

A few projects I'm proud of...

Oussama Zekri¹

A selection of projects I am proud of, and a few points of interest

¹ Research Student, oussama.zekri@ens-paris-saclay.fr

June 30, 2024



A few projects I'm proud of...

Summary

① Maths project : Consistent error bounds

Projection algorithm

Douglas-Rachford splitting algorithm

Main results

② Code project : Convolutional Kernel Networks

CKN but... from scratch !

Blogpost incoming

③ Both : SGD through LLMs ICL

LLMs understand the convergence of SGD

Estimating the trans. kernel of SGD

Experiments

④ My affinities, interests and usefulness

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Maths project : Consistent error bounds

Math-oriented project : Consistent error bounds

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm
Douglas-Rachford
splitting algorithm
Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

The Convex Feasibility Problem

The (pretty simple) problem of interest

Let C_1, \dots, C_m be closed convex sets included in a finite-dimensional v.s. \mathcal{E} and $C = \cap_{i=1}^m C_i \neq \emptyset$.

Find $x \in C$ (CFP)

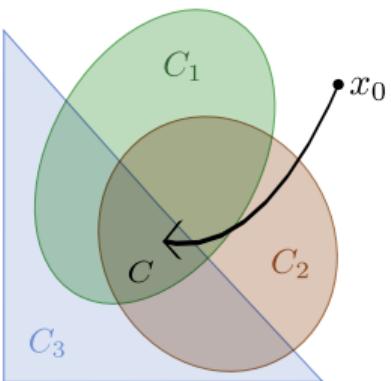


Figure: Illustration of CFP in the case of $m = 3$ sets

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Projection algorithm

The first idea that comes to mind : Cyclic projections.

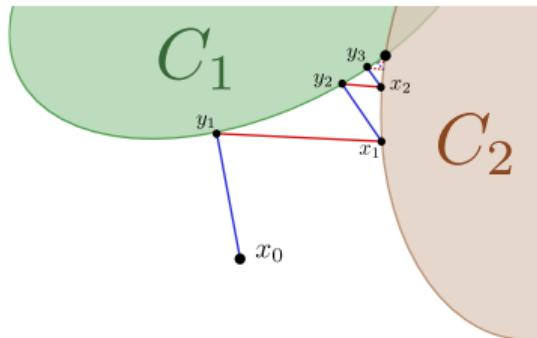
A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm
Douglas-Rachford
splitting algorithm

Classic idea for solving the problem : We project alternatively between the C_i sets.



Cyclic projection algorithm (1933)

Require: x_0 and N

- 1: **for** $k = 0$ to N **do**
- 2: $y_{k+1} = p_{C_1}(x_k)$
- 3: $x_{k+1} = p_{C_2}(y_{k+1})$
- 4: **end for**

Douglas-Rachford splitting algorithm

A slightly more sophisticated algorithm...

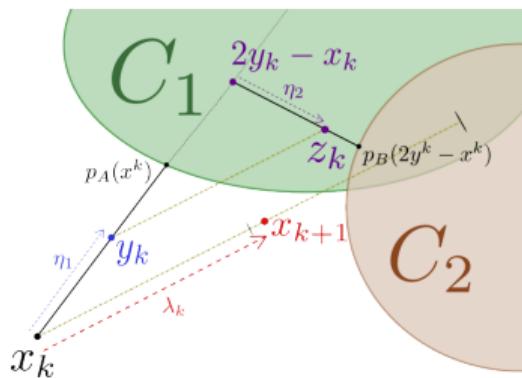
A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds
Projection algorithm
Douglas-Rachford
splitting algorithm

We introduce $f(x) = \text{dist}^2(x, C_1)$ and $g(x) = \text{dist}^2(x, C_2)$. We replace projections by proximal operators :

$$\text{prox}_{\eta f}(x) = \frac{1}{2\eta + 1}x + \frac{2\eta}{2\eta + 1}p_{C_1}(x)$$



(Damped)Douglas-Rachford algorithm (1956)

Require: x_0, N, η_1, η_2 , and $(\lambda_k)_k \in (0, 1]$

- 1: **for** $k = 0$ to N **do**
- 2: $y_k = \text{prox}_{\eta_1 f}(x_k)$
- 3: $z_k = \text{prox}_{\eta_2 g}(2y_k - x_k)$
- 4: $x_{k+1} = x_k + 2\lambda_k(z_k - y_k)$
- 5: **end for**

Damped Douglas-Rachford

One iteration of the Damped Douglas-Rachford algorithm.

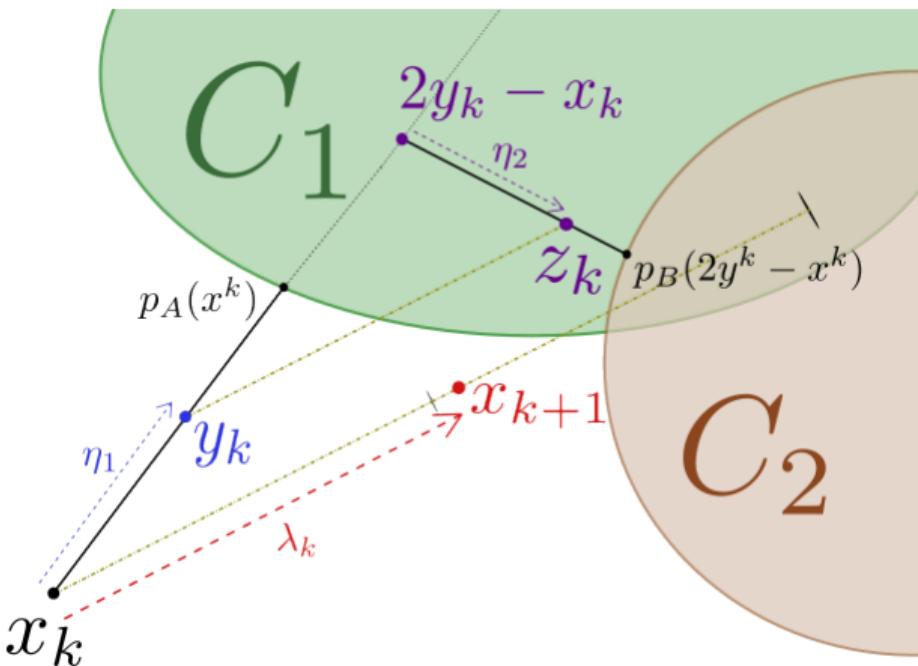


Figure: One iteration of the Damped Douglas-Rachford algorithm

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Main results

Consistent error bounds

A few projects
I'm proud of...

ZEKRI

- Let C_1, \dots, C_m be closed convex sets included in a finite-dimensional e.v. \mathcal{E} and $C = \cap_{i=1}^m C_i \neq \emptyset$. Find $x \in C$ (CFP)
- (strict) Consistent error bound function for C_1, \dots, C_m ¹ : $\Phi : \mathbb{R}^+ \times \mathbb{R}^+ \mapsto \mathbb{R}^+$ such that
 - ① $\forall x \in \mathcal{E}, \text{dist}(x, C) \leq \Phi\left(\max_{1 \leq i \leq m} \text{dist}(x, C_i), \|x\|\right)$
 - ② $\forall b \in \mathbb{R}^+, \Phi(., b)$ is monotone (**increasing**) nondecreasing on \mathbb{R}^+ , right-continuous at 0 and satisfies $\Phi(0, b) = 0$
 - ③ $\forall a \in \mathbb{R}^+, \Phi(a, .)$ is monotone nondecreasing on \mathbb{R}^+

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

¹T. Liu & B.F. Lourenço, Convergence analysis under consistent error bounds

Main results

Validity assumption

Assumption 1

Let $\{x^k\} \subseteq \mathcal{E}$ be a sequence such that the following conditions hold.

i) *Fejér monotonicity condition.* For any fixed $c \in C$, it holds that

$$\|x^{k+1} - c\| \leq \|x^k - c\| \quad \forall k. \quad (1)$$

ii) *Sufficient decrease condition.* There exist some positive integer ℓ and nonnegative sequence $\{a_k\}$ with $\sum_{k=0}^{\infty} a_k = \infty$ such that

$$\text{dist}^2(x^k, C) \geq \text{dist}^2(x^{k+\ell}, C) + a_k \max_{1 \leq i \leq m} \text{dist}^2(x^k, C_i) \quad \forall k. \quad (2)$$

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Main results

The big (and ugly ?) Theorem !

Proposition

Let Assumption 1 holds. Then $\{x^k\}$ converges to some point in C .

Theorem

Suppose that Assumption 1 holds. Let Φ be a strict consistent error bound function for C_1, \dots, C_m . Let $\Phi_{\widehat{\kappa}}^{\spadesuit}$ be defined as in Definition 2 b with $\widehat{\kappa}$ such that $\widehat{\kappa} \geq \|x^0\| + 2 \operatorname{dist}(0, C)$. Then, the convergence of $\{x^k\}$ is either finite or

$$\operatorname{dist}(x^k, C) \leq \sqrt{(\Phi_{\widehat{\kappa}}^{\spadesuit})^{-1} \left(\Phi_{\widehat{\kappa}}^{\spadesuit}(\operatorname{dist}^2(x^0, C)) - \sum_{i=0}^{b_k-1} a_{k_0+i\ell} \right)} \quad \forall k \geq 2\ell \quad (3)$$

holds for any integer $k_0 \in [0, \ell - 1]$ and $b_k := \frac{k-\ell-(k \bmod \ell)}{\ell}$.

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm
Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

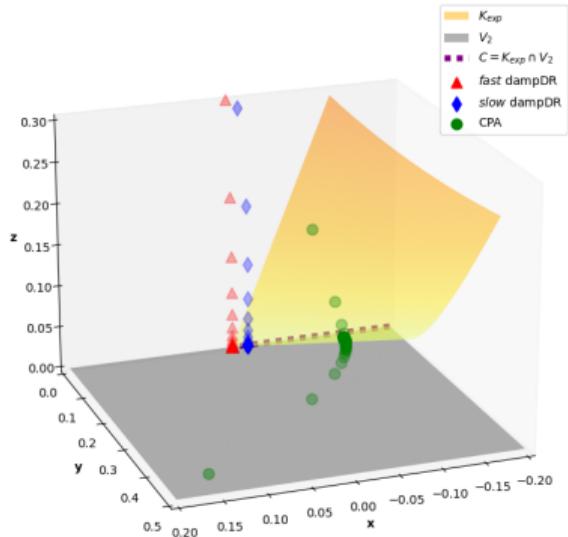
Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Paper incoming

to be submitted really soon to Optimization Letters



-Title still to be determined-

Oussema Zekri^a Ellen H. Fukuda^a Tianxiang Liu^b Bruno F. Lourenço^b

June 6, 2024

Abstract

The new notion of constraint error bounds provides a general framework for the study of error bounds for convex feasibility problems (CFP), including linear-convex cases. In particular, it provides an upper bound on the convergence rates of multiple optimization algorithms, suitable for solving the CFP. The first one, the fast dampDR, is a well-known algorithm for solving the CFP, which we are now seeking to establish this result for two more complex, significantly more efficient and numerically used algorithms. The fast one, Dykstra's algorithm, is a fairly old one, but has recently regained interest due to its efficiency in solving the CFP. The second one, CPA, is a very recent algorithm for the CFP, by giving the closest possible solution to the initial point of the algorithm. The second one is a refinement version of the proximal point splitting method. While used, this algorithm is not popular yet, probably because of the need of being initialized in some way. Once we have established the results for these two algorithms, we will turn our attention to two very specific cases of the CFP, namely the regression case. The error bound underlying these specific cases are not classical, but are given by the notion of constraint error bounds. We also obtain an upper bound on the convergence rates of algorithms for these specific cases, which we will observe through various numerical experiments.

1 Introduction

In this paper, we consider the following convex feasibility problem (CFP)

$$\text{find } x \in C := \bigcap_{i=1}^m C_i \quad (\text{CFP})$$

where C_1, \dots, C_m are closed convex sets contained in a finite dimensional real vector space E with $C \neq \emptyset$.

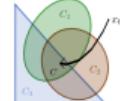


Figure 1: Illustration of CFP in the case of $m = 3$ sets.

Given some fixed algorithm for solving (CFP), the following two questions are of natural interest.

^aDepartment of Teaching and Research of Mathematics, ENS Paris-Saclay, France. enseignants.ens-paris-saclay.fr

^bGraduate School of Informatics, Kyoto University, Kyoto, Japan. informatics.kyoto-u.ac.jp

^cResearch Institute of Mathematics, Kyoto University, Kyoto, Japan. math.kyoto-u.ac.jp

1

Joint work with [Ellen H. Fukuda](#), [Bruno F. Lourenço](#) and [Tianxiang Liu](#), realized during the 2nd year internship at the Maths department of ENS Paris-Saclay.

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Bibliography



T. Liu and B. F. Lourenço.

Convergence analysis under consistent error bounds.

Found Comput Math, 2022.

arXiv:2008.12968.



J. Douglas and H. H. Rachford.

On the Numerical Solution of Heat Conduction Problems in Two and Three Space Variables.

JSTOR, 1956.



B. F. Svaiter.

On weak convergence of the Douglas-Rachford method.

SIAM Journal on Control and Optimization, 2011.

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Code project : Convolutional Kernel Networks

Code-oriented project : Convolutional Kernel Networks

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

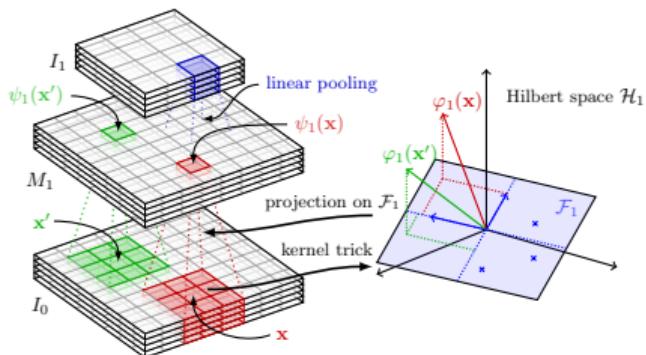
Experiments

My affinities,
interests and
usefulness

CKN but... from scratch !

Recoding the whole architecture from scratch

To win the challenge of J. Mairal's course, we decided to implement the CKN. The only requirement of the course was to use Kernel Methods and to code everything from scratch (**no ML libraries !**).



This implies many challenges, such as

- Recoding automatic differentiation (track correctly the gradients, computational graph etc...)
- In-depth understanding of how the architecture works.
- Harder : Parallelizing the code (with CUDA !)

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Blog : logB project.

Our new blog, with Ambroise Odonnat.

A blogpost about this work is being finalized.

The blogpost : <https://logb-research.github.io/blog/2024/cnk/>

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Bibliography

A few projects
I'm proud of...

ZEKRI



J. Mairal, P. Koniusz, Z. Harchaoui and C. Schmid.
Convolutional Kernel Networks.
NIPS, 2014.
arXiv:1406.3332.

Maths project :
Consistent error
bounds

Projection algorithm
Douglas-Rachford
splitting algorithm
Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness



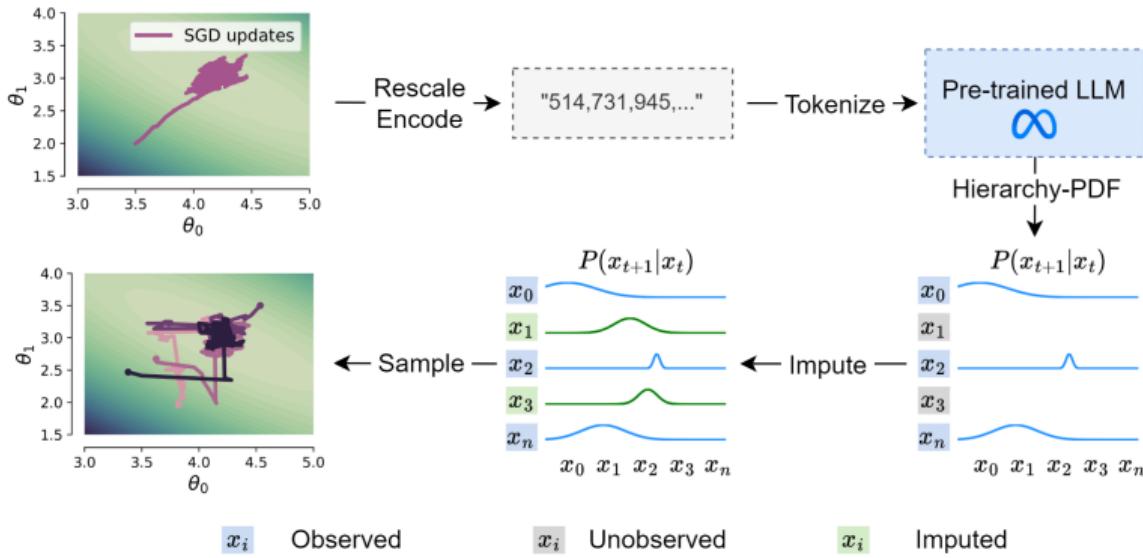
J. Mairal.
End-to-End Kernel Learning with Supervised Convolutional Kernel Networks.
NIPS, 2016.



A. Bietti.
Foundations of deep convolutional models through kernel methods.
PhD Thesis, Université Grenoble Alpes, 2019.

A clever mix : Understanding SGD through LLMs ICL abilities

A clever mix : Understanding SGD through LLMs ICL abilities



In Context Learning

Procedure: ICL for dynamics learning

Procedure: ICL for dynamics learning

Input: time serie $(x_i)_{i \leq t}$, LLM M , precision k

1. Rescale and encode the time serie with k digits

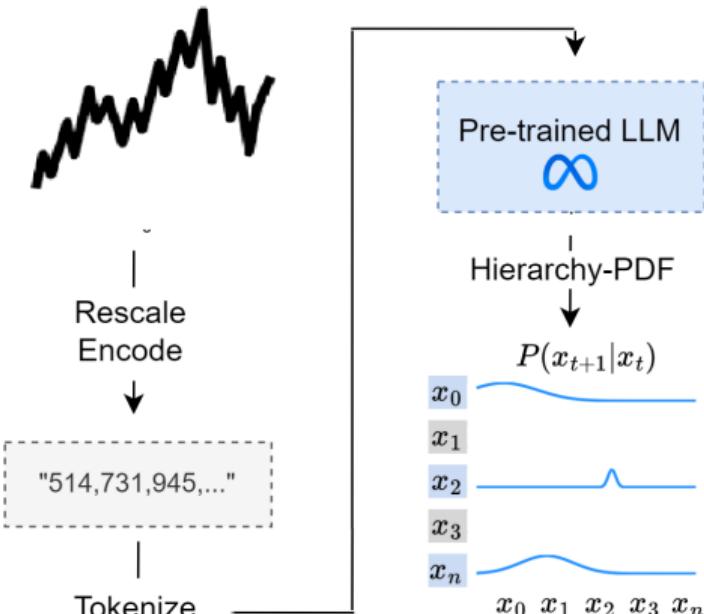
$$\hat{x}_t = "x_1^1 x_1^2 \dots x_1^k, \dots"$$

2. Call $M(\hat{x}_t)$

3. Extract the digits logits $(0, 1, 2, 3, \dots, 9)$

4. Build the next state probability distribution using the Hierarchy-PDF algorithm in [Liu et al., 2024]

Return: predicted transition rules for the observed states: $\{P(X_{i+1}|X_i = x_i)\}_{i \leq t}$



A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm
Douglas-Rachford
splitting algorithm
Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD
Estimating the trans.
kernel of SGD
Experiments

My affinities,
interests and
usefulness

LLMs understand the convergence of SGD

Problem setup

Training set $x = (x_1, \dots, x_N)$ of N i.i.d samples,

$$\min_{\theta} F(\theta), \quad F(\theta) = \frac{1}{N} \sum_{i=1}^N f(x_i, \theta), \quad (4)$$

where $\theta \in \mathbb{R}^d$. Minibatch SGD updates :

$$\theta^{t+1} = \theta^t - \gamma_t \nabla \tilde{f}_t(\theta^t) \quad (5)$$

where θ^t denotes the parameters after t iterations, and $\nabla \tilde{f}_t(\theta^t) = \frac{1}{m} \sum_{x \in B_t} \nabla_{\theta} f(x, \theta^t)$ where B_t is a minibatch of size m of training examples selected randomly.

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

LLMs understand the convergence of SGD

Overparametrized vs. underparametrized regime

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

$$\theta^{t+1} = \theta^t - \gamma \nabla \tilde{f}_t(\theta^t) \quad (6)$$

If $\gamma_t = \gamma$, θ^t from (6) form a homogeneous Markov chain that converge to a unique stationary distribution π_γ .

- In the *overparametrized* regime(i.e. when $d >> N$), $\pi_\gamma = \delta_{\tilde{\theta}^*}$ where $\tilde{\theta}^*$ is a specific optimum.
- In the *underparametrized* regime (i.e. when $d << N$), π_γ is a distribution with a strictly positive variance, e.g. $\mathcal{N}(\theta^*, \gamma^{1/2})$ where θ^* is an optimum.

LLMs understand the convergence of SGD

Overparametrized vs. underparametrized regime

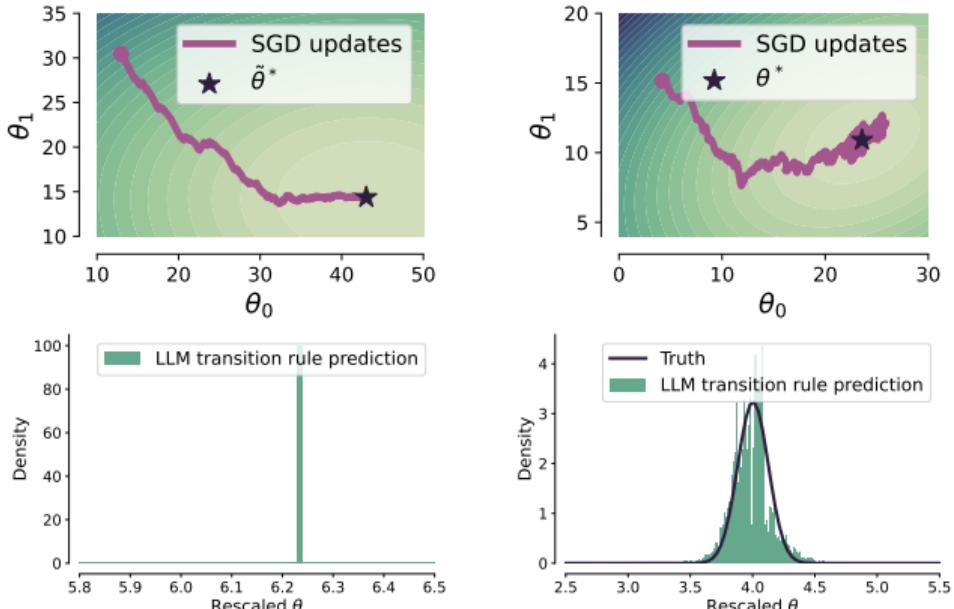


Figure: Top Left and Top Right, a run of SGD in the overparameterized and underparameterized regimes. Bottom Left and Bottom Right, transition probabilities predicted by LLM in overparameterized and underparameterized regimes.

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm
Douglas-Rachford
splitting algorithm
Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD
Experiments

My affinities,
interests and
usefulness

Estimating the transition kernel of SGD

For each parameter $\theta_i, i \in \{1, \dots, d\}$, we consider the one-hot discretized state vector Θ_i^t at time t . Then, we can write

$$\Theta_i^{t+1} = \sum_{j=1}^d \lambda_{i,j} P^{(i,j)} \Theta_j^t$$

where $\forall i, j, \lambda_{i,j} \geq 0, \sum_{j=1}^d \lambda_{i,j} = 1$ and $P^{(i,j)} = P(\theta_j | \theta_i)$.

Then, the discretized transition kernel of SGD can be seen as a matrix

$$Q = \begin{pmatrix} \lambda_{1,1} P^{(1,1)} & \dots & \lambda_{1,d} P^{(1,d)} \\ \vdots & \ddots & \vdots \\ \lambda_{d,1} P^{(d,1)} & \dots & \lambda_{d,d} P^{(d,d)} \end{pmatrix}$$

which satisfies $\Theta^{t+1} = Q\Theta^t$.

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Algorithm

Estimating $P^{(i,i)}$

Input: time serie $(\theta_i^{t+1})_{t \geq 0}$, LLM M , precision k , regularization ε

1. Fill $s < 10^k$ rows of the 10^k rows of $P^{(i,i)}$ with Procedure(θ_i^{t+1}, M, k), denoted as $(P_1^{(i,i)}, \dots, P_s^{(i,i)})$

2. Fill the remaining $10^k - s$ rows of $P^{(i,i)}$ with debiased Sinkhorn barycenter of regularization parameter ε :

for $j = 1$ **to** $s - 1$ **do**

if empty rows between $P_j^{(i,i)}$ and $P_{j+1}^{(i,i)}$ **then**

 Compute debiased Sinkhorn barycenter between

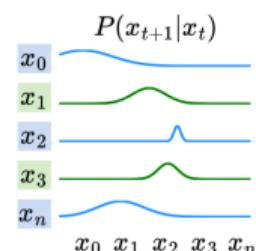
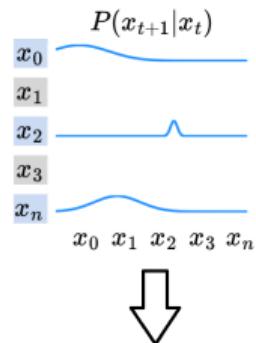
$P_j^{(i,i)}$ and $P_{j+1}^{(i,i)}$, with regularization parameter ε

 Fill the empty rows

end if

end for

Return: Estimated matrix $P^{(i,i)}$



A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds
Projection algorithm
Douglas-Rachford
splitting algorithm
Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !
Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Experiments

Convex case

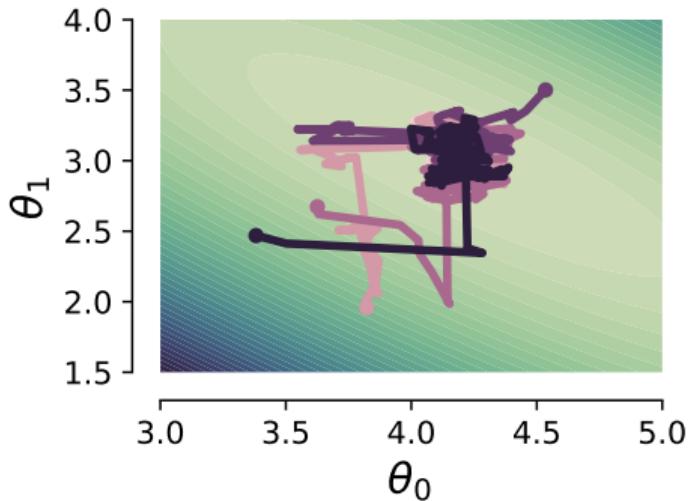
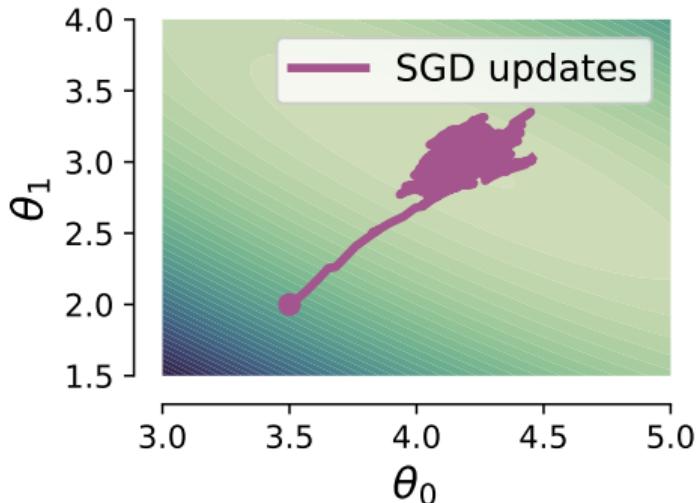


Figure: We optimize F defined in (4) with $f(x_i, \theta) = \frac{1}{2}(\langle x_i, \theta \rangle_{\mathbb{R}^2} - y)^2$ for $d = 2$ and $N = 100$. **Left.** A full SGD run used to learn Q . **Right.** Starting from different initial points, simulating the SGD thanks to Q .

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CNN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Experiments

Non-convex case

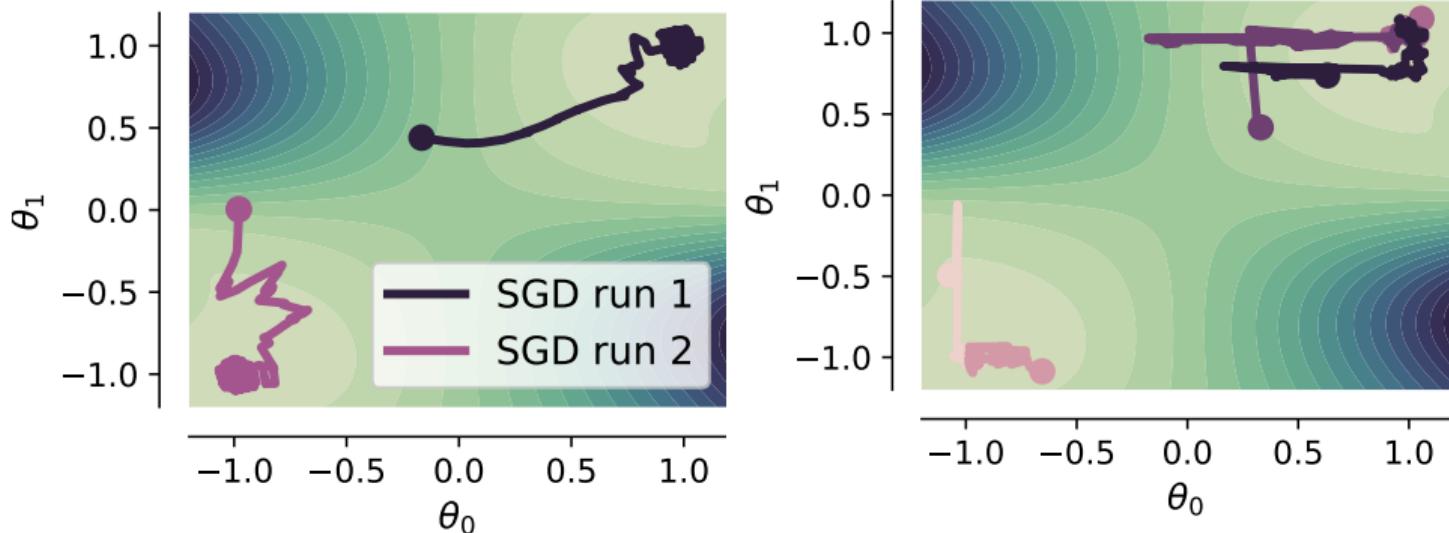


Figure: We optimize F defined in (4) with $f(x_i, \theta) = \frac{1}{2}(\theta_0 \sin(\theta_1 x_i) - y)^2$ for $d = 2$ and $N = 100$.
Left. A full SGD run used to learn Q . **Right.** Starting from different initial points, simulating the SGD thanks to Q .

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CNN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Paper submitted

submitted really soon to ICL workshop at ICML 2024

Can LLMs predict the convergence of Stochastic Gradient Descent?

Anonymous Authors[†]

Abstract

Large-language models are notoriously famous for their impressive performance across a wide range of tasks. One surprising example of such impressive performance is a recently identified capacity of LLMs to understand the governing principles of dynamical systems satisfying the Markovian property. In this paper, we seek to explore this direction further by studying the dynamics of stochastic gradient descent in convex and non-convex optimization. By leveraging the theoretical link between the SGD and Markov chains, we show a remarkable zero-shot performance of LLMs in predicting the local minima to which SGD converges for previously unseen starting points. On a more general level, we inquire about the possibility of using LLMs to perform zero-shot randomized trials for larger deep learning models used in practice.

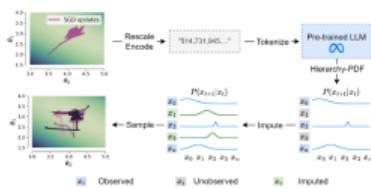


Figure 1. Overview of the proposed approach. After having run SGD on a given optimization problem, we tokenize the obtained iterates and feed them to an LLM of choice. We further use the logits to fill the transition kernel of the Markov chain underlying the SGD with probabilities $P(x_{t+1}|x_t)$, while imputing those of its elements that were not observed. Finally, we use the estimate transition kernel to do forecasting for previously unseen inputs.

Joint work with [Abdelhakim Benechehab](#) and [Ievgen Redko](#), started 2 months ago during the 3rd year internship at the Maths department of ENS Paris-Saclay.

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

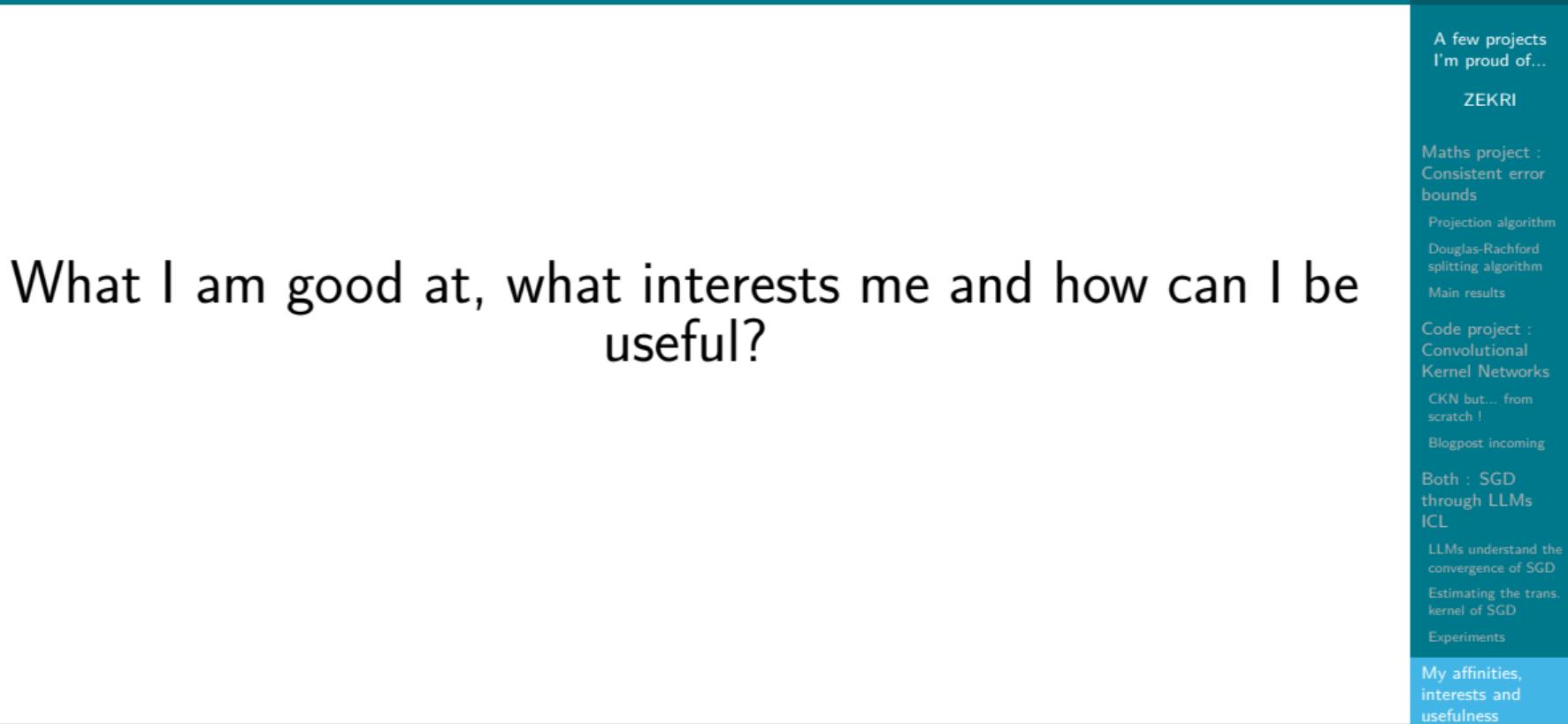
LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

What I am good at, what interests me and how can I be useful?



Affinity and interests

Because of my background at ENS Saclay, and my personal preferences, **I'm more competent and confident when I'm doing maths**, with some theory. But I also really enjoy coding for my research.

Topics of interest include :

- **Theoretical fields with application to ML** especially Optimization, Kernel Methods and Optimal Transport.
- Differential programming, theory of deep learning and neural nets in general...
- **"Applied" Large-scale ML** : LLMs, ICL...

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm
Douglas-Rachford
splitting algorithm
Main results

Code project :
Convolutional
Kernel Networks
CKN but... from
scratch !

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD
Estimating the trans.
kernel of SGD
Experiments

My affinities,
interests and
usefulness

How can I be useful

I did 2 years of oral exams in preparatory class at LLG, in maths. I am also, for this year 2024, a **member of the HEC entrance exam jury** : I have corrected 300 papers in the written exams, and I will provide 2 weeks of oral exams, on my own exercises.

This could be useful, as I could be a **teaching assistant**, which could help finance part of the potential year.

A few projects
I'm proud of...

ZEKRI

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Thank you for your attention !

A few projects
I'm proud of...

ZEKRI

Thank you for your attention !



@oussamazekri_



My website : www.oussamazekri.fr

Maths project :
Consistent error
bounds

Projection algorithm

Douglas-Rachford
splitting algorithm

Main results

Code project :
Convolutional
Kernel Networks

CKN but... from
scratch !

Blogpost incoming

Both : SGD
through LLMs
ICL

LLMs understand the
convergence of SGD

Estimating the trans.
kernel of SGD

Experiments

My affinities,
interests and
usefulness

Appendix

A few projects
I'm proud of...

ZEKRI

Appendix - CEB

Appendix - LLM
ICL

Definition (Inverse smoothing function)

Let Φ be a *strict* consistent error bound function.

- a for $\kappa > 0$, we define $\phi_{\kappa, \Phi}$ as $\phi_{\kappa, \Phi}(t) := (\Phi(\sqrt{t}, \kappa))^2$, $t \geq 0$.
- b For $\kappa > 0$ and for $\phi_{\kappa, \Phi}$ we define Φ_κ^\spadesuit as $\Phi_\kappa^\spadesuit(t) := \int_\delta^t \frac{1}{\phi_{\kappa, \Phi}^-(s)} ds$, $t \in (0, \sup \phi_{\kappa, \Phi})$, where $\delta \in (0, \sup \phi_{\kappa, \Phi})$ is some fixed number.

Dykstra's projection algorithm

Case of two sets

Dykstra's algorithm solves not only the CFP but also the BAP, i.e. it finds,

$$x \in C := \bigcap_{i=1}^m C_i, \text{ s.t. } x = p_C(x_0)$$

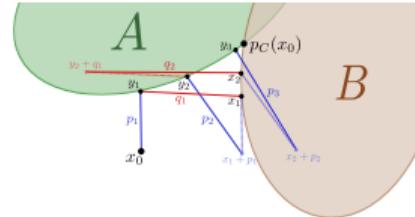


Figure: First iterations of DPA

Dykstra's projection algorithm for two sets

Require: x_0, N

- 1: $p_0 = q_0 = 0$
- 2: **for** $k = 0$ to N **do**
- 3: $y_k = p_A(x_{k-1} + p_{k-1})$
- 4: $p_k = x_{k-1} + p_{k-1} - y_k$
- 5: $x_k = p_B(y_k + q_{k-1})$
- 6: $q_k = y_k + q_{k-1} - x_k$
- 7: **end for**

A few projects
I'm proud of...

ZEKRI

Appendix - CEB

Appendix - LLM
ICL

Dykstra's projection algorithm

Case of two sets

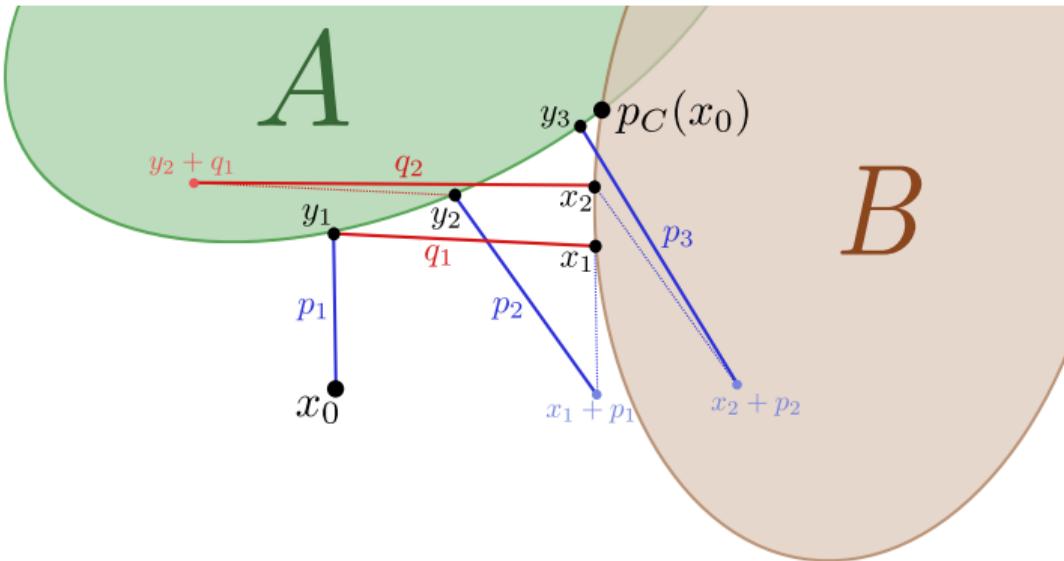


Figure: First iterations of DPA

A few projects
I'm proud of...

ZEKRI

Appendix - CEB

Appendix - LLM
ICL

Dykstra's projection algorithm

Properties

Lemma (Iterates inequality)

For each k , $\langle x_k - b, q_k \rangle \geq 0, \forall b \in B$ and $\langle y_k - a, p_k \rangle \geq 0, \forall a \in A$.

A few projects
I'm proud of...

ZEKRI

Appendix - CEB

Appendix - LLM
ICL

Dykstra's projection algorithm

Properties

Lemma (Iterates inequality)

For each k , $\langle x_k - b, q_k \rangle \geq 0, \forall b \in B$ and $\langle y_k - a, p_k \rangle \geq 0, \forall a \in A$.

Lemma (Norm equality)

For each k , for all $c \in C$, define

$$\begin{aligned} c_k &= \|x_k - y_{k+1}\|^2 + \|x_{k+1} - y_{k+1}\|^2 + 2\langle y_{k+1} - c, p_{k+1} \rangle \\ &\quad + 2\langle x_{k+1} - c, q_{k+1} \rangle + 2\langle y_k - y_{k+1}, p_k \rangle + 2\langle x_k - x_{k+1}, q_k \rangle \\ &\quad - 2\langle y_k - c, p_k \rangle - 2\langle x_k - c, q_k \rangle. \end{aligned}$$

Then, the following holds:

$$\|x_k - c\|^2 = \|x_{k+1} - c\|^2 + c_k. \tag{7}$$

Dykstra's projection algorithm

Main result

A few projects
I'm proud of...

ZEKRI

Lemma (Assumption ?? holds for Dykstra's algorithm)

Let the sequences $\{x_k\}$ and $\{y_k\}$ be generated by the algorithm scheme 2, $\{c_k\}$ defined as in (7) and define

$$d_k := 1 + \frac{c_k}{\|x_k - y_{k+1}\|^2}. \quad (8)$$

If c_k is nonnegative sequence, then the following statements hold,

- (i) The sequences $\{x_k\}$ is Fejér monotone with respect to $A \cap B$.
- (ii) For any k , it holds that

$$\text{dist}^2(x_k, A \cap B) \geq \text{dist}^2(x_{k+1}, A \cap B) + d_k \max \{\text{dist}^2(x_k, A), \text{dist}^2(x_k, B)\}. \quad (9)$$

Appendix - CEB

Appendix - LLM
ICL

Dykstra's projection algorithm

Main result

A few projects
I'm proud of...

ZEKRI

Appendix - CEB

Appendix - LLM
ICL

Theorem (Convergence of Dykstra with $m = 2$)

If Φ is a strict consistent e.b. function for C_1, C_2 , with $\hat{\kappa} = \|x_0\| + 2 \operatorname{dist}(0, C)$. Then either the sequence $\{x_k\}$ has finite convergence or converges with rate

$$\operatorname{dist}(x_k, C) \leq \sqrt{(\Phi_{\kappa}^{\blacklozenge})^{-1} \left(\Phi_{\kappa}^{\blacklozenge}(\operatorname{dist}^2(x_0, C)) - \sum_{i=0}^{k-2} d_i \right)} \quad \forall k \geq 2.$$

Applications : Image restoration

Now... We can unblur faster!

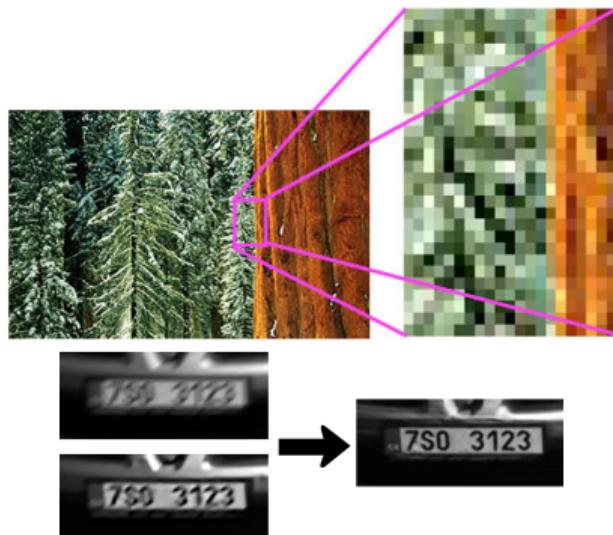


Figure: Image restoration procedure applied to the *LPD* dataset samples

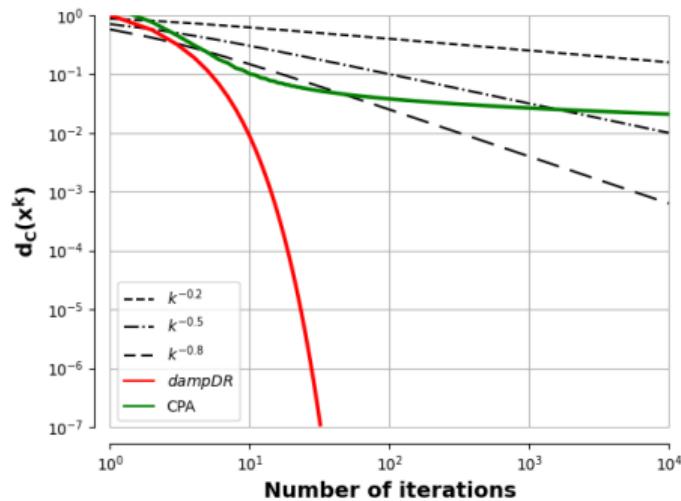


Figure: dampDR (red) vs CPA (green) convergence rates

A few projects
I'm proud of...

ZEKRI

Appendix - CEB

Appendix - LLM
ICL