

# Learning Kernels -Tutorial

Part III: Theoretical Guarantees.

Corinna Cortes

Google Research

[corinna@google.com](mailto:corinna@google.com)

Mehryar Mohri

Courant Institute &

Google Research

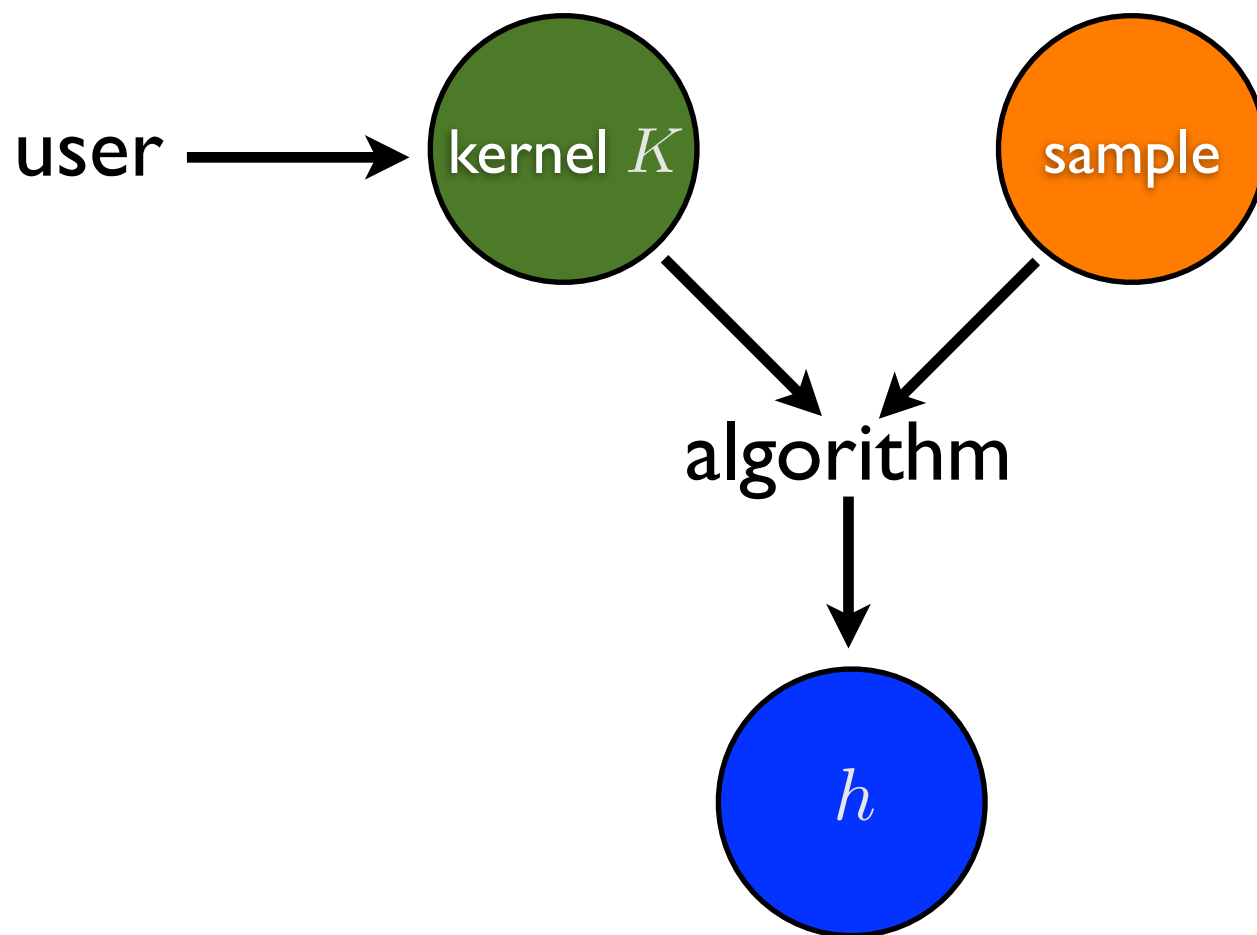
[mohri@cims.nyu.edu](mailto:mohri@cims.nyu.edu)

Afshin Rostami

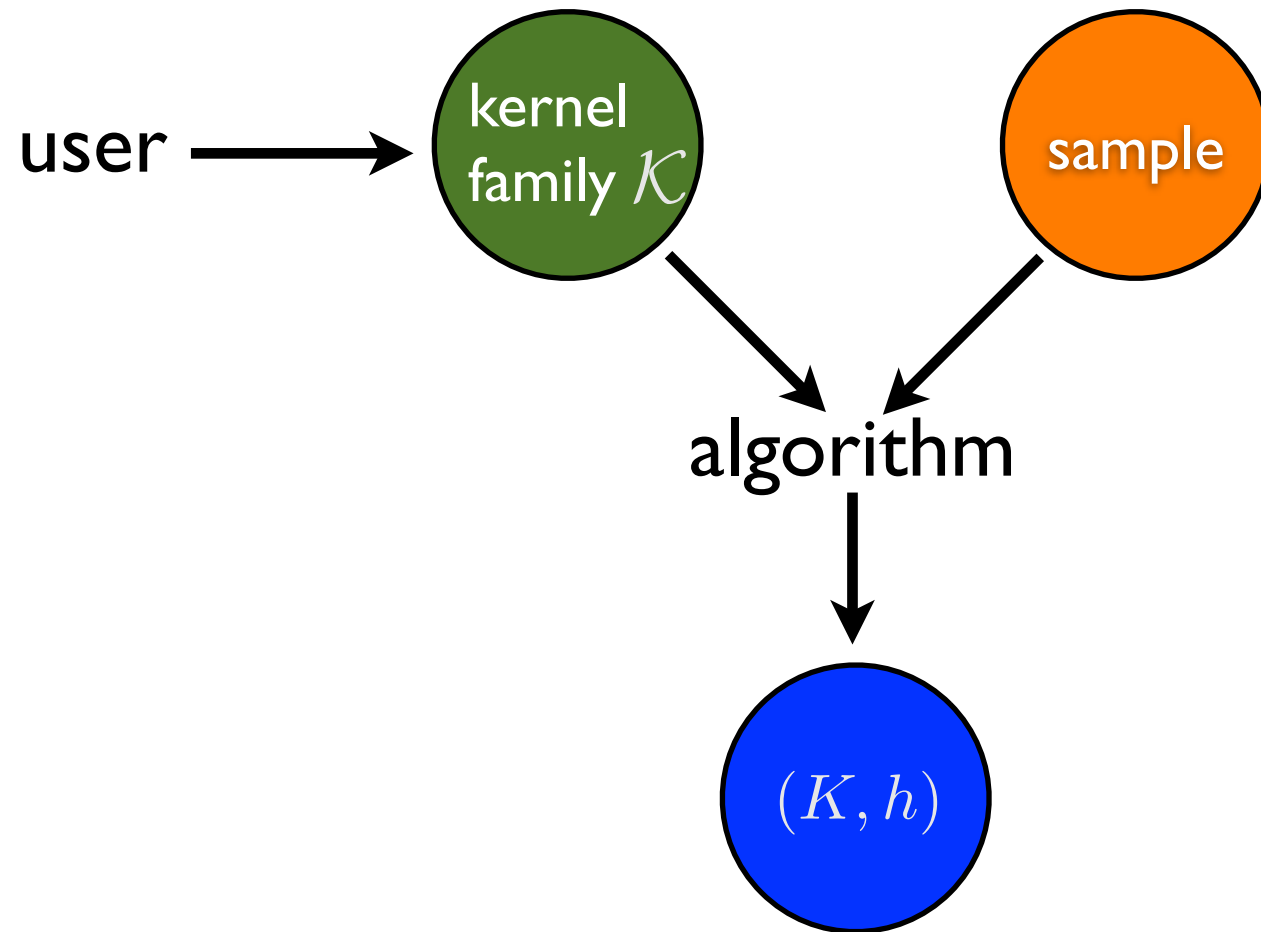
UC Berkeley

[arostami@eecs.  
berkeley.edu](mailto:arostami@eecs.berkeley.edu)

# Standard Learning with Kernels



# Learning Kernel Framework



# Learning Kernels

## ■ Theoretical questions:

- what is the price to pay for relaxing the requirement from the user to specify a kernel?
- how does the choice of the kernel family affect generalization?

# Part III

- Non-negative combinations.
- General case.

# Kernel Families

- Most frequently used kernel families,  $q \geq 1$ ,

$$\mathcal{K}_q = \left\{ \sum_{k=1}^p \mu_k K_k : \boldsymbol{\mu} \in \Delta_q \right\}$$

with  $\Delta_q = \left\{ \boldsymbol{\mu} : \mu_k \geq 0, \|\boldsymbol{\mu}\|_q = 1 \right\}.$

- Hypothesis sets:

$$H_q = \left\{ h \in \mathbb{H}_K : K \in \mathcal{K}_q, \|h\|_{\mathbb{H}_K} \leq 1 \right\}.$$

# Rademacher Complexity

- Empirical Rademacher complexity of  $H$ : for a sample  $S = (x_1, \dots, x_m)$ ,

$$\hat{\mathfrak{R}}_S(H) = \mathbb{E}_{\sigma} \left[ \sup_{h \in H} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right],$$

where  $\sigma_i$ s are independent uniform random variables taking values in  $\{-1, +1\}$ .

- Rademacher complexity of  $H$ :

$$\mathfrak{R}_m(H) = \mathbb{E}_{S \sim D^m} [\hat{\mathfrak{R}}_S(H)].$$

# Single Kernel Margin Bound

- **Theorem** (Koltchinskii and Panchenko, 2002): fix  $\rho > 0$ . Assume that  $K(x, x) \leq R^2$  for all  $x$ , then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , for any  $h \in H_1$ ,

$$R(h) \leq \hat{R}_\rho(h) + 2\sqrt{\frac{R^2/\rho^2}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}.$$



# Early Learning Kernel Bounds

(Bousquet and Herrmann 2003; Lanckriet et al., 2004)

- For any  $\delta > 0$ , with probability at least  $1 - \delta$ , for any  $h \in H_1$ ,

$$R(h) \leq \hat{R}_\rho(h) + \frac{1}{\sqrt{m}} \left[ \sqrt{\frac{\max_{k=1}^p \text{Tr}(\mathbf{K}_k) \max_{k=1}^p \frac{\|\mathbf{K}_k\|}{\text{Tr}(\mathbf{K}_k)}}{\rho^2}} + 4 + \sqrt{2 \log \frac{1}{\delta}} \right].$$

- but, bound always greater than one (Srebro and Ben-David, 2006)!
- other bound of (Lanckriet et al., 2004) for linear combination case also always greater than one!

# Multiplicative Learning Bound

(Lanckriet et al., 2004)

- Assume that for all  $k \in [1, p]$ ,  $K_k(x, x) \leq R^2$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , for any  $h \in H_1$ ,

$$R(h) \leq \hat{R}_\rho(h) + O\left(\sqrt{\frac{p R^2 / \rho^2}{m}}\right).$$

- bound multiplicative in  $p$  (number of kernels).

# Additive Learning Bound

(Srebro and Ben-David, 2006)

- Assume that for all  $k \in [1, p]$ ,  $K_k(x, x) \leq R^2$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , for any  $h \in H_1$ ,

$$R(h) \leq \hat{R}_\rho(h) + \sqrt{8 \frac{2 + p \log \frac{128em^3 R^2}{\rho^2 p} + 256 \frac{R^2}{\rho^2} \log \frac{\rho em}{8R} \log \frac{128mR^2}{\rho^2} + \log(1/\delta)}{m}}.$$

- bound additive in  $p$  (modulo log terms).
- not informative for  $p > m$ .
- based on pseudo-dimension of kernel family.
- similar guarantees for other families.

# New Data-Dependent Bound

(CC, MM, and AR, 2010)

■ **Theorem:** for any sample  $S$  of size  $m$ , and positive integer  $r$ ,

$$\hat{\mathfrak{R}}_S(H_1) \leq \frac{\sqrt{\frac{23}{22} r \|\mathbf{u}\|_r}}{m},$$

with  $\mathbf{u} = (\text{Tr}[\mathbf{K}_1], \dots, \text{Tr}[\mathbf{K}_p])^\top$ .

- similarity with single kernel bound.
- can be used directly to derive an algorithm.

# New Data-Dependent Bound

■ **Proof:** Let  $q, r \geq 1$  with  $\frac{1}{q} + \frac{1}{r} = 1$ .

$$\begin{aligned}
 \hat{\mathfrak{R}}_S(H_q) &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{h \in H_q} \sum_{i=1}^m \sigma_i h(x_i) \right] \\
 &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\boldsymbol{\mu} \in \Delta_q, \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \leq 1} \sum_{i,j=1}^m \sigma_i \alpha_j K_{\boldsymbol{\mu}}(x_i, x_j) \right] \\
 &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\boldsymbol{\mu} \in \Delta_q, \boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha} \leq 1} \boldsymbol{\sigma}^\top \mathbf{K}_{\boldsymbol{\mu}} \boldsymbol{\alpha} \right] = \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\boldsymbol{\mu} \in \Delta_q, \|\boldsymbol{\alpha}\|_{\mathbf{K}^{1/2}} \leq 1} \langle \boldsymbol{\sigma}, \boldsymbol{\alpha} \rangle_{\mathbf{K}^{1/2}} \right] \\
 &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\boldsymbol{\mu} \in \Delta_q} \sqrt{\boldsymbol{\sigma}^\top \mathbf{K}_{\boldsymbol{\mu}} \boldsymbol{\sigma}} \right] \quad \text{(Cauchy-Schwarz)} \\
 &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\boldsymbol{\mu} \in \Delta_q} \sqrt{\boldsymbol{\mu} \cdot \mathbf{u}_{\boldsymbol{\sigma}}} \right] \quad [\mathbf{u}_{\boldsymbol{\sigma}} = (\boldsymbol{\sigma}^\top \mathbf{K}_1 \boldsymbol{\sigma}, \dots, \boldsymbol{\sigma}^\top \mathbf{K}_p \boldsymbol{\sigma})^\top] \\
 &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sqrt{\|\mathbf{u}_{\boldsymbol{\sigma}}\|_r} \right]. \quad \text{(definition of dual norm)}
 \end{aligned}$$

# New Data-Dependent Bound

■ **Proof:** in the following,  $r \geq 1$  is arbitrary integer.

$$\begin{aligned}\widehat{\mathfrak{R}}_S(H_1) &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} [\sqrt{\|\mathbf{u}_{\boldsymbol{\sigma}}\|_{\infty}}] \\ &\leq \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} [\sqrt{\|\mathbf{u}_{\boldsymbol{\sigma}}\|_r}] & (\forall r \geq 1, \|\mathbf{u}_{\boldsymbol{\sigma}}\|_{\infty} \leq \|\mathbf{u}_{\boldsymbol{\sigma}}\|_r) \\ &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \left[ \sum_{k=1}^p (\boldsymbol{\sigma}^{\top} \mathbf{K}_k \boldsymbol{\sigma})^r \right]^{\frac{1}{2r}} \right] \\ &\leq \frac{1}{m} \left[ \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sum_{k=1}^p (\boldsymbol{\sigma}^{\top} \mathbf{K}_k \boldsymbol{\sigma})^r \right] \right]^{\frac{1}{2r}} \quad (\text{Jensen's inequality}) \\ &= \frac{1}{m} \left[ \sum_{k=1}^p \mathbb{E}_{\boldsymbol{\sigma}} \left[ (\boldsymbol{\sigma}^{\top} \mathbf{K}_k \boldsymbol{\sigma})^r \right] \right]^{\frac{1}{2r}} \\ &\leq \frac{1}{m} \left[ \sum_{k=1}^p \left( \frac{23}{22} r \operatorname{Tr}[\mathbf{K}_k] \right)^r \right]^{\frac{1}{2r}} = \frac{\sqrt{\frac{23}{22} r \|\mathbf{u}\|_r}}{m}. & (\text{lemma})\end{aligned}$$

# Key Lemma

- **Lemma:** Let  $\mathbf{K}$  be a kernel matrix for a finite sample. Then, for any integer  $r$ ,

$$\mathbb{E}_{\boldsymbol{\sigma}} \left[ (\boldsymbol{\sigma}^\top \mathbf{K} \boldsymbol{\sigma})^r \right] \leq \left( \frac{23}{22} r \operatorname{Tr}[\mathbf{K}] \right)^r.$$

- proof based on combinatorial argument.

# New Learning Bound - LI

(CC, MM, and AR, 2010)

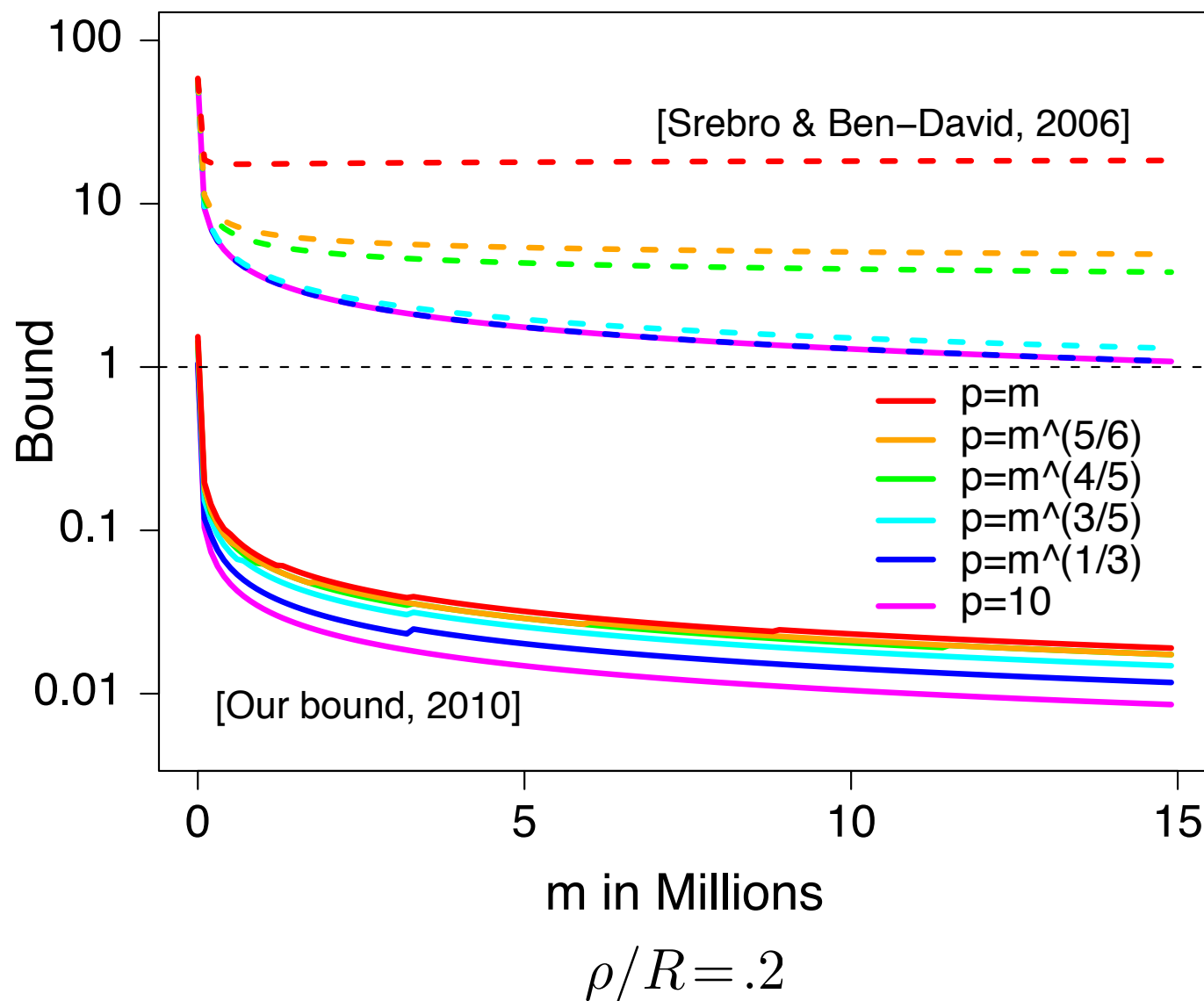
- **Theorem:** assume that for all  $k \in [1, p]$ ,  $K_k(x, x) \leq R^2$ .  
Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , for any  $h \in H_1$ ,

$$R(h) \leq \hat{R}_\rho(h) + 2\sqrt{\frac{\frac{23}{22}e \lceil \log p \rceil R^2 / \rho^2}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}.$$

- very weak dependency on  $p$ , no extra log terms.
- analysis based on Rademacher complexity.
- bound valid for  $p \gg m$ .
- see also (Kakade et al., 2010).



# Comparison



# Lower Bound

## ■ Tight bound:

- dependency  $\sqrt{\log p}$  cannot be improved.
- argument based on VC dimension or example.

## ■ Observations: case $\mathcal{X} = \{-1, +1\}^p$ .

- canonical projection kernels  $K_k(\mathbf{x}, \mathbf{x}') = x_k x'_k$  .
- $H_1$  contains  $J_p = \{\mathbf{x} \mapsto s x_k : k \in [1, p], s \in \{-1, +1\}\}$  .
- $\text{VCdim}(J_p) = \Omega(\log p)$  .
- for  $\rho = 1$  and  $h \in J_p$ ,  $\hat{R}_\rho(h) = \hat{R}(h)$  .
- VC lower bound:  $\Omega(\sqrt{\text{VCdim}(J^p)/m})$ .

# New Paper

- **Recent claim** (Hussain and Shawe-Taylor, AISTATS 2011): additive bound in terms of  $\log p$ , instead of multiplicative.
- main proof incorrect: probabilistic bound on Rademacher complexity, but slack term left out of proof of theorem 8. Adding it → **multiplicative bound**.
- however: authors are preparing new version (private communication: J. Shawe-Taylor).

# New Learning Bound - Lq

(CC, MM, and AR, 2010)

- **Theorem:** let  $q, r \geq 1$  with  $\frac{1}{q} + \frac{1}{r} = 1$  and  $r$  integer. Assume that for all  $k \in [1, p]$ ,  $K_k(x, x) \leq R^2$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , for any  $h \in H_q$ ,

$$R(h) \leq \hat{R}_\rho(h) + 2p^{\frac{1}{2r}} \sqrt{\frac{\frac{23}{22} r R^2 / \rho^2}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}.$$

- mild dependency on  $p$ .
- analysis based on Rademacher complexity.

# Lower Bound

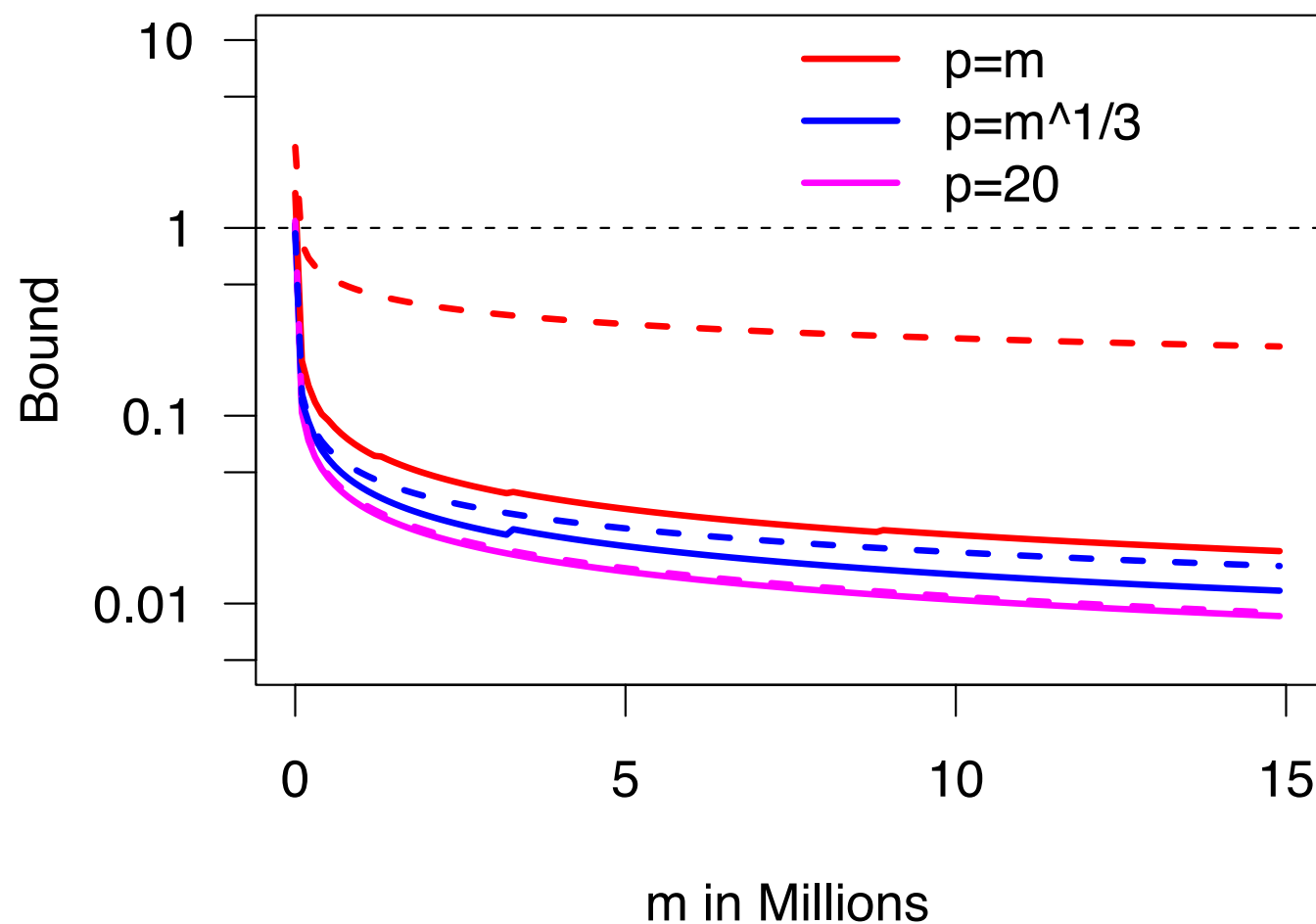
## ■ Tight bound:

- dependency  $p^{\frac{1}{2r}}$  cannot be improved.
- in particular  $p^{\frac{1}{4}}$  tight for  $L_2$  regularization.

## ■ Observations: equal kernels.

- $\sum_{k=1}^p \mu_k K_k = \left( \sum_{k=1}^p \mu_k \right) K_1$  .
- thus,  $\|h\|_{\mathbb{H}_{K_1}}^2 = \left( \sum_{k=1}^p \mu_k \right) \|h\|_{\mathbb{H}_K}^2$  for  $\sum_{k=1}^p \mu_k \neq 0$  .
- $\sum_{k=1}^p \mu_k \leq p^{\frac{1}{r}} \|\boldsymbol{\mu}\|_q = p^{\frac{1}{r}}$  (Hölder's inequality).
- $H_q$  coincides with  $\{h \in \mathbb{H}_{K_1} : \|h\|_{\mathbb{H}_{K_1}} \leq p^{\frac{1}{2r}}\}$  .

# Comparison L1 vs L2



# Conclusion

- **Theory:** tight generalization bounds for learning kernels with  $L_1$  or  $L_q$  regularization ( $p$  dependency).
  - mild dependency on  $p$ .
  - similar proof and analysis for other regularizations.
- **Applications:**
  - results suggest using large number of kernels.
  - recent results show significant improvements (CC, MM, AR, ICML 2010).

# Part III

- Non-negative combinations.
- General case.



# Kernel Family

## ■ General case:

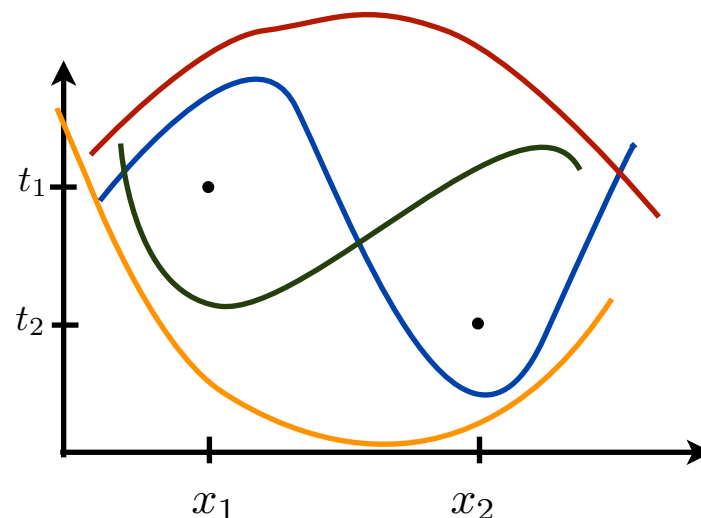
- $\mathcal{K}$  a family of kernels bounded by  $R$ .
- finite pseudo-dimension:  $\text{Pdim}(\mathcal{K}) < \infty$ .
- general hypothesis set:

$$H_{\mathcal{K}} = \left\{ h \in \mathbb{H}_K : K \in \mathcal{K}, \|h\|_{\mathbb{H}_K} \leq 1 \right\}.$$

# Shattering

■ **Definition:** Let  $H$  be a hypothesis set of functions from  $X$  to  $\mathbb{R}$ .  $A = \{x_1, \dots, x_m\}$  is **shattered** by  $H$  if there exist  $t_1, \dots, t_m \in \mathbb{R}$  such that

$$\left| \left\{ \begin{bmatrix} \text{sgn} (L(h(x_1), f(x_1)) - t_1) \\ \vdots \\ \text{sgn} (L(h(x_m), f(x_m)) - t_m) \end{bmatrix} : h \in H \right\} \right| = 2^m.$$



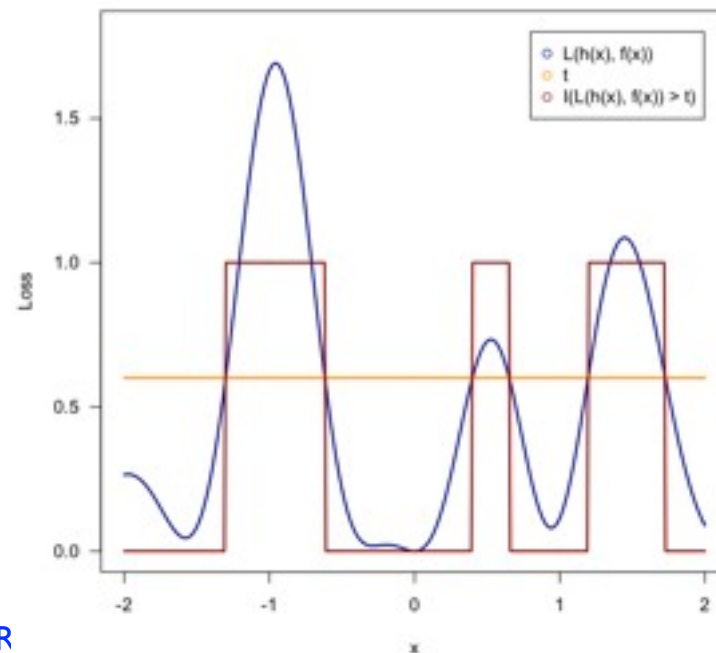
# Pseudo-Dimension

(Pollard, 1984)

■ **Definition:** Let  $H$  be a hypothesis set of functions from  $X$  to  $\mathbb{R}$ . The pseudo-dimension of  $H$ ,  $\text{Pdim}(H)$ , is the size of the largest set shattered by  $H$ .

■ **Definition** (equivalent, see also (Vapnik, 1995)):

$$\text{Pdim}(H) = \text{VCdim}\left(\left\{(x, t) \mapsto 1_{(h(x) - t) > 0} : h \in H\right\}\right).$$



# Pseudo-Dimension - Properties

- **Theorem:** Pseudo-dimension of hyperplanes.

$$\text{Pdim}(\mathbf{x} \mapsto \mathbf{w} \cdot \mathbf{x} + b : \mathbf{w} \in \mathbb{R}^N, b \in \mathbb{R}) = N + 1.$$

- **Theorem:** Pseudo-dimension of a vector space of real-valued functions  $H$ .

$$\text{Pdim}(H) = \dim(H).$$

- **Theorem:** Pseudo-dimension of  $\phi(H) = \{\phi \circ h : h \in H\}$  where  $\phi$  is a monotone function:

$$\text{Pdim}(\phi(H)) \leq \dim(H).$$

# General Pdim Learning Bound

(Srebro and Ben-David, 2006)

- Let  $\mathcal{K}$  a family of kernel functions bounded by  $R$ .  
Let  $d = \text{Pdim}(\mathcal{K})$ , then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , for any  $h \in H_{\mathcal{K}}$ ,

$$R(h) \leq \hat{R}_{\rho}(h) + \sqrt{8 \frac{2 + d \log \frac{128em^3 R^2}{\rho^2 d} + 256 \frac{R^2}{\rho^2} \log \frac{\rho em}{8R} \log \frac{128mR^2}{\rho^2} + \log(1/\delta)}{m}}.$$

- bound additive in  $d$  (modulo log terms).
- not informative for  $d > m$ .

# Application: Linear Combinations

- Linear and non-negative combination of base kernels (previous section):

$$\mathcal{K}_{\text{lin}} = \left\{ K_{\boldsymbol{\mu}} = \sum_{k=1}^p \mu_k K_k : \left( \sum_{k=1}^p \mu_k = 1 \right) \wedge (\mathbf{K}_{\boldsymbol{\mu}} \succeq \mathbf{0}) \right\}$$

$$\mathcal{K}_1 = \left\{ K_{\boldsymbol{\mu}} = \sum_{k=1}^p \mu_k K_k : \left( \sum_{k=1}^p \mu_k = 1 \right) \wedge (\boldsymbol{\mu} \geq \mathbf{0}) \right\}$$

- Since  $\mathcal{K}_{\text{lin}} \subseteq \mathcal{K}_1 \subseteq \left\{ \sum_{k=1}^p \mu_k K_k \right\}$ ,

$$\text{Pdim}(\mathcal{K}_1) \leq \text{Pdim}(\mathcal{K}_{\text{lin}}) \leq \dim \left( \left\{ \sum_{k=1}^p \mu_k K_k \right\} \right) = p.$$

# Application: Gaussian Kernels

- Gaussian kernels with a fixed covariance matrix:

$$\mathcal{K}_{\text{Gaussian}} = \left\{ (\mathbf{x}_1, \mathbf{x}_2) \mapsto \exp(-(\mathbf{x}_2 - \mathbf{x}_1)^\top \mathbf{A} (\mathbf{x}_2 - \mathbf{x}_1)) : \mathbf{A} \in \mathbb{S}_+^N \right\}.$$

- since  $\exp$  is monotone and since

$$\begin{aligned} & \left\{ (\mathbf{x}_1, \mathbf{x}_2) \mapsto (\mathbf{x}_2 - \mathbf{x}_1)^\top \mathbf{A} (\mathbf{x}_2 - \mathbf{x}_1) : \mathbf{A} \in \mathbb{S}_+^N \right\} \\ &= \left\{ (\mathbf{x}_1, \mathbf{x}_2) \mapsto \sum_{i,j=1}^n \mathbf{A}_{ij} (\mathbf{x}_2 - \mathbf{x}_1)_i (\mathbf{x}_2 - \mathbf{x}_1)_j : \mathbf{A} \in \mathbb{S}_+^N \right\} \\ &\subseteq \text{span} \left\{ (\mathbf{x}_1, \mathbf{x}_2) \mapsto (\mathbf{x}_2 - \mathbf{x}_1)_i (\mathbf{x}_2 - \mathbf{x}_1)_j : 1 \leq i \leq j \leq N \right\}, \end{aligned}$$

- $\text{Pdim}(\mathcal{K}_{\text{Gaussian}}) \leq \frac{N(N-1)}{2}.$
- Similar for  $\mathbf{A}$  diagonal,  $\text{Pdim}(\mathcal{K}_{\text{Gaussian}}) \leq N.$

# References

- Bousquet, Olivier and Herrmann, Daniel J. L. On the complexity of learning the kernel matrix. In NIPS, 2002.
- Cortes, Corinna, Mohri, Mehryar, and Rostamizadeh, Afshin. Generalization Bounds for Learning Kernels. In ICML, 2010.
- Cortes, Corinna, Mohri, Mehryar, and Rostamizadeh, Afshin. Two-stage learning kernel methods. In ICML, 2010.
- Zakria Hussain, John Shawe-Taylor. Improved Loss Bounds For Multiple Kernel Learning. In AISTATS, 2011.
- Kakade, Sham M., Shalev-Shwartz, Shai, and Tewari, Ambuj. Applications of strong convexity–strong smoothness duality to learning with matrices, 2010. arXiv:0910.0610v1.
- Koltchinskii, V. and Panchenko, D. Empirical margin distributions and bounding the generalization error of combined classifiers. Annals of Statistics, 30, 2002.
- Koltchinskii, Vladimir and Yuan, Ming. Sparse recovery in large ensembles of kernel machines on-line learning and bandits. In COLT, 2008.



# References

- Lanckriet, Gert, Cristianini, Nello, Bartlett, Peter, Ghaoui, Laurent El, and Jordan, Michael. Learning the kernel matrix with semidefinite programming. JMLR, 5, 2004.
- Srebro, Nathan and Ben-David, Shai. Learning bounds for support vector machines with learned kernels. In COLT, 2006.
- Ying, Yiming and Campbell, Colin. Generalization bounds for learning the kernel problem. In COLT, 2009.