

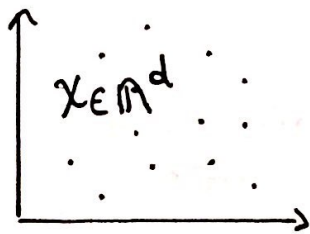
Metric learning (1)

Le choix de la similarité / distance est cruciale pour les performances.

=> Comment apprendre une similarité appropriée pour une tâche donnée ?

Le **metric learning** a pour objectif de **construire une métrique de distance spécifiques à une tâche donnée** à partir de données supervisées. La distance apprise peut ensuite être utilisée pour \neq tâches (ex: KNN, clustering...)

ex: vérification faciale.



Recette de metric learning :

1) choisir une **distance paramétrique** en amont

$D_H(x, x')$ paramétrisé par une matrice H

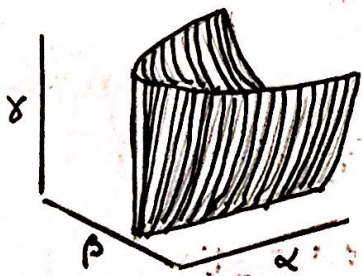
2) Collecter des **jugements** sur des paires d'images (couple (S, D) , triplet (S, D, R))

3) Estimer les paramètres (métrique qui s'accorde le + avec les jugements).

problème d'optimisation

$$H^* = \underset{H}{\operatorname{argmin}} \left[\underbrace{\ell(H, S, D, R)}_{\text{loss}} + \lambda \underbrace{\operatorname{reg}(H)}_{\substack{\text{regul} \\ (\text{Frobenius}, \text{L}_1 \text{ norm}, \text{Trace})}} \right]$$

En général, on utilise la **distance de Mahalanobis**.



$$D_H(x, x') = \sqrt{(x - x')^T H (x - x')}$$

pseudo-distance

H est PSD
 $H \in \mathbb{R}^{d \times d}$

c'est une distance euclidienne après transformation linéaire

L'ensemble des matrices PSD est un **cône convexe** donc cela va faciliter les problèmes d'optimisation.

Si $\operatorname{rg}(H) = R \leq d$ alors $L \in \mathbb{R}^{R \times d} \Rightarrow$ **réduction de dimension**.
On a $x^T H x \geq 0 \forall x \in \mathbb{R}^d$ et $H = L L^T$ (Choleski)

* Pairwise constraints.

=> tâche : clustering

Problème d'optimisation

$$\max_{H \in S_+^d} \sum_{(x_i, x_j) \in D} D_H(x_i, x_j)$$

concave car composée de racine + linéaire

$$\text{sous contraintes } \sum_{(x_i, x_j) \in S} D_H^2(x_i, x_j) \leq 1$$

change rien au pb.

convexe car linéaire en H

=> toujours faisable ($\Pi = 0$)

=> résolu avec une méthode de gradient projeté

complexité $O(d^2)$ si projeté sur les contraintes de distance.

$O(d^3)$ si = sur S_+^d

* Triplet constraints

=> tâche : recherche d'information.

norme de Frobenius (regul)

$$\min_{H \in S_+^d, \xi \geq 0} \|H\|_F^2 + \lambda \sum_{i,j,k} \xi_{ij,k}$$

Formulation convexe

$$\text{sous contraintes } D_H^2(x_i, x_k) - D_H^2(x_i, x_j) \geq 1 - \xi_{ij,k}$$

1 contrainte par triplet.

=> on veut que les x_i soient plus près de x_j que x_k si possible avec une certaine marge. Si on n'arrive pas à avoir cette marge, on va payer un coût dans la fonction objectif.

* Large Margin Nearest Neighbor

=> tâche : RNN classification

trade off parameter $\mu \in [0, 1]$

Formulation convexe

$$\min_{H \in S_+^d, \xi \geq 0} (1-\mu) \sum_{(x_i, x_j) \in S} D_H^2(x_i, x_j) + \mu \sum_{i,j,k} \xi_{ij,k}$$

$$\text{sous contraintes } D_H^2(x_i, x_k) - D_H^2(x_i, x_j) \geq 1 - \xi_{ij,k}$$



on veut rapprocher les points d'une même classe et "éjecter" les autres avec une marge de sécurité.

Netic learning (2)

Extension non-linéaire

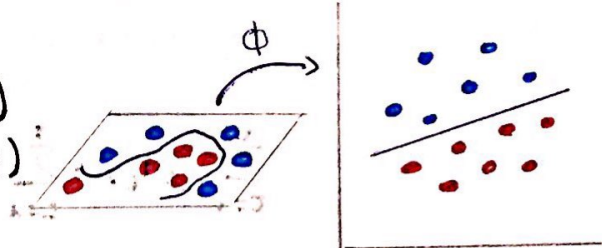
* Approche par noyau

on fixe une transformation riche $K(x, x') = \langle \phi(x), \phi(x') \rangle$
 \Rightarrow on apprend une métrique linéaire dans l'espace noyau.

distance de Mahalanobis

$$D_H^2(\phi_i, \phi_j) = (\phi_i - \phi_j)^T H (\phi_i - \phi_j) \\ = (\phi_i - \phi_j)^T L^T L (\phi_i - \phi_j)$$

en posant $L^T = \Phi U^T$



$$D_H^2(\phi(x), \phi(x')) = (R - R')^T H (R - R')$$

les paramètres ne dépendent pas de la taille de l'espace ϕ
 on peut calculer ces distances que par évaluation du noyau.

\Rightarrow problème de passage à l'échelle pour cette méthode. $O(n^2)$

* Approche peu apprentissage d'une transformation linéaire

on fixe un mapping ϕ

\Rightarrow on apprend une distance euclidienne dans ce mapping.

$$D_\phi(x, x') = \|\phi(x) - \phi(x')\|_2$$

Marche bien mais ne conserve pas la convexité.

Exemple d'application : image avec réseau siamois (on peut plager la fonction objectif au bout) par les.

* Apprentissage multiple de métrique locale (LMNN)

\Rightarrow grouper les données en C clusters

Pour chaque groupe, on apprend une métrique :

$$\min_{\xi \geq 0} (1 - \mu) \sum_{(x_i, x_j) \in S} D_{H_C(x_j)}^2(x_i, x_j) + \mu \sum_{i,j,k} \xi_{ijk}$$

$$\text{sous contraintes } D_{H_C(x_k)}^2(x_i, x_k) - D_{H_C(x_j)}^2(x_i, x_j) \geq 1 - \xi_{ijk}$$

garde la convexité, mais overfitting.

Large-Scale Netric learning

Deux problèmes:

- Large datasets $\Rightarrow n$ grand ($O(n^2)$, $O(n^3)$)
- Grande dimension $\Rightarrow d$ grand ($O(d^2)$, $O(d^3)$)

* Cas de n large

1- Online learning

\Rightarrow des données arrivent séquentiellement (on a pas besoin de gérer tous les points à la fois)

On cherche à minimiser le regret:

$$\sum_{t=1}^T \ell_t(H_t) - \sum_{t=1}^T \ell_t(H^*) \leq f(T) \quad (T \text{ étapes})$$

H^* optimal

(A quel point j'ai été moins bon que si j'avais travaillé avec toutes les données en batch).

\Rightarrow On peut montrer qu'on a des garanties: pour certains algorithmes, le regret est en \sqrt{T} ou $\log(T)$

Algorithme OASIS.

A chaque étapes t ,

$$H^t = \underset{H, \xi}{\operatorname{argmin}} \frac{1}{2} \|H - H^{t-1}\|_F^2 + C \xi$$

paramètre d'agressivité

sous contraintes $1 - S_H(x_i, x_j) + S_H(x_i, x_k) \leq \xi \quad \xi \geq 0$

puis $S_H(x, x') = x^T H x' \quad C \geq 0$

Solution donnée par $H^t = H^{t-1} + \beta V$

pas
(si C faible, β faible)

$$V = x_i(x_j - x_k)^T$$

à quel point on est prêt à changer la métrique si on n'est trompé que un triplet

2- Méthodes de gradient stochastique

$$\min_H \frac{1}{|R|} \sum_{(x_i, x_j, x_k) \in R} \ell(H, x_i, x_j, x_k)$$

\Rightarrow utilise un échantillon aléatoire pour estimer le gradient

\Rightarrow mieux que le gradient classique quand n est grand.

\Rightarrow peut être combiné avec de l'optimisation distribuée.

* Cas de d large

1- Apprendre une matrice diagonale

=> on apprend d paramètres

=> on apprend qu'une pondération des features (pas très riche).

2- Apprendre la métrique après réduction de dimension (PCA)

=> on perd de l'information

=> difficile à interpréter.

3- Notions de décompositions

Décomposition de faible rang

$$M = L^T L \quad L \in \mathbb{R}^{r \times d}$$

=> on apprend $r \times d$ paramètres

=> en général non convexe, on doit tuner r

Décomposition rang-1

$$M = \sum_{i=1}^K w_i b_i b_i^T$$

=> on apprend K paramètres

=> on doit choisir un bon set de base.

Régularisation

=> pour éviter l'overfitting.

* norme de Frobenius $\|M\|_F^2 = \sum_{i,j=1}^d m_{ij}^2$ convexe, facile à optimiser

* norme mixte $L_{2,1}$ $\|M\|_{2,1} = \sum_{i=1}^d \|m_i\|_2$ amène certaines colonnes à 0 (sélection features).
convexe

* norme trace $\|M\|_* = \sum_{i=1}^d \sigma_i(M)$ favorise les matrices de rang faible (réduction dimension), convexe.

-
- Composantes qu'on retrouve dans plein d'algo
 - Solutions existes pour passer à l'échelle
 - Gap entre recherche et software

=> librairie **metric-learn** sur Python.
compatible avec **scikit-learn**.