

# ANALYSIS OF NYSTRÖM METHOD WITH SEQUENTIAL RIDGE LEVERAGE SCORE SAMPLING

DANIELE CALANDRIELLO, ALESSANDRO LAZARIC, MICHAL VALKO



## MOTIVATION

- Kernel regression is *versatile* and *accurate*
- Strong accuracy guarantees but *poor scalability*

$$\mathcal{O}(n^3) \text{ time } \mathcal{O}(n^2) \text{ space } (n \text{ number of samples})$$

- Current limitation: Many approximate schemes are either *not scalable* or *not accurate*

⇒ We propose an incremental approximation scheme for kernel regression with *complexity and error guarantees* depending on the *kernel structure*

## KERNEL RIDGE REGRESSION (KRR)

### The setting (fixed-design)

- Dataset  $\mathcal{D} = \{\mathbf{x}_t, y_t\}_{t=1}^n$

- *arbitrary*  $\mathbf{x}_t \in \mathcal{X}$
- $y_t = f^*(\mathbf{x}_t) + \eta_t$

- Kernel function  $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

- Kernel matrix  $\mathbf{K}_t \in \mathbb{R}^{t \times t}$ , with  $[\mathbf{K}_t]_{i,j} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$ ,  $i, j \leq t$

### Kernel regression

- Objective (after  $t$  samples)

$$\hat{\mathbf{w}}_t = \arg \min_{\mathbf{w}} \|\mathbf{y}_t - \mathbf{K}_t \mathbf{w}\|^2 + \mu \|\mathbf{w}\|^2.$$

- Closed-form solution

$$\hat{\mathbf{w}}_t = (\mathbf{K}_t + \mu \mathbf{I})^{-1} \mathbf{y}_t$$

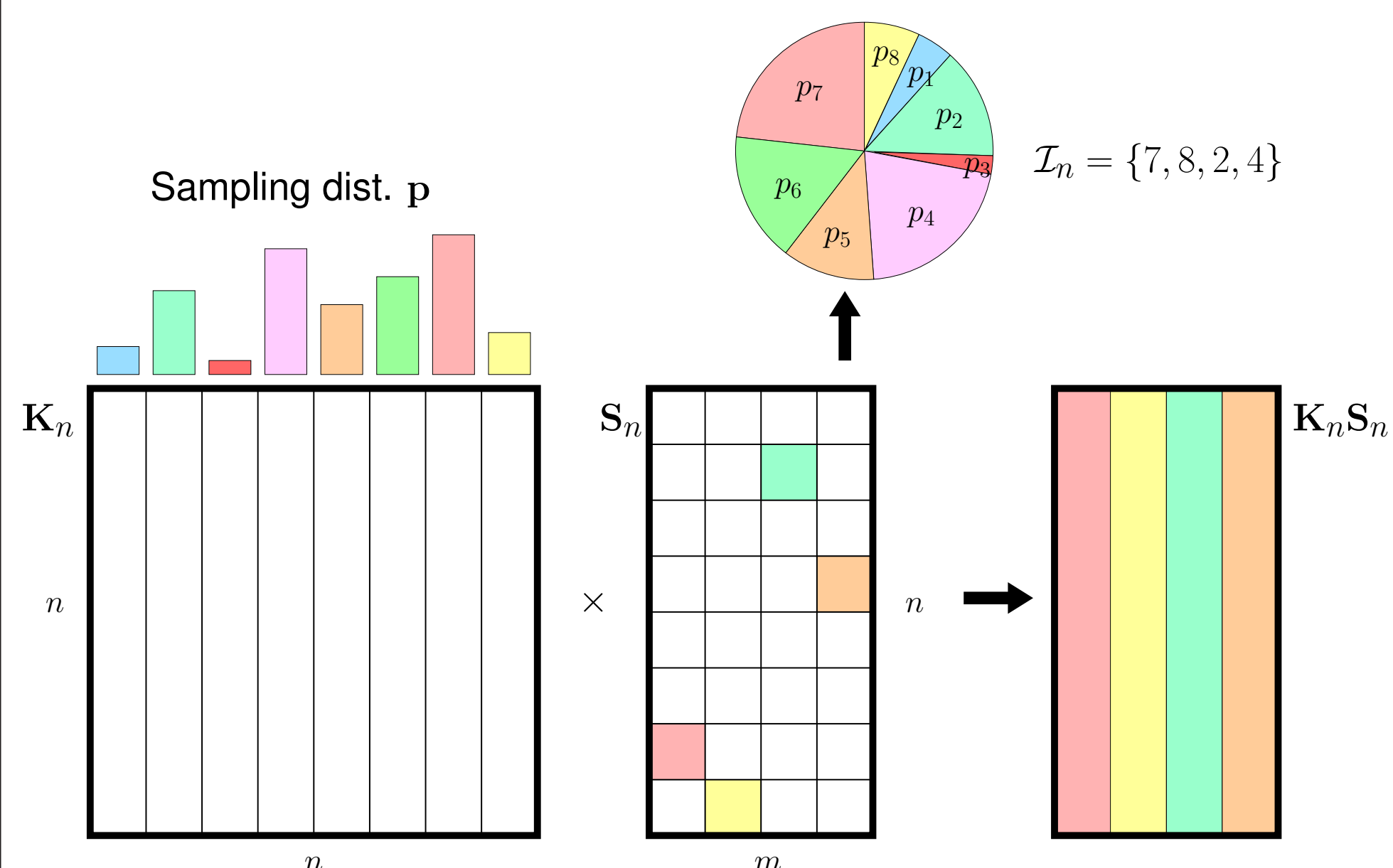
- On-sample risk

$$\mathcal{R}(\hat{\mathbf{w}}_t) = \mathbb{E}_{\eta} [\|\mathbf{f}_t^* - \mathbf{K}_t \hat{\mathbf{w}}_t\|^2]$$

## NYSTRÖM APPROXIMATION

### Subsampling

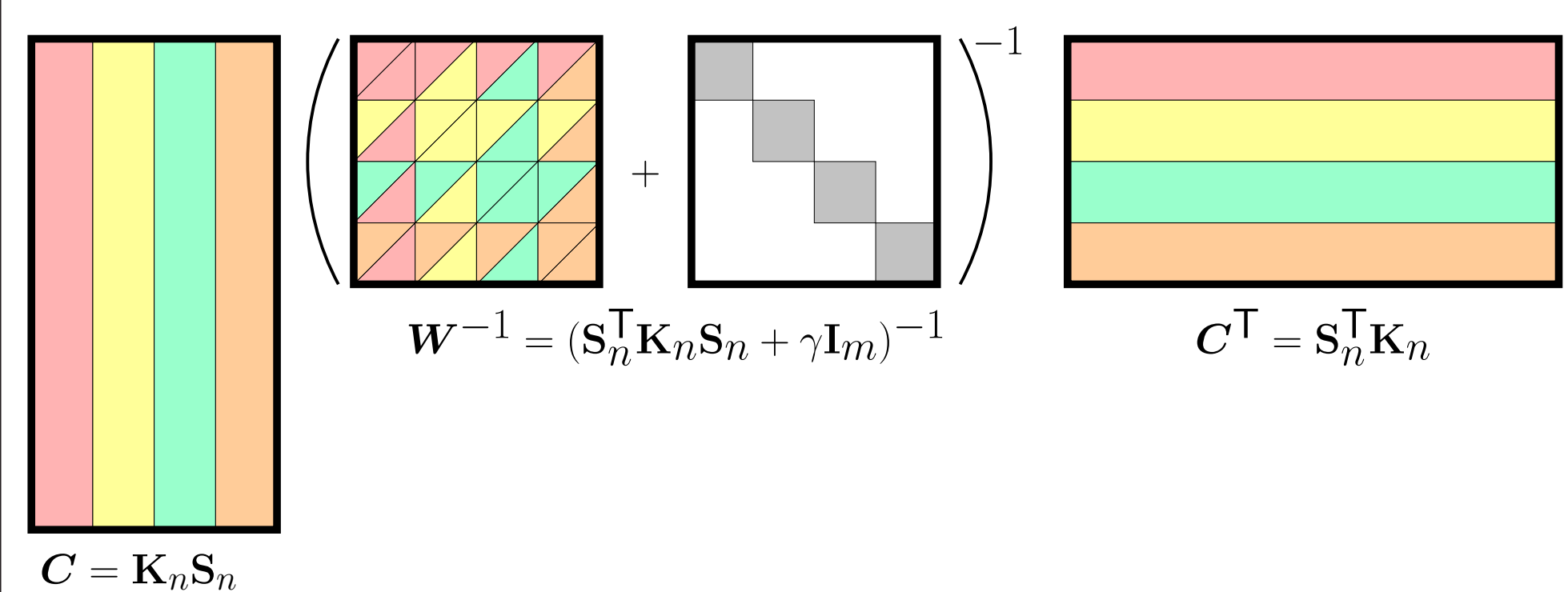
- Select a subset (dictionary)  $\mathcal{I}_n$  of  $m$  representative samples
- Constructs a sparse matrix  $\mathbf{S}_n$  to select and reweight the columns associated with the points in  $\mathcal{I}_n$



### 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

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

### Scalability

now depends on  $m$

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

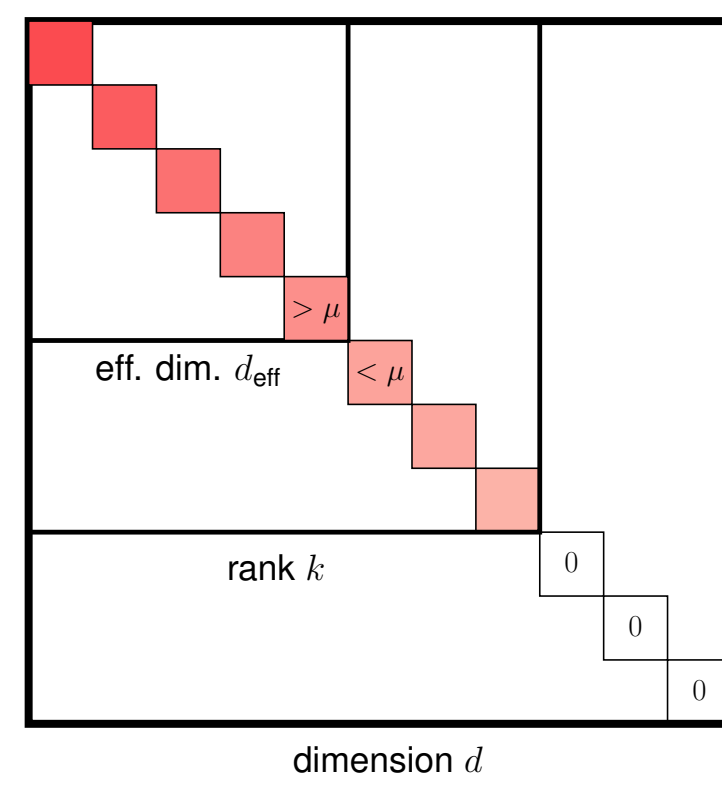
### Problems:

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

## REFERENCES

- [Alaoui and Mahoney(2015)] Ahmed El Alaoui and Michael W. Mahoney. Fast randomized kernel methods with statistical guarantees. In *Neural Information Processing Systems*, 2015.
- [Bach(2013)] Francis Bach. Sharp analysis of low-rank kernel matrix approximations. In *International Conference on Learning Theory*, 2013.
- [Pachocki(2016)] Jakub Pachocki. Analysis of reparsification. *arXiv preprint arXiv:1605.08194*, 2016.
- [Rudi et al.(2015)] Alessandro Rudi, Raffaello Camoriano, and Lorenzo Rosasco. Less is more: Nyström computational regularization. In *Neural Information Processing Systems*, 2015.

## KERNEL RIDGE LEVERAGE SCORES (RLS) SAMPLING FOR KRR



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

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

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

$$\text{sampling distribution } [\mathbf{p}_n]_i = p_{i,n} = \frac{\tau_{i,n}(\gamma)}{\sum_{j=1}^n \tau_{j,n}(\gamma)} = \frac{\tau_{i,n}}{d_{\text{eff}}(\gamma)} \quad (3)$$

**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 the sampling distribution  $\{p_{i,t}\}$ , 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  $1 - \delta$

$$0 \leq \mathbf{K}_n - \tilde{\mathbf{K}}_n \leq \frac{\gamma}{1 - \epsilon} \mathbf{I}_n, \quad \mathcal{R}(\tilde{\mathbf{w}}_n) \leq \left( 1 + \frac{\gamma}{\mu} \frac{1}{1 - \epsilon} \right)^2 \mathcal{R}(\hat{\mathbf{w}}_n)$$

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

**Pros:** +  $m$  scales with the effective dimension  
 + the risk for  $\tilde{\mathbf{w}}_n$  is almost the same as for the exact solution

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

## INCREMENTAL ESTIMATES OF RLS AND EFFECTIVE DIMENSION

For any column  $i$  in  $\mathcal{I}_t$  and  $\mathbf{k}_{t+1}$  compute the ridge leverage score estimator ( $\alpha = \frac{2-\epsilon}{1-\epsilon}$ )

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

- $\tilde{\tau}_{i,t+1} = \mathbf{e}_i^T \tilde{\mathbf{K}}_t (\tilde{\mathbf{K}}_t + \gamma \mathbf{I})^{-1} \mathbf{e}_i$  would fail

- $\tilde{\tau}_{i,t+1}$  is computed only for columns stored in  $\mathcal{I}_t$  (accurate)

- $\tilde{\tau}_{i,t+1}$  can be computed in a space/time efficient way

- $\alpha$  trades off accuracy of the estimator and space/time cost

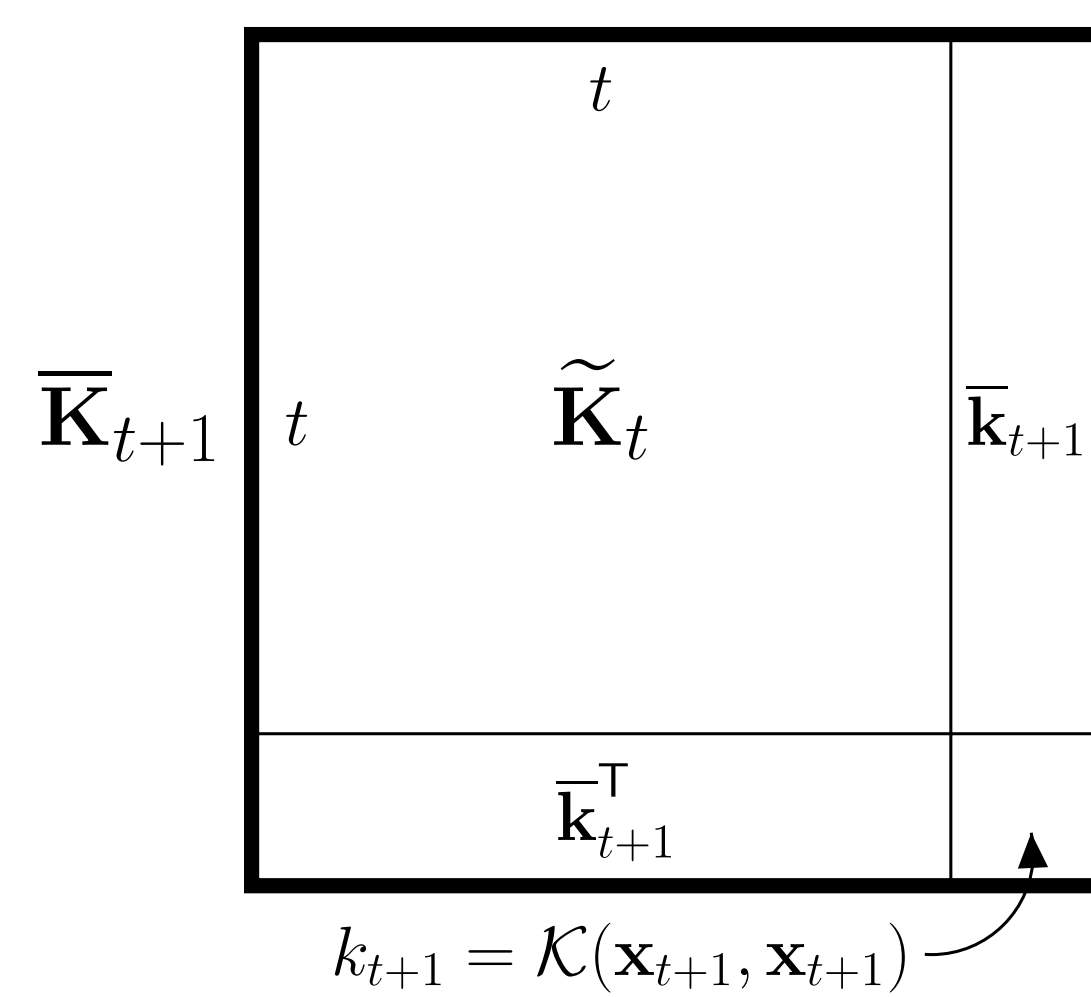
Compute the effective dimension estimator  $\tilde{d}_{\text{eff}}(\gamma)_{t+1} = \tilde{d}_{\text{eff}}(\gamma)_t + \alpha \tilde{\Delta}_t$  with

$$\tilde{\Delta}_t = \frac{(k_{t+1} - \bar{\mathbf{k}}_{t+1}^T (\tilde{\mathbf{K}}_t + \alpha \gamma \mathbf{I})^{-1} \bar{\mathbf{k}}_{t+1}) - \frac{(1-\epsilon)^2}{4} \gamma \bar{\mathbf{k}}_{t+1}^T (\tilde{\mathbf{K}}_t + \gamma \mathbf{I})^{-2} \bar{\mathbf{k}}_{t+1}}{k_{t+1} + \gamma - \bar{\mathbf{k}}_{t+1}^T (\tilde{\mathbf{K}}_t + \alpha \gamma \mathbf{I})^{-1} \bar{\mathbf{k}}_{t+1}}$$

- $\tilde{d}_{\text{eff}}(\gamma)_{t+1} = \sum_{i=1}^{t+1} \tilde{\tau}_{i,t+1}$  requires  $\tilde{\tau}_{i,t+1}$  for  $i \notin \mathcal{I}_t$  (not accurate)

- $\tilde{\Delta}_t$  captures the interaction between the new and past samples

- $\tilde{\Delta}_t$  requires approximating "second order" terms for which first order reconstruction guarantees ( $0 \leq \mathbf{K}_t - \tilde{\mathbf{K}}_t \leq \frac{\gamma}{1-\epsilon} \mathbf{I}$ ) are not enough



**Lemma 1.** Let  $\epsilon$  be the accuracy and  $\rho = \lambda_{\max}(\mathbf{K}_n)/\gamma$  a soft condition number. If after  $t$  samples  $\tilde{\mathbf{K}}_t$  is such that  $0 \leq \mathbf{K}_t - \tilde{\mathbf{K}}_t \leq \frac{\gamma}{1-\epsilon} \mathbf{I}$ , then for  $\alpha = \frac{2-\epsilon}{1-\epsilon}$  and  $\beta = \left( \frac{2-\epsilon}{1-\epsilon} \right)^2 (1 + \rho)$ , the estimators satisfy for any  $i \in \{\mathcal{I}_t \cup t+1\}$

$$\frac{1}{\alpha} \tau_{i,t+1}(\gamma) \leq \tilde{\tau}_{i,t+1} \leq \mathbf{1} \cdot \tau_{i,t+1}(\gamma), \quad \mathbf{1} \cdot d_{\text{eff}}(\gamma)_{t+1} \leq \tilde{d}_{\text{eff}}(\gamma)_{t+1} \leq \beta d_{\text{eff}}(\gamma)_{t+1}.$$

and the estimated probabilities satisfy

$$\frac{1}{\alpha \beta} p_{i,t+1} \leq \tilde{p}_{i,t+1} \leq \mathbf{1} \cdot p_{i,t+1}$$

## INK-ESTIMATE

### INK-ESTIMATE

**Input:** Dataset  $\mathcal{D}$ , regularization  $\gamma$ , sampling budget  $\bar{q}$

**Output:**  $\tilde{\mathbf{K}}_n, \mathbf{S}_n$

- Initialize  $\mathcal{I}_0$  as empty,  $\tilde{p}_{1,0} = 1, b_{1,0} = 1$ , budget  $\bar{q}$
- for  $t = 0, \dots, n-1$  do
- Receive new column  $\bar{\mathbf{k}}_{t+1}$  and scalar  $k_{t+1}$
- Compute *approximate leverage scores*  $\{\tilde{\tau}_{i,t+1} : i \in \mathcal{I}_t \cup \{t+1\}\}$
- Compute *approximate effective dimension*  $\tilde{d}_{\text{eff}}(\gamma)_{t+1}$
- Set  $\tilde{p}_{i,t+1} = \min\{\tilde{\tau}_{i,t+1}/\tilde{d}_{\text{eff}}(\gamma)_{t+1}, \tilde{p}_{i,t}\}$
- $\mathcal{I}_{t+1}, \mathbf{b}_{t+1} = \text{SHRINK-EXPAND}(\mathcal{I}_t, \tilde{\mathbf{p}}_{t+1}, \mathbf{b}_t, \bar{q})$
- Compute  $\mathbf{S}_{t+1}$  using  $\mathcal{I}_{t+1}$  and weights  $\sqrt{b_{i,t+1}}$
- Compute  $\tilde{\mathbf{K}}_{t+1}$  using  $\mathbf{S}_{t+1}$
- end for
- Return  $\tilde{\mathbf{K}}_n$  and  $\mathbf{S}_n$

**Theorem 1.** Let  $\epsilon$  be the desired accuracy and  $\rho = \lambda_{\max}(\mathbf{K}_n)/\gamma$  a soft condition number. If INK-ESTIMATE is run with

$$\bar{q} \geq \left( \frac{28\alpha\beta d_{\text{eff}}(\gamma)_t}{\epsilon^2} \right) \log \left( \frac{4t}{\delta} \right),$$

then the approximate kernel solution  $\tilde{\mathbf{w}}_n$  satisfies

$$\mathcal{R}(\tilde{\mathbf{w}}_n) \leq \left( 1 + \frac{\gamma}{\mu} \frac{1}{1 - \epsilon} \right)^2 \mathcal{R}(\hat{\mathbf{w}}_n)$$

and INK-ESTIMATE runs in at most

$$\mathcal{O}(n\bar{q}) \leq \tilde{\mathcal{O}}(n\rho d_{\text{eff}}(\gamma)_n) \quad \text{space,}$$

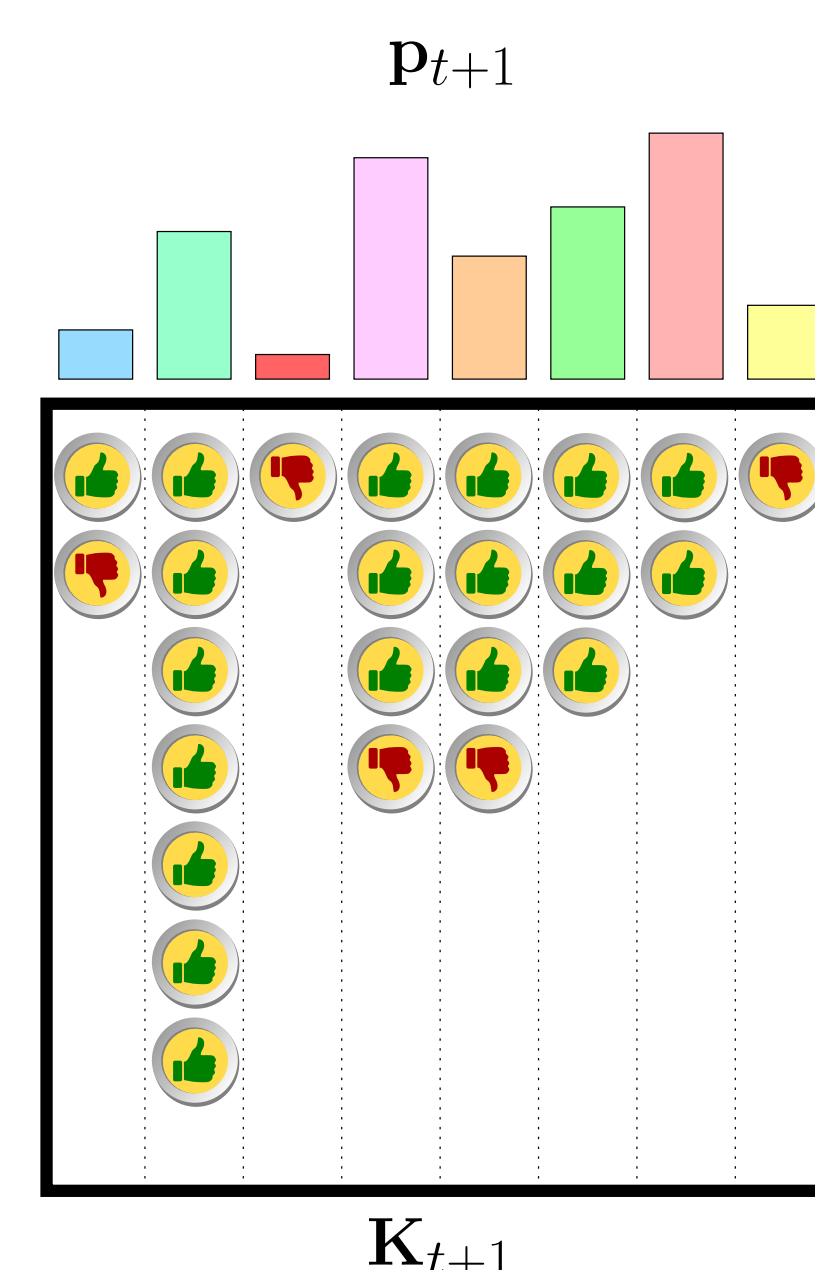
$$\mathcal{O}(n^2\bar{q}^2 + n\bar{q}^3) \leq \tilde{\mathcal{O}}(n^2\rho^2 d_{\text{eff}}(\gamma)_n^2) \quad \text{time}$$

### SHRINK-EXPAND (Pachocki, 2016)

**Input:**  $\mathcal{I}_t, \{\tilde{p}_{i,t+1}, b_{i,t}\} : i \in \mathcal{I}_t, \tilde{p}_{t+1,t+1}, \bar{q}$

**Output:**  $\mathcal{I}_{t+1}$ , the set of all columns with  $b_{i,t+1} \neq 0$

- $b_{i,t+1} = b_{i,t}$  for all  $i \in \mathcal{I}_t, b_{t+1,t+1} = 1$
- for all  $i \in \{1, \dots, t\} : b_{i,t} \neq 0$  do  $\triangleright$ SHRINK
- while  $b_{i,t+1} \tilde{p}_{i,t+1} \leq 1/\bar{q}$  do
- Sample a random Bernoulli  $\mathcal{B} \left( \frac{b_{i,t+1}}{b_{i,t+1}+1} \right)$
- On success set  $b_{i,t+1} = b_{i,t+1} + 1$
- On failure set  $b_{i,t+1} = 0$ , break
- end while
- end for
- while  $b_{t+1,t+1} \tilde{p}_{t+1,t+1} \leq 1/\bar{q}$  do  $\triangleright$ EXPAND
- Sample a random Bernoulli  $\mathcal{B} \left( \frac{b_{t+1,t+1}}{b_{t+1,t+1}+1} \right)$
- On success set  $b_{t+1,t+1} = b_{t+1,t+1} + 1$
- On failure set  $b_{t+1,t+1} = 0$ , break
- end while



**Lemma 2.** For any kernel matrix  $\mathbf{K}_t$  at time  $t$ , and its bordering  $\mathbf{K}_{t+1}$  at time  $t+1$ ,

$$\frac{\tau_{i,t+1}}{d_{\text{eff},t+1}} = \mathbf{p}_{i,t+1} \leq \mathbf{p}_{i,t} = \frac{\tau_{i,t}}{d_{\text{eff},t}}$$

### Pros:

- + Accuracy *and* space/time guarantees
- + In the worst case, only  $\sqrt{n}$  space overhead (wrt exact method)
- + Anytime risk guarantees

### Cons:

- The time complexity is not fully satisfactory
- The current formulation of the estimators is not "fully" incremental

### Open questions:

- ? Removing the dependency on  $\rho$
- ? Random design (Rudi et al., 2015)
- ? Online learning

	Time	Space	Acc. loss	Inc.
EXACT	$n^3$	$n^2$	1	/
Bach'13	$\frac{nd_{\max}^2 + d_{\max}^3}{\epsilon}$	$\frac{nd_{\max}}{\epsilon}$	$(1 + 4\epsilon)$	No
A&M'15	$n(\text{space})^2$	$\left( \frac{\lambda_{\min} + n\mu\epsilon}{\lambda_{\min} - n\mu\epsilon} \right) nd_{\text{eff}} + \frac{\text{Tr}(\mathbf{K}_n)}{\mu\epsilon}$	$(1 + 2\epsilon)^2$	No
INK-EST	$\frac{\rho^2 n^2 d_{\text{eff}}^2}{\epsilon}$	$\frac{\rho n d_{\text{eff}}}{\epsilon}$	$(1 + 2\epsilon)^2$	Yes