



Pack only the essentials: distributed sequential sampling for adaptive kernel DL

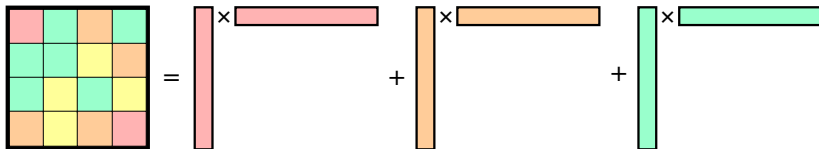
with **Daniele Calandriello** and **Alessandro Lazaric**

SequeL team, Inria Lille - Nord Europe, France

appeared in AISTATS 2017

Distributed sequential sampling for adaptive kernel DL

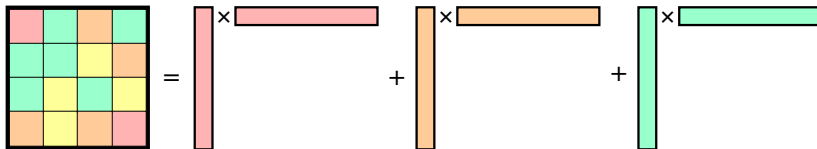
What is Dictionary Learning (DL)?



Finding an **accurate** representation of the input data as a linear combination of a **small** set of basic elements (**atoms**)

Distributed sequential sampling for adaptive kernel DL

What is Dictionary Learning (DL)?

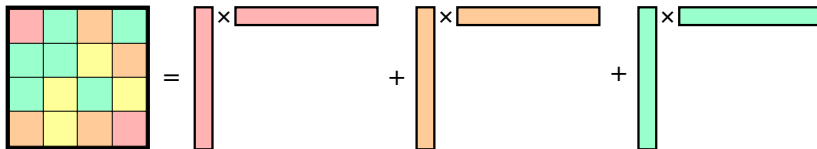


Finding an **accurate** representation of the input data as a linear combination of a **small** set of basic elements (**atoms**)

Representation/Unsupervised learning

Distributed sequential sampling for adaptive kernel DL

What is Dictionary Learning (DL)?



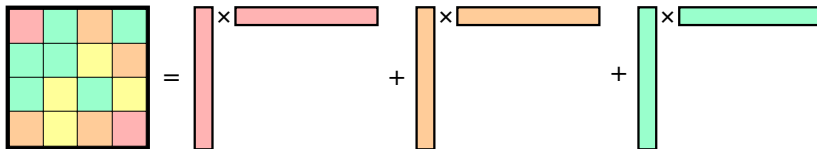
Finding an **accurate** representation of the input data as a linear combination of a **small** set of basic elements (**atoms**)

Representation/Unsupervised learning

“Most important open problem in ML” **Y. LeCun**, NIPS 2016

Distributed sequential sampling for adaptive kernel DL

What is Dictionary Learning (DL)?



Finding an **accurate** representation of the input data as a linear combination of a **small** set of basic elements (**atoms**)

Representation/Unsupervised learning

“Most important open problem in ML” **Y. LeCun**, NIPS 2016

“Already solved” **J. Schmidhuber**, NIPS 2016

Distributed sequential sampling for adaptive kernel DL

Why DL for kernel problems?

Kernel methods have huge scalability problem

Problem: for a dataset \mathcal{D} with n samples

$\mathcal{O}(n^2)$ time to construct kernel matrix \mathbf{K}

$\mathcal{O}(n^3)$ time to compute solution

$\mathcal{O}(n^2)$ space to store it

Distributed sequential sampling for adaptive kernel DL

Why DL for kernel problems?

Kernel methods have huge scalability problem

Problem: for a dataset \mathcal{D} with n samples

$\mathcal{O}(n^2)$ time to construct kernel matrix \mathbf{K}

$\mathcal{O}(n^3)$ time to compute solution

$\mathcal{O}(n^2)$ space to store it

Solution:

compute accurate, small dictionary \mathcal{I} to represent \mathcal{D}

compute approximate solution on \mathcal{I} efficiently

Distributed sequential sampling for adaptive kernel DL

Why DL for kernel problems?

Problem: Existing DL methods guarantee either scalability or accuracy

Distributed sequential sampling for adaptive kernel DL

Why DL for kernel problems?

Problem: Existing DL methods guarantee either scalability or accuracy

we want both

Distributed sequential sampling for adaptive kernel DL

Why DL for kernel problems?

Problem: Existing DL methods guarantee either scalability or accuracy

we want both

We present SQUEAK — a dictionary learning algorithm that guarantees

Distributed sequential sampling for adaptive kernel DL

Why DL for kernel problems?

Problem: Existing DL methods guarantee either scalability or accuracy

we want both

We present SQUEAK — a dictionary learning algorithm that guarantees
In all cases accurate reconstruction of the input

Distributed sequential sampling for adaptive kernel DL

Why DL for kernel problems?

Problem: Existing DL methods guarantee either scalability or accuracy

we want both

We present SQUEAK — a dictionary learning algorithm that guarantees

In all cases accurate reconstruction of the input

Adapts to the data:

on “easy” problems small $\mathcal{O}(n)$ space/time requirements

on “hard” problems not worse than storing whole input

Distributed sequential sampling for adaptive kernel DL

Why DL for kernel problems?

Problem: Existing DL methods guarantee either scalability or accuracy

we want both

We present SQUEAK — a dictionary learning algorithm that guarantees

In all cases accurate reconstruction of the input

Adapts to the data:

on “easy” problems small $\mathcal{O}(n)$ space/time requirements

on “hard” problems not worse than storing whole input

Only local data access, distributed version with $\mathcal{O}(\log(n))$ runtime

Distributed sequential sampling for adaptive kernel DL

We consider Positive Semi-Definite matrices

$$\mathbf{A} = \mathbf{A}^{1/2}(\mathbf{A}^{1/2})^T = \sum_{i=1}^n \mathbf{a}_i \mathbf{a}_i^T \quad \tilde{\mathbf{A}} = \sum_{i=1}^m w_i \mathbf{x}_i \mathbf{x}_i^T$$

Method	w_i	\mathbf{x}_i	Accuracy	Space	Time
Whole Input	1	\mathbf{a}_i	★★★★★		

Distributed sequential sampling for adaptive kernel DL

We consider Positive Semi-Definite matrices

$$\mathbf{A} = \mathbf{A}^{1/2}(\mathbf{A}^{1/2})^T = \sum_{i=1}^n \mathbf{a}_i \mathbf{a}_i^T \quad \tilde{\mathbf{A}} = \sum_{i=1}^m w_i \mathbf{x}_i \mathbf{x}_i^T$$

Method	w_i	\mathbf{x}_i	Accuracy	Space	Time
Whole Input	1	\mathbf{a}_i	★★★★★		
Empty dictionary	0	$\mathbf{0}$		★★★★★	★★★★★

Distributed sequential sampling for adaptive kernel DL

We consider Positive Semi-Definite matrices

$$\mathbf{A} = \mathbf{A}^{1/2}(\mathbf{A}^{1/2})^T = \sum_{i=1}^n \mathbf{a}_i \mathbf{a}_i^T \quad \tilde{\mathbf{A}} = \sum_{i=1}^m w_i \mathbf{x}_i \mathbf{x}_i^T$$

Method	w_i	\mathbf{x}_i	Accuracy	Space	Time
Whole Input	1	\mathbf{a}_i	★★★★★		
PCA	λ_i	\mathbf{u}_i	★★★★★	★★★★★	★
Empty dictionary	0	$\mathbf{0}$		★★★★★	★★★★★

Distributed sequential sampling for adaptive kernel DL

We consider Positive Semi-Definite matrices

$$\mathbf{A} = \mathbf{A}^{1/2}(\mathbf{A}^{1/2})^T = \sum_{i=1}^n \mathbf{a}_i \mathbf{a}_i^T \quad \tilde{\mathbf{A}} = \sum_{i=1}^m w_i \mathbf{x}_i \mathbf{x}_i^T$$

Method	w_i	\mathbf{x}_i	Accuracy	Space	Time
Whole Input	1	\mathbf{a}_i	★★★★★		
PCA	λ_i	\mathbf{u}_i	★★★★★	★★★★★	★
RLS (this)	$1/\tau_i$	\mathbf{a}_i	★★★★	★★★★	★★★★
Empty dictionary	0	$\mathbf{0}$		★★★★★	★★★★★

Distributed sequential sampling for **adaptive** kernel DL

We consider Positive Semi-Definite matrices

$$\mathbf{A} = \mathbf{A}^{1/2}(\mathbf{A}^{1/2})^T = \sum_{i=1}^n \mathbf{a}_i \mathbf{a}_i^T$$

$$\tilde{\mathbf{A}} = \sum_{i=1}^m w_i \mathbf{x}_i \mathbf{x}_i^T$$

Method	w_i	\mathbf{x}_i	Accuracy	Space	Time
Whole Input	1	\mathbf{a}_i	★★★★★		
PCA	λ_i	\mathbf{u}_i	★★★★★	★★★★★	★
RLS (this)	$1/\tau_i$	\mathbf{a}_i	★★★★★	★★★★★	★★★★
Uniform	n/m	\mathbf{a}_i	★★	★★	★★★★★
Empty dictionary	0	$\mathbf{0}$		★★★★★	★★★★★

Preliminaries: Setting and Kernels

Indexing $[t] = \{1, \dots, t\}$, notation \mathbf{K} matrices, \mathbf{k} vectors, k scalar

Dataset $\mathcal{D}_n = \{\mathbf{x}_i\}_{i=1}^n$, samples $\mathbf{x}_i \in \mathcal{X}$ (e.g., \mathbb{R}^d)

Kernel function $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

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

Kernel trick

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

Feature matrix $\Phi_t = [\phi_1, \phi_2, \dots, \phi_t] : \mathbb{R}^t \rightarrow \mathcal{H}$

Empirical kernel matrix $\mathbf{K}_t \in \mathbb{R}^{t \times t} = \mathbf{K}_{[t],[t]} = \Phi_t^\top \Phi_t$

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

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

Find a dictionary $\mathcal{I} = \{(w_j, \phi_j)\}_{j=1}^m$ such that $\tilde{\mathbf{K}} = f(\mathcal{I})$ close to \mathbf{K}

Preliminaries: Linear Algebra

(Full) Singular Value Decomposition $\Phi = \mathbf{V}\Sigma\mathbf{U}^T$, Σ rectangular
Eigendecomposition $\Phi^T\Phi = \mathbf{U}\Sigma^T\Sigma\mathbf{U}^T = \mathbf{U}\Lambda\mathbf{U}^T = \mathbf{K}$

Matrix norms (if omitted, ℓ_2 norm)

$$\ell_2 \text{ norm} \quad \|\mathbf{A}\|_2 = \sup_{\|\mathbf{x}\|_2=1} \|\mathbf{A}\mathbf{x}\|_2 = \max \lambda_i$$

$$\text{Frobenius norm} \quad \|\mathbf{A}\|_F^2 = \sum a_{i,j}^2 = \sum \lambda_i^2$$

Useful equality for arbitrary $n \times m$ matrix (or operator)

$$\Phi\Phi^T(\Phi\Phi^T + \gamma\mathbf{I}_n)^{-1} = \Phi(\Phi^T\Phi + \gamma\mathbf{I}_m)^{-1}\Phi^T$$

Example: Kernel Ridge Regression

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

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

Example: Kernel Ridge Regression

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

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

If we can have accurate low-rank approximations ...

$$\tilde{\mathbf{K}}_n \preceq \mathbf{K}_n \preceq \tilde{\mathbf{K}}_n + \frac{\gamma}{1 - \varepsilon} \mathbf{I}$$

Example: Kernel Ridge Regression

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

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

If we can have accurate low-rank approximations ...

$$\tilde{\mathbf{K}}_n \preceq \mathbf{K}_n \preceq \tilde{\mathbf{K}}_n + \frac{\gamma}{1 - \varepsilon} \mathbf{I}$$

... then we can use them for to get good approximate solutions:

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

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

Example: Kernel Ridge Regression

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

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

If we can have accurate low-rank approximations ...

$$\tilde{\mathbf{K}}_n \preceq \mathbf{K}_n \preceq \tilde{\mathbf{K}}_n + \frac{\gamma}{1 - \varepsilon} \mathbf{I}$$

... then we can use them for to get good approximate solutions:

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

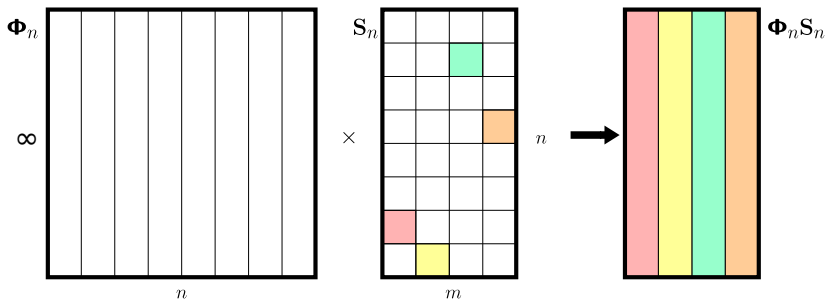
$$R(\tilde{\mathbf{w}}_n) \leq \left(1 + \frac{1}{1 - \varepsilon} \right) R(\hat{\mathbf{w}}_n)$$

~~$\mathcal{O}(n^3)$~~ \Rightarrow $\mathcal{O}(nm + m^3)$ time to compute the approx. solution

~~$\mathcal{O}(n^2)$~~ \Rightarrow $\mathcal{O}(nm)$ space to store dictionary

Reconstruction Guarantees

Given dataset \mathcal{D}_n and dictionary \mathcal{I}_n , the selection matrix \mathbf{S}_n is defined as



$$\sum_{i=1}^m w_i \phi_i \phi_i^T = \sum_{i=1}^m (\sqrt{w_i} \phi_i)(\sqrt{w_i} \phi_i)^T = \Phi_n \mathbf{S}_n \mathbf{S}_n^T \Phi_n^T$$

Reconstruction guarantees

Consider the regularized projection Ψ_n

$$\begin{aligned}\Psi_n &= \Phi_n \Phi_n^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1} = (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \Phi_n \Phi_n^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \\ &= \sum_{i=1}^n (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \phi_i \phi_i^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} = \sum_{i=1}^n \psi_i \psi_i^T\end{aligned}$$

$$\tilde{\Psi}_n = (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \Phi_n \mathbf{S}_n \mathbf{S}_n^T \Phi_n^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} = \sum_{j=1}^m w_j \psi_j \psi_j^T$$

Reconstruction guarantees

Consider the regularized projection Ψ_n

$$\begin{aligned}\Psi_n &= \Phi_n \Phi_n^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1} = (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \Phi_n \Phi_n^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \\ &= \sum_{i=1}^n (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \phi_i \phi_i^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} = \sum_{i=1}^n \psi_i \psi_i^T\end{aligned}$$

$$\tilde{\Psi}_n = (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \Phi_n \mathbf{S}_n \mathbf{S}_n^T \Phi_n^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} = \sum_{j=1}^m w_j \psi_j \psi_j^T$$

An accurate dictionary satisfies

$$\|\Psi_n - \tilde{\Psi}_n\|_2 \leq \varepsilon$$

Reconstruction guarantees

Consider the regularized projection Ψ_n

$$\begin{aligned}\Psi_n &= \Phi_n \Phi_n^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1} = (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \Phi_n \Phi_n^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \\ &= \sum_{i=1}^n (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \phi_i \phi_i^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} = \sum_{i=1}^n \psi_i \psi_i^T\end{aligned}$$

$$\tilde{\Psi}_n = (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \Phi_n \mathbf{S}_n \mathbf{S}_n^T \Phi_n^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} = \sum_{j=1}^m w_j \psi_j \psi_j^T$$

An accurate dictionary satisfies

$$\|\Psi_n - \tilde{\Psi}_n\|_2 \leq \varepsilon$$

equivalent to mixed additive/multiplicative error in quadratic form

$$(1 - \varepsilon) \Phi_n \Phi_n^T - \varepsilon \gamma \mathbf{I} \preceq \Phi_n \mathbf{S}_n \mathbf{S}_n^T \Phi_n^T \preceq (1 + \varepsilon) \Phi_n \Phi_n^T + \varepsilon \gamma \mathbf{I}$$

Reconstruction guarantees

Why would bounding $\|\Psi_n - \tilde{\Psi}_n\|_2$ be useful?

Reconstruction guarantees

Why would bounding $\|\Psi_n - \tilde{\Psi}_n\|_2$ be useful?

$$\begin{aligned}\|\Psi_n - \tilde{\Psi}_n\|_2 &= \|(\Phi_n \Phi_n^\top + \gamma \mathbf{I})^{-1/2} \Phi_n (\mathbf{I} - \mathbf{S}_n \mathbf{S}_n^\top) \Phi_n (\Phi_n \Phi_n^\top + \gamma \mathbf{I})^{-1/2}\|_2 \\ &= \|(\Sigma \Sigma^\top + \gamma \mathbf{I})^{-1/2} \Sigma \mathbf{U}^\top (\mathbf{I} - \mathbf{S}_n \mathbf{S}_n^\top) \mathbf{U} \Sigma^\top (\Sigma \Sigma^\top + \gamma \mathbf{I})^{-1/2}\|_2 \\ &= \|(\mathbf{K}_n + \gamma \mathbf{I})^{-1/2} \mathbf{K}_n^{1/2} (\mathbf{I} - \mathbf{S}_n \mathbf{S}_n^\top) \mathbf{K}_n^{1/2} (\mathbf{K}_n + \gamma \mathbf{I})^{-1/2}\|_2 \\ &= \|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2\end{aligned}$$

Reconstruction guarantees

Why would bounding $\|\Psi_n - \tilde{\Psi}_n\|_2$ be useful?

$$\begin{aligned}\|\Psi_n - \tilde{\Psi}_n\|_2 &= \|(\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2} \Phi_n (\mathbf{I} - \mathbf{S}_n \mathbf{S}_n^T) \Phi_n (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1/2}\|_2 \\ &= \|(\Sigma \Sigma^T + \gamma \mathbf{I})^{-1/2} \Sigma \mathbf{U}^T (\mathbf{I} - \mathbf{S}_n \mathbf{S}_n^T) \mathbf{U} \Sigma^T (\Sigma \Sigma^T + \gamma \mathbf{I})^{-1/2}\|_2 \\ &= \|(\mathbf{K}_n + \gamma \mathbf{I})^{-1/2} \mathbf{K}_n^{1/2} (\mathbf{I} - \mathbf{S}_n \mathbf{S}_n^T) \mathbf{K}_n^{1/2} (\mathbf{K}_n + \gamma \mathbf{I})^{-1/2}\|_2 \\ &= \|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2\end{aligned}$$

with

$$\begin{aligned}\mathbf{P}_n &= \mathbf{K}_n (\mathbf{K}_n + \gamma \mathbf{I})^{-1} \\ \tilde{\mathbf{P}}_n &= (\mathbf{K}_n + \gamma \mathbf{I})^{-1/2} \mathbf{K}_n^{1/2} \mathbf{S}_n \mathbf{S}_n^T \mathbf{K}_n^{1/2} (\mathbf{K}_n + \gamma \mathbf{I})^{-1/2}\end{aligned}$$

Reconstruction guarantees

Why would bounding $\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2$ be useful?

Reconstruction guarantees

Why would bounding $\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2$ be useful?

It appears in many problems e.g., Kernel Ridge Regression

Reconstruction guarantees

Why would bounding $\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2$ be useful?

It appears in many problems e.g., Kernel Ridge Regression

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

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

Reconstruction guarantees

Why would bounding $\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2$ be useful?

We can compute accurate low rank approximations. Let

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

then

$$\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2 \leq \varepsilon \Rightarrow \tilde{\mathbf{K}}_n \preceq \mathbf{K}_n \preceq \tilde{\mathbf{K}}_n + \frac{\gamma}{1 - \varepsilon} \mathbf{I}$$

Reconstruction guarantees

Why would bounding $\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2$ be useful?

We can compute accurate low rank approximations. Let

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

then

$$\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2 \leq \varepsilon \Rightarrow \tilde{\mathbf{K}}_n \preceq \mathbf{K}_n \preceq \tilde{\mathbf{K}}_n + \frac{\gamma}{1-\varepsilon} \mathbf{I}$$

e.g., Kernel Ridge Regression

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

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

~~$\mathcal{O}(n^3)$~~ \Rightarrow $\mathcal{O}(nm + m^3)$ time to compute the approx. solution

~~$\mathcal{O}(n^2)$~~ \Rightarrow $\mathcal{O}(nm)$ space to store dictionary

Reconstruction guarantees

Why would bounding $\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2$ be useful?

We can compute accurate low rank approximations. Let

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

then

$$\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2 \leq \varepsilon \Rightarrow \tilde{\mathbf{K}}_n \preceq \mathbf{K}_n \preceq \tilde{\mathbf{K}}_n + \frac{\gamma}{1 - \varepsilon} \mathbf{I}$$

e.g., Kernel Ridge Regression*

*Gaussian Processes

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

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

~~$\mathcal{O}(n^3)$~~ \Rightarrow $\mathcal{O}(nm + m^3)$ time to compute the approx. solution

~~$\mathcal{O}(n^2)$~~ \Rightarrow $\mathcal{O}(nm)$ space to store dictionary

Reconstruction guarantees

Why would bounding $\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2$ be useful?

We can compute accurate low rank approximations. Let

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

then

$$\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2 \leq \varepsilon \Rightarrow \tilde{\mathbf{K}}_n \preceq \mathbf{K}_n \preceq \tilde{\mathbf{K}}_n + \frac{\gamma}{1 - \varepsilon} \mathbf{I}$$

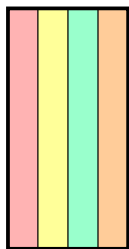
e.g., Kernel PCA, \mathbf{K}_n and $\tilde{\mathbf{K}}_n$ have close leading eigenvalues/vectors

e.g., Kernel K -means can be formulated as a quadratic form

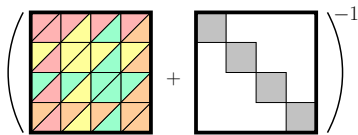
$$\min_{\mathbf{C}} \text{Tr}(\mathbf{K}_n - \mathbf{C}\mathbf{C}^T \mathbf{K}_n \mathbf{C}\mathbf{C}^T) \sim \min_{\tilde{\mathbf{C}}} \text{Tr}(\tilde{\mathbf{K}}_n - \tilde{\mathbf{C}}\tilde{\mathbf{C}}^T \tilde{\mathbf{K}}_n \tilde{\mathbf{C}}\tilde{\mathbf{C}}^T)$$

Regularized Nyström reconstruction

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



$$\mathbf{C} = \mathbf{K}_n \mathbf{S}_n$$



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



$$\mathbf{C}^\top = \mathbf{S}_n^\top \mathbf{K}_n$$

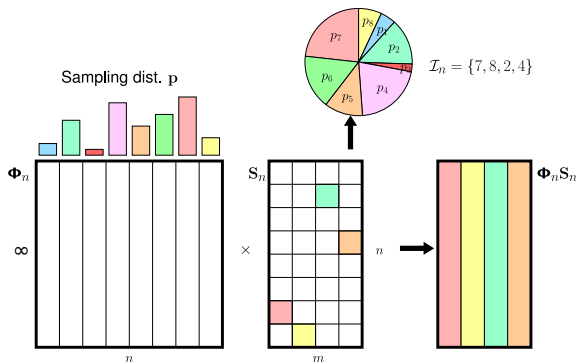
Distributed sequential **sampling** for adaptive kernel DL

How do we compute an accurate ($\|\Psi_n - \tilde{\Psi}_n\|_2 \leq \varepsilon$) dictionary?

Distributed sequential sampling for adaptive kernel DL

How do we compute an accurate ($\|\Psi_n - \tilde{\Psi}_n\|_2 \leq \varepsilon$) dictionary?

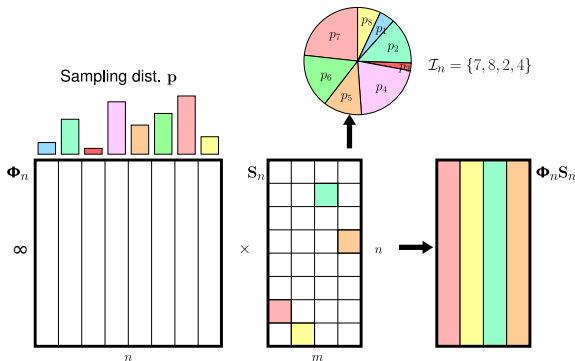
Sample m points w.p. $p_{n,i}$, add to \mathcal{I} with weight $1/p_{n,i}$ (unbiased)



Distributed sequential sampling for adaptive kernel DL

How do we compute an accurate ($\|\Psi_n - \tilde{\Psi}_n\|_2 \leq \varepsilon$) dictionary?

Sample m points w.p. $p_{n,i}$, add to \mathcal{I} with weight $1/p_{n,i}$ (unbiased)



? How to choose the sampling distribution?

? How to choose m ?

Ridge Leverage Scores and Effective Dimension

Definition

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

$$\begin{aligned} \gamma\text{-RLS} \quad \tau_{n,i} &= \mathbf{e}_{n,i} \mathbf{K}_n^T (\mathbf{K}_n + \gamma \mathbf{I}_n)^{-1} \mathbf{e}_{n,i} \\ &= \phi_i^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1} \phi_i \end{aligned} \quad (1)$$

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

Ridge Leverage Scores

Intuitively, RLS capture orthogonality

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

If all ϕ_i are orthogonal, we have

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

If all ϕ_i are identical (collinear), we have

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

Ridge Leverage Scores

Intuitively, RLS capture orthogonality

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

If all ϕ_i are orthogonal, we have

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

If all ϕ_i are identical (collinear), we have

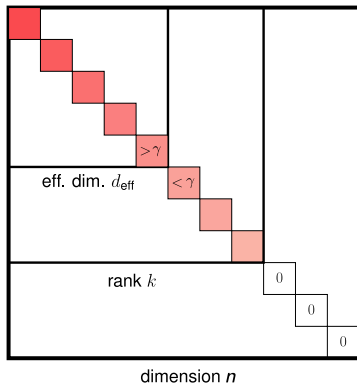
$$\tau_{n,i} = \phi_i^T (\Phi_n \Phi_n^T + \gamma \mathbf{I})^{-1} \phi_i = \phi_i^T (n \phi_i \phi_i^T + \gamma \mathbf{I})^{-1} \phi_i = \frac{\phi_i^T \phi_i}{n \phi_i^T \phi_i + \gamma} \sim \frac{1}{n}$$

Given Φ_{t-1} , adding a new column to it can only reduce the RLS of columns already in Φ_{t-1}

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

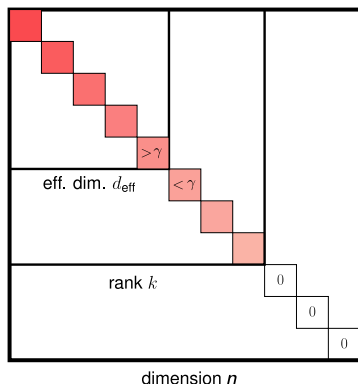
Effective Dimension

Intuitively, the **effective dimension** is a **soft version of matrix rank**



Effective Dimension

Intuitively, the **effective dimension** is a **soft version of matrix rank**



Given $d_{\text{eff}}(\gamma)_{t-1}$, adding a new column to Φ_{t-1} can only increase $d_{\text{eff}}(\gamma)_t$

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

Nyström Sampling

Theorem (Alaoui, Mahoney, 2015)

Given γ be the Nyström regularization, ε the accuracy, δ the confidence.
If the dictionary \mathcal{I}_n is computed using the sampling distribution $p_{n,i} \propto \tau_{n,i}$ and using at least m columns

$$m \geq \left(\frac{2\mathbf{d}_{\text{eff}}(\gamma)_n}{\varepsilon^2} \right) \log \left(\frac{n}{\delta} \right),$$

then with probability $1 - \delta$

$$\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2 \leq \varepsilon$$

Nyström Sampling

Theorem (Alaoui, Mahoney, 2015)

Given γ be the Nyström regularization, ε the accuracy, δ the confidence.
If the dictionary \mathcal{I}_n is computed using the sampling distribution $p_{n,i} \propto \tau_{n,i}$ and using at least m columns

$$m \geq \left(\frac{2\mathbf{d}_{\text{eff}}(\gamma)_n}{\varepsilon^2} \right) \log \left(\frac{n}{\delta} \right),$$

then with probability $1 - \delta$

$$\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2 \leq \varepsilon$$

Done!

Nyström Sampling

Theorem (Alaoui, Mahoney, 2015)

Given γ be the Nyström regularization, ε the accuracy, δ the confidence.
If the dictionary \mathcal{I}_n is computed using the sampling distribution $p_{n,i} \propto \tau_{n,i}$ and using at least m columns

$$m \geq \left(\frac{2\mathbf{d}_{\text{eff}}(\gamma)_n}{\varepsilon^2} \right) \log \left(\frac{n}{\delta} \right),$$

then with probability $1 - \delta$

$$\|\mathbf{P}_n - \tilde{\mathbf{P}}_n\|_2 \leq \varepsilon$$

Done!

If someone gave us the RLS

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

Estimating RLS

Idea 1: Instead of computing exact RLS, compute good approximations

Estimating RLS

Idea 1: Instead of computing **exact** RLS, compute **good approximations**

Idea 2: When all you have is a dictionary, you use the dictionary

Estimating RLS

Idea 1: Instead of computing exact RLS, compute good approximations

Idea 2: When all you have is a dictionary, you use the dictionary

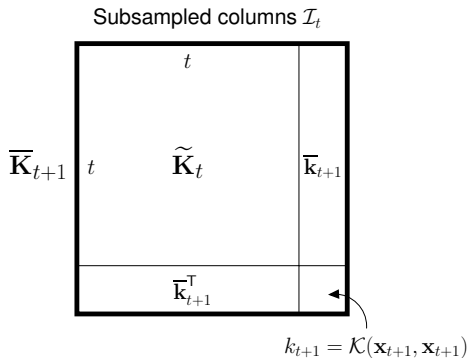
Lemma

Assume that the dictionary \mathcal{I}_{t-1} is accurate, and let $\bar{\mathbf{S}}_t$ be constructed by adding $(1, \phi_t)$ to \mathcal{I}_{t-1} . Then, denoting $\alpha = (1 + \varepsilon)/(1 - \varepsilon)$, for all i such that $i \in \{\mathcal{I}_{t-1} \cup \{t\}\}$,

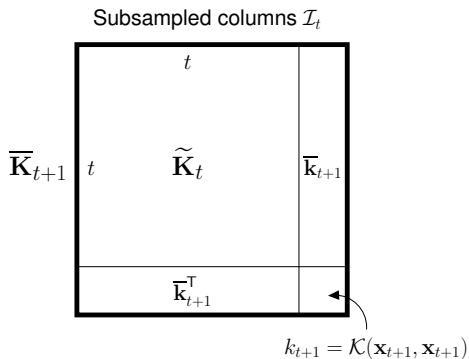
$$\tilde{\tau}_{t,i} = \frac{1 + \varepsilon}{\alpha \gamma} \left(k_{i,i} - \mathbf{k}_{t,i} \bar{\mathbf{S}} \left(\bar{\mathbf{S}}^T \mathbf{K}_t \bar{\mathbf{S}} + \gamma \mathbf{I} \right)^{-1} \bar{\mathbf{S}}^T \mathbf{k}_{t,i} \right), \quad (3)$$

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

The problem of estimating RLS

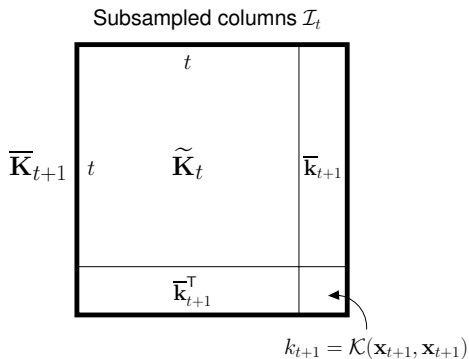


The problem of estimating RLS



Approximate sampling distribution \mathbf{p}_{t+1}

The problem of estimating RLS



Approximate sampling distribution \mathbf{p}_{t+1}

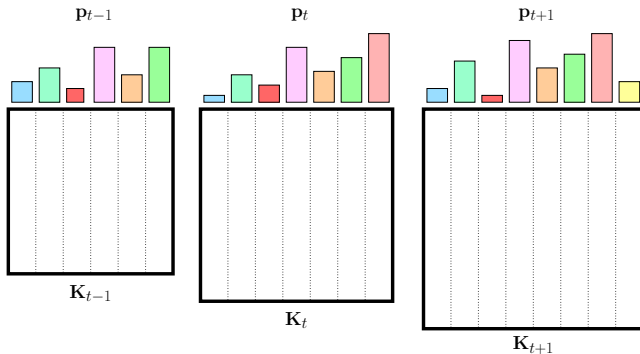
\Rightarrow since $p_{i,t+1} \propto \tau_{i,t+1}$, *approximate* $\tau_{i,t+1}$

Estimating RLS

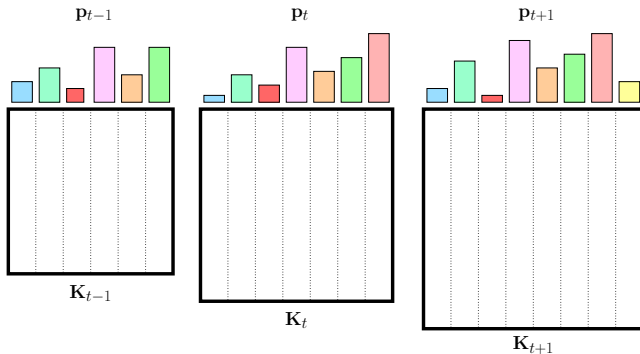
$$\tilde{\tau}_{t,i} = \frac{1 + \varepsilon}{\alpha\gamma} \left(k_{i,i} - \mathbf{k}_{t,i}^T \bar{\mathbf{S}} \left(\bar{\mathbf{S}}^T \mathbf{K}_t \bar{\mathbf{S}} + \gamma \mathbf{I} \right)^{-1} \bar{\mathbf{S}}^T \mathbf{k}_{t,i} \right),$$

- ▶ $\tilde{\tau}_{t,i} = \mathbf{e}_i^T \tilde{\mathbf{K}}_t (\tilde{\mathbf{K}}_t + \gamma \mathbf{I})^{-1} \mathbf{e}_i$ would fail
- ▶ Instead, approximate $\tau_{t,i}$ directly in \mathcal{H} , and then reformulate using kernel trick $\tilde{\tau}_{t,i} = \phi_i^T (\Phi \bar{\mathbf{S}} \bar{\mathbf{S}}^T \Phi^T + \gamma \mathbf{I})^{-1} \phi_i$
- ▶ $\tilde{\tau}_{t,i}$ can be computed in $\mathcal{O}(|\mathcal{I}_t|^2)$ space and $\mathcal{O}(|\mathcal{I}_t|^3)$ time
↳ independent from t
- ▶ $\tilde{\tau}_{t,i}$ for $i \in \mathcal{I}_t$ can be computed using only samples contained in \mathcal{I}_t .

Estimating RLS incrementally

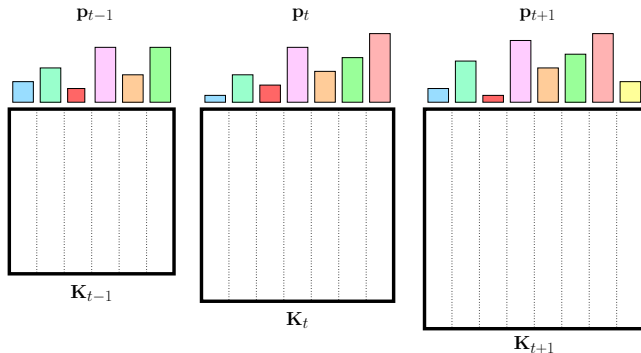


Estimating RLS incrementally



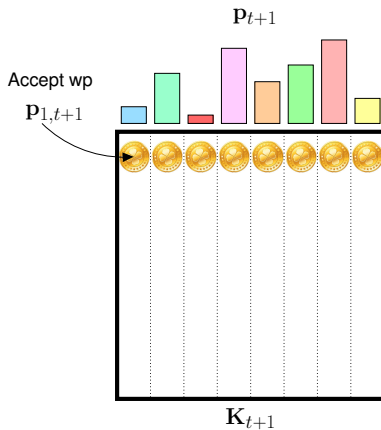
At each time step t construct $\tilde{\mathbf{K}}_t$ **as if** it was drawn from \mathbf{p}_t

Estimating RLS incrementally

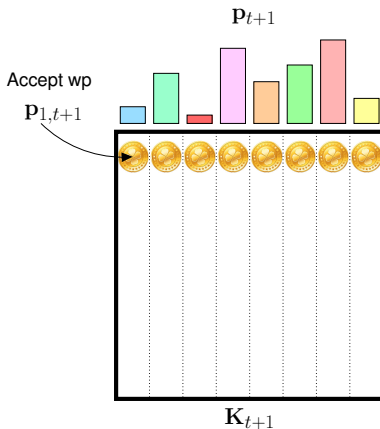


At each time step t construct \tilde{K}_t *as if* it was drawn from p_t
 \Rightarrow update the sampling set \mathcal{I}_t incrementally as p_t changes

Estimating RLS incrementally by rejection sampling



Estimating RLS incrementally by rejection sampling



m calls to a multinomial \mathbf{p}_{t+1}

\approx calls to $t + 1$ binomials each with probability $p_{i,t+1}$

Distributed **sequential** sampling for adaptive kernel DL

Instead of sampling from multinomial consider the sampling process

$$q_{i,i} \sim \mathcal{B}(\tilde{p}_{i,i}, \bar{q})$$

$$q_{t,i} \sim \mathcal{B}(\tilde{p}_{t,i}/\tilde{p}_{t-1,i}, q_{t-1,i})$$

Distributed sequential sampling for adaptive kernel DL

Instead of sampling from multinomial consider the sampling process

$$\begin{aligned}q_{i,i} &\sim \mathcal{B}(\tilde{p}_{i,i}, \bar{q}) \\q_{t,i} &\sim \mathcal{B}(\tilde{p}_{t,i}/\tilde{p}_{t-1,i}, q_{t-1,i})\end{aligned}$$

Similar to importance sampling. If the $\tilde{p}_{t,i}$ were fixed in advance

$$\begin{aligned}\mathbb{P}(z_{t,i,j} = 1) &= \mathbb{P}(\mathcal{B}(\tilde{p}_{t,i}/\tilde{p}_{t-1,i}) = 1)z_{t-1,i,j} \\&= \mathbb{P}(\mathcal{B}(\tilde{p}_{t,i}/\tilde{p}_{t-1,i}) = 1)\mathbb{P}(\mathcal{B}(\tilde{p}_{t-1,i}/\tilde{p}_{t-2,i}) = 1)z_{t-2,i,j} \\&= \frac{\tilde{p}_{t,i}}{\tilde{p}_{t-1,i}} \frac{\tilde{p}_{t-1,i}}{\tilde{p}_{t-2,i}} \dots \frac{\tilde{p}_{i+1,i}}{\tilde{p}_{i,i}} \frac{\tilde{p}_{i,i}}{1} = \tilde{p}_{t,i}\end{aligned}$$

Distributed sequential sampling for adaptive kernel DL

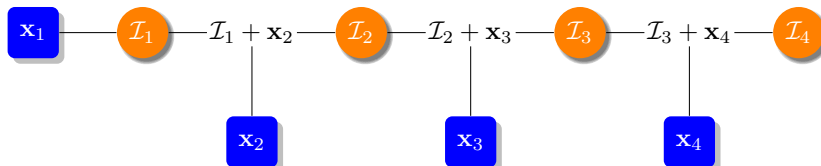
Instead of sampling from multinomial consider the sampling process

$$q_{i,i} \sim \mathcal{B}(\tilde{p}_{i,i}, \bar{q})$$

$$q_{t,i} \sim \mathcal{B}(\tilde{p}_{t,i}/\tilde{p}_{t-1,i}, q_{t-1,i})$$

Similar to importance sampling. If the $\tilde{p}_{t,i}$ were fixed in advance

$$\begin{aligned} \mathbb{P}(z_{t,i,j} = 1) &= \mathbb{P}(\mathcal{B}(\tilde{p}_{t,i}/\tilde{p}_{t-1,i}) = 1)z_{t-1,i,j} \\ &= \mathbb{P}(\mathcal{B}(\tilde{p}_{t,i}/\tilde{p}_{t-1,i}) = 1)\mathbb{P}(\mathcal{B}(\tilde{p}_{t-1,i}/\tilde{p}_{t-2,i}) = 1)z_{t-2,i,j} \\ &= \frac{\tilde{p}_{t,i}}{\tilde{p}_{t-1,i}} \frac{\tilde{p}_{t-1,i}}{\tilde{p}_{t-2,i}} \dots \frac{\tilde{p}_{i+1,i}}{\tilde{p}_{i,i}} \frac{\tilde{p}_{i,i}}{1} = \tilde{p}_{t,i} \end{aligned}$$



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 α -app. RLS $\{\tilde{\tau}_{t,i} : i \in \mathcal{I}_{t-1} \cup \{t\}\}$, using \mathcal{I}_{t-1} , \mathbf{x} , and Eq. 3
 - 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

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

(1) $\|\mathbf{P}_t - \tilde{\mathbf{P}}_t\|_2 \leq \varepsilon$.

(2) $|\mathcal{I}_t| = \sum_i q_{t,i} \leq \mathcal{O}(\bar{q} d_{\text{eff}}(\gamma)_t) \leq \mathcal{O}\left(\frac{\alpha}{\varepsilon^2} \mathbf{d}_{\text{eff}}(\gamma)_n \log\left(\frac{n}{\delta}\right)\right)$.

SQUEAK

Theorem

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

(1) $\|\mathbf{P}_t - \tilde{\mathbf{P}}_t\|_2 \leq \varepsilon$.

(2) $|\mathcal{I}_t| = \sum_i q_{t,i} \leq \mathcal{O}(\bar{\mathbf{q}} d_{\text{eff}}(\gamma)_t) \leq \mathcal{O}\left(\frac{\alpha}{\varepsilon^2} \mathbf{d}_{\text{eff}}(\gamma)_n \log\left(\frac{n}{\delta}\right)\right)$.

- ▶ Accuracy and space/time guarantees
- ▶ Anytime risk 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 $\tilde{p}_{t,i}$, no need for appr. $d_{\text{eff}}(\gamma)_t$

SQUEAK

Theorem

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

(1) $\|\mathbf{P}_t - \tilde{\mathbf{P}}_t\|_2 \leq \varepsilon$.

(2) $|\mathcal{I}_t| = \sum_i q_{t,i} \leq \mathcal{O}(\bar{\mathbf{q}} d_{\text{eff}}(\gamma)_t) \leq \mathcal{O}\left(\frac{\alpha}{\varepsilon^2} \mathbf{d}_{\text{eff}}(\gamma)_n \log\left(\frac{n}{\delta}\right)\right)$.

- ▶ 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 $\mathcal{O}(n^3) \Rightarrow \mathcal{O}(n|\mathcal{I}_n|^3) \leq \tilde{\mathcal{O}}(n \mathbf{d}_{\text{eff}}(\gamma)_n^3)$
- ▶ Store points directly in the dictionary
 - ↳ $\tilde{\mathcal{O}}(\mathbf{d}_{\text{eff}}(\gamma)_n^2 + \mathbf{d}_{\text{eff}}(\gamma)_n \mathbf{d})$ space constant in n
 - ↳ single pass over the dataset (streaming)

Proof sketch

Need to bound

$$\mathbb{P}\left(\exists t \in \{1, \dots, n\} : \|\mathbf{P}_t - \tilde{\mathbf{P}}_t\|_2 \geq \varepsilon \cup |\mathcal{I}_t| \geq 3\bar{q}d_{\text{eff}}(\gamma)_t\right)$$

Proof sketch

Need to bound

$$\mathbb{P}\left(\exists t \in \{1, \dots, n\} : \|\mathbf{P}_t - \tilde{\mathbf{P}}_t\|_2 \geq \varepsilon \cup |\mathcal{I}_t| \geq 3\bar{q}d_{\text{eff}}(\gamma)_t\right)$$

After a union bound

$$\begin{aligned} & \sum_{t=1}^n \mathbb{P}\left(\|\mathbf{P}_t - \tilde{\mathbf{P}}_t\|_2 \geq \varepsilon\right) \\ & + \sum_{t=1}^n \mathbb{P}\left(|\mathcal{I}_t| \geq 3\bar{q}d_{\text{eff}}(\gamma)_t \cap \left\{\forall t' \in \{1, \dots, t\} : \|\mathbf{P}_{t'} - \tilde{\mathbf{P}}_{t'}\|_2 \leq \varepsilon\right\}\right) \end{aligned}$$

Proof sketch

We start by bounding $\mathbb{P} \left(\|\mathbf{P}_t - \tilde{\mathbf{P}}_t\|_2 \geq \varepsilon \right)$. Let

$$z_{s,i,j} = \mathbb{I} \left\{ u_{s,i,j} \leq \frac{\tilde{p}_{s,i}}{\tilde{p}_{s-1,i}} \right\} z_{s-1,i,j}, \quad \mathbf{v}_i = (\mathbf{K}_t + \gamma \mathbf{I})^{-1} \mathbf{K}_t^{1/2} \mathbf{e}_{t,i}$$

with $u_{s,i,j} \sim \mathcal{U}(0, 1)$. Then

$$\mathbf{Y}_t = \mathbf{P}_t - \tilde{\mathbf{P}}_t = \frac{1}{\bar{q}} \sum_{i=1}^t \sum_{j=1}^{\bar{q}} \left(1 - \frac{z_{t,i,j}}{\tilde{p}_{t,i}} \right) \mathbf{v}_i \mathbf{v}_i^\top$$

Proof sketch

We start by bounding $\mathbb{P} \left(\|\mathbf{P}_t - \tilde{\mathbf{P}}_t\|_2 \geq \varepsilon \right)$. Let

$$z_{s,i,j} = \mathbb{I} \left\{ u_{s,i,j} \leq \frac{\tilde{p}_{s,i}}{\tilde{p}_{s-1,i}} \right\} z_{s-1,i,j}, \quad \mathbf{v}_i = (\mathbf{K}_t + \gamma \mathbf{I})^{-1} \mathbf{K}_t^{1/2} \mathbf{e}_{t,i}$$

with $u_{s,i,j} \sim \mathcal{U}(0, 1)$. Then

$$\mathbf{Y}_t = \mathbf{P}_t - \tilde{\mathbf{P}}_t = \frac{1}{\bar{q}} \sum_{i=1}^t \sum_{j=1}^{\bar{q}} \left(1 - \frac{z_{t,i,j}}{\tilde{p}_{t,i}} \right) \mathbf{v}_i \mathbf{v}_i^\top$$

Cannot use concentrations for independent r.v., because $z_{t,i,j}$ and $z_{t,i',j'}$ are both dependent on $z_{t-1,i'',j''}$ through the estimates.

Proof sketch

Build the martingale

$$\mathbf{X}_{\{s,i,j\}} = \left(\frac{Z_{s-1,i,j}}{\tilde{p}_{s-1,i}} - \frac{Z_{t,i,j}}{\tilde{p}_{s,i}} \right) \mathbf{v}_i \mathbf{v}_i^\top$$

We can use variants of Bernstein's inequality for matrix martingales, we need a bound on the range

$$\begin{aligned} \|\mathbf{X}_{\{s,i,j\}}\| &= \frac{1}{\tilde{q}} \left\| \left(\frac{Z_{s-1,i,j}}{\tilde{p}_{s-1,i}} - \frac{Z_{t,i,j}}{\tilde{p}_{s,i}} \right) \right\| \|\mathbf{v}_i \mathbf{v}_i^\top\| \leq \frac{1}{\tilde{q}} \frac{1}{\tilde{p}_{s,i}} \|\mathbf{v}_i\|^2 \\ &\leq \frac{1}{\tilde{q}} \frac{1}{\tilde{p}_{s,i}} \mathbf{v}_i^\top \mathbf{v}_i = \frac{1}{\tilde{q}} \frac{1}{\tilde{p}_{s,i}} \mathbf{e}_i^\top \mathbf{K}_t^{1/2} (\mathbf{K}_t + \gamma \mathbf{I})^{-1} \mathbf{K}_t^{1/2} \mathbf{e}_i \\ &= \frac{1}{\tilde{q}} \frac{1}{\tilde{p}_{s,i}} \mathbf{e}_i^\top \mathbf{P}_t \mathbf{e}_i = \frac{1}{\tilde{q}} \frac{\tau_{t,i}}{\tilde{p}_{s,i}} \leq \frac{\alpha}{\tilde{q}} \frac{\tau_{t,i}}{p_{s,i}} = \frac{\alpha}{\tilde{q}} \frac{\tau_{t,i}}{\tau_{s,i}} \leq \frac{\alpha}{\tilde{q}} := R, \end{aligned}$$

Proof sketch

Build the martingale

$$\mathbf{X}_{\{s,i,j\}} = \left(\frac{z_{s-1,i,j}}{\tilde{p}_{s-1,i}} - \frac{z_{t,i,j}}{\tilde{p}_{s,i}} \right) \mathbf{v}_i \mathbf{v}_i^\top$$

We can use variants of Bernstein's inequality for matrix martingales, we need a bound on the range

$$\begin{aligned} \|\mathbf{X}_{\{s,i,j\}}\| &= \frac{1}{q} \left\| \left(\frac{z_{s-1,i,j}}{\tilde{p}_{s-1,i}} - \frac{z_{t,i,j}}{\tilde{p}_{s,i}} \right) \right\| \|\mathbf{v}_i \mathbf{v}_i^\top\| \leq \frac{1}{q} \frac{1}{\tilde{p}_{s,i}} \|\mathbf{v}_i\|^2 \\ &\leq \frac{1}{q} \frac{1}{\tilde{p}_{s,i}} \mathbf{v}_i^\top \mathbf{v}_i = \frac{1}{q} \frac{1}{\tilde{p}_{s,i}} \mathbf{e}_i^\top \mathbf{K}_t^{1/2} (\mathbf{K}_t + \gamma \mathbf{I})^{-1} \mathbf{K}_t^{1/2} \mathbf{e}_i \\ &= \frac{1}{q} \frac{1}{\tilde{p}_{s,i}} \mathbf{e}_i^\top \mathbf{P}_t \mathbf{e}_i = \frac{1}{q} \frac{\tau_{t,i}}{\tilde{p}_{s,i}} \leq \frac{\alpha}{q} \frac{\tau_{t,i}}{p_{s,i}} = \frac{\alpha}{q} \frac{\tau_{t,i}}{\tau_{s,i}} \leq \frac{\alpha}{q} := R, \end{aligned}$$

RLS normalize our r.v.

Proof sketch

Now bound the total variation

$$\begin{aligned} \mathbf{W} &= \sum \mathbb{E} \left[\mathbf{x}_{\{s,i,j\}}^2 \mid \{\mathbf{x}_r\}_{r=0}^{\{s,i,j\}-1} \right] \\ &= \frac{1}{q^2} \sum_{j=1}^{\bar{q}} \sum_{i=1}^t \sum_{s=1}^t \frac{z_{s-1,i,j}}{\tilde{p}_{s-1,i}} \left(\frac{1}{\tilde{p}_{s,i}} - \frac{1}{\tilde{p}_{s-1,i}} \right) \mathbf{v}_i \mathbf{v}_i^\top \mathbf{v}_i \mathbf{v}_i^\top \end{aligned}$$

Proof sketch

Now bound the total variation

$$\begin{aligned}\mathbf{W} &= \sum \mathbb{E} \left[\mathbf{X}_{\{s,i,j\}}^2 \mid \{\mathbf{X}_r\}_{r=0}^{\{s,i,j\}-1} \right] \\ &= \frac{1}{q^2} \sum_{j=1}^{\bar{q}} \sum_{i=1}^t \sum_{s=1}^t \frac{z_{s-1,i,j}}{\tilde{\rho}_{s-1,i}} \left(\frac{1}{\tilde{\rho}_{s,i}} - \frac{1}{\tilde{\rho}_{s-1,i}} \right) \mathbf{v}_i \mathbf{v}_i^\top \mathbf{v}_i \mathbf{v}_i^\top\end{aligned}$$

Deterministically

$$\begin{aligned}\|\mathbf{W}\| &= \left\| \frac{1}{q^2} \sum_{j=1}^{\bar{q}} \sum_{i=1}^t \sum_{s=1}^t \frac{z_{s-1,i,j}}{\tilde{\rho}_{s-1,i}} \left(\frac{1}{\tilde{\rho}_{s,i}} - \frac{1}{\tilde{\rho}_{s-1,i}} \right) \mathbf{v}_i \mathbf{v}_i^\top \mathbf{v}_i \mathbf{v}_i^\top \right\| \\ &\leq \left\| \frac{1}{q^2} \sum_{j=1}^{\bar{q}} \sum_{i=1}^t \frac{\mathbf{v}_i^\top \mathbf{v}_i}{\tilde{\rho}_{t,i}^2} \mathbf{v}_i \mathbf{v}_i^\top \right\| \leq \left\| \frac{\alpha}{q} \sum_{i=1}^t \frac{1}{\tilde{\rho}_{t,i}} \mathbf{v}_i \mathbf{v}_i^\top \right\| \\ &\leq \left\| \frac{\alpha^2}{q} \sum_{i=1}^t \mathbf{I} \right\| = \frac{\alpha^2}{q} t\end{aligned}$$

Proof sketch

Now bound the total variation

$$\begin{aligned}\mathbf{W} &= \sum \mathbb{E} \left[\mathbf{X}_{\{s,i,j\}}^2 \mid \{\mathbf{X}_r\}_{r=0}^{\{s,i,j\}-1} \right] \\ &= \frac{1}{q^2} \sum_{j=1}^{\bar{q}} \sum_{i=1}^t \sum_{s=1}^t \frac{z_{s-1,i,j}}{\tilde{\rho}_{s-1,i}} \left(\frac{1}{\tilde{\rho}_{s,i}} - \frac{1}{\tilde{\rho}_{s-1,i}} \right) \mathbf{v}_i \mathbf{v}_i^\top \mathbf{v}_i \mathbf{v}_i^\top\end{aligned}$$

Deterministically

$$\begin{aligned}\|\mathbf{W}\| &= \left\| \frac{1}{q^2} \sum_{j=1}^{\bar{q}} \sum_{i=1}^t \sum_{s=1}^t \frac{z_{s-1,i,j}}{\tilde{\rho}_{s-1,i}} \left(\frac{1}{\tilde{\rho}_{s,i}} - \frac{1}{\tilde{\rho}_{s-1,i}} \right) \mathbf{v}_i \mathbf{v}_i^\top \mathbf{v}_i \mathbf{v}_i^\top \right\| \\ &\leq \left\| \frac{1}{q^2} \sum_{j=1}^{\bar{q}} \sum_{i=1}^t \frac{\mathbf{v}_i^\top \mathbf{v}_i}{\tilde{\rho}_{t,i}^2} \mathbf{v}_i \mathbf{v}_i^\top \right\| \leq \left\| \frac{\alpha}{q} \sum_{i=1}^t \frac{1}{\tilde{\rho}_{t,i}} \mathbf{v}_i \mathbf{v}_i^\top \right\| \\ &\leq \left\| \frac{\alpha^2}{q} \sum_{i=1}^t \mathbf{I} \right\| = \frac{\alpha^2}{q} t\end{aligned}$$

Deterministic bound on variance too large

Proof sketch

This looks **too pessimistic**. When $\frac{1}{\bar{p}_{s,i}}$ is large, $z_{s,i,j}$ should be zero.
We should take advantage of that.

Proof sketch

This looks **too pessimistic**. When $\frac{1}{\bar{p}_{s,i}}$ is large, $z_{s,i,j}$ should be zero.
We should take advantage of that.

We can use a finer concentration, Freedman's inequality, that treats \mathbf{W} itself as a random variable.

$$\mathbb{P}(\|\mathbf{Y}_t\| \geq \varepsilon \cap \|\mathbf{W}\| \leq \sigma^2) \leq t \exp\{-\dots\}$$

Proof sketch

This looks **too pessimistic**. When $\frac{1}{\tilde{p}_{s,i}}$ is large, $z_{s,i,j}$ should be zero. We should take advantage of that.

We can use a finer concentration, Freedman's inequality, that treats \mathbf{W} itself as a random variable.

$$\mathbb{P}(\|\mathbf{Y}_t\| \geq \varepsilon \cap \|\mathbf{W}\| \leq \sigma^2) \leq t \exp\{-\dots\}$$

Starting from an upper bound on \mathbf{W} that is still a r.v.

$$\mathbf{W} \preceq \frac{1}{\bar{q}^2} \sum_{j=1}^{\bar{q}} \sum_{i=1}^t \max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\tilde{p}_{s,i}^2} \right\} \mathbf{v}_i \mathbf{v}_i^\top \mathbf{v}_i \mathbf{v}_i^\top$$

Proof sketch

This looks **too pessimistic**. When $\frac{1}{\bar{\rho}_{s,i}}$ is large, $z_{s,i,j}$ should be zero. We should take advantage of that.

We can use a finer concentration, Freedman's inequality, that treats \mathbf{W} itself as a random variable.

$$\mathbb{P}(\|\mathbf{Y}_t\| \geq \varepsilon \cap \|\mathbf{W}\| \leq \sigma^2) \leq t \exp\{-\dots\}$$

Starting from an upper bound on \mathbf{W} that is still a r.v.

$$\mathbf{W} \preceq \frac{1}{\bar{q}^2} \sum_{j=1}^{\bar{q}} \sum_{i=1}^t \max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\bar{\rho}_{s,i}^2} \right\} \mathbf{v}_i \mathbf{v}_i^\top \mathbf{v}_i \mathbf{v}_i^\top$$

This still has high variance: cannot simply apply martingale Bernstein

Proof sketch

$\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{p_{s,i}^2} \right\}$ is still hard to analyze, since it is the
maximum of dependent variables

Proof sketch

$\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\tilde{p}_{s,i}^2} \right\}$ is still hard to analyze, since it is the
maximum of dependent variables

Moreover $\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\tilde{p}_{s,i}^2} \right\}$ depends on $\max_{s=0}^{t-1} \left\{ \frac{z_{s,i',j'}}{\tilde{p}_{s,i'}^2} \right\}$

Proof sketch

$\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\tilde{p}_{s,i}^2} \right\}$ is still hard to analyze, since it is the
maximum of dependent variables

Moreover $\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\tilde{p}_{s,i}^2} \right\}$ depends on $\max_{s=0}^{t-1} \left\{ \frac{z_{s,i',j'}}{\tilde{p}_{s,i'}^2} \right\}$

We will find another set of dominating r.v. $1/w_{i,j}$, indep. from each other
Then apply Bernstein for indep. r.v.

Proof sketch

$\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\tilde{p}_{s,i}^2} \right\}$ is still hard to analyze, since it is the
maximum of dependent variables

Moreover $\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\tilde{p}_{s,i}^2} \right\}$ depends on $\max_{s=0}^{t-1} \left\{ \frac{z_{s,i',j'}}{\tilde{p}_{s,i'}^2} \right\}$

We will find another set of dominating r.v. $1/w_{i,j}$, indep. from each other
Then apply Bernstein for indep. r.v.

Random variable A stochastically dominates random variable B , if for all values a the two equivalent conditions are verified

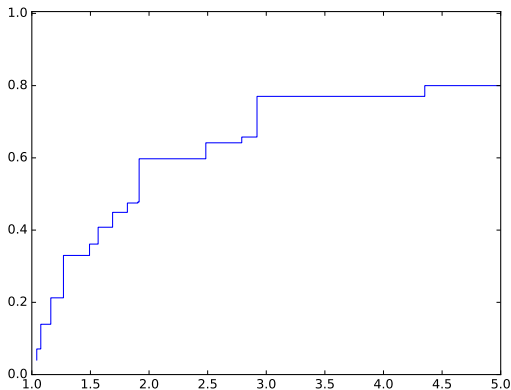
$$\mathbb{P}(A \geq a) \geq \mathbb{P}(B \geq a) \Leftrightarrow \mathbb{P}(A \leq a) \leq \mathbb{P}(B \leq a).$$

Proof sketch

Imagine the sequence $\tilde{p}_{s,i}$ was fixed in advance. I can compute exactly the distribution of all $z_{s,i,j}$.

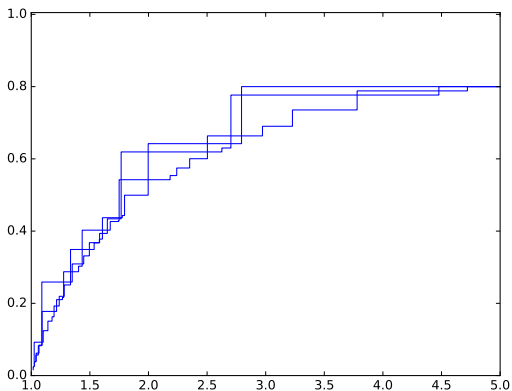
Proof sketch

Imagine the sequence $\tilde{p}_{s,i}$ was fixed in advance. I can compute exactly the distribution of all $z_{s,i,j}$.



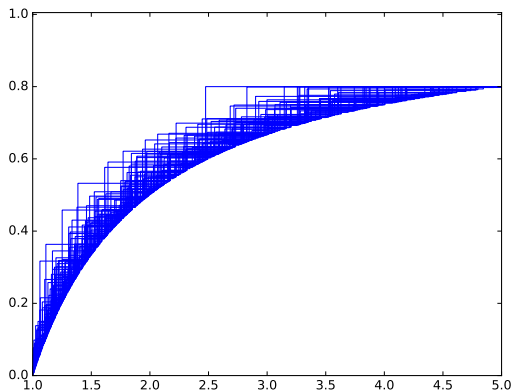
Proof sketch

Imagine the sequence $\tilde{p}_{s,i}$ was fixed in advance. I can compute exactly the distribution of all $z_{s,i,j}$.



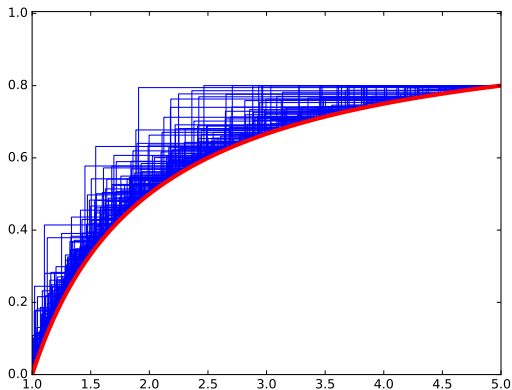
Proof sketch

Imagine the sequence $\tilde{p}_{s,i}$ was fixed in advance. I can compute exactly the distribution of all $z_{s,i,j}$.



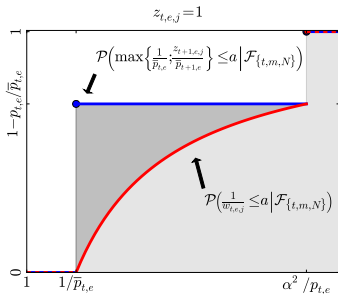
Proof sketch

Imagine the sequence $\tilde{p}_{s,i}$ was fixed in advance. I can compute exactly the distribution of all $z_{s,i,j}$.



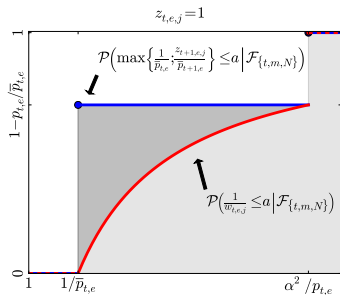
Proof sketch

Imagine the sequence $\tilde{p}_{s,i}$ was fixed in advance. I can compute exactly the distribution of all $z_{s,i,j}$.



Proof sketch

Imagine the sequence $\tilde{p}_{s,i}$ was fixed in advance. I can compute exactly the distribution of all $z_{s,i,j}$.



$$\mathbb{P}\left(\frac{1}{w_{0,i,j}} \leq a\right) = \begin{cases} 0 & \text{for } a < 1 \\ 1 - \frac{1}{a} & \text{for } 1 \leq a < \alpha/p_{t,i} \\ 1 & \text{for } \alpha/p_{t,i} \leq a \end{cases}$$

Proof sketch

We can now unwind the proof

$$5 \text{ dominate } \max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\bar{\rho}_{s,i}^2} \right\} \text{ with } 1/w_{i,j}$$

Proof sketch

We can now unwind the proof

5 dominate $\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\bar{\rho}_{s,i}^2} \right\}$ with $1/w_{i,j}$

4 apply Bernstein inequality for indep. r.v. to bound $\mathbb{P}(\|\mathbf{W}\| \geq \sigma^2)$

Proof sketch

We can now unwind the proof

5 dominate $\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\bar{\rho}_{s,i}^2} \right\}$ with $1/w_{i,j}$

4 apply Bernstein inequality for indep. r.v. to bound $\mathbb{P}(\|\mathbf{W}\| \geq \sigma^2)$

3 apply Freedman inequality to bound $\mathbb{P}(\|\mathbf{Y}\| \geq \varepsilon \cap \|\mathbf{W}\| \leq \sigma^2)$

Proof sketch

We can now unwind the proof

5 dominate $\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\bar{\rho}_{s,i}^2} \right\}$ with $1/w_{i,j}$

4 apply Bernstein inequality for indep. r.v. to bound $\mathbb{P}(\|\mathbf{W}\| \geq \sigma^2)$

3 apply Freedman inequality to bound $\mathbb{P}(\|\mathbf{Y}\| \geq \varepsilon \cap \|\mathbf{W}\| \leq \sigma^2)$

2 apply another stochastic dominance argument to bound

$$\mathbb{P}\left(|\mathcal{I}_t| \geq 3\bar{q}d_{\text{eff}}(\gamma)_t \cap \left\{ \forall t' \in \{1, \dots, t\} : \|\mathbf{P}_{t'} - \tilde{\mathbf{P}}_{t'}\|_2 \leq \varepsilon \right\}\right)$$

Proof sketch

We can now unwind the proof

5 dominate $\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\bar{\rho}_{s,i}^2} \right\}$ with $1/w_{i,j}$

4 apply Bernstein inequality for indep. r.v. to bound $\mathbb{P}(\|\mathbf{W}\| \geq \sigma^2)$

3 apply Freedman inequality to bound $\mathbb{P}(\|\mathbf{Y}\| \geq \varepsilon \cap \|\mathbf{W}\| \leq \sigma^2)$

2 apply another stochastic dominance argument to bound
 $\mathbb{P} \left(|\mathcal{I}_t| \geq 3\bar{q}d_{\text{eff}}(\gamma)_t \cap \left\{ \forall t' \in \{1, \dots, t\} : \|\mathbf{P}_{t'} - \tilde{\mathbf{P}}_{t'}\|_2 \leq \varepsilon \right\} \right)$

1 union bound

Proof sketch

We can now unwind the proof

5 dominate $\max_{s=0}^{t-1} \left\{ \frac{z_{s,i,j}}{\bar{\rho}_{s,i}^2} \right\}$ with $1/w_{i,j}$

4 apply Bernstein inequality for indep. r.v. to bound $\mathbb{P}(\|\mathbf{W}\| \geq \sigma^2)$

3 apply Freedman inequality to bound $\mathbb{P}(\|\mathbf{Y}\| \geq \varepsilon \cap \|\mathbf{W}\| \leq \sigma^2)$

2 apply another stochastic dominance argument to bound
 $\mathbb{P}\left(|\mathcal{I}_t| \geq 3\bar{q}d_{\text{eff}}(\gamma)_t \cap \left\{ \forall t' \in \{1, \dots, t\} : \|\mathbf{P}_{t'} - \tilde{\mathbf{P}}_{t'}\|_2 \leq \varepsilon \right\}\right)$

1 union bound

0 Q.E.D.

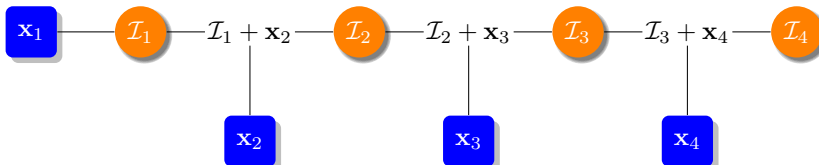
Distributed sequential sampling for adaptive kernel DL

SQUEAK is a strictly sequential algorithm

Distributed sequential sampling for adaptive kernel DL

SQUEAK is a strictly sequential algorithm

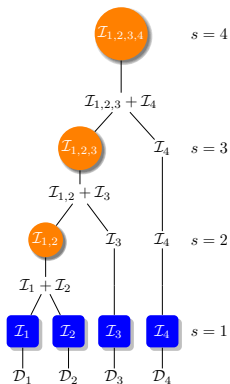
We just did a sequential analysis



Distributed sequential sampling for adaptive kernel DL

SQUEAK is a strictly sequential algorithm

We just did a sequential analysis



Distributed sequential sampling for adaptive kernel DL

SQUEAK is a strictly sequential algorithm

Distributed sequential sampling for adaptive kernel DL

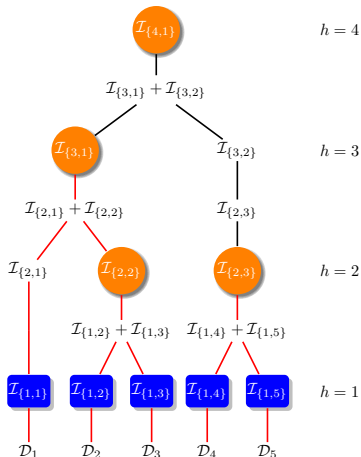
SQUEAK is a strictly sequential algorithm

DISQUEAK is the distributed equivalent

Distributed sequential sampling for adaptive kernel DL

SQUEAK is a strictly sequential algorithm

DISQUEAK is the distributed equivalent



DISQUEAK

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

- 1: Partition \mathcal{D} into disjoint sub-datasets \mathcal{D}_i
- 2: Run SQUEAK on each \mathcal{D}_i , build set $\mathcal{S}_1 = \{\mathcal{I}_{\mathcal{D}_i}\}_{i=1}^k$
- 3: **for** $h = 1, \dots, k - 1$ **do**
- 4: **if** $|\mathcal{S}_h| > 1$ **then** ▷DICT-MERGE
- 5: Pick two dictionaries $\mathcal{I}_{\mathcal{D}}, \mathcal{I}_{\mathcal{D}'}$ from \mathcal{S}_h
- 6: $\bar{\mathcal{I}} = \mathcal{I}_{\mathcal{D}} \cup \mathcal{I}_{\mathcal{D}'}$
- 7: $\mathcal{I}_{\mathcal{D}, \mathcal{D}'} = \text{DICT-UPDATE}(\bar{\mathcal{I}})$ using Eq. (4)
- 8: Place $\mathcal{I}_{\mathcal{D}, \mathcal{D}'}$ back into \mathcal{S}_{h+1}
- 9: **else**
- 10: $\mathcal{S}_{h+1} = \mathcal{S}_h$
- 11: **end if**
- 12: **end for**
- 13: Return $\mathcal{I}_{\mathcal{D}}$, the last dictionary in \mathcal{S}_k

$$\tilde{\tau}_{\mathcal{D} \cup \mathcal{D}', i} = \frac{1 - 2\varepsilon}{\gamma} (k_{i,j} - \mathbf{k}_i^T \bar{\mathbf{S}} (\bar{\mathbf{S}}^T \mathbf{K} \bar{\mathbf{S}} + \gamma \mathbf{I})^{-1} \bar{\mathbf{S}}^T \mathbf{k}_i), \quad (4)$$

DISQUEAK

Theorem

Let $\alpha = \left(\frac{1+2\varepsilon}{1-2\varepsilon}\right)$ and $\gamma > 1$. For any $0 \leq \varepsilon \leq 1$, and $0 \leq \delta \leq 1$, if we run DISQUEAK with $\bar{\mathbf{q}} = \mathcal{O}\left(\frac{\alpha}{\varepsilon^2} \log\left(\frac{n}{\delta}\right)\right)$, then w.p. $1 - \delta$, for all nodes $\{h, l\}$ in the merge tree

$$(1) \quad \|\mathbf{P}_{\{h,l\}} - \tilde{\mathbf{P}}_{\{h,l\}}\|_2 \leq \varepsilon.$$

$$(2) \quad |\mathcal{I}_{\{h,l\}}| \leq \mathcal{O}(\bar{\mathbf{q}} d_{\text{eff}}(\gamma)_{\{h,l\}}) \leq \mathcal{O}\left(\frac{\alpha}{\varepsilon^2} \mathbf{d}_{\text{eff}}(\gamma)_{\mathbf{n}} \log\left(\frac{n}{\delta}\right)\right).$$

- ▶ Same accuracy as SQUEAK but much faster
- ▶ Space/accuracy guarantees for all nodes
- ▶ Much more space used, but spread across many machines
- ▶ Runtime depends on exact merge tree
 - ↳ Fully unbalanced tree: $\mathcal{O}(n|\mathcal{I}_n|^3)$, same as SQUEAK
 - ↳ Fully balanced tree: $\mathcal{O}(\log(n)|\mathcal{I}_n|^3)$ time, $\mathcal{O}(n|\mathcal{I}_n|^3)$ work!

Comparison

	Time	$ \mathcal{I}_n $	Increm.
EXACT	n^3	n	-
Bach'13	$\frac{nd_{\max,n}^2}{\epsilon}$	$\frac{d_{\max,n}}{\epsilon}$	No
A&M'15	$n(\mathcal{I}_n)^2$	$\left(\frac{\lambda_{\min} + n\gamma\epsilon}{\lambda_{\min} - n\gamma\epsilon}\right) d_{\text{eff},n} + \frac{\text{Tr}(\mathbf{K}_n)}{\gamma\epsilon}$	No
Cal&al'16	$\frac{\lambda_{\max}^2 n^2 d_{\text{eff},n}^3}{\gamma^2 \epsilon^2}$	$\frac{\lambda_{\max} d_{\text{eff},n}}{\gamma \epsilon^2}$	Yes
SQUEAK	$\frac{nd_{\text{eff},n}^3}{\epsilon^2}$	$\frac{d_{\text{eff},n}}{\epsilon^2}$	Yes
RLS-SAMPLING	$\frac{nd_{\text{eff},n}^2}{\epsilon^2}$	$\frac{d_{\text{eff},n}}{\epsilon^2}$	-
M&M'16	$\frac{nd_{\text{eff},n}^3}{\epsilon^2}$	$\frac{d_{\text{eff},n}}{\epsilon^2}$	No

Conclusions

SQUEAK and DISQUEAK

First method (with guarantees) to break $\mathcal{O}(n)$ time barrier using DISQUEAK, with M&M'16 first to break $\mathcal{O}(n^2)$ barrier

Strong reconstruction guarantees, suitable for many downstream kernel (and not) tasks

Conclusions

SQUEAK and DISQUEAK

First method (with guarantees) to break $\mathcal{O}(n)$ time barrier using DISQUEAK, with M&M'16 first to break $\mathcal{O}(n^2)$ barrier

Strong reconstruction guarantees, suitable for many downstream kernel (and not) tasks

Final dictionary can be updated if new samples arrive

Conclusions

SQUEAK and DISQUEAK

First method (with guarantees) to break $\mathcal{O}(n)$ time barrier using DISQUEAK, with M&M'16 first to break $\mathcal{O}(n^2)$ barrier

Strong reconstruction guarantees, suitable for many downstream kernel (and not) tasks

Final dictionary can be updated if new samples arrive

Novel analysis, potentially useful for general importance sampling

Conclusions

SQUEAK and DISQUEAK

First method (with guarantees) to break $\mathcal{O}(n)$ time barrier using DISQUEAK, with M&M'16 first to break $\mathcal{O}(n^2)$ barrier

Strong reconstruction guarantees, suitable for many downstream kernel (and not) tasks

Final dictionary can be updated if new samples arrive

Novel analysis, potentially useful for general importance sampling

Conclusions

SQUEAK and DISQUEAK

First method (with guarantees) to break $\mathcal{O}(n)$ time barrier using DISQUEAK, with M&M'16 first to break $\mathcal{O}(n^2)$ barrier

Strong reconstruction guarantees, suitable for many downstream kernel (and not) tasks

Final dictionary can be updated if new samples arrive

Novel analysis, potentially useful for general importance sampling

Future work

Experiments

- ↳ Trivial to implement: 328 lines of python, single file, including distributed task queue

Conclusions

SQUEAK and DISQUEAK

First method (with guarantees) to break $\mathcal{O}(n)$ time barrier using DISQUEAK, with M&M'16 first to break $\mathcal{O}(n^2)$ barrier

Strong reconstruction guarantees, suitable for many downstream kernel (and not) tasks

Final dictionary can be updated if new samples arrive

Novel analysis, potentially useful for general importance sampling

Future work

Experiments

↳ Trivial to implement: 328 lines of python, single file, including distributed task queue

Preliminary results promising, easily scales to 100k of samples

Conclusions

SQUEAK and DISQUEAK

First method (with guarantees) to break $\mathcal{O}(n)$ time barrier using DISQUEAK, with M&M'16 first to break $\mathcal{O}(n^2)$ barrier

Strong reconstruction guarantees, suitable for many downstream kernel (and not) tasks

Final dictionary can be updated if new samples arrive

Novel analysis, potentially useful for general importance sampling

Future work

Experiments

↳ Trivial to implement: 328 lines of python, single file, including distributed task queue

Preliminary results promising, easily scales to 100k of samples

Beyond closed formulas: SQUEAK for gradient based methods