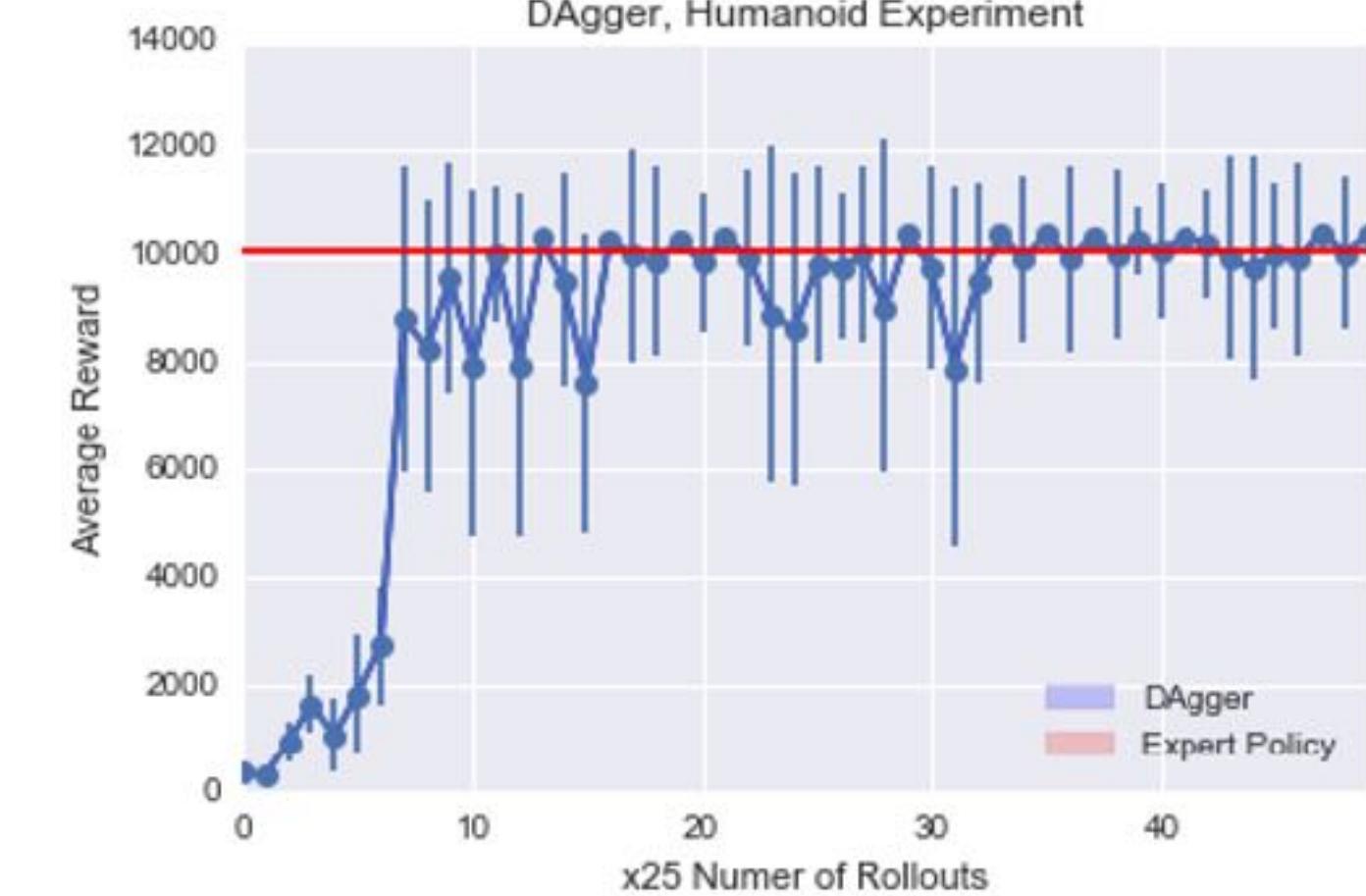
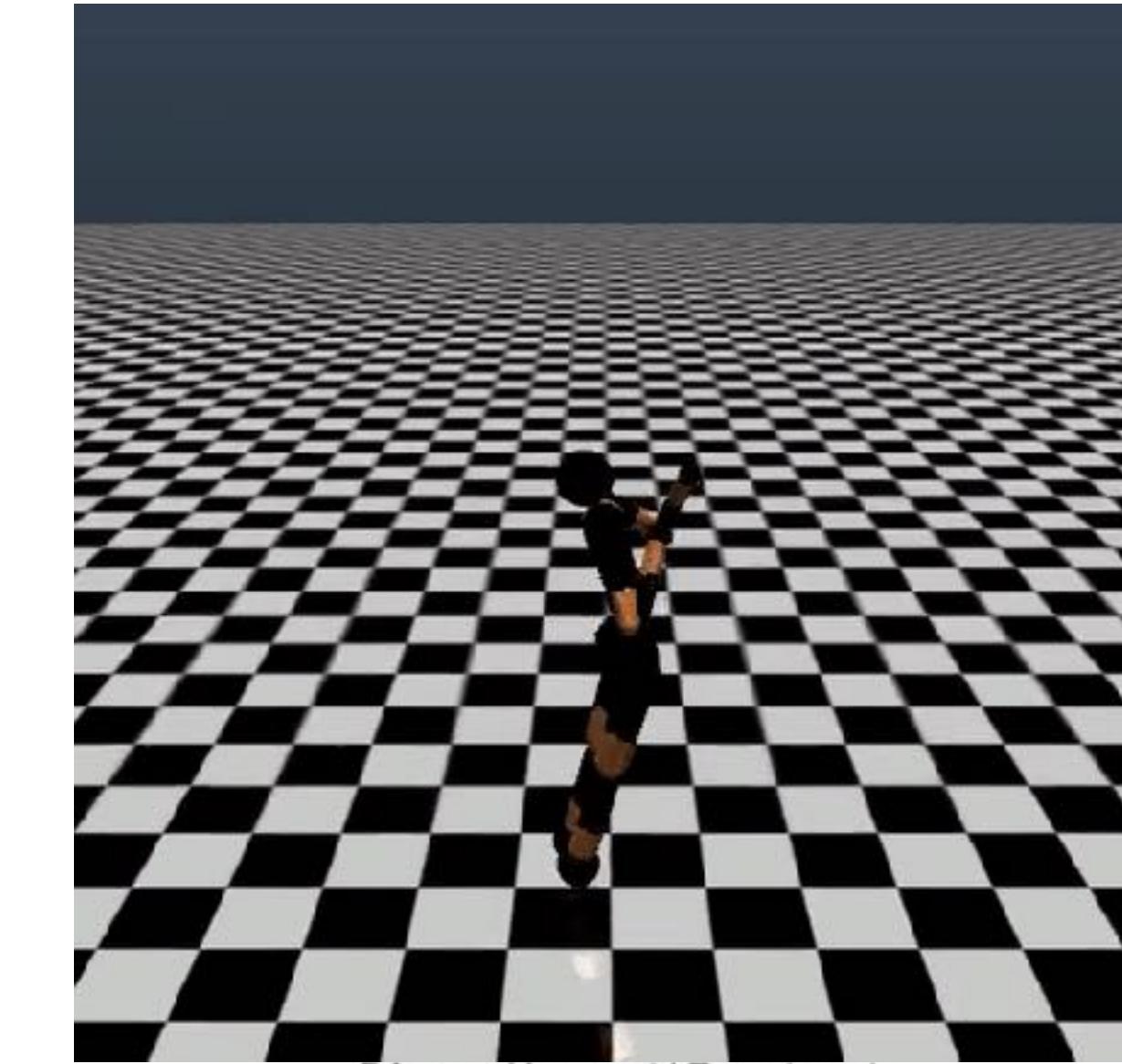


# Fix: covariate shift in imitation learning

DAgger (Data Aggregation)

1. Train a policy  $\pi_\theta$  from expert data  $D$
2. Run  $\pi_\theta$  (on the robot) to obtain new states
3. Ask the expert to label the states with actions and obtain new data  $D'$
4. Add the labeled data to the training set  $D = D \cup D'$  (aggregation)

**Want:**  $\min_{\theta \in \Theta} \int \|\pi_\theta(x) - \pi_e(x)\|^2 \mu^{\pi_\theta}(dx)$



Other related topics: off-policy/offline RL

# Conclusions

# Conclusions

- We have introduced distributional metrics such as the MMD, Wasserstein distance as tools for comparing distributions.

# Conclusions

- We have introduced distributional metrics such as the MMD, Wasserstein distance as tools for comparing distributions.
- We learned about RO and DRO, and the principles of robustness against distribution shift in optimization and machine learning.

# Conclusions

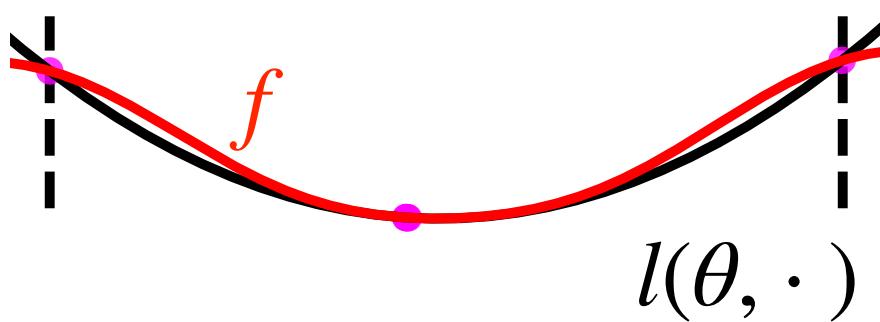
- We have introduced distributional metrics such as the MMD, Wasserstein distance as tools for comparing distributions.
- We learned about RO and DRO, and the principles of robustness against distribution shift in optimization and machine learning.
- We looked at the robustness and distribution shift of multi-stage decision-making and control.

# Conclusions

- We have introduced distributional metrics such as the MMD, Wasserstein distance as tools for comparing distributions.
- We learned about RO and DRO, and the principles of robustness against distribution shift in optimization and machine learning.
- We looked at the robustness and distribution shift of multi-stage decision-making and control.
- We have applied the robustness principle for decision-making :)

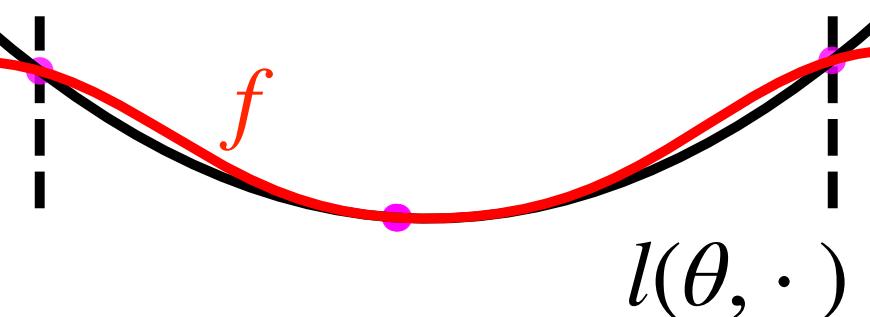
# Conclusions

- We have introduced distributional metrics such as the MMD, Wasserstein distance as tools for comparing distributions.
- We learned about RO and DRO, and the principles of robustness against distribution shift in optimization and machine learning.
- We looked at the robustness and distribution shift of multi-stage decision-making and control.
- We have applied the robustness principle for decision-making :)
- Takeaway
  - **Flatten the curve, smooth is robust**



# Conclusions

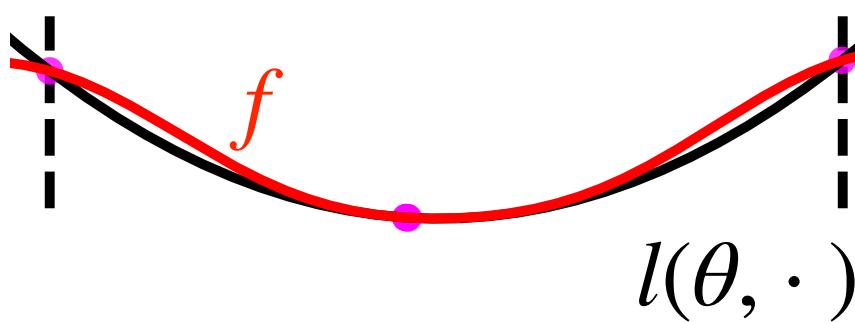
- We have introduced distributional metrics such as the MMD, Wasserstein distance as tools for comparing distributions.
- We learned about RO and DRO, and the principles of robustness against distribution shift in optimization and machine learning.
- We looked at the robustness and distribution shift of multi-stage decision-making and control.
- We have applied the robustness principle for decision-making :)
- Takeaway
  - **Flatten the curve, smooth is robust**



# Other topics

# Conclusions

- We have introduced distributional metrics such as the MMD, Wasserstein distance as tools for comparing distributions.
- We learned about RO and DRO, and the principles of robustness against distribution shift in optimization and machine learning.
- We looked at the robustness and distribution shift of multi-stage decision-making and control.
- We have applied the robustness principle for decision-making :)
- Takeaway
  - **Flatten the curve, smooth is robust**



# Other topics

- Generalizaiton and statistical bounds of DRO, minimax rate of non-parametric statistics
- Distributionally robust MPC
- Distributionally robust dynamics learning
- Off-policy RL, off-line RL
- Multi-stage DRO
- Adversarial attacks
- Causal inference
- Double descent
- Robust generative modeling
- You tell me!

# Thank you!

J.J. (Jia-Jie) Zhu

[jj-zhu.github.io](https://jj-zhu.github.io)

[zhu@wias-berlin.de](mailto:zhu@wias-berlin.de)

Weierstrass-Institute, Berlin &  
Max-Planck-Institute, Tübingen  
Germany