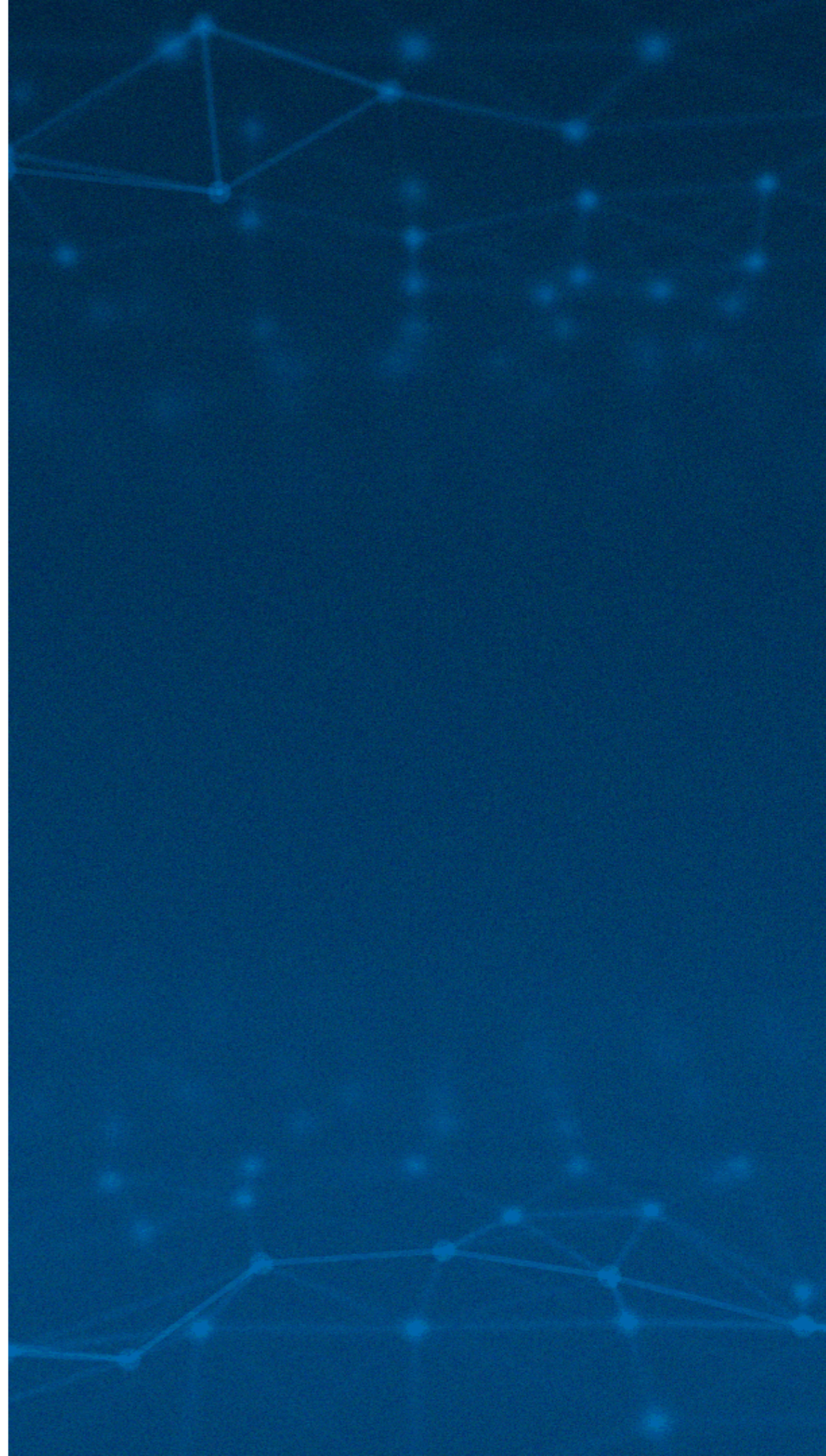


# How negative dependence broke the quadratic barrier for learning with graphs and kernels

Michal Valko

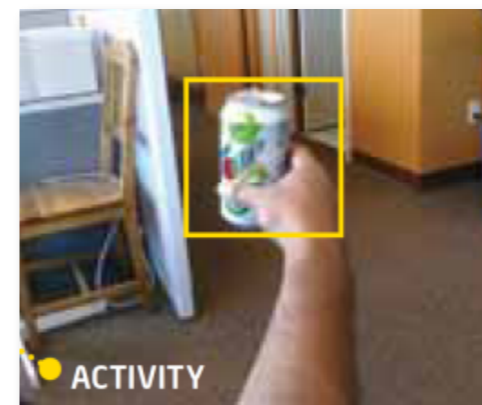
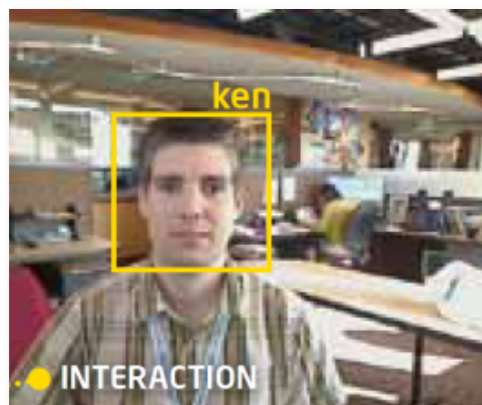
# ONLINE LEARNING

.....  
when we reason on the fly



# IN 2007 IT ALL STARTED WITH AN IDEA...

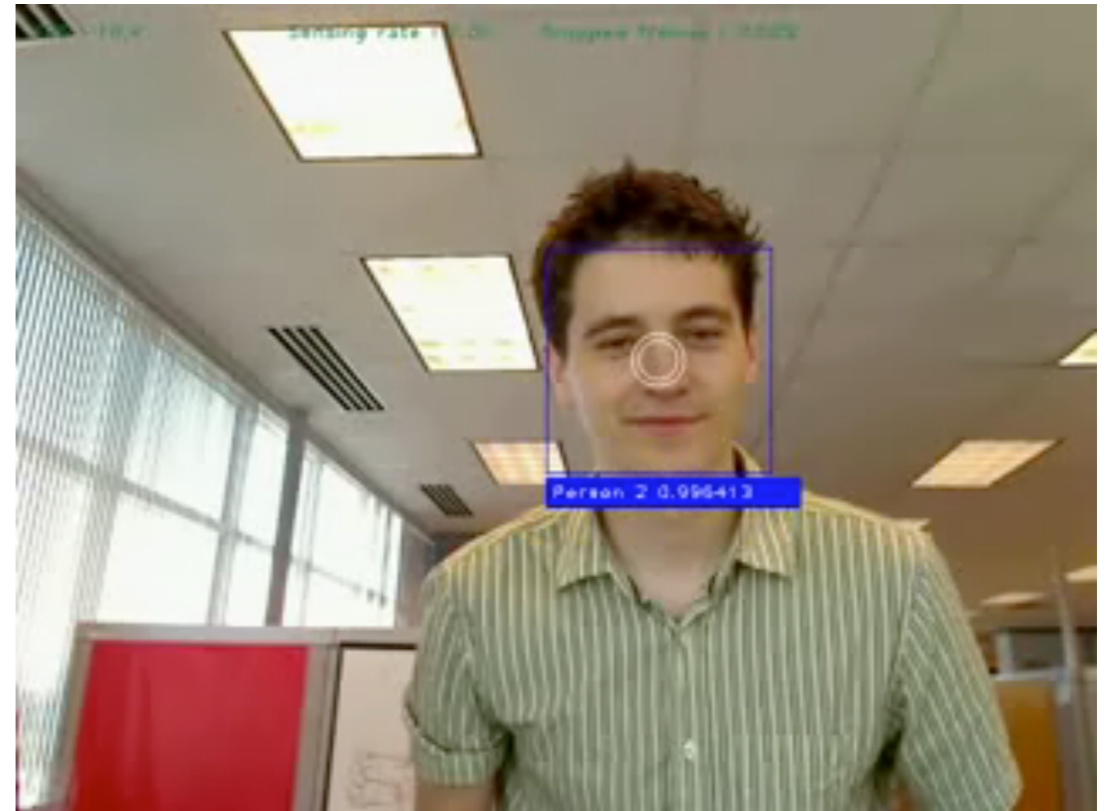
- Develop **sequential machine learning** recognition system
- System with **minimal feedback**
- 90% accurate over 90% of time
- With **theory** that guarantee's its performance
- **Efficient** (e.g., mobile device)



from B. Kveton

# ... AND RESULTED IN A REAL SYSTEM IN 2009

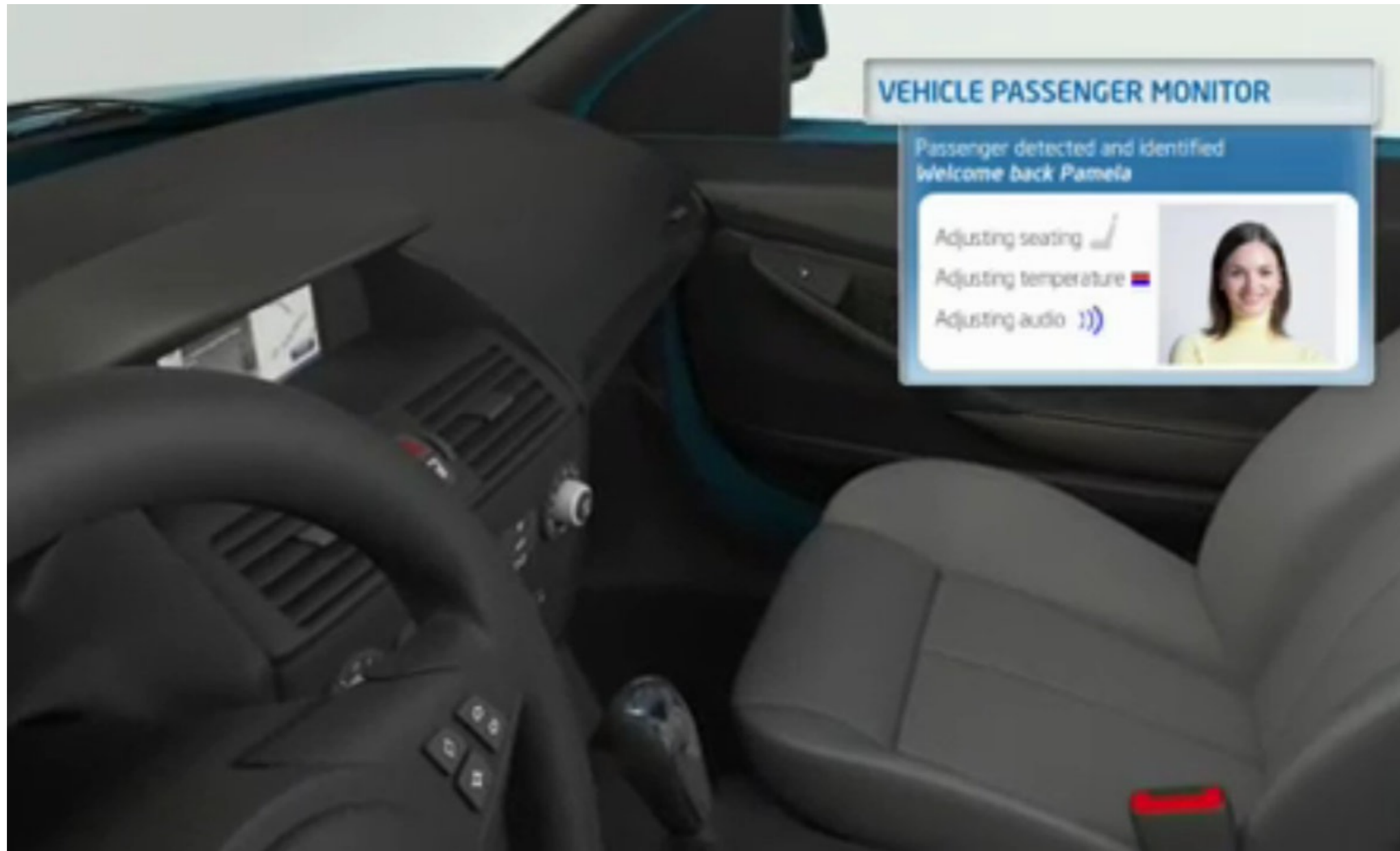
- **adaptive graph-based** recognition system
  - highly accurate
  - trained from a **small amount** of labeled data
  - real-time running time
  - robust to outliers
  - theoretical analysis



from B. Kveton

$$\frac{1}{n} \sum_t (\ell_t^q[t] - y_t)^2 \leq \frac{1}{n_l} \sum_{i \in I} (l_i^* - y_i)^2 + O(n^{-\frac{1}{2}})$$

# THIS CAN'T SCALE: CONNECTED CAR



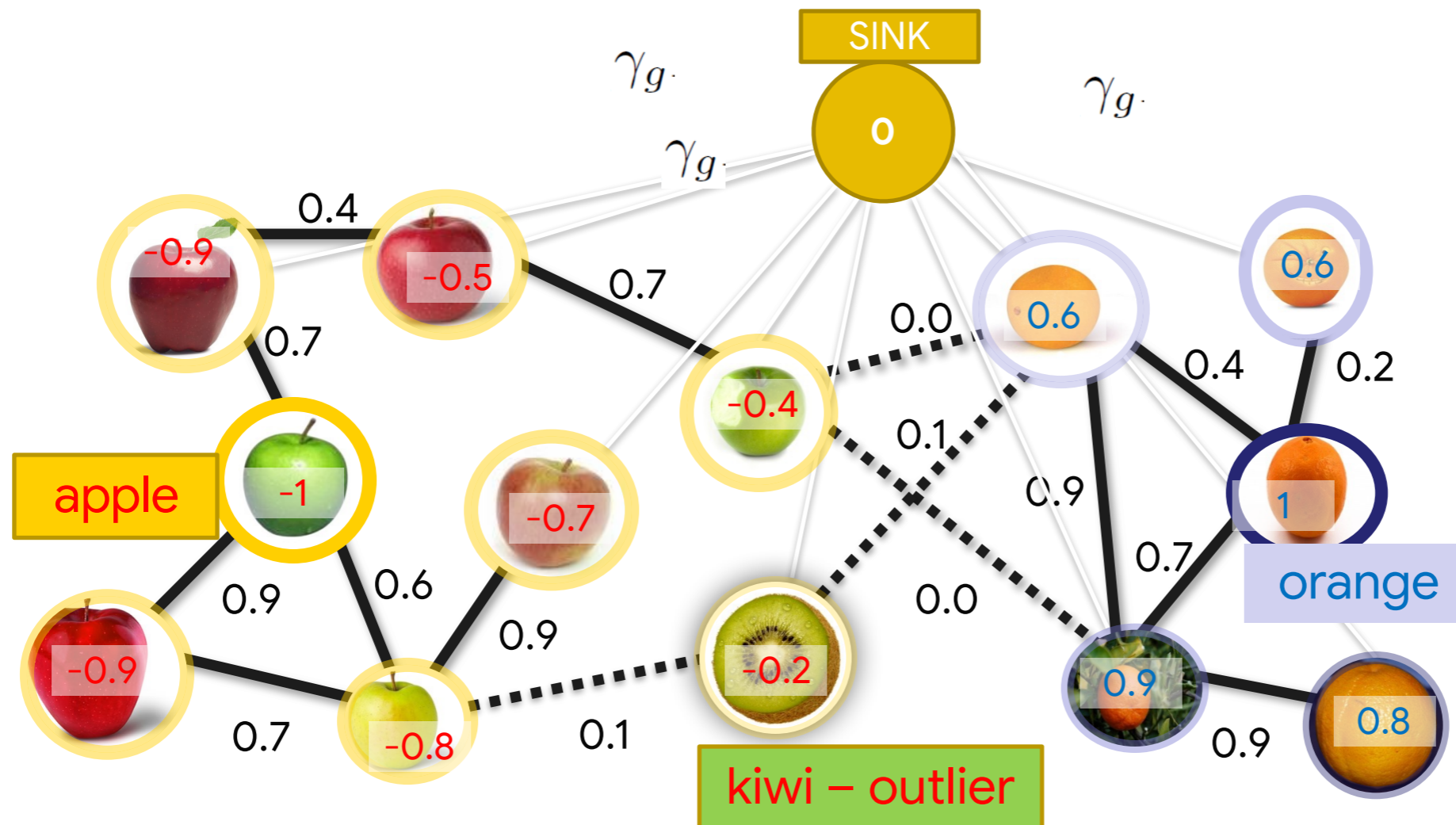
Personalization

# 2 BIG REAL-WORLD ISSUES

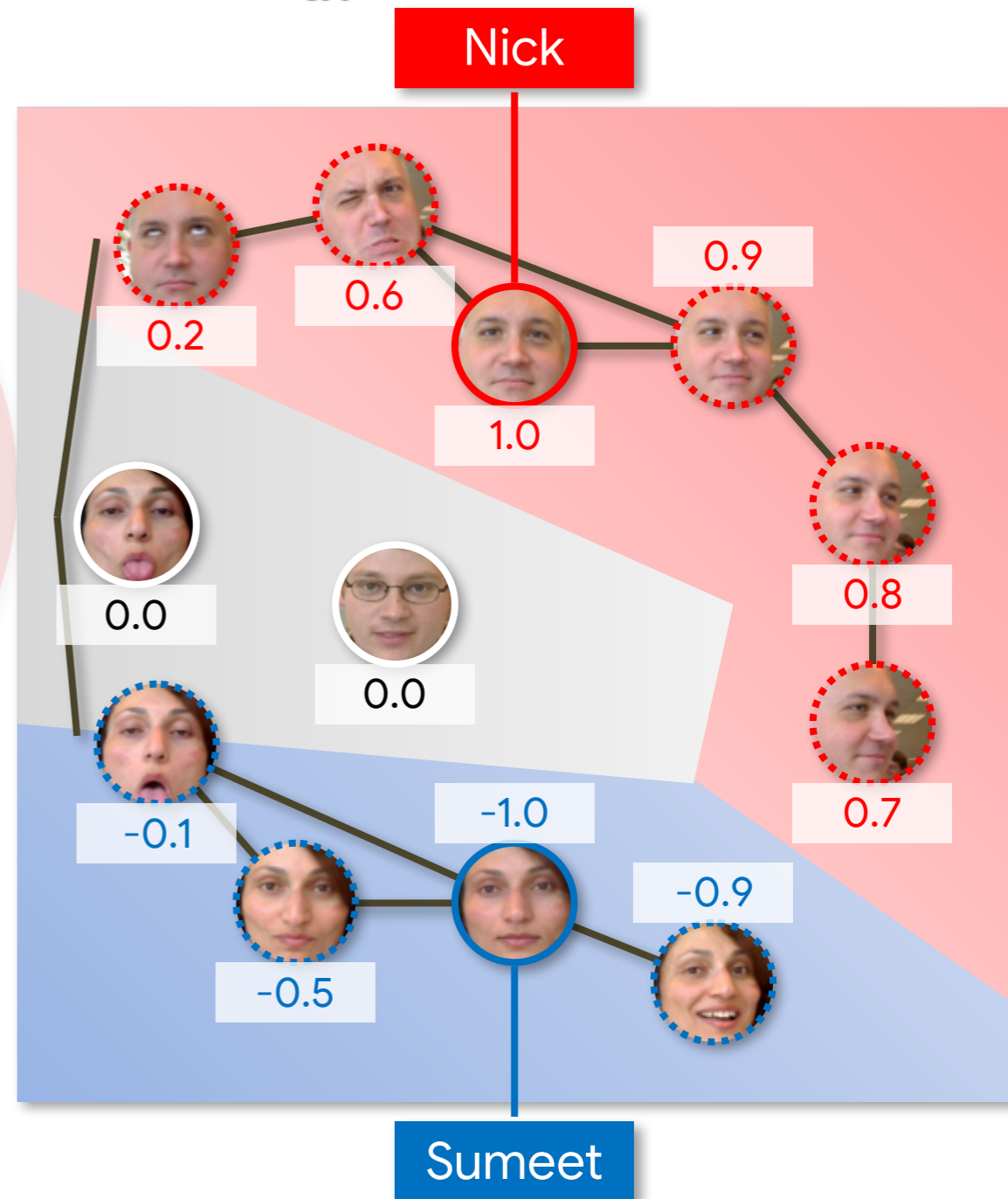
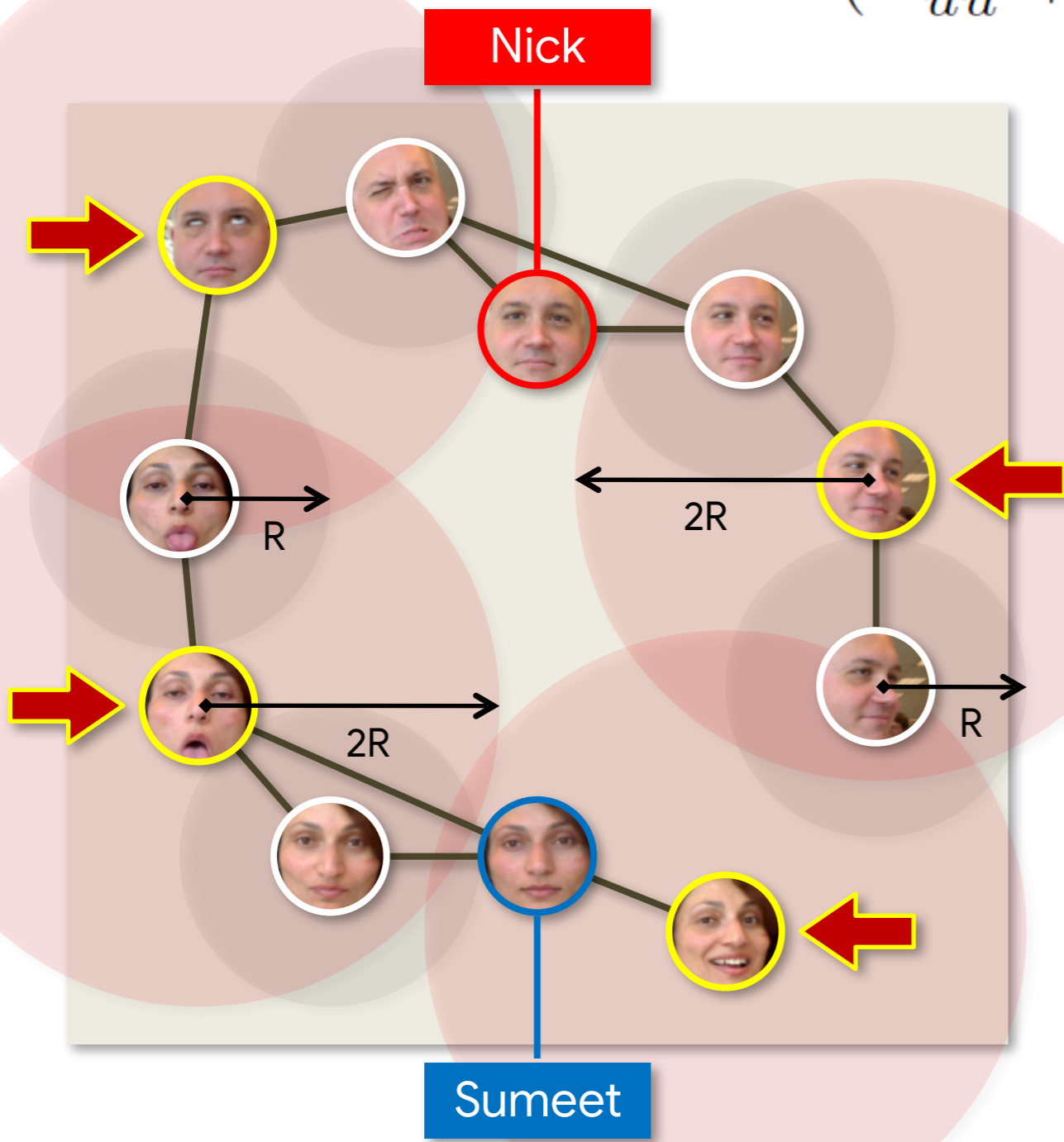
► SIZE and SPEED

► ANOMALIES

$$\mathbf{f}_u = (\mathbf{L}_{uu} + \gamma \mathbf{g} \mathbf{I})^{-1} (\mathbf{W}_{ul} \mathbf{f}_l)$$



$$\ell^q = (L_{uu}^q + \gamma_g V)^{-1} W_{ul}^q \ell_l$$



$$\frac{1}{n} \sum_t (\ell_t^q[t] - y_t)^2 \leq \frac{3}{n} \sum_t (\ell_t^* - y_t)^2 + \frac{3}{n} \sum_t (\ell_t^o[t] - \ell_t^*)^2 + \frac{3}{n} \sum_t (\ell_t^q[t] - \ell_t^o[t])^2$$

Error of our solution

Offline learning error

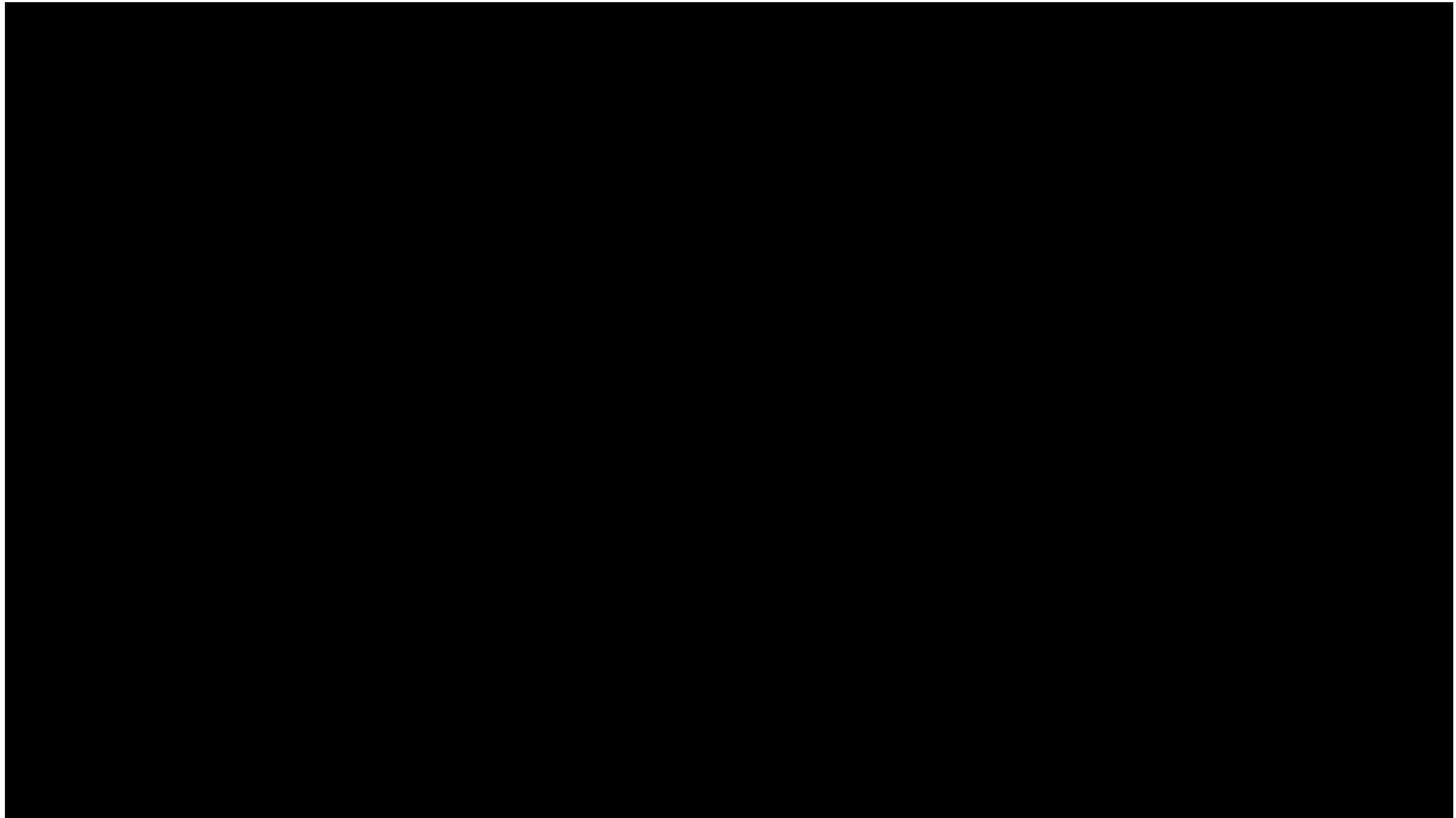
Online learning error

Quantization error

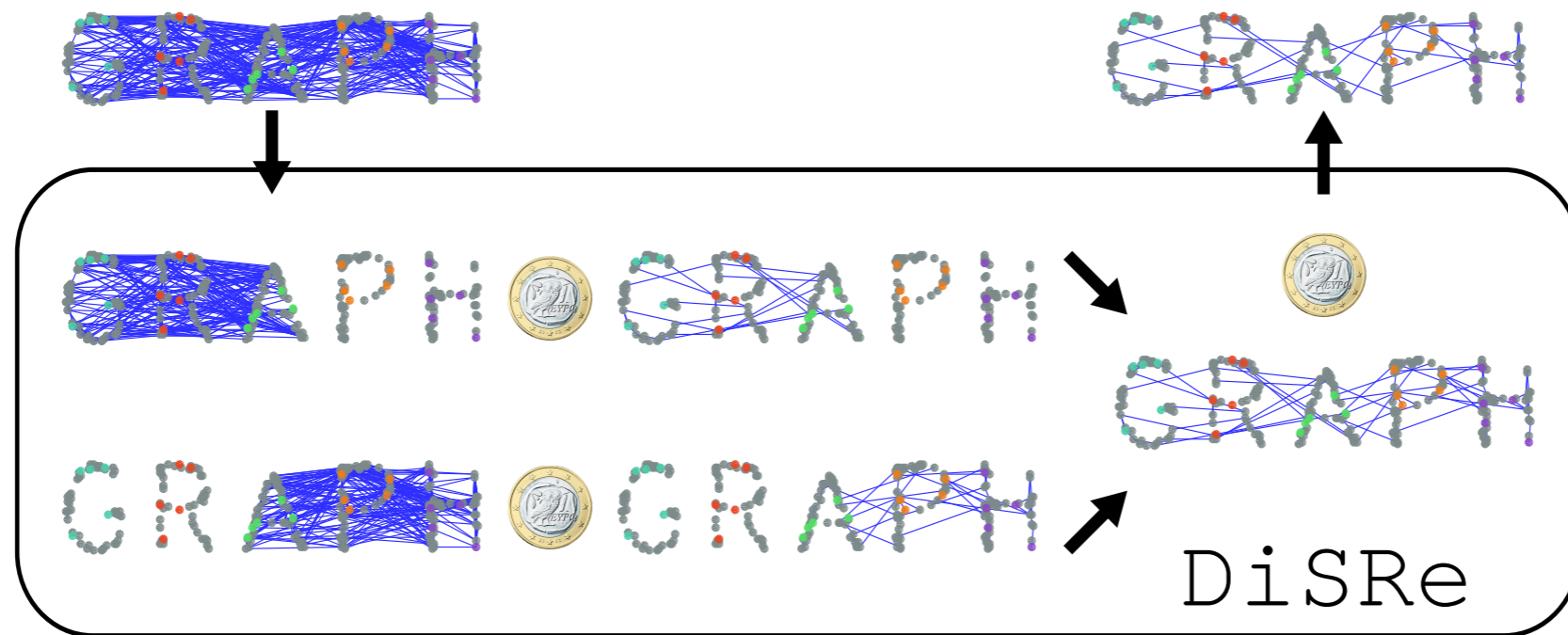
# FACE-RECOGNITION FOR INTEL



DeepMind







MV, Kveton, Huang, Ting: **Online Semi-Supervised Learning on Quantized Graphs** UAI 2010

Kveton, MV, Rahimi, Huang: **Semi-Supervised Learning with Max-Margin Graph Cuts** AISTATS 2010

Calandriello, Lazaric, MV: **Distributed sequential sampling for kernel matrix approximation** AISTATS **2017**

Calandriello, Lazaric, MV: **Second-order kernel online convex optimization with adaptive sketching**, ICML 2017

Calandriello, Lazaric, MV: **Efficient second-order online kernel learning with adaptive embedding**, NIPS 2017

Calandriello, Koutis, Lazaric, MV: **Improved large-scale graph learning through ridge spectral sparsification**, ICML 2018

Calandriello, Carratino, Lazaric, MV, Rosasco: **Gaussian process optimization with adaptive sketching: Scalable and no regret**, COLT 2019 and [NEGDEP@ICML2019](#)

Dereziński\*, Calandriello\*, MV: **Exact sampling of determinantal point processes with sublinear time preprocessing**, [NEGDEP@ICML2019](#)

**code:** <http://researchers.lille.inria.fr/~valko/hp/publications/squeak.py>

# COMING UP..

---

- ▶ **Sparsification**
- ▶ **Resistance distance**
- ▶ **Leverage scores**
- ▶ **1-pass is a must**
- ▶ **Online leverage scores**
- ▶ **Negative dependence!**
- ▶ **SQUEAK**
- ▶ **Back to the beginning**
  - **Spectral sparsifiers**
- ▶ **Back to the future**
  - **GP-UCB & DPPs**

# JOINT WORK WITH...



**Alessandro  
Lazaric**  
FAIR Paris



**Ali Rahimi**  
Google  
Research



**Branislav  
Kveton**  
Google  
Research



**Daniel Ting**  
Tableau  
Research



**Daniele  
Calandriello**  
IIT, Genova



**Ling Huang**  
AHI Fintech



**Lorenzo  
Rosasco**  
IIT, Genova



**Luigi  
Carratino**  
IIT, Genova



**Michał  
Dereziński**  
UC Berkeley



**Yiannis  
Koutis**  
NJIT & CMU

## Laplacians and kernels

*Reproducing kernel Hilbert space\**

Vector space  $\mathcal{H}$  with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$

Feature map  $\varphi(\mathbf{x}) : \mathcal{X} \rightarrow \mathcal{H}$

Kernel function  $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \langle \mathcal{K}(\mathbf{x}, \cdot), \mathcal{K}(\mathbf{x}', \cdot) \rangle_{\mathcal{H}} = \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle_{\mathcal{H}}$

*Kernels evaluated at the dataset*

Features  $\varphi(\mathbf{x}_i) = \phi_i$

Kernel  $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \langle \varphi(\mathbf{x}_i), \varphi(\mathbf{x}_j) \rangle_{\mathcal{H}} = \phi_i^{\top} \phi_j$

Feature map  $\Phi_n = [\phi_1, \phi_2, \dots, \phi_n] : \mathbb{R}^n \rightarrow \mathcal{H}$

Empirical kernel matrix  $\mathbf{K}_n \in \mathbb{R}^{n \times n}$ , s.t.  $[\mathbf{K}]_{i,j} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$

Column  $\mathbf{k}_{[t-1],t} \in \mathbb{R}^{t-1} = \Phi_{t-1}^{\top} \phi_t$

Kernel at a point  $k_{i,i} \in \mathbb{R} = \phi_t^{\top} \phi_t$

\*Not entering into formal details

## Laplacians and kernels

*Reproducing kernel Hilbert space\**

Vector space  $\mathcal{H}$  with inner product  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$

Feature map  $\varphi(\mathbf{x}) : \mathcal{X} \rightarrow \mathcal{H}$

Kernel function  $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \langle \mathcal{K}(\mathbf{x}, \cdot), \mathcal{K}(\mathbf{x}', \cdot) \rangle_{\mathcal{H}} = \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle_{\mathcal{H}}$

*Kernels evaluated at the dataset*

Features  $\varphi(\mathbf{x}_i) = \phi_i$

Kernel  $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) = \langle \varphi(\mathbf{x}_i), \varphi(\mathbf{x}_j) \rangle_{\mathcal{H}} = \phi_i^{\top} \phi_j$

Feature map  $\Phi_n = [\phi_1, \phi_2, \dots, \phi_n] : \mathbb{R}^n \rightarrow \mathcal{H}$

Empirical kernel matrix  $\mathbf{K}_n \in \mathbb{R}^{n \times n}$ , s.t.  $[\mathbf{K}]_{i,j} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$

Column  $\mathbf{k}_{[t-1],t} \in \mathbb{R}^{t-1} = \Phi_{t-1}^{\top} \phi_t$

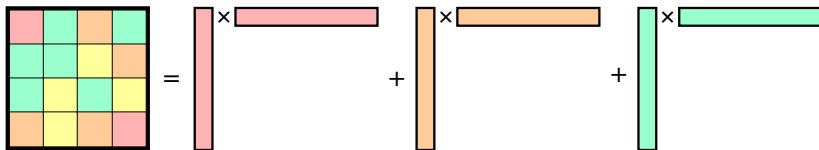
Kernel at a point  $k_{i,i} \in \mathbb{R} = \phi_t^{\top} \phi_t$

\*Not entering into formal details

# Part 1: Kernel Dictionary Learning - The Hammer

# Dictionary Learning

Covariance operator:  $\Phi_n \Phi_n^T = \sum_{i=1}^n \phi_i \phi_i^T$



*Dictionary learning*\*: find an **accurate** representation of the input data as a linear combination of a **small** set of basic elements (**atoms**)

\*other people may give other definitions...

## Singular Value Decomposition – Learning Atoms

SVD of  $\Phi_n = \mathbf{V}\Sigma\mathbf{U}^T$  (with rank  $r$ )\*

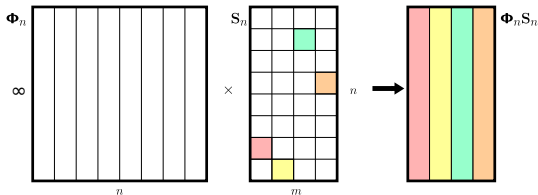
$$\Phi_n \Phi_n^T = \sum_{i=1}^n \phi_i \phi_i^T = \sum_{j=1}^r \sigma_j^2 \mathbf{v}_j \mathbf{v}_j^T$$

\*With kernels (e.g., Gaussian),  $r$  is often as large as  $n$



# Dataset Subsampling – Learning Weights

Dictionary  $\mathcal{I} = \{(w_j, \phi_j)\}_{j=1}^m$



$$\sum_{j=1}^m w_j \phi_j \phi_j^T = \sum_{j=1}^m (\sqrt{w_j} \phi_j)(\sqrt{w_j} \phi_j)^T = \Phi_n S_n S_n^T \Phi_n^T$$

which points? ( $\phi_j$ )  
how many? ( $m$ )

which weights? ( $w_j$ )  
which guarantees?

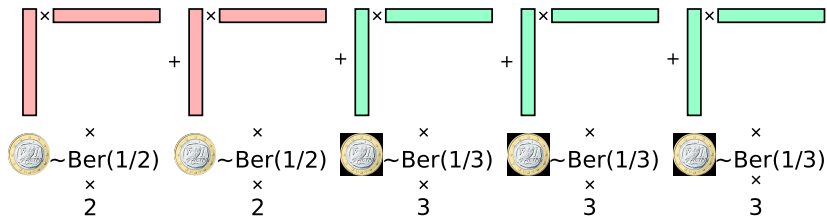
\*Remark: we do not reduce the vectors  $\phi_j$

## Nyström Sampling – Intuition

Sample points  $\mathbf{x}_i$  w.p.  $p_i$  and add it to  $\mathcal{I}$  with weight  $\propto 1/p_i$

# Nyström Sampling – Intuition

Sample points  $\mathbf{x}_i$  w.p.  $p_i$  and add it to  $\mathcal{I}$  with weight  $\propto 1/p_i$



## Nyström Sampling – Formally

**Input:** budget  $\bar{q}$ , probabilities  $\{p_i\}_i$  (*not necessarily normalized!*)

**Init:**  $\mathcal{I} = \emptyset$

**For all**  $i = 1, \dots, n$

Draw  $q_i \sim \mathcal{B}(p_i, \bar{q})$

Compute weight  $w_i = \frac{1}{p_i} \frac{q_i}{\bar{q}}$

Add  $(w_i, \mathbf{x}_i)$  to  $\mathcal{I}$

**Output:**  $\mathcal{I}$

$q_i$  may be seen as adding  $q_i$  copies of  $\mathbf{x}_i$  with weight  $1/(p_i \bar{q})$

## Nyström Sampling – Formally

### Lemma

The Nyström estimator ( $z_{i,j}$ : one out of  $\bar{q}$  Bernoulli trials of probability  $p_i$ )

$$\Phi_n \mathbf{S}_n \mathbf{S}_n^\top \Phi_n^\top = \sum_{i=1}^n \sum_{j=1}^{\bar{q}} \frac{1}{p_i} \frac{z_{i,j}}{\bar{q}} \phi_i \phi_i^\top$$

is unbiased

$$\mathbb{E}_{\mathbf{S}_n} \left[ \Phi_n \mathbf{S}_n \mathbf{S}_n^\top \Phi_n^\top \right] = \Phi_n \Phi_n^\top$$

and its dictionary has size

$$\mathbb{P} \left( |\mathcal{I}| \geq 3\bar{q} \sum_{i=1}^n p_i \right) \leq \exp \left( -\bar{q} \sum_{i=1}^n p_i \right)$$

E.g., uniform sampling  $p_i = 1/n$ ,  $|\mathcal{I}| \leq 3\bar{q}$  w.h.p.

## Nyström Sampling – Formally

### Lemma

The Nyström estimator ( $z_{i,j}$ : one out of  $\bar{q}$  Bernoulli trials of probability  $p_i$ )

$$\Phi_n \mathbf{S}_n \mathbf{S}_n^T \Phi_n^T = \sum_{i=1}^n \sum_{j=1}^{\bar{q}} \frac{1}{p_i} \frac{z_{i,j}}{\bar{q}} \phi_i \phi_i^T$$

is unbiased

$$\mathbb{E}_{\mathbf{S}_n} [\Phi_n \mathbf{S}_n \mathbf{S}_n^T \Phi_n^T] = \Phi_n \Phi_n^T$$

and its dictionary has size

$$\mathbb{P}\left(|\mathcal{I}| \geq 3\bar{q} \sum_{i=1}^n p_i\right) \leq \exp\left(-\bar{q} \sum_{i=1}^n p_i\right)$$

E.g., uniform sampling  $p_i = 1/n$ ,  $|\mathcal{I}| \leq 3\bar{q}$  w.h.p.

But is the approximate covariance good?

## Reconstruction Guarantees

An  $(\varepsilon, \gamma)$ -accurate dictionary  $\mathcal{I}$  satisfies\*

$$\overbrace{(1 - \varepsilon)\Phi_n\Phi_n^T}^{\text{multiplicative error}} - \overbrace{\varepsilon\gamma\mathbf{I}}^{\text{additive error}} \preceq \Phi_n\mathbf{S}\mathbf{S}^T\Phi_n^T \preceq \overbrace{(1 + \varepsilon)\Phi_n\Phi_n^T}^{\text{multiplicative error}} + \overbrace{\varepsilon\gamma\mathbf{I}}^{\text{additive error}}$$

### Remarks

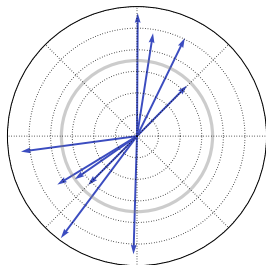
If  $\gamma = 0$ , all spectrum of  $\Phi_n\Phi_n^T$  is preserved up to  $1 \pm \varepsilon$  multiplicative error

If  $\gamma > 0$  only eigenvalues larger than  $\varepsilon\gamma$  are preserved

\*If a dictionary is accurate in this sense, then it is accurate to build many other things

## Nyström Sampling Guarantees – Intuition

Uniform sampling  $p_i = 1/n$



$$\Phi_n \Phi_n^T - \Phi_n \mathbf{S} \mathbf{S}^T \Phi_n^T = \sum_{i=1}^n \phi_i \phi_i^T - \sum_{j=1}^m w_j \phi_j \phi_j^T$$

“Important” directions may have probability **too small** to be selected

“Redundant” directions may have probability **too large** to be selected



## Ridge Leverage Scores\*

$$\tau_{n,i} = \mathbf{e}_{n,i} \mathbf{K}_n^\top (\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1} \mathbf{e}_{n,i} = \phi_i^\top (\Phi_n \Phi_n^\top + \gamma \mathbf{I})^{-1} \phi_i$$

\*leverage scores evaluate the “relevance” of a point in statistics

## Ridge Leverage Scores\*

$$\tau_{n,i} = \mathbf{e}_{n,i} \mathbf{K}_n^T (\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1} \mathbf{e}_{n,i} = \phi_i^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1} \phi_i$$

RLS capture “soft” orthogonality

- ▶ If all  $\phi_i$  are *orthogonal*

$$\tau_{n,i} = \phi_i^T (\phi_i \phi_i^T + \gamma \mathbf{I})^{-1} \phi_i = \frac{\phi_i^T \phi_i}{\phi_i^T \phi_i + \gamma} \sim \mathbf{1}$$

- ▶ If all  $\phi_i$  are *collinear*

$$\tau_{n,i} = \phi_i^T (n \phi_i \phi_i^T + \gamma \mathbf{I})^{-1} \phi_i = \frac{\phi_i^T \phi_i}{n \phi_i^T \phi_i + \gamma} \sim \frac{1}{n}$$

\*leverage scores evaluate the “relevance” of a point in statistics

## Ridge Leverage Scores\*

$$\tau_{n,i} = \mathbf{e}_{n,i} \mathbf{K}_n^T (\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1} \mathbf{e}_{n,i} = \phi_i^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1} \phi_i$$

RLS capture “soft” orthogonality

- ▶ If all  $\phi_i$  are *orthogonal*

$$\tau_{n,i} = \phi_i^T (\phi_i \phi_i^T + \gamma \mathbf{I})^{-1} \phi_i = \frac{\phi_i^T \phi_i}{\phi_i^T \phi_i + \gamma} \sim \mathbf{1}$$

- ▶ If all  $\phi_i$  are *collinear*

$$\tau_{n,i} = \phi_i^T (n \phi_i \phi_i^T + \gamma \mathbf{I})^{-1} \phi_i = \frac{\phi_i^T \phi_i}{n \phi_i^T \phi_i + \gamma} \sim \frac{1}{n}$$

Given  $\Phi_{t-1}$ , adding columns reduce previous RLS

$$\tau_{t,i} \leq \tau_{t-1,i}$$

\*leverage scores evaluate the “relevance” of a point in statistics

## Ridge Leverage Scores\*

$$\tau_{n,i} = \mathbf{e}_{n,i} \mathbf{K}_n^T (\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1} \mathbf{e}_{n,i} = \phi_i^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1} \phi_i$$

RLS capture “soft” orthogonality

- ▶ If all  $\phi_i$  are *orthogonal*

$$\tau_{n,i} = \phi_i^T (\phi_i \phi_i^T + \gamma \mathbf{I})^{-1} \phi_i = \frac{\phi_i^T \phi_i}{\phi_i^T \phi_i + \gamma} \sim \mathbf{1}$$

- ▶ If all  $\phi_i$  are *collinear*

$$\tau_{n,i} = \phi_i^T (n \phi_i \phi_i^T + \gamma \mathbf{I})^{-1} \phi_i = \frac{\phi_i^T \phi_i}{n \phi_i^T \phi_i + \gamma} \sim \frac{1}{n}$$

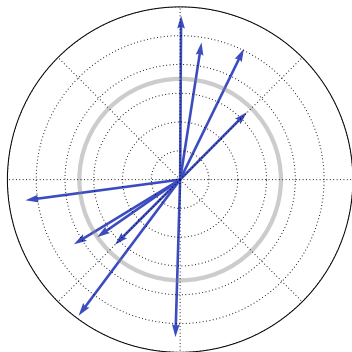
Given  $\Phi_{t-1}$ , adding columns reduce previous RLS

$$\tau_{t,i} \leq \tau_{t-1,i}$$

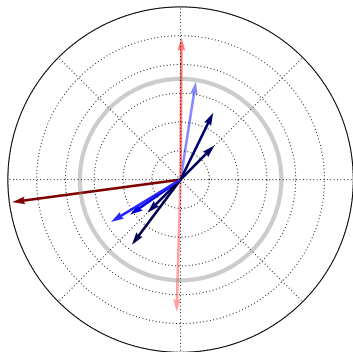
RLS decrease with  $\gamma$

\*leverage scores evaluate the “relevance” of a point in statistics

## Ridge Leverage Scores

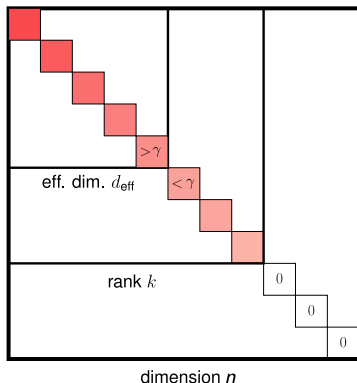


## Ridge Leverage Scores



## Effective Dimension

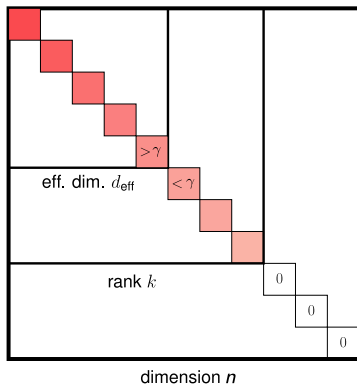
The **effective dimension** is the **number of relevant directions in the data**



$$d_{\text{eff}}^n(\gamma) = \sum_{i=1}^n \tau_{n,i} = \text{Tr}(\mathbf{K}_n(\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1}) = \sum_{i=1}^n \frac{\lambda_i(\mathbf{K}_n)}{\lambda_i(\mathbf{K}_n) + \gamma} \leq \text{Rank}(\mathbf{K}_n)$$

## Effective Dimension

The **effective dimension** is the **number of relevant directions in the data**



Given  $d_{\text{eff}}^{t-1}(\gamma)$ , adding a new column to  $\Phi_{t-1}$  may increase  $d_{\text{eff}}^t(\gamma)$

$$\mathbf{d}_{\text{eff}}^t(\gamma) \geq \mathbf{d}_{\text{eff}}^{t-1}(\gamma)$$



## Nyström Sampling – Formally

**Input:** budget  $\bar{q}$ , probabilities  $\{p_i\}_i$  (*not necessarily normalized!*)

**Init:**  $\mathcal{I} = \emptyset$

**For all**  $i = 1, \dots, n$

Set  $p_i = \tau_{n,i}$

Draw  $q_i \sim \mathcal{B}(p_i, \bar{q})$

Compute weight  $w_i = \frac{1}{p_i} \frac{q_i}{\bar{q}}$

Add  $(w_i, \mathbf{x}_i)$  to  $\mathcal{I}$

**Output:**  $\mathcal{I}$

$q_i$  may be seen as adding  $q_i$  copies of  $\mathbf{x}_i$  with weight  $1/(p_i \bar{q})$

## Oracle RLS Sampling

### Theorem (Alaoui and Mahoney, 2014)

Consider the Nyström estimator with oracle RLS sampling  $p_i = \tau_{n,i}$ . If

$$\bar{q} \geq \frac{4 \log(n/\delta)}{\varepsilon^2}$$

then  $\mathcal{I}$  is an  $(\varepsilon, \gamma)$ -accurate dictionary w.p.  $1 - \delta$  and

$$|\mathcal{I}| \leq 3\bar{q}d_{\text{eff}}^n(\gamma)$$

## Oracle RLS Sampling

### Theorem (Alaoui and Mahoney, 2014)

Consider the Nyström estimator with oracle RLS sampling  $p_i = \tau_{n,i}$ . If

$$\bar{q} \geq \frac{4 \log(n/\delta)}{\varepsilon^2}$$

then  $\mathcal{I}$  is an  $(\varepsilon, \gamma)$ -accurate dictionary w.p.  $1 - \delta$  and

$$|\mathcal{I}| \leq 3\bar{q}d_{\text{eff}}^n(\gamma)$$

Small and accurate dictionary adapting to the “complexity” of the data

$$d_{\text{eff}}^n(\gamma) = \sum_{i=1}^n \tau_{i,n} \ll n\tau_{\max}$$

Given the RLS as input

## Oracle RLS Sampling

### Theorem (Alaoui and Mahoney, 2014)

Consider the Nyström estimator with oracle RLS sampling  $p_i = \tau_{n,i}$ . If

$$\bar{q} \geq \frac{4 \log(n/\delta)}{\varepsilon^2}$$

then  $\mathcal{I}$  is an  $(\varepsilon, \gamma)$ -accurate dictionary w.p.  $1 - \delta$  and

$$|\mathcal{I}| \leq 3\bar{q}d_{\text{eff}}^n(\gamma)$$

Small and accurate dictionary adapting to the “complexity” of the data

$$d_{\text{eff}}^n(\gamma) = \sum_{i=1}^n \tau_{i,n} \ll n\tau_{\max}$$

Given the RLS as input

Computing  $\tau_{n,i} = \mathbf{e}_{n,i} \mathbf{K}_n^T (\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1} \mathbf{e}_{n,i}$  requires storing and inverting  $\mathbf{K}_n$

## Estimating RLS from a Dictionary

Approximate the kernel matrix directly

$$\tau_{n,i} = \mathbf{e}_{n,i} \mathbf{K}_n^T (\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1} \mathbf{e}_{n,i}$$

$$\tilde{\tau}_{n,i} = \mathbf{e}_i^T \tilde{\mathbf{K}}_n (\tilde{\mathbf{K}}_n + \gamma \mathbf{I})^{-1} \mathbf{e}_i$$

## Estimating RLS from a Dictionary

Approximate the kernel matrix directly

$$\tau_{n,i} = \mathbf{e}_{n,i} \mathbf{K}_n^T (\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1} \mathbf{e}_{n,i}$$

$$\tilde{\tau}_{n,i} = \mathbf{e}_i^T \tilde{\mathbf{K}}_n (\tilde{\mathbf{K}}_n + \gamma \mathbf{I})^{-1} \mathbf{e}_i$$

Instead, approximate  $\tau_{n,i}$  directly in  $\mathcal{H}$

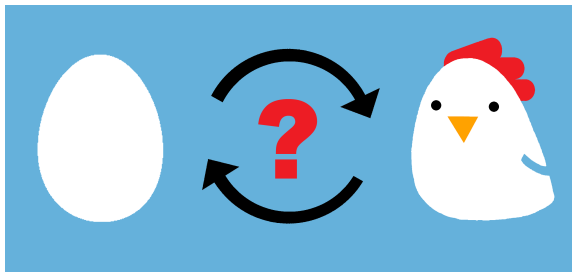
$$\tau_{n,i} = \phi_i^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1} \phi_i$$

$$\tilde{\tau}_{n,i} = \phi_i^T (\Phi_n \mathbf{S}_n \mathbf{S}_n^T \Phi_n^T + \gamma \mathbf{I})^{-1} \phi_i$$

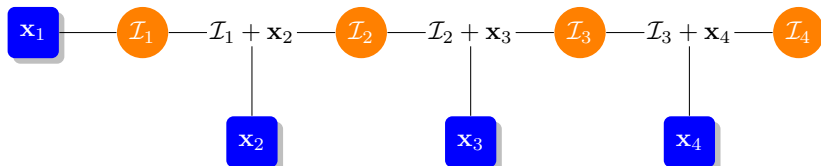
# Chicken and Egg problem

Given accurate  $\tilde{\tau}_{n,i}$   $\Rightarrow$  compute accurate dictionary

Given accurate dictionary  $\Rightarrow$  compute accurate  $\tilde{\tau}_{n,i}$



## Sequential RLS Sampling – Intuition





# SQUEAK

Dictionary  $\mathcal{I}_t = \{(j, \phi_j, q_{t,j}, \tilde{p}_{t,j})\}$ , weights  $w_i = \frac{q_{t,j}}{\tilde{p}_{t,j} \bar{q}}$

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

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

# SQUEAK

## Theorem

Consider the Nyström estimator built using SQUEAK . If

$$\bar{q} \geq \frac{4\alpha \log(n/\delta)}{\varepsilon^2} \quad \text{where } \alpha = \left(\frac{1+\varepsilon}{1-\varepsilon}\right),$$

then for all  $t = 1, \dots, n$   $\mathcal{I}_t$  is an  $(\varepsilon, \gamma)$ -accurate dictionary w.p.  $1 - \delta$  and

$$|\mathcal{I}_t| \leq 3\bar{q}d_{\text{eff}}^t(\gamma)$$

# SQUEAK

## Theorem

Consider the Nyström estimator built using SQUEAK . If

$$\bar{q} \geq \frac{4\alpha \log(n/\delta)}{\varepsilon^2} \quad \text{where } \alpha = \left(\frac{1+\varepsilon}{1-\varepsilon}\right),$$

then for all  $t = 1, \dots, n$   $\mathcal{I}_t$  is an  $(\varepsilon, \gamma)$ -accurate dictionary w.p.  $1 - \delta$  and

$$|\mathcal{I}_t| \leq 3\bar{q}d_{\text{eff}}^t(\gamma)$$

- ▶ Accuracy and space/time guarantees
- ▶ Anytime guarantees
- ▶ In worst case, no space gain (stores full  $\mathbf{K}_n$ )
- ▶ In worst case, no space overhead (stores full  $\mathbf{K}_n$ )
- ▶ RLS estimator not incremental, not easy because of changing weights
- ▶ Unnormalized  $\tilde{\rho}_{t,i}$

# SQUEAK

## Theorem

Consider the Nyström estimator built using SQUEAK . If

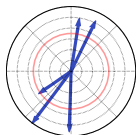
$$\bar{q} \geq \frac{4\alpha \log(n/\delta)}{\varepsilon^2} \quad \text{where } \alpha = \left(\frac{1+\varepsilon}{1-\varepsilon}\right),$$

then for all  $t = 1, \dots, n$   $\mathcal{I}_t$  is an  $(\varepsilon, \gamma)$ -accurate dictionary w.p.  $1 - \delta$  and

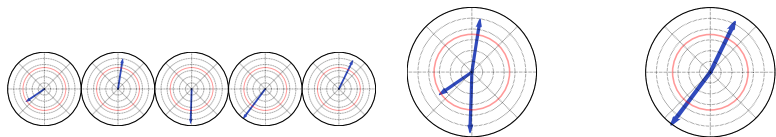
$$|\mathcal{I}_t| \leq 3\bar{q}d_{\text{eff}}^t(\gamma)$$

- ▶ Only need to compute  $\tilde{\tau}_{t,i}$  if  $i \in \mathcal{I}_t$ , never recompute after dropping
  - ↳ Never construct the whole  $\mathbf{K}_n$
  - ↳ subquadratic runtime  $\mathcal{O}(n^3) \Rightarrow \mathcal{O}(n|\mathcal{I}_n|^3) \leq \tilde{\mathcal{O}}(nd_{\text{eff}}^n(\gamma)^3)$
- ▶ Store points directly in the dictionary
  - ↳  $\tilde{\mathcal{O}}(\mathbf{d}_{\text{eff}}^n(\gamma)^2 + \mathbf{d}_{\text{eff}}^n(\gamma)\mathbf{D})$  space “constant” in  $n$
  - ↳ single pass over the dataset (streaming)

# Sequential RLS sampling – Distributed Version

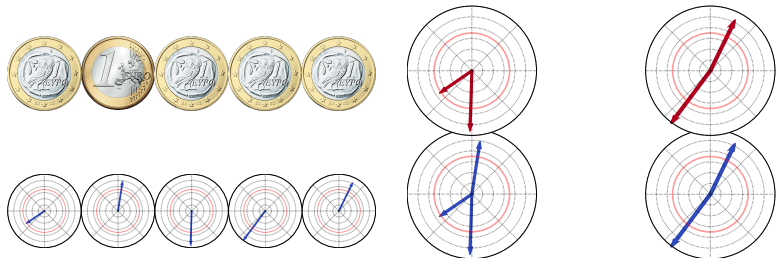


## Sequential RLS sampling – Distributed Version



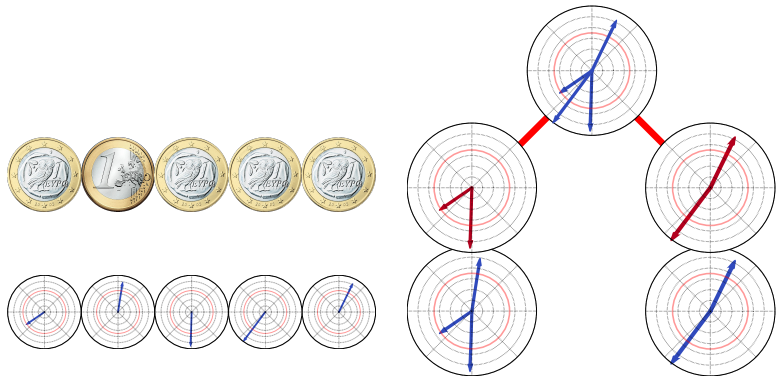
## Sequential RLS sampling – Distributed Version

$$\tilde{p}_{1,i} \propto \tilde{\tau}_{1,i},$$
$$z_{1,i} = \mathbb{I}\{\text{Ber}(\tilde{p}_{1,i})\}$$



## Sequential RLS sampling – Distributed Version

$$\tilde{p}_{1,i} \propto \tilde{\tau}_{1,i},$$
$$z_{1,i} = \mathbb{I}\{\text{Ber}(\tilde{p}_{1,i})\}$$

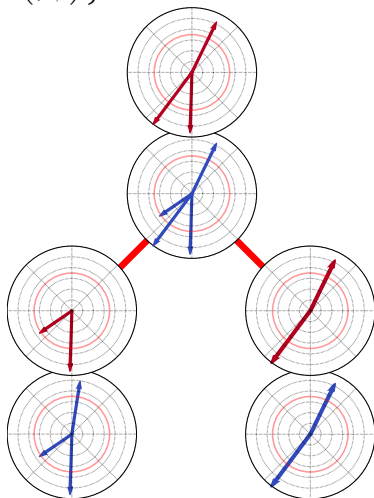
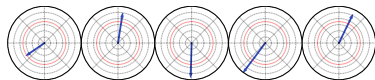




## Sequential RLS sampling – Distributed Version

$$\tilde{p}_{1,i} \propto \tilde{\tau}_{1,i},$$
$$z_{1,i} = \mathbb{I}\{\text{Ber}(\tilde{p}_{1,i})\}$$

$$\tilde{p}_{2,i} \propto \tilde{\tau}_{2,i}$$
$$z_{2,i} = \mathbb{I}\left\{\text{Ber}\left(\frac{\tilde{p}_{2,i}}{\tilde{p}_{1,i}}\right)\right\} z_{1,i}$$

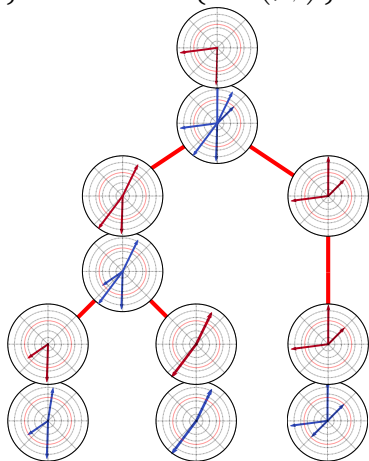


## Sequential RLS sampling – Distributed Version

$$\begin{aligned}\tilde{p}_{1,i} &\propto \tilde{\tau}_{1,i}, \\ z_{1,i} &= \mathbb{I}\{\text{Ber}(\tilde{p}_{1,i})\}\end{aligned}$$

$$\begin{aligned}\tilde{p}_{2,i} &\propto \tilde{\tau}_{2,i} \\ z_{2,i} &= \mathbb{I}\left\{\text{Ber}\left(\frac{\tilde{p}_{2,i}}{\tilde{p}_{1,i}}\right)\right\} z_{1,i}\end{aligned}$$

$$\begin{aligned}\tilde{p}_{3,i} &\propto \tilde{\tau}_{3,i} \\ z_{3,i} &= \mathbb{I}\left\{\text{Ber}\left(\frac{\tilde{p}_{3,i}}{\tilde{p}_{2,i}}\right)\right\} z_{2,i}\end{aligned}$$



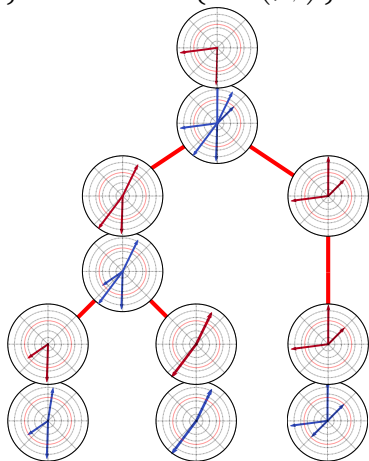
## Sequential RLS sampling – Distributed Version

$$\begin{aligned}\tilde{p}_{1,i} &\propto \tilde{\tau}_{1,i}, \\ z_{1,i} &= \mathbb{I}\{\text{Ber}(\tilde{p}_{1,i})\}\end{aligned}$$

$$\begin{aligned}\tilde{p}_{2,i} &\propto \tilde{\tau}_{2,i}, \\ z_{2,i} &= \mathbb{I}\left\{\text{Ber}\left(\frac{\tilde{p}_{2,i}}{\tilde{p}_{1,i}}\right)\right\} z_{1,i}\end{aligned}$$

$$\begin{aligned}\tilde{p}_{3,i} &\propto \tilde{\tau}_{3,i}, \\ z_{3,i} &= \mathbb{I}\left\{\text{Ber}\left(\frac{\tilde{p}_{3,i}}{\tilde{p}_{2,i}}\right)\right\} z_{2,i}\end{aligned}$$

- ▶ Dataset is distributed over multiple machines



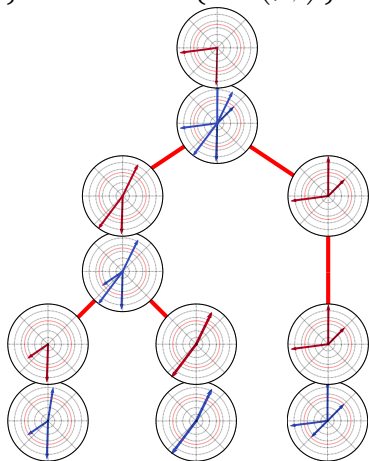
## Sequential RLS sampling – Distributed Version

$$\tilde{p}_{1,i} \propto \tilde{\tau}_{1,i}, \\ z_{1,i} = \mathbb{I}\{\text{Ber}(\tilde{p}_{1,i})\}$$

$$\tilde{p}_{2,i} \propto \tilde{\tau}_{2,i}, \\ z_{2,i} = \mathbb{I}\left\{\text{Ber}\left(\frac{\tilde{p}_{2,i}}{\tilde{p}_{1,i}}\right)\right\} z_{1,i}$$

$$\tilde{p}_{3,i} \propto \tilde{\tau}_{3,i}, \\ z_{3,i} = \mathbb{I}\left\{\text{Ber}\left(\frac{\tilde{p}_{3,i}}{\tilde{p}_{2,i}}\right)\right\} z_{2,i}$$

- ▶ Dataset is distributed over multiple machines
- ▶ Communication is limited to samples in the dictionaries



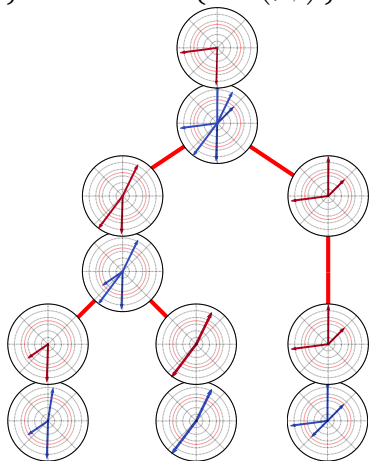
## Sequential RLS sampling – Distributed Version

$$\tilde{p}_{1,i} \propto \tilde{\tau}_{1,i}, \\ z_{1,i} = \mathbb{I}\{\text{Ber}(\tilde{p}_{1,i})\}$$


$$\tilde{p}_{2,i} \propto \tilde{\tau}_{2,i}, \\ z_{2,i} = \mathbb{I}\left\{\text{Ber}\left(\frac{\tilde{p}_{2,i}}{\tilde{p}_{1,i}}\right)\right\} z_{1,i}$$

$$\tilde{p}_{3,i} \propto \tilde{\tau}_{3,i}, \\ z_{3,i} = \mathbb{I}\left\{\text{Ber}\left(\frac{\tilde{p}_{3,i}}{\tilde{p}_{2,i}}\right)\right\} z_{2,i}$$

- ▶ Dataset is distributed over multiple machines
- ▶ Communication is limited to samples in the dictionaries
- ▶ Runtime depends on merge tree




## Comparison

 = oracle,  $\mu(\gamma) = \max_i \tau_{n,i}(\gamma) \leq 1/\gamma$  regularized coherence

	$\tilde{\mathcal{O}}(\text{Runtime})$	$\mathcal{O}( \mathcal{I}_n )$	Passes
Bach, 2013 (Uniform)	$n\mu(\gamma) + \text{oracle}$	$n\mu(\gamma)$	1

## Comparison

 = oracle,  $\mu(\gamma) = \max_i \tau_{n,i}(\gamma) \leq 1/\gamma$  regularized coherence

	$\tilde{\mathcal{O}}(\text{Runtime})$	$\mathcal{O}( \mathcal{I}_n )$	Passes
Bach, 2013 (Uniform)	$n\mu(\gamma) + \text{oracle}$	$n\mu(\gamma)$	1
Oracle RLS sampling	$n + \text{oracle}$	$d_{\text{eff}}^n(\gamma) \log(n)$	Many


## Comparison

 = oracle,  $\mu(\gamma) = \max_i \tau_{n,i}(\gamma) \leq 1/\gamma$  regularized coherence

	$\tilde{\mathcal{O}}(\text{Runtime})$	$\mathcal{O}( \mathcal{I}_n )$	Passes
Bach, 2013 (Uniform)	$n\mu(\gamma) + \text{oracle}$	$n\mu(\gamma)$	1
Oracle RLS sampling	$n + \text{oracle}$	$d_{\text{eff}}^n(\gamma) \log(n)$	Many
Exact RLS sampling	$n^3$	$d_{\text{eff}}^n(\gamma) \log(n)$	Many



## Comparison

 = oracle,  $\mu(\gamma) = \max_i \tau_{n,i}(\gamma) \leq 1/\gamma$  regularized coherence


	$\tilde{\mathcal{O}}(\text{Runtime})$	$\mathcal{O}( \mathcal{I}_n )$	Passes
Bach, 2013 (Uniform)	$n\mu(\gamma) + \text{oracle}$	$n\mu(\gamma)$	1
Oracle RLS sampling	$n + \text{oracle}$	$d_{\text{eff}}^n(\gamma) \log(n)$	Many
Exact RLS sampling	$n^3$	$d_{\text{eff}}^n(\gamma) \log(n)$	Many
Alaoui and Mahoney, 2015	$n^3 \mu(\gamma)^2$	$n\mu(\gamma) + d_{\text{eff}}^n(\gamma) \log(n)$	3

## Comparison

 = oracle,  $\mu(\gamma) = \max_i \tau_{n,i}(\gamma) \leq 1/\gamma$  regularized coherence


	$\tilde{\mathcal{O}}(\text{Runtime})$	$\mathcal{O}( \mathcal{I}_n )$	Passes
Bach, 2013 (Uniform)	$n\mu(\gamma) + \text{oracle}$	$n\mu(\gamma)$	1
Oracle RLS sampling	$n + \text{oracle}$	$d_{\text{eff}}^n(\gamma) \log(n)$	Many
Exact RLS sampling	$n^3$	$d_{\text{eff}}^n(\gamma) \log(n)$	Many
Alaoui and Mahoney, 2015	$n^3 \mu(\gamma)^2$	$n\mu(\gamma) + d_{\text{eff}}^n(\gamma) \log(n)$	3
SQUEAK Calandriello et al., 2017a	$(n/k)d_{\text{eff}}^n(\gamma)^3$	$d_{\text{eff}}^n(\gamma) \log(n)$	1

## Comparison

 = oracle,  $\mu(\gamma) = \max_i \tau_{n,i}(\gamma) \leq 1/\gamma$  regularized coherence

	$\tilde{\mathcal{O}}(\text{Runtime})$	$\mathcal{O}( \mathcal{I}_n )$	Passes
Bach, 2013 (Uniform)	$n\mu(\gamma) + \text{Oracle}$	$n\mu(\gamma)$	1
Oracle RLS sampling	$n + \text{Oracle}$	$d_{\text{eff}}^n(\gamma) \log(n)$	Many
Exact RLS sampling	$n^3$	$d_{\text{eff}}^n(\gamma) \log(n)$	Many
Alaoui and Mahoney, 2015	$n^3 \mu(\gamma)^2$	$n\mu(\gamma) + d_{\text{eff}}^n(\gamma) \log(n)$	3
SQUEAK Calandriello et al., 2017a	$(n/k)d_{\text{eff}}^n(\gamma)^3$	$d_{\text{eff}}^n(\gamma) \log(n)$	1
KORS Calandriello et al., 2017b	$nd_{\text{eff}}^n(\gamma)^2$	$d_{\text{eff}}^n(\gamma) \log^2(n)$	1

## Comparison

 = oracle,  $\mu(\gamma) = \max_i \tau_{n,i}(\gamma) \leq 1/\gamma$  regularized coherence

	$\tilde{\mathcal{O}}(\text{Runtime})$	$\mathcal{O}( \mathcal{I}_n )$	Passes
Bach, 2013 (Uniform)	$n\mu(\gamma) + \text{Oracle}$	$n\mu(\gamma)$	1
Oracle RLS sampling	$n + \text{Oracle}$	$d_{\text{eff}}^n(\gamma) \log(n)$	Many
Exact RLS sampling	$n^3$	$d_{\text{eff}}^n(\gamma) \log(n)$	Many
Alaoui and Mahoney, 2015	$n^3 \mu(\gamma)^2$	$n\mu(\gamma) + d_{\text{eff}}^n(\gamma) \log(n)$	3
SQUEAK Calandriello et al., 2017a	$(n/k)d_{\text{eff}}^n(\gamma)^3$	$d_{\text{eff}}^n(\gamma) \log(n)$	1
KORS Calandriello et al., 2017b	$nd_{\text{eff}}^n(\gamma)^2$	$d_{\text{eff}}^n(\gamma) \log^2(n)$	1
Musco and Musco, 2017	$nd_{\text{eff}}^n(\gamma)^2$	$d_{\text{eff}}^n(\gamma) \log(n)$	$\log(n)$

## Recap

Construct a small, provably accurate dictionary in near-linear time

SQUEAK and DISQUEAK

Sub-linear time using multiple machines

Final dictionary can be updated if new samples arrive

## Recap

Construct a small, provably accurate dictionary in near-linear time

SQUEAK and DISQUEAK

Sub-linear time using multiple machines

Final dictionary can be updated if new samples arrive

Novel analysis, potentially useful for general importance sampling

## Recap

Construct a small, provably accurate dictionary in near-linear time

SQUEAK and DISQUEAK

Sub-linear time using multiple machines

Final dictionary can be updated if new samples arrive

Novel analysis, potentially useful for general importance sampling

Future work

Experiments

- ↳ Easy to implement: distributed task queue
- Preliminary results promising, easily scales to 1M+ samples

## Recap

Construct a small, provably accurate dictionary in near-linear time

SQUEAK and DISQUEAK

Sub-linear time using multiple machines

Final dictionary can be updated if new samples arrive

Novel analysis, potentially useful for general importance sampling

Future work

Experiments

- ↳ Easy to implement: distributed task queue
- Preliminary results promising, easily scales to 1M+ samples

Beyond passive processing: SQUEAK for active learning



## Part 2: Applications - The Nails

# Kernel Regression

*Kernel ridge regression*

$$\hat{\omega}_n = \arg \min_{\omega} \|\mathbf{y}_n - \mathbf{K}_n \omega\|^2 + \lambda \|\omega\|^2 = (\mathbf{K}_n + \lambda \mathbf{I})^{-1} \mathbf{y}_n$$

*Regularized Nyström kernel approximation*

$$\tilde{\mathbf{K}}_n = \mathbf{K}_n \mathbf{S}_n (\mathbf{S}_n^T \mathbf{K}_n \mathbf{S}_n + \gamma \mathbf{I}_{\mathcal{I}_n})^{-1} \mathbf{S}_n^T \mathbf{K}_n = \Phi_n^T \Phi_n \mathbf{S}_n (\mathbf{S}_n^T \mathbf{K}_n \mathbf{S}_n + \gamma \mathbf{I}_{\mathcal{I}_n})^{-1} \mathbf{S}_n^T \Phi_n^T \Phi_n$$

$$\begin{aligned} \tilde{\omega}_n &= (\tilde{\mathbf{K}}_n + \lambda \mathbf{I}_n)^{-1} \mathbf{y}_n \\ &= \frac{1}{\lambda} (\mathbf{y}_n - \mathbf{K}_n \mathbf{S}_n (\mathbf{S}_n^T \mathbf{K}_n \mathbf{S}_n + \lambda (\mathbf{S}_n^T \mathbf{K}_n \mathbf{S}_n + \gamma \mathbf{I}_{\mathcal{I}_n}))^{-1} \mathbf{S}_n^T \mathbf{K}_n \mathbf{y}_n) \end{aligned}$$

*Efficient computation*

- ▶ Construct the matrix  $\mathcal{O}(n|\mathcal{I}_n|^2)$
- ▶ Invert the matrix  $\mathcal{O}(|\mathcal{I}_n|^3)$
- ▶ Time  $\mathcal{O}(n^3) \Rightarrow \mathcal{O}(n|\mathcal{I}_n|^2 + |\mathcal{I}_n|^3)$
- ▶ Space  $\mathcal{O}(n^2) \Rightarrow \mathcal{O}(n|\mathcal{I}_n|)$

# Kernel Regression

## Theorem (Alaoui and Mahoney, 2014)

Consider the regularized Nyström kernel approximation generated by an  $(\varepsilon, \gamma)$ -accurate dictionary. Then

$$\mathbf{0} \preceq \mathbf{K}_n - \tilde{\mathbf{K}}_n \preceq \frac{\gamma}{1 - \varepsilon} \mathbf{K}_n (\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1} \preceq \frac{\gamma}{1 - \varepsilon} \mathbf{I}_n.$$

and

$$\mathcal{R}_{\mathcal{D}}(\tilde{\omega}) \leq \left(1 + \frac{\gamma}{\lambda} \frac{\varepsilon}{1 - \varepsilon}\right)^2 \mathcal{R}_{\mathcal{D}}(\hat{\omega}),$$

If  $\gamma = \lambda$  (i.e., additive error of the same order of the regularization)

- ▶ SQUEAK can be used to compute  $\tilde{\omega}$  in  $\mathcal{O}(nd_{\text{eff}}^n(\lambda)^2 + d_{\text{eff}}^n(\lambda)^3)$  time
- ▶ with a prediction error  $1/(1 - \varepsilon)^2$  larger than the exact solution

## Online Kernel Learning (OKL)

**Online** game between learner and adversary, at each round  $t \in [T]$

- 1 the **adversary** reveals a new point  $\varphi(\mathbf{x}_t) = \phi_t \in \mathcal{H}$
- 2 the learner chooses a function  $f_{\mathbf{w}_t}$  and predicts  $f_{\mathbf{w}_t}(\mathbf{x}_t) = \varphi(\mathbf{x}_t)^\top \mathbf{w}_t$ ,
- 3 the adversary reveals the **curved** loss  $\ell_t$ ,
- 4 the learner suffers  $\ell_t(\phi_t^\top \mathbf{w}_t)$  and observes the associated gradient  $\mathbf{g}_t$ .

# Online Kernel Learning (OKL)

**Online** game between learner and adversary, at each round  $t \in [T]$

- 1 the **adversary** reveals a new point  $\varphi(\mathbf{x}_t) = \phi_t \in \mathcal{H}$
- 2 the learner chooses a function  $f_{\mathbf{w}_t}$  and predicts  $f_{\mathbf{w}_t}(\mathbf{x}_t) = \varphi(\mathbf{x}_t)^\top \mathbf{w}_t$ ,
- 3 the adversary reveals the **curved** loss  $\ell_t$ ,
- 4 the learner suffers  $\ell_t(\phi_t^\top \mathbf{w}_t)$  and observes the associated gradient  $\mathbf{g}_t$ .

**Kernel** flexible but **curse of kernelization**

$t$  parameters  $\Rightarrow \mathcal{O}(t)$  per-step prediction cost

$$\mathbf{g}_t = \ell'_t(\phi_t^\top \mathbf{w}_t) \phi_t := \dot{g}_t \phi_t$$

# Online Kernel Learning (OKL)

**Online** game between learner and adversary, at each round  $t \in [T]$

- 1 the **adversary** reveals a new point  $\varphi(\mathbf{x}_t) = \phi_t \in \mathcal{H}$
- 2 the learner chooses a function  $f_{\mathbf{w}_t}$  and predicts  $f_{\mathbf{w}_t}(\mathbf{x}_t) = \varphi(\mathbf{x}_t)^\top \mathbf{w}_t$ ,
- 3 the adversary reveals the **curved** loss  $\ell_t$ ,
- 4 the learner suffers  $\ell_t(\phi_t^\top \mathbf{w}_t)$  and observes the associated gradient  $\mathbf{g}_t$ .

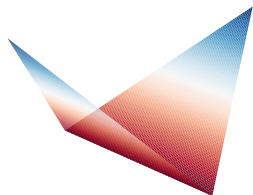
**Kernel** flexible but **curse of kernelization**

$t$  parameters  $\Rightarrow \mathcal{O}(t)$  per-step prediction cost

$$\mathbf{g}_t = \ell'_t(\phi_t^\top \mathbf{w}_t) \phi_t := \dot{g}_t \phi_t$$

**Learning** to minimize **regret**  $R(\mathbf{w}) = \sum_{t=1}^T \ell_t(\phi_t^\top \mathbf{w}_t) - \ell_t(\phi_t^\top \mathbf{w}^*)$   
and **compete** with **best-in-hindsight**  $\mathbf{w}^* := \arg \min_{\mathbf{w} \in \mathcal{H}} \sum_{t=1}^T \ell_t(\phi_t \mathbf{w})$

## OGD and losses

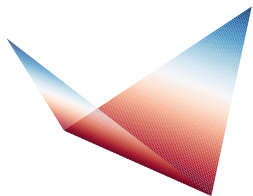


**convex**

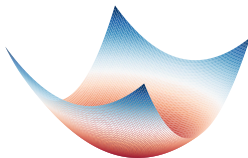
First order (GD) [Kivinen et al., 2004; Zinkevich, 2003]

$\sqrt{T}$  regret,  $\mathcal{O}(d)/\mathcal{O}(t)$  time/space per-step

## OGD and losses



**convex**



**strongly convex**

First order (GD) [Kivinen et al., 2004; Zinkevich, 2003]

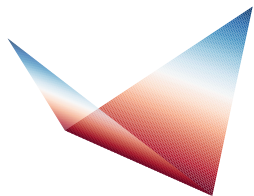
$\sqrt{T}$  regret,  $\mathcal{O}(d)/\mathcal{O}(t)$  time/space per-step

First order (GD) [Hazan et al., 2008]

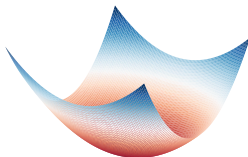
$\log(T)$  regret,



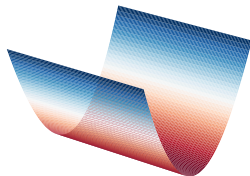
## OGD and losses



convex



strongly convex



First order (GD) [Kivinen et al., 2004; Zinkevich, 2003]

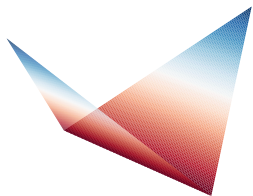
$\sqrt{T}$  regret,  $\mathcal{O}(d)/\mathcal{O}(t)$  time/space per-step

First order (GD) [Hazan et al., 2008]

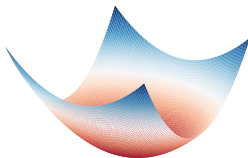
$\log(T)$  regret, but often not satisfied in practice

↳ (e.g.  $(y_t - \phi_t^T \mathbf{w}_t)^2$ )

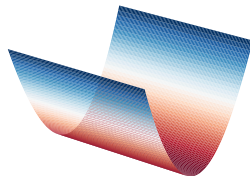
## OGD and losses



convex



strongly convex

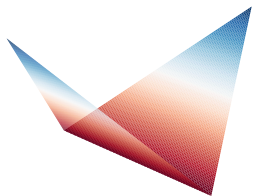


$\sigma$ -curved

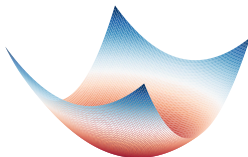
Second order (Newton-like) [Hazan et al., 2006; Zhdanov and Kalnishkan, 2010]

$\log(T)$  regret,  $\mathcal{O}(d^2)/\mathcal{O}(t^2)$  time/space per-step

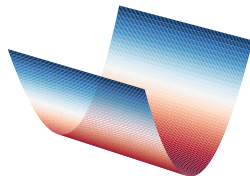
## OGD and losses



convex



strongly convex



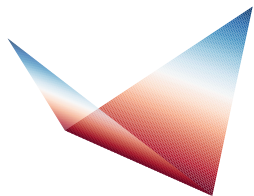
$\sigma$ -curved

Second order (Newton-like) [Hazan et al., 2006; Zhdanov and Kalnishkan, 2010]

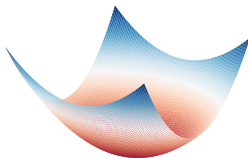
$\log(T)$  regret,  $\mathcal{O}(d^2)/\mathcal{O}(t^2)$  time/space per-step

Weaker than strong convexity

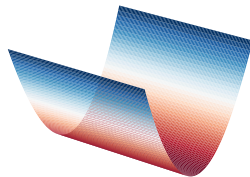
## OGD and losses



convex



strongly convex



$\sigma$ -curved

Second order (Newton-like) [Hazan et al., 2006; Zhdanov and Kalnishkan, 2010]

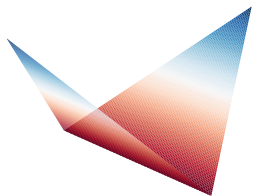
$\log(T)$  regret,  $\mathcal{O}(d^2)/\mathcal{O}(t^2)$  time/space per-step

Weaker than strong convexity

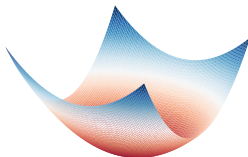
Satisfied by exp-concave losses:

↳ squared loss, squared hinge-loss, logistic loss

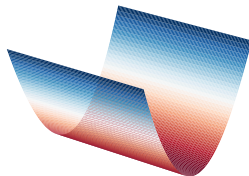
## OGD and losses



convex



strongly convex



$\sigma$ -curved

Second order (Newton-like) [Hazan et al., 2006; Zhdanov and Kalnishkan, 2010]  
 $\log(T)$  regret,  $\mathcal{O}(d^2)/\mathcal{O}(t^2)$  time/space per-step

Weaker than strong convexity

Satisfied by exp-concave losses:

↳ squared loss, squared hinge-loss, logistic loss

### Assumptions:

$\ell_t$  are  $\sigma$ -curved and  $|\ell'_t(z)| \leq L$  whenever  $|z| \leq C$  (scalar Lipschitz)

## Second-Order OKL (Kernel Online Newton Step)

Second-Order Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t, \quad \mathbf{A}_t = \sum_{s=1}^t \sigma \mathbf{g}_s \mathbf{g}_s^T + \alpha \mathbf{I} = \mathbf{G}_t \mathbf{G}_t^T + \alpha \mathbf{I}$$

## Second-Order OKL (Kernel Online Newton Step)

Second-Order Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t, \quad \mathbf{A}_t = \sum_{s=1}^t \sigma \mathbf{g}_s \mathbf{g}_s^T + \alpha \mathbf{I} = \mathbf{G}_t \mathbf{G}_t^T + \alpha \mathbf{I}$$

Regret [Hazan et al., 2006; Luo et al., 2016]

$$R(\mathbf{w}^*) \leq \overbrace{\alpha \|\mathbf{w}^* - \mathbf{w}_0\|_2^2}^{\text{initial error}} + \mathcal{O} \left( \sum_{t=1}^T \mathbf{g}_t^T (\mathbf{G}_t \mathbf{G}_t^T + \alpha \mathbf{I})^{-1} \mathbf{g}_t \right)$$

## Second-Order OKL (Kernel Online Newton Step)

Second-Order Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t, \quad \mathbf{A}_t = \sum_{s=1}^t \sigma \mathbf{g}_s \mathbf{g}_s^\top + \alpha \mathbf{I} = \mathbf{G}_t \mathbf{G}_t^\top + \alpha \mathbf{I}$$

Regret [Hazan et al., 2006; Luo et al., 2016]

$$\begin{aligned} R(\mathbf{w}^*) &\leq \overbrace{\alpha \|\mathbf{w}^* - \mathbf{w}_0\|_2^2}^{\text{initial error}} + \mathcal{O} \left( \sum_{t=1}^T \mathbf{g}_t^\top (\mathbf{G}_t \mathbf{G}_t^\top + \alpha \mathbf{I})^{-1} \mathbf{g}_t \right) \\ &\leq \alpha \|\mathbf{w}^* - \mathbf{w}_0\|^2 + \mathcal{O} \left( L \sum_{t=1}^T \phi_t^\top (\Phi_t \Phi_t^\top + \alpha \mathbf{I})^{-1} \phi_t \right) \end{aligned}$$



## Second-Order OKL (Kernel Online Newton Step)

Second-Order Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t, \quad \mathbf{A}_t = \sum_{s=1}^t \sigma \mathbf{g}_s \mathbf{g}_s^\top + \alpha \mathbf{I} = \mathbf{G}_t \mathbf{G}_t^\top + \alpha \mathbf{I}$$

Regret [Hazan et al., 2006; Luo et al., 2016]

$$\begin{aligned} R(\mathbf{w}^*) &\leq \overbrace{\alpha \|\mathbf{w}^* - \mathbf{w}_0\|_2^2}^{\text{initial error}} + \mathcal{O} \left( \sum_{t=1}^T \mathbf{g}_t^\top (\mathbf{G}_t \mathbf{G}_t^\top + \alpha \mathbf{I})^{-1} \mathbf{g}_t \right) \\ &\leq \alpha \|\mathbf{w}^* - \mathbf{w}_0\|_2^2 + \mathcal{O} \left( \overbrace{L \sum_{t=1}^T \phi_t^\top (\Phi_t \Phi_t^\top + \alpha \mathbf{I})^{-1} \phi_t}^{\text{online effective dimension}} \right) \end{aligned}$$

## Second-Order OKL (Kernel Online Newton Step)

Second-Order Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t, \quad \mathbf{A}_t = \sum_{s=1}^t \sigma \mathbf{g}_s \mathbf{g}_s^\top + \alpha \mathbf{I} = \mathbf{G}_t \mathbf{G}_t^\top + \alpha \mathbf{I}$$

Regret [Hazan et al., 2006; Luo et al., 2016]

$$\begin{aligned} R(\mathbf{w}^*) &\leq \overbrace{\alpha \|\mathbf{w}^* - \mathbf{w}_0\|_2^2}^{\text{initial error}} + \mathcal{O} \left( \sum_{t=1}^T \mathbf{g}_t^\top (\mathbf{G}_t \mathbf{G}_t^\top + \alpha \mathbf{I})^{-1} \mathbf{g}_t \right) \\ &\leq \alpha \|\mathbf{w}^* - \mathbf{w}_0\|^2 + \mathcal{O} \left( \overbrace{L \sum_{t=1}^T \phi_t^\top (\Phi_t \Phi_t^\top + \alpha \mathbf{I})^{-1} \phi_t}^{\text{online effective dimension}} \right) \\ &\leq \alpha \|\mathbf{w}^* - \mathbf{w}_0\|^2 + \mathcal{O}(\log \text{Det}(\mathbf{K}_T / \alpha + \mathbf{I}_n)) \end{aligned}$$

## Second-Order OKL (Kernel Online Newton Step)

Second-Order Gradient Descent

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{A}_t^{-1} \mathbf{g}_t, \quad \mathbf{A}_t = \sum_{s=1}^t \sigma \mathbf{g}_s \mathbf{g}_s^\top + \alpha \mathbf{I} = \mathbf{G}_t \mathbf{G}_t^\top + \alpha \mathbf{I}$$

Regret [Hazan et al., 2006; Luo et al., 2016]

$$\begin{aligned} R(\mathbf{w}^*) &\leq \overbrace{\alpha \|\mathbf{w}^* - \mathbf{w}_0\|_2^2}^{\text{initial error}} + \mathcal{O} \left( \sum_{t=1}^T \mathbf{g}_t^\top (\mathbf{G}_t \mathbf{G}_t^\top + \alpha \mathbf{I})^{-1} \mathbf{g}_t \right) \\ &\leq \alpha \|\mathbf{w}^* - \mathbf{w}_0\|^2 + \mathcal{O} \left( \overbrace{L \sum_{t=1}^T \phi_t^\top (\Phi_t \Phi_t^\top + \alpha \mathbf{I})^{-1} \phi_t}^{\text{online effective dimension}} \right) \\ &\leq \alpha \|\mathbf{w}^* - \mathbf{w}_0\|^2 + \mathcal{O}(\log \text{Det}(\mathbf{K}_T / \alpha + \mathbf{I}_n)) \\ &\leq \alpha \|\mathbf{w}^* - \mathbf{w}_0\|^2 + \mathcal{O}(d_{\text{eff}}^T(\alpha) \log(T)) \text{ [Calandriello et al., 2017b]} \end{aligned}$$

## Effective Dimension in online learning

$$R(\mathbf{w}^*) \leq \alpha \|\mathbf{w}^* - \mathbf{w}_0\|^2 + \mathcal{O}(d_{\text{eff}}^T(\alpha) \log(T))$$

$d_{\text{eff}}^T(\alpha)$  number of relevant orthogonal directions played by the adversary.

Every new orthogonal direction causes some regret.

↳ if it is played often enough (i.e.,  $\geq \alpha/(L\sigma)$ )

## Effective Dimension in online learning

$$R(\mathbf{w}^*) \leq \alpha \|\mathbf{w}^* - \mathbf{w}_0\|^2 + \mathcal{O}(d_{\text{eff}}^T(\alpha) \log(T))$$

$d_{\text{eff}}^T(\alpha)$  number of relevant orthogonal directions played by the adversary.

Every new orthogonal direction causes some regret.

↳ if it is played often enough (i.e.,  $\geq \alpha/(L\sigma)$ )

If all  $\phi_t$  are orthogonal

$$d_{\text{eff}}^T(\sqrt{T}) \sim \sqrt{T}$$

and

$$R(\mathbf{w}^*) \leq \sqrt{T} + \sqrt{T} \log(T) \sim \sqrt{T}$$

## Effective Dimension in online learning

$$R(\mathbf{w}^*) \leq \alpha \|\mathbf{w}^* - \mathbf{w}_0\|^2 + \mathcal{O}(d_{\text{eff}}^T(\alpha) \log(T))$$

$d_{\text{eff}}^T(\alpha)$  number of relevant orthogonal directions played by the adversary.

Every new orthogonal direction causes some regret.

↳ if it is played often enough (i.e.,  $\geq \alpha/(L\sigma)$ )

If all  $\phi_t$  are orthogonal

$$d_{\text{eff}}^T(\sqrt{T}) \sim \sqrt{T}$$

and

$$R(\mathbf{w}^*) \leq \sqrt{T} + \sqrt{T} \log(T) \sim \sqrt{T}$$

If  $\phi_t$  from finite subspace

$$d_{\text{eff}}^T(1) \sim \mathcal{O}(1) \leq r$$

is constant in  $T$  and

$$R(\mathbf{w}^*) \leq \mathcal{O}(1) + \mathcal{O}(1) \log(T) \sim \log T$$

## Approximating KONS

KONS:  $d_{\text{eff}}^T(\alpha) \log(T)$  regret

↳ large  $\mathcal{H} \Rightarrow \mathcal{O}(t)$  prediction  $\phi_t^\top \mathbf{w}_t$ ,  $\mathcal{O}(t^2)$  updates  $\mathbf{g}_t - \mathbf{A}_t^{-1} \mathbf{g}_t$

## Approximating KONS

KONS:  $d_{\text{eff}}^T(\alpha) \log(T)$  regret

↳ large  $\mathcal{H} \Rightarrow \mathcal{O}(t)$  prediction  $\phi_t^\top \mathbf{w}_t$ ,  $\mathcal{O}(t^2)$  updates  $\mathbf{g}_t - \mathbf{A}_t^{-1} \mathbf{g}_t$

Use approximate second order updates in large  $\mathcal{H}$  [Calandriello et al., 2017b]

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  regret, but prediction still depends on  $t$



## Approximating KONS

KONS:  $d_{\text{eff}}^T(\alpha) \log(T)$  regret

↳ large  $\mathcal{H} \Rightarrow \mathcal{O}(t)$  prediction  $\phi_t^T \mathbf{w}_t$ ,  $\mathcal{O}(t^2)$  updates  $\mathbf{g}_t - \mathbf{A}_t^{-1} \mathbf{g}_t$

Use approximate second order updates in large  $\mathcal{H}$  [Calandriello et al., 2017b]

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  regret, but prediction still depends on  $t$

Use exact second order updates in small approximate  $\tilde{\mathcal{H}}$

↳ replace  $\varphi$  with approximate map  $\tilde{\varphi}$  (random features, embeddings)

## Approximating KONS

KONS:  $d_{\text{eff}}^T(\alpha) \log(T)$  regret

↳ large  $\mathcal{H} \Rightarrow \mathcal{O}(t)$  prediction  $\phi_t^\top \mathbf{w}_t$ ,  $\mathcal{O}(t^2)$  updates  $\mathbf{g}_t - \mathbf{A}_t^{-1} \mathbf{g}_t$

Use approximate second order updates in large  $\mathcal{H}$  [Calandriello et al., 2017b]

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  regret, but prediction still depends on  $t$

Use exact second order updates in small approximate  $\tilde{\mathcal{H}}$

↳ replace  $\varphi$  with approximate map  $\tilde{\varphi}$  (random features, embeddings)  
finite  $\tilde{\mathcal{H}} \Rightarrow$  constant per-step prediction/update cost

## Approximating KONS

KONS:  $d_{\text{eff}}^T(\alpha) \log(T)$  regret

↳ large  $\mathcal{H} \Rightarrow \mathcal{O}(t)$  prediction  $\phi_t^\top \mathbf{w}_t$ ,  $\mathcal{O}(t^2)$  updates  $\mathbf{g}_t - \mathbf{A}_t^{-1} \mathbf{g}_t$

Use approximate second order updates in large  $\mathcal{H}$  [Calandriello et al., 2017b]

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  regret, but prediction still depends on  $t$

Use exact second order updates in small approximate  $\tilde{\mathcal{H}}$

↳ replace  $\varphi$  with approximate map  $\tilde{\varphi}$  (random features, embeddings)

finite  $\tilde{\mathcal{H}} \Rightarrow$  constant per-step prediction/update cost

$$\sum_{t=1}^T \ell_t(\tilde{\phi}_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \mathbf{w}^*) = \sum_{t=1}^T \underbrace{\ell_t(\tilde{\phi}_t \tilde{\mathbf{w}}_t) - \ell_t(\tilde{\phi}_t \bar{\mathbf{w}})}_a + \underbrace{\ell_t(\phi_t \bar{\mathbf{w}}) - \ell_t(\phi_t \mathbf{w}^*)}_b$$

## Approximating KONS

KONS:  $d_{\text{eff}}^T(\alpha) \log(T)$  regret

↳ large  $\mathcal{H} \Rightarrow \mathcal{O}(t)$  prediction  $\phi_t^\top \mathbf{w}_t$ ,  $\mathcal{O}(t^2)$  updates  $\mathbf{g}_t - \mathbf{A}_t^{-1} \mathbf{g}_t$

Use approximate second order updates in large  $\mathcal{H}$  [Calandriello et al., 2017b]

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  regret, but prediction still depends on  $t$

Use exact second order updates in small approximate  $\tilde{\mathcal{H}}$

↳ replace  $\varphi$  with approximate map  $\tilde{\varphi}$  (random features, embeddings)  
finite  $\tilde{\mathcal{H}} \Rightarrow$  constant per-step prediction/update cost

$$\sum_{t=1}^T \ell_t(\tilde{\phi}_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \mathbf{w}^*) = \sum_{t=1}^T \underbrace{\ell_t(\tilde{\phi}_t \tilde{\mathbf{w}}_t) - \ell_t(\tilde{\phi}_t \bar{\mathbf{w}})}_a + \underbrace{\ell_t(\phi_t \bar{\mathbf{w}}) - \ell_t(\phi_t \mathbf{w}^*)}_b$$

(a) Exact KONS in  $\tilde{\mathcal{H}}$ :  $d_{\text{eff}}^T(\alpha) \log(T)$

## Approximating KONS

KONS:  $d_{\text{eff}}^T(\alpha) \log(T)$  regret

↳ large  $\mathcal{H} \Rightarrow \mathcal{O}(t)$  prediction  $\phi_t^\top \mathbf{w}_t$ ,  $\mathcal{O}(t^2)$  updates  $\mathbf{g}_t - \mathbf{A}_t^{-1} \mathbf{g}_t$

Use approximate second order updates in large  $\mathcal{H}$  [Calandriello et al., 2017b]

↳  $d_{\text{eff}}^T(\alpha) \log(T)$  regret, but prediction still depends on  $t$

Use exact second order updates in small approximate  $\tilde{\mathcal{H}}$

↳ replace  $\varphi$  with approximate map  $\tilde{\varphi}$  (random features, embeddings)  
finite  $\tilde{\mathcal{H}} \Rightarrow$  constant per-step prediction/update cost

$$\sum_{t=1}^T \ell_t(\tilde{\phi}_t \tilde{\mathbf{w}}_t) - \ell_t(\phi_t \mathbf{w}^*) = \sum_{t=1}^T \underbrace{\ell_t(\tilde{\phi}_t \tilde{\mathbf{w}}_t) - \ell_t(\tilde{\phi}_t \bar{\mathbf{w}})}_a + \underbrace{\ell_t(\phi_t \bar{\mathbf{w}}) - \ell_t(\phi_t \mathbf{w}^*)}_b$$

(a) Exact KONS in  $\tilde{\mathcal{H}}$ :  $d_{\text{eff}}^T(\alpha) \log(T)$

(b) error between  $\bar{\mathbf{w}}$  best in  $\tilde{\mathcal{H}}$  and  $\mathbf{w}^*$  best in  $\mathcal{H}$ : bound how?

## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this

## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this same for fixed budget (e.g.,  $k$ -rank approx [Luo et al., 2016])

## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this same for fixed budget (e.g.,  $k$ -rank approx [Luo et al., 2016])

Use Nyström approximation instead and adapt it online



## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this same for fixed budget (e.g.,  $k$ -rank approx [Luo et al., 2016])

Use Nyström approximation instead and adapt it online

↳ if the adversary plays a “sufficiently orthogonal”  $\phi_t$ , add it to  $\mathcal{I}_{t+1}$

## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this same for fixed budget (e.g.,  $k$ -rank approx [Luo et al., 2016])

Use Nyström approximation instead and adapt it online

↳ if the adversary plays a “sufficiently orthogonal”  $\phi_t$ , add it to  $\mathcal{I}_{t+1}$   
 $\tilde{\mathcal{H}}_t = \text{Span}(\mathcal{I}_t)$  defined using  $m_t$  inducing points  $\mathcal{I}_t = \{\phi_s\}_{s=1}^{m_t}$

## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this same for fixed budget (e.g.,  $k$ -rank approx [Luo et al., 2016])

Use Nyström approximation instead and adapt it online

↳ if the adversary plays a “sufficiently orthogonal”  $\phi_t$ , add it to  $\mathcal{I}_{t+1}$

$\tilde{\mathcal{H}}_t = \text{Span}(\mathcal{I}_t)$  defined using  $m_t$  inducing points  $\mathcal{I}_t = \{\phi_s\}_{s=1}^{m_t}$

Use RLS (KORS) to select inducing points

## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this same for fixed budget (e.g.,  $k$ -rank approx [Luo et al., 2016])

Use Nyström approximation instead and adapt it online

↳ if the adversary plays a “sufficiently orthogonal”  $\phi_t$ , add it to  $\mathcal{I}_{t+1}$   
 $\tilde{\mathcal{H}}_t = \text{Span}(\mathcal{I}_t)$  defined using  $m_t$  inducing points  $\mathcal{I}_t = \{\phi_s\}_{s=1}^{m_t}$

Use RLS (KORS) to select inducing points

↳ SQUEAK without removal ( $\mathcal{I}_t \subseteq \mathcal{I}_{t+1}$ ,  $\tilde{\mathcal{H}}_t \subseteq \tilde{\mathcal{H}}_{t+1}$ )

## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this same for fixed budget (e.g.,  $k$ -rank approx [Luo et al., 2016])

Use Nyström approximation instead and adapt it online

↳ if the adversary plays a “sufficiently orthogonal”  $\phi_t$ , add it to  $\mathcal{I}_{t+1}$   
 $\tilde{\mathcal{H}}_t = \text{Span}(\mathcal{I}_t)$  defined using  $m_t$  inducing points  $\mathcal{I}_t = \{\phi_s\}_{s=1}^{m_t}$

Use RLS (KORS) to select inducing points

↳ SQUEAK without removal ( $\mathcal{I}_t \subseteq \mathcal{I}_{t+1}$ ,  $\tilde{\mathcal{H}}_t \subseteq \tilde{\mathcal{H}}_{t+1}$ )

w.h.p. accurate and maximum size  $|\tilde{\mathcal{H}}_t| \leq \mathcal{O}(d_{\text{eff}}^T(\gamma) \log^2(T))$

## Subspace approximation error

$\tilde{\mathcal{H}}$  cannot be fixed

↳ the adversary will find orthogonal points and exploit this same for fixed budget (e.g.,  $k$ -rank approx [Luo et al., 2016])

Use Nyström approximation instead and adapt it online

↳ if the adversary plays a “sufficiently orthogonal”  $\phi_t$ , add it to  $\mathcal{I}_{t+1}$   
 $\tilde{\mathcal{H}}_t = \text{Span}(\mathcal{I}_t)$  defined using  $m_t$  inducing points  $\mathcal{I}_t = \{\phi_s\}_{s=1}^{m_t}$

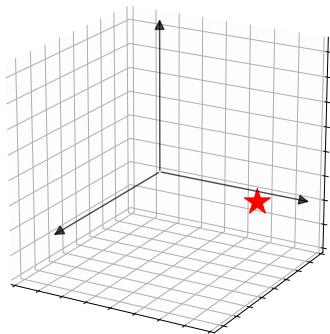
Use RLS (KORS) to select inducing points

↳ SQUEAK without removal ( $\mathcal{I}_t \subseteq \mathcal{I}_{t+1}$ ,  $\tilde{\mathcal{H}}_t \subseteq \tilde{\mathcal{H}}_{t+1}$ )

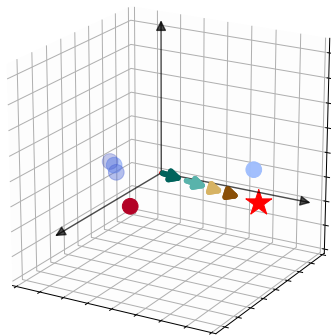
w.h.p. accurate and maximum size  $|\tilde{\mathcal{H}}_t| \leq \mathcal{O}(d_{\text{eff}}^T(\gamma) \log^2(T))$

$\tilde{\mathcal{O}}(d_{\text{eff}}^T(\gamma)^2)$  time/space cost to run exact KONS in  $\tilde{\mathcal{H}}_t$

# PROS-N-KONS

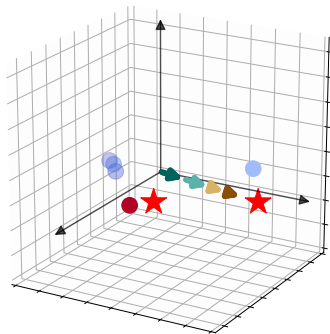


# PROS-N-KONS

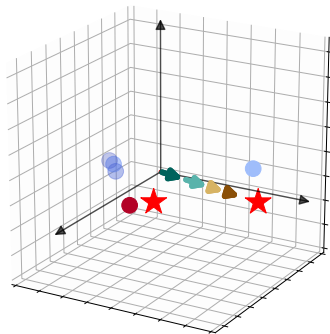




# PROS-N-KONS

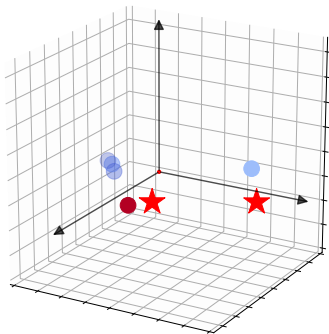


# PROS-N-KONS



Every time we change  $\tilde{\mathcal{H}}$  we pay  $\alpha \|\bar{\mathbf{w}}_j - \mathbf{w}_{t_j}\|_2^2$  (initial error in GD)  
↳ the **adversary** can influence  $\mathbf{w}_{t_j}$  and make it large

# PROS-N-KONS



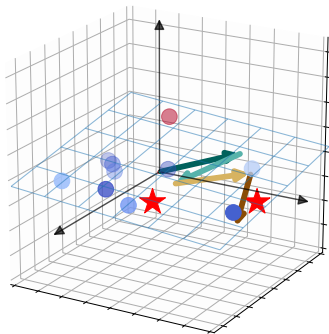
Reset  $\tilde{\mathbf{w}}_t$  and  $\tilde{\mathbf{A}}_t$  when  $\tilde{\mathcal{H}}_t$  changes

↳ **wasteful**, but **not too often**. At most  $J \leq d_{\text{eff}}^T(\gamma)$  times.

**learning** is preserved through  $\tilde{\mathcal{H}}_t$  that always improves

**adaptive** doubling trick

# PROS-N-KONS



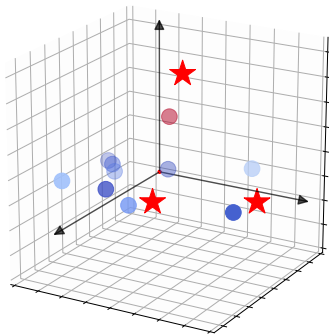
Reset  $\tilde{\mathbf{w}}_t$  and  $\tilde{\mathbf{A}}_t$  when  $\tilde{\mathcal{H}}_t$  changes

↳ **wasteful**, but **not too often**. At most  $J \leq d_{\text{eff}}^T(\gamma)$  times.

**learning** is preserved through  $\tilde{\mathcal{H}}_t$  that always improves

**adaptive** doubling trick

# PROS-N-KONS



Reset  $\tilde{\mathbf{w}}_t$  and  $\tilde{\mathbf{A}}_t$  when  $\tilde{\mathcal{H}}_t$  changes

↳ **wasteful**, but **not too often**. At most  $J \leq d_{\text{eff}}^T(\gamma)$  times.

**learning** is preserved through  $\tilde{\mathcal{H}}_t$  that always improves

**adaptive** doubling trick

## Experiments - regression

$\alpha = 1, \gamma = 1$						
Algorithm	cadata $n = 20k, d = 8$			casp $n = 45k, d = 9$		
	Avg. Squared Loss	#SV	Time	Avg. Squared Loss	#SV	Time
FOGD	$0.04097 \pm 0.00015$	30	—	$0.08021 \pm 0.00031$	30	—
NOGD	$0.03983 \pm 0.00018$	30	—	$0.07844 \pm 0.00008$	30	—
<b>PROS-N-KONS</b>	$0.03095 \pm 0.00110$	20	18.59	<b><math>0.06773 \pm 0.00105</math></b>	21	40.73
CON-KONS	<b><math>0.02850 \pm 0.00174</math></b>	19	18.45	<b><math>0.06832 \pm 0.00315</math></b>	20	40.91
B-KONS	$0.03095 \pm 0.00118$	19	18.65	<b><math>0.06775 \pm 0.00067</math></b>	21	41.13
BATCH	$0.02202 \pm 0.00002$	—	—	$0.06100 \pm 0.00003$	—	—

Algorithm	slice $n = 53k, d = 385$			year $n = 463k, d = 90$		
	Avg. Squared Loss	#SV	Time	Avg. Squared Loss	#SV	Time
FOGD	<b><math>0.00726 \pm 0.00019</math></b>	30	—	$0.01427 \pm 0.00004$	30	—
NOGD	$0.02636 \pm 0.00460$	30	—	$0.01427 \pm 0.00004$	30	—
DUAL-SGD	—	—	—	$0.01440 \pm 0.00000$	100	—
<b>PROS-N-KONS</b>	did not complete	—	—	$0.01450 \pm 0.00014$	149	884.82
CON-KONS	did not complete	—	—	$0.01444 \pm 0.00017$	147	889.42
B-KONS	<b><math>0.00913 \pm 0.00045</math></b>	100	60	<b><math>0.01302 \pm 0.00006</math></b>	100	505.36
BATCH	$0.00212 \pm 0.00001$	—	—	$0.01147 \pm 0.00001$	—	—

## Experiments - binary classification

$\alpha = 1, \gamma = 1$						
Algorithm	ijcnn1 $n = 141,691, d = 22$			cod-rna $n = 271,617, d = 8$		
	accuracy	#SV	time	accuracy	#SV	time
FOGD	9.06 $\pm$ 0.05	400	—	10.30 $\pm$ 0.10	400	—
NOGD	9.55 $\pm$ 0.01	100	—	13.80 $\pm$ 2.10	100	—
DUAL-SGD	<b>8.35</b> $\pm$ 0.20	100	—	<b>4.83</b> $\pm$ 0.21	100	—
PROS-N-KONS	9.70 $\pm$ 0.01	100	211.91	13.95 $\pm$ 1.19	38	270.81
CON-KONS	9.64 $\pm$ 0.01	101	215.71	18.99 $\pm$ 9.47	38	271.85
B-KONS	9.70 $\pm$ 0.01	98	206.53	13.99 $\pm$ 1.16	38	274.94
BATCH	8.33 $\pm$ 0.03	—	—	3.781 $\pm$ 0.01	—	—

$\alpha = 0.01, \gamma = 0.01$						
Algorithm	ijcnn1 $n = 141,691, d = 22$			cod-rna $n = 271,617, d = 8$		
	accuracy	#SV	time	accuracy	#SV	time
FOGD	9.06 $\pm$ 0.05	400	—	10.30 $\pm$ 0.10	400	—
NOGD	9.55 $\pm$ 0.01	100	—	13.80 $\pm$ 2.10	100	—
DUAL-SGD	8.35 $\pm$ 0.20	100	—	4.83 $\pm$ 0.21	100	—
PROS-N-KONS	10.73 $\pm$ 0.12	436	1003.82	4.91 $\pm$ 0.04	111	459.28
CON-KONS	6.23 $\pm$ 0.18	432	987.33	5.81 $\pm$ 1.96	111	458.90
B-KONS	<b>4.85</b> $\pm$ 0.08	100	147.22	<b>4.57</b> $\pm$ 0.05	100	333.57
BATCH	5.61 $\pm$ 0.01	—	—	3.61 $\pm$ 0.01	—	—

## PROS-N-KONS - recap

PROS-N-KONS: avoid **curse of kernelization**, constant per-step cost



## PROS-N-KONS - recap

PROS-N-KONS: avoid **curse of kernelization**, constant per-step cost  
First approximate method with **logarithmic regret**

## PROS-N-KONS - recap

PROS-N-KONS: avoid **curse of kernelization**, constant per-step cost

First approximate method with **logarithmic regret**

Future work

## PROS-N-KONS - recap

PROS-N-KONS: avoid **curse of kernelization**, constant per-step cost

First approximate method with **logarithmic regret**

Future work

Restarts really necessary?

## PROS-N-KONS - recap

PROS-N-KONS: avoid **curse of kernelization**, constant per-step cost

First approximate method with **logarithmic regret**

Future work

Restarts really necessary?

Adaptive  $\alpha$  and  $\gamma$ ?

## PROS-N-KONS - recap

PROS-N-KONS: avoid **curse of kernelization**, constant per-step cost

First approximate method with **logarithmic regret**

Future work

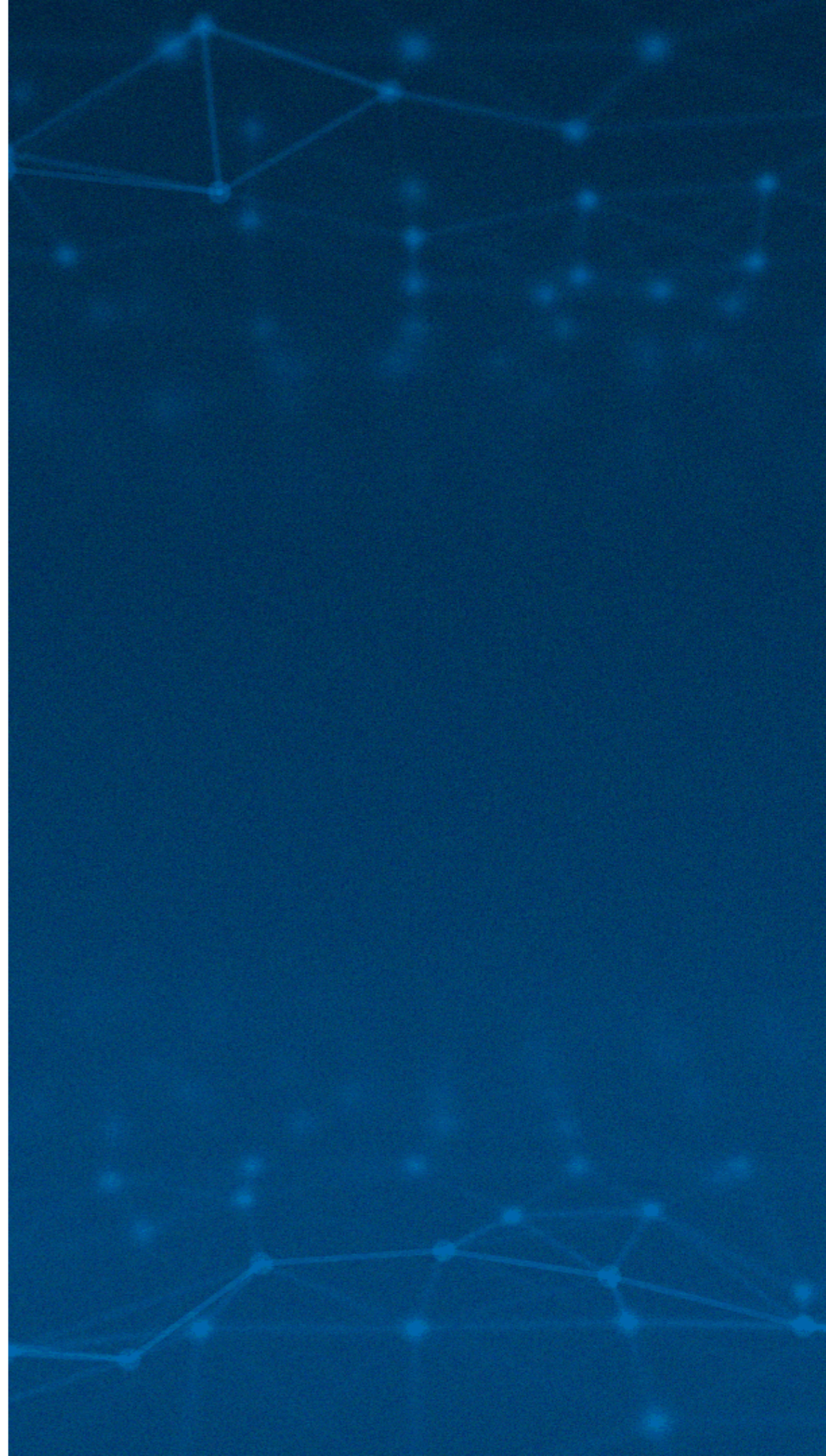
Restarts really necessary?

Adaptive  $\alpha$  and  $\gamma$ ?

... and now, back to the beginning!

# BACK TO THE BEGINNING: GRAPH SPARSIFICATION

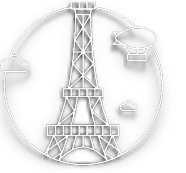
---



# SCALING UP GRAPH LEARNING



DeepMind

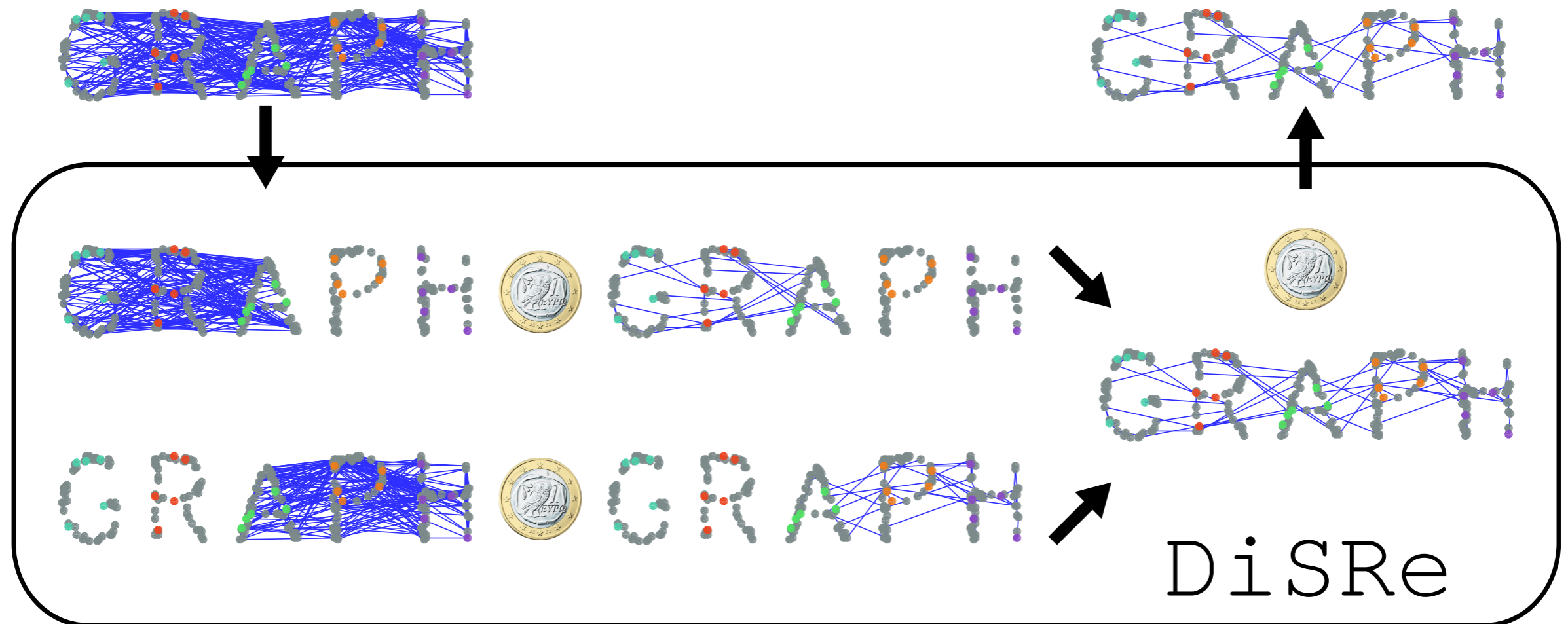


- ▶ Large graphs do not fit in a single machine memory
- ▶ multiple passes slow, distribution has communication costs
- ▶ removing edges impacts structure/accuracy
- ▶ Make the graph sparse, while preserving its structure for learning

$$(1 - \varepsilon)\mathbf{L}_G \preceq \mathbf{L}_H \preceq (1 + \varepsilon)\mathbf{L}_G$$

$$(1 - \varepsilon)\mathbf{L}_G - \varepsilon\gamma\mathbf{I} \preceq \mathbf{L}_H \preceq (1 + \varepsilon)\mathbf{L}_G + \varepsilon\gamma\mathbf{I}$$

# DISRE GUARANTEES



## Theorem

Given an arbitrary graph  $\mathcal{G}$  w.h.p. DISRE satisfies

- (1) each sub-graphs is an  $(\varepsilon, \gamma)$ -sparsifier
- (2) with at most  $\mathcal{O}(d_{\text{eff}}(\gamma) \log(n))$  edges.



# DISRE EXPERIMENTS

**Dataset:** Amazon co-purchase graph [Yang and Leskovec 2015]

↳ **natural**, artificially sparse (true graph known only to Amazon)

↳ we compute 4-step random walk to recover removed co-purchases  
[Gleich and Mahoney 2015]

**Target:** eigenvector  $\mathbf{v}$  associated with  $\lambda_2(\mathbf{L}_{\mathcal{G}})$  [Sadhanala et al. 2016]

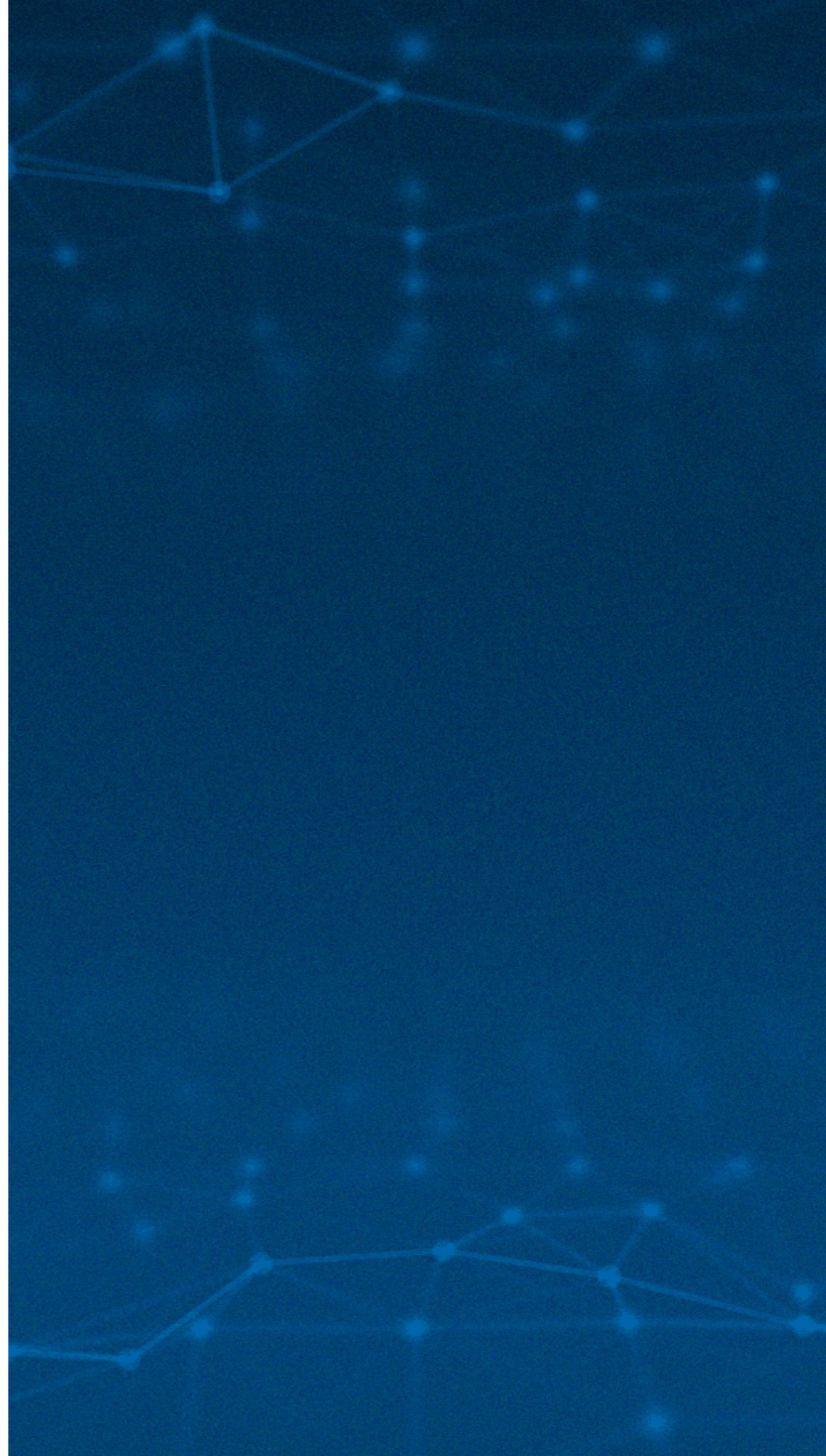
$n = 334,863$  nodes,  $m = 98,465,352$  edges (294 avg. degree)

<i>Alg.</i>	<i>Parameters</i>	$ \mathcal{E} $ ( $\times 10^6$ )	$\ \tilde{\mathbf{f}} - \mathbf{v}\ _2^2$ ( $\sigma = 10^{-3}$ )	$\ \hat{\mathbf{f}} - \mathbf{v}\ _2^2$ ( $\sigma = 10^{-2}$ )
EXACT		98.5	$0.067 \pm 0.0004$	$0.756 \pm 0.006$
kN	$k = 60$	15.7	$0.172 \pm 0.0004$	$0.822 \pm 0.002$
DISRE	$\gamma = 0$	22.8	$0.068 \pm 0.0004$	<b><math>0.756 \pm 0.005</math></b>
DISRE	$\gamma = 10^2$	11.8	<b><math>0.068 \pm 0.0002</math></b>	$0.772 \pm 0.004$

**Time:** Loading  $\mathcal{G}$  from disk 90sec, DISRE 120sec( $k = 4 \times 32$  CPU),  
computing  $\tilde{\mathbf{f}}$  120sec, computing  $\hat{\mathbf{f}}$  720sec

**AFTER 12  
YEARS?**

**THIS IS JUST THE  
BEGINNING!**



## ▶ SPARSIFYING GP-UCB RIGHT

- More than 20 years of heuristics
- Even 2019 results on sparsifying LinUCB can go wrong
- BKB - **adaptive** dictionary, guarantees regret and is fast

## ▶ BATCHED GP-UCB SPARSIFICATION - stay tuned!

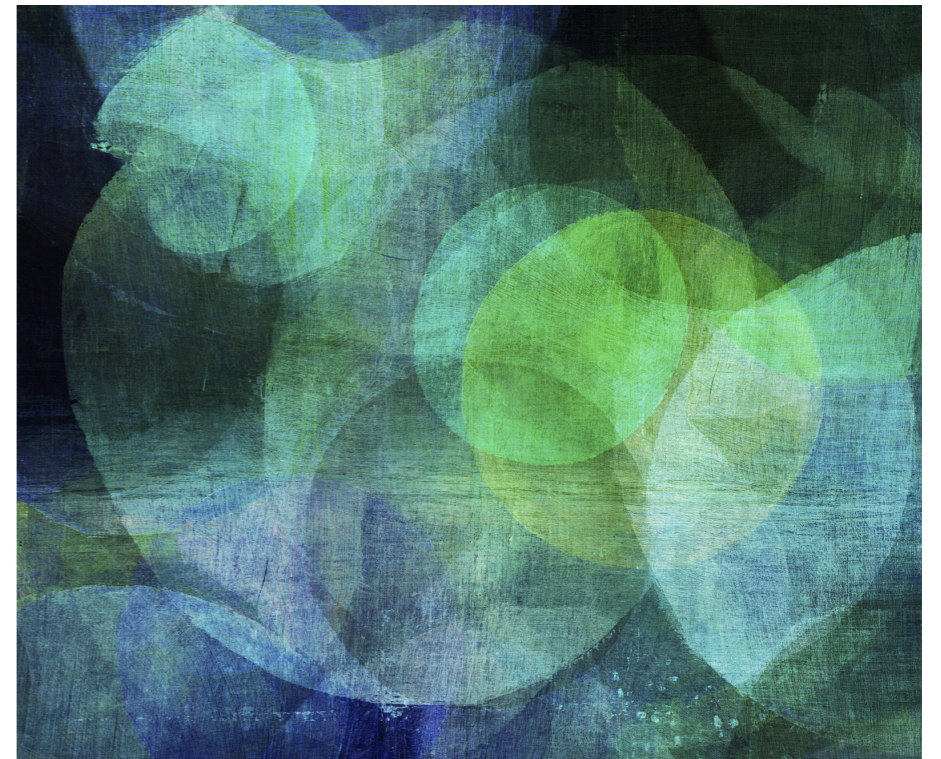
## ▶ *Negative dependence/online leverages scores/DPPs*

## ▶ FAST SAMPLING OF DPPs - repulsion for the sets!

- w/Michał Dereziński and Daniele Calandriello
- online lev. Scores + R-DPP + downsampling → **perfect**

Michał Valko, [valkom@google.com](mailto:valkom@google.com)

<http://researchers.lille.inria.fr/~valko/hp/>



Michal Valko, [valkom@google.com](mailto:valkom@google.com)  
<http://researchers.lille.inria.fr/~valko/hp/>