# The Stochastic Bregman Primal-Dual Algorithm

Antonio Silveti-Falls, Cesare Molinari, and Jalal Fadili

# Main Idea

Solving structured convex optimization problems:

$$\min_{x \in \mathcal{C}} f(x) + g(Tx).$$

## Question

How to take advantage of properties of the individual terms?

# Main Idea

Solving structured convex optimization problems:

$$\min_{x \in \mathcal{C}} f(x) + g(Tx).$$

### Question

How to take advantage of properties of the individual terms?

- Lipschitz-smoothness - $\nabla f(x)$.

# Main Idea

Solving structured convex optimization problems:

$$\min_{x \in \mathcal{C}} f(x) + g(Tx).$$

## Question

How to take advantage of properties of the individual terms?

- Lipschitz-smoothness - $\nabla f(x)$.
- prox-friendliness - $\operatorname{prox}_g(x) \stackrel{\text{def}}{=} \operatorname*{argmin}_u \left\{ g(u) + \frac{1}{2} \|x - u\|_2^2 \right\}$.

Solving structured convex optimization problems:

$$\min_{x \in \mathcal{C}} f(x) + g(Tx).$$

### Question

How to take advantage of properties of the individual terms?

- Lipschitz-smoothness - $\nabla f(x)$.
- prox-friendliness - $\operatorname{prox}_g(x) \stackrel{\text{def}}{=} \operatorname*{argmin}_u \left\{ g(u) + \frac{1}{2} \|x - u\|_2^2 \right\}$.
- Projection onto $\mathcal{C}$ - $P_{\mathcal{C}}(x) \stackrel{\text{def}}{=} \operatorname*{argmin}_{u \in \mathcal{C}} \|x - u\|_2^2$.

Solving structured convex optimization problems:

$$\min_{x \in \mathcal{C}} f(x) + g(Tx).$$

## Question

How to take advantage of properties of the individual terms?

- Lipschitz-smoothness - $\nabla f(x)$.
- prox-friendliness - $\operatorname{prox}_g(x) \overset{\text{def}}{=} \underset{u}{\operatorname{argmin}} \left\{ g(u) + \frac{1}{2} \|x - u\|_2^2 \right\}$.
- Projection onto $\mathcal{C}$ - $P_{\mathcal{C}}(x) \overset{\text{def}}{=} \underset{u \in \mathcal{C}}{\operatorname{argmin}} \|x - u\|_2^2$.

Changing the geometry?

Solving structured convex optimization problems:

$$\min_{x \in \mathcal{C}} f(x) + g(Tx).$$

## Question

How to take advantage of properties of the individual terms?

- ~~Lipschitz~~-smoothness - $\nabla f(x)$.
- prox-friendliness - $\text{prox}_g(x) \stackrel{\text{def}}{=} \underset{u}{\text{argmin}} \left\{ g(u) + \frac{1}{2} \|x - u\|_2^2 \right\}$.
- Projection onto $\mathcal{C}$ - $P_{\mathcal{C}}(x) \stackrel{\text{def}}{=} \underset{u \in \mathcal{C}}{\text{argmin}} \|x - u\|_2^2$.

Changing the geometry?

Consider a matrix $Y \stackrel{\text{def}}{=} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}_{++}^{n \times p}$ with $y_i \in \Delta^p$ and matrices $A_1, \ldots, A_n \in \mathbb{R}_+^{p \times m}$ without any zero rows.

Consider a matrix $Y \overset{\text{def}}{=} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}_{++}^{n \times p}$ with $y_i \in \Delta^p$ and matrices

$A_1, \ldots, A_n \in \mathbb{R}_+^{p \times m}$ without any zero rows.

We examine the *trend filtering* problem of recovering a matrix

$X \overset{\text{def}}{=} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}_+^{n \times m}$ with $x_i \in \Delta^m$ under the model

$$\mathcal{A}X \approx Y \quad \text{with} \quad \mathcal{A}X \overset{\text{def}}{=} \begin{pmatrix} A_1 x_1 \\ \vdots \\ A_n x_n \end{pmatrix} \in \mathbb{R}_+^{n \times p}$$

and assuming that the columns of $X$ are piecewise constant.

## The Kullback-Leibler divergence

For $u, v \in \mathbb{R}_+$,

$$\mathrm{KL}(u, v) \stackrel{\text{def}}{=} \begin{cases} u \log \left( \frac{u}{v} \right) - u + v & \text{if } u, v > 0, \\ v & \text{if } u = 0, \\ +\infty & \text{otherwise.} \end{cases}$$

# Trend Filtering - Notation

## The Kullback-Leibler divergence

For $u, v \in \mathbb{R}_+$,

$$\mathrm{KL}\left(u, v\right) \stackrel{\text{def}}{=} \begin{cases} u \log\left(\frac{u}{v}\right) - u + v & \text{if } u, v > 0, \\ v & \text{if } u = 0, \\ +\infty & \text{otherwise.} \end{cases}$$

## The row gradient

$\nabla_{\mathrm{row}} : \mathbb{R}^{n \times m} \to \mathbb{R}^{m(n-1)}$. For a matrix $X \in \mathbb{R}^{n \times m}$,

$$\nabla_{\mathrm{row}} X \stackrel{\text{def}}{=} \begin{pmatrix} x_2 - x_1 \\ \vdots \\ x_n - x_{n-1} \end{pmatrix}.$$

## Trend filtering

$$\min_{\substack{X \in \mathbb{R}_+^{n \times m} \\ X \mathbb{1}_m = \mathbb{1}_n}} \underbrace{\sum_{i=1}^{n} \mathrm{KL}\left(A_i x_i, y_i\right)}_{f(X)} + \underbrace{\beta \left\| \nabla_{\mathrm{row}} X \right\|_1}_{g \circ \nabla_{\mathrm{row}}(X)}$$

## Trend filtering

$$\min_{\substack{X \in \mathbb{R}_+^{n \times m} \\ X \mathbb{1}_m = \mathbb{1}_n}} \underbrace{\sum_{i=1}^{n} \mathrm{KL}\left(A_i x_i, y_i\right)}_{f(X)} + \underbrace{\beta \left\| \nabla_{\mathrm{row}} X \right\|_1}_{g \circ \nabla_{\mathrm{row}}(X)}$$

Obstacles:

- $f$ is not Lipschitz-smooth at the origin.

## Trend filtering

$$\min_{\substack{X\in\mathbb{R}_+^{n\times m} \\ X\mathbb{1}_m=\mathbb{1}_n}} \underbrace{\sum_{i=1}^n \mathrm{KL}\left(A_i x_i, y_i\right)}_{f(X)} \quad + \quad \underbrace{\beta\left\|\nabla_{\mathrm{row}}X\right\|_1}_{g\circ\nabla_{\mathrm{row}}(X)}$$

Obstacles:

- $f$ is not Lipschitz-smooth at the origin.
- $\nabla_{\mathrm{row}}$ makes computing $\mathrm{prox}_{g\circ\nabla_{\mathrm{row}}}$ difficult.

## Trend filtering

$$\min_{\substack{X \in \mathbb{R}^{n \times m}_+ \\ X \mathbb{1}_m = \mathbb{1}_n}} \underbrace{\sum_{i=1}^n \mathrm{KL}\left(A_i x_i, y_i\right)}_{f(X)} \quad + \quad \underbrace{\beta \left\| \nabla_{\mathrm{row}} X \right\|_1}_{g \circ \nabla_{\mathrm{row}}(X)}$$

Obstacles:

- $f$ is not Lipschitz-smooth at the origin.
- $\nabla_{\mathrm{row}}$ makes computing $\mathrm{prox}_{g \circ \nabla_{\mathrm{row}}}$ difficult.
- $\mathrm{prox}_f$ is computable when the $A_i$ are nice but requires special functions (Lambert $W$-function).

## Trend filtering

$$\min_{\substack{X \in \mathbb{R}_+^{n \times m} \\ X \mathbb{1}_m = \mathbb{1}_n}} \underbrace{\sum_{i=1}^n \mathrm{KL}\left(A_i x_i, y_i\right)}_{f(X)} + \underbrace{\beta \left\| \nabla_{\mathrm{row}} X \right\|_1}_{g \circ \nabla_{\mathrm{row}}(X)}$$

Obstacles:

- $f$ is not Lipschitz-smooth at the origin.
- $\nabla_{\mathrm{row}}$ makes computing $\mathrm{prox}_{g \circ \nabla_{\mathrm{row}}}$ difficult.
- $\mathrm{prox}_f$ is computable when the $A_i$ are nice but requires special functions (Lambert $W$-function).
- Projecting (in the euclidean norm) onto the constraint set requires sorting.

Using a primal-dual formulation of the problem, we have

$$\min_{\substack{X \in \mathbb{R}_+^{n \times m} \\ X \mathbb{1}_m = \mathbb{1}_n}} \max_{\mu \in \mathbb{R}^{m(n-1)}} \quad \sum_{i=1}^{n} \mathrm{KL}\left(A_i x_i, y_i\right) + \left\langle \nabla_{\mathrm{row}} X, \mu \right\rangle - \iota_{\mathcal{B}_\infty^\beta}\left(\mu\right).$$

This formulation solves the issue of computing $\mathrm{prox}_{g \circ \nabla_{\mathrm{row}}}$.

Using a primal-dual formulation of the problem, we have

$$\min_{\substack{X \in \mathbb{R}_+^{n \times m} \\ X \mathbb{1}_m = \mathbb{1}_n}} \max_{\mu \in \mathbb{R}^{m(n-1)}} \quad \sum_{i=1}^{n} \mathrm{KL}\left(A_i x_i, y_i\right) + \langle \nabla_{\mathrm{row}} X, \mu \rangle - \iota_{\mathcal{B}_\infty^\beta}\left(\mu\right).$$

This formulation solves the issue of computing $\mathrm{prox}_{g \circ \nabla_{\mathrm{row}}}$.
Remaining obstacles:

- Projection onto the constraint set.

Using a primal-dual formulation of the problem, we have

$$\min_{\substack{X \in \mathbb{R}_+^{n \times m} \\ X \mathbb{1}_m = \mathbb{1}_n}} \max_{\mu \in \mathbb{R}^{m(n-1)}} \quad \sum_{i=1}^{n} \mathrm{KL}\left(A_i x_i, y_i\right) + \langle \nabla_{\mathrm{row}} X, \mu \rangle - \iota_{\mathcal{B}_\infty^\beta}\left(\mu\right).$$

This formulation solves the issue of computing $\mathrm{prox}_{g \circ \nabla_{\mathrm{row}}}$.
Remaining obstacles:

- Projection onto the constraint set.
- Utilize differentiability of $\sum_{i=1}^{n} K\left(A_i x_i, y_i\right)$.

Let $\mathcal{X}_p$ and $\mathcal{X}_d$ be reflexive Banach spaces.

### Primal-dual problem

$$\min_{x \in \mathcal{X}_p} \max_{\mu \in \mathcal{X}_d} \quad \underbrace{f(x) + g(x) + \langle Tx, \mu \rangle - h^*(\mu) - \ell^*(\mu) + \iota_{\mathcal{C}_p}(x) - \iota_{\mathcal{C}_d}(\mu)}_{\mathcal{L}(x,\mu)}$$

Let $\mathcal{X}_p$ and $\mathcal{X}_d$ be reflexive Banach spaces.

### Primal-dual problem

$$\min_{x \in \mathcal{X}_p} \max_{\mu \in \mathcal{X}_d} \quad \underbrace{f(x) + g(x) + \langle Tx, \mu \rangle - h^*(\mu) - \ell^*(\mu) + \iota_{\mathcal{C}_p}(x) - \iota_{\mathcal{C}_d}(\mu)}_{\mathcal{L}(x, \mu)}$$

- $\mathcal{C}_p$ and $\mathcal{C}_d$ are nonempty closed convex subsets.

Let $\mathcal{X}_p$ and $\mathcal{X}_d$ be reflexive Banach spaces.

## Primal-dual problem

$$\min_{x \in \mathcal{X}_p} \max_{\mu \in \mathcal{X}_d} \underbrace{f(x) + g(x) + \langle Tx, \mu \rangle - h^*(\mu) - \ell^*(\mu) + \iota_{\mathcal{C}_p}(x) - \iota_{\mathcal{C}_d}(\mu)}_{\mathcal{L}(x,\mu)}$$

- $\mathcal{C}_p$ and $\mathcal{C}_d$ are nonempty closed convex subsets.
- $f$ and $h^*$ are relatively smooth with respect to $\phi_p$ and $\phi_d$, respectively.

Let $\mathcal{X}_p$ and $\mathcal{X}_d$ be reflexive Banach spaces.

### Primal-dual problem

$$\min_{x \in \mathcal{X}_p} \max_{\mu \in \mathcal{X}_d} \underbrace{f(x) + g(x) + \langle Tx, \mu \rangle - h^*(\mu) - \ell^*(\mu) + \iota_{\mathcal{C}_p}(x) - \iota_{\mathcal{C}_d}(\mu)}_{\mathcal{L}(x,\mu)}$$

- $\mathcal{C}_p$ and $\mathcal{C}_d$ are nonempty closed convex subsets.
- $f$ and $h^*$ are relatively smooth with respect to $\phi_p$ and $\phi_d$, respectively.
- $T$ is a bounded linear operator.

# A Different Kind of Distance

## Bregman divergence

Let $\mathcal{X}$ be a Banach space and define the *Bregman divergence* of a differentiable function $f : \mathcal{C} \subset \mathcal{X} \to \mathbb{R}$, for any $u, v \in \mathcal{C}$,

$$D_f(u, v) \stackrel{\text{def}}{=} f(u) - f(v) - \langle \nabla f(v), u - v \rangle.$$

## Bregman divergence

Let $\mathcal{X}$ be a Banach space and define the *Bregman divergence* of a differentiable function $f : \mathcal{C} \subset \mathcal{X} \to \mathbb{R}$, for any $u, v \in \mathcal{C}$,

$$D_f(u, v) \overset{\text{def}}{=} f(u) - f(v) - \langle \nabla f(v), u - v \rangle.$$

- $D_f(u, v)$ is a sort of distance between $u$ and $v$. For the euclidean squared norm $f(x) = \frac{1}{2} \|x\|_2^2$, it holds

$$D_f(u, v) = \frac{1}{2} \|u - v\|_2^2.$$

## Euclidean prox operator

Given a function $f : \mathcal{H} \to \mathbb{R} \cup \{+\infty\}$, we define the proximal operator

$$\mathrm{prox}_f\left(u\right) \stackrel{\mathrm{def}}{=} \underset{v \in \mathcal{H}}{\mathrm{argmin}} \left\{ f\left(v\right) + \frac{1}{2} \left\| v - u \right\|_2^2 \right\}.$$

# $D$-prox Operators

## Relative smoothness

$f$ is *relatively smooth* [Bauschke et al. 2017], [Lu et al. 2018] with respect to a differentiable function $\phi : \mathcal{C} \subset \mathcal{X} \to \mathbb{R}$ if there exists $L > 0$ such that, for any $u, v \in \mathcal{X}$,

$$D_f(u, v) \leq L D_\phi(u, v)$$

(equivalently, if $L\phi - f$ is a convex function).

# Going Beyond Lipschitz-smoothness

## Relative smoothness

$f$ is *relatively smooth* [Bauschke et al. 2017], [Lu et al. 2018] with respect to a differentiable function $\phi : \mathcal{C} \subset \mathcal{X} \to \mathbb{R}$ if there exists $L > 0$ such that, for any $u, v \in \mathcal{X}$,

$$D_f(u, v) \leq L D_\phi(u, v)$$

(equivalently, if $L\phi - f$ is a convex function).

- Lipschitz-smooth functions in $\Gamma_0(\mathcal{X})$ are relatively smooth with respect to the euclidean squared norm $\frac{1}{2} \|\cdot\|_2^2$:

$$D_f(u, v) \leq L \|u - v\|_2^2$$
$$\implies f(u) \leq f(v) + \langle \nabla f(v), u - v \rangle + L \|u - v\|_2^2$$
$$\implies f \text{ is } L\text{-smooth (Baillon-Haddad Theorem)}.$$

| Algorithm: | Bregman Primal-Dual ( BPD) |
|---|---|

Input: $x_0 \in \mathcal{C}_p,\ \mu_0 \in \mathcal{C}_d,\ (\lambda_k)_{k \in \mathbb{N}},\ (\nu_k)_{k \in \mathbb{N}},$
   $\phi_p : \mathcal{X}_p \to \mathbb{R} \cup \{+\infty\},\ \phi_d : \mathcal{X}_d \to \mathbb{R} \cup \{+\infty\}.$

$k = 0$

repeat

$$x_{k+1} = \operatorname*{argmin}_{x \in \mathcal{C}_p} \Big\{ g(x) + \big\langle \nabla f(x_k) \qquad, x \big\rangle$$

$$+ \big\langle x, T^*\mu_k \big\rangle + \tfrac{1}{\lambda_k} D_{\phi_p}(x, x_k) \Big\}$$

$$\mu_{k+1} = \operatorname*{argmin}_{\mu \in \mathcal{C}_d} \Big\{ \ell^*(\mu) + \big\langle \nabla h^*(\mu_k) \qquad, \mu \big\rangle$$

$$- \big\langle T(2x_{k+1} - x_k), \mu \big\rangle + \tfrac{1}{\nu_k} D_{\phi_d}(\mu, \mu_k) \Big\}$$

   $k \leftarrow k + 1$

until *convergence*;

Output: $x_k, \mu_k.$

| Algorithm: | Bregman Primal-Dual ( BPD) |
|---|---|

Input: $x_0 \in \mathcal{C}_p$, $\mu_0 \in \mathcal{C}_d$, $(\lambda_k)_{k\in\mathbb{N}}$, $(\nu_k)_{k\in\mathbb{N}}$,
$\quad\quad \phi_p : \mathcal{X}_p \rightarrow \mathbb{R} \cup \{+\infty\}$, $\phi_d : \mathcal{X}_d \rightarrow \mathbb{R} \cup \{+\infty\}$.
$k = 0$
repeat

$\quad x_{k+1} = \underset{x\in\mathcal{C}_p}{\operatorname{argmin}} \left\{ g(x) + \left\langle \nabla f(x_k) \quad , x \right\rangle \right.$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad + \left\langle x, T^*\mu_k \right\rangle + \frac{1}{\lambda_k} D_{\phi_p}(x, x_k) \Big\}$

$\quad \mu_{k+1} = \underset{\mu\in\mathcal{C}_d}{\operatorname{argmin}} \left\{ \ell^*(\mu) + \left\langle \nabla h^*(\mu_k) \quad , \mu \right\rangle \right.$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad - \left\langle T(2x_{k+1} - x_k), \mu \right\rangle + \frac{1}{\nu_k} D_{\phi_d}(\mu, \mu_k) \Big\}$

$\quad k \leftarrow k + 1$
until *convergence*;
Output: $x_k, \mu_k$.

---

**Algorithm: Stochastic Bregman Primal-Dual (SBPD)**

---

Input: $x_0 \in \mathcal{C}_p$, $\mu_0 \in \mathcal{C}_d$, $(\lambda_k)_{k \in \mathbb{N}}$, $(\nu_k)_{k \in \mathbb{N}}$,

$\qquad \phi_p : \mathcal{X}_p \to \mathbb{R} \cup \{+\infty\}$, $\phi_d : \mathcal{X}_d \to \mathbb{R} \cup \{+\infty\}$.

$k = 0$

repeat

$\quad x_{k+1} = \underset{x \in \mathcal{C}_p}{\mathrm{argmin}} \left\{ g\left(x\right) + \left\langle \nabla f\left(x_k\right) + \delta_k^p, x \right\rangle \right.$

$\qquad\qquad\qquad\qquad\qquad\qquad \left. + \left\langle x, T^* \mu_k \right\rangle + \frac{1}{\lambda_k} D_{\phi_p}\left(x, x_k\right) \right\}$

$\quad \mu_{k+1} = \underset{\mu \in \mathcal{C}_d}{\mathrm{argmin}} \left\{ \ell^*\left(\mu\right) + \left\langle \nabla h^*\left(\mu_k\right) + \delta_k^d, \mu \right\rangle \right.$

$\qquad\qquad\qquad\qquad\qquad\qquad \left. - \left\langle T\left(2x_{k+1} - x_k\right), \mu \right\rangle + \frac{1}{\nu_k} D_{\phi_d}\left(\mu, \mu_k\right) \right\}$

$\quad k \leftarrow k + 1$

until *convergence*;

Output: $x_k, \mu_k$.

---

# Interpretation of the Algorithm

Alternatively,

$$x_{k+1} = \underbrace{[\nabla\phi_p + \lambda_k \partial g]^{-1}}_{\text{Backward step}} \underbrace{(\nabla\phi_p(x_k) - \lambda_k \nabla f(x_k) - \lambda_k T^* \mu_k)}_{\text{Forward step}};$$

$$\mu_{k+1} = \underbrace{[\nabla\phi_d + \nu_k \partial \ell^*]^{-1}}_{\text{Backward step}} \underbrace{(\nabla\phi_d(\mu_k) - \nu_k \nabla h^*(\mu_k) + \nu_k T(2x_{k+1} - x_k))}_{\text{Forward step}}$$

Alternatively,

$$x_{k+1} = \underbrace{[\nabla\phi_p + \lambda_k\partial g]^{-1}}_{\text{Backward step}} \underbrace{(\nabla\phi_p(x_k) - \lambda_k\nabla f(x_k) - \lambda_k T^*\mu_k)}_{\text{Forward step}};$$

$$\mu_{k+1} = \underbrace{[\nabla\phi_d + \nu_k\partial\ell^*]^{-1}}_{\text{Backward step}} \underbrace{(\nabla\phi_d(\mu_k) - \nu_k\nabla h^*(\mu_k) + \nu_k T(2x_{k+1} - x_k))}_{\text{Forward step}}$$

- $\phi_p = \frac{1}{2}\|\cdot\|_2^2 \implies \nabla\phi_p = \text{Id}$ (likewise for $\phi_d$).

Alternatively,

$$x_{k+1} = \underbrace{[\nabla\phi_p + \lambda_k\partial g]^{-1}}_{\text{Backward step}} \underbrace{(\nabla\phi_p(x_k) - \lambda_k\nabla f(x_k) - \lambda_k T^*\mu_k)}_{\text{Forward step}};$$

$$\mu_{k+1} = \underbrace{[\nabla\phi_d + \nu_k\partial\ell^*]^{-1}}_{\text{Backward step}} \underbrace{(\nabla\phi_d(\mu_k) - \nu_k\nabla h^*(\mu_k) + \nu_k T(2x_{k+1} - x_k))}_{\text{Forward step}}$$

- $\phi_p = \frac{1}{2}\|\cdot\|_2^2 \implies \nabla\phi_p = \mathrm{Id}$ (likewise for $\phi_d$).

- Flavor of mirror descent [Nemirovsky et al. 83], Chambolle-Pock [Chambolle et al. 2011], [Chambolle et al., 2016], NoLips [Bauschke et al. 2017], Bregman Forward-Backward [Nguyen, 2017], etc.

- $f$ is $L_p$ relatively smooth with respect to $\phi_p$ on $\operatorname{int}(\mathcal{C}_p)$.

- $f$ is $L_p$ relatively smooth with respect to $\phi_p$ on $\operatorname{int}(\mathcal{C}_p)$.
- $h^*$ is $L_d$ relatively smooth with respect to $\phi_d$ on $\operatorname{int}(\mathcal{C}_d)$.

## Matching the Geometries

- $f$ is $L_p$ relatively smooth with respect to $\phi_p$ on $\operatorname{int}(\mathcal{C}_p)$.
- $h^*$ is $L_d$ relatively smooth with respect to $\phi_d$ on $\operatorname{int}(\mathcal{C}_d)$.
- $\phi_p$ and $\phi_d$ are Legendre functions with domains $\mathcal{C}_p$ and $\mathcal{C}_d$ and the mappings $[\nabla\phi_p + \lambda_k\partial g]^{-1}$ and $[\nabla\phi_d + \nu_k\partial\ell^*]^{-1}$ are well-defined.

# Matching the Geometries

- $f$ is $L_p$ relatively smooth with respect to $\phi_p$ on $\text{int}\,(\mathcal{C}_p)$.
- $h^*$ is $L_d$ relatively smooth with respect to $\phi_d$ on $\text{int}\,(\mathcal{C}_d)$.
- $\phi_p$ and $\phi_d$ are Legendre functions with domains $\mathcal{C}_p$ and $\mathcal{C}_d$ and the mappings $[\nabla\phi_p + \lambda_k \partial g]^{-1}$ and $[\nabla\phi_d + \nu_k \partial \ell^*]^{-1}$ are well-defined.

## Note

The geometry of $\phi_p$ and $\phi_d$ must match the problem!

- $f$ is $L_p$ relatively smooth with respect to $\phi_p$ on $\operatorname{int}(\mathcal{C}_p)$.
- $h^*$ is $L_d$ relatively smooth with respect to $\phi_d$ on $\operatorname{int}(\mathcal{C}_d)$.
- $\phi_p$ and $\phi_d$ are Legendre functions with domains $\mathcal{C}_p$ and $\mathcal{C}_d$ and the mappings $[\nabla\phi_p + \lambda_k \partial g]^{-1}$ and $[\nabla\phi_d + \nu_k \partial\ell^*]^{-1}$ are well-defined.

### Note

The geometry of $\phi_p$ and $\phi_d$ must match the problem!

### Note

We do not assume strong convexity of $\phi_p$ or $\phi_d$ (cf. [Chambolle et al., 2016], [Van Dung et al., 2021]).

Define

$$M\left(x_1, x_2, \mu_1, \mu_2\right) \stackrel{\text{def}}{=} \left\langle T\left(x_1 - x_2\right), \mu_1 - \mu_2 \right\rangle$$

We require existence of two functions $r_p : \mathcal{X}_p^2 \to \mathbb{R}_+$ and $r_d : \mathcal{X}_d^2 \to \mathbb{R}_+$ and $\epsilon \geq 0$ such that, for all $x_1 \in \mathcal{C}_p \cap \mathrm{dom}\partial g$, $x_2 \in \mathrm{int}\left(\mathcal{C}_p\right) \cap \mathrm{dom}\partial g$, $\mu_1 \in \mathcal{C}_d \cap \mathrm{dom}\partial\ell$, and $\mu_2 \in \mathrm{int}\left(\mathcal{C}_d\right) \cap \mathrm{dom}\partial\ell$,

$$\frac{\left(\lambda_\infty^{-1} - L_p\right) D_{\phi_p}\left(x_1, x_2\right) + \left(\nu_\infty^{-1} - L_d\right) D_{\phi_d}\left(\mu_1, \mu_2\right) - M\left(x_1, x_2, \mu_1, \mu_2\right)}{r_p\left(x_1, x_2\right) + r_d\left(\mu_1, \mu_2\right)}$$

$$\geq \epsilon$$

$$\frac{\left(\lambda_\infty^{-1} - L_p\right) D_{\phi_p}\left(x_1, x_2\right) + \left(\nu_\infty^{-1} - L_d\right) D_{\phi_d}\left(\mu_1, \mu_2\right) - M\left(x_1, x_2, \mu_1, \mu_2\right)}{r_p\left(x_1, x_2\right) + r_d\left(\mu_1, \mu_2\right)}$$

$$\geq \epsilon$$

Consider $\phi_p = \frac{1}{2}\left\|\cdot\right\|_2^2$ and $\phi_d = \frac{1}{2}\left\|\cdot\right\|_2^2$,

$$\frac{\lambda_\infty^{-1} - L_p}{2}\left\|x_1 - x_2\right\|_2^2 + \frac{\nu_\infty^{-1} - L_d}{2}\left\|\mu_1 - \mu_2\right\|_2^2 - \left\langle T\left(x_1 - x_2\right), \mu_1 - \mu_2\right\rangle$$

$$\geq \underbrace{\frac{\lambda_\infty^{-1} - L_p - \|T\|}{2}\left\|x_1 - x_2\right\|_2^2}_{r_p(x_1, x_2)} + \underbrace{\frac{\nu_\infty^{-1} - L_d - 1}{2}\left\|\mu_1 - \mu_2\right\|_2^2}_{r_d(\mu_1, \mu_2)}$$

## Theorem (Ergodic Convergence Rate)

Define $\bar{x}_k \overset{\text{def}}{=} \frac{1}{k} \sum\limits_{i=0}^{k} x_i$, $\bar{\mu}_k \overset{\text{def}}{=} \frac{1}{k} \sum\limits_{i=0}^{k} \mu_i$, and, for $w \overset{\text{def}}{=} (x, \mu)$,
$M(w, w') = \langle T(x - x'), \mu - \mu' \rangle$. Under [assumptions], for each
$k \in \mathbb{N}$, for every $w \in \mathcal{C}_p \times \mathcal{C}_d$,

$$\mathcal{L}(\bar{x}_k, \mu) - \mathcal{L}(x, \bar{\mu}_k) \leq \frac{\Lambda_0^{-1} D_{\phi_p, \phi_d}(w, w_0) - M(w, w_0)}{k}.$$

In particular, every weak cluster point of the sequence $(\bar{x}_k, \bar{\mu}_k)_{k \in \mathbb{N}}$
is a solution to the primal-dual problem.

## Theorem (Ergodic Convergence Rate)

*Under [the same assumptions], if the errors $\delta_k^p$ and $\delta_k^d$ are unbiased conditioned on the the previous iterates, for each $k \in \mathbb{N}$, for every $w \in \mathcal{C}_p \times \mathcal{C}_d$,*

$$\mathbb{E}\left[\mathcal{L}\left(\bar{x}_k, \mu\right) - \mathcal{L}\left(x, \bar{\mu}_k\right)\right] \leq \frac{\Lambda_0^{-1} D_{\phi_p, \phi_d}\left(w, w_0\right) - M\left(w, w_0\right)}{k}$$

$$+ \frac{\sum\limits_{i=0}^{k-1} \mathbb{E}\left[\langle \Delta_i, w - w_{i+1}\rangle\right]}{k}.$$

*In particular, every almost sure weak cluster point of the sequence $\left(\bar{x}_k, \bar{\mu}_k\right)_{k\in\mathbb{N}}$ is a solution to the primal-dual problem in expectation ($\mathbb{E}\left[\left(x_\infty, \mu_\infty\right)\right]$ is a saddle-point).*

# Trend Filtering

## Trend filtering problem - primal-dual formulation

$$\min_{\substack{X \in \mathbb{R}_+^{n \times m} \\ X \mathbb{1}_m = \mathbb{1}_n}} \max_{\mu \in \mathbb{R}^{m(n-1)}} \sum_{i=1}^{n} \mathrm{KL}\left(A_i x_i, y_i\right) + \left\langle \nabla_{\mathrm{row}} X, \mu \right\rangle - \iota_{\mathcal{B}_\infty^\beta}\left(\mu\right).$$

## Trend filtering problem - primal-dual formulation

$$\min_{\substack{X \in \mathbb{R}_+^{n \times m} \\ X\mathbb{1}_m = \mathbb{1}_n}} \max_{\mu \in \mathbb{R}^{m(n-1)}} \quad \sum_{i=1}^n \mathrm{KL}\left(A_i x_i, y_i\right) + \langle \nabla_{\mathrm{row}} X, \mu \rangle - \iota_{\mathcal{B}_\infty^\beta}\left(\mu\right).$$

Apply SBPD with

$$f(X) = \sum_{i=1}^n \mathrm{KL}\left(A_i x_i, y_i\right), \quad g(X) = \iota_{\mathbb{1}_n}\left(X\mathbb{1}_m\right), \quad \mathcal{C}_p = \mathbb{R}_+^{n \times m},$$

$$T = \nabla_{\mathrm{row}} \quad h^*\left(\mu\right) = 0, \quad \ell^*\left(\mu\right) = \iota_{\mathcal{B}_\infty^\beta}\left(\mu\right) \quad \text{and} \quad \mathcal{C}_d = \mathbb{R}^{m(n-1)}$$

# Choosing $\phi_p$ and $\phi_d$

## Primal entropy $\phi_p$

- $\mathcal{C}_p = \mathbb{R}_+^{n \times m}$

$$\phi_p(X) = \sum_{i=1}^{n} \sum_{j=1}^{m} X_{i,j} \log(X_{i,j}).$$

- Must show $\exists L_p > 0$ such that $L_p \phi_p - f$ is convex.
- Must compute $\operatorname{prox}_{\lambda_k g}^{D_{\phi_p}}(X)$.

# Choosing $\phi_p$ and $\phi_d$

## Primal entropy $\phi_p$

- $\mathcal{C}_p = \mathbb{R}_+^{n \times m}$

$$\phi_p(X) = \sum_{i=1}^{n} \sum_{j=1}^{m} X_{i,j} \log(X_{i,j}).$$

- Must show $\exists L_p > 0$ such that $L_p \phi_p - f$ is convex.
- Must compute $\text{prox}_{\lambda_k g}^{D_{\phi_p}}(X)$.

## Dual entropy $\phi_d$

- $\mathcal{C}_d = \mathbb{R}^{m(n-1)}$ (trivial constraint)

$$\phi_d(\mu) = \frac{1}{2} \|\mu\|_2^2.$$

- Euclidean prox of $\ell^*(\mu) = \iota_{\mathcal{B}_\infty^\beta}$ is accessible.

# New Geometry of $\phi_p$

## Relative smoothness

For each $i \in \{1, \dots, n\}$, let $L_i \geq \max\limits_{1 \leq q \leq m} \sum\limits_{j=1}^{p} A_i(j, q)$ and let $L_p = \max\limits_{1 \leq i \leq n} L_i$. Then $L_p \phi_p - f$ is convex on $\operatorname{int}(\mathcal{C}_p)$.

# New Geometry of $\phi_p$

## Relative smoothness

For each $i \in \{1, \ldots, n\}$, let $L_i \geq \max\limits_{1 \leq q \leq m} \sum\limits_{j=1}^{p} A_i(j, q)$ and let $L_p = \max\limits_{1 \leq i \leq n} L_i$. Then $L_p \phi_p - f$ is convex on $\text{int}(\mathcal{C}_p)$.
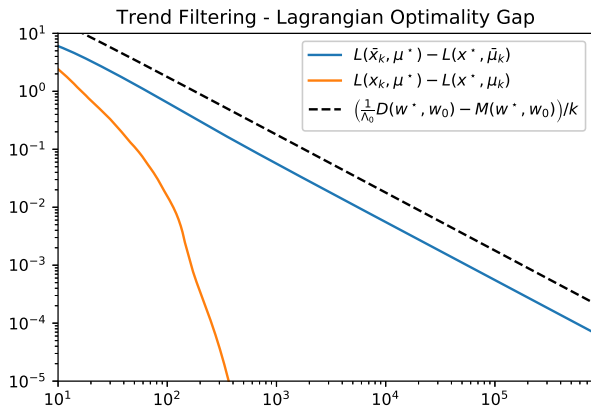
## $D$-prox under $\phi_p$

For each $X \in \mathcal{C}_p$,

$$\text{prox}_{\lambda_k g}^{D_{\phi_p}}(X) = \underset{\substack{U \in \mathbb{R}_+^{n \times m} \\ U^T \mathbb{1}_m = \mathbb{1}_n}}{\text{argmin}} \left\{ D_{\phi_p}(U, X) \right\} = \left( \frac{\exp(X_{i,j})}{\sum\limits_{q=1}^{m} \exp(X_{i,q})} \right)_{i,j}$$

i.e., project each row onto the simplex under $D_{\phi_p}$.

We take $n = 100$, $m = 3$ and $\beta = 1$ with synthetic (randomly generated) data $Y$ and $A_i = \mathrm{Id}$.



Trend Filtering - Lagrangian Optimality Gap

Legend:
- $L(\bar{x}_k, \mu^\star) - L(x^\star, \bar{\mu}_k)$
- $L(x_k, \mu^\star) - L(x^\star, \mu_k)$
- $\left(\frac{1}{\Lambda_0} D(w^\star, w_0) - M(w^\star, w_0)\right)/k$

Simplest case: discrete measures $\rho$ and $\theta$ with ground cost matrix $C \in \mathbb{R}_+^{n \times m}$.

## Entropically regularized Wasserstein distance

$$W_\gamma \left( \rho, \theta \right) = \inf_{\pi \in \Pi(\rho, \theta)} \left\{ \gamma \mathrm{KL} \left( \pi, \exp \left( -\gamma^{-1} C \right) \right) \right\}.$$

where $\Pi \left( \rho, \theta \right) \stackrel{\text{def}}{=} \left\{ \pi \in \mathbb{R}_+^{n \times m} : \pi \mathbb{1}_m = \rho, \pi^T \mathbb{1}_n = \theta \right\}$

# Entropically Regularized Wasserstein Inverse Problems

Simplest case: discrete measures $\rho$ and $\theta$ with ground cost matrix $C \in \mathbb{R}_+^{n \times m}$.

## Entropically regularized Wasserstein distance

$$W_\gamma (\rho, \theta) = \inf_{\pi \in \Pi(\rho, \theta)} \left\{ \gamma \mathrm{KL} \left( \pi, \exp \left( -\gamma^{-1} C \right) \right) \right\}.$$

where $\Pi (\rho, \theta) \stackrel{\text{def}}{=} \left\{ \pi \in \mathbb{R}_+^{n \times m} : \pi \mathbb{1}_m = \rho, \pi^T \mathbb{1}_n = \theta \right\}$

## Inverse problem

$$\min_{\rho \in \Delta^n} \underbrace{\inf_{\pi \in \Pi(F\rho, \theta)} \left\{ \gamma \mathrm{KL} \left( \pi, \exp \left( -\gamma^{-1} C \right) \right) \right\}}_{W_\gamma(F\rho, \theta)} + J \circ A (\rho),$$

where $J \in \Gamma_0 \left( \mathbb{R}^p \right)$, $F : \Delta^n \to \Delta^m$ is linear, and $A \in \mathbb{R}^{n \times p}$.

Since $W_\gamma(F\rho, \theta)$ is itself a minimization problem, we introduce a dual variable $\tau$ and use Lagrangian duality to have

$$\min_{\rho \in \Delta^n} \max_{\tau \in \mathbb{R}^m} \langle \tau, F\rho \rangle - \gamma \sum_{j=1}^m \theta_j \log \left( \sum_{i=1}^m \exp \left( \frac{\tau_i - C_{i,j}}{\gamma} \right) \right) + J \circ A(\rho)$$

# Primal-Dual Splitting for OT

Since $W_\gamma (F\rho, \theta)$ is itself a minimization problem, we introduce a dual variable $\tau$ and use Lagrangian duality to have

$$\min_{\rho \in \Delta^n} \max_{\tau \in \mathbb{R}^m} \langle \tau, F\rho \rangle - \gamma \sum_{j=1}^m \theta_j \log \left( \sum_{i=1}^m \exp \left( \frac{\tau_i - C_{i,j}}{\gamma} \right) \right) + J \circ A (\rho)$$

We can further dualize to split $J \circ A$,

$$\min_{\rho \in \Delta^n} \max_{\substack{\tau \in \mathbb{R}^m \\ \zeta \in \mathbb{R}^p}} \left\langle \begin{pmatrix} \tau \\ \zeta \end{pmatrix}, \begin{pmatrix} F\rho \\ A\rho \end{pmatrix} \right\rangle - \gamma \sum_{j=1}^m \theta_j \log \left( \sum_{i=1}^m \exp \left( \frac{\tau_i - C_{i,j}}{\gamma} \right) \right) - J^* (\zeta)$$

and now we can apply SBPD.

# Splitting the Inverse Problem

## Inverse problem - primal-dual formulation

$$\min_{\rho \in \Delta^n} \max_{\substack{\tau \in \mathbb{R}^m \\ \zeta \in \mathbb{R}^p}} \left\langle \begin{pmatrix} \tau \\ \zeta \end{pmatrix}, \begin{pmatrix} F\rho \\ A\rho \end{pmatrix} \right\rangle - \gamma \sum_{j=1}^{m} \theta_j \log \left( \sum_{i=1}^{m} \exp \left( \frac{\tau_i - C_{i,j}}{\gamma} \right) \right) - J^* (\zeta)$$

> **Inverse problem - primal-dual formulation**
>
> $$\min_{\rho \in \Delta^n} \max_{\substack{\tau \in \mathbb{R}^m \\ \zeta \in \mathbb{R}^p}} \left\langle \begin{pmatrix} \tau \\ \zeta \end{pmatrix}, \begin{pmatrix} F\rho \\ A\rho \end{pmatrix} \right\rangle - \gamma \sum_{j=1}^{m} \theta_j \log \left( \sum_{i=1}^{m} \exp \left( \frac{\tau_i - C_{i,j}}{\gamma} \right) \right) - J^* \left( \zeta \right)$$

Apply SBPD with

$$f\left(\rho\right) = 0, \quad g\left(\rho\right) = \iota_{\{1\}}\left(\rho^T \mathbb{1}_n\right), \quad \mathcal{C}_p = \mathbb{R}_+^n,$$

$$T\left(\rho\right) = \begin{pmatrix} F\rho \\ A\rho \end{pmatrix}, \quad h^*\left(\mu\right) = h^*\left(\tau\right) = \gamma \sum_{j=1}^{m} \theta_j \log \left( \sum_{i=1}^{m} \exp \frac{\tau_i - C_{i,j}}{\gamma} \right),$$

$$\ell^*\left(\mu\right) = \ell^*\left(\zeta\right) = J^*\left(\zeta\right), \quad \text{and} \quad \mathcal{C}_d = \mathbb{R}^{m+p}.$$

## Primal entropy $\phi_p$

- $\mathcal{C}_p = \mathbb{R}^n_+$

$$\phi_p(\rho) = \sum_{i=1}^{n} \rho_i \log(\rho_i).$$

- $\mathrm{prox}_{\lambda_k g}^{D_{\phi_p}}$ is same as in trend filtering (consider 1 row).

# Choosing $\phi_p$ and $\phi_d$

## Primal entropy $\phi_p$

- $\mathcal{C}_p = \mathbb{R}_+^n$

$$\phi_p(\rho) = \sum_{i=1}^{n} \rho_i \log(\rho_i).$$

- $\mathrm{prox}_{\lambda_k g}^{D_{\phi_p}}$ is same as in trend filtering (consider 1 row).
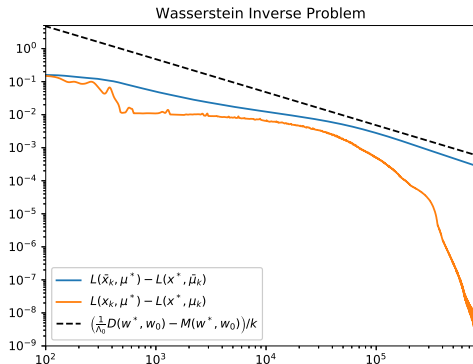
## Dual entropy $\phi_d$

- $\mathcal{C}_d = \mathbb{R}^{m+p}$ (trivial constraint)
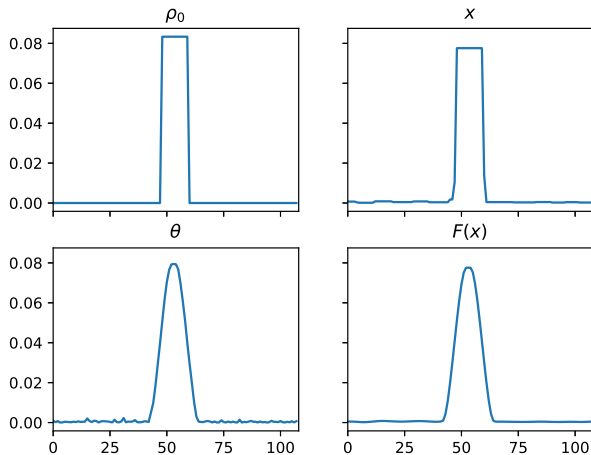
$$\phi_d(\mu) = \frac{1}{2} \|\mu\|_2^2$$

- Must show that $h^*$ is Lipschitz-smooth (straightforward).

# An Example Problem

- $n = 108$,
- $C_{i,j} = \frac{1}{2} \|i - j\|_2^2$,
- $F$ - convolution operator (bump function),
- $J \circ A = \|\cdot\|_1 \circ \nabla$.



Wasserstein Inverse Problem

## Contributions

- No Lipschitz-smoothness assumptions: $\nabla\mathrm{KL}$ vs $\mathrm{prox}_{\mathrm{KL}}$.

## Contributions

- No Lipschitz-smoothness assumptions: $\nabla \mathrm{KL}$ vs $\mathrm{prox}_{\mathrm{KL}}$.
- Improved constants: $\|x - x_0\|_2^2$ vs $D_{\phi_p}(x, x_0)$.

## Contributions

- No Lipschitz-smoothness assumptions: $\nabla \mathrm{KL}$ vs $\mathrm{prox}_{\mathrm{KL}}$.
- Improved constants: $\|x - x_0\|_2^2$ vs $D_{\phi_p}(x, x_0)$.
- Improved complexities: sorting vs softmax.

# Let's Recap

## Contributions

- No Lipschitz-smoothness assumptions: $\nabla \text{KL}$ vs $\text{prox}_{\text{KL}}$.
- Improved constants: $\|x - x_0\|_2^2$ vs $D_{\phi_p}(x, x_0)$.
- Improved complexities: sorting vs softmax.
- Stochasticity permitted.

## Contributions

- No Lipschitz-smoothness assumptions: $\nabla \text{KL}$ vs $\text{prox}_{\text{KL}}$.
- Improved constants: $\|x - x_0\|_2^2$ vs $D_{\phi_p}(x, x_0)$.
- Improved complexities: sorting vs softmax.
- Stochasticity permitted.
- Infinite-dimensional problems (Reflexive Banach spaces) permitted.

## Note

Code (NumPy) is available on github soon.

Thanks for listening.

Full paper available on arxiv: https://arxiv.org/abs/ 2112.11928

"A Stochastic Bregman Primal-Dual Splitting Algorithm for Composite Optimization" - Antonio Silveti-Falls, Cesare Molinari, Jalal Fadili.