Regularized Optimal Transport

Marco Cuturi



book with Gabriel Peyré https://optimaltransport.github.io/

Joint works with G. Peyré, F. Bach, N. Bonneel, A. Genevay, L. Chizat, A. Rolet, J. Solomon, G. Carlier, JD Benamou, L. Nenna, M. Heitz, and many others.









Kantorovich Problem





Kantorovich Problem



Kantorovich Problem



Kantorovich Problem à la française





Wasserstein on Discrete Measures

Consider
$$\boldsymbol{\mu} = \sum_{i=1}^{n} a_i \delta_{x_i}$$
 and $\boldsymbol{\nu} = \sum_{j=1}^{m} b_j \delta_{y_j}$.
 $M_{\boldsymbol{X}\boldsymbol{Y}} \stackrel{\text{def}}{=} [D(\boldsymbol{x}_i, \boldsymbol{y}_j)^p]_{ij}$
 $U(\boldsymbol{a}, \boldsymbol{b}) \stackrel{\text{def}}{=} \{ \boldsymbol{P} \in \mathbb{R}^{n \times m}_+ | \boldsymbol{P} \boldsymbol{1}_m = \boldsymbol{a}, \boldsymbol{P}^T \boldsymbol{1}_n = \boldsymbol{b} \}$

Def. Optimal Transport Problem $W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$

Solving the OT Problem



Solving the OT Problem



Early application: Earth Mover's





Early application: Earth Mover's





Early application: Earth Mover's



[Rubner'98] dist $(I_1, I_2) = W_1(\mu, \nu)$

Word Mover's Distance



word2vec embedding

Word Mover's Distance



word2vec embedding

[Kusner'15]

 $\operatorname{dist}(D_1, D_2) = W_2(\boldsymbol{\mu}, \boldsymbol{\nu})$

Variational OT Problems in ML

Up to 2010: OT solvers used mostly for retrieval in databases of histograms

 $W_p(\boldsymbol{\mu}, \boldsymbol{\nu}) = ?$ $W_p(\boldsymbol{\mu}, \boldsymbol{\nu}) \leq \cdots ?$

The field has now transitioned to OT as a **loss or fidelity** term

 $\operatorname{argmin}_{\boldsymbol{\mu} \in \mathcal{P}(\Omega)} F(W_p(\boldsymbol{\mu}, \boldsymbol{\nu_1}), W_p(\boldsymbol{\mu}, \boldsymbol{\nu_2}), \dots, \boldsymbol{\mu}) =?$ $``\nabla_{\boldsymbol{\mu}}"W_p(\boldsymbol{\mu}, \boldsymbol{\nu_1}) =?$

Recent spike in interest for [Ambrosio'05]

"Wasserstein + Data" Problems

- Quantization: k-means problem [Lloyd'82] $\min_{\substack{\mu \in \mathcal{P}(\mathbb{R}^d) \\ |\operatorname{supp} \mu| = k}} W_2^2(\mu, \nu_{data})$
- [McCann'95] Interpolant

$$\min_{\boldsymbol{\mu}\in\mathcal{P}(\Omega)}(1-t)W_2^2(\boldsymbol{\mu},\boldsymbol{\nu_1})+tW_2^2(\boldsymbol{\mu},\boldsymbol{\nu_2})$$

• [JKO'98] PDE's as "gradient" flows in $(\mathcal{P}(\Omega), W)$.

$$\mu_{t+1} = \operatorname*{argmin}_{\boldsymbol{\mu} \in \mathcal{P}(\Omega)} J(\boldsymbol{\mu}) + \lambda_t W_p^p(\boldsymbol{\mu}, \mu_t)$$



N $\min_{\boldsymbol{\mu}\in\mathcal{P}(\Omega)}\sum_{i=1}^{\infty}\lambda_i W_p^p(\boldsymbol{\mu},\boldsymbol{\nu_i})$ ${\cal V}_1$ Wasserstein $\mathcal{P}(\Omega)$ Barycenter [Agueh'11] ν_2 $\overline{
u}_3$





Ex: Barycenters for shapes

Graphics: simple testing ground for relevance of Wasserstein geometry





Ex: Barycenters for shapes

Graphics: simple testing ground for relevance of Wasserstein geometry





Ex: Barycenters for shapes



[**PC'18**]

Dataset $\{(x_i, y_i)\}, x_i \in \mathbb{R}^p, y_i \in \mathbb{R}^n_+$



Goal is to find f_{θ} : Images \mapsto Labels

N $\min_{\boldsymbol{\theta}\in\Theta}\sum_{i=1}\mathcal{L}(f_{\boldsymbol{\theta}}(x_i),y_i)$

SNOW

sled

men

 y_i





dog driver winter ice

 $f_{\boldsymbol{\theta}}(x_i)$

husky snow sled slope men





Use for \mathcal{L} a Wasserstein type loss. [Frogner'15] 17

Example: Generative Models



Example: Generative Models



Statistics 0.1: Density Fitting



Maximum Likelihood Estimation

ON AN ABSOLUTE CRITERION FOR FITTING FREQUENCY CURVES.

By R. A. Fisher, Gonville and Caius College, Cambridge.

1. IF we set ourselves the problem, in its frequent occurrence, of finding the arbitrary function of known form, which best suit a observations, we are met at the outset by an which appears to invalidate any results we ma



 $\max_{\boldsymbol{\theta}\in\Theta}\frac{\mathbf{I}}{N}\sum_{i=1}\log \boldsymbol{p}_{\boldsymbol{\theta}}(\boldsymbol{x}_{i})$

 $\nu_{\rm data}$

Maximum Likelihood Estimation

ON AN ABSOLUTE CRITERION FOR FITTING FREQUENCY CURVES.

By R. A. Fisher, Gonville and Caius College, Cambridge.

1. IF we set ourselves the problem, in its frequent occurrence, of finding the arbitrary function of known form, which best suit a observations, we are met at the outset by an which appears to invalidate any results we ma



 $\max_{\boldsymbol{\theta}\in\Theta}\frac{1}{N}\sum_{i}\log \boldsymbol{p}_{\boldsymbol{\theta}}(\boldsymbol{x}_{i})$ i=1

 $\log 0 = -\infty$ $p_{\theta}(x_i) \text{ must be } > 0$

PH_A

 $\nu_{\rm data}$
Maximum Likelihood Estimation



Maximum Likelihood Estimation



In higher dimensional spaces...



























- Formulation as adversarial problem [GPM...'14]
 - $\min_{\boldsymbol{\theta}\in\Theta} \max_{\text{classifiers } \boldsymbol{g}} \operatorname{Accuracy}_{\boldsymbol{g}} \left((\boldsymbol{f}_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}, +1), (\boldsymbol{\nu}_{\text{data}}, -1) \right)$



- Formulation as adversarial problem [GPM...'14]
 - $\min_{\boldsymbol{\theta} \in \Theta} \max_{\text{classifiers } \boldsymbol{g}} \operatorname{Accuracy}_{\boldsymbol{g}} \left((\boldsymbol{f}_{\boldsymbol{\theta} \sharp} \boldsymbol{\mu}, +1), (\boldsymbol{\nu}_{\text{data}}, -1) \right)$



- Formulation as adversarial problem [GPM...'14]
 - $\min_{\boldsymbol{\theta}\in\Theta} \max_{\text{classifiers } \boldsymbol{g}} \operatorname{Accuracy}_{\boldsymbol{g}} \left((\boldsymbol{f}_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}, +1), (\boldsymbol{\nu}_{\text{data}}, -1) \right)$



• Formulation as adversarial problem [GPM...'14]



• Formulation as adversarial problem [GPM...'14]

$\min_{\boldsymbol{\theta} \in \Theta \text{ classifiers } \boldsymbol{g}} \max_{\boldsymbol{\theta} \in \Theta \text{ classifiers } \boldsymbol{g}} \operatorname{Accuracy}_{\boldsymbol{g}} \left((\boldsymbol{f}_{\boldsymbol{\theta} \sharp} \boldsymbol{\mu}, +1), (\boldsymbol{\nu}_{\text{data}}, -1) \right)$



• Formulation as adversarial problem [GPM...'14]

$\min_{\boldsymbol{\theta} \in \Theta \text{ classifiers } \boldsymbol{g}} \max_{\boldsymbol{\theta} \in \Theta \text{ classifiers } \boldsymbol{g}} \operatorname{Accuracy}_{\boldsymbol{g}} \left((\boldsymbol{f}_{\boldsymbol{\theta} \sharp} \boldsymbol{\mu}, +1), (\boldsymbol{\nu}_{\text{data}}, -1) \right)$



- Formulation as adversarial problem [GPM...'14]
- $\min_{\boldsymbol{\theta} \in \Theta \text{ classifiers } \boldsymbol{g}} \max_{\boldsymbol{\theta} \in \Theta \text{ classifiers } \boldsymbol{g}} \operatorname{Accuracy}_{\boldsymbol{g}} \left((\boldsymbol{f}_{\boldsymbol{\theta} \sharp} \boldsymbol{\mu}, +1), (\boldsymbol{\nu}_{\text{data}}, -1) \right)$



Another idea?



• Use a metric Δ for probability measures, that can handle measures with non-overlapping supports:

$$\min_{\boldsymbol{\theta}\in\Theta} \Delta(\boldsymbol{\nu}_{data}, \boldsymbol{p}_{\boldsymbol{\theta}}), \quad \min_{\boldsymbol{\theta}\in\Theta} \mathrm{KL}(\boldsymbol{\nu}_{data} \| \boldsymbol{p}_{\boldsymbol{\theta}})$$

• The original GAN paper can be interpreted in that light using the Jensen-Shannon divergence.

Minimum Δ Estimation

The Annals of Statistics 1980, Vol. 8, No. 3, 457-487

MINIMU 1 CHI-SQUARE, NOT MAXIMUM LIKELIHOOD!

By JOSEPH BERKSON Mavo Clinic, Rochester, Minnesota



COMPUTATIONAL STATISTICS & DATA ANALYSIS

ELSEVIER Computational Statistics & Data Analysis 29 (1998) 81-103



Minimum Hellinger listance estimation for Poisson mixtures

Dimitris Karlis, Evdokia Xekalaki* Department of Statistics. Athens University of Economics and Business, 76 Patissian Str., 104 34 Athens, Greece



Available online at www.sciencedirect.com

SCIENCE DIRECT.



Statistics & Probability Letters 76 (2006) 1298-1302

www.elsevier.com/locate/stapro

On minimum Kantorovich listance estimators

Federico Bassetti^a, Antonella Bodini^b, Eugenio Regazzini^{a,*}

Δ Generative Model Estimation

Generative Moment Matching Networks

Training generative neural networks via Maximum Mean Discrepancy optimization

Yujia Li¹ Kevin Swersky¹ Richard Zemel^{1,2}

YUJIALI@CS.TORONTO.EDU KSWERSKY@CS.TORONTO.EDU ZEMEL@CS.TORONTO.EDU

¹Department of Computer Science, University of Toronto, Toronto, ON, CANADA ²Canadian Institute for Advanced Research, Toronto, ON, CANADA



Gintare Karolina Dziugaite University of Cambridge **Daniel M. Roy** University of Toronto Zoubin Ghahramani University of Cambridge

Chun-Liang Li^{1,*} Wei-Cheng Chang^{1,*} Yu Cheng² Yiming Yang¹ Barnabás Póczos¹ ¹ Carnegie Mellon University, ²IBM Research {chunlial,wchang2,yiming,bapoczos}@cs.cmu.edu chengyu@us.ibm.com

Δ Generative Model Estimation

Generative Moment Matching Networks Training generative neural networks via Maximum Mean Discrepancy optimization Yuiia Li¹ YUJIALI@CS.TORONTO.EDU Kevin Swersky¹ KSWERSKY @CS.TORONTO.EDU Richard Zemel^{1,2} ZEMEL@CS.TORONTO.EDU ¹Department of Computer Science, University of Toronto, Toronto, ON, CANADA. ²Canadian Institute for Advanced Research, Toronto, ON, CANADA **Gintare Karolina Dziugaite** Daniel M. Roy Zoubin Ghahramani University of Cambridge University of Toronto University of Cambridge **MMD** GAN: Towards Deeper Understanding of **Moment Matching Network** Wasserstein Fraining of **Restricted Boltzmann Machines** Chun-Liang Li^{1,*} Wei-Cheng Chang^{1,*} Yu Cheng² Yiming Yang¹ Barnabás Póczos¹ ¹ Carnegie Mellon University, ²IBM Research {chunlial,wchang2,yiming,bapoczos}@cs.cmu.edu chengyu@us.ibm.com **Grégoire Montavon** Klaus-Robert Müller* Technische Universität Berlin Technische Universität Berlin gregoire.montavon@tu-berlin.de klaus-robert.mueller@tu-berlin.de Inference in generative models using the Wasserstein distance Marco Cuturi CREST, ENSAE, Université Paris-Saclay marco.cuturi@ensae.fr Espen Bernton, Mathieu Gerber, Pierre E. Jacob, Christian P. Robert Wasserstein GAN Martin Arjovsky¹, Soumith Chintala², and Léon Bottou^{1,2} ¹Courant Institute of Mathematical Sciences

²Facebook AI Research

Δ Generative Model Estimation



33

Rutgers University dnm@cs.rutgers.edu

Optimal Transport in ML

OT is establishing itself as a generic toolbox to handle probability measures in ML tasks



OT Computations



OT Computations














Discrete OT Problem

```
c emd.c
   Image: Second c.6.1 + <No selected symbol > +
                                                                                                2, -, C, #, E
    1*
1
2
         end.c
3
4
        Last update: 3/14/98
5
        An implementation of the Earth Movers Distance.
G
         Based of the solution for the Transportation problem as described in
7
         "Introduction to Mathematical Programming" by F. S. Hillier and
g
        G. J. Lieberman, McGraw-Hill, 1990.
9
10
        Copyright (C) 1998 Yossi Rubner
11
         Computer Science Department, Stanford University
12
13
         E-Mail: rupher@cs.stanford.edu URL: http://vision.stanford.edu/~rupher
14
    *1
15
    /##include <stdio.h>
16
    #include <stdlib.h>+/
17
    #include <math.h>
18
19
    finclude "end.h"
20
21
22
    #define DEBUG_LEVEL 0
23
    1+
24
     DEBUG_LEVEL:
25
       0 = NO MESSAGES
        1 = PRINT THE NUMBER OF ITERATIONS AND THE FINAL RESULT
26
27
       2 = PRINT THE RESULT AFTER EVERY ITERATION
28
       3 = PRINT ALSO THE FLOW AFTER EVERY ITERATION
29
        4 - PRINT A LOT OF INFORMATION (PROBABLY USEFUL ONLY FOR THE AUTHOR)
30
    41
31
32
33
    #define MAX_SIG_SIZE1 (MAX_SIG_SIZE+1) /* FOR THE POSIBLE DUMMY FEATURE */
34
35
    /* NEW TYPES DEFINITION */
36
77
     /* node1_t IS USED FOR SINGLE-LINKED LISTS */
385
    typedef struct node1_t {
49
      int i:
40
      double val;
41
      struct node1_t *Next;
42
    } node1_t;
43
    /* node1_t IS USED FOR DOUBLE-LINKED LISTS */
44
45
   typedef struct node2_t {
46
      int i, j;
47
      double val;
48
      struct node2_t *NextC;
                                            /* NEXT COLUMN */
49
       struct node7_t *NextR;
                                             /* NEXT ROW */
50
    } node2_t;
51
52
53
    /* GLOBAL VARIABLE DECLARATION */
54
    static int _n1, _n2; /* SIGNATURES SIZES */
static float _C[MAX_SIG_SIZE1][MAX_SIG_SIZE1];/* THE COST MATRIX */
55
                                                    /* SIGNATURES SIZES */
56
    static node2_t _X[MAX_SIC_SIZE1+2]; /* THE EASIC VARIABLES VECTOR +/
57
     58
```

Discrete OT Problem

```
c emd.c
                                                                                                 2, , C, #, E
   Image: Second c.6.1 + <No selected symbol > +
    1*
1
2
         end.c
3
4
        Last update: 3/14/98
5
G
         An implementation of the Earth Movers Distance.
         Based of the solution for the Transportation problem as described in
7
         "Introduction to Mathematical Programming" by F. S. Hillier and
g
         G. J. Lieberman, McGraw-Hill, 1990.
9
10
11
         Copyright (C) 1998 Yossi Rubner
12
         Computer Science Department, Stanford University
13
         E-Mail: rupher@cs.stanford.edu URL: http://vision.stanford.edu/~rupher
14
    *1
15
16
    /##include <stdio.h>
    #include <stdlib.h>+/
17
    #include <math.h>
18
19
     finclude "emd.h"
20
21
22
    #define DEBUG_LEVEL 0
23
    1+
24
     DEBUG_LEVEL:
25
        0 = NO MESSAGES
        1 = PRINT THE NUMBER OF ITERATIONS AND THE FINAL RESULT
26
27
        2 = PRINT THE RESULT AFTER EVERY ITERATION
28
       3 = PRINT ALSO THE FLOW AFTER EVERY ITERATION
29
        4 - PRINT A LOT OF INFORMATION (PROBABLY USEFUL ONLY FOR THE AUTHOR)
    41
30
31
32
33
    #define MAX_SIG_SIZE1 (MAX_SIG_SIZE+1) /* FOR THE POSIBLE DUMMY FEATURE */
34
35
    /* NEW TYPES DEFINITION */
36
77
     /* node1_t IS USED FOR SINGLE-LINKED LISTS */
385
    typedef struct node1_t {
49
      int 1:
40
      double val;
      struct model_t *Next;
41
42
    } node1_t;
43
     /* node1_t IS USED FOR DOUBLE-LINKED LISTS */
44
15
    typedef struct node2_t {
46
      int i, j;
47
      double val;
48
       struct node2_t *NextC;
                                            /* NEXT COLUMN */
49
       struct node7_t *NextR;
                                             /* NEXT ROW */
    } node2_t;
50
51
52
53
    /* GLOBAL VARIABLE DECLARATION */
54
    static int _n1, _n2; /* SIGNATURES SIZES */
static float _C[MAX_SIG_SIZE1][MAX_SIG_SIZE1];/* THE COST MATRIX */
55
                                                    /* SIGNATURES SIZES */
56
57
    static node2_t _X[MAX_SIC_SIZE1+2]; /* THE EASIC VARIABLES VECTOR */
      58
```

Discrete OT Problem

```
c emd.c
                                                                                                 2, , C, #, E
   Image: Second c.6.1 + <No selected symbol > +
    1*
1
2
         end.c
3
4
        Last update: 3/14/98
5
G
         An implementation of the Earth Movers Distance.
         Based of the solution for the Transportation problem as described in
7
         "Introduction to Mathematical Programming" by F. S. Hillier and
g
         G. J. Lieberman, McGraw-Hill, 1990.
9
10
11
         Copyright (C) 1998 Yossi Rubner
12
         Computer Science Department, Stanford University
13
         E-Mail: rupher@cs.stanford.edu URL: http://vision.stanford.edu/~rupher
14
    *1
15
16
    /##include <stdio.h>
    #include <stdlib.h>+/
17
    #include <math.h>
18
19
     finclude "emd.h"
20
21
22
    #define DEBUG_LEVEL 0
23
    1+
24
     DEBUG_LEVEL:
25
        0 = NO MESSAGES
        1 = PRINT THE NUMBER OF ITERATIONS AND THE FINAL RESULT
26
27
        2 = PRINT THE RESULT AFTER EVERY ITERATION
28
       3 = PRINT ALSO THE FLOW AFTER EVERY ITERATION
29
        4 - PRINT A LOT OF INFORMATION (PROBABLY USEFUL ONLY FOR THE AUTHOR)
    41
30
31
32
33
    #define MAX_SIG_SIZE1 (MAX_SIG_SIZE+1) /* FOR THE POSIBLE DUMMY FEATURE */
34
35
    /* NEW TYPES DEFINITION */
36
77
     /* node1_t IS USED FOR SINGLE-LINKED LISTS */
385
    typedef struct node1_t {
49
      int 1:
40
      double val;
      struct model_t *Next;
41
42
    } node1_t;
43
     /* node1_t IS USED FOR DOUBLE-LINKED LISTS */
44
15
    typedef struct node2_t {
46
      int i, j;
47
      double val;
48
       struct node2_t *NextC;
                                            /* NEXT COLUMN */
49
       struct node7_t *NextR;
                                             /* NEXT ROW */
    } node2_t;
50
51
52
53
    /* GLOBAL VARIABLE DECLARATION */
54
    static int _n1, _n2; /* SIGNATURES SIZES */
static float _C[MAX_SIG_SIZE1][MAX_SIG_SIZE1];/* THE COST MATRIX */
55
                                                    /* SIGNATURES SIZES */
56
57
    static node2_t _X[MAX_SIC_SIZE1+2]; /* THE EASIC VARIABLES VECTOR */
      58
```

Solution: Regularization



Def. Regularized Wasserstein,
$$\gamma \ge 0$$

 $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$

$$E(P) \stackrel{\text{def}}{=} - \sum_{i,j=1}^{nm} P_{ij} (\log P_{ij} - 1)$$

Note: Unique optimal solution because of strong concavity of entropy

Def. Regularized Wasserstein,
$$\gamma \ge 0$$

 $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$



Note: Unique optimal solution because of strong concavity of entropy

Def. Regularized Wasserstein,
$$\gamma \ge 0$$

 $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$



Note: Unique optimal solution because of strong concavity of entropy

Def. Regularized Wasserstein,
$$\gamma \ge 0$$

 $W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \min_{\boldsymbol{P} \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$



"static" problem associated with Schrödinger problem

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle P, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

$$L(P,\alpha,\beta) = \sum_{ij} P_{ij}M_{ij} + \gamma P_{ij}(\log P_{ij} - 1) + \alpha^T (P\mathbf{1} - \mathbf{a}) + \beta^T (P^T\mathbf{1} - \mathbf{b})$$

 $\partial L/\partial P_{ij} = M_{ij} + \gamma \log P_{ij} + \alpha_i + \beta_j$ $(\partial L/\partial P_{ij} = 0) \Rightarrow P_{ij} = e^{\frac{\alpha_i}{\gamma}} e^{-\frac{M_{ij}}{\gamma}} e^{\frac{\beta_j}{\gamma}} = u_i K_{ij} v_j$

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

$$P_{\gamma} \in U(\boldsymbol{a}, \boldsymbol{b}) \Leftrightarrow \begin{cases} \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}) \boldsymbol{1}_{m} &= \boldsymbol{a} \\ \operatorname{diag}(\boldsymbol{v}) K^{T} \operatorname{diag}(\boldsymbol{u}) \boldsymbol{1}_{n} &= \boldsymbol{b} \end{cases}$$

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

$$P_{\gamma} \in U(\boldsymbol{a}, \boldsymbol{b}) \Leftrightarrow \begin{cases} \operatorname{diag}(\boldsymbol{u}) \boldsymbol{K} \operatorname{diag}(\boldsymbol{v}) \boldsymbol{1}_{m} &= \boldsymbol{a} \\ \operatorname{diag}(\boldsymbol{v}) \boldsymbol{K}^{T} \operatorname{diag}(\boldsymbol{u}) \boldsymbol{1}_{n} &= \boldsymbol{b} \end{cases}$$

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

$$P_{\gamma} \in U(\boldsymbol{a}, \boldsymbol{b}) \Leftrightarrow \begin{cases} \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}) \mathbf{1}_{m} &= \boldsymbol{a} \\ \operatorname{diag}(\boldsymbol{v}) K^{T} \operatorname{diag}(\boldsymbol{u}) \mathbf{1}_{n} &= \boldsymbol{b} \end{cases}$$

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

$$P_{\gamma} \in U(\boldsymbol{a}, \boldsymbol{b}) \Leftrightarrow \begin{cases} \operatorname{diag}(\boldsymbol{u}) K \boldsymbol{v} &= \boldsymbol{a} \\ \operatorname{diag}(\boldsymbol{v}) K^{T} \boldsymbol{u} &= \boldsymbol{b} \end{cases}$$

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

$$P_{\gamma} \in U(\boldsymbol{a}, \boldsymbol{b}) \Leftrightarrow \begin{cases} \boldsymbol{u} \odot \boldsymbol{K} \boldsymbol{v} &= \boldsymbol{a} \\ \boldsymbol{v} \odot \boldsymbol{K}^{T} \boldsymbol{u} &= \boldsymbol{b} \end{cases}$$

Prop. If
$$P_{\gamma} \stackrel{\text{def}}{=} \operatorname{argmin}_{P \in U(\boldsymbol{a}, \boldsymbol{b})} \langle \boldsymbol{P}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle - \gamma E(\boldsymbol{P})$$

then $\exists ! \boldsymbol{u} \in \mathbb{R}^{n}_{+}, \boldsymbol{v} \in \mathbb{R}^{m}_{+}$, such that
 $P_{\gamma} = \operatorname{diag}(\boldsymbol{u}) K \operatorname{diag}(\boldsymbol{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{\boldsymbol{X}\boldsymbol{Y}}/\gamma}$

$$P_{\gamma} \in U(\boldsymbol{a}, \boldsymbol{b}) \Leftrightarrow \begin{cases} \boldsymbol{u} = \boldsymbol{a}/K\boldsymbol{v} \\ \boldsymbol{v} = \boldsymbol{b}/K^{T}\boldsymbol{u} \end{cases}$$

Sinkhorn's Algorithm : Repeat

1.
$$\boldsymbol{u} = \boldsymbol{a}/K\boldsymbol{v}$$

2. $\boldsymbol{v} = \boldsymbol{b}/K^T\boldsymbol{u}$

Sinkhorn's Algorithm : Repeat

1.
$$\boldsymbol{u} = \boldsymbol{a}/K\boldsymbol{v}$$

2. $\boldsymbol{v} = \boldsymbol{b}/K^T\boldsymbol{u}$

- [Sinkhorn'64] proved convergence for the first time.
- [Lorenz'89] linear convergence, see [Altschuler'17]
- O(nm) complexity, GPGPU parallel [Cuturi'13].
- $O(n \log n)$ on gridded spaces using convolutions. [Solomon'15]





$$\mu = \sum_{i=1}^{n} a_{i} \delta_{x_{i}} \quad \nu = \sum_{j=1}^{m} b_{j} \delta_{y_{j}}$$

$$\mathcal{E}(\mu, \nu) = \langle ab^{T}, M_{XY} \rangle$$

$$W_{\gamma}(\mu, \nu) = \langle P_{\gamma}, M_{XY} \rangle$$

$$W^{p}(\mu, \nu) = \langle P^{\star}, M_{XY} \rangle$$

$$M_{XY} P^{\star}$$

$$\mu = \sum_{i=1}^{n} a_{i} \delta_{x_{i}} \quad \nu = \sum_{j=1}^{m} b_{j} \delta_{y_{j}}$$

$$\mathcal{E}(\mu, \nu) = \langle ab^{T}, M_{XY} \rangle$$

$$W_{\gamma}(\mu, \nu) = \langle P_{\gamma}, M_{XY} \rangle$$

$$W^{p}(\mu, \nu) = \langle P^{\star}, M_{XY} \rangle$$

$$M_{XY} P^{\star}$$

$$\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle \boldsymbol{a}\boldsymbol{b}^{T}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$
$$\mathcal{M}\mathcal{M}\mathcal{D}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2}(\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$
$$W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle P_{\gamma}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$
$$\bar{W}_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2}(W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\mu}) + W_{\gamma}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

$$W^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle \boldsymbol{P}^{\star}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

How to compare them?

i.i.d samples $x_1, \ldots, x_n \sim \mu, y_1, \ldots, y_m \sim \nu$, $\hat{\boldsymbol{\mu}}_{\boldsymbol{n}} \stackrel{\text{def}}{=} \frac{1}{n} \sum \delta_{\boldsymbol{x}_{\boldsymbol{i}}}, \hat{\boldsymbol{\nu}}_{\boldsymbol{m}} \stackrel{\text{def}}{=} \frac{1}{m} \sum \delta_{\boldsymbol{y}_{\boldsymbol{j}}}$ Computational properties Effort to compute/approximate $\Delta(\hat{\mu}_n, \hat{\nu}_m)$? Statistical properties $|\Delta(\boldsymbol{\mu}, \boldsymbol{\nu}) - \Delta(\hat{\boldsymbol{\mu}}_{\boldsymbol{n}}, \hat{\boldsymbol{\nu}}_{\boldsymbol{n}})| \leq f(n)?$

Sinkhorn in between W and MMD

$$\mathcal{MMD}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2}(\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

$$(n+m)^2$$

$$O(1/\sqrt{n})$$

$$W^{p}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle \boldsymbol{P}^{\star}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$
$$O((n+m)nm\log(n+m)) \qquad O(1/n^{1/d})$$

$$\mathcal{MMD}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2}(\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

$$(n+m)^{2} \qquad O(1/\sqrt{n})$$

$$\overline{W}_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) = W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2}(W_{\gamma}(\boldsymbol{\mu}, \boldsymbol{\mu}) + W_{\gamma}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

$$O((n+m)^{2}) \qquad O\left(\frac{1}{\gamma^{d/2}\sqrt{n}}\right) \qquad [\textbf{GCBCP'18}]$$

$$[\textbf{FSVATP'18}]$$

$$W^{p}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle \boldsymbol{P^{\star}}, M_{\boldsymbol{XY}} \rangle$$

$$O((n+m)nm\log(n+m) \qquad O(1/n^{1/d})$$

Differentiability of W

 $W((\boldsymbol{a}, \boldsymbol{X}), (\boldsymbol{b}, \boldsymbol{Y}))$



Differentiability of W

 $W((a + \Delta a, X), (b, Y)) = W((a, X), (b, Y)) + ??$



Differentiability of W

 $W((a + \Delta a, X), (b, Y)) = W((a, X), (b, Y)) + ??$



Sinkhorn ----> Differentiability

 $W((a, X + \Delta X), (b, Y)) = W((a, X), (b, Y)) + ??$



Sinkhorn ----> Differentiability

 $W((a, X + \Delta X), (b, Y)) = W((a, X), (b, Y)) + ??$



Sinkhorn: A Programmer View

Def. For $L \geq 1$, define $W_L(\boldsymbol{\mu},\boldsymbol{\nu}) \stackrel{\text{def}}{=} \langle \boldsymbol{P}_L, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle,$ where $P_L \stackrel{\text{def}}{=} \operatorname{diag}(\boldsymbol{u}_L) K \operatorname{diag}(\boldsymbol{v}_L)$, $\boldsymbol{v_0} = \boldsymbol{1}_m; l \ge 0, \boldsymbol{u_l} \stackrel{\text{def}}{=} \boldsymbol{a}/K\boldsymbol{v_l}, \boldsymbol{v_{l+1}} \stackrel{\text{def}}{=} \boldsymbol{b}/K^T\boldsymbol{u_l}.$ **Prop.** $\frac{\partial W_L}{\partial X}, \frac{\partial W_L}{\partial a}$ can be computed recursively, in O(L) kernel $K \times$ vector products.

Sinkhorn: A Programmer View

Def. For
$$L \ge 1$$
, define
 $W_L(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \langle \boldsymbol{P_L}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$,



Sinkhorn $\ell = 1, \ldots, L-1$

Sinkhorn: A Programmer View

Def. For
$$L \ge 1$$
, define
 $W_L(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \langle \boldsymbol{P_L}, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$,

Prop. $\frac{\partial W_L}{\partial \mathbf{X}}, \frac{\partial W_L}{\partial \mathbf{a}}$ can be computed recursively, in O(L) kernel $K \times \text{vector products.}$

[Hashimoto'16] [Bonneel'16][Shalit'16]

Primal Descent on Regularized W



[Cuturi'14]
Primal Descent on Regularized W



[Cuturi'14]

Primal Descent on Regularized W



[Cuturi'14]

On Regularizing or Not





On Regularizing or Not



[Schmitzer'16]

On Regularizing or Not



Dictionary Learning

 $\min_{\boldsymbol{A}\in(\Sigma_{R})^{K},\boldsymbol{\Lambda}\in(\Sigma_{K})^{N}}\sum_{i=1}^{N}W\left(\boldsymbol{b_{i}},\sum_{k=1}^{K}\boldsymbol{\Lambda_{k}^{i}a_{k}}\right)$



[Sandler'11] [Zen'14] [Rolet'16] $_{57}$

Dictionary Learning

 $\min_{\boldsymbol{A}\in(\Sigma_{n})^{K},\boldsymbol{\Lambda}\in(\Sigma_{K})^{N}}\sum_{i=1}^{N}W\left(\boldsymbol{b}_{\boldsymbol{i}},\sum_{k=1}^{K}\boldsymbol{\Lambda}_{\boldsymbol{k}}^{\boldsymbol{i}}\boldsymbol{a}_{\boldsymbol{k}}\right)$



OT Dictionary Learning

• [Hoffman'98] proposed to learn dictionaries (topics) for text, seen as histograms-of-words.

$$\Omega = \{ \text{words} \}, \quad |\Omega| \approx 13,000$$

Vector embeddings for words [Mikolov'13]
[Pennington'14] defines geometry:

$$\boldsymbol{D}(\text{public}, \text{car}) = \|x_{\text{public}} - x_{\text{car}}\|^2$$

• Data: 7,034 Reuters, 737 BBC sports news articles

Topic Models



[**Rolet'16**]

Inverse Wasserstein Problems

• consider Barycenter operator:

$$\boldsymbol{b}(\lambda) \stackrel{\text{def}}{=} \operatorname{argmin}_{\boldsymbol{a}} \sum_{i=1}^{N} \lambda_i W_{\gamma}(\boldsymbol{a}, \boldsymbol{b}_i)$$

• address now Wasserstein inverse problems:

Given \boldsymbol{a} , find $\operatorname*{argmin}_{\lambda \in \Sigma_N} \mathcal{E}(\lambda) \stackrel{\text{def}}{=} \operatorname{Loss}(\boldsymbol{a}, \boldsymbol{b}(\lambda))$

Wasserstein Inverse Problems



Barycenters = Fixed Points

Prop. [BCCNP'15] Consider
$$\boldsymbol{B} \in \Sigma_d^N$$

and let $\boldsymbol{U_0} = \boldsymbol{1_{d \times N}}$, and then for $l \ge 0$:
 $\boldsymbol{b}^{l \text{ def}} \exp\left(\log\left(K^T \boldsymbol{U_l}\right)\lambda\right); \begin{cases} \boldsymbol{V_{l+1}} \stackrel{\text{def}}{=} \frac{\boldsymbol{b}^{l} \boldsymbol{1}_N^T}{K^T \boldsymbol{U_l}}, \\ \boldsymbol{U_{l+1}} \stackrel{\text{def}}{=} \frac{\boldsymbol{B}}{K \boldsymbol{V_{l+1}}}. \end{cases}$

Using Truncated Barycenters

- instead of using the exact barycenter $\operatorname{argmin} \mathcal{E}(\lambda) \stackrel{\text{def}}{=} \operatorname{Loss}(\boldsymbol{a}, \boldsymbol{b}(\lambda))$ $\lambda \in \Sigma_N$
- use instead the L-iterate barycenter

$$\operatorname{argmin}_{\lambda \in \Sigma_N} \mathcal{E}^{(L)}(\lambda) \stackrel{\text{def}}{=} \operatorname{Loss}(\boldsymbol{a}, \boldsymbol{b}^{(L)}(\lambda))$$

• Differente using the chain rule.

$$\nabla \mathcal{E}^{(L)}(\lambda) = [\partial \boldsymbol{b}^{(L)}]^T(\boldsymbol{g}), \ \boldsymbol{g} \stackrel{\text{def}}{=} \nabla \text{Loss}(\boldsymbol{a}, \cdot)|_{\boldsymbol{b}^{(L)}(\lambda)}.$$

Gradient / Barycenter Computation

$$\begin{aligned} & \text{function SINKHORN-DIFFERENTIATE}((p_s)_{s=1}^S, q, \lambda) \\ & \forall s, b_s^{(0)} \leftarrow 1 \\ & (w, r) \leftarrow (0^S, 0^{S \times N}) \\ & \text{for } \ell = 1, 2, \dots, L \quad // Sinkhorn \ loop \\ & \forall s, \varphi_s^{(\ell)} \leftarrow K^\top \frac{p_s}{Kb_s^{(\ell-1)}} \\ & p \leftarrow \prod_s \left(\varphi_s^{(\ell)}\right)^{\lambda_s} \\ & \forall s, b_s^{(\ell)} \leftarrow \frac{p}{\varphi_s^{(\ell)}} \\ & g \leftarrow \nabla \mathcal{L}(p, q) \odot p \\ & \text{for } \ell = L, L - 1, \dots, 1 \quad // Reverse \ loop \\ & \forall s, w_s \leftarrow w_s + \langle \log \varphi_s^{(\ell)}, g \rangle \\ & \forall s, r_s \leftarrow -K^\top \left(K(\frac{\lambda_s g - r_s}{\varphi_s^{(\ell)}}) \odot \frac{p_s}{(Kb_s^{(\ell-1)})^2}\right) \odot b_s^{(\ell-1)} \\ & g \leftarrow \sum_s r_s \\ & \text{return } P^{(L)}(\lambda) \leftarrow p, \nabla \mathcal{E}_L(\lambda) \leftarrow w \end{aligned}$$

Application: Volume Reconstruction



Shape database (p_1, \ldots, p_5)

Input shape q

Projection $P(\lambda)$

Iso-surface

[Bonneel'16]







 $\lambda_0 = 0.03$

 $\lambda_1 = 0.12$



 $\lambda_2 = 0.40$



 $\lambda_{3} = 0.43$





Wasserstein Barycentric Coordinates: Histogram Regression using Optimal Transport, **SIGGRAPH'16**

[**BPC'16**]

Application: Brain Mapping















Euclidean projection







Wasserstein projection

Application: Brain Mapping







end-to-end W Dictionary Learning

N $\min_{\boldsymbol{A} \in (\Sigma_{n})^{K} \boldsymbol{\Lambda} \in (\Sigma_{K})^{N}} \sum_{i=1}^{K} \mathcal{L}\left(\boldsymbol{b}_{i}, \boldsymbol{a}(\boldsymbol{\lambda}_{i})\right)$

[Schmitz'18]



end-to-end W Dictionary Learning



Minimum Kantorovich Estimators

$$\min_{\boldsymbol{\theta}\in\Theta} W(\boldsymbol{\nu}_{\text{data}}, f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu})$$

[Bassetti'06] 1st reference discussing this approach.

Challenge:
$$\nabla_{\boldsymbol{\theta}} W(\boldsymbol{\nu}_{\text{data}}, f_{\boldsymbol{\theta} \sharp} \boldsymbol{\mu})$$
?

[Montavon'16] use regularized OT in a finite setting.

[**Arjovsky'17**] (WGAN) uses a NN to approximate dual solutions and recover gradient w.r.t. parameter

[Bernton'17] (Wasserstein ABC)

[Genevay'17, Salimans'17] (Sinkhorn approach)

Proposal: Autodiff OT using Sinkhorn

Approximate W loss by the transport cost \overline{W}_L after L Sinkhorn iterations.



[GPC'17]

Example: MNIST, Learning f_{θ}



Example: MNIST, Learning f_{θ}

| | 0 - | 5 | 5 | 5 | 5 | 8 | 8 | 8 | 8 | 8 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
|-------------------|--------------------------|---|-----|---|---|---|---|-----|----|---|----|-----|---|---|-----|---|---|---|-----|---|---|
| | | 5 | 5 | 5 | 8 | 8 | 8 | 8 | 8 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| | 100 - | 5 | 5 | 5 | 8 | 8 | 8 | 8 | 3 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| | | 5 | 5 | 5 | 8 | 8 | 8 | 2 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| | | 5 | 5 | 5 | 8 | 8 | 3 | 2 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| | | 5 | 5 | 5 | 8 | 8 | 3 | 2 | 2 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| Tataat | | 5 | 5 | 5 | 5 | 3 | 3 | 2 | 2 | 2 | 7 | 1 | 1 | 1 | T | 1 | 1 | 1 | 1 | 1 | 1 |
| Latent | 200 - | 3 | 5 | 5 | 3 | 3 | 3 | 2 | 2 | 2 | J. | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 010000 | | 3 | 3 | 3 | 3 | 3 | 3 | 2 | 2 | 2 | 5 | 5 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| space | | з | 3 | 3 | 3 | 3 | 3 | 3 | 2 | 4 | 6 | £ | 4 | 1 | Ŧ | 1 | 1 | 1 | 7 | 7 | 1 |
| - | 300 - 3 3 3 | 3 | 3 | 3 | 3 | 3 | 5 | 6 | 6 | 6 | 6 | 6 | 4 | 4 | 9 | 7 | 7 | 7 | 7 | 7 | 1 |
| $[0 \ 1]^2$ | | 3 | 3 | 3 | 3 | 8 | 6 | 6 | 6 | 6 | 6 | 4 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 7 | 1 |
| $[0, \mathbf{I}]$ | | 3 | 3 | 3 | В | Б | 6 | 6 | 6 | 6 | 4 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 |
| | | 3 | 0 | 0 | 0 | 0 | 6 | 6 | 6 | 6 | 4 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 |
| | 400 - | 0 | 0 | 0 | 0 | 0 | 6 | 6 | 6 | 4 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 |
| | | 0 | 0 | 0 | 0 | 0 | 6 | 6 | 6 | 4 | 9 | 9 | 9 | 9 | 9 | 4 | 4 | 9 | 9 | 9 | 9 |
| | | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 6 | 4 | 9 | 9 | 9 | 9 | 4 | 4 | 9 | 9 | 7 | 7 | 7 |
| | 500 - | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 6 | 4 | 9 | 9 | 9 | 9 | 7 | 7 | 7 | 7 | 7 | 7 | 7 |
| | (| 0 | 0 | 0 | 0 | 0 | 0 | 6 | 6 | 9 | 9 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 |
| | | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 64 | q | ٦ | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 |
| | (| ò | 100 | | | | | 200 | | | | 300 | | | 400 | | | | 500 | | |

75

Example: Generation of Images



arxiv.org/pdf/1710.05488

[Salimans'18]

Example: Generation of Images



arxiv.org/pdf/1710.05488

[Salimans'18]

Concluding Remarks

- *Regularized* OT is much faster than OT when handled in full generality.
- *Regularized* OT can interpolate between W and the MMD / *Energy distance* (MMD) metrics.
- The solution of *regularized OT* is *"auto-differentiable"*.
- Many open problems remain in ML that can be addressed with OT.