

# Distributed Adaptive Sampling for Kernel Matrix Approximation

Daniele Calandriello, Alessandro Lazaric, Michal Valko



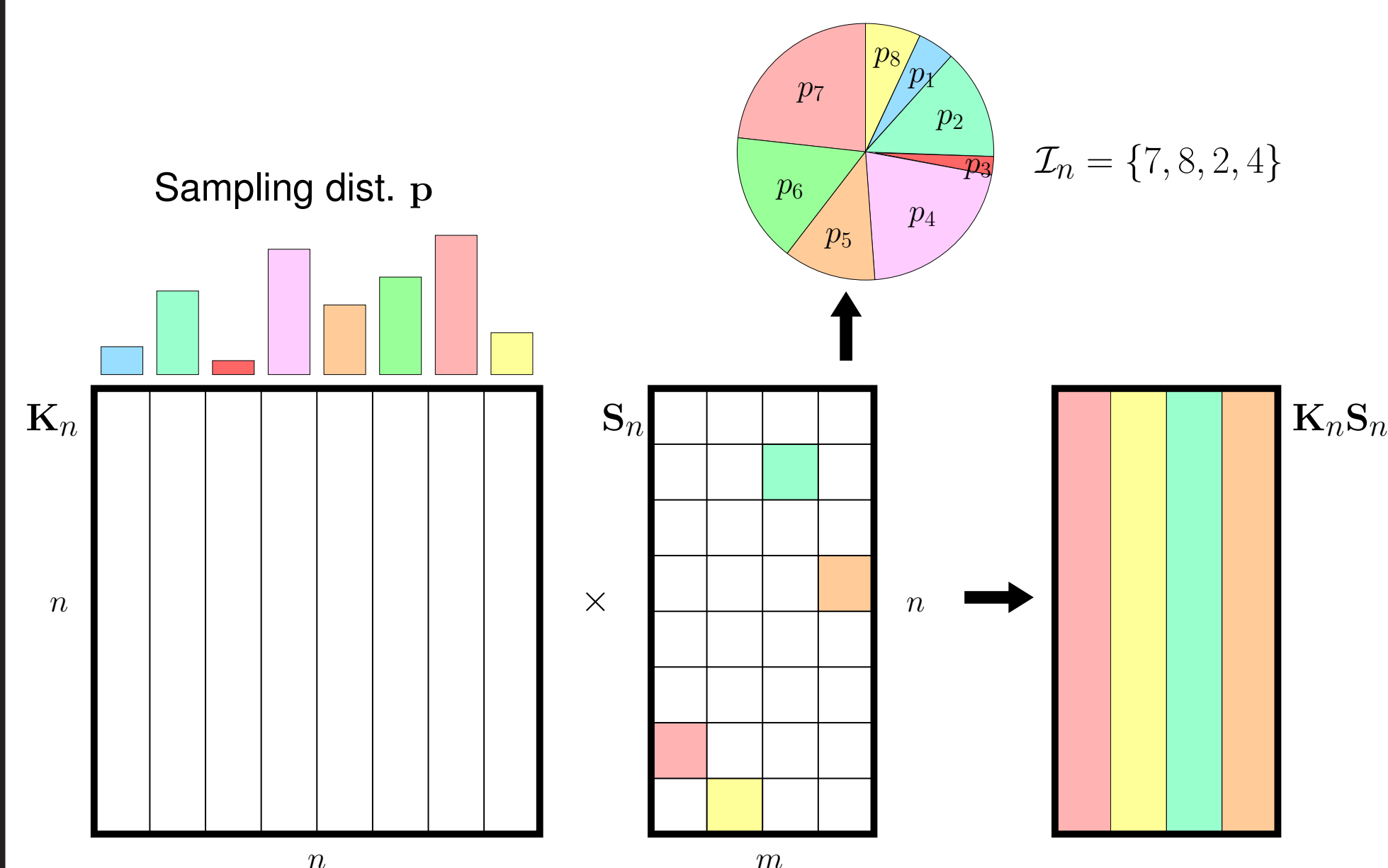
## Motivation

- Kernel methods are *versatile* and *accurate*
- Strong generalization guarantees but *poor scalability*
- $\mathcal{O}(n^3)$  time  $\mathcal{O}(n^2)$  space ( $n$  number of samples)
- Current limitation: Many approximate schemes are either *not scalable* or *not accurate*
- ⇒ We propose a **parallel distributed incremental approximation** scheme for kernel methods with **complexity and error guarantees adaptive to the dataset and kernel structure**
- ⇒ Runs in a **single pass** over the dataset, and can **update its solution** if **new data** arrives
- ⇒ On a single machine, computes a solution in **subquadratic**  $\mathcal{O}(n)$  time, avoids constructing the whole  $\mathbf{K}_n$
- ⇒ On multiple machines, computes a solution in **logarithmic**  $\mathcal{O}(\log(n))$  time, without increase in total work
- ⇒ **Black-box applicability** to many downstream tasks

## Nyström Approximation

### Subsampling

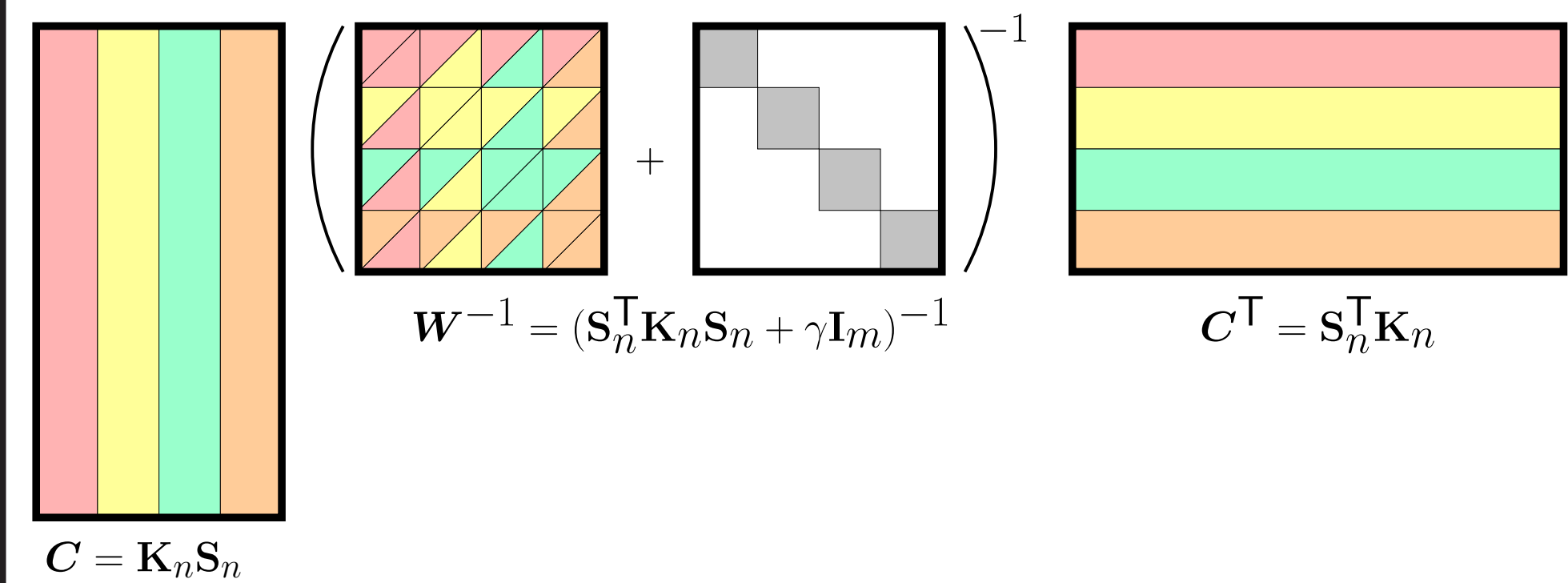
- Select a subset (dictionary)  $\mathcal{I}_n$  of  $m$  representative samples
- Constructs a sparse matrix  $\mathbf{S}_n$  to select and reweight the columns associated with the points in  $\mathcal{I}_n$



### Low-Rank Approximation

- Compute approximate, low-rank matrix  $\tilde{\mathbf{K}}_n = \mathbf{C}\mathbf{W}^{-1}\mathbf{C}^T$  as

$$\tilde{\mathbf{K}}_n = \mathbf{C}\mathbf{W}^{-1}\mathbf{C}^T = \mathbf{K}_n \mathbf{S}_n (\mathbf{S}_n^T \mathbf{K}_n \mathbf{S}_n + \gamma \mathbf{I}_m)^{-1} \mathbf{S}_n^T \mathbf{K}_n$$



### Efficient Solution

- Compute approximate solution Kernel Ridge Regression

$$\tilde{\mathbf{w}}_n = (\tilde{\mathbf{K}}_n + \mu \mathbf{I})^{-1} \mathbf{y}_n = \frac{1}{\mu} (\mathbf{y}_n - \mathbf{C} (\mathbf{C}^T \mathbf{C} + \mu \mathbf{W})^{-1} \mathbf{C}^T \mathbf{y}_n)$$

### Kernel K-Means

$$\min_{\mathbf{A}} \text{Tr}(\mathbf{K}_n - \mathbf{A}\mathbf{A}^T \mathbf{K}_n \mathbf{A}\mathbf{A}^T) \sim \min_{\mathbf{A}} \text{Tr}(\tilde{\mathbf{K}}_n - \mathbf{A}\mathbf{A}^T \tilde{\mathbf{K}}_n \mathbf{A}\mathbf{A}^T)$$

### Kernel PCA

$$\min_{\mathbf{Z}} \|\mathbf{K}_n - \mathbf{Z}\mathbf{Z}^T \mathbf{K}_n\|_F \sim \min_{\mathbf{Z}} \|\tilde{\mathbf{K}}_n - \mathbf{Z}\mathbf{Z}^T \tilde{\mathbf{K}}_n\|_F$$

Also Kernel CCA, Kernel [Your downstream problem here]

### Scalability

 now depends on  $m$ 

$$\text{Space: } \mathcal{O}(n^2) \Rightarrow \mathcal{O}(nm), \quad \text{Time: } \mathcal{O}(n^3) \Rightarrow \mathcal{O}(nm^2 + m^3)$$

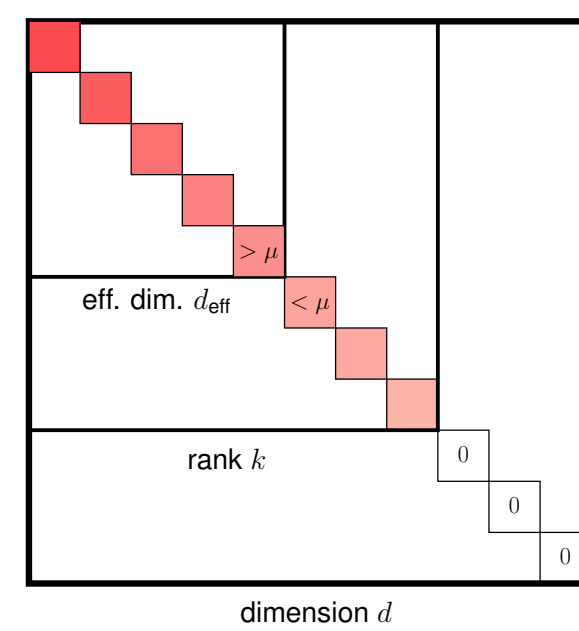
### Problems:

- ? How to choose the sampling distribution?
- ? How to choose  $m$ ?

## References

- [Alaoui and Mahoney (2015)] A. El Alaoui and M. W. Mahoney. Fast randomized kernel methods with statistical guarantees. In *NIPS*, 2015.
- [Bach (2013)] F. Bach. Sharp analysis of low-rank kernel matrix approximations. In *COLT*, 2013.
- [Calandriello et al. (2016)] D. Calandriello, A. Lazaric, and M. Valko. Analysis of Nyström method with sequential ridge leverage scores. In *UAI*, 2016.
- [Rudi et al. (2015)] A. Rudi, R. Camoriano, and L. Rosasco. Less is more: Nyström computational regularization. In *NIPS*, 2015.
- [Musco and Musco (2016)] C. Musco and C. Musco. Provably useful kernel matrix approximation in linear time. In arXiv, 2016.

## Kernel Ridge Leverage Scores (RLS) Sampling



**Definition 1.** Given a kernel matrix  $\mathbf{K}_n \in \mathbb{R}^{n \times n}$ , define

$$\gamma\text{-ridge leverage score} \quad \tau_{n,i}(\gamma) = \mathbf{e}_{n,i}^T \mathbf{K}_n (\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1} \mathbf{e}_{n,i} = \phi(\mathbf{x}_i)^T (\phi(\mathbf{X}_n) \phi(\mathbf{X}_n)^T + \gamma \mathbf{I})^{-1} \phi(\mathbf{x}_i) \quad (1)$$

$$\text{effective dimension} \quad d_{\text{eff}}(\gamma)_n = \sum_{i=1}^n \tau_{n,i}(\gamma) = \text{Tr}(\mathbf{K}_n (\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1}) \quad (2)$$

**Proposition 1** (Alaoui, Mahoney, 2015). Let  $\epsilon$  be the accuracy,  $\delta$  the confidence. If the regularized Nyström approximation  $\tilde{\mathbf{K}}_n$  is computed using a sampling distribution proportional to  $\tau_{n,i}$ , and at least

$$m \geq \left( \frac{2d_{\text{eff}}(\gamma)_n}{\epsilon^2} \right) \log \left( \frac{n}{\delta} \right)$$

columns, then with probability  $1 - \delta$ ,  $\mathbf{0} \preceq \mathbf{K}_n - \tilde{\mathbf{K}}_n \preceq \frac{\gamma}{1-\epsilon} \mathbf{I}_n$ .

**Intuitively:**  $\tau_{n,i}$  sensitivity of prediction on point  $\mathbf{x}_i$   
 $\Rightarrow \hat{\mathbf{y}}_{n,i} = \mathbf{e}_i^T (\mathbf{K}_n \hat{\mathbf{w}}_n) = \mathbf{e}_i^T \mathbf{K}_n (\mathbf{K}_n + \mu \mathbf{I})^{-1} \mathbf{y}_n$

**Pros:** +  $m$  scales with the effective dimension

**Cons:** - computing  $\tau_{n,i}(\mu)$  is as difficult as solving the original problem  
 - the probabilities need to be **recomputed at any new sample** (=multipass)

## SQUEAK

**Lemma 1.** Assume that the dictionary  $\mathcal{I}_{t-1}$  induces a  $\gamma$ -approx.  $\tilde{\mathbf{K}}_{t-1}$ , and let  $\tilde{\mathbf{S}}_t$  be constructed by adding  $\bar{q}$  copies of  $(\bar{q})^{-1/2} \mathbf{e}_{t,t}$  to the selection matrix. Then, denoting  $\alpha = (1 + \epsilon)/(1 - \epsilon)$ , for all  $i$  such that  $i \in \mathcal{I}_{t-1} \cup \{t\}$ ,

$$\tilde{\tau}_{t,i} = \frac{1 + \epsilon}{\alpha \gamma} \left( k_{i,i} - \mathbf{k}_{t,i} \tilde{\mathbf{S}} (\tilde{\mathbf{S}}^T \mathbf{K}_t \tilde{\mathbf{S}} + \gamma \mathbf{I})^{-1} \tilde{\mathbf{S}}^T \mathbf{k}_{t,i} \right), \quad (3)$$

is an  $\alpha$ -approximation of the RLS  $\tau_{t,i}$ , that is  $\tau_{t,i}(\gamma)/\alpha \leq \tilde{\tau}_{t,i} \leq \tau_{t,i}(\gamma)$ .

### SQUEAK

**Input:**  $\mathcal{D}$ , regularization  $\gamma, \bar{q}, \epsilon$ , **Output:**  $\mathcal{I}_n$

- Initialize  $\mathcal{I}_0$  as empty,  $\tilde{p}_{1,0} = 1$
- for**  $t = 1, \dots, n$  **do**
- Receive new sample  $\mathbf{x}_t$
- Compute  $\alpha$ -app. RLS  $\{\tilde{\tau}_{t,i} : i \in \mathcal{I}_{t-1} \cup \{t\}\}$ , using  $\mathcal{I}_{t-1}$ ,  $\mathbf{x}$ , and Eq. 3
- Set  $\tilde{p}_{t,i} = \min\{\tilde{\tau}_{t,i}, \tilde{p}_{t-1,i}\}$
- Initialize  $\mathcal{I}_t = \emptyset$
- for all**  $j \in \{1, \dots, t-1\}$  **do**
- if**  $q_{t-1,j} \neq 0$  **then**
- $\mathbf{q}_{t,j} \sim \mathcal{B}(\tilde{p}_{t,j}/\tilde{p}_{t-1,j}, q_{t-1,j})$
- Add  $(j, \phi_j, q_{t,j}, \tilde{p}_{t,j})$  to  $\mathcal{I}_t$ . SHRINK
- end if**
- end for** DICT-UPDATE
- $\mathbf{q}_{t,t} \sim \mathcal{B}(\tilde{p}_{t,t}, \bar{q})$
- Add  $q_{t,t}$  copies of  $(t, \phi_t, q_{t,t}, \tilde{p}_{t,t})$  to  $\mathcal{I}_t$  EXPAND
- end for**

**Theorem 1.** Let  $\alpha = (1+\epsilon)/(1-\epsilon)$  and  $\gamma > 1$ . For any  $0 \leq \epsilon \leq 1$ , and  $0 \leq \delta \leq 1$ , if we run SQUEAK with  $\bar{q} = \mathcal{O}(\frac{\alpha}{\epsilon^2} \log(\frac{n}{\delta}))$ , then w.p.  $1 - \delta$ , for all  $t \in [n]$

- $\tilde{\mathbf{K}}_t$  computed with  $\mathcal{I}_t$  is a  $\gamma$ -approximation of  $\mathbf{K}_t$ .
- $|\mathcal{I}_t| = \sum_i q_{t,i} \leq \mathcal{O}(\bar{q} d_{\text{eff}}(\gamma)_t) \leq \mathcal{O}(\frac{\alpha}{\epsilon^2} d_{\text{eff}}(\gamma)_n \log(\frac{n}{\delta}))$ .

Accuracy and space/time anytime guarantees, matches exact RLS sampling.

Using unnormalized  $\tilde{p}_{t,i}$ , no need for appr.  $d_{\text{eff}}(\gamma)_t$

Only need to compute RLS for points in  $\mathcal{I}_t$ , never recompute after dropping  
 ↳ Never construct the whole  $\mathbf{K}_n$ , subquadratic runtime  $\mathcal{O}(n^2 |\mathcal{I}_n|^2) \Rightarrow \mathcal{O}(n |\mathcal{I}_n|^3)$

Store points directly in the dictionary  
 ↳  $\mathcal{O}(d_{\text{eff}}(\gamma)_n^2 + d_{\text{eff}}(\gamma)_n d)$  space constant in  $n$   
 Single pass over the dataset (streaming)

Dictionary changes a lot between iteration, total runtime  $\mathcal{O}(n |\mathcal{I}_n|^3)$

Extend DICT-UPDATE (point + dict.) to DICT-MERGE (dict. + dict.)  
 ↳ Distributed SQUEAK, multiple workers in parallel, without sharing memory  
 Recursive merging to build dictionary,  $\mathcal{O}(\log(n) |\mathcal{I}_n|^3)$  time,  $\mathcal{O}(n |\mathcal{I}_n|^3)$  work

	Time	$ \mathcal{I}_n $	Incr.
EXACT	$n^3$	$n$	-
Bach'13	$\frac{n d_{\text{max},n}^2}{\epsilon} + \frac{d_{\text{max},n}^3}{\epsilon}$	$\frac{d_{\text{max},n}}{\epsilon}$	No
A&M'15	$n( \mathcal{I}_n )^2$	$\left( \frac{\lambda_{\min} + n\mu\epsilon}{\lambda_{\min} - n\mu\epsilon} \right) d_{\text{eff}}(\gamma)_n$	No
INK (C&a'16)	$\frac{\lambda_{\text{max}}^2}{\gamma^2} n^2 d_{\text{eff}}(\gamma)_n^2$	$\frac{\lambda_{\text{max}}}{\gamma} d_{\text{eff}}(\gamma)_n$	Yes
SQUEAK	$\frac{n^2 d_{\text{eff}}(\gamma)_n^2}{\epsilon^2}$	$d_{\text{eff}}(\gamma)_n$	Yes

## Downstream guarantees (Musco & Musco 2016)

RLS sampling preserves well the projection on  $\mathbf{K}_n$ 's range  $\mathbf{P} = \mathbf{K}_n^{1/2} (\mathbf{K}_n + \gamma \mathbf{I})^{-1} \mathbf{K}_n^{1/2} = \phi(\mathbf{X}_n)^T (\phi(\mathbf{X}_n) \phi(\mathbf{X}_n)^T + \gamma \mathbf{I})^{-1} \phi(\mathbf{X}_n)$

Kernel ridge regression:

$$\hat{\mathbf{y}}_{n,i} = \mathbf{e}_i^T \mathbf{K}_n (\mathbf{K}_n + \mu \mathbf{I})^{-1} \mathbf{y}_n = \mathbf{e}_i^T \mathbf{P} \mathbf{y}_n$$

$$\tilde{\mathbf{w}}_n = (\tilde{\mathbf{K}}_n + \gamma \mathbf{I})^{-1} \mathbf{y}_n$$

Kernel PCA:

$$\mathbf{K}_n = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T, \quad \mathbf{P} = \mathbf{U} \mathbf{\Lambda} (\mathbf{\Lambda} + \gamma \mathbf{I})^{-1} \mathbf{U}^T$$

$$\tilde{\mathbf{Z}} \text{ computed using } \tilde{\mathbf{K}}_n = \tilde{\mathbf{U}} \tilde{\mathbf{\Lambda}} \tilde{\mathbf{U}}^T$$

$$R(\tilde{\mathbf{w}}_n) \leq \left(1 + \frac{1}{1-\epsilon}\right) R(\hat{\mathbf{w}}_n)$$

$$\|\mathbf{K}_n - \tilde{\mathbf{Z}} \tilde{\mathbf{Z}}^T \mathbf{K}_n\|_F \leq (1 + 2\epsilon) \|\mathbf{K}_n - \mathbf{Z}^* \mathbf{Z}^{*T} \mathbf{K}_n\|_F$$

Kernel K-Means:

$$\tilde{\mathbf{A}} \text{ } \rho\text{-optimal cluster assignment for } \tilde{\mathbf{K}}_n$$

$$\mathbf{A}^* \text{ optimal cluster assignment for } \mathbf{K}_n$$

$$\xi = (1 + \epsilon)(1 + \rho)$$

$$\text{Tr}(\mathbf{K}_n - \mathbf{A} \mathbf{A}^T \mathbf{K}_n \mathbf{A} \mathbf{A}^T) \leq \xi \text{Tr}(\mathbf{K}_n - \mathbf{A}^* \mathbf{A}^{*T} \mathbf{K}_n \mathbf{A}^* \mathbf{A}^{*T})$$

