

# Kernel approaches to covariate shift

*Arthur Gretton*

Carnegie Mellon University  
Max Planck Institute for Biological Cybernetics

December 2009

# Transfer learning and covariate shift

---

- Patterns  $\mathcal{X}$ , labels  $\mathcal{Y}$
- **Training:** get  $Z_{\text{tr}}$  are  $n_{\text{tr}}$  pairs  $(x^{\text{tr}}, y^{\text{tr}})$  from  $\mathbf{P}_{\text{tr}}$
- **Test:** get  $Z_{\text{te}}$  are  $n_{\text{te}}$  pairs  $(x^{\text{te}}, y^{\text{te}})$  from  $\mathbf{P}_{\text{te}}$
- Predict on  $\mathbf{P}_{\text{te}}$  given data from  $\mathbf{P}_{\text{tr}}$
- **Examples:**
  - Medical diagnosis
  - Brain computer interfaces
  - Gene expression profiles

# Transfer learning and covariate shift

---

- Patterns  $\mathcal{X}$ , labels  $\mathcal{Y}$
- **Training:** get  $Z_{\text{tr}}$  are  $n_{\text{tr}}$  pairs  $(x^{\text{tr}}, y^{\text{tr}})$  from  $\mathbf{P}_{\text{tr}}$
- **Test:** get  $Z_{\text{te}}$  are  $n_{\text{te}}$  pairs  $(x^{\text{te}}, y^{\text{te}})$  from  $\mathbf{P}_{\text{te}}$
- Predict on  $\mathbf{P}_{\text{te}}$  given data from  $\mathbf{P}_{\text{tr}}$
- **Examples:**
  - Medical diagnosis
  - Brain computer interfaces
  - Gene expression profiles

Does this make sense?

# Transfer learning and covariate shift

---

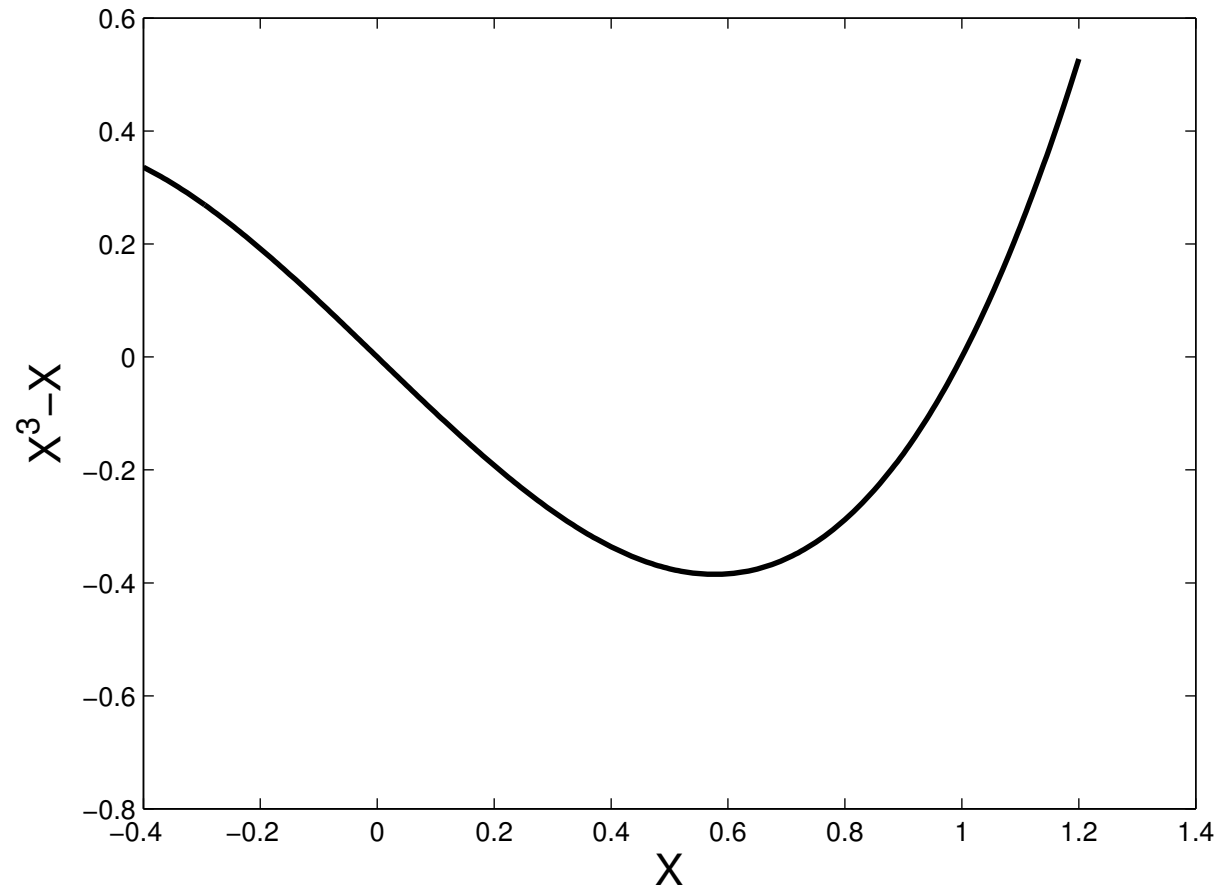
- Patterns  $\mathcal{X}$ , labels  $\mathcal{Y}$
- **Training:** get  $Z_{\text{tr}}$  are  $n_{\text{tr}}$  pairs  $(x^{\text{tr}}, y^{\text{tr}})$  from  $\mathbf{P}_{\text{tr}}$
- **Test:** get  $Z_{\text{te}}$  are  $n_{\text{te}}$  pairs  $(x^{\text{te}}, y^{\text{te}})$  from  $\mathbf{P}_{\text{te}}$
- Predict on  $\mathbf{P}_{\text{te}}$  given data from  $\mathbf{P}_{\text{tr}}$
- **Examples:**
  - Medical diagnosis
  - Brain computer interfaces
  - Gene expression profiles
- **Assumption:**  $\mathbf{P}_{\text{tr}}(x, y) = \mathbf{P}(y|x)\mathbf{P}_{\text{tr}}(x)$  and  $\mathbf{P}_{\text{te}}(x, y) = \mathbf{P}(y|x)\mathbf{P}_{\text{te}}(x)$

Conditional probs unchanged: **covariate shift**

# A toy example

---

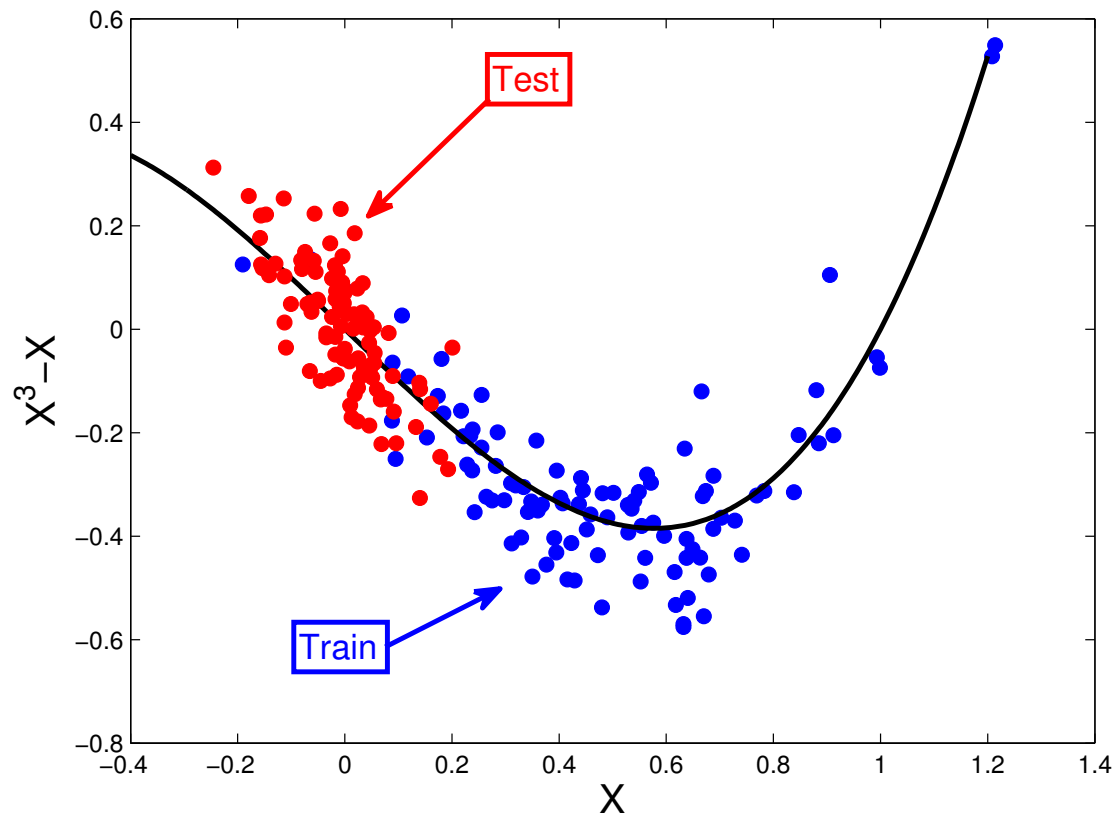
- Toy data [Shimodaira, 2000]
  - $\mathbf{P}_{\text{tr}}(x) \sim \mathcal{N}(0.5, 0.5^2)$ ,
  - $\mathbf{P}_{\text{te}}(x) \sim \mathcal{N}(0, 0.3^2)$
- $y = -x + x^3 + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, 0.3^2)$
- Linear regression



# A toy example

---

- Toy data [Shimodaira, 2000]
  - $\mathbf{P}_{\text{tr}}(x) \sim \mathcal{N}(0.5, 0.5^2)$ ,
  - $\mathbf{P}_{\text{te}}(x) \sim \mathcal{N}(0, 0.3^2)$
- $y = -x + x^3 + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, 0.3^2)$
- Linear regression



# A toy example

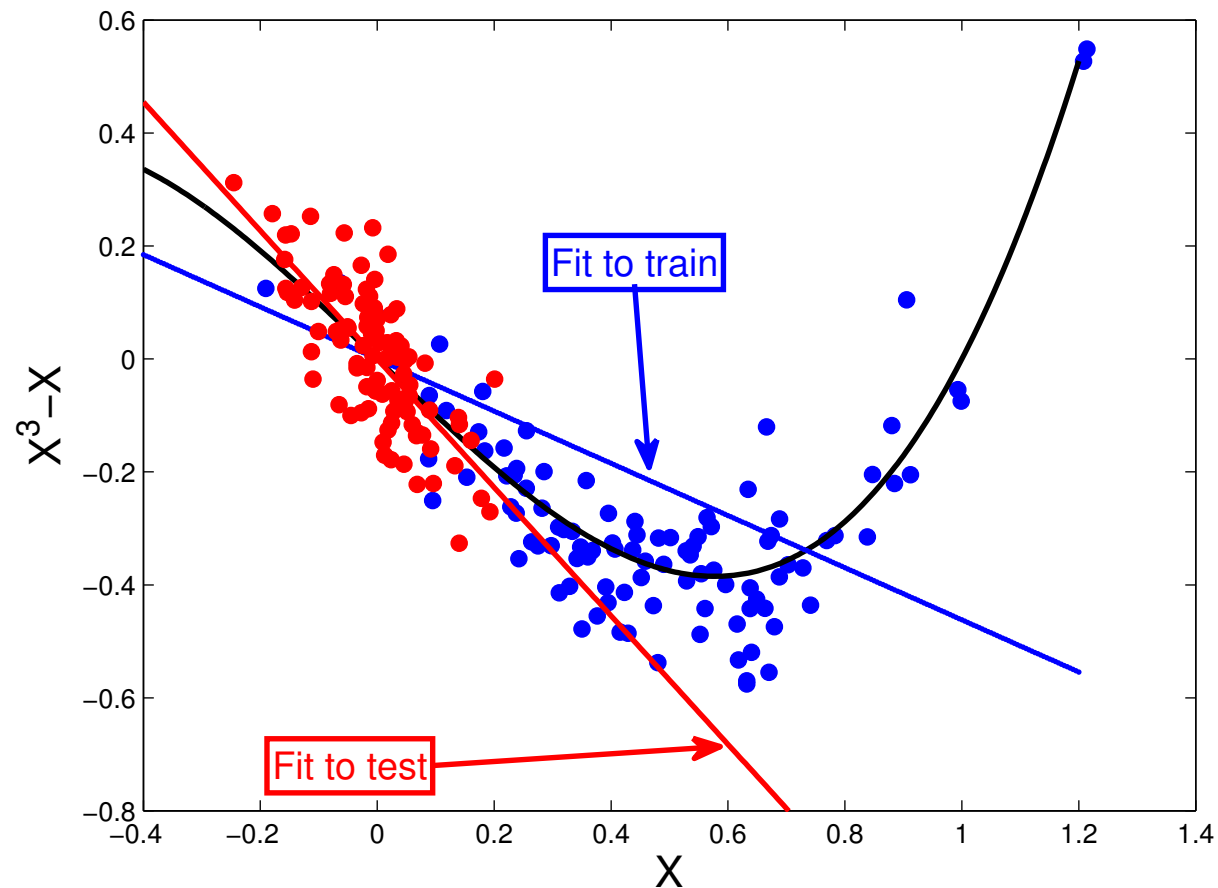
- Toy data [Shimodaira, 2000]

- $\mathbf{P}_{\text{tr}}(x) \sim \mathcal{N}(0.5, 0.5^2)$ ,

- $\mathbf{P}_{\text{te}}(x) \sim \mathcal{N}(0, 0.3^2)$

- $y = -x + x^3 + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, 0.3^2)$

- Linear regression



# The solution procedure

---

- Classical setting: (regularized) expected risk

$$R[\mathbf{P}, l(x, y, \theta)] = \mathbf{E} [l(x, y, \theta)] + \lambda \Omega[\theta]$$

- Loss  $l(x, y, \theta)$ , eg  $-\log \mathbf{P}(y|x, \theta)$
- Minimize over  $\theta$



# The solution procedure

---

- Classical setting: (regularized) expected risk

$$R[\mathbf{P}, l(x, y, \theta)] = \mathbf{E} [l(x, y, \theta)] + \lambda\Omega[\theta]$$

– Loss  $l(x, y, \theta)$ , eg  $-\log \mathbf{P}(y|x, \theta)$

– Minimize over  $\theta$

- Covariate shift setting:

$$\begin{aligned} R[\mathbf{P}_{\text{te}}, l(x, y, \theta)] &= \mathbf{E}_{\mathbf{P}_{\text{te}}} [l(x, y, \theta)] + \lambda\Omega[\theta] \\ &= \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x, y)l(x, y, \theta)] + \lambda\Omega[\theta] \end{aligned}$$

# The solution procedure

---

- Classical setting: (regularized) expected risk

$$R[\mathbf{P}, l(x, y, \theta)] = \mathbf{E} [l(x, y, \theta)] + \lambda\Omega[\theta]$$

- Loss  $l(x, y, \theta)$ , eg  $-\log \mathbf{P}(y|x, \theta)$
- Minimize over  $\theta$

- Covariate shift setting:

$$\begin{aligned} R[\mathbf{P}_{\text{te}}, l(x, y, \theta)] &= \mathbf{E}_{\mathbf{P}_{\text{te}}} [l(x, y, \theta)] + \lambda\Omega[\theta] \\ &= \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x, y)l(x, y, \theta)] + \lambda\Omega[\theta] \end{aligned}$$

- Importance weighting:

$$\mathbf{E}_{\mathbf{P}_{\text{te}}} [l(x, y, \theta)] = \mathbf{E}_{\mathbf{P}_{\text{tr}}} \left[ \underbrace{\frac{\mathbf{P}_{\text{te}}(x, y)}{\mathbf{P}_{\text{tr}}(x, y)}}_{:=\beta_{\text{imp}}(x, y)} l(x, y, \theta) \right] \quad \text{provided } \mathbf{P}_{\text{te}} \ll \mathbf{P}_{\text{tr}}$$

# Importance weighting

---

- **Variance** of importance weighted risk [Robert and Casella, 2004]

$$\begin{aligned} & \text{var} \left( l(x, y, \theta) \frac{\mathbf{P}_{\text{te}}(x, y)}{\mathbf{P}_{\text{tr}}(x, y)} \right) \\ &= \mathbf{E}_{\mathbf{P}_{\text{tr}}} \left[ l^2(x, y, \theta) \frac{\mathbf{P}_{\text{te}}^2(x, y)}{\mathbf{P}_{\text{tr}}^2(x, y)} \right] - (\mathbf{E}_{\mathbf{P}_{\text{te}}} [l(x, y, \theta)])^2 \end{aligned}$$

# Importance weighting

---

- **Variance** of importance weighted risk [Robert and Casella, 2004]

$$\begin{aligned} & \text{var} \left( l(x, y, \theta) \frac{\mathbf{P}_{\text{te}}(x, y)}{\mathbf{P}_{\text{tr}}(x, y)} \right) \\ &= \mathbf{E}_{\mathbf{P}_{\text{tr}}} \left[ l^2(x, y, \theta) \frac{\mathbf{P}_{\text{te}}^2(x, y)}{\mathbf{P}_{\text{tr}}^2(x, y)} \right] - R^2[\mathbf{P}_{\text{te}}, \theta, l(x, y, \theta)] \\ &= \mathbf{E}_{\mathbf{P}_{\text{te}}} \left[ l^2(x, y, \theta) \frac{\mathbf{P}_{\text{te}}(x, y)}{\mathbf{P}_{\text{tr}}(x, y)} \right] - R^2[\mathbf{P}_{\text{te}}, \theta, l(x, y, \theta)] \end{aligned}$$

# Importance weighting

---

- **Variance** of importance weighted risk [Robert and Casella, 2004]

$$\begin{aligned} & \text{var} \left( l(x, y, \theta) \frac{\mathbf{P}_{\text{te}}(x, y)}{\mathbf{P}_{\text{tr}}(x, y)} \right) \\ &= \mathbf{E}_{\mathbf{P}_{\text{tr}}} \left[ l^2(x, y, \theta) \frac{\mathbf{P}_{\text{te}}^2(x, y)}{\mathbf{P}_{\text{tr}}^2(x, y)} \right] - R^2[\mathbf{P}_{\text{te}}, \theta, l(x, y, \theta)] \\ &= \mathbf{E}_{\mathbf{P}_{\text{te}}} \left[ l^2(x, y, \theta) \underbrace{\frac{\mathbf{P}_{\text{te}}(x, y)}{\mathbf{P}_{\text{tr}}(x, y)}}_{\leq B} \right] - R^2[\mathbf{P}_{\text{te}}, \theta, l(x, y, \theta)] \end{aligned}$$

# Importance weighting

---

- **Variance** of importance weighted risk [Robert and Casella, 2004]

$$\begin{aligned} & \text{var} \left( l(x, y, \theta) \frac{\mathbf{P}_{\text{te}}(x, y)}{\mathbf{P}_{\text{tr}}(x, y)} \right) \\ &= \mathbf{E}_{\mathbf{P}_{\text{tr}}} \left[ l^2(x, y, \theta) \frac{\mathbf{P}_{\text{te}}^2(x, y)}{\mathbf{P}_{\text{tr}}^2(x, y)} \right] - R^2[\mathbf{P}_{\text{te}}, \theta, l(x, y, \theta)] \\ &= \mathbf{E}_{\mathbf{P}_{\text{te}}} \left[ l^2(x, y, \theta) \underbrace{\frac{\mathbf{P}_{\text{te}}(x, y)}{\mathbf{P}_{\text{tr}}(x, y)}}_{\leq B} \right] - R^2[\mathbf{P}_{\text{te}}, \theta, l(x, y, \theta)] \end{aligned}$$

- $\mathbf{P}_{\text{tr}}$  should have **heavier tails** than  $\mathbf{P}_{\text{te}}$

# Importance weighting

---

- **Example:** kernel ridge regression
- Loss  $l(x, y, \theta) = (y - \langle \Phi(x), \theta \rangle)^2$

# Importance weighting

---

- **Example:** kernel ridge regression
- Loss  $l(x, y, \theta) = (y - \langle \Phi(x), \theta \rangle)^2$
- Solve

$$\underset{\theta}{\text{minimize}} \sum_{i=1}^{n_{\text{tr}}} \beta_i (y_i^{\text{tr}} - \langle \Phi(x_i^{\text{tr}}), \theta \rangle)^2 + \lambda \|\theta\|^2. \quad (2)$$



# Importance weighting

---

- **Example:** kernel ridge regression
- Loss  $l(x, y, \theta) = (y - \langle \Phi(x), \theta \rangle)^2$
- Solve

$$\text{minimize}_{\theta} \sum_{i=1}^{n_{\text{tr}}} \beta_i (y_i^{\text{tr}} - \langle \Phi(x_i^{\text{tr}}), \theta \rangle)^2 + \lambda \|\theta\|^2. \quad (3)$$

- Equivalently:

$$\text{minimize}_{\alpha} (y - K\alpha)^{\top} \bar{\beta} (y - K\alpha) + \lambda \alpha^{\top} K \alpha$$

- $\bar{\beta} = \text{diag}(\beta_1, \dots, \beta_{n_{\text{tr}}})$
- $K_{ij} = k(x_i^{\text{tr}}, x_j^{\text{tr}}) = \langle \Phi(x_i^{\text{tr}}), \Phi(x_j^{\text{tr}}) \rangle$

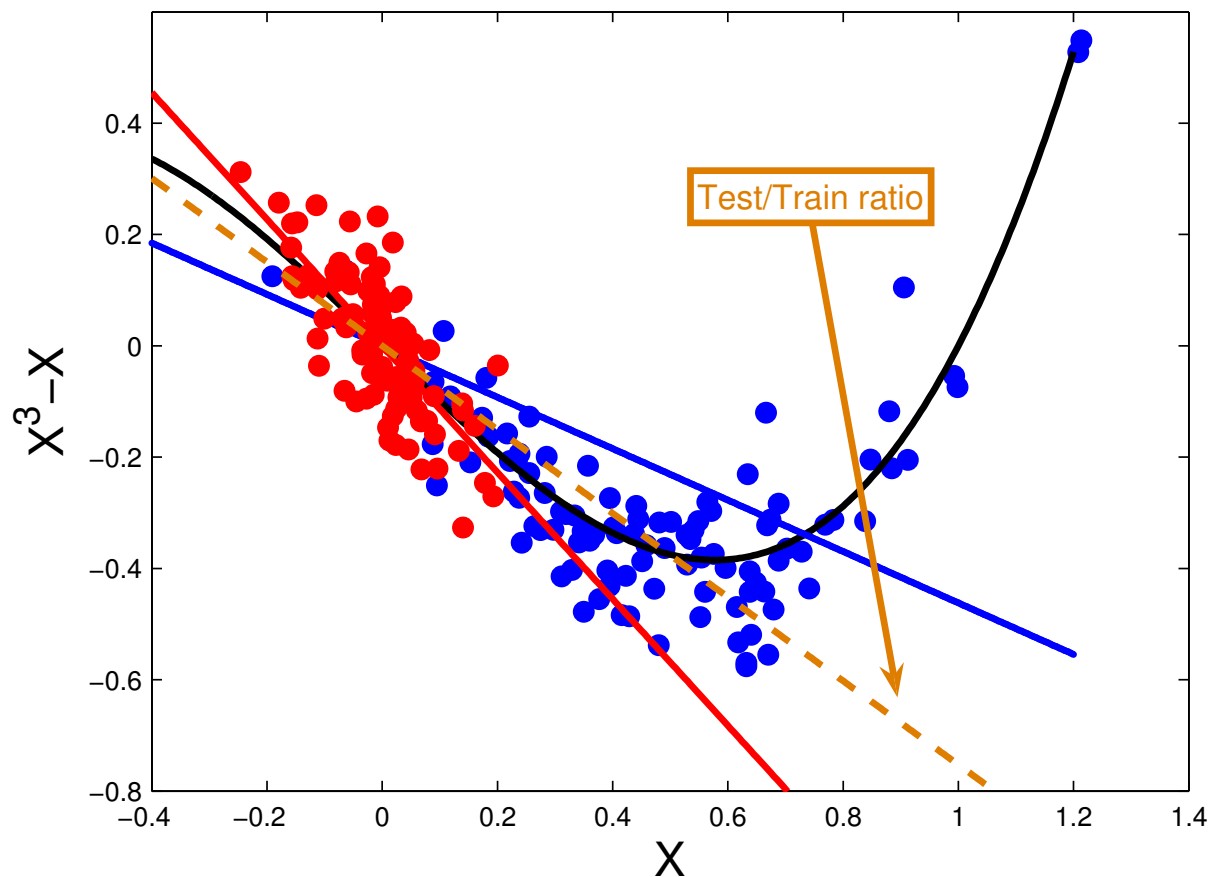
- **Solution**

$$\alpha = (\lambda \bar{\beta}^{-1} + K)^{-1} y$$

# Importance weighting

---

- Ridge regression, linear kernel
- Importance weighting improves performance



# Alternatives to density estimation

---

- Difficulties with direct density estimation
  - Empirical  $\mathbf{P}_{tr}$  and  $\mathbf{P}_{te}$  difficult for structured/high dimensional data
  - Variance can be large if empirical  $\mathbf{P}_{te}/\mathbf{P}_{tr}$  large

# Alternatives to density estimation

---

- Difficulties with direct density estimation
  - Empirical  $\mathbf{P}_{tr}$  and  $\mathbf{P}_{te}$  difficult for structured/high dimensional data
  - Variance can be large if empirical  $\mathbf{P}_{te}/\mathbf{P}_{tr}$  large
- Some other reweighting approaches:
  - Minimize classification error of  $\mathbf{P}_{tr}$  vs  $\mathbf{P}_{te}$  [Qin, 1998, Cheng and Chu, 2004, Bickel et al., 2009]
  - Minimize KL divergence between  $\mathbf{P}_{tr}$  and  $\mathbf{P}_{te}$  (KLIEP) [Sugiyama et al., 2008]
  - Ratio  $\mathbf{P}_{te}/\mathbf{P}_{tr}$  via least-squares function fitting [Kanamori et al., 2009]
  - Minimize Maximum Mean Discrepancy (MMD) between  $\mathbf{P}_{tr}$  and  $\mathbf{P}_{te}$  [Huang et al., 2007, Gretton et al., 2008]

Maximum mean discrepancy

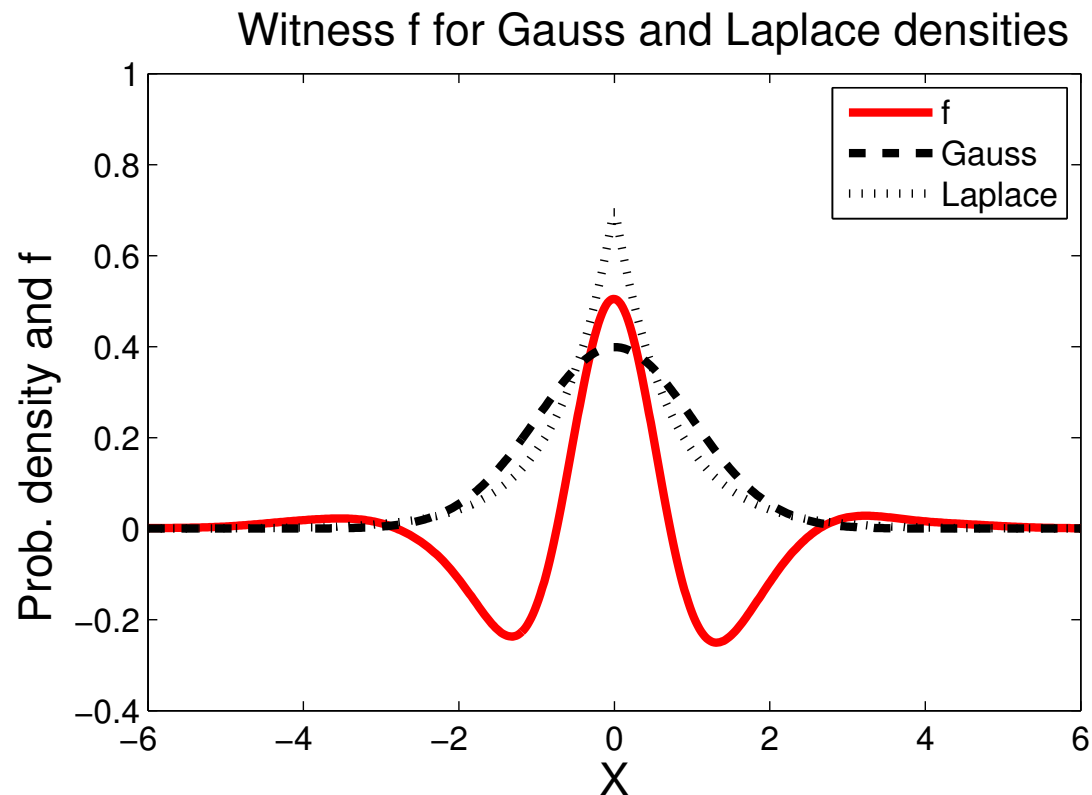
# Function Showing Difference in Distributions

---

- Idea: **avoid density estimation** when comparing distributions **P** and **Q**

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)].$$

- **Example:** Gauss **P** vs Laplace **Q**



# Function Showing Difference in Distributions

---

- Idea: **avoid density estimation** when comparing distributions  $\mathbf{P}$  and  $\mathbf{Q}$

[Fortet and Mourier, 1953]

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)].$$

- **Classical results:**  $\text{MMD}(\mathbf{P}, \mathbf{Q}; F) = 0$  iff  $\mathbf{P} = \mathbf{Q}$ , when
  - $F =$  bounded continuous [Dudley, 2002]
  - $F =$  bounded variation 1 (Kolmogorov metric) [Müller, 1997]
  - $F =$  bounded Lipschitz (Earth mover's distances) [Dudley, 2002]

# Function Showing Difference in Distributions

---

- Idea: **avoid density estimation** when comparing distributions  $\mathbf{P}$  and  $\mathbf{Q}$

[Fortet and Mourier, 1953]

$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)].$$

- **Classical results**:  $\text{MMD}(\mathbf{P}, \mathbf{Q}; F) = 0$  iff  $\mathbf{P} = \mathbf{Q}$ , when
  - $F =$  bounded continuous [Dudley, 2002]
  - $F =$  bounded variation 1 (Kolmogorov metric) [Müller, 1997]
  - $F =$  bounded Lipschitz (Earth mover's distances) [Dudley, 2002]
- $\text{MMD}(\mathbf{P}, \mathbf{Q}; F) = 0$  iff  $\mathbf{P} = \mathbf{Q}$  when  $F =$  the unit ball in a **characteristic RKHS**  $\mathcal{F}$  [Fukumizu et al., 2008, Sriperumbudur et al., 2008]



## Function Showing Difference in Distributions (2)

---

- $\mathcal{F}$  RKHS from  $\mathcal{X}$  to  $\mathbb{R}$  with positive definite kernel  $k(x_i, x_j)$
- $\mathcal{F} = \overline{\text{span}\{k(x, \cdot) \mid x \in \mathcal{X}\}}$ 
  - Example:  $f(\cdot) = \sum_{i=1}^m \alpha_i k(x_i, \cdot)$  for arbitrary  $m \in \mathbb{N}$ ,  $\alpha_i \in \mathbb{R}$ ,  $x_i \in \mathcal{X}$ .

# Function Showing Difference in Distributions (2)

---

- $\mathcal{F}$  RKHS from  $\mathcal{X}$  to  $\mathbb{R}$  with positive definite kernel  $k(x_i, x_j)$
- $\mathcal{F} = \overline{\text{span}\{k(x, \cdot) \mid x \in \mathcal{X}\}}$ 
  - Example:  $f(\cdot) = \sum_{i=1}^m \alpha_i k(x_i, \cdot)$  for arbitrary  $m \in \mathbb{N}$ ,  $\alpha_i \in \mathbb{R}$ ,  $x_i \in \mathcal{X}$ .
- Kernel is inner product between two feature maps:

$$\langle \Phi(x_1), \Phi(x_2) \rangle_{\mathcal{F}} = k(x_1, x_2)$$

# Function Showing Difference in Distributions (2)

---

- $\mathcal{F}$  RKHS from  $\mathcal{X}$  to  $\mathbb{R}$  with positive definite kernel  $k(x_i, x_j)$
- $\mathcal{F} = \overline{\text{span}\{k(x, \cdot) \mid x \in \mathcal{X}\}}$ 
  - Example:  $f(\cdot) = \sum_{i=1}^m \alpha_i k(x_i, \cdot)$  for arbitrary  $m \in \mathbb{N}$ ,  $\alpha_i \in \mathbb{R}$ ,  $x_i \in \mathcal{X}$ .
- Kernel is inner product between two feature maps:

$$\langle \Phi(x_1), \Phi(x_2) \rangle_{\mathcal{F}} = k(x_1, x_2)$$

- Evaluating functions at  $x$

$$f(x) = \langle f, \Phi(x) \rangle_{\mathcal{F}}$$

- $\Phi(x)$  feature map

# Function Showing Difference in Distributions

---

- **The (kernel) MMD:** [Gretton et al., 2007]

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F)$$

$$= \left( \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)] \right)^2$$

# Function Showing Difference in Distributions

---

- **The (kernel) MMD:** [Gretton et al., 2007]

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F)$$

$$= \left( \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)] \right)^2$$

using

$$\begin{aligned} \mathbf{E}_{\mathbf{P}}(f(x)) &= \mathbf{E}_{\mathbf{P}} [\langle \Phi(x), f \rangle_{\mathcal{F}}] \\ &=: \langle \mu_x, f \rangle_{\mathcal{F}} \end{aligned}$$

# Function Showing Difference in Distributions

---

- **The (kernel) MMD:** [Gretton et al., 2007]

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F)$$

$$= \left( \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)] \right)^2$$

$$= \left( \sup_{f \in F} \langle f, \mu_x - \mu_y \rangle_{\mathcal{F}} \right)^2$$

using

$$\mathbf{E}_{\mathbf{P}}(f(x)) = \mathbf{E}_{\mathbf{P}} [\langle \Phi(x), f \rangle_{\mathcal{F}}]$$

$$=: \langle \mu_x, f \rangle_{\mathcal{F}}$$

# Function Showing Difference in Distributions

---

- **The (kernel) MMD:** [Gretton et al., 2007]

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F)$$

$$= \left( \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)] \right)^2$$

using

$$= \left( \sup_{f \in F} \langle f, \mu_x - \mu_y \rangle_{\mathcal{F}} \right)^2$$

$$\|\mu\|_{\mathcal{F}} = \sup_{f \in F} \langle f, \mu \rangle_{\mathcal{F}}$$

$$= \|\mu_x - \mu_y\|_{\mathcal{F}}^2$$

# Function Showing Difference in Distributions

---

- **The (kernel) MMD:** [Gretton et al., 2007]

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F)$$

$$= \left( \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(\mathbf{x}) - \mathbf{E}_{\mathbf{Q}} f(\mathbf{y})] \right)^2$$

$$= \left( \sup_{f \in F} \langle f, \mu_x - \mu_y \rangle_{\mathcal{F}} \right)^2$$

$$= \|\mu_x - \mu_y\|_{\mathcal{F}}^2$$

$$= \langle \mu_x - \mu_y, \mu_x - \mu_y \rangle_{\mathcal{F}}$$

$$= \mathbf{E}_{\mathbf{P}, \mathbf{P}} k(\mathbf{x}, \mathbf{x}') + \mathbf{E}_{\mathbf{Q}, \mathbf{Q}} k(\mathbf{y}, \mathbf{y}') - 2\mathbf{E}_{\mathbf{P}, \mathbf{Q}} k(\mathbf{x}, \mathbf{y})$$

- $\mathbf{x}'$  is a R.V. independent of  $\mathbf{x}$  with distribution **P**
- $\mathbf{y}'$  is a R.V. independent of  $\mathbf{y}$  with distribution **Q**.



# Function Showing Difference in Distributions

---

- **The (kernel) MMD:** [Gretton et al., 2007]

$$\text{MMD}^2(\mathbf{P}, \mathbf{Q}; F)$$

$$= \left( \sup_{f \in F} [\mathbf{E}_{\mathbf{P}} f(x) - \mathbf{E}_{\mathbf{Q}} f(y)] \right)^2$$

$$= \left( \sup_{f \in F} \langle f, \mu_x - \mu_y \rangle_{\mathcal{F}} \right)^2$$

$$= \|\mu_x - \mu_y\|_{\mathcal{F}}^2$$

$$= \langle \mu_x - \mu_y, \mu_x - \mu_y \rangle_{\mathcal{F}}$$

$$= \mathbf{E}_{\mathbf{P}, \mathbf{P}} k(x, x') + \mathbf{E}_{\mathbf{Q}, \mathbf{Q}} k(y, y') - 2\mathbf{E}_{\mathbf{P}, \mathbf{Q}} k(x, y)$$

- $x'$  is a R.V. independent of  $x$  with distribution  $\mathbf{P}$
- $y'$  is a R.V. independent of  $y$  with distribution  $\mathbf{Q}$ .

- **Mean map:**

$$\mu_x := \mathbf{E}_{\mathbf{P}} \Phi(x) = \int k(\cdot, x) d\mathbf{P}(x)$$

Transfer learning using maximum mean discrepancy

# Transfer learning by KMM

---

Kernel mean matching (KMM)

# Transfer learning by KMM

---

## Kernel mean matching (KMM)

- Reweight training points so feature means match

$$\underset{\beta}{\text{minimize}} \quad \|\mu(\mathbf{P}_{\text{te}}) - \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x)\Phi(x)]\|$$

$$\text{subject to } \beta(x) \geq 0 \text{ and } \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x)] = 1.$$

- If  $\mathbf{P}_{\text{te}} \ll \mathbf{P}_{\text{tr}}$ , characteristic kernel, solution is  $\mathbf{P}_{\text{te}}(x) = \beta_{\text{imp}}(x)\mathbf{P}_{\text{tr}}(x)$

# Transfer learning by KMM

---

## Kernel mean matching (KMM)

- Reweight training points so feature means match

$$\underset{\beta}{\text{minimize}} \quad \|\mu(\mathbf{P}_{\text{te}}) - \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x)\Phi(x)]\|$$

$$\text{subject to } \beta(x) \geq 0 \text{ and } \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x)] = 1.$$

- If  $\mathbf{P}_{\text{te}} \ll \mathbf{P}_{\text{tr}}$ , characteristic kernel, solution is  $\mathbf{P}_{\text{te}}(x) = \beta_{\text{imp}}(x)\mathbf{P}_{\text{tr}}(x)$
- What about non-characteristic?

# Transfer learning by KMM

---

## Kernel mean matching (KMM)

- Reweight training points so feature means match

$$\text{minimize}_{\beta} \quad \|\mu(\mathbf{P}_{\text{te}}) - \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x)\Phi(x)]\|$$

$$\text{subject to } \beta(x) \geq 0 \text{ and } \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x)] = 1.$$

- If  $\mathbf{P}_{\text{te}} \ll \mathbf{P}_{\text{tr}}$ , characteristic kernel, solution is  $\mathbf{P}_{\text{te}}(x) = \beta_{\text{imp}}(x)\mathbf{P}_{\text{tr}}(x)$
- Empirical:

$$\min_{\beta} \left\| \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \beta_i \Phi(x_i^{\text{tr}}) - \frac{1}{n_{\text{te}}} \sum_{i=1}^{n_{\text{te}}} \Phi(x_i^{\text{te}}) \right\|^2 = \frac{1}{n_{\text{tr}}^2} \beta^{\top} K \beta - \frac{2}{n_{\text{tr}}^2} \kappa^{\top} \beta + \text{const.}$$

# Transfer learning by KMM

---

## Kernel mean matching (KMM)

- Reweight training points so feature means match

$$\text{minimize}_{\beta} \quad \|\mu(\mathbf{P}_{\text{te}}) - \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x)\Phi(x)]\|$$

$$\text{subject to } \beta(x) \geq 0 \text{ and } \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x)] = 1.$$

- If  $\mathbf{P}_{\text{te}} \ll \mathbf{P}_{\text{tr}}$ , characteristic kernel, solution is  $\mathbf{P}_{\text{te}}(x) = \beta_{\text{imp}}(x)\mathbf{P}_{\text{tr}}(x)$
- Empirical:

$$\min_{\beta} \left\| \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \beta_i \Phi(x_i^{\text{tr}}) - \frac{1}{n_{\text{te}}} \sum_{i=1}^{n_{\text{te}}} \Phi(x_i^{\text{te}}) \right\|^2 = \frac{1}{n_{\text{tr}}^2} \beta^{\top} K \beta - \frac{2}{n_{\text{tr}}^2} \kappa^{\top} \beta + \text{const.}$$

$$\text{subject to } \beta_i \in [0, B] \quad \text{and} \quad \left| \sum_{i=1}^{n_{\text{tr}}} \beta_i - n_{\text{tr}} \right| \leq \sqrt{n_{\text{tr}}} \epsilon.$$

# Transfer learning by KMM

---

## Kernel mean matching (KMM)

- Reweight training points so feature means match

$$\text{minimize}_{\beta} \quad \|\mu(\mathbf{P}_{\text{te}}) - \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x)\Phi(x)]\|$$

$$\text{subject to } \beta(x) \geq 0 \text{ and } \mathbf{E}_{\mathbf{P}_{\text{tr}}} [\beta(x)] = 1.$$

- If  $\mathbf{P}_{\text{te}} \ll \mathbf{P}_{\text{tr}}$ , characteristic kernel, solution is  $\mathbf{P}_{\text{te}}(x) = \beta_{\text{imp}}(x)\mathbf{P}_{\text{tr}}(x)$
- Empirical:

$$\min_{\beta} \left\| \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \beta_i \Phi(x_i^{\text{tr}}) - \frac{1}{n_{\text{te}}} \sum_{i=1}^{n_{\text{te}}} \Phi(x_i^{\text{te}}) \right\|^2 = \frac{1}{n_{\text{tr}}^2} \beta^{\top} K \beta - \frac{2}{n_{\text{tr}}^2} \kappa^{\top} \beta + \text{const.}$$

$$\text{subject to } \beta_i \in [0, B] \quad \text{and} \quad \underbrace{\left| \sum_{i=1}^{n_{\text{tr}}} \beta_i - n_{\text{tr}} \right|}_{\leq \sqrt{n_{\text{tr}}}\epsilon} \leq \sqrt{n_{\text{tr}}}\epsilon.$$

$$\left[ \frac{1}{\sqrt{n_{\text{tr}}}} \sum_i \beta_{\text{imp}}(x_i^{\text{tr}}) - \sqrt{n_{\text{tr}}} \right] \xrightarrow{D} \mathcal{N}(0, \sigma^2)$$



# Transfer learning by KMM

---

- What if **given**  $\beta_{\text{imp}}$ : finite sample effects?
- Assume  $k(x, x) \leq R^2$  for all  $x \in \mathcal{X}$ .

# Transfer learning by KMM

---

- What if **given**  $\beta_{\text{imp}}$ : finite sample effects?
- Assume  $k(x, x) \leq R^2$  for all  $x \in \mathcal{X}$ .
- With probability at least  $1 - \delta$ ,

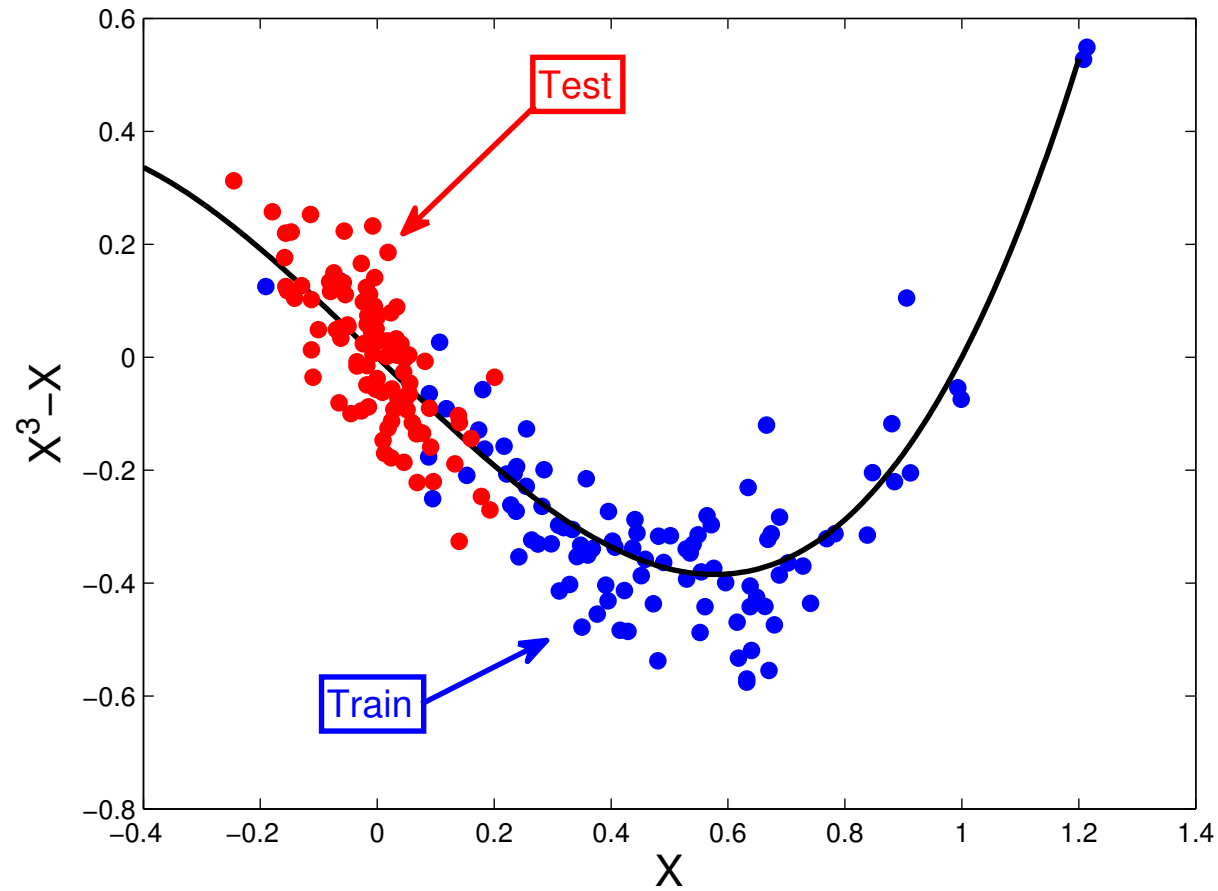
$$\left\| \frac{1}{n_{\text{tr}}} \sum_{i=1}^{n_{\text{tr}}} \beta_{\text{imp}}(x_i^{\text{tr}}) \Phi(x_i^{\text{tr}}) - \frac{1}{n_{\text{te}}} \sum_{i=1}^{n_{\text{te}}} \Phi(x_i^{\text{te}}) \right\| \leq \left( 1 + \sqrt{2 \log 2 / \delta} \right) R \sqrt{B^2 / n_{\text{tr}} + 1 / n_{\text{te}}}.$$

- Still (potentially) **high variance** for large  $B$ .
- **Convergence** of KMM procedure: [Cortes et al., 2008]

# Transfer learning by KMM

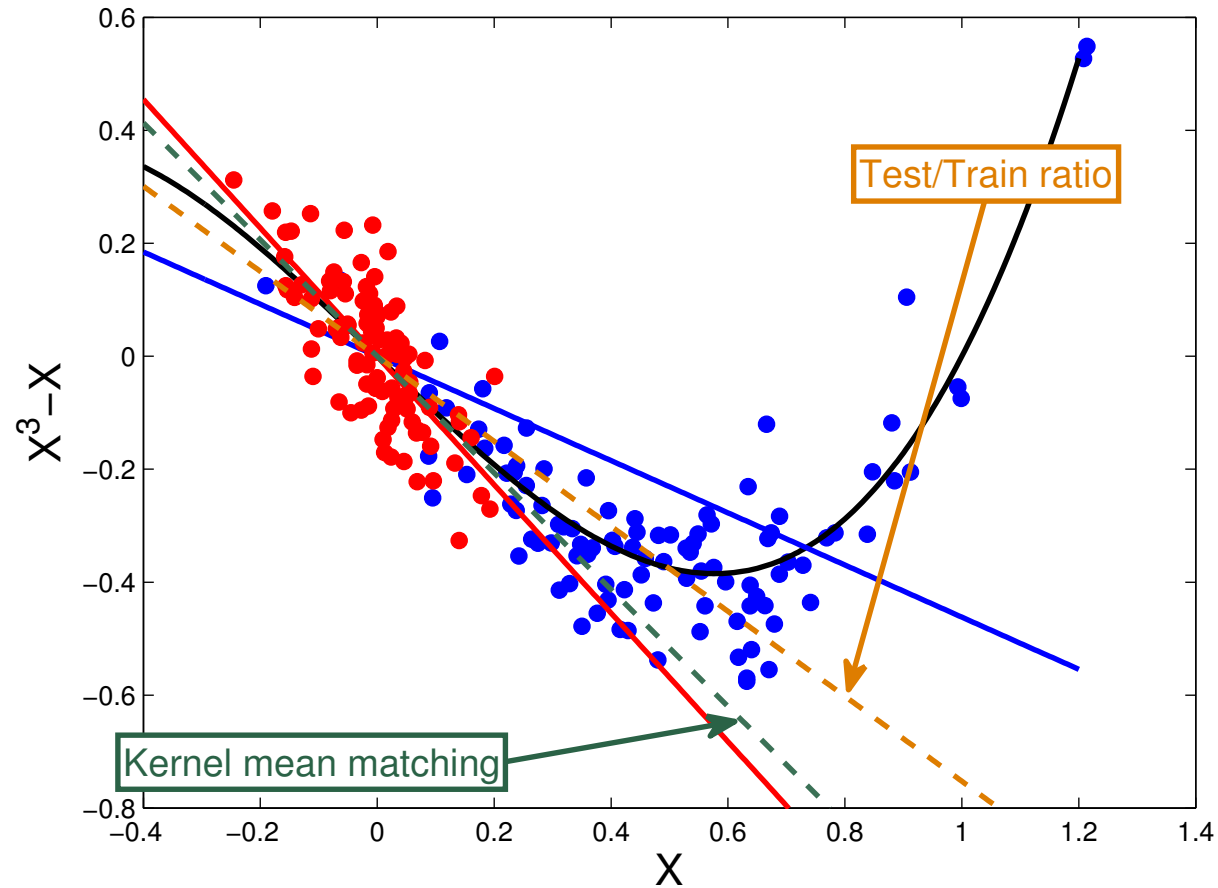
---

- Compare KMM and importance sampling



# Transfer learning by KMM

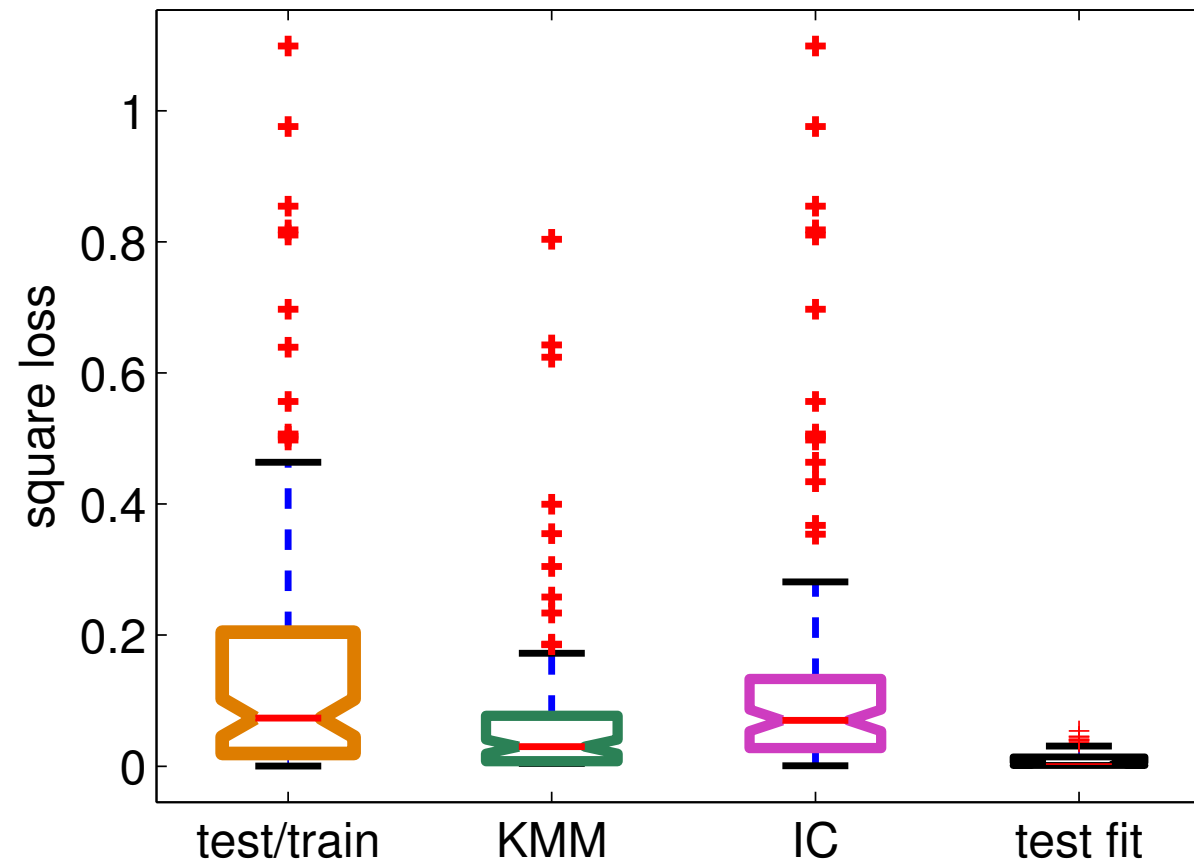
- Compare KMM and importance sampling



# Transfer learning by KMM

---

- Compare KMM and importance sampling



IC method due to [Shimodaira, 2000]

# Reweighting by classification

---

- Use train/test classification error to reweight [Qin, 1998, Cheng and Chu, 2004, Bickel et al., 2009]
- $\mathbf{P}(S|x^{\text{tr}}, x^{\text{te}}, \theta_{\text{shift}})$  classifies training ( $s = 1$ ) vs test ( $s = 0$ )

# Reweighting by classification

---

- Use train/test classification error to reweight [Qin, 1998, Cheng and Chu, 2004, Bickel et al., 2009]
- $\mathbf{P}(S|x^{\text{tr}}, x^{\text{te}}, \theta_{\text{shift}})$  classifies training ( $s = 1$ ) vs test ( $s = 0$ )
- Estimate **importance ratio**:

$$\frac{\mathbf{P}_{\text{te}}(x_i^{\text{tr}})}{\mathbf{P}_{\text{tr}}(x_i^{\text{tr}})} = \frac{\mathbf{P}(s = 1)}{\mathbf{P}(s = 0)} (\mathbf{P}^{-1}(s = 1|x_i^{\text{tr}}, \theta_{\text{shift}}) - 1)$$

- Learn **two** classifiers: train vs test and covariate to label

# Reweighting by classification

---

- Use train/test classification error to reweight [Qin, 1998, Cheng and Chu, 2004, Bickel et al., 2009]
- $\mathbf{P}(S|x^{\text{tr}}, x^{\text{te}}, \theta_{\text{shift}})$  classifies training ( $s = 1$ ) vs test ( $s = 0$ )
- Estimate **importance ratio**:

$$\frac{\mathbf{P}_{\text{te}}(x_i^{\text{tr}})}{\mathbf{P}_{\text{tr}}(x_i^{\text{tr}})} = \frac{\mathbf{P}(s = 1)}{\mathbf{P}(s = 0)} (\mathbf{P}^{-1}(s = 1|x_i^{\text{tr}}, \theta_{\text{shift}}) - 1)$$

- Learn **two** classifiers: train vs test and covariate to label
- Single joint optimization? [Bickel et al., 2009]

$$\max_{\theta_{\text{shift}}, \theta_{\text{learn}}} \mathbf{P}(y^{\text{tr}}|S, x^{\text{tr}}, \theta_{\text{shift}}, \theta_{\text{learn}}) \mathbf{P}(S|x^{\text{tr}}, x^{\text{te}}, \theta_{\text{shift}}) \mathbf{P}(\theta_{\text{shift}}) \mathbf{P}(\theta_{\text{learn}})$$



# Experiments

# Breast Cancer data

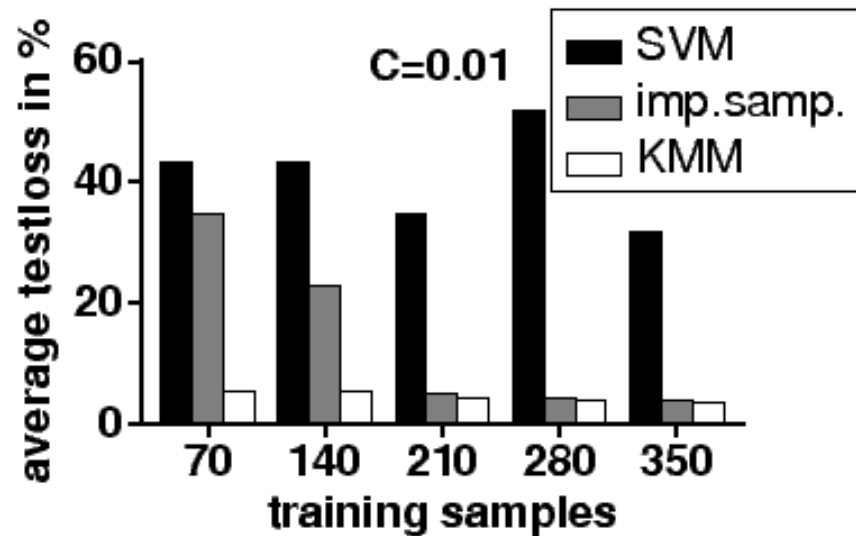
---

- Gaussian kernel  $\exp(-|x_i - x_j|^2/(2\sigma))$  for **KMM and SVN**,  $\sigma = 5$
- Performance vs  $C$ 
  - Small  $C$   $\rightarrow$  prioritize smoothness
- Selection procedure:
  - Random training/test split
  - Training set from 10% - 50% of test
  - $P(s_i = 1|x_i) \propto \exp(-0.05\|x_i - \bar{x}\|^2)$

## Breast Cancer data

---

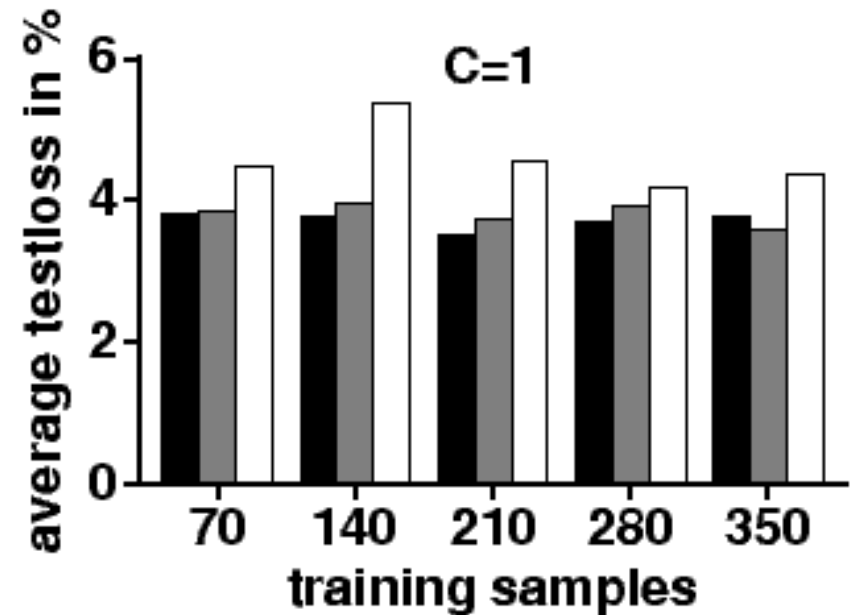
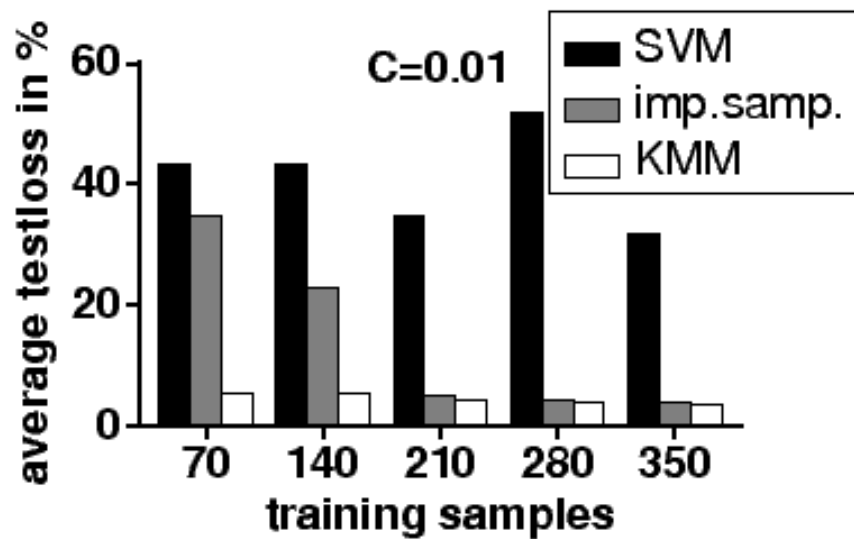
- Reweighting greatly improves performance
- KMM outperforms IS at small sample sizes



# Breast Cancer data

---

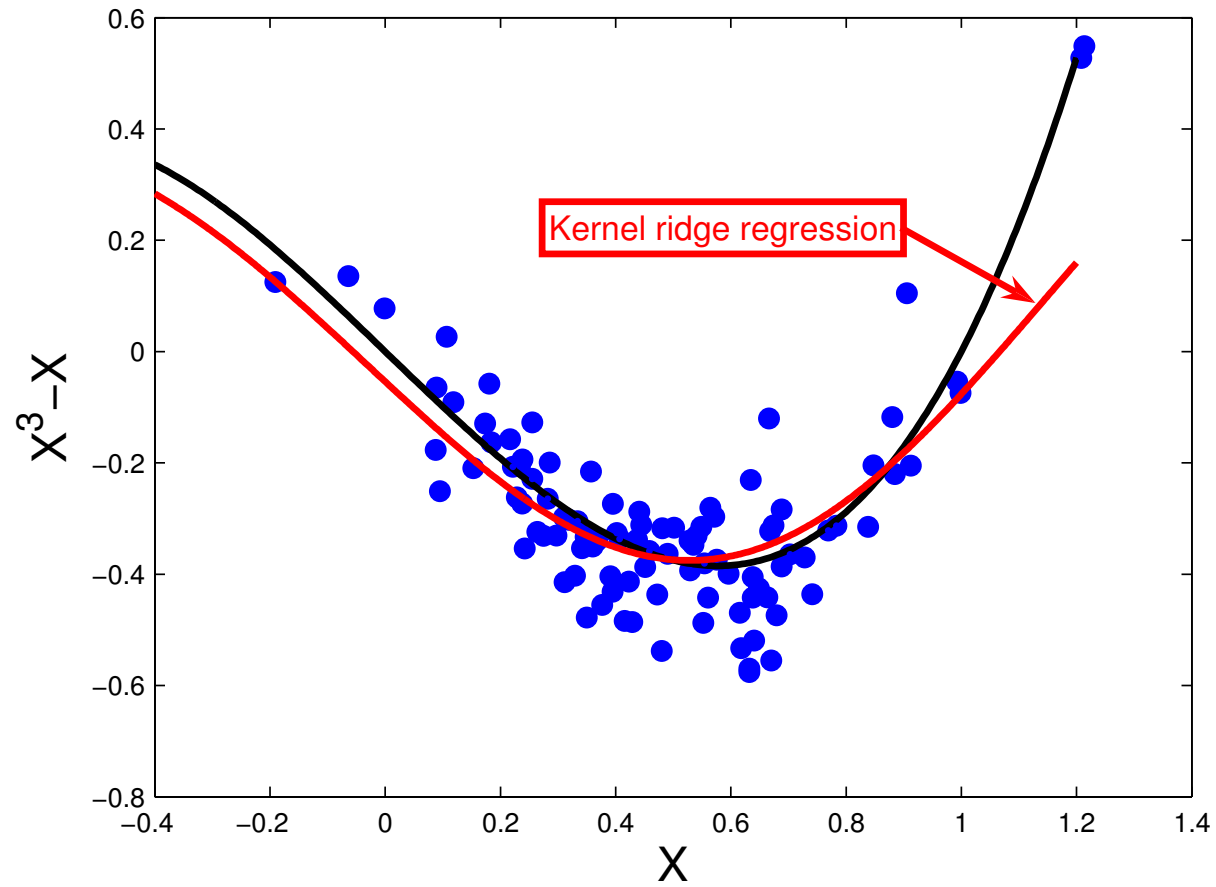
- KMM slightly decreases performance
- IS does not help



# Toy example revisited

---

- Kernel ridge regression result



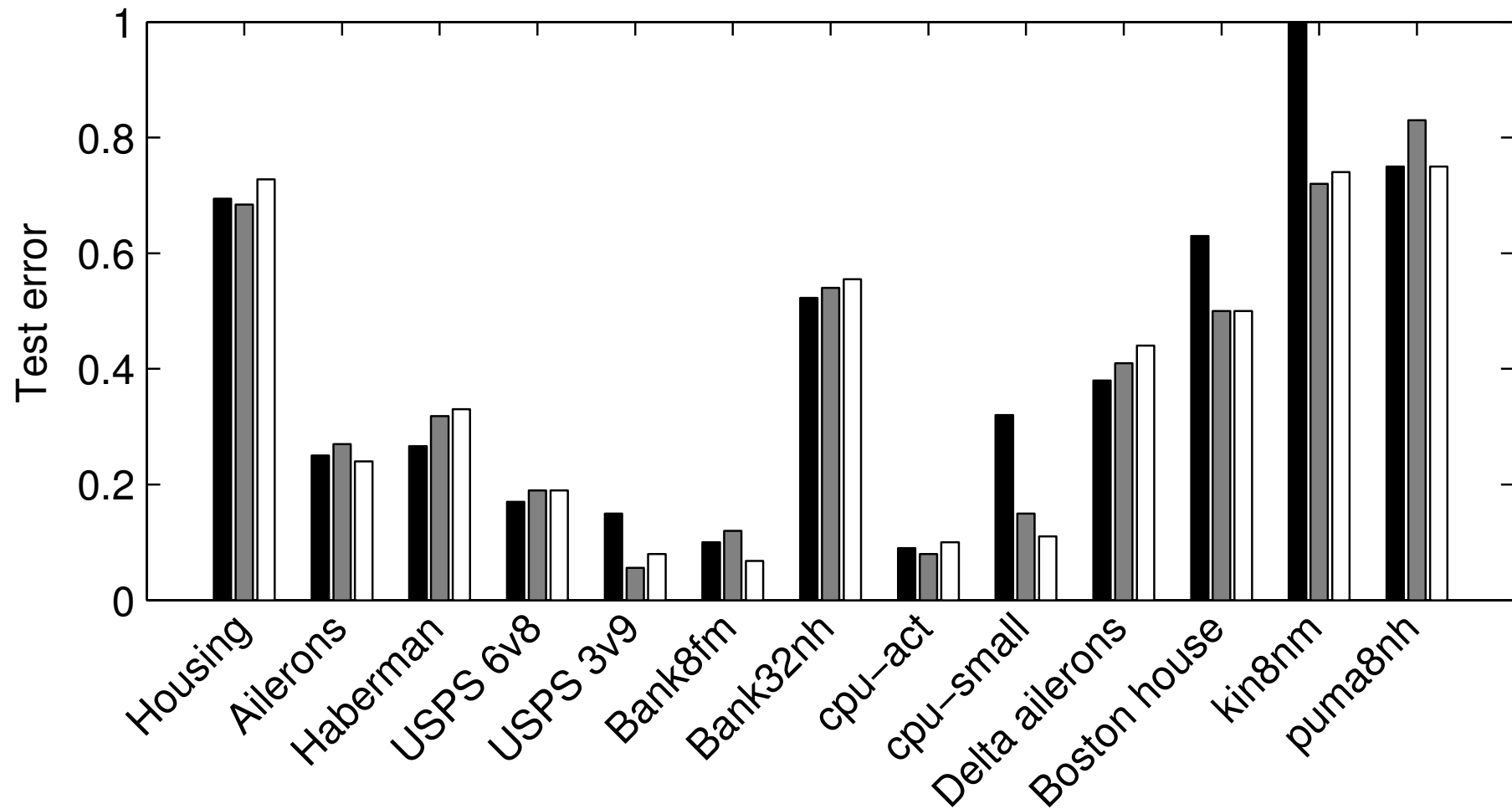
# Large scale experiments

---

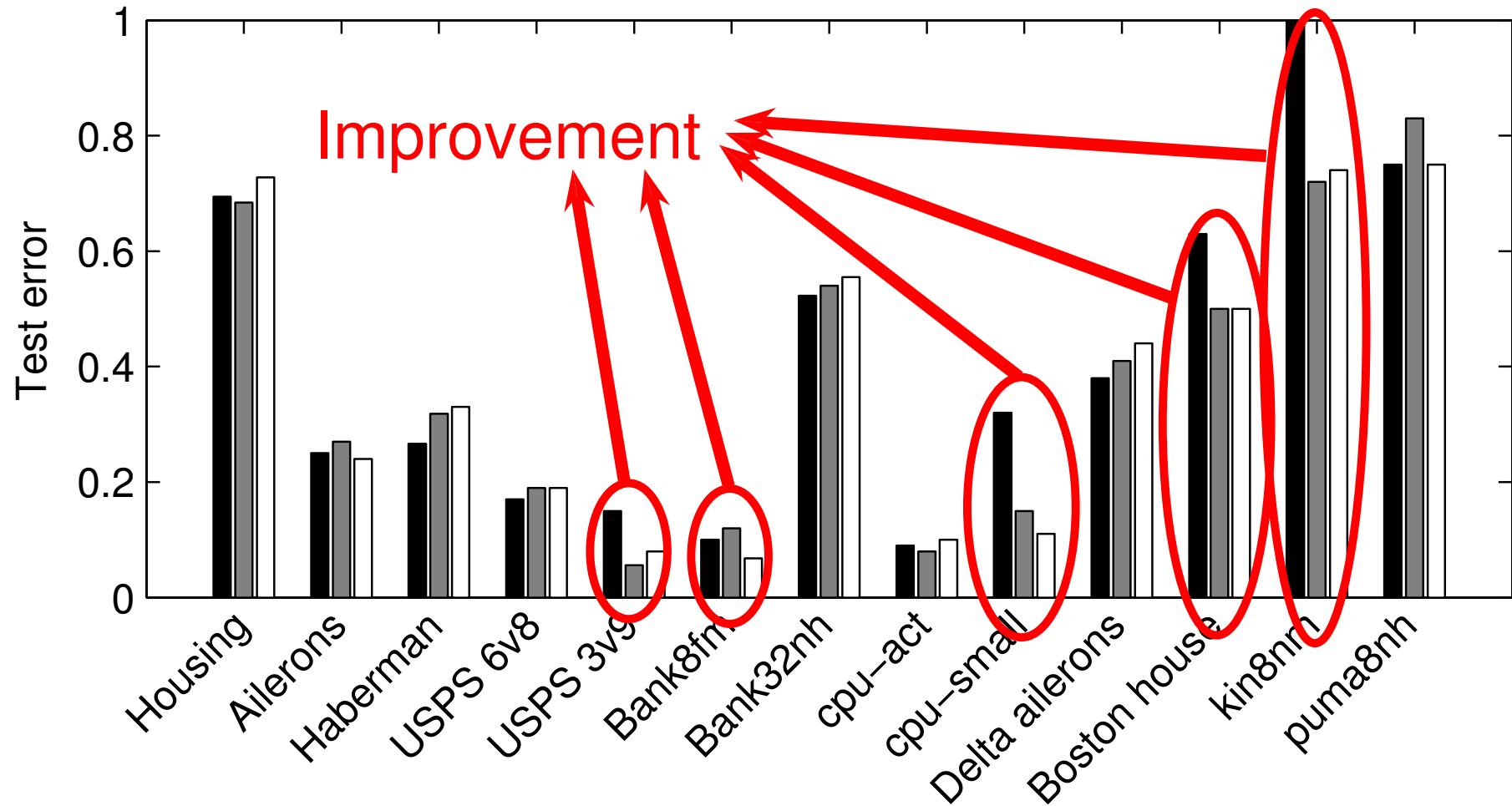
- Regression and classification
- Sampling scheme: training data missing at random
  - Sampling by Gaussian distribution on first principal component
- Cross validate on unweighted training set for  $C$  and  $\sigma$
- Same  $\sigma$  for classifier/regressor and KMM

# Large scale experiments

---

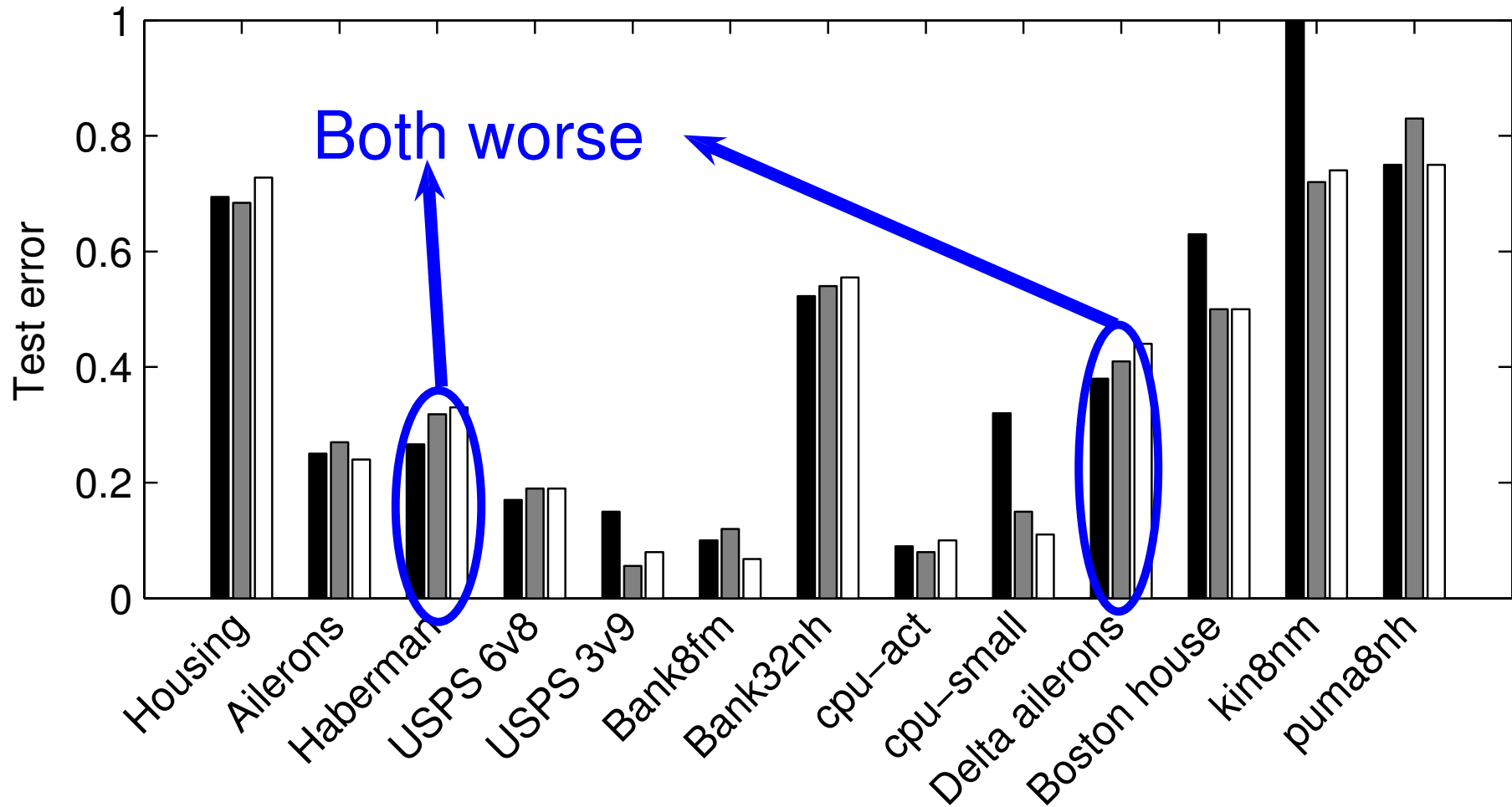


# Large scale experiments

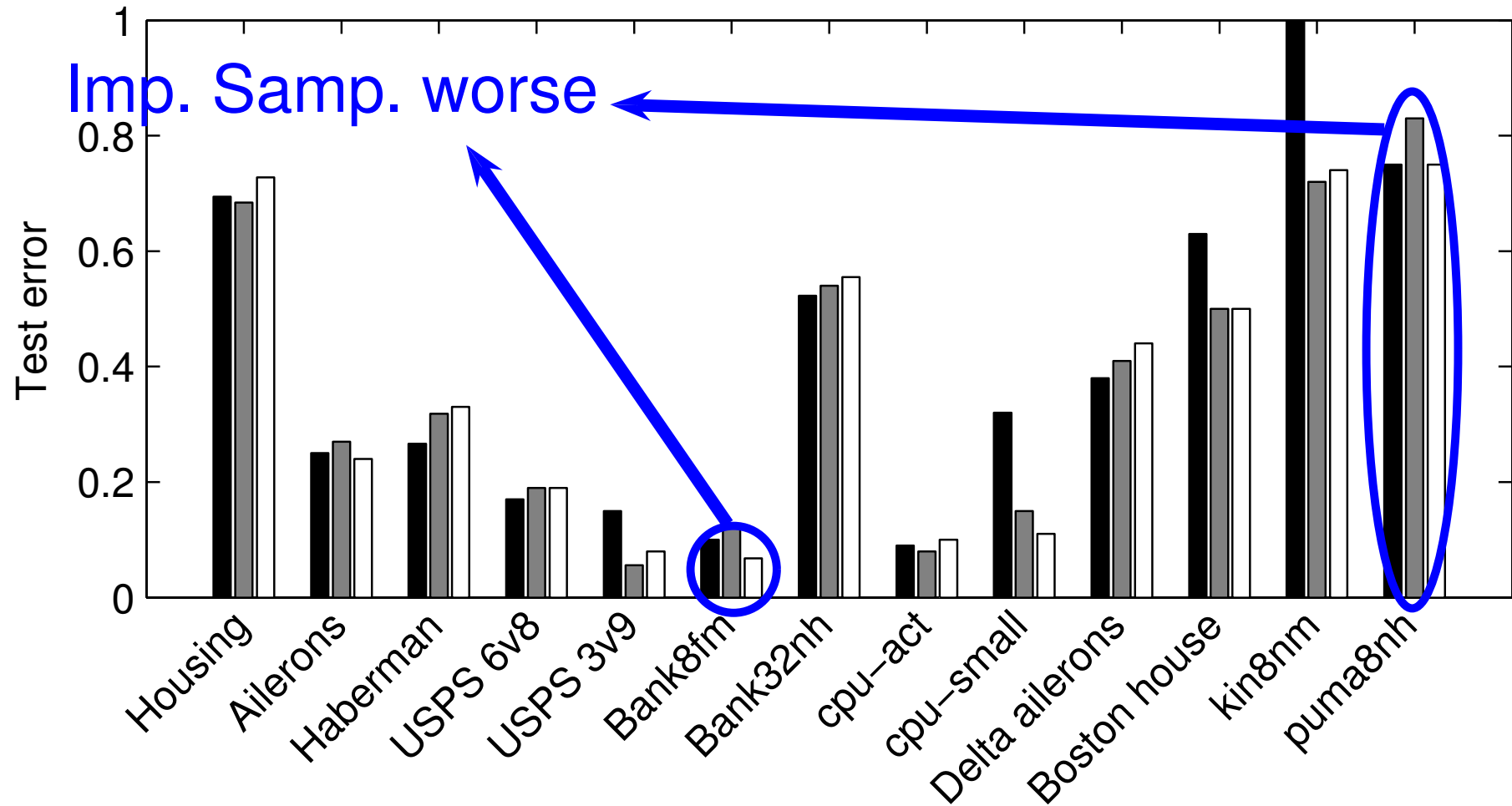




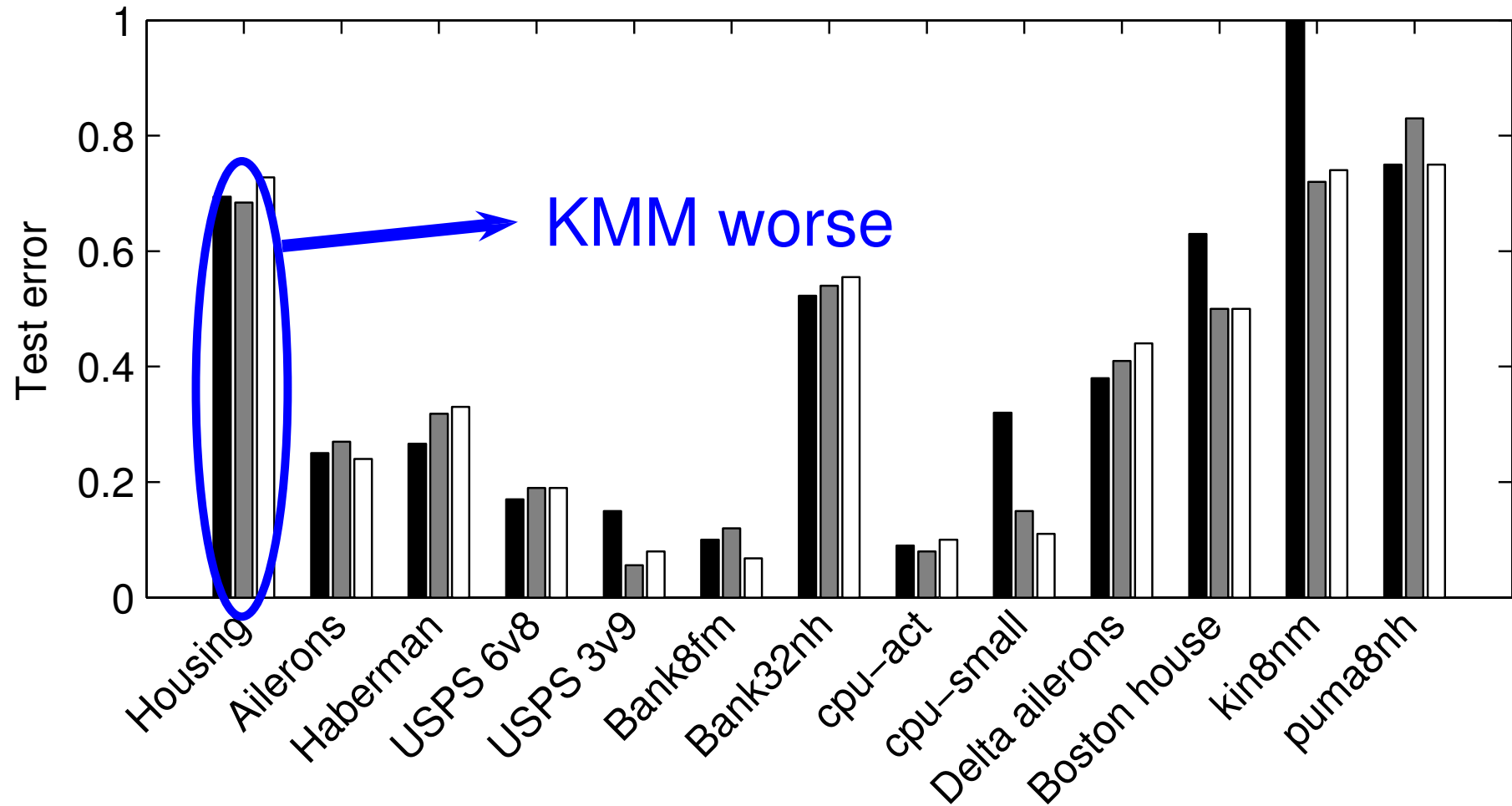
# Large scale experiments



# Large scale experiments

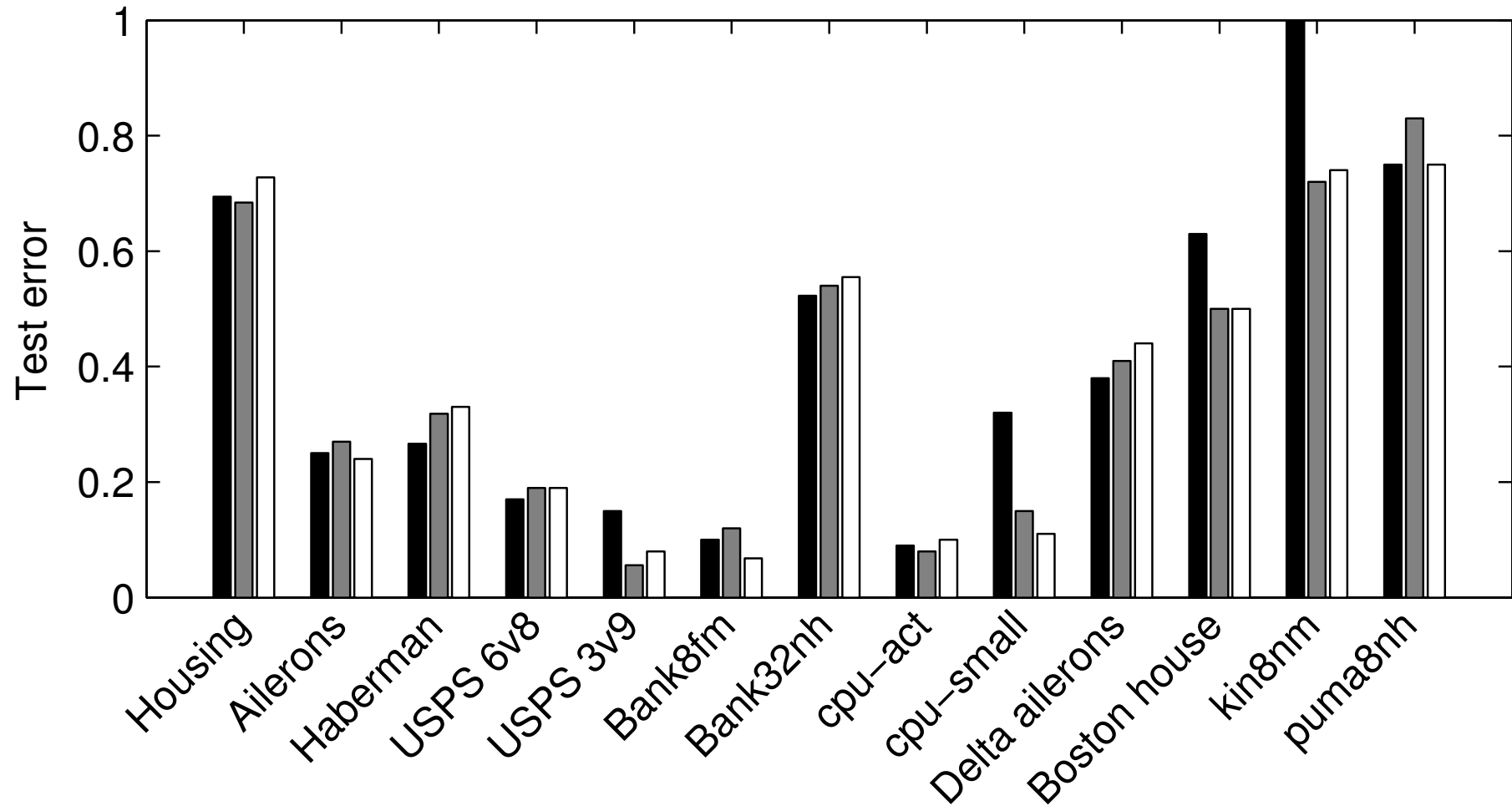


# Large scale experiments



# Large scale experiments

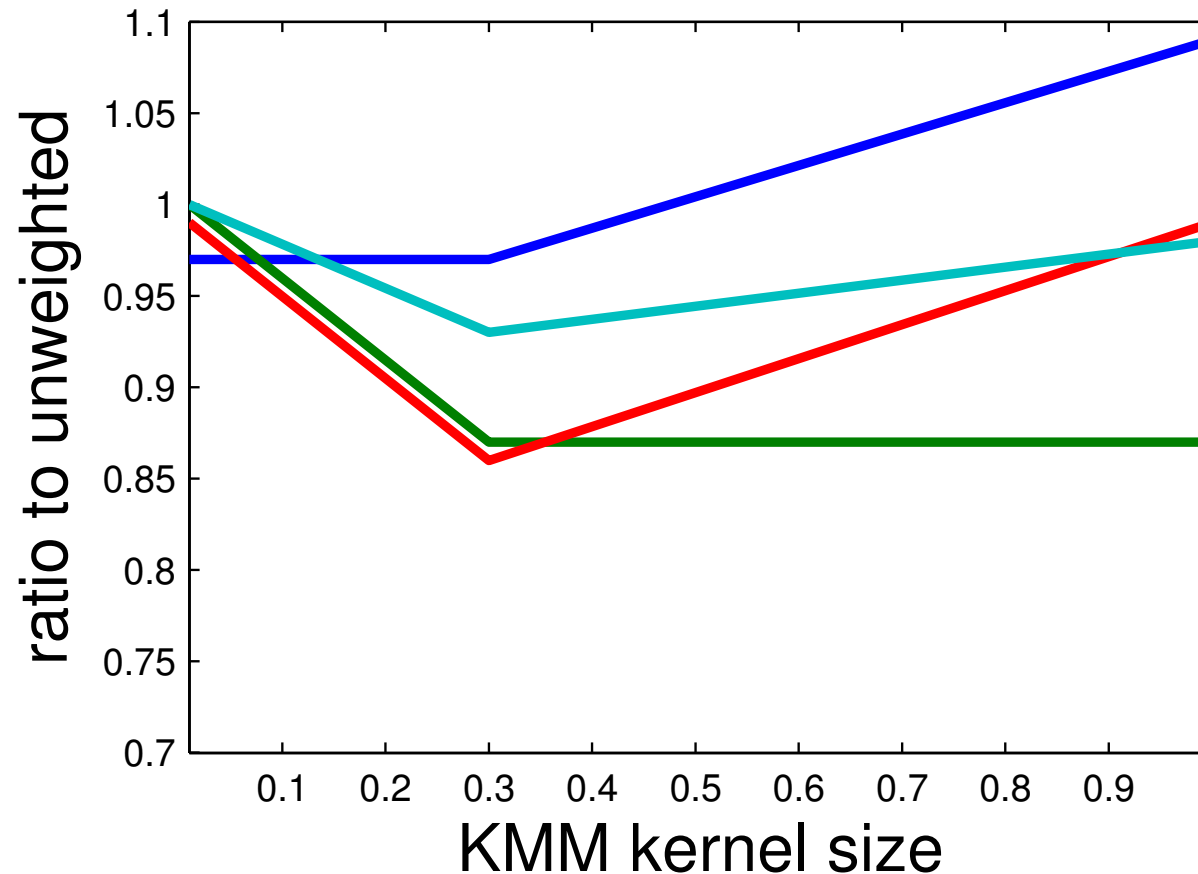
---



# Further work: model selection

---

- Model selection for covariate shift
- Results from [Sugiyama et al., 2008]
- Data have 18-21 dimensions



# Further work: model selection

---

- Model selection **for covariate shift**
- Some strategies [Bickel et al., 2009]
  - **Systematic drift**: can be learned [Bickel et al., 2009]
  - **Cross validation** to obtain error for current  $\beta$  estimate [Sugiyama et al., 2008, Kanamori et al., 2009]
  - **Classifier** of training vs test: again, **cross-validate** [Bickel et al., 2009]
  - **Supremum** of MMD over set of kernels? **(this NIPS)** [Sriperumbudur et al., 2010]
- Does knowing something about the **learning problem** help?

# Further work: model selection

---

- Model selection **for covariate shift**
- Some strategies [Bickel et al., 2009]
  - **Systematic drift**: can be learned [Bickel et al., 2009]
  - **Cross validation** to obtain error for current  $\beta$  estimate [Sugiyama et al., 2008, Kanamori et al., 2009]
  - **Classifier** of training vs test: again, **cross-validate** [Bickel et al., 2009]
  - **Supremum** of MMD over set of kernels? **(this NIPS)** [Sriperumbudur et al., 2010]
- Does knowing something about the **learning problem** help?
- Model selection **for weighted learning**: bias for unweighted? [Kanamori et al., 2009]

# Summary

---

- **Kernel mean matching**: perform **covariate shift**...
  - ...**without** density estimation
  - ...using only particular covariate **features**
  - ...on **structured domains**
- **Large** performance advantage for “**simple**” learning algorithms
- **Mixed** results for **powerful** learning algorithms
- **Model selection** remains an issue



# Acknowledgements

---

- Co-authors on KMM papers:
  - Karsten Borgwardt
  - Jiayuan Huang
  - Marcel Schmittful
  - Bernhard Schölkopf
  - Alex Smola
- Discussions
  - Paul von Büнау
  - Corinna Cortes
  - Klaus-Robert Müller
  - Masashi Sugiyama

Questions?



# Bibliography

---

## References

- S. Bickel, M. Brückner, and T. Scheffer. Discriminative learning under covariate shift. *JMLR*, 10:2137–2155, 2009.
- K. F. Cheng and C. K. Chu. Semiparametric density estimation under a two-sample density ratio model. *Bernoulli*, 10(4):583–604, 2004.
- C. Cortes, M. Mohri, M. Riley, and A. Rostamizadeh. Sample selection bias correction theory. In *ALT*, 2008.
- R. M. Dudley. *Real analysis and probability*. Cambridge University Press, Cambridge, UK, 2002.
- R. Fortet and E. Mourier. Convergence de la réparation empirique vers la réparation théorique. *Ann. Scient. École Norm. Sup.*, 70:266–285, 1953.
- K. Fukumizu, A. Gretton, X. Sun, and B. Schölkopf. Kernel measures of conditional dependence. In *Advances in Neural Information Processing Systems 20*, pages 489–496, Cambridge, MA, 2008. MIT Press.
- A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, and A. Smola. A kernel method for the two-sample-problem. In *Advances in Neural Information Processing Systems 19*, pages 513–520, Cambridge, MA, 2007. MIT Press.
- A. Gretton, A. Smola, J. Huang, M. Schmittfull, K. Borgwardt, and B. Schölkopf. Dataset shift in machine learning. In J. Quiñero-Candela, M. Sugiyama, A. Schwaighofer, and N. Lawrence, editors, *Covariate Shift and Local Learning by Distribution Matching*, pages 131–160, Cambridge, MA, 2008. MIT Press.
- J. Huang, A. Smola, A. Gretton, K. Borgwardt, and B. Schölkopf. Correcting sample selection bias by unlabeled data. In *Advances in Neural Information Processing Systems 19*, Cambridge, MA, 2007. MIT Press.
- T. Kanamori, S. Hido, , and M Sugiyama. A least-squares approach to direct importance estimation. *Journal of*

# Characteristic kernels

# Characteristic Kernels (1)

---

- **Characteristic:** MMD a **metric** (MMD = 0 iff **P** = **Q**) [NIPS07b, COLT08]

# Characteristic Kernels (1)

---

- **Characteristic:** MMD a **metric** (MMD = 0 iff **P** = **Q**) [NIPS07b, COLT08]
- **Translation invariant** kernels:  $k(x, y) = k(x - y)$

# Characteristic Kernels (1)

---

- **Characteristic**: MMD a **metric** (MMD = 0 iff **P** = **Q**) [NIPS07b, COLT08]
- **Translation invariant** kernels:  $k(x, y) = k(x - y)$
- **Bochner's theorem**:

$$k(x) = \int_{\mathbb{R}^d} e^{-ix^\top \omega} d\Lambda(\omega)$$

- $\Lambda$  finite non-negative Borel measure



# Characteristic Kernels (1)

---

- **Characteristic:** MMD a **metric** (MMD = 0 iff  $\mathbf{P} = \mathbf{Q}$ ) [NIPS07b, COLT08]
- **Translation invariant** kernels:  $k(x, y) = k(x - y)$

- **Bochner's theorem:**

$$k(x) = \int_{\mathbb{R}^d} e^{-ix^\top \omega} d\Lambda(\omega)$$

- $\Lambda$  finite non-negative Borel measure

- **Fourier representation of MMD:**

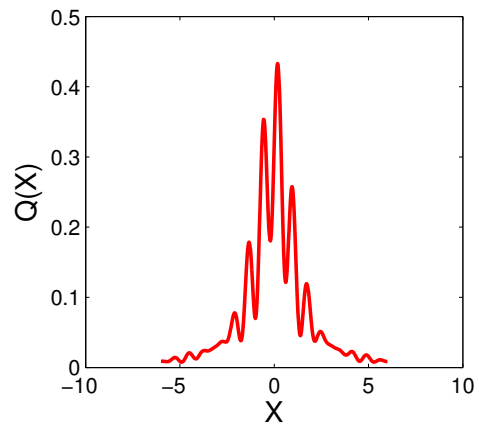
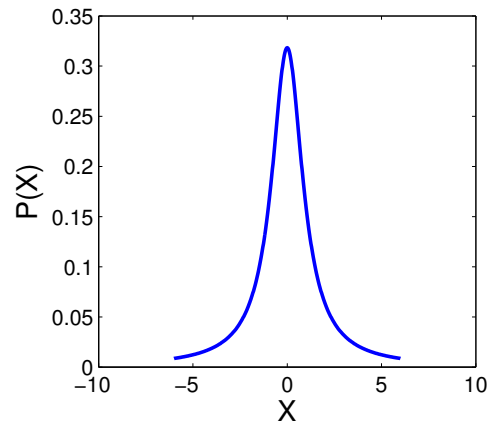
$$\text{MMD}(\mathbf{P}, \mathbf{Q}; F) := \left\| [(\bar{\phi}_{\mathbf{P}} - \bar{\phi}_{\mathbf{Q}}) \Lambda]^\vee \right\|_{\mathcal{F}}$$

- $\phi_{\mathbf{P}}$  characteristic function of  $\mathbf{P}$
- $f^\wedge$  is Fourier transform,  $f^\vee$  is inverse Fourier transform
- $\mu_x := \int k(\cdot, x) d\mathbf{P}(x)$

# Characteristic Kernels (2)

---

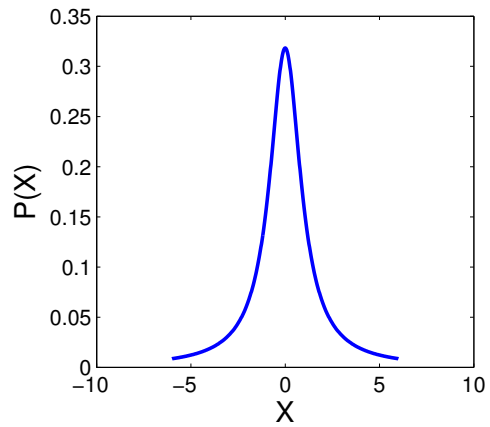
- Example: **P** differs from **Q** at (roughly) one frequency



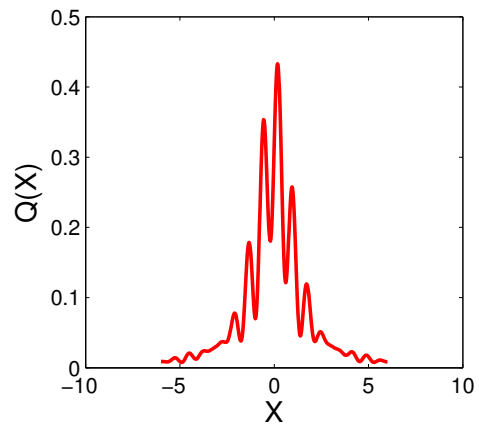
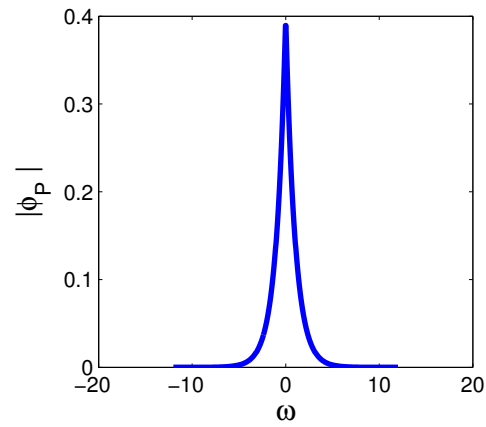
# Characteristic Kernels (2)

---

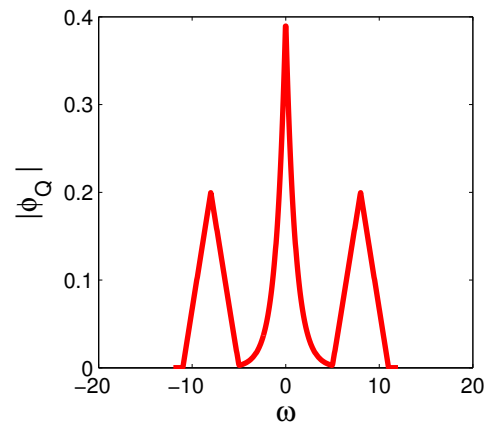
- Example: **P** differs from **Q** at (roughly) one frequency



$F \rightarrow$

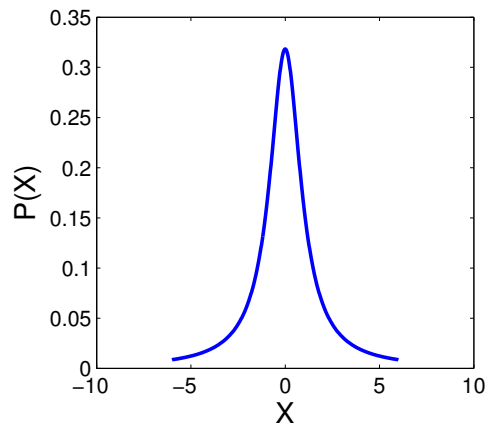


$F \rightarrow$

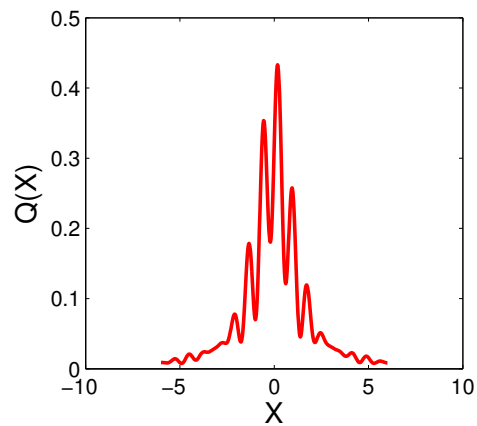
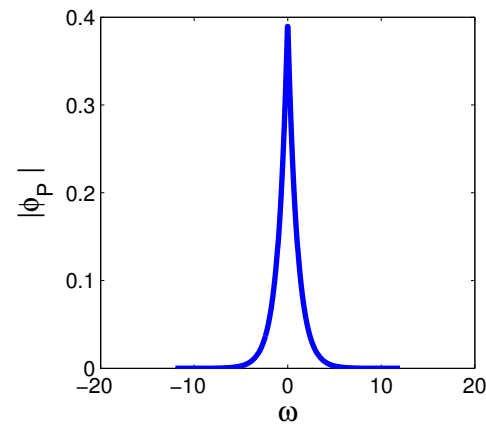


# Characteristic Kernels (2)

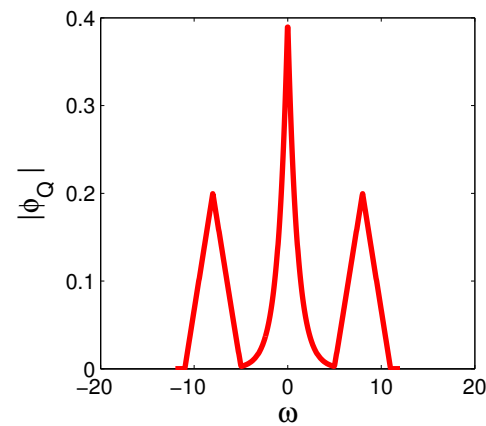
- Example: **P** differs from **Q** at (roughly) one frequency



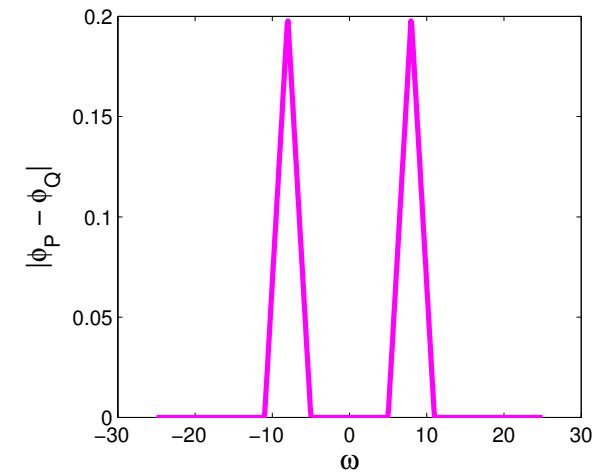
$F \rightarrow$



$F \rightarrow$



Characteristic function difference



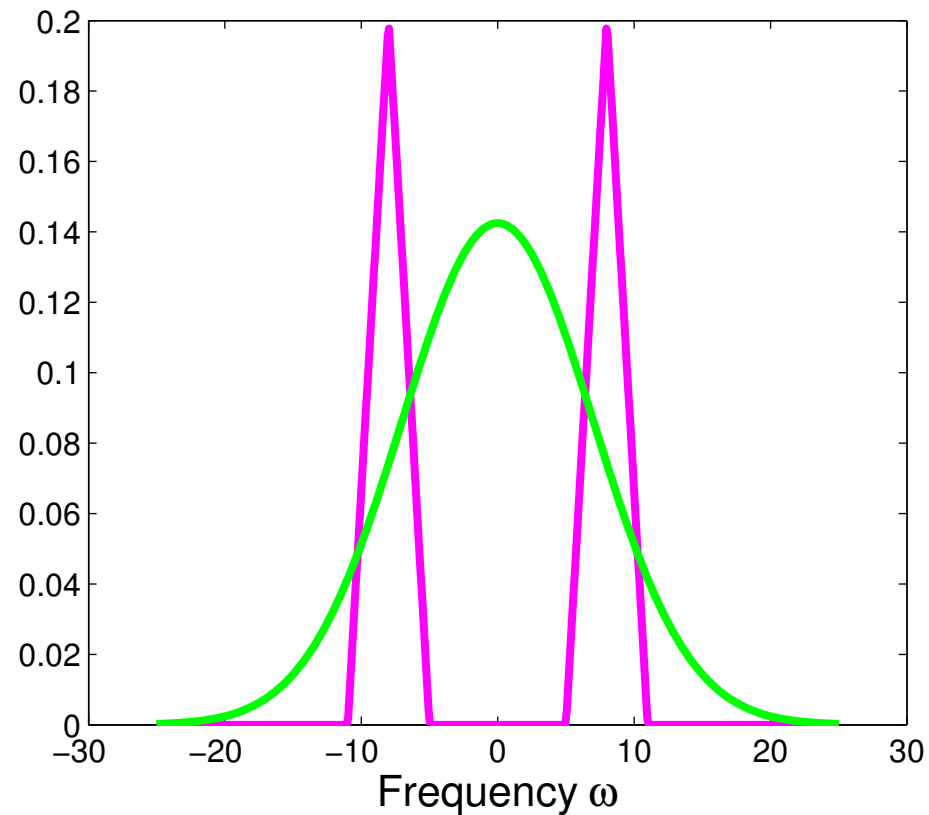
# Characteristic Kernels (3)

---

- Example: **P** differs from **Q** at (roughly) one frequency

Gaussian kernel

Difference  $|\phi_P - \phi_Q|$

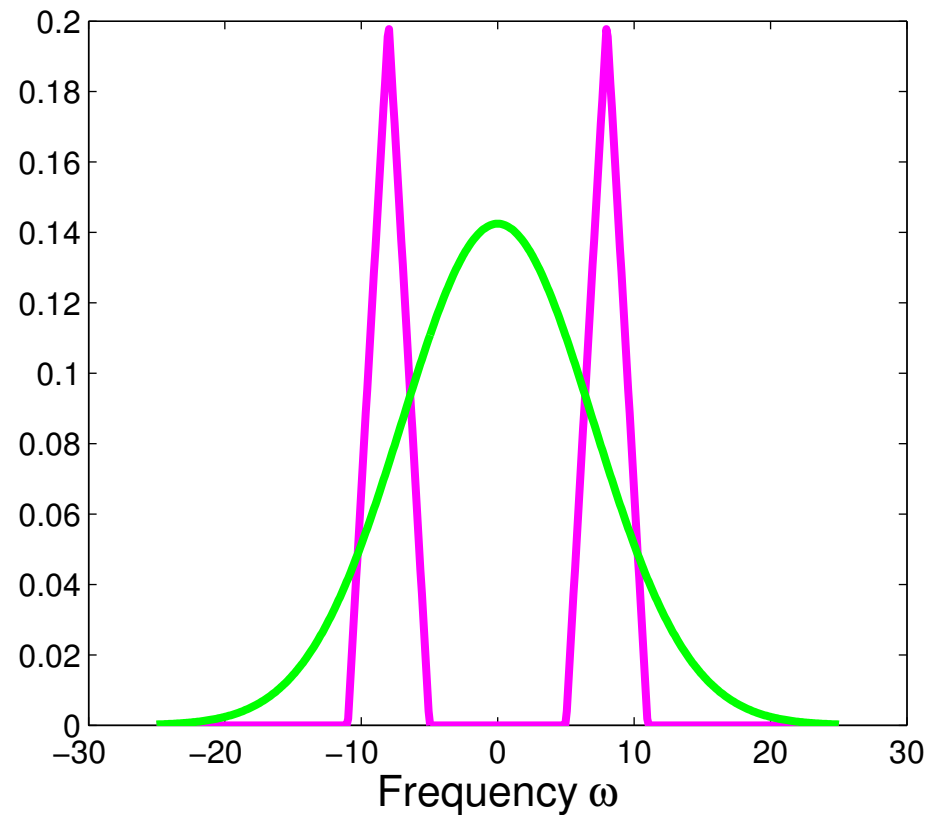


# Characteristic Kernels (3)

---

- Example: **P** differs from **Q** at (roughly) one frequency

## Characteristic



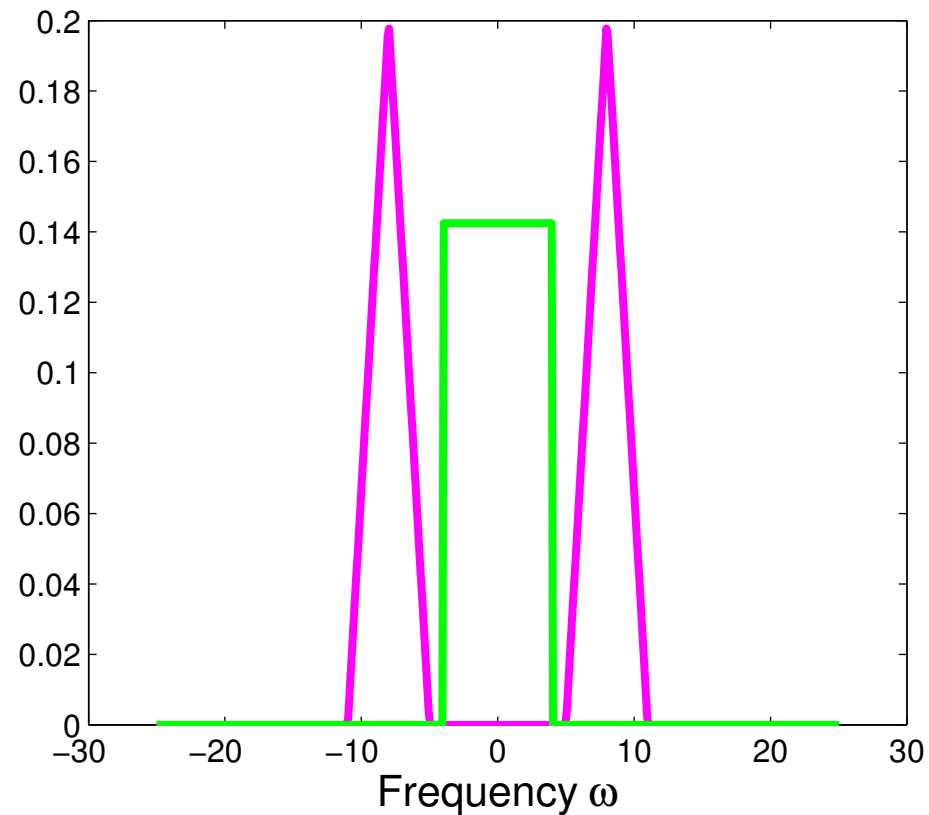
# Characteristic Kernels (4)

---

- Example: **P** differs from **Q** at (roughly) one frequency

Sinc kernel

Difference  $|\phi_P - \phi_Q|$

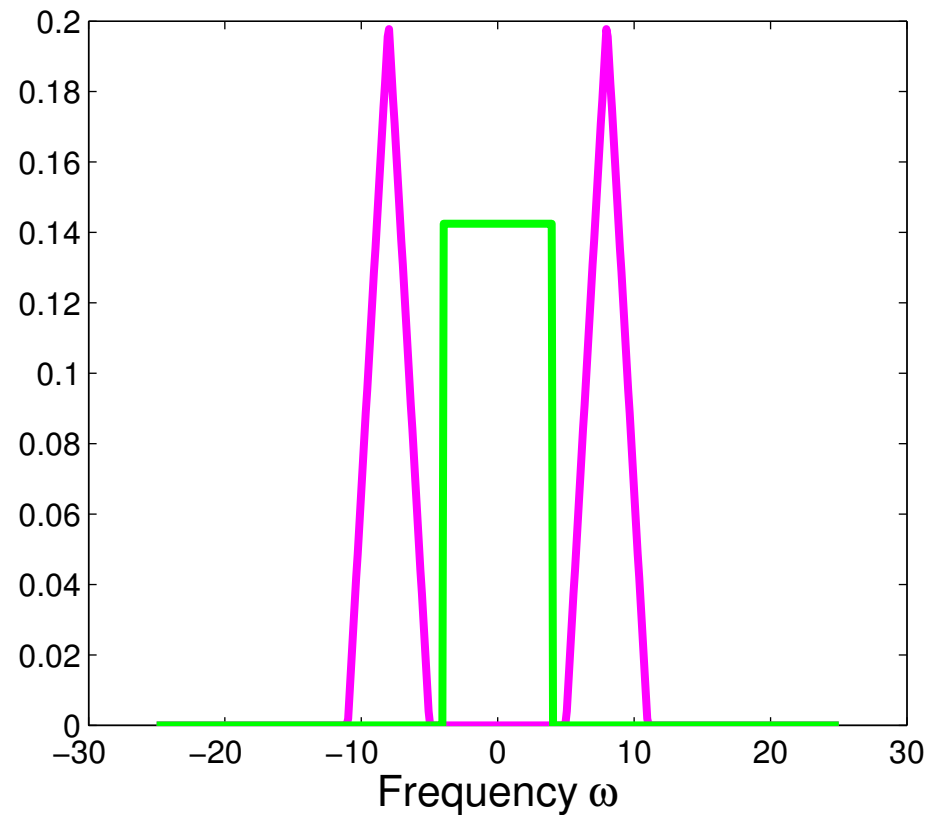


# Characteristic Kernels (4)

---

- Example: **P** differs from **Q** at (roughly) one frequency

NOT characteristic





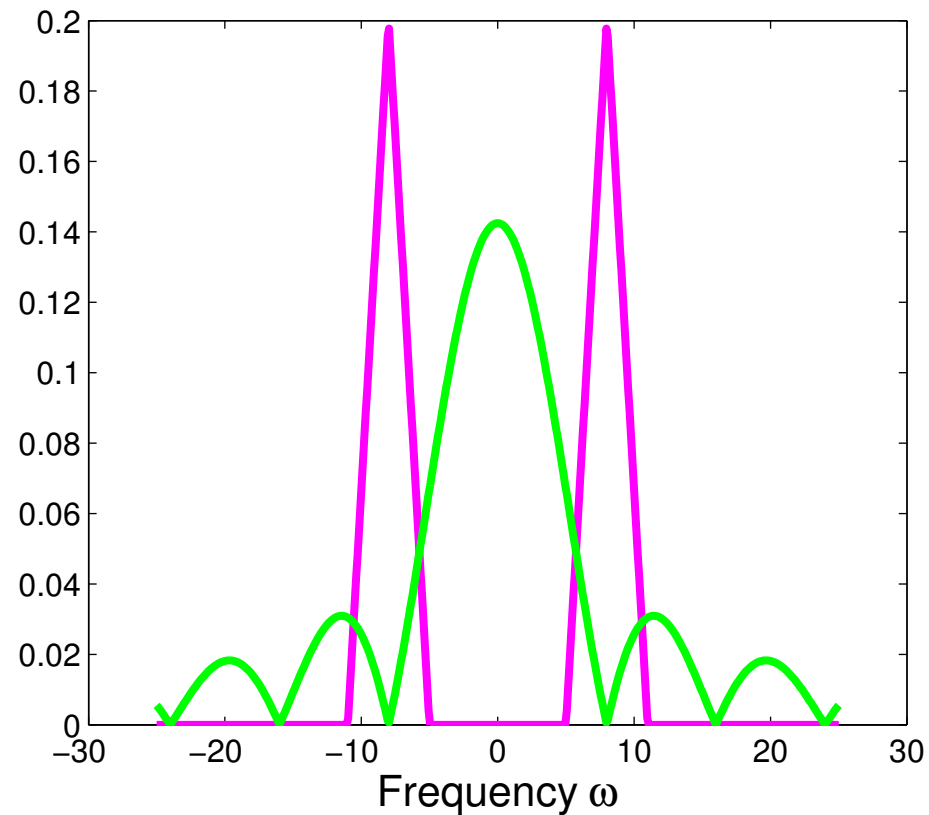
# Characteristic Kernels (5)

---

- Example: **P** differs from **Q** at (roughly) one frequency

**B-Spline kernel**

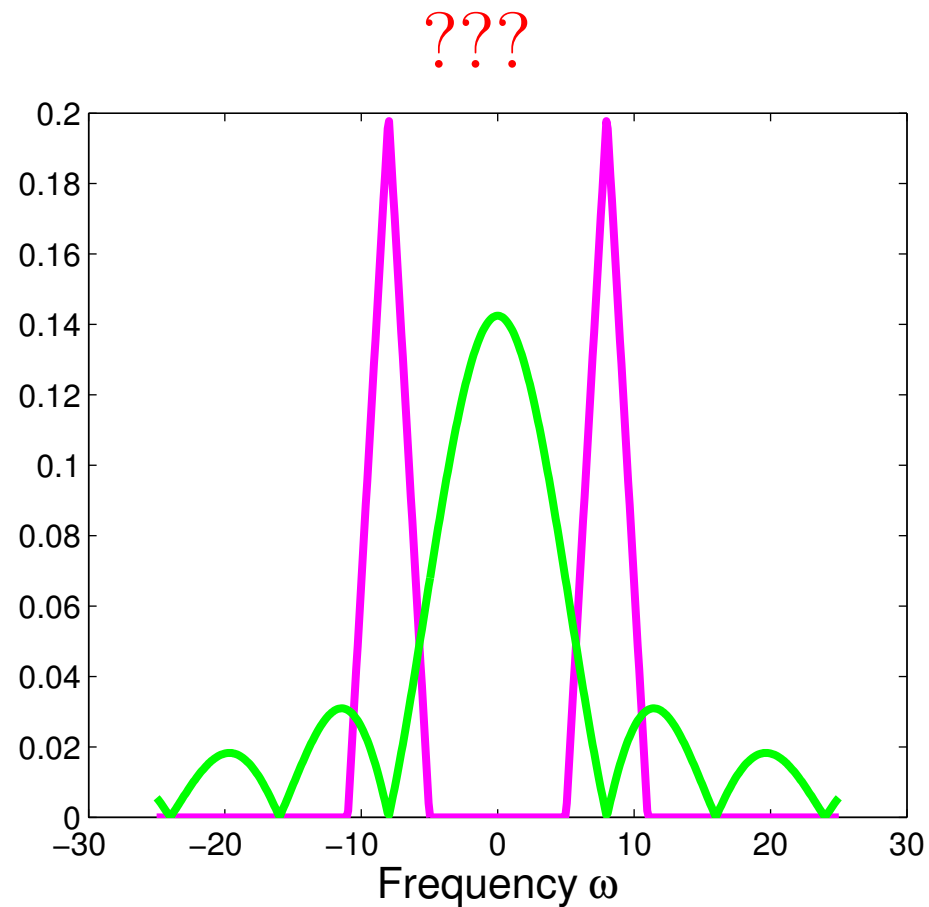
Difference  $|\phi_P - \phi_Q|$



# Characteristic Kernels (5)

---

- Example: **P** differs from **Q** at (roughly) one frequency

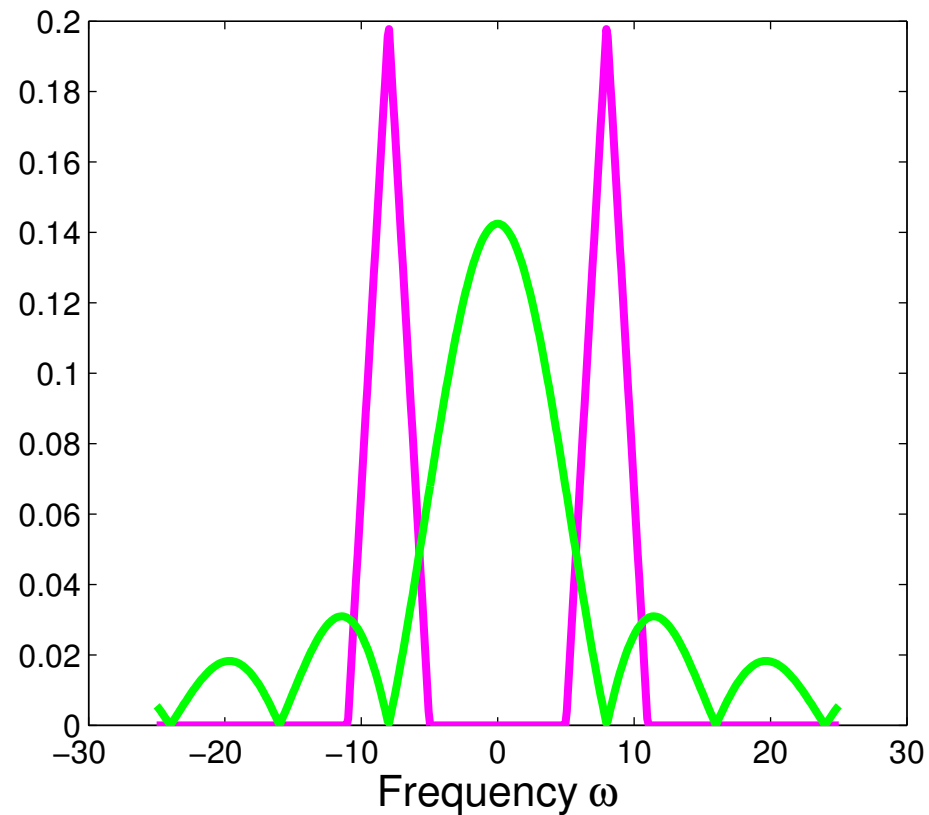


# Characteristic Kernels (5)

---

- Example: **P** differs from **Q** at (roughly) one frequency

## Characteristic



# Summary: Characteristic Kernels

---

- **Characteristic kernel:** (MMD = 0 iff  $\mathbf{P} = \mathbf{Q}$ ) [NIPS07b, COLT08]
- **Main theorem:**  $k$  characteristic if and only if  $\text{supp}(\Lambda) = \mathbb{R}^d$  [COLT08]

# Summary: Characteristic Kernels

---

- **Characteristic kernel:** (MMD = 0 iff  $\mathbf{P} = \mathbf{Q}$ ) [NIPS07b, COLT08]
- **Main theorem:**  $k$  characteristic if and only if  $\text{supp}(\Lambda) = \mathbb{R}^d$  [COLT08]
  - Corollary: continuous, compactly supported  $k$  characteristic

# Summary: Characteristic Kernels

---

- **Characteristic kernel:** (MMD = 0 iff  $\mathbf{P} = \mathbf{Q}$ ) [NIPS07b, COLT08]
- **Main theorem:**  $k$  characteristic if and only if  $\text{supp}(\Lambda) = \mathbb{R}^d$  [COLT08]
  - Corollary: continuous, compactly supported  $k$  characteristic
- **Alternative property:** continuous, strictly P.D., includes NON-translation invariant [COLT09?]

$$k(x, y) = e^{\sigma x^\top y}, \sigma > 0$$

# Summary: Characteristic Kernels

---

- **Characteristic kernel:** (MMD = 0 iff  $\mathbf{P} = \mathbf{Q}$ ) [NIPS07b, COLT08]
- **Main theorem:**  $k$  characteristic if and only if  $\text{supp}(\Lambda) = \mathbb{R}^d$  [COLT08]
  - Corollary: continuous, compactly supported  $k$  characteristic
- **Alternative property:** continuous, strictly P.D., includes NON-translation invariant [COLT09?]
- **Similar reasoning** wherever extensions of **Bochner's theorem** exist: [NIPS08a]
  - Locally compact Abelian groups (periodic domains)
  - Compact, non-Abelian groups (orthogonal matrices)
  - The semigroup  $\mathbb{R}_n^+$  (histograms)