

Lecture 2:
Optimal Transport,
Sinkhorn's Algorithm,
SMART

Gabriele Steidl
Applied Mathematics - Imaging Sciences
TU Berlin

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Outline

1. Basic Notation: Spaces of Measures
2. Monge and Kantorovich Problem of OT
3. Discrete OT and its Dual
4. Regularized Optimal Transport
5. Mirror Descent Algorithm - SMART - Sinkhorn Algorithm

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

1. Basic Notation: Spaces of Measures

- ◆ $X = \mathbb{R}^d$ or $X \subset \mathbb{R}^d$ compact subset of \mathbb{R}^d
 (General: Polish space = separable, completely metrizable topol. space)
- ◆ $\mathcal{B}(X)$ Borel σ -algebra on X : smallest σ -algebra which contains all open sets
 = σ -algebra generated by the open sets
 (Remember: σ -algebra Σ : $X \in \Sigma$, $A \in \Sigma \Rightarrow \bar{A} \in \Sigma$, $A_k \in \Sigma \Rightarrow \bigcup_{k \in \mathbb{N}} A_k \in \Sigma$)
- ◆ $\mu : \mathcal{B}(X) \rightarrow \mathbb{R}$ is a finite, signed measure on $\mathcal{B}(X)$ (Borel measure)
 i.e. for pairwise disjoint sets $A_k \in \mathcal{B}(X)$, $k \in \mathbb{N}$,

$$\mu\left(\bigcup_{k \in \mathbb{N}} A_k\right) = \sum_{k \in \mathbb{N}} \mu(A_k) \quad \text{--- } \sigma\text{-additivity}$$

- ◆ $\mathcal{M}(X)$ linear space of Borel measures
 (General: for Radon measures (fulfill inner/outer regularity condition); coincide with Borel measures on Polish spaces, in particular \mathbb{R}^d)
- ◆ $\mathcal{M}(X)$ is a Banach space with **total variation norm**

$$\|\mu\|_{TV} := |\mu|(X), \quad \text{where}$$

$$|\mu|(A) := \sup_{A = \bigcup A_k, A_i \cap A_j = \emptyset} \sum_{k=1}^n |\mu(A_k)|$$

1	2
---	---

3	4
---	---

5	6
---	---

7	8
---	---

9	10
---	----

11	12
----	----

13	14
----	----

15	16
----	----

17	18
----	----

19	20
----	----

21	22
----	----

23	24
----	----

25	26
----	----

1. Basic Notation: Spaces of Measures

- ◆ pre-dual space of $\mathcal{M}(X)$ is $C_0(X)$, i.e. $C_0(X)' = \mathcal{M}(X)$
- ◆ weak-* convergence of measures $\mu_n \rightharpoonup \mu$ if

$$\int_X \varphi(x) d\mu_n(x) \rightarrow \int_X \varphi(x) d\mu(x) \quad \text{for all } \varphi \in C_0(X).$$



$$\|\mu\|_{TV} = \sup_{\|\varphi\|_\infty \leq 1} |\langle \varphi, \mu \rangle|$$

with dual pairing

$$\langle x, \mu \rangle := \int_X \varphi(x) d\mu(x)$$

Further, let

- ◆ $\mathcal{M}_+(X)$ subset of non-negative measures on X
- ◆ $\mathcal{P}(X)$ subset of probability measures on X ,
i.e. $\mu \in \mathcal{M}_+(X)$ and $\mu(X) = 1$

1	2
---	---

3	4
---	---

5	6
---	---

7	8
---	---

9	10
---	----

11	12
----	----

13	14
----	----

15	16
----	----

17	18
----	----

19	20
----	----

21	22
----	----

23	24
----	----

25	26
----	----

Remark on Convergence of Measures

Let X locally compact Polish space (separable, complete metric space) and $\text{rba}(X)$ space of finitely additive set functions

- ◆ $C_c(X)' \cong \mathcal{M}(X)$ with $T_\mu(\varphi) = \langle \mu, \varphi \rangle := \int_X \varphi d\mu$ for all $\varphi \in C_c(X)$.
- ◆ $C_0(X)' \cong \mathcal{M}(X)$ with $T_\mu(\varphi) = \langle \mu, \varphi \rangle := \int_X \varphi d\mu$ for all $\varphi \in C_0(X)$.
- ◆ $C_b(X)' \cong \text{rba}(X)$

Note that $C_0(X) \subset C_b(X)$, but $C_b(X)' \cong \text{rba}(X) \not\subseteq \mathcal{M}(X) \cong C_0(X)'$

- ◆ $(\mu_n)_{n \in \mathbb{N}} \subset \mathcal{M}(X)$ bounded Then

$$\mu_n \xrightarrow{*} \mu \text{ in } C_c(X)' \quad \text{if and only if} \quad \mu_n \xrightarrow{*} \mu \text{ in } C_0(X)'.$$

Counterexample: $\mu_n = n\delta_n$ test against $\varphi(x) = \sin(x)/x \in C_0(\mathbb{R})$

- ◆ $(\mu_n)_{n \in \mathbb{N}} \subset \mathcal{M}(X)$ be bounded and tight, then

$$\mu_n \xrightarrow{*} \mu \text{ in } C_c(X)' \quad \text{if and only if} \quad \mu_n \xrightarrow{*} \mu \text{ in } C_b(X)'.$$

Counterexample: $\mu_n = \delta_n$ test against $\varphi(x) = \sin(x) \in C_b(\mathbb{R})$

$(\mu_n)_n$ is *tight* if for all $\epsilon > 0$, there exists a compact set $K \subset X$ such that $|\mu_n|(X \setminus K) < \epsilon$ for all $n \in \mathbb{N}$.

Theorem (Prokhorov) Let $(\mu_n)_{n \in \mathbb{N}}$ be a tight sequence of probability measures. Then there exists a subsequence $(\mu_{n_k})_{k \in \mathbb{N}}$ which converges weakly to a probability measure μ .

Ref: Plonka, Potts, Steidl, Tasche, Numerical Fourier Analysis, Springer 2023,
 Chapter: Fourier Analysis of Measures

1	2
---	---

3	4
---	---

5	6
---	---

7	8
---	---

9	10
---	----

11	12
----	----

13	14
----	----

15	16
----	----

17	18
----	----

19	20
----	----

21	22
----	----

23	24
----	----

25	26
----	----

Examples

1. Atomic measures (empirical measures if same weights)

$$\mu = \sum_{i=1}^m \mu_i \delta_{x_i}, \quad \delta_x(A) := \begin{cases} 1 & \text{if } x \in A, \\ 0 & \text{otherwise} \end{cases}$$

If $\mu \in \mathcal{P}(X)$, then $\sum_{i=1}^m \mu_i = 1$ and $\mu_i > 0$, $i = 1, \dots, m$. Then

$$\|\mu\|_{TV} = |\mu_1| + \dots + |\mu_m|$$

2. Absolutely continuous measures $\mu \in \mathcal{M}(X)$ with density $\varphi \in L^1(X)$

$$\mu(A) = \int_A \varphi(x) dx$$

(Remember: $\mu \ll \lambda$ with Lebesgue measure λ , if $\lambda(A) = 0 \Rightarrow \mu(A) = 0 \forall A \in \mathcal{B}(X)$)

Then

$$\|\mu\|_{TV} = \int_X |\varphi(x)| dx$$

TV-norm is not a good measure for our purposes.

1	2
---	---

3	4
---	---

5	6
---	---

7	8
---	---

9	10
---	----

11	12
----	----

13	14
----	----

15	16
----	----

17	18
----	----

19	20
----	----

21	22
----	----

23	24
----	----

25	26
----	----



2. Monge Problem (1781)

Given $\mu \in \mathcal{P}(X)$, $\nu \in \mathcal{P}(Y)$

Find an optimal **transport map** $\hat{T} : X \rightarrow Y$ such that

$$\hat{T} \in \operatorname{argmin}_{T \text{ measurable}} \int_X c(x, T(x)) \, d\mu(x) \quad \text{subject to} \quad \nu = T_{\#}\mu$$

with the **push forward measure**

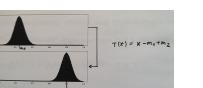
$$T_{\#}\mu := \mu \circ T^{-1}$$

Note that

$$\int_{T^{-1}(A)} h(T(x)) d\mu(x) = \int_A h(y) d\underbrace{(T_{\#}\mu)}_{\nu}(y)$$

and in case of existing densities and a diffeomorphism T :

$$p_{T_{\#}\mu}(y) = p_{\mu}(T^{-1}(y)) |\det \nabla T^{-1}(y)|$$



Example: law of random variables, sliced OT, measures on curves

1 2

3 4

5 6

7 8

9 10

11 12

13 14

15 16

17 18

19 20

21 22

23 24

25 26

Discrete Monge Problem

Given

$$\mu = \sum_{i=1}^m \mu_i \delta_{x_i}, \quad \nu = \sum_{j=1}^n \nu_j \delta_{y_j}$$

Find an optimal transport map $\hat{T} : \{x_1, \dots, x_m\} \rightarrow \{y_1, \dots, y_n\}$ such that

$$\hat{T} \in \operatorname{argmin}_T \sum_{i=1}^N c(x_i, T(x_i)) \mu_i \quad \text{s.t.} \quad \nu_j = \sum_{T(x_i)=y_j} \mu_i$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26



Kantorovich Problem (1942)

Given $\mu \in \mathcal{P}(X)$, $\nu \in \mathcal{P}(Y)$

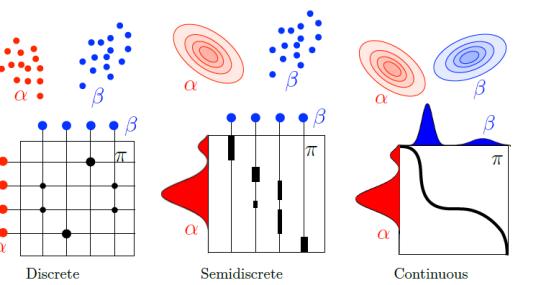
Find an optimal **transport plan** $\hat{\pi} \in \mathcal{P}(X \times Y)$ with given marginals μ and ν

$$\hat{\pi} \in \operatorname{argmin}_{\pi \in \Pi(\mu, \nu)} \int_{X \times X} c(x, y) d\pi(x, y)$$

$$\text{OT}(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \int_{X \times X} c(x, y) d\pi(x, y)$$

where

- ◆ $\Pi(\mu, \nu) := \{\pi \in \mathcal{P}(X \times X) : (P_1)_\# \pi = \mu, (P_2)_\# \pi = \nu$
- ◆ $P_1(x_1, x_2) = x_1, P_2(x_1, x_2) = x_2$



: Book Peyré/Cuturi: Computational optimal transport

- ◆ **Existence** of a minimizer is ensured if c is lsc and bounded from below

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Discrete Kantorovich Problem

Given

$$\mu = \sum_{i=1}^m \mu_i \delta_{x_i}, \quad \nu = \sum_{j=1}^n \nu_j \delta_{y_j}$$

Find an optimal **transport plan** $\hat{\pi} = \sum_{i,j=1}^{m,n} \hat{\pi}_{ij} \delta_{x_i, y_j}$, i.e. $(\pi_{ij})_{i,j} \in \mathbb{R}_{\geq 0}^{m,n}$,

$$\hat{\pi} \in \operatorname{argmin}_{\pi} \sum_{i=1}^m \sum_{j=1}^n c(x_i, y_j) \pi_{ij}$$

subject to $\pi_{ij} \geq 0$

$$\sum_{i=1}^m \pi_{ij} = \nu_j, \quad j = 1, \dots, n,$$

$$\sum_{j=1}^n \pi_{ij} = \mu_i, \quad i = 1, \dots, m$$

Shorter notation:

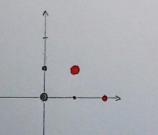
$$\hat{\pi} \in \operatorname{argmin}_{\pi} \langle c, \pi \rangle \quad \text{subject to} \quad \pi \mathbf{1}_n = \mu, \quad \pi^\top \mathbf{1}_m = \nu, \quad \pi \geq 0$$

$$P_1 \pi = \mu, \quad P_2 \pi = \nu, \quad \pi \geq 0$$

$$(1_n \otimes I_m) \pi = \mu, \quad (I_n \otimes 1_m) \pi = \nu, \quad \pi \geq 0$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Examples

<u>EXAMPLE ①</u>	$x_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$	$\mu = 3 \delta_{x_1} + 1 \delta_{x_2} + 2 \delta_{x_3}$
	$x_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$	
	$x_3 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$	
	$y_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$	$v = 4 \delta_{y_1} + 2 \delta_{y_2}$
	$y_2 = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$	
		
		$C(x, y) = \ x - y\ ^2$
		The only possible transport map is
		$T(x_1) = y_1, T(x_2) = y_1, T(x_3) = y_2$
		The optimal transport plan is
	$\begin{array}{c cc} & y_1 & y_2 \\ \hline x_1 & 2 & 4 \\ x_2 & 1 & 1 \\ x_3 & 1 & 5 \\ \hline C & \underbrace{}_{\frac{1}{2}} & \end{array}$	$\begin{array}{c cc} & 4 & 2 \\ \hline 3 & 2 & 1 \\ 1 & 0 & 1 \\ 2 & 2 & 0 \\ \hline \hat{T} & \underbrace{}_{\frac{1}{2}} & \end{array}$
		$\begin{array}{c cc} & 4 & 2 \\ \hline 3 & 3 & 0 \\ 1 & 1 & 0 \\ 2 & 0 & 2 \\ \hline \bar{T} = (\text{id}, T) \# \mu & \underbrace{}_{\frac{1}{2}} & \end{array}$
		(does not correspond to)

<u>②</u>	$\mu = 2 \delta_{y_1} + 4 \delta_{y_2}$; not the same
	Map:	$T(x_1) = y_2, T(x_2) = y_2, T(x_3) = y_1$
	$\begin{array}{c cc} & y_1 & y_2 \\ \hline x_1 & 4 & 2 \\ x_2 & 1 & 1 \\ x_3 & 5 & 1 \\ \hline C & \underbrace{}_{\frac{1}{2}} & \end{array}$	$\begin{array}{c cc} & 2 & 4 \\ \hline 3 & 0 & 3 \\ 1 & 0 & 1 \\ 2 & 2 & 0 \\ \hline \hat{T} & \underbrace{}_{\frac{1}{2}} & \end{array} = \bar{T} = (\text{id}, T) \# \mu$

<u>③</u>	$x_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$	$\mu = 2 \delta_{x_1} + 4 \delta_{x_2}$
	$x_2 = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$	
	$y_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$	
	$y_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$	$v = 3 \delta_{y_1} + 1 \delta_{y_2} + 2 \delta_{y_3}$
	$y_3 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$	
		There does not exist a transport map!
	$\begin{array}{c cc} & y_1 & y_2 & y_3 \\ \hline x_1 & 4 & 1 & 5 \\ x_2 & 2 & 1 & 1 \\ \hline C & \underbrace{}_{\frac{1}{2}} & & \end{array}$	$\begin{array}{c ccc} & 3 & 1 & 2 \\ \hline 2 & 0 & 0 & 2 \\ 4 & 3 & 1 & 0 \\ \hline \hat{T} & \underbrace{}_{\frac{1}{2}} & & \end{array}$
		optimal transport plan

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Brenier's Theorem

Theorem Let $\mu, \nu \in \mathcal{P}_2(X)$, where $\mu \ll \lambda$ and $c(x, y) = \|x - y\|^2$. Then

- ◆ Kantorovich problem has a unique solution $\hat{\pi}$
- ◆ $\hat{\pi} = (I, \hat{T})_{\#}\mu$, where $\hat{T} \in L^2_{\mu}(X, X)$ is the optimal transport map
- ◆ If ν has bounded support, then

$$\hat{T}(x) = x - \nabla \varphi(x) = \nabla \psi(x) \quad \text{for } \mu - \text{a.e. } x,$$

for some lower semi-continuous, convex, differentiable μ -a.e function ψ .

- ◆ Conversely, if ψ is lower semi-continuous, convex and differentiable μ -a.e. with $|\nabla \psi| \in L^2_{\mu}(X, X)$, then

$$T = \nabla \psi$$

is the optimal transport map from μ to $\nu = T_{\#}\mu \in \mathcal{P}_2(X)$.

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Wasserstein Spaces

For $p \in [1, \infty)$, the space of measures with finite p -th moment is defined by

$$\mathcal{P}_p(\mathcal{X}) := \left\{ \mu \in \mathcal{P}(\mathbb{R}^d) : \int_{\mathbb{R}^d} \|x\|^p d\mu(x) < \infty \right\}.$$

For $\mu, \nu \in \mathcal{P}_p(\mathbb{R}^d)$, the **Wasserstein p -distance** is given by

$$W_p^p(\mu, \nu) := \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^p d\pi(x, y),$$

which actually defines a **metric**.

The metric space $(\mathcal{P}_p(\mathcal{X}), W_p)$ is called the p -th **Wasserstein space**.

Convergence of $(\mu_n)_n$, $\mu_n \in \mathbb{P}_2(\mathbb{R}^d)$: $W_2(\mu_n, \mu) \rightarrow 0$ as $n \rightarrow \infty$ if and only if we have weak/narrow convergence

$$\mu_n \xrightarrow{*} \mu \text{ in } C_b(X)'$$

and

$$\int_{\mathbb{R}^d} \|x\|^2 d\mu_n(x) \rightarrow \int_{\mathbb{R}^d} \|x\|^2 d\mu(x) \text{ as } n \rightarrow \infty$$

Example: $\mu_1 \sim \mathcal{N}(m_1, \Sigma_1)$, $\mu_2 \sim \mathcal{N}(m_2, \Sigma_2)$ and $\Sigma_{1,2} := \left(\Sigma_2^{\frac{1}{2}} \Sigma_1 \Sigma_2^{\frac{1}{2}} \right)^{\frac{1}{2}}$

Then

$$W_2^2(\mu_1, \mu_2) = \|m_1 - m_2\|^2 + \text{tr}(\Sigma_1 - 2\Sigma_{1,2} + \Sigma_2)$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Outline

1. Basic Notation: Spaces of Measures
2. Monge and Kantorovich Problem of OT
3. Discrete OT and its Dual
4. Regularized Optimal Transport
5. Mirror Descent Algorithm - SMART - Sinkhorn Algorithm

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Discrete OT and its Dual

<u>Discrete OT:</u>	$\begin{array}{c cccc} & v_1 & \dots & v_n \\ \mu_1 & T_{11} & \dots & T_{1n} \\ \vdots & \vdots & & \vdots \\ \mu_m & T_{m1} & \dots & T_{mn} \end{array}$	<u>somewhat double notation</u>
<u>Matrixform</u>	$\min_{\pi \in \mathbb{R}^{m,n}} \langle C, \pi \rangle$	$\text{st } \pi^T 1_n = \underline{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_m \end{pmatrix}$
		$\pi^T 1_m = \underline{v}$
		$\pi \geq 0$
<u>Vectorform</u> \downarrow	$\min_{\pi \in \mathbb{R}^N} \langle C, \pi \rangle$	$\text{st } (\underbrace{I_m \dots I_m}_m) \pi = \underline{\mu}$
columnwise reshaped $N = m \cdot n$		$\begin{pmatrix} \underbrace{1 \dots 1}_m & & \\ & \ddots & \\ & & \underbrace{1 \dots 1}_m \end{pmatrix} \pi = \underline{v}$
		$\pi \geq 0$
<u>Show:</u>	$\min_{\pi \in \mathbb{R}^N} \langle C, \pi \rangle$	$\text{st } \underbrace{P_1}_{\pi^T} = (1_n \otimes I_m) \pi = \underline{\mu}$
		$P_2 = (I_m \otimes 1_m) \pi = \underline{v}$
		$\pi \geq 0$
<u>Linear optimization problem</u> , dimension $m \cdot m$		
(P)	$\min_{V \geq 0} \max_{\psi, \underline{\psi}} \langle C, \pi \rangle + \langle \underline{\mu} - P_1 \pi, \underline{\psi} \rangle + \langle \underline{v} - P_2 \pi, \underline{\psi} \rangle$	
(D)	$\max_{\psi, \underline{\psi}} \min_{\pi \geq 0} \langle C, \pi \rangle - \langle P_1 \pi, \underline{\psi} \rangle - \langle P_2 \pi, \underline{\psi} \rangle + \langle \underline{\mu}, \underline{\psi} \rangle + \langle \underline{v}, \underline{\psi} \rangle$	
	$= \max_{\psi, \underline{\psi}} \min_{\pi \geq 0} \underbrace{\langle C - P_1^T \underline{\psi} - P_2^T \underline{\psi}, \pi \rangle}_{(C - P_1^T \underline{\psi} - P_2^T \underline{\psi})_{ij} < 0 \Rightarrow -\infty \text{ if for some } i, j} + \langle \underline{\mu}, \underline{\psi} \rangle + \langle \underline{v}, \underline{\psi} \rangle$	
	$= \max_{\psi, \underline{\psi}} \min_{\pi \geq 0} \underbrace{(C - P_1^T \underline{\psi} - P_2^T \underline{\psi})_{ij} \geq 0}_{\text{for all } i, j} \Rightarrow 0 \text{ and } \pi_{ij} = 0 \text{ if } C_{ij} \geq \psi_i - \underline{\psi}_j$	
	$= \max_{\psi, \underline{\psi}} \langle \underline{\mu}, \underline{\psi} \rangle + \langle \underline{v}, \underline{\psi} \rangle$	<u>dual problem</u> , dimension $m+m$
		$\Rightarrow \text{support cone}, \text{ fm } \pi$

OT and its Dual

Primal OT:

$$\text{OT}(\mu, \nu) := \min_{\pi \in \Pi(\mu, \nu)} \int_{X^2} c \, d\pi,$$

Dual OT:

$$\text{OT}(\mu, \nu) = \max_{\substack{(\varphi, \psi) \in C(X)^2 \\ \varphi(x) + \psi(y) \leq c(x, y)}} \int_X \varphi \, d\mu + \int_X \psi \, d\nu.$$

$(\hat{\varphi}, \hat{\psi}) = (\hat{\varphi}, \hat{\varphi}^c)$ with *c-transformed function*

$$\varphi^c(y) = \min_{x \in X} \{c(x, y) - \varphi(x)\}.$$

Example: Set $X := [0, 1]$, $c(x, y) = |x - y|$, $\mu = \delta_0/2 + \delta_1/2$, $\nu = \delta_{0.1}/2 + \delta_{0.9}/2$. Then, $\text{OT}(\mu, \nu) = 0.1$ with unique optimal transport plan $\hat{\pi} = \frac{1}{2}\delta_{0,0.1} + \frac{1}{2}\delta_{1,0.9}$.

Optimal dual potentials

$$\hat{\varphi}_1(x) = \begin{cases} 0.1 - x & \text{for } x \in [0, 0.1], \\ x - 0.9 & \text{for } x \in [0.9, 1], \\ 0 & \text{else,} \end{cases} \quad \text{and} \quad \hat{\varphi}_2(x) = \begin{cases} 0.2 - x & \text{for } x \in [0, 0.2], \\ x - 0.9 & \text{for } x \in [0.9, 1], \\ 0 & \text{else.} \end{cases}$$

Clearly, these potentials do not differ only by a constant.

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Regularized Optimal Transport

Primal Problem:

$$\text{OT}_\varepsilon = \min_{\pi \in \Pi(\mu, \nu)} \int_{X \times X} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi, \mu \otimes \nu)$$

Discrete setting:

$$\begin{aligned} \text{OT}_\varepsilon &= \min_{\pi \in \Pi(\mu, \nu)} \langle c, \pi \rangle + \varepsilon \left(\sum_{i,j} \pi_{i,j} \log \pi_{i,j} - \pi_{i,j} \log(\mu_i \nu_j) - \pi_{i,j} \right) \\ &= \varepsilon \min_{\pi \in \Pi(\mu, \nu)} \sum_{i,j} \pi_{i,j} \log \pi_{i,j} - \pi_{i,j} \log(\mu_i \nu_j e^{-c/\varepsilon}) - \pi_{i,j} \\ &= \varepsilon \min_{\pi \in \Pi(\mu, \nu)} \text{KL}(\pi, \underbrace{\text{diag}(\mu) e^{-c/\varepsilon} \text{diag}(\nu)}_K) \end{aligned}$$

Lagrangian:

$$\min_{\pi} \max_{\varphi, \psi} \text{KL}(\pi, K) + \langle P_1 \pi - \mu, \varphi \rangle + \langle P_2 \pi - \nu, \psi \rangle$$

Dual Problem:

$$\min_{\phi, \psi} \langle K, e^{-P_1^\top \varphi - P_2^\top \psi} \rangle + \langle \mu, \varphi \rangle + \langle \nu, \psi \rangle$$

1 2

3 4

5 6

7 8

9 10

11 12

13 14

15 16

17 18

19 20

21 22

23 24

25 26

Regularized Optimal Transport

Alternating minimization: **Sinkhorn algorithm** with $u := e^{-\varphi}$, $v := e^{-\psi}$

$$\nabla_\varphi = 0 : \text{diag}(e^{-\varphi}) K \text{diag}(e^{-\psi^{(r-1)}}) \mathbf{1} = \mu \quad \rightarrow \quad u^{(r)} = \frac{\mu}{K v^{(r-1)}}$$

$$\nabla_\psi = 0 : (\text{diag}(e^{-\varphi^{(r)}}) K \text{diag}(e^{-\psi}))^\top \mathbf{1} = \nu \quad \rightarrow \quad v^{(r)} = \frac{\nu}{K^\top u^{(r)}}$$

Relation to π : (set gradient of Lagrangian wrt π to 0)

$$0 = \log \pi - \log K + P_1^\top \varphi + P_2^\top \psi \quad \iff \quad \pi = \text{diag}(u) K \text{diag}(v)$$

We can also consider the Sinkhorn algorithm wrt π :

$$\pi^{(0)} := K,$$

for $r = 0, 1, \dots$

$$\pi^{(2r+1)} := \text{diag} \left(\frac{\mu}{\pi^{(2r)} \mathbf{1}} \right) \pi^{(2r)}$$

$$\pi^{(2r+2)} := \pi^{(2r+1)} \text{diag} \left(\frac{\nu}{(\pi^{(2r+1)})^\top \mathbf{1}} \right)$$

1 2

3 4

5 6

7 8

9 10

11 12

13 14

15 16

17 18

19 20

21 22

23 24

25 26

Convergence of Sinkhorn Algorithm

Linear convergence: in $\mathbb{R}_{>0,*}^d$, $x \sim y$ if $x = \alpha y$ for some $\alpha > 0$

$$d_H(u^{(r+1)}, u^*) \leq \lambda(K)^2 d_H(u^{(r)}, u^*)$$

with

$$\lambda(K) := \sup \left\{ \frac{d_H(Kx, Ky)}{d_H(x, y)} : x \not\sim y \right\}$$

and the Hilbert norm

$$d_H(x, y) := \|\log x - \log y\|_V,$$

$$\|x\|_V := \max_i x_i - \min_i x_i$$

1 2

3 4

5 6

7 8

9 10

11 12

13 14

15 16

17 18

19 20

21 22

23 24

25 26

Is there a relation to proximal algorithms?

- ◆ Mirror Descent Algorithm = linearized proximal algorithm with a general f -divergence instead of the special squared norm
- ◆ SMART
- ◆ BI-SMART

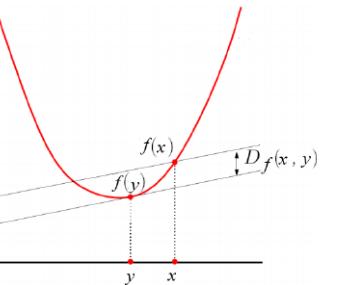
1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Bregman Distances

Let $f : \mathbb{R}^d \rightarrow R \cup \{+\infty\}$ convex lsc and $\text{dom } f \cap (0, +\infty) \neq \emptyset$ with nonempty subdifferential ∂f

Bregman distance:

$$D_f(x, y) = \{f(x) - f(y) + \langle p, x - y \rangle : p \in \partial f(y)\}$$



Examples:

$$1. \quad f(x) = \frac{1}{2}\|x\|^2$$

$$D_f(x, y) = \frac{1}{2}\|x\|^2 - \frac{1}{2}\|y\|^2 - \langle y, x - y \rangle = \frac{1}{2}\|x - y\|^2$$

$$2. \quad f(x) = x \log x$$

$$D_f(x, y) = x \log x - y \log y - \langle \log y + 1, x - y \rangle = x \log x - x \log y - x + y = \text{KL}(x, y)$$

Function Name	$f(x)$	$\text{dom } f$	$D_f(x, y)$
Squared Norm	$\frac{1}{2}x^2$	$(-\infty, \infty)$	$\frac{1}{2}(x - y)^2$
Shannon Entropy	$x \log x$	$[0, \infty)$	$x \log \frac{x}{y} - x + y$
Bit Entropy	$x \log x + (1 - x) \log(1 - x)$	$[0, 1]$	$x \log \frac{x}{y} + (1 - x) \log \frac{1-x}{1-y}$
Burg Entropy	$-\log x$	$(0, \infty)$	$\frac{x}{y} - \log \frac{x}{y} - 1$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Properties of Bregman Distances

Properties:

- ◆ $D_f(x, y) \geq 0$
- ◆ $D_f(x, y) = 0$ iff $x = y$ in case f is strictly convex.
- ◆ In general not symmetric and does not fulfill a triangular inequality
- ◆ Jointly convex, lsc.
- ◆ If f is strictly convex, D_f is strictly convex in the first argument.
- ◆ If expressions exist

$$\nabla_x D_f(x, y) = \nabla f(x) - \nabla f(y)$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Is there a relation to proximal algorithms?

- ◆ Mirror Descent Algorithm = linearized proximal algorithm with a general f -divergence instead of the special squared norm
- ◆ SMART (simultaneous multiplicative algebraic reconstruction technique)
- ◆ BI-SMART (block iterative SMART)

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Blackboard

GENERAL TASK: $A^T 1 = 1$, $A \geq 0$ stochastic matrix, $b \geq 0$

(1) $\min_x KL(x|y)$ st. $Ax=b$

$$P_C^{KL}(y), \quad C := \{x : Ax=b\} \quad \text{KL projection onto affine set}$$

(2) $\min_{x \geq 0} KL(Ax, b)$

derivation of alg: $\nabla_{Ax} KL(Ax, b) = A^T \log \frac{Ax}{b}$

$$0 = A^T \log \frac{Ax}{b}$$

$$1 = e^{-\tau A^T \log \frac{Ax}{b}}$$

$$x = x_0 e^{-\tau A^T \log \frac{Ax_0}{b}} \quad \text{fixed point eq.}$$

$$\boxed{x^{(0)} = x^{(1)} \circ e^{-\tau A^T \log \frac{Ax^{(0)}}{b}}} \quad \text{Picard it.}$$

SMART (Byrne 1996, ...)

relation to (1)

$$x^{(r)} \xrightarrow{\text{def}} x^* \in C_{\geq 0} := \{x \geq 0 : Ax=b\} \neq \emptyset$$

$$x^* = \underset{x \in C_{\geq 0}}{\operatorname{argmin}} KL(x, x^{(r)})$$

\uparrow

(3) $\min_{x \in C} f(x)$

$x \in C \leftarrow$ convex, closed

- gradient descent alg.: ⊕ not Lipschitz and gradient of f
- $x^{(r+1)} = \underset{x \in C}{\operatorname{argmin}} \frac{1}{\tau} KL(x, x^{(r)}) + f(x)$ (BPG)
- ⊕ $0 = \nabla f(x) - \tau \nabla f(x^{(r)})$ (more general)

Generalization of f :

$$x^{(r+1)} = \underset{x \in C}{\operatorname{argmin}} \frac{1}{\tau} KL(x, x^{(r)}) + f(x^{(r)}) + \langle \nabla f(x^{(r)}), x - x^{(r)} \rangle$$

$$= \underset{x \in C}{\operatorname{argmin}} KL(x, x^{(r)}) + \mathbb{E}_x \langle \nabla f(x^{(r)}), x \rangle$$

Mirror descent alg. (**MDA**) (Nesterov/Yudin 1983; Beck/Teboulle 2003)

$$x^{(r+1)} = \underset{x \in C}{\operatorname{argmin}} D_\alpha(x, x^{(r)}) + \mathbb{E}_x \langle \nabla f(x^{(r)}), x \rangle \quad \text{or equivalent:}$$

$$= \underset{x \in C}{\operatorname{argmin}} h(x) + \langle \mathbb{E}_x \nabla f(x^{(r)}) - \nabla h(x^{(r)}), x \rangle$$

$D_\alpha = KL$: $x^{(r+1)} = x^{(r)} \circ e^{-\tau \nabla f(x^{(r)})}$

$f(x) = KL(Ax, b)$: $x^{(r+1)} = x^{(r)} \circ e^{-\tau A^T \log \frac{Ax^{(r)}}{b}}$

Convergence: - general MDA $O(\frac{1}{\tau^2})$
 - accelerated BPG: $O(\frac{1}{\tau^2}) \quad \tau \in [2, \infty)$ for special h
 (biangular scaling prop. of D_α)

• SMART for OT_E : $A = \begin{pmatrix} p_1 \\ p_2 \end{pmatrix}$ recall by $\frac{1}{2}$ to get $A^T 1 = 1$

$$T^{(r+1)} = \operatorname{diag}\left(\frac{1}{(\pi^{(r)})^T 1}\right)^{1/2} \pi^{(r)} \operatorname{diag}\left(\frac{\mu}{\pi^{(r)} 1}\right)^{1/2}$$

- this alg is $\frac{1}{2}$ slower than Sinkhorn alg.
- it works also with $\alpha < 1$ (not $\alpha=1$!) instead of $\alpha=\frac{1}{2}$;
 from $\alpha=1/2$ slower than Sinkhorn alg. Question 1: Why?

GENERALIZATION

$C = C_0 \cap C_1 \cap \dots \cap C_N$, $C_k := \{x : A_k x = b_k\}, \quad A_k x_k \geq 0, \quad A^T 1 = 1$

(4) $\underset{x \in C}{\operatorname{argmin}} KL(x|y)$

Iterative uniformization projection: $x^{(0)} = y$ $\xrightarrow{\text{see (1) and (2)}}$
 $x^{(r+1)} = P_{C_r}^{KL}(x^{(r)})$ $\xrightarrow{\text{C} \rightarrow \mathbb{R}^N = C_r \text{ (periodic projection)}}$
 (Caesar 1995)

- convergence for affine sets; not for general convex sets (dystopian alg. for convex sets...)

Generalized iterative scaling (**GIS**): $x^{(0)} = y$
 $B1-SMART: \quad x^{(r+1)} = x^{(r)} \circ e^{-\tau A^T \log \frac{A_r x^{(r)}}{b_r}}$ one SMART step

- B1-SMART for OT_E : $\pi^{(r+1)} = \pi^{(r)} \circ \operatorname{diag} \frac{\mu}{\pi^{(r)} 1}$
- Sinkhorn alg.: $\pi^{(r+2)} = \operatorname{diag} \frac{\nu}{(\pi^{(r+1)})^T 1} \quad T^{(r+1)}$

Question 2: B1-SMART for unbalanced OT: $(1-\lambda) KL(x|y) + \lambda KL(Ax, b)$
 does not converge in the obvious form
 Alternatives?

B1-SMART in case $C = \emptyset$. (old problem, e.g. of Byrne)

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Comparison

SMART solely performs the multiplicative update specified in (2) with its step-size fixed to $\tau_k = \frac{1}{L}$, where again f is L-smooth relative to φ .

FSMART is based on the iteration suggested in [2], where initially $\theta_0 = 1$ is chosen, which is then subsequently updated via $\theta_{k+1} = \frac{\sqrt{\theta_k^4 + 4\theta_k^2} - \theta_k^2}{2}$, as suggested in [5].

ABPG-e as described in [6, Algorithm 2] was applied to (1) with parameters $\gamma_{\min} = 1$, $\gamma_0 = 5$ and $\delta = 0.05$. The choices for γ_0 and δ deviate slightly from the recommendation in [6], but were chosen to facilitate fastest possible convergence on the selected problem instances. To ensure comparability restarting mechanisms and stopping criteria based on the divergence of iterates were foregone. Updates for θ were conducted via Newton's method.

ABPG-g specified in [6, Algorithm 3] to (1) is used with parameters: $\rho = 1.2$, $\gamma = 2$ and $G_{\min} = 10^{-3}$. Restarting, stopping criteria and updating θ was handled analogously to ABPG-e.

RG is a SMART iteration with Armijo line search for choosing the step size τ_k via the retraction in (29) to iterate according to (37). The line search parameters are $\sigma = 0.5$, $\beta = 0.8$, $\alpha = 5.0$.

PD is the Chambolle-Pock primal dual algorithm [17, Algorithm 1] for solving convex composite structured optimization problems of the form $f(x) = g(x) + h(Ax)$. For $h(y) := \text{KL}(y, b)$ with $y = Ax$ and $g \equiv 0$ we obtain

$$x^{k+1} = x^k e^{-\tau A^\top y^k} \quad (\text{primal-step}) \quad (39)$$

$$y^{k+1} = \log \left(\frac{e^{y^k} + \sigma A(2x^{k+1} - x^k)}{1 + \sigma b} \right), \quad (\text{dual-step}) \quad (40)$$

Chambolle-Pock 2015:

$$x^{(r+1)} = \operatorname{argmin}_x \left\{ \langle Ax, y^n \rangle + g(x) + \frac{1}{\tau} D_\varphi(x, x^{(r)}) \right\}$$

$$y^{(r+1)} = \operatorname{argmin}_y \left\{ h^*(y) - \langle A(2x^{(r+1)} - x^{(r)}), y \rangle + \frac{1}{\sigma} D_{\varphi^*}(y, y^{(r)}) \right\}$$

with $\varphi(x) = \langle x, \log x \rangle - \langle 1, x \rangle$ and $\varphi^*(y) = \langle 1, e^y \rangle$

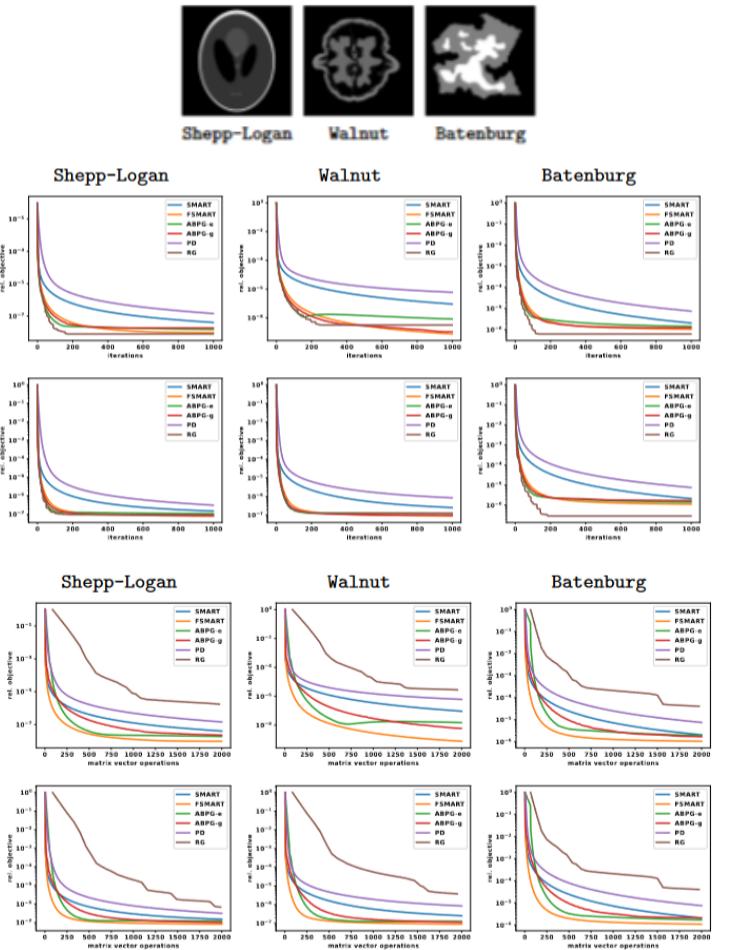
1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

Comparison

CT: parallel beam geometry and equidistant angles in the range $[0, \pi]$.

Undersampling rate was chosen to be 20

Setting without noise and with Poisson noise of SNR = 20 db.



1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26

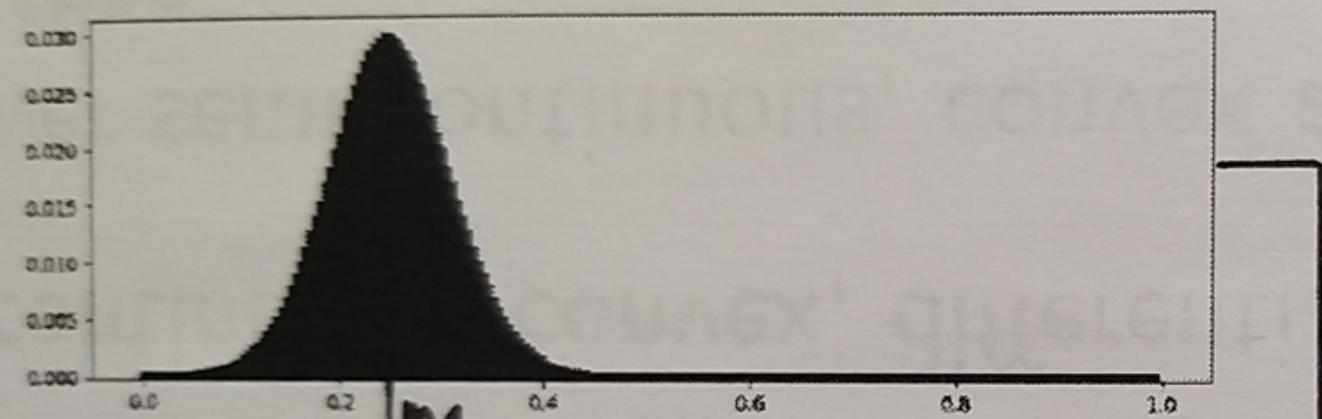
Berlin Mathematics Research Center



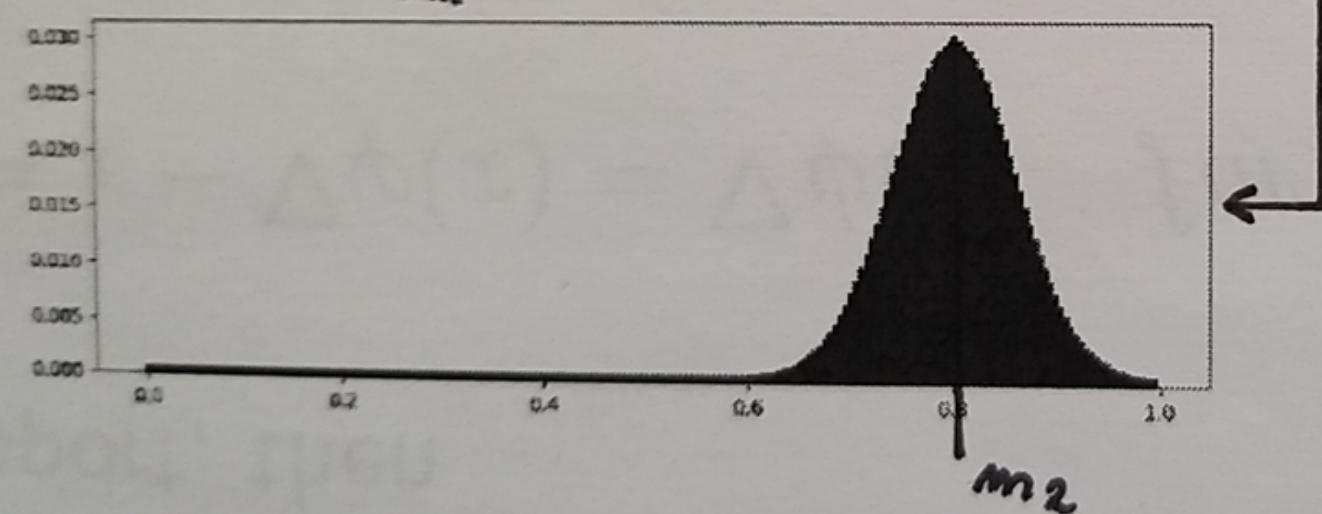
Funded under Germany's Excellence Strategy by



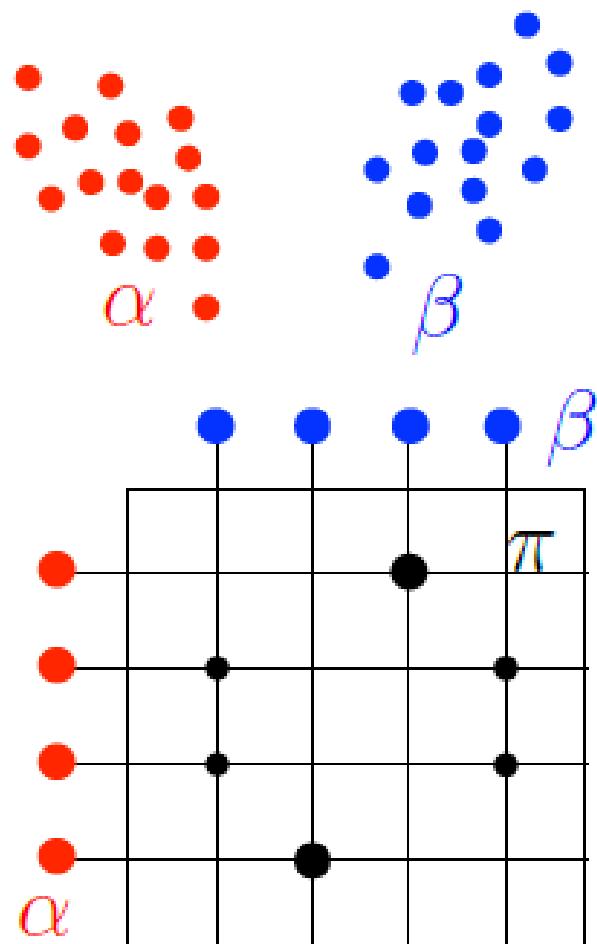
Deutsche
Forschungsgemeinschaft



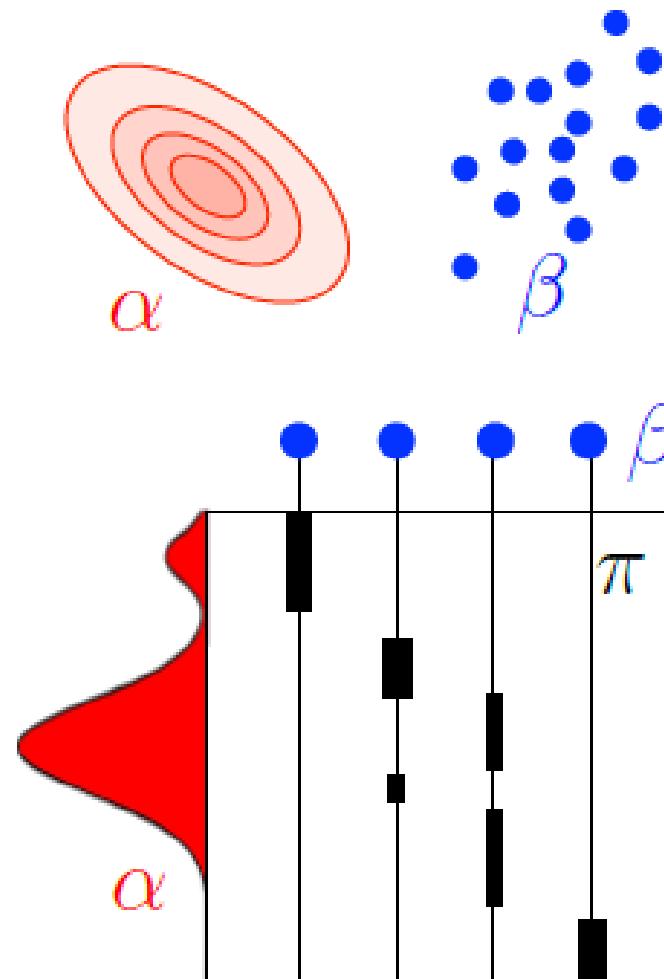
$$\tau(x) = x - m_1 + m_2$$



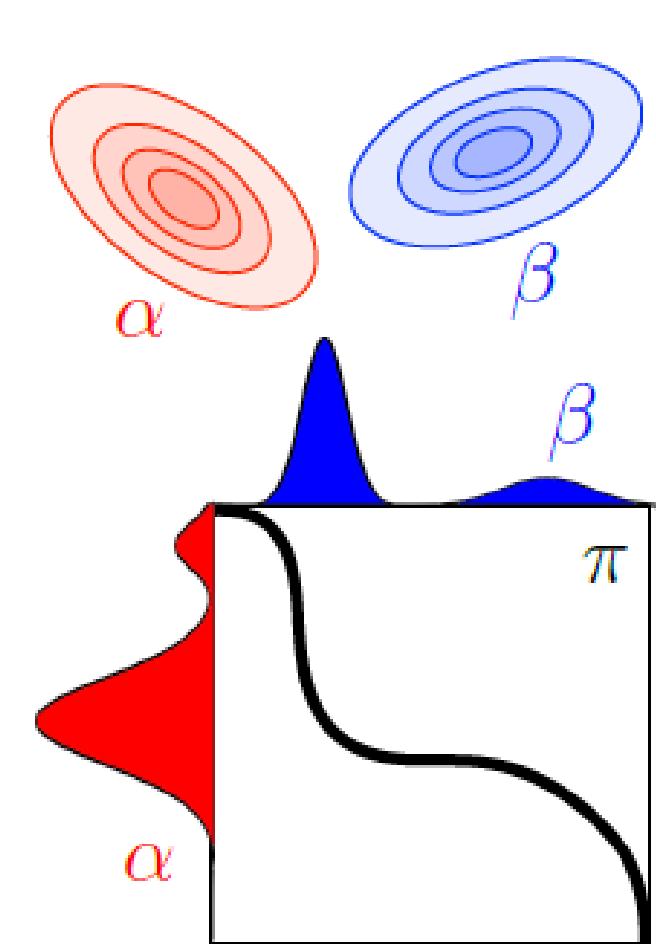




Discrete



Semidiscrete



Continuous



$$\text{EXAMPLE } ① \quad x_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

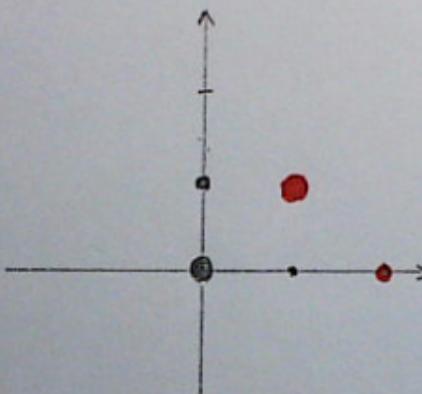
$$x_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$x_3 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$y_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$y_2 = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$$

$$v = 4d_{y_1} + 2d_{y_2}$$



The only possible transport map is

$$T(x_1) = y_1, \quad T(x_2) = y_1, \quad T(x_3) = y_2$$

The optimal transport plan is

	y_1	y_2
x_1	2	4
x_2	1	1
x_3	1	5

C

	4	2
3	2	1
1	0	1
2	2	0

$\hat{\pi}$

optimal plan

(does not correspond to

	4	2
3	3	0
1	1	0
2	0	2

$$\pi = (\text{id}, T) \# \mu$$

② $D = 2\partial_{y_1} + 4\partial_{y_2}$; not the same

Map: $T(x_1) = y_2, T(x_2) = y_2, T(x_3) = y_1$

	y_1	y_2
x_1	4	2
x_2	1	1
x_3	5	1

$\underbrace{}_c$

$$\begin{array}{c|cc}
& 2 & 4 \\
\hline
3 & 0 & 3 \\
1 & 0 & 1 \\
2 & \underbrace{2 & 0}_{\hat{\pi}}
\end{array}$$

$$= \bar{\pi} = (\text{id}, T) \# \mu$$

$$\begin{array}{c|cc}
& 2 & 4 \\
\hline
3 & 0 & 3 \\
1 & 0 & 1 \\
2 & \underbrace{2 & 0}_{\hat{\pi}}
\end{array}$$

(3)

$$x_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$x_2 = \begin{pmatrix} 2 \\ 0 \end{pmatrix}$$

$$\mu = 2 \delta_{x_1} + 4 \delta_{x_2}$$

$$y_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$y_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$\nu = 3 \delta_{y_1} + 1 \delta_{y_2} + 2 \delta_{y_3}$$

$$y_3 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

There does not exist a transport map !

	y_1	y_2	y_3
x_1	4	1	5
x_2	2	1	1
C			

	3	1	2
2	0	0	2
4	3	1	0

 $\hat{\pi}$

optimal transport plan

Discrete OT:

	v_1	\dots	v_n
M_1	T_{11}	\dots	T_{1n}
\vdots	\vdots	\vdots	\vdots
M_m	T_{m1}	\dots	T_{mn}

somewhat double
notation

Matrixform

$$\min_{\pi \in \mathbb{R}^{m,n}} \langle c, \pi \rangle \quad \text{s.t.} \quad \pi^T 1_n = \underline{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_m \end{pmatrix}$$

$$\pi^T 1_m = \underline{\nu}$$

$$\pi \geq 0$$

Vectorform

$$\min_{\pi \in \mathbb{R}^N} \langle c, \pi \rangle \quad \text{s.t.} \quad \underbrace{(I_m \ \dots \ I_m)}_m \pi = \underline{\mu}$$

columnwise
reshaped

$N = m \cdot n$

$$\begin{pmatrix} \underbrace{1 \dots 1}_m & & \\ & \ddots & \\ & & \underbrace{1 \dots 1}_m \end{pmatrix} \pi = \underline{\nu}$$

$$\pi \geq 0$$

Shorter:

$$\min_{\pi \in \mathbb{R}^N} \langle c, \pi \rangle \quad \text{s.t.} \quad \begin{array}{l} P_1 \pi = (1_n \otimes I_m) \pi = \underline{\mu} \\ P_2 \pi = (I_m \otimes 1_m) \pi = \underline{\nu} \\ \pi \geq 0 \end{array}$$

Linear optimization problem, dimension $n \cdot m$

(P) $\min_{\pi \geq 0} \max_{\underline{\nu}, \underline{\mu}} \langle c, \pi \rangle + \langle \underline{\mu} - P_1 \pi, \underline{\nu} \rangle + \langle \underline{\nu} - P_2 \pi, \underline{\mu} \rangle$

(D) $\max_{\underline{\nu}, \underline{\mu}} \min_{\pi \geq 0} \langle c, \pi \rangle - \langle P_1 \pi, \underline{\nu} \rangle - \langle P_2 \pi, \underline{\mu} \rangle + \langle \underline{\mu}, \underline{\nu} \rangle + \langle \underline{\nu}, \underline{\mu} \rangle$

$$= \max_{\underline{\nu}, \underline{\mu}} \min_{\pi \geq 0} \underbrace{\langle c - P_1^T \underline{\nu} - P_2^T \underline{\mu}, \pi \rangle}_{(c - P_1^T \underline{\nu} - P_2^T \underline{\mu})_{ij} < 0 \Rightarrow -\infty \text{ if } \text{for some } i, j} + \langle \underline{\mu}, \underline{\nu} \rangle + \langle \underline{\nu}, \underline{\mu} \rangle$$

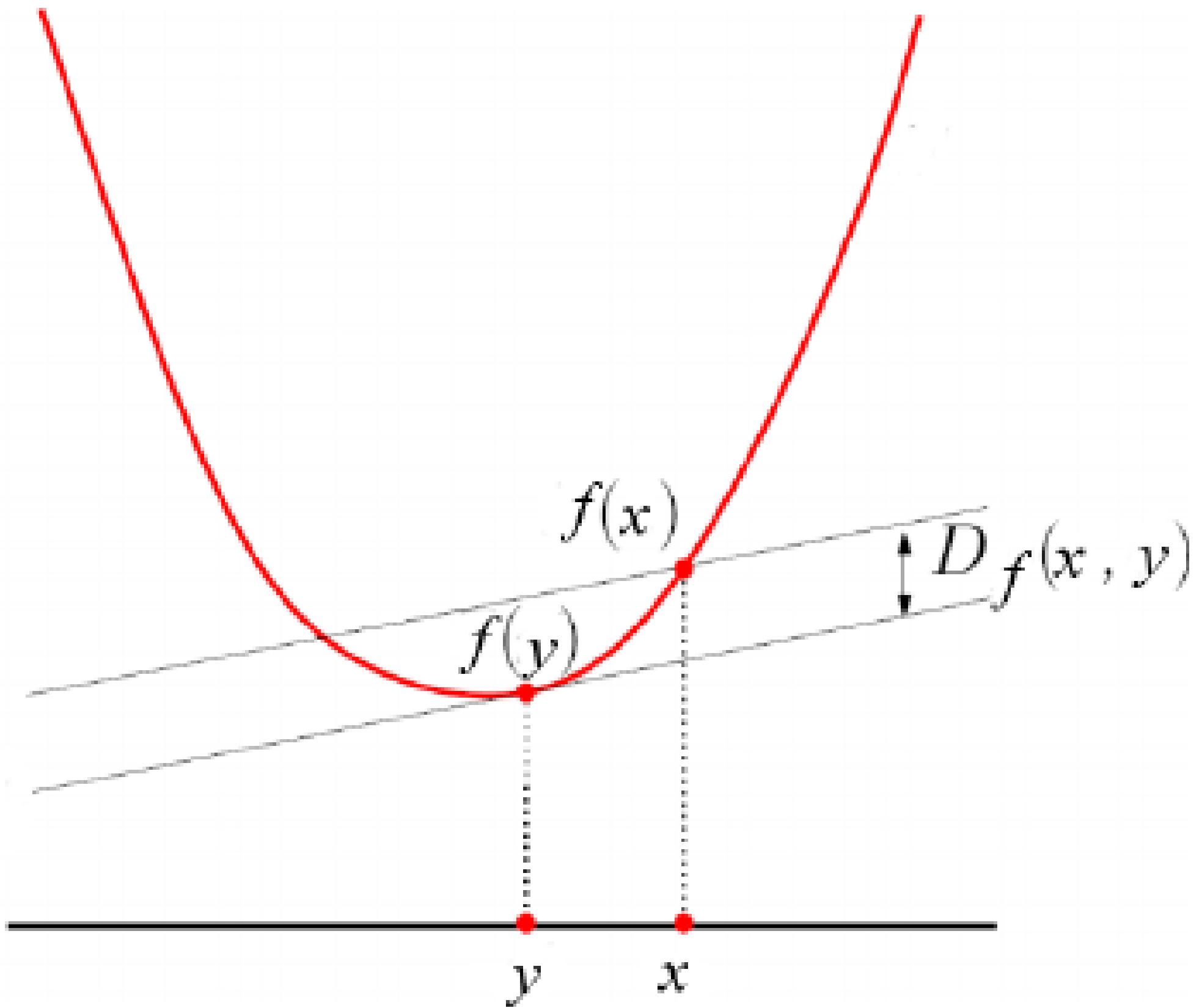
$$= \max_{\underline{\nu}, \underline{\mu}} \min_{\pi \geq 0} \underbrace{\langle c - P_1^T \underline{\nu} - P_2^T \underline{\mu}, \pi \rangle}_{(c - P_1^T \underline{\nu} - P_2^T \underline{\mu})_{ij} \geq 0 \Rightarrow 0 \text{ and } \pi_{ij} = 0 \text{ if } c_{ij} > \nu_i - \mu_j} + \langle \underline{\mu}, \underline{\nu} \rangle + \langle \underline{\nu}, \underline{\mu} \rangle$$

$$= \max_{\underline{\nu}, \underline{\mu}} \min_{\pi \geq 0} \underbrace{\langle c - P_1^T \underline{\nu} - P_2^T \underline{\mu}, \pi \rangle}_{\nu_i + \mu_j \leq c_{ij}} + \langle \underline{\mu}, \underline{\nu} \rangle + \langle \underline{\nu}, \underline{\mu} \rangle$$

dual problem,

dimension $m+n$

\Rightarrow support cond. for π



Function Name	$f(x)$	$\text{dom } f$	$D_f(x, y)$
Squared Norm	$\frac{1}{2}x^2$	$(-\infty, \infty)$	$\frac{1}{2}(x - y)^2$
Shannon Entropy	$x \log x$	$[0, \infty)$	$x \log \frac{x}{y} - x + y$
Bit Entropy	$x \log x + (1 - x) \log(1 - x)$	$[0, 1]$	$x \log \frac{x}{y} + (1 - x) \log \frac{1-x}{1-y}$
Burg Entropy	$-\log x$	$(0, \infty)$	$\frac{x}{y} - \log \frac{x}{y} - 1$

$$\boxed{OT_{\varepsilon}(\mu, \nu) = \min_{\pi} KL(\pi, \kappa) \quad \text{st} \quad \begin{pmatrix} p_1 \\ p_2 \end{pmatrix} \pi = \begin{pmatrix} \mu \\ \nu \end{pmatrix} \\ \begin{pmatrix} \Delta & \Delta \\ - & - \end{pmatrix}}$$

GENERAL TASK: $A^T 1 = 1$, $A \geq 0$ stochastic matrix, $b \geq 0$

$$(1) \min_x KL(x, y) \quad \text{st} \quad Ax = b$$

$P_C^{KL}(y)$, $C := \{x : Ax = b\}$ KL projection onto affine set

$$(2) \min_{x \geq 0} KL(Ax, b)$$

derivation of alg: $\nabla KL(Ax, b) = A^T \log \frac{Ax}{b}$

$$0 = A^T \log \frac{Ax}{b}$$

$$1 = e^{-\tau A^T \log \frac{Ax}{b}}$$

$$x = x^0 e^{-\tau A^T \log \frac{Ax}{b}} \quad \text{fixed point eq.}$$

$$\begin{aligned} x^{(0)} \\ x^{(\tau+1)} = x^{(\tau)} \circ e^{-\tau A^T \log \frac{Ax^{(\tau)}}{b}} \end{aligned}$$

Picard it.

SMART

(Byrne 1996, ...)

relation to (1)

$$x^{(\tau)} \xrightarrow{\infty} x^* \in C_{\geq 0} := \{x \geq 0 : Ax = b, y \neq \emptyset\}$$

$$x^* = \underset{x \in C_{\geq 0}}{\operatorname{argmin}} KL(x, x^{(0)})$$

↑
y

$$(3) \min_x f(x)$$

$x \in C \leftarrow \text{convex, closed}$

- gradient descent typ. \circlearrowleft not Lipschitz cont. gradient of f

$$x^{(\tau+1)} = \underset{x \in C}{\operatorname{argmin}} \frac{1}{\tau} KL(x, x^{(\tau)}) + f(x) \quad (\text{BPP})$$

$$\circlearrowleft 0 = \log x - \log(x^{(\tau)}) + \tau \nabla f(x) \quad \downarrow D_h(x, x^{(\tau)}) \\ (\text{more general})$$

linearization of f :

$$x^{(\tau+1)} = \underset{x \in C}{\operatorname{argmin}} \frac{1}{\tau} KL(x, x^{(\tau)}) + f(x^{(\tau)}) + \langle \nabla f(x^{(\tau)}), x - x^{(\tau)} \rangle$$

$$= \underset{x \in C}{\operatorname{argmin}} KL(x, x^{(\tau)}) + \tau_x \langle \nabla f(x^{(\tau)}), x \rangle$$

Mirror descent alg (**MDA**)

(Nesterov/Yudin 1983, Beck/Teboulle 2003)

$$x^{(\tau+1)} = \underset{x \in C}{\operatorname{argmin}} D_h(x, x^{(\tau)}) + \tau_x \langle \nabla f(x^{(\tau)}), x \rangle \quad \text{or equivalent}$$

$$= \underset{x \in C}{\operatorname{argmin}} h(x) + \langle \tau_x \nabla f(x^{(\tau)}) - \nabla h(x^{(\tau)}), x \rangle$$

$$D_h = KL : \quad x^{(\tau+1)} = x^{(\tau)} \circ e^{-\tau \nabla f(x^{(\tau)})}$$

$$f(x) = KL(Ax, b) : \quad x^{(\tau+1)} = x^{(\tau)} \circ e^{-\tau A^T \log \frac{Ax^{(\tau)}}{b}}$$

Convergence: - general MDA $\mathcal{O}\left(\frac{1}{\tau^2}\right)$

- accelerated BPP: $\mathcal{O}\left(\frac{1}{\tau^2}\right) \quad \tau \in [1, 2] \quad \text{for special } h$
 (triangular scaling prop. of D_h)

• SMART for OTE : $A = \begin{pmatrix} P_1 \\ P_2 \end{pmatrix}$ rescale by $\frac{1}{2}$ to get $A^T 1 = 1$

$$\pi^{(\tau+1)} = \text{diag}\left(\frac{\nu}{(\pi^{(\tau)})^T 1}\right)^{1/2} \pi^{(\tau)} \text{diag}\left(\frac{\mu}{\pi^{(\tau)} 1}\right)^{1/2}$$

- This alg. is $^{1/4}$ slower than Sinkhorn alg.

- it works also with $\alpha < 1$ (not $\alpha=1$!) instead of $\alpha = \frac{1}{2}$;
then it is $^{1/2}$ slower than Sinkhorn alg. Question 1: Why?

GENERALIZATION

$$C = C_1 \cap C_2 \cap \dots \cap C_N, \quad C_k := \{x : A_k x = b_k\}, \quad A_k^T b_k \geq 0, \quad A_k^T 1 = 1$$

$$(4) \quad \underset{x \in C}{\operatorname{argmin}} \text{KL}(x, y)$$

→ see (1) and (2)

Iterative information projection : $x^{(0)} = y$
 (Csiszar 1989) $x^{(\tau+1)} = P_{C_\tau}^{\text{KL}}(x^{(\tau)})$ $C_{t+eN} = C_t$ (periodic projection)

- convergence for affine sets ; not for general convex sets (Dijkstra alg. for convex sets...)

Generalized iterative scaling (GIS) : $x^{(0)} = y$
BI-SMART : $x^{(\tau+1)} = x^{(\tau)} \circ e^{-\tau A_\tau^T \log \frac{A_\tau x^{(\tau)}}{b_\tau}}$ one SMART step

• BI-SMART for OTE : $\pi^{(\tau+1)} = \pi^{(\tau)} \circ \text{diag} \frac{\mu}{\pi^{(\tau)} 1}$

= Sinkhorn alg.

$$\pi^{(\tau+2)} = \text{diag} \frac{\nu}{(\pi^{(\tau+1)})^T 1} \pi^{(\tau+1)}$$

Question 2: BI-SMART for unbalanced OT : $(1-\lambda) \text{KL}(x, y) + \lambda \text{KL}(Ax, b)$
 does not converge in the obvious form

Alternatives ?

BI-SMART in case $C = \emptyset$ (odd problem, e.g. of Byrne)

SMART solely performs the multiplicative update specified in (2) with its step-size fixed to $\tau_k = \frac{1}{L}$, where again f is L-smooth relative to φ .

FSMART is based on the iteration suggested in [2], where initially $\theta_0 = 1$ is chosen, which is then subsequently updated via $\theta_{k+1} = \frac{\sqrt{\theta_k^4 + 4\theta_k^2} - \theta_k^2}{2}$, as suggested in [5].

ABPG-e as described in [6, Algorithm 2] was applied to (1) with parameters $\gamma_{\min} = 1$, $\gamma_0 = 5$ and $\delta = 0.05$. The choices for γ_0 and δ deviate slightly from the recommendation in [6], but were chosen to facilitate fastest possible convergence on the selected problem instances. To ensure comparability restarting mechanisms and stopping criteria based on the divergence of iterates were foregone. Updates for θ were conducted via Newton's method.

ABPG-g specified in [6, Algorithm 3] to (1) is used with parameters: $\rho = 1.2$, $\gamma = 2$ and $G_{\min} = 10^{-3}$. Restarting, stopping criteria and updating θ was handled analogously to ABPG-e.

RG is a SMART iteration with Armijo line search for choosing the step size τ_k via the retraction in (29) to iterate according to (37). The line search parameters are $\sigma = 0.5$, $\beta = 0.8$, $\alpha = 5.0$.

PD is the Chambolle-Pock primal dual algorithm [17, Algorithm 1] for solving convex composite structured optimization problems of the form $f(x) = g(x) + h(Ax)$. For $h(y) := \text{KL}(y, b)$ with $y = Ax$ and $g \equiv 0$ we obtain

$$x^{k+1} = x^k e^{-\tau A^\top y^k} \quad (\text{primal-step}) \quad (39)$$

$$y^{k+1} = \log \left(\frac{e^{y^k} + \sigma A(2x^{k+1} - x^k)}{1 + \sigma b} \right), \quad (\text{dual-step}) \quad (40)$$

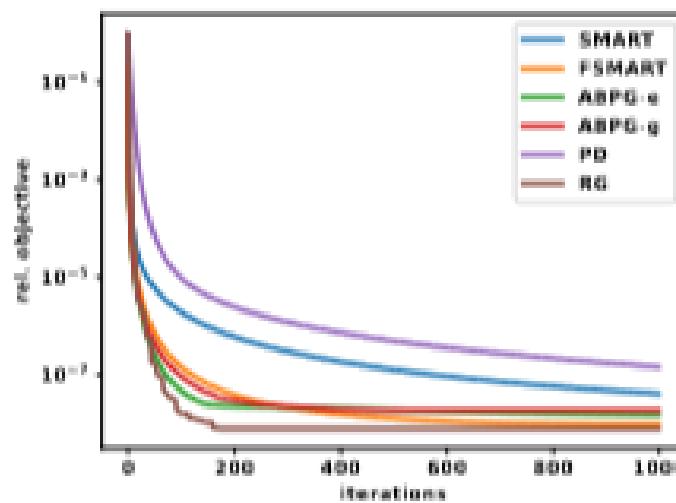


Stepp-Logos

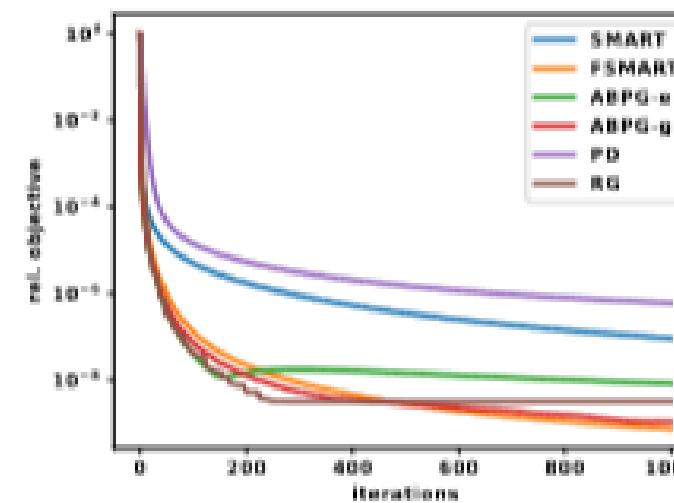
Stepp-Logos

Stepp-Logos

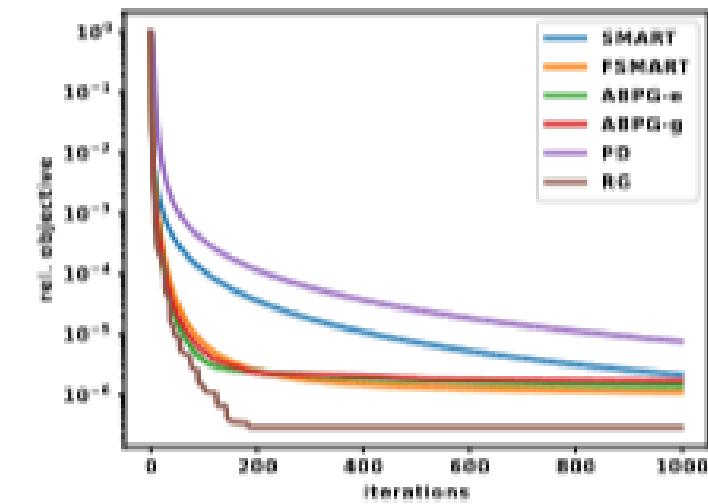
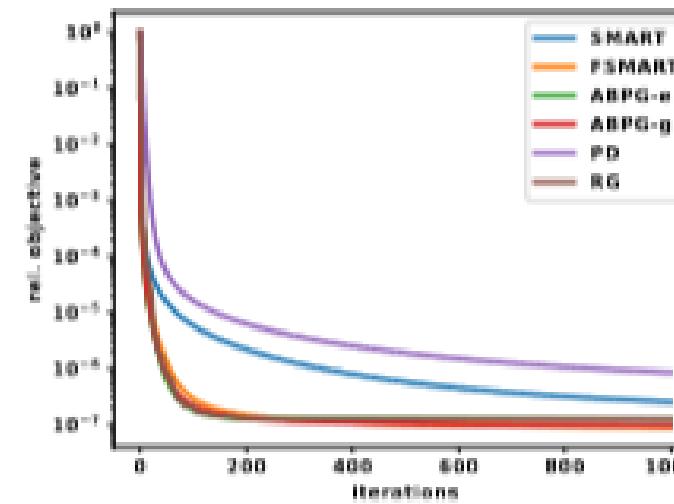
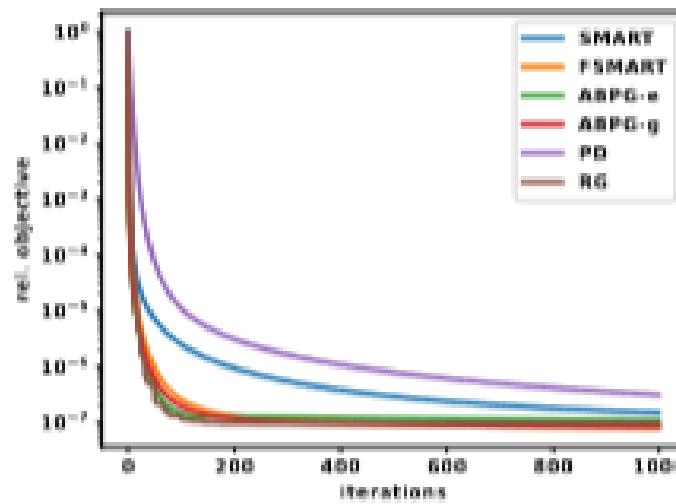
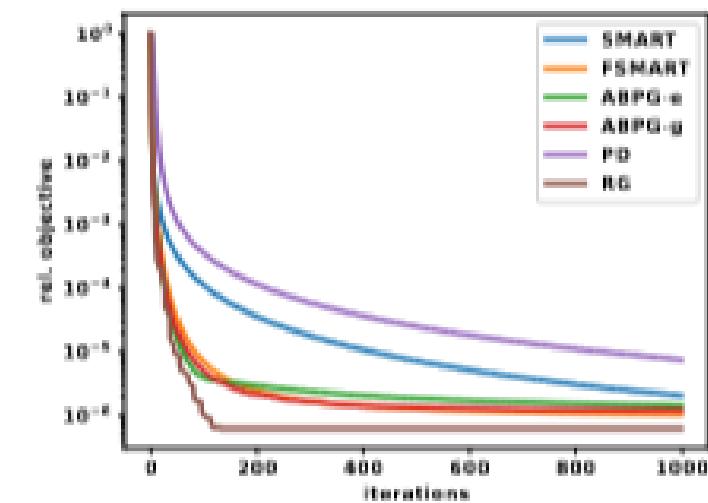
Shepp-Logan



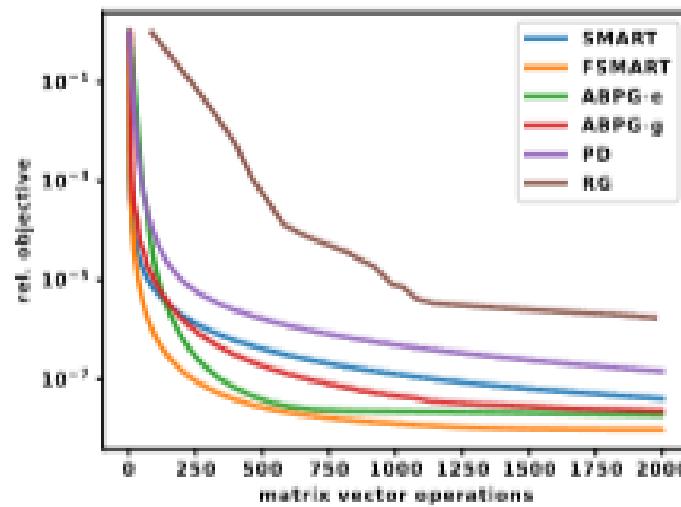
Walnut



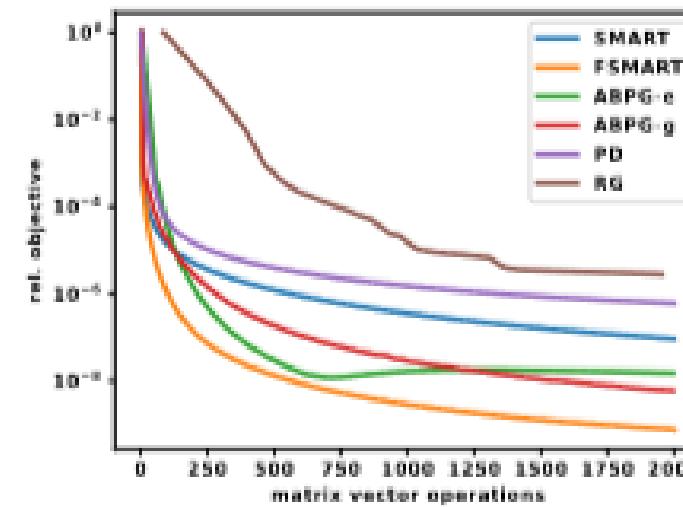
Batenburg



Shepp-Logan



Walnut



Batenburg

