

RAPPORT: TRANSPORT OPTIMAL

ANDRIAMAHERISOA Liantsoa

Juin 2020

1 Introductions

Le principe du chemin le plus court guide la plupart des décisions dans les domaines de la vie et des sciences: lorsqu'un produit, une personne ou un seul élément d'information est disponible à un point donné et doit être accessible à envoyer à un point cible. Il convient de privilégier le moins d'effort possible. La théorie du transport optimal (TO) généralise cette intuition dans le cas où, au lieu de déplacer un seul élément à la fois, on se préoccupe du problème de déplacement simultanément de plusieurs éléments (ou une distribution continue de ceux-ci) d'un espace à un autre.

Introduite au 18ème siècle, ce n'est que récemment que cette théorie constitue un terrain fertile pour les chercheurs en informatique, en imagerie et plus généralement en sciences des données car cette théorie fournissait des outils très puissants pour étudier les distributions dans un contexte différent et plus abstrait, celui de comparer les distributions, facilement disponibles sous la forme d'un sac de caractéristiques ou d'informations. Maintenant que le TO s'est progressivement imposé comme un outil informatique, centré sur les applications à la science des données, notamment, l'imagerie et l'apprentissage automatique.

2 Fondements Théoriques

Ici, nous allons introduire les notions relatives d'appariements et de couplage optimaux entre vecteurs de probabilité (a, b), généralisant progressivement ce calcul pour un transporter des mesures discrètes (α, β) , pour couvrir enfin le cadre général des mesures arbitraires.

2.1 Histogrammes et Mesures

On utilise de manière interchangeable les termes histogramme et vecteur de probabilité pour tout élément $a \in \sum_n$ qui appartient au simplexe de probabilité.

$$\Sigma_n \stackrel{\text{def.}}{=} \left\{ \mathbf{a} \in \mathbb{R}_+^n : \sum_{i=1}^n \mathbf{a}_i = 1 \right\}.$$

Pour des mesures discrètes avec des poids a et les emplacements $x_1, \dots, x_n \in X$. Noté:

$$\alpha = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i},$$

où δ_x est le Dirac à la position x, intuitivement une unité de masse est concentrée à l'infini au position x. Cette mesure décrit une mesure de probabilité si, en plus, $\mathbf{a} \in \sum_n$ et plus généralement une mesure positive si tous les éléments du vecteur sont non négatifs.

Pour les mesures en générales: discrètes ou continues. On s'appuie sur l'ensemble des mesure de Randon $M(X)$ sur l'espace X. X doit être accompagné d'une distance d, car on ne peut accéder à une mesure qu'en l'intégrant contre les fonctions continues, notés $f \in C(X)$:

$$\int_X f(x) d\alpha(x) = \sum_{i=1}^n \mathbf{a}_i f(x_i).$$

2.2 Affectation et problème de Monge

Dans ce cas, on va avoir une matrice de coût

$$(C_{i,j})_{i \in [n], j \in [m]}$$

En supposant $n = m$, le problème d'affectation optimal cherche une bijection σ dans l'ensemble $\text{Perm}(n)$ des permutations d'éléments. On cherche alors à faire:

$$\min_{\sigma \in \text{Perm}(n)} \frac{1}{n} \sum_{i=1}^n C_{i,\sigma(i)}.$$

Une méthode naïve est d'évaluer cette fonction coût en utilisant toutes les permutations dans l'ensemble $\text{Perm}(n)$. Cependant, cet ensemble a une taille $n!$ Qui est énorme même pour des petits n. On va devoir utiliser des algorithmes efficaces pour optimiser le calcul. Le problème d'affectation optimale peut avoir plusieurs solutions optimales, comme par exemple dans le cas de deux éléments identiques, on peut avoir deux affectations possibles qui sont tout les deux optimales.

Dans le cas de mesures discrètes, le problème Monge cherche une carte qui associe à chaque point x_i un seul point y_j et qui va pousser la masse de α vers la masse de β . Une telle carte T doit vérifier:

$$\forall j \in [m], \quad \mathbf{b}_j = \sum_{i: T(x_i)=y_j} \mathbf{a}_i,$$

Dans le cas de mesures arbitraires (α, β) , supposant deux espaces (X, Y) peuvent être liés par une carte $T: X \rightarrow Y$ qui minimise:

$$\min_T \left\{ \int_X c(x, T(x)) d\alpha(x) : T_\# \alpha = \beta \right\}.$$

La contrainte $T_\# \alpha = \beta$ veut dire que T pousse en avant la masse de α vers la masse de β , en utilisant la méthode de push-forward. La nature combinatoire du problème dans le cas discret et la complexité de la contrainte $T_\# \alpha$ qui est convexe font que le problème de Monge est très difficile.

2.3 La formulation relaxée de Kantorovich

Puisque le problème de Monge est trop ardu, il faut nous résigner à en résoudre un plus simple. C'est la stratégie proposée par Kantorovich. Plutôt que d'exiger que l'élément x d'une distribution α soit transporté de manière déterministe vers un point T(x) dans la distribution cible β , Kantorovich propose d'autoriser la répartition des éléments sur différents points dans la cible. Un plan de transport γ entre deux distributions discrètes α et β . La masse d'un point dans α peut être répartie sur plusieurs points dans β . $\gamma(x,y)$ est une distribution de probabilités jointe qui indique la proportion d'élément au voisinage de x dans α que l'on va transporter au voisinage de y dans β . Pour que ce plan de transport soit compatible avec α et β il faut exiger les contraintes de marginalité. Cette flexibilité est alors codée en utilisant, à la place d'une permutation σ ou une carte T, une matrice de couplage $P \in \mathbb{R}(n \times m)_+$, où $P_{i,j}$ décrit la quantité de masse circulant de i vers j, ou de la masse trouvée en x_i vers y_j . Les couplages admissibles admettent une caractérisation beaucoup plus simple que les cartes Monge, où, on utilise la matrice vectorielle suivante:

$$\mathbf{P} \mathbf{1}_m = \left(\sum_j P_{i,j} \right)_i \in \mathbb{R}^n \quad \text{and} \quad \mathbf{P}^T \mathbf{1}_n = \left(\sum_i P_{i,j} \right)_j \in \mathbb{R}^m.$$

L'ensemble des matrices $U(a,b)$ est borné et défini par $n+m$ contraintes d'égalité, et est donc un polytope convexe (la coque convexe d'un ensemble fini de matrices). De plus, alors que la formulation de Monge était intrinsèquement asymétrique, la formulation détendue de Kantorovich est toujours symétrique, dans le sens où un couplage P est dans $U(a,b)$ si et seulement si P^T est dans $U(b,a)$, d'où:

$$L_C(a, b) \stackrel{\text{def}}{=} \min_{P \in U(a,b)} \langle C, P \rangle \stackrel{\text{def}}{=} \sum_{i,j} C_{i,j} P_{i,j}.$$

C'est une problème linéaire, et comme c'est généralement le cas avec de telles problèmes, ses solutions optimales ne sont pas forcément uniques.

Matrices de permutation comme couplages: Pour une permutation $\sigma \in \text{Perm}(n)$, on écrit P_σ pour la matrice de permutation correspondantes:

$$\forall (i, j) \in [n]^2, \quad (P_\sigma)_{i,j} = \begin{cases} 1/n & \text{if } j = \sigma_i, \\ 0 & \text{otherwise.} \end{cases}$$

Pour le cas des mesures discrètes, nous stockons dans la matrice C tous les coûts par paire entre les points dans les supports de α, β , à savoir $C_{i,j} = c(x_i, y_j)$, pour définir:

$$L_C(\alpha, \beta) \stackrel{\text{def}}{=} L_C(a, b).$$

Par conséquent, la formulation de Kantorovich du TO entre les mesures discrètes est la même que le problème entre leurs vecteurs de probabilité de poids a, b sauf que la matrice de coût C dépend du support de α et β . Dans le cas des mesures arbitraires, L_c est étendu à des mesures arbitraires en considérant les couplages $\pi \in M_+^1(X \times Y)$ qui sont des distributions conjointes sur l'espace produit. Le cas discret est une situation particulière où l'on impose que cette mesure de produit soit de la forme $\pi = \sum_{i,j} P_{i,j} \delta(x_i, y_j)$. Dans le cas général, la contrainte de conservation de la masse doit être réécrite en tant que contrainte marginale sur la distribution de probabilité conjointe.

$$U(\alpha, \beta) \stackrel{\text{def}}{=} \left\{ \pi \in M_+^1(X \times Y) : P_{X\#} \pi = \alpha \quad \text{and} \quad P_{Y\#} \pi = \beta \right\}.$$

2.4 Propriétés métriques du TO : Wasserstein

Pour définir une notion de distance entre deux distributions α et β , on utilise la distance de p -Wasserstein entre α et β . Pour la fonction de coût de transport C on choisit $C(x,y) = D^p(x,y)$ où $D(x,y)$ est une distance entre x et y et où $p \geq 1$.

$$W_p(\alpha, \beta) \equiv [L_D^p(\alpha, \beta)]^{1/p} \text{ pour } p \geq 1$$

On doit vérifier: $W_p(\alpha, \gamma) \geq W_p(\alpha, \beta) + W_p(\beta, \gamma)$. Dans certains cas spécifiques on peut néanmoins W^p explicitement:

- distance entre deux distributions $\alpha = \delta_x$ et $\beta = \delta_y$, on vérifie immédiatement que $W_p(\delta_x, \delta_y) = D(x,y)$

- dans le cas $p = 2$, si l'on définit α' et β' comme les versions centrées de moyennes nulles α et β , qu'on définit m_α et m_β comme les moyennes respectives de α et β alors on a la décomposition $W_2(\alpha, \beta)^2 = W_2(\alpha', \beta')^2 + \|m_\alpha - m_\beta\|^2$.

- dans le cas $p = 2$ avec deux gaussiennes $\alpha = N(m_\alpha, \Sigma_\alpha)$ et $\beta = N(m_\beta, \Sigma_\beta)$, on a: $W_2(\alpha, \beta)^2 = \|m_\alpha - m_\beta\|^2 + B(\Sigma_\alpha, \Sigma_\beta)$ où B est une métrique entre matrices de covariances que l'on sait calculer explicitement.

2.5 Problème Dual

Le problème de Kantorovich est un problème de minimisation convexe avec contraintes et, par conséquent, il peut naturellement être associé à un problème dual, qui est un problème de maximisation concave avec contraintes. Le problème de Kantorovich admet:

$$L_C(a, b) = \max_{(f,g) \in R(C)} \langle f, a \rangle + \langle g, b \rangle,$$

où, l'ensemble des variables duals est:

$$R(C) \stackrel{\text{def}}{=} \{(f, g) \in \mathbb{R}^n \times \mathbb{R}^m : \forall (i, j) \in [n] \times [m], f_i \oplus g_j \leq C_{i,j}\}.$$

appelées potentiels de Kantorovich. Ce résultat est une conséquence directe du résultat plus général sur la forte dualité des programmes linéaires. Dédoublons notre résultat en utilisant la dualité Langrangienne. D'où:

$$\min_{P \geq 0} \max_{(f,g) \in \mathbb{R}^n \times \mathbb{R}^m} \langle C, P \rangle + \langle a - P \mathbf{1}_m, f \rangle + \langle b - P^T \mathbf{1}_n, g \rangle.$$

En échangeant le min et le max, on obtient:

$$\max_{(f,g) \in \mathbb{R}^n \times \mathbb{R}^m} \langle a, f \rangle + \langle b, g \rangle + \min_{P \geq 0} \langle C - f \mathbf{1}_m^T - \mathbf{1}_n g^T, P \rangle.$$

On en conclut grâce à ce remarque:

$$\min_{P \geq 0} \langle Q, P \rangle = \begin{cases} 0 & \text{if } Q \geq 0, \\ -\infty & \text{otherwise} \end{cases}$$

de sorte que la contrainte se lit

$$C - f \mathbf{1}_m^T - \mathbf{1}_n g^T = C - f \oplus g \geq 0.$$

Ainsi, la relation d'optimalité primale-dual pour le lagrangien permet de localiser le support du plan de transport optimal.