

Minimizing Quadratic Functions in Constant Time

Kohei Hayashi [National Institute of Advanced Industrial Science and Technology] and Yuichi Yoshida [National Institute of Informatics, Preferred Infrastructure, Inc.]

What We Solve

Let $A \in \mathbb{R}^{n \times n}$ be a matrix and $\mathbf{d}, \mathbf{b} \in \mathbb{R}^n$ be vectors. Then, we consider the following n -dimensional quadratic problem:

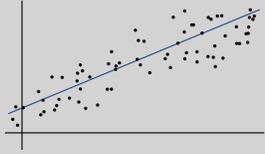
$$z^* = \min_{\mathbf{v} \in \mathbb{R}^n} p_{n,A,\mathbf{d},\mathbf{b}}(\mathbf{v}), \quad (1)$$

$$\text{where } p_{n,A,\mathbf{d},\mathbf{b}}(\mathbf{v}) = \langle \mathbf{v}, A\mathbf{v} \rangle + n\langle \mathbf{v}, \text{diag}(\mathbf{d})\mathbf{v} \rangle + n\langle \mathbf{b}, \mathbf{v} \rangle. \quad (2)$$

Caveat: not **argmin** but **min**!

Applications

- Least square distance (a.k.a. linearity check): $\min_{\mathbf{w}} \|\mathbf{y} - X\mathbf{w}\|^2$



- Kernel approximation of the Pearson divergence [Yamada+ NIPS'11]:

$$-\frac{1}{2} - \min_{\mathbf{v} \in \mathbb{R}^n} \frac{1}{2} \langle \mathbf{v}, H\mathbf{v} \rangle - \langle \mathbf{h}, \mathbf{v} \rangle + \frac{\lambda}{2} \langle \mathbf{v}, \mathbf{v} \rangle \quad (3)$$

- $H \in \mathbb{R}^{n \times n}, \mathbf{h} \in \mathbb{R}^n$: defined by a kernel function
- $\lambda \in \mathbb{R}$: regularization coefficient

How Problem (1) Has Been Solved

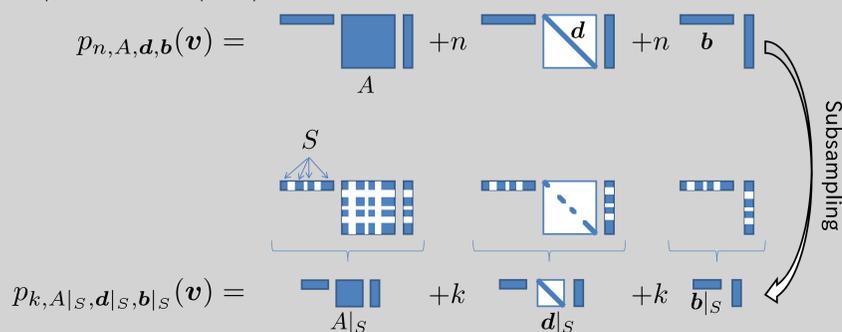
- Quadratic programming
 - Stochastic gradient
 - Nystrom's method
- Problem:** All of them need $\Omega(n)$ time
- How can we solve ultra-dimensional problem, e.g. $n \sim 10^{15}$?

Contributions

Goal: Approximately solve (1) in $O(1)$ time

Method: Solve subsampled problem $p_{k,A|_S,\mathbf{d}|_S,\mathbf{b}|_S}(\mathbf{v})$ instead of (1), where

- $k = O(1)$: sampling size
- $S \subset \{1, \dots, n\}$: sampled indices ($|S| = k$)
- $A|_S \in \mathbb{R}^{k \times k}, \mathbf{d}|_S, \mathbf{b}|_S \in \mathbb{R}^k$: subsamples of $A, \mathbf{d}, \mathbf{b}$, resp.



Main Theorem

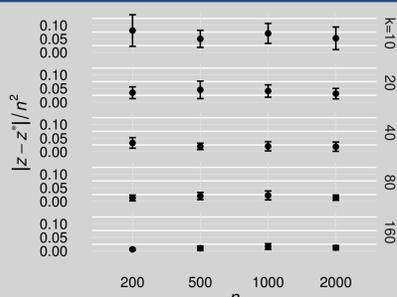
Assume $\forall_{i,j} |A_{ij}|, |b_i|, |d_i| = O(1)$. With parameters $\epsilon, \delta \in (0, 1)$, an approximate minimum $z = \frac{n^2}{k^2} \min_{\mathbf{v} \in \mathbb{R}^k} p_{k,A|_S,\mathbf{d}|_S,\mathbf{b}|_S}(\mathbf{v})$ in which $k = k(\delta, \epsilon)$ satisfies, with probability at least $1 - \delta$,

$$|z - z^*| = O(\epsilon n^2) \quad (4)$$

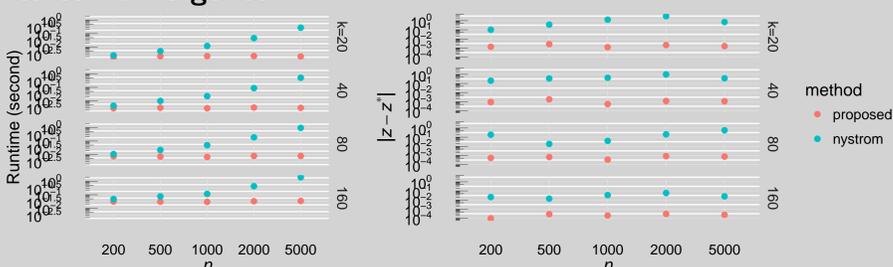
Experiments

Synthetic Data

- $A_{ij}, b_i \sim \text{unif}(-1, 1)$
- $d_i \sim \text{unif}(0, 1)$
- Minimize $p_{k,A|_S,\mathbf{d}|_S,\mathbf{b}|_S}(\mathbf{v})$ by QP



Pearson Divergence



Proof Idea

Rewrite $p_{n,A,\mathbf{d},\mathbf{b}}(\mathbf{v})$ as

$$p_{A,D,B}(\mathbf{v}) = \langle \mathbf{v}, A\mathbf{v} \rangle + \langle \mathbf{v}^2, D\mathbf{1} \rangle + \langle \mathbf{v}, B\mathbf{1} \rangle,$$

where $(v^2)_i = v_i^2$, $D = \mathbf{d}\mathbf{1}^\top$, and $B = \mathbf{b}\mathbf{1}^\top$.

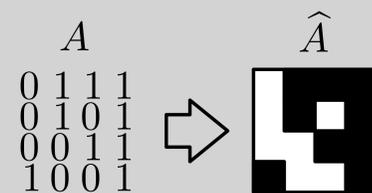
Goal: Show $\min_{\mathbf{v}} p_{A,D,B}(\mathbf{v}) \approx \frac{n^2}{k^2} \min_{\mathbf{v}} p_{A|_S,D|_S,B|_S}(\mathbf{v})$.

- Want to say $A \approx A|_S$, $D \approx D|_S$, and $B \approx B|_S$.
- How can we measure the distance between matrices of different sizes?
- Embed matrices to the same space: exploit the **graph limit theory**.

Dikernel

Dikernel: a measurable function $W : [0, 1]^2 \rightarrow \mathbb{R}$.

Any matrix $A \in \mathbb{R}^{n \times n}$ has a corresponding dikernel $\hat{A} : [0, 1]^2 \rightarrow \mathbb{R}$.



Reduction to Dikernels

For a function $f : [0, 1] \rightarrow \mathbb{R}$, consider a dikernel analogue of $p_{A,D,B}(f)$:

$$\hat{p}_{A,D,B}(f) = \langle f, \hat{A}f \rangle + \langle f^2, \hat{D}\mathbf{1} \rangle + \langle f, \hat{B}\mathbf{1} \rangle,$$

where $\langle f, Wg \rangle = \int_{[0,1]} \int_{[0,1]} W(x,y) f(x)g(y) dx dy$ and $f^2(x) = f(x)^2$.

New goal: Show $\min_f \hat{p}_{A,D,B}(f) \approx \min_f \hat{p}_{A|_S,D|_S,B|_S}(f)$.

Key Lemma

Lemma: $|\hat{p}_{A,D,B}(f) - \hat{p}_{A|_S,D|_S,B|_S}(f)|$ is small for any bounded $f : [0, 1] \rightarrow \mathbb{R}$.

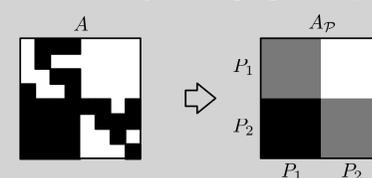
Proof of the new goal: Let $f^* = \text{argmin}_f \hat{p}_{A,D,B}(f)$ and $f' = \text{argmin}_f \hat{p}_{A|_S,D|_S,B|_S}(f)$. Then,

$$\begin{aligned} \hat{p}_{A,D,B}(f^*) &\leq \hat{p}_{A,D,B}(f') \approx \hat{p}_{A|_S,D|_S,B|_S}(f'), \\ \hat{p}_{A|_S,D|_S,B|_S}(f') &\leq \hat{p}_{A|_S,D|_S,B|_S}(f^*) \approx \hat{p}_{A,D,B}(f^*). \end{aligned}$$

Szemerédi's (Weak) Regularity Lemma

Any matrix $A \in \mathbb{R}^{n \times n}$ with $|A_{ij}| = O(1)$ has a partition $\mathcal{P} = (P_1, \dots, P_k)$ of $\{1, 2, \dots, n\}$ for constant k with the following property:

Let $A_{\mathcal{P}}$ be the matrix obtained by averaging each part $P_i \times P_j$ of A :



Then, $\|\hat{A} - \hat{A}_{\mathcal{P}}\|_{\square}$ is small.

Cut norm: $\|W\|_{\square} = \sup_{S,T \subseteq [0,1]} |\int_S \int_T W(x,y) dx dy|$.

Proof of the Key Lemma

Claim: If the cut norm $\|W\|_{\square}$ is small, then $|\langle f, Wg \rangle|$ is also small.

By the claim,

$$\begin{aligned} &|\hat{p}_{A,D,B}(f) - \hat{p}_{A|_S,D|_S,B|_S}(f)| \\ &\leq |\langle f, (\hat{A} - \hat{A}|_S)f \rangle| + |\langle f^2, (\hat{D} - \hat{D}|_S)\mathbf{1} \rangle| + |\langle f, (\hat{B} - \hat{B}|_S)\mathbf{1} \rangle| \end{aligned}$$

is small when $\|\hat{A} - \hat{A}|_S\|_{\square}$, $\|\hat{D} - \hat{D}|_S\|_{\square}$, and $\|\hat{B} - \hat{B}|_S\|_{\square}$ are small.

Let \mathcal{P} be the partition given by Szemerédi's regularity lemma. Then,

$$\|\hat{A} - \hat{A}|_S\|_{\square} \leq \|\hat{A} - \hat{A}_{\mathcal{P}}\|_{\square} + \|\hat{A}|_S - \hat{A}_{\mathcal{P}}\|_{\square},$$

which is small because $A|_S$ has enough information to approximate $A_{\mathcal{P}}$. (The arguments for D and B are almost identical.)