

L'IA au LMAC, Sciences de Données avec Transport Optimal

Mokhtar Z. Alaya



Journée Scientifique, Chaire SAFE AI

UTC, Octobre 2022

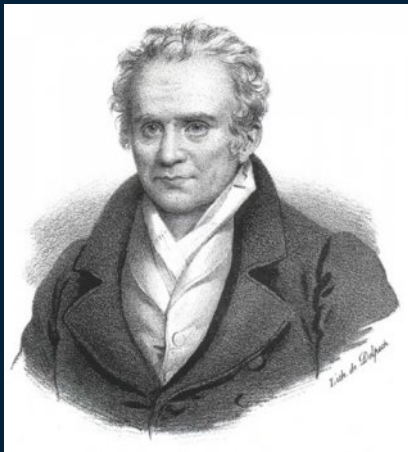
donnons un sens à l'innovation



I. What is optimal transport
(OT)?

OT is ...

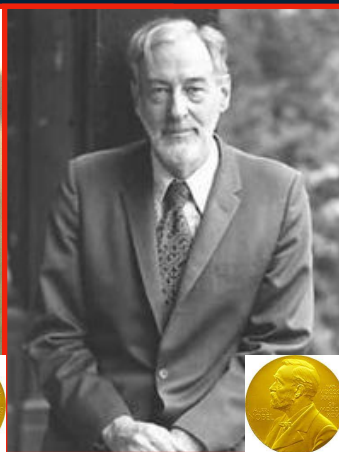
A method for comparing probability distributions with the ability to incorporate spatial information.



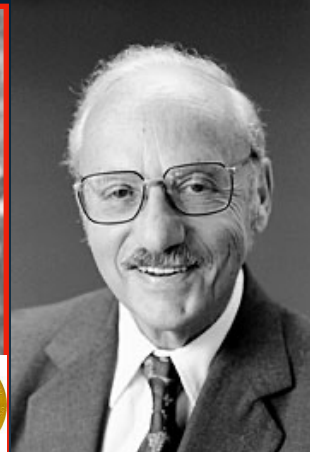
G. Monge
(1746 - 1818)



L. Kantorovich
(1912 - 1986)



T. Koopmans
(1910 - 1985)



G. Dantzig
(1914 - 2005)



Y. Brenier



F. Otto



R. McCann

Nobel Prize '75



C. Villani

Fields' 10



A. Figalli

Fields' 18

Origin: Monge Problem (1781)

666. MÉMOIRES DE L'ACADÉMIE ROYALE



M É M O I R E

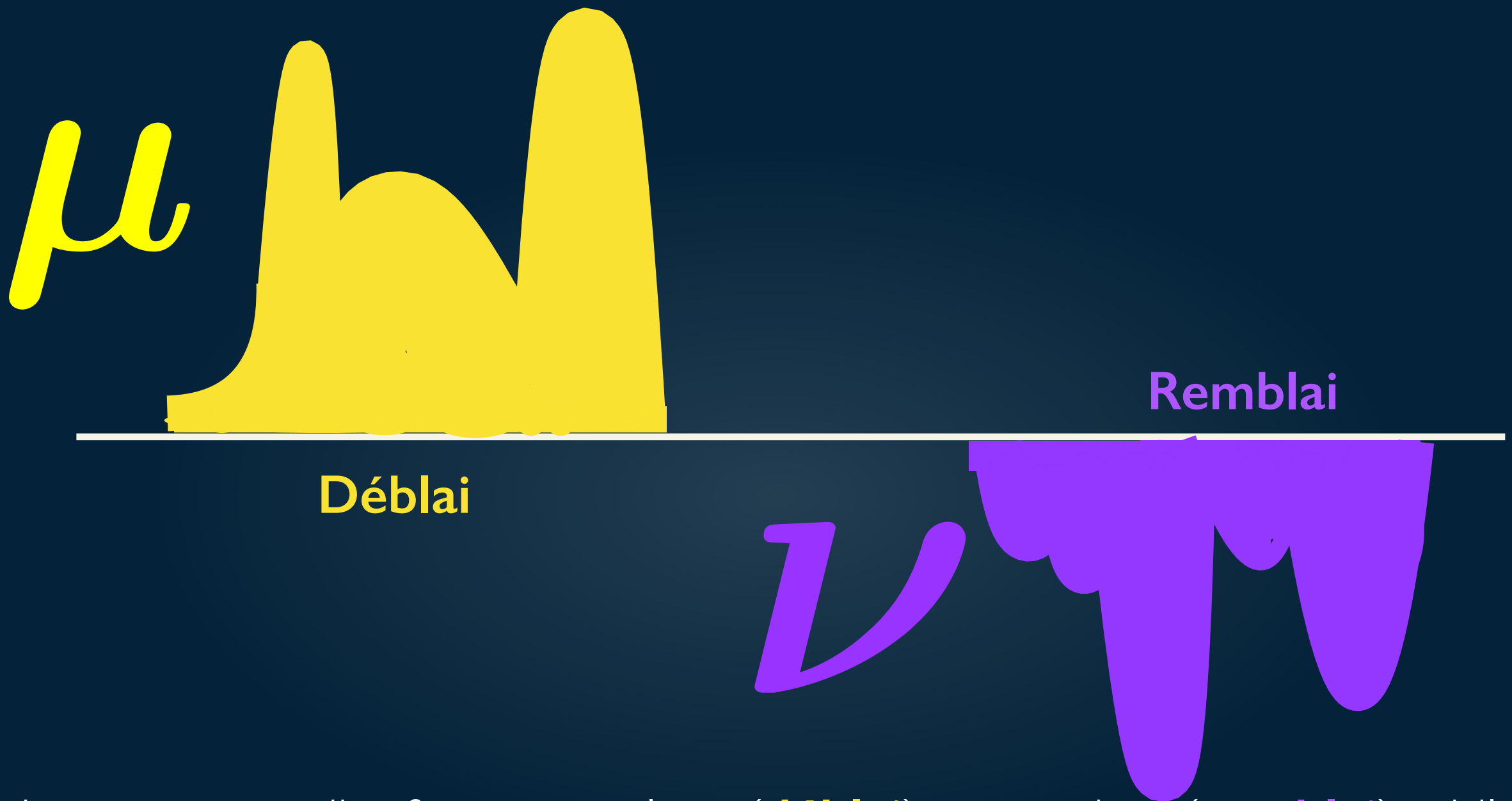
S U R L A

T H É O R I E D E S D É B L A I S
E T D E S R E M B L A I S.

Par M. M O N G E.

LORSQU'ON doit transporter des terres d'un lieu dans un autre, on a coutume de donner le nom de *Déblai* au volume des terres que l'on doit transporter, & le nom de *Remblai* à l'espace qu'elles doivent occuper après le transport.

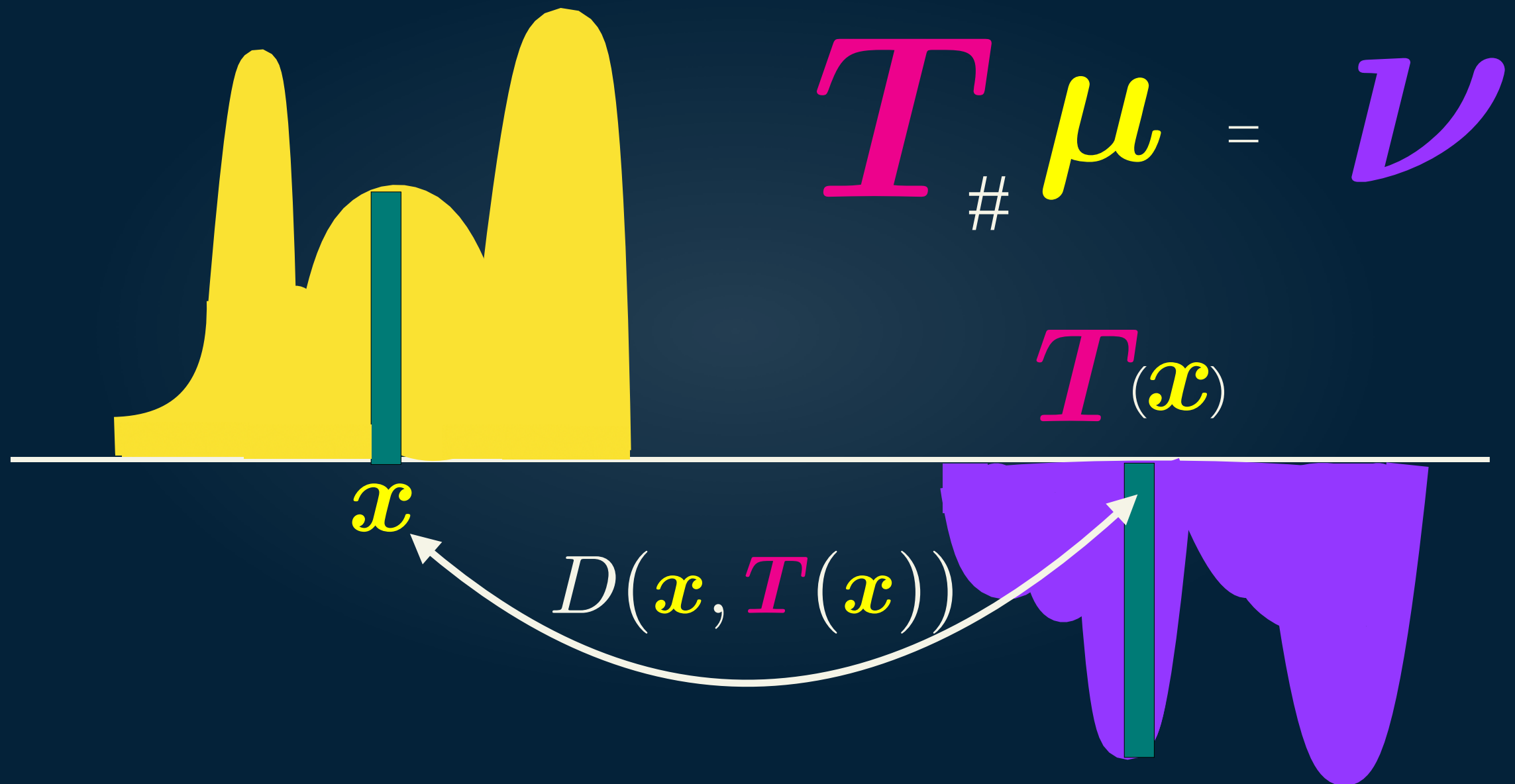
Monge Problem (1781)



- How to move dirt from one place (**déblai**) to another (**remblai**) while minimizing the effort?
- Find a mapping T between the two distributions of mass (**transport**).
- Optimize with respect to a displacement cost (**optimal**).

Monge Problem (1781)

- The mapping T must **push-forward** the “**déblai**” measure towards the “**remblai**”.



Monge Problem (1781)

- Monge formulation aim at finding a mapping T such that:

$$\inf_{T_{\# \mu = \nu}} \int C(x, T(x)) \mu(x) dx$$

- Mapping T does not exist in the general case.
- Brenier, 1991 proved existence and unicity of the Monge map for Euclidean cost and distributions with densities.

Discrete OT Framework

$$\nu = \sum_{j=1}^m \nu_j \delta_{y_j}$$

$$\mu = \sum_{i=1}^n \mu_i \delta_{x_i}$$

$$D^p(x_i, y_j)$$

C
Cost Matrix

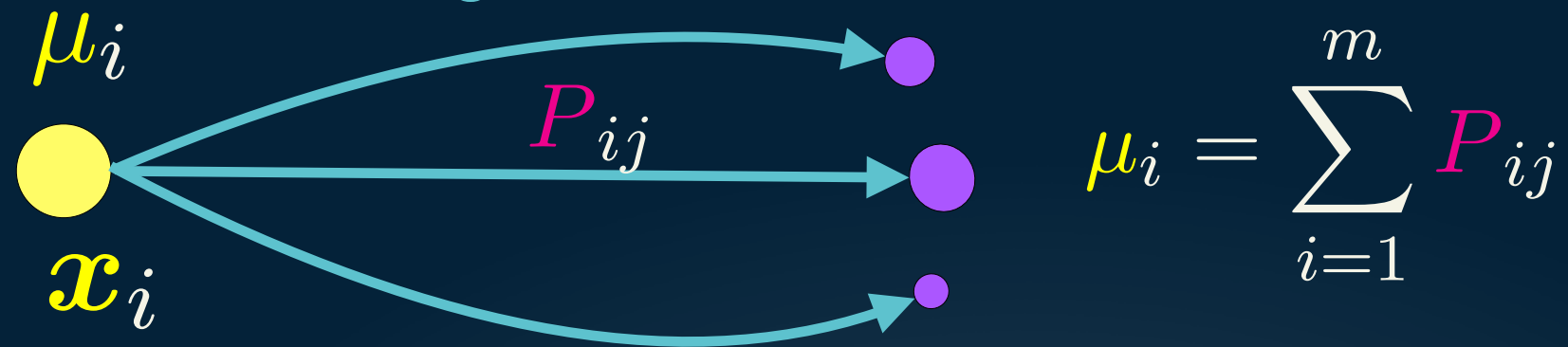
$$C_{ij} = (D^p(x_i, y_j))$$

$$\min_{T_{\#} \mu = \nu} \sum_i C(x_i, T(x_i)) \mu_i$$

$$\forall j \in \{1, \dots, m\}, \nu_j = \sum_{i: T(x_i) = y_j} \mu_i$$

Kantorovich's Formula (1942)

Relaxed: Fractional Assignments



- Focus on where the mass goes, allow splitting.
- Applications mainly for resource allocation problems.

$$\min_{\gamma \in \Pi_{\text{con}}(\mu, \nu)} \int_{\mathbb{R}^n \times \mathbb{R}^m} C(x, y) \gamma(x, y) dx dy$$

Probabilistic couplings set (Transport Polytope)

$$\Pi_{\text{con}}(\mu, \nu) = \left\{ \gamma \geq 0, \int_{\mathbb{R}^m} \gamma(x, y) dy = \mu, \int_{\mathbb{R}^n} \gamma(x, y) dx = \nu \right\}$$

Mass conservation constraints

Discrete Version

$$\Pi_{\text{dis}}(\mu, \nu) = \left\{ P \in \mathbb{R}_+^{n \times m}, P \mathbf{1}_m = \mu, P^\top \mathbf{1}_n = \nu \right\}$$

A first simple examples

Matching words embeddings



Word Mover's Distance avec Word2vec embeddings
[Kusner et al, 2015 (ICML)]

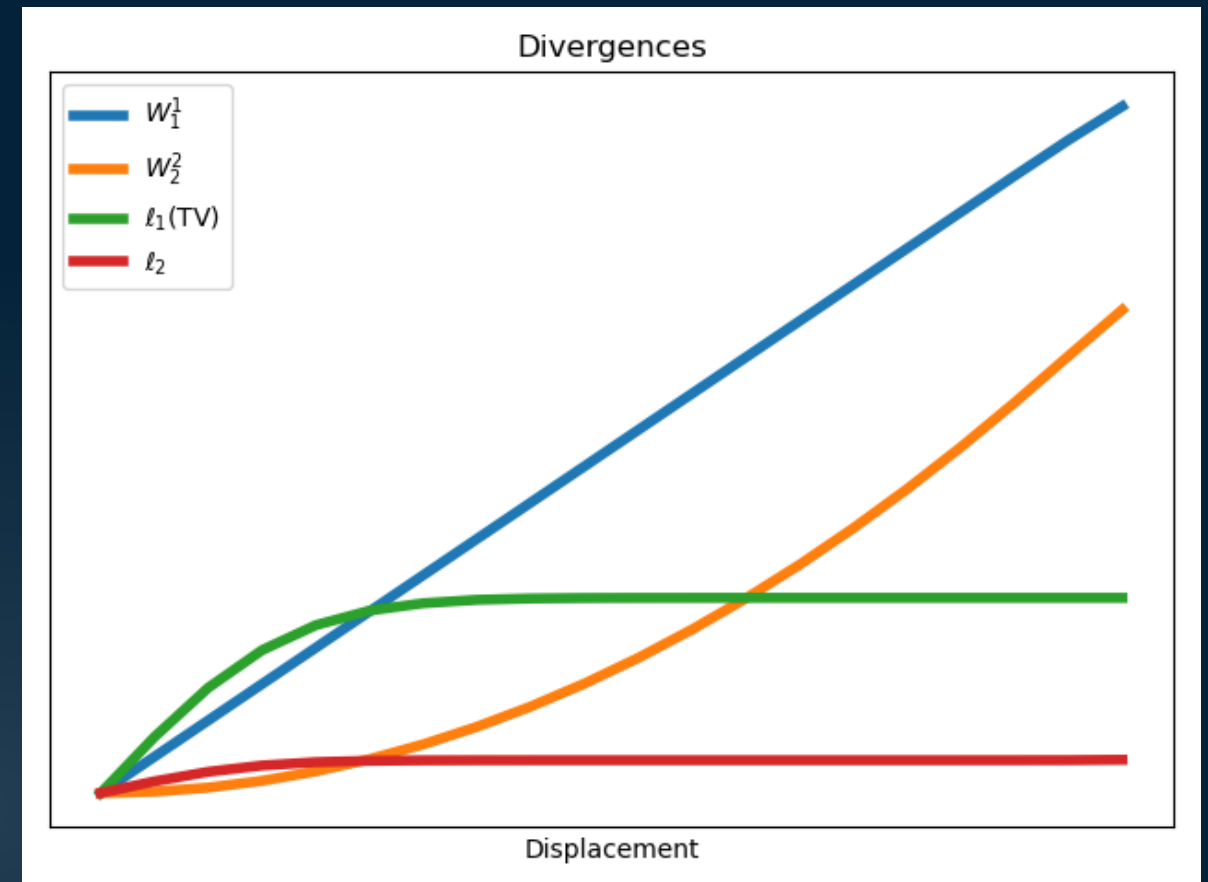
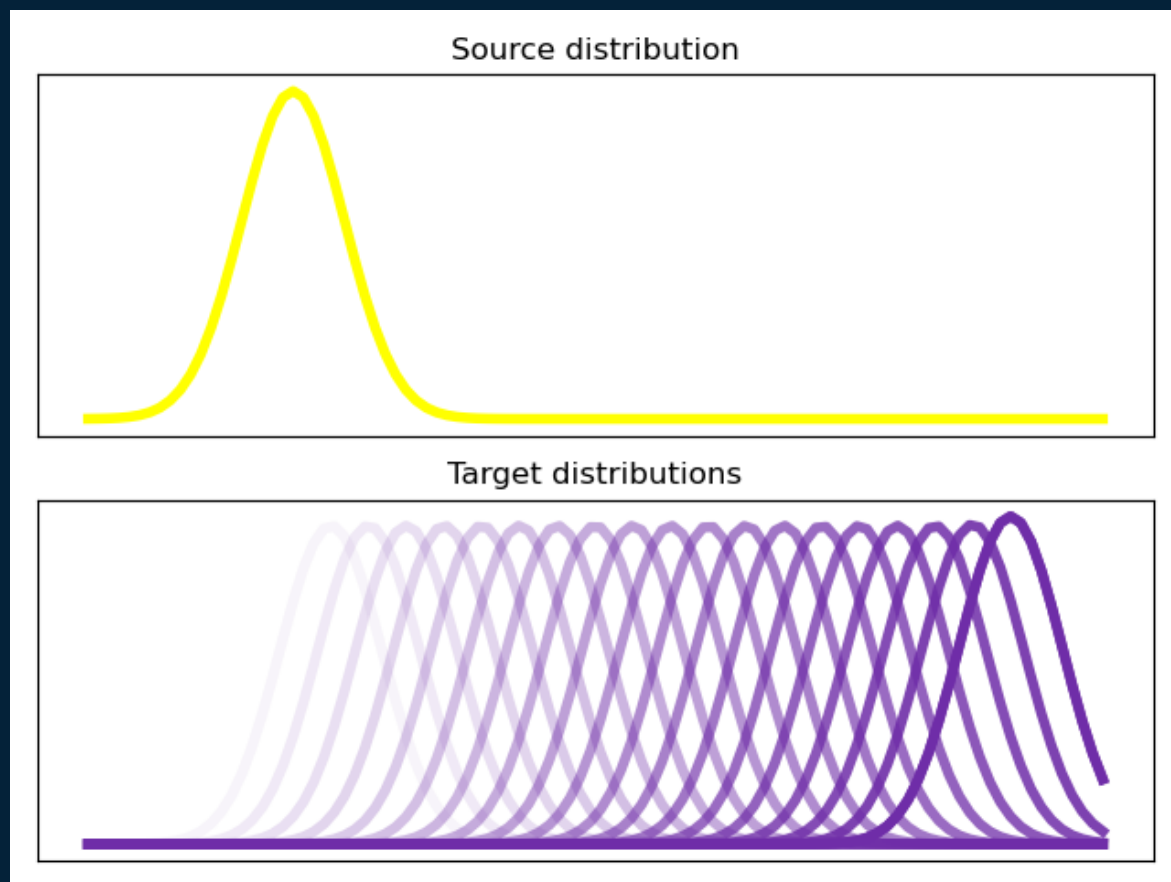
- Words are embedded in a high-dimensional space with deep neural networks.
- Matching two documents in an OT problem, with the Euclidean distance in the embedded space.

Color transfer



Wasserstein distance

Wasserstein distance



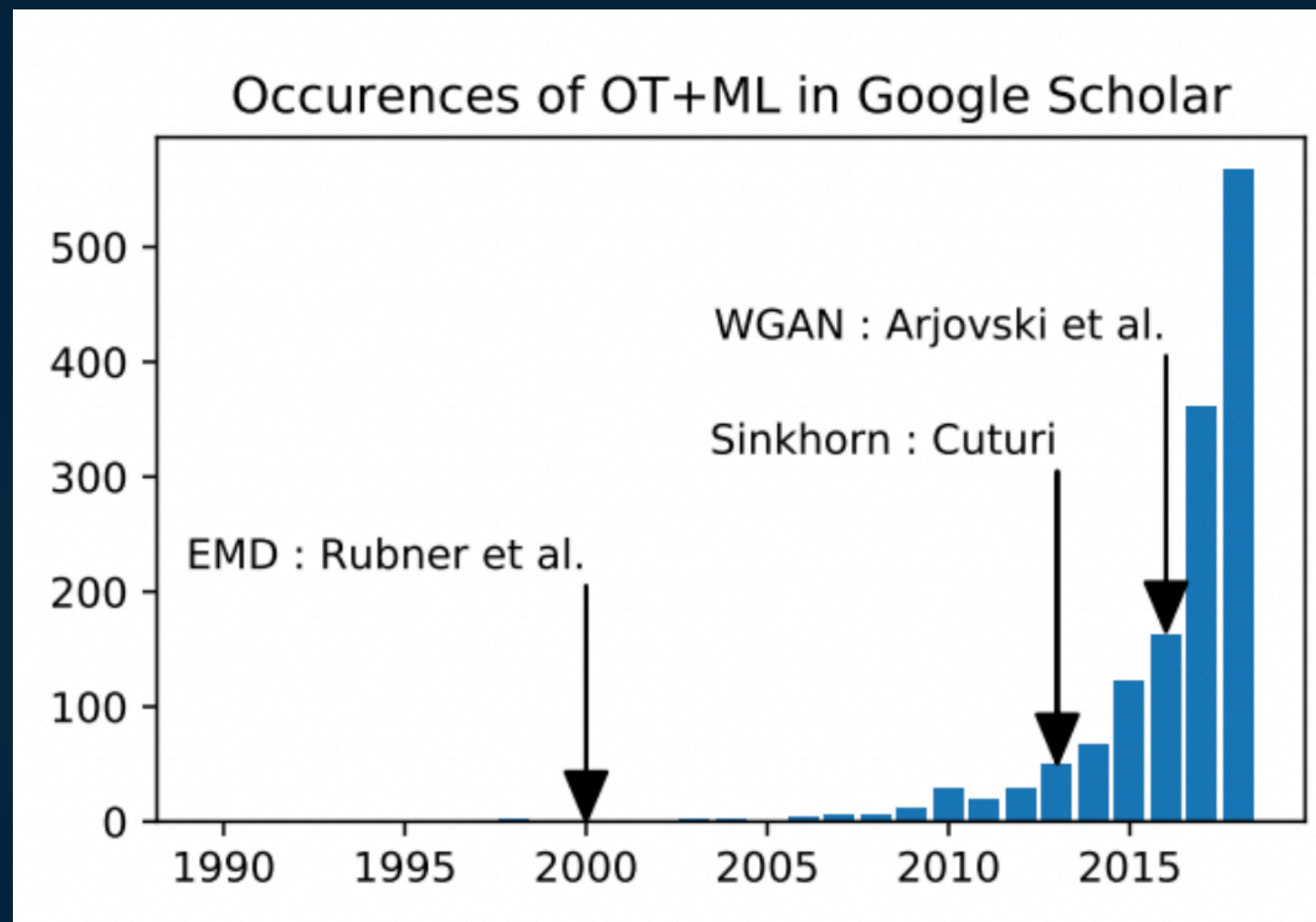
Monge-Kantorovich / Wasserstein Distance

$$\mathcal{W}_p^p(\mu, \nu) = \min_{\gamma \in \Pi_{\text{con}}(\mu, \nu)} \int_{\mathbb{R}^n \times \mathbb{R}^m} C(x, y) \gamma(x, y) dx dy = \mathbb{E}_{(x, y) \sim \gamma} [C(x, y)]$$

- Do not need the distributions to have overlapping support.
- Works for continuous and discrete distributions (histograms, empirical).

2. How can it be used in data science?

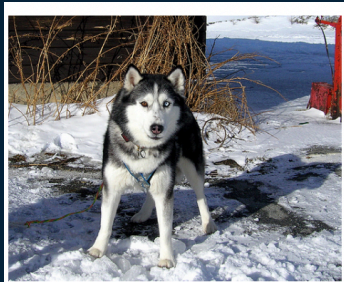
History of OT for machine learning



[R. Flamary, 2019 (HDR)]

- Recently introduced to ML (well know in image processing since 2000).
- Computational OT allows numerous applications (regularization).
- Deep learning boost (numerical optimisation and GAN).

Wasserstein distance as a multi-label loss



Siberian husky



Eskimo dog



Flickr user tags: street, parade, dragon

Predictions: people, protest, parade



Flickr user tags: water, boat, reflection, sunshine

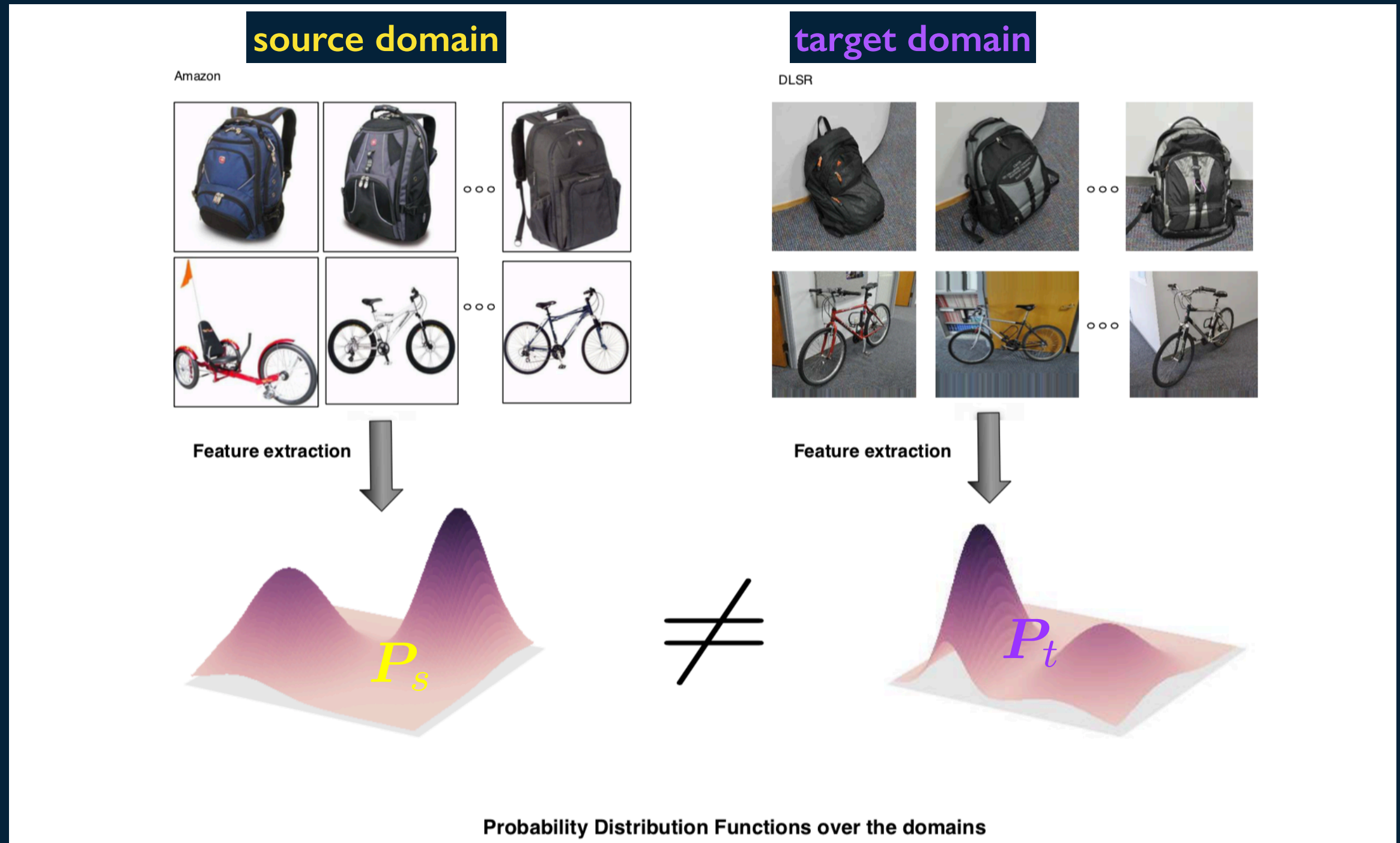
Predictions: water, river, lake, summer

Leveraging output space structure [Frogner et al., 2015, (NeurIPS)]

- Classes of a multiclass (multi-label) problem have structure.
- Takes into account semantic of classes in the output distribution probability.
- Error in "similar" class is less penalized than to dissimilar one .
- Can be represented as a Wasserstein distance between true label and output a model.
- Ground metric represent the distance between classes

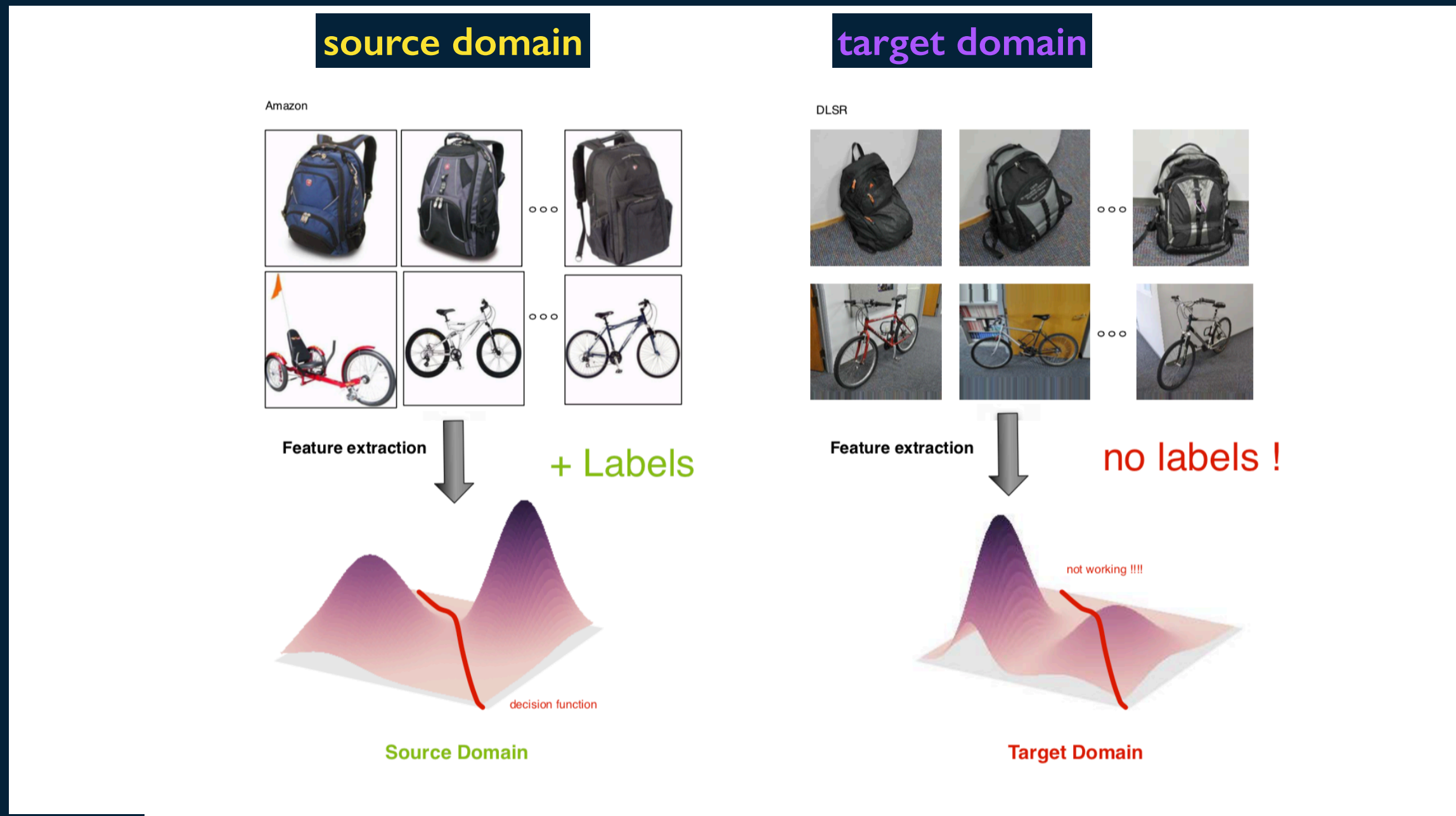
$$\min_{f_{\theta}} \frac{1}{n} \sum_{i=1}^n W_{17}^1 \left(f_{\theta}(x_i), y_i \right)$$

Domain Adaptation Problem



- We have a Classification problem with data coming from different source (domains).
- Distributions are different but related.

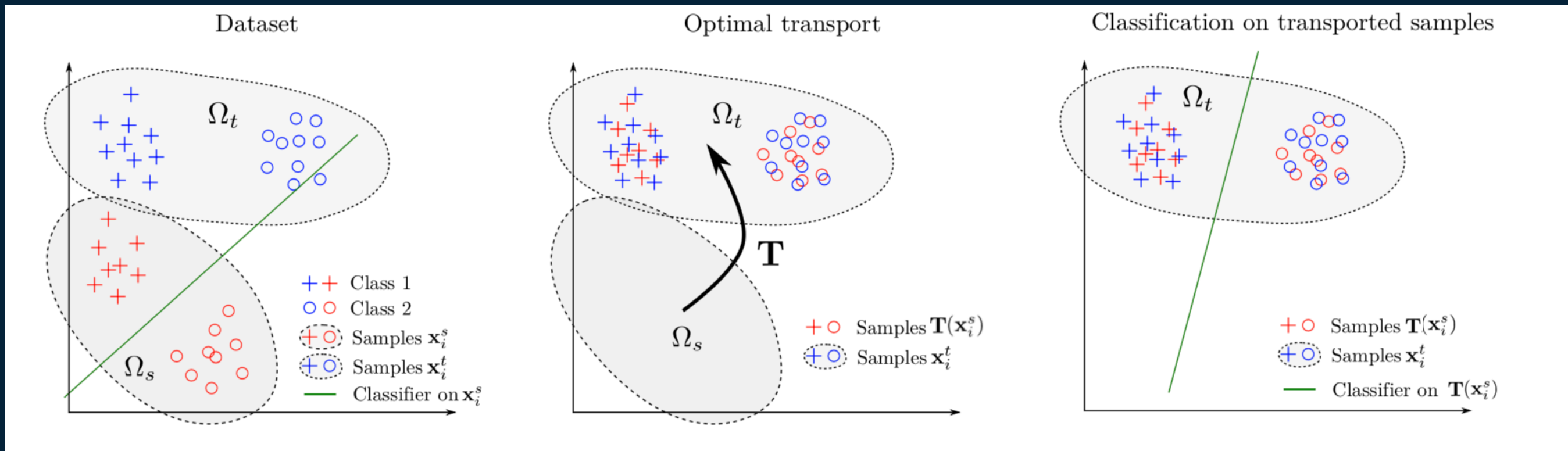
Domain Adaptation Problem



Problems

- Labels only available in the source domain, and classification is conducted in the target domain.
- Classifier on the source domain data performs badly in the target domain.

OT for Domain Adaptation



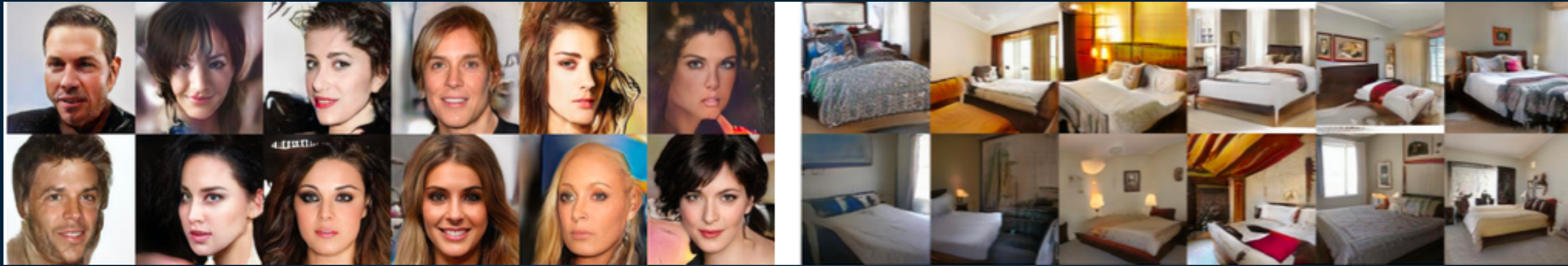
- There exist a transport T in the feature space between the two domains.
- The transport preserves the conditional distributions:

$$P_s[y|\mathbf{x}_s] = P_t[y|T(\mathbf{x}_s)]$$

3-step strategy [Courty et al., 2017]

1. Estimate optimal transport between distributions.
2. Transport the training samples onto the target distribution using barycentric mapping [Ferradans et al., 2013].
3. Learn a classifier on the transported training samples.

Wasserstein loss for generative modelling



Generative modelling as a matching distribution problem

- Learn a model that maps random vector to target space.
- Distribution of the model is targeted to be similar to the learning samples.
- Similarity as Wasserstein sense [Arjovsky et al. 2017, Deshpande et al. 2018, Nguyen et al. 2020).

$$\min_{f_{\theta}} W_p^p \left(\{f_{\theta}(z_i)\}_{i=1}^K, \{x_j\}_{j=1}^K \right)$$

$\{z_i\}$ some random vectors, $\{x_j\}$ some samples from the target distribution.

3. Conclusion

Take Home Message

- A powerful tool, well theoretically grounded, for manipulating distributions in machine learning.
- Despite its initial computational complexity, a lot of applications, even in large scale/deep learning settings.
- Others OT aspects (out the scope of the presentation): unbalanced OT, Gromov-Wasserstein distance (working with structured data), and many more

Some References

- G. Peyré and M. Cuturi,
Computational Optimal Transport with Applications to Data Sciences, 2019
- N. country, R. Flamary, D. Tuia and A. Rakotomamonjy.
Optimal Transport for Domain Adaptation, *PAMI* 2017
- R. Flamary et al. POT: Python Optimal Transport Library, 2017.

POT: Python Optimal Transport

Contents

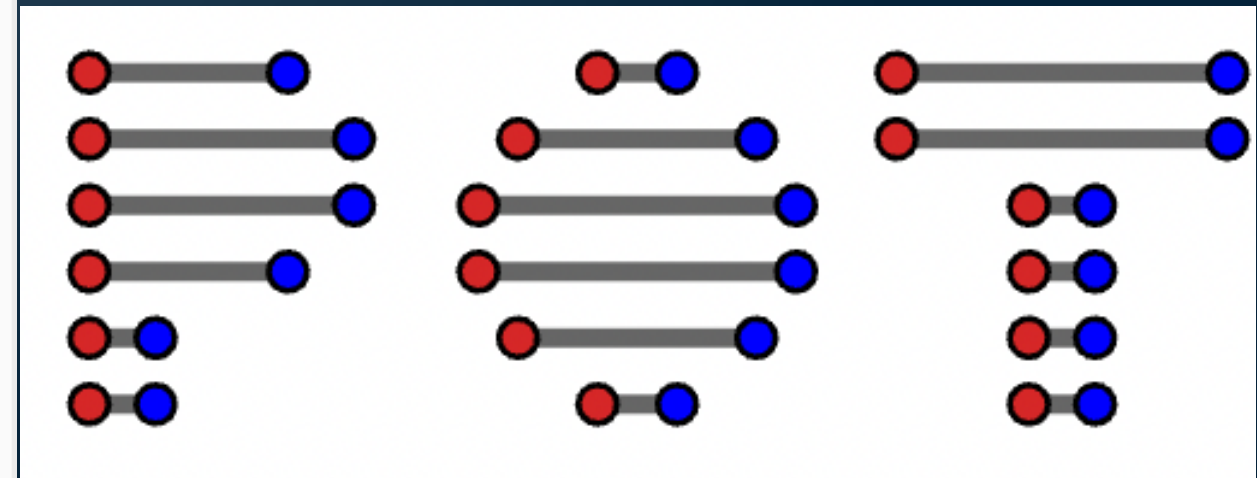
- [POT: Python Optimal Transport](#)
- [Quick start guide](#)
- [API and modules](#)
- [Examples gallery](#)
- [Releases](#)

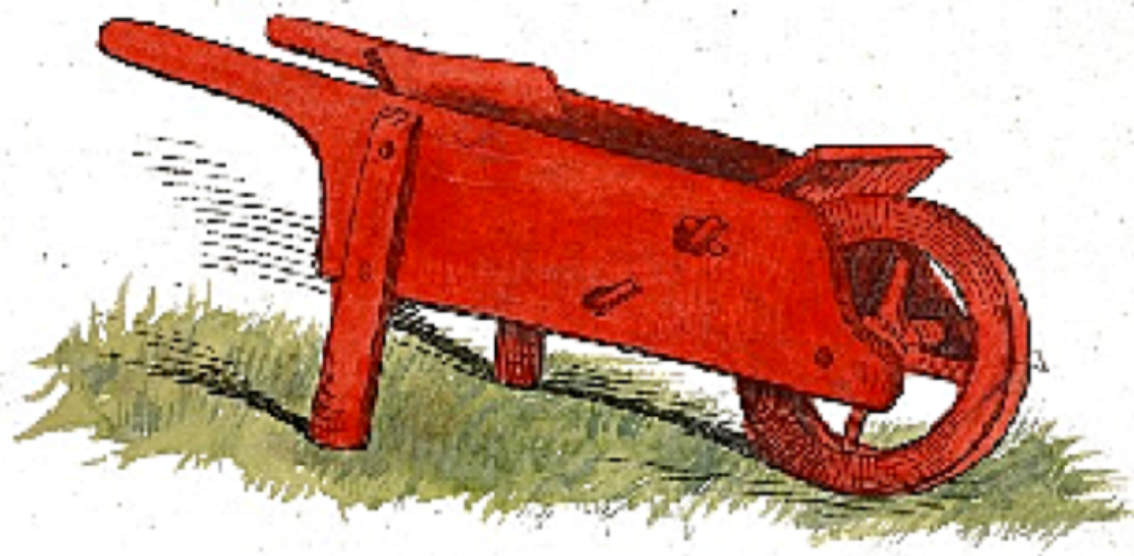
pypi package **0.7.0** Anaconda Cloud **0.7.0** build **passing** codecov **92%** downloads **177k**
downloads **86k total** license **MIT**

This open source Python library provide several solvers for optimization problems related to Optimal Transport for signal, image processing and machine learning.

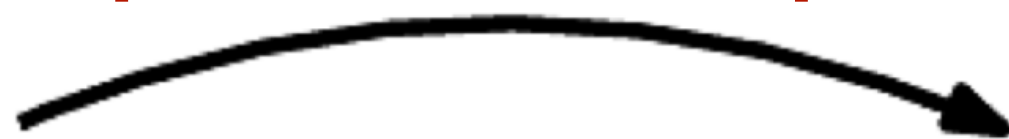
Website and documentation: <https://PythonOT.github.io/>

Source Code (MIT): <https://github.com/PythonOT/POT>

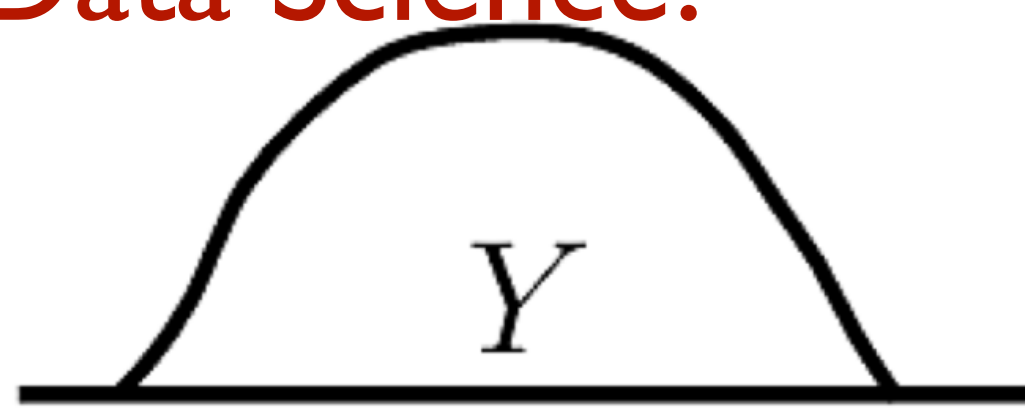
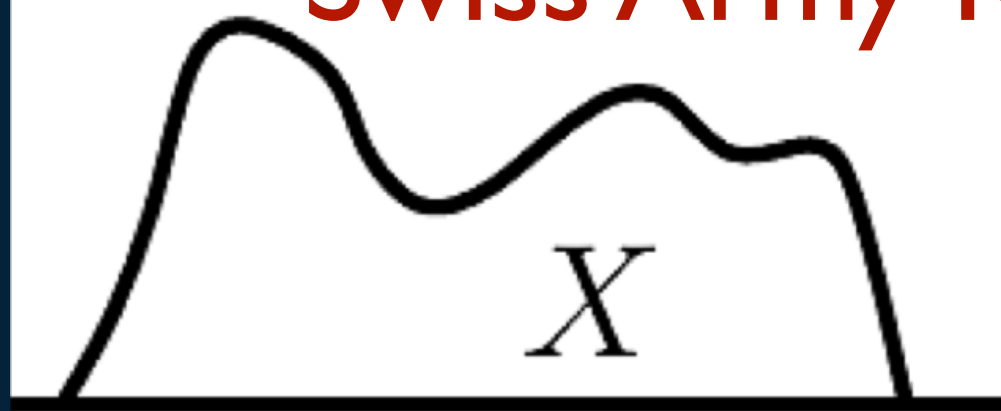




Optimal Transport



Swiss Army Knife for Data Science!



Thank You!