

# Article InfoGraph (ICLR 2020)

F.Y. Sun, J. Hoffmann, V. Verra, J. Tiang

National Taiwan Univ, MQILA, Aalto Univ, Harvard, HEC Montréal, CIFAR

Apprentissage de rep. d'un graphe dans un ensemble

Approche non supervisée (mais peut être étendue au cas  $\frac{1}{2}$  sup)

Basé sur Deep InfoMax (DIM) de Hjelm et al. (ICLR 2019)

Préliminaires : papiers DIM et Norouzi et al. (NIPS 2016)

DIM maximise la MI entre l'entrée ( $X$ ) et la représentation latente  $E_\psi(X)$

L'objectif est alors de maximiser  $I(X, E_\psi(X))$

↳ paramètre du modèle

tout en imposant certaines contraintes sur la loi marginale  $\mathbb{P}_{\psi, \mathbb{P}}$

$\mathbb{P}$  est la distribution empirique de  $X$

Or, maximiser directement  $I$  est difficile et pas forcément nécessaire -

On peut passer par des estimateurs, ce qui nous permet de recourir à des fonctions de la famille des  $f$ -divergences en plus de la KL sur laquelle se base la MI

cf. papier de Norouzi et al.

En effet, on a :

$$D_f(P \parallel Q) \geq \sup_{T \in \mathcal{T}} (\underbrace{\mathbb{E}_{x \sim P}[T(x)]}_{\text{terme 1}} - \underbrace{\mathbb{E}_{x \sim Q}[f^*(T(x))]}_{\text{terme 2}}) \quad (\text{eq. 4})$$

ce qui correspond à une borne inférieure de la divergence

où :  $f^*$  est la fonction duale de  $f$  (conjuguée)

$\mathcal{T}$  est une classe de fonctions arbitraires

Il peut être prouvé que la borne est très proche pour  $T^* = f'(\frac{f(x)}{q(x)})$

Si  $f$  = Entropie Shannon (JS) alors on peut calculer :

$$D_{JS} = \frac{1}{2} \int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx$$

Le "générateur"  $f$  est :  $f(u) = -(u+1) \log \frac{1+u}{2} + u \log u$

$$\text{et } T^* \text{ est : } T^*(x) = \log \frac{2p(x)}{p(x)+q(x)}$$

DIM

Reprenons le terme 1 de l'eq. 4 pour le cas où  $p := p(x, y)$  et  $q := p(x) p(y)$

$$\mathbb{E}_{x \sim P} [T^*(x, E_\psi(X))] = \mathbb{E}_P \left[ \underbrace{-\log \frac{2p(x, y)}{p(x, y) + p(x)p(y)}}_{\text{A}} \right]$$

$\uparrow$  objet  $x$        $\uparrow$  rep. latente  $y$

$$\text{A se réécrit : } -\log \frac{p(x, y) + p(x)p(y)}{2p(x, y)}$$

$$\propto -\log(1 + e^{-M}) \quad \text{si on pose : } M := \log \left( \frac{p(x, y)}{p(x)p(y)} \right)$$

$$= -\text{sp}(-M) \quad \text{avec } \text{sp}(z) = \log(1 + e^z)$$

Même travail avec le terme 2 de l'eq. 4 :

$$\mathbb{E}_{x \sim Q} [f^*(T^*(x, E_\psi(x)))] = \mathbb{E}_Q \left[ -\log \left( 2 - \exp \left( \log \frac{2p(x, y)}{p(x, y) + p(x)p(y)} \right) \right) \right]$$

$\uparrow$   $T^*$

car  $f^*(z) = -\log(2 - \exp(z))$  pour JS

$$\textcircled{B} \text{ se réécrit : } -\log\left(2 - \frac{2p(x,y) + 2p(x)p(y) - 2p(x)p(y)}{p(x,y) + p(x)p(y)}\right)$$

$$= -\log\left(\frac{2p(x)p(y)}{p(x,y) + p(x)p(y)}\right)$$

$$\propto -\log \frac{1}{e^M + 1} = \log(e^M + 1) = \text{sp}(M)$$

On peut enfin réécrire la borne en suivant l'estimateur JS :

$$\hat{I}_{\theta}^{(JS)}(X, E_{\psi}(X)) := \mathbb{E}_{\mathbb{P}}[-\text{sp}(-M(x, E_{\psi}(x)))] - \mathbb{E}_{\tilde{\mathbb{P}}}[\text{sp}(M(x', E_{\psi}(x)))]$$

$\tilde{\mathbb{P}} \times \mathbb{P}$   
 lecture "local" d'un objet distinct de  $x$       lecture "global" associée à  $x$

Une autre interprétation :

$$\hat{I}_{\theta}^{(JS)} := \mathbb{E}_{\mathbb{P}}[-\text{sp}(-M(x, E_{\psi}(x)))] - \mathbb{E}_{\tilde{\mathbb{P}} \times \mathbb{P}}[\text{sp}(M(x', E_{\psi}(x)))]$$

↑  
 « discriminateur »  
 { valeur ↑ si MI forte entre  $x$  et  $E_{\psi}(x)$   
 { valeur ↓ si MI faible entre  $x'$  et  $E_{\psi}(x)$

On remplace  $M$  par une fonction paramétrique (un réseau de neurones) :  $M_{\omega}$   
 donc  $\theta = \{\psi, \omega\}$       cf. Belghazi et al., cf. Pennington

Au final, on cherche  $\hat{\theta}$  t.q. :

$$\hat{\theta} = \arg \max_{\theta} \hat{I}_{\theta}^{(JS)}(X, E_{\psi}(X))$$

DIM repose sur le principe **InfoMax** (Linzen 1988; Bell & Sejnowski, 1995) qui cherche à maximiser la MI entre l'entrée  $X$  et la sortie d'une fonction comme un réseau calculant une rep. latente  $Z$ , ici  $E_{\psi}(X)$  en suivant un principe d'apprentissage adverse (adversarial training) inspiré des GAN et AAE (Makhadmeh et al., 2015) pour contraindre l'espace latent. L'adv. training agit comme une forme de régularisation dans l'espace appris par  $E$ . L'important est d'être capable d'échantillonner :

$\mathbb{P} : (x, E_{\psi}(x))$ , avec  $x$  pouvant être "une partie" (feature local) de l'objet ; et  $E_{\psi}(x)$  la représentation globale de  $X$

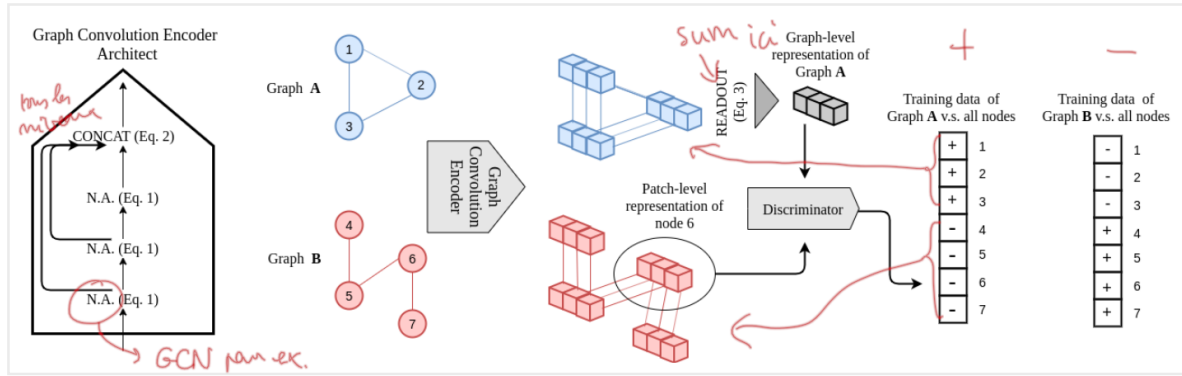
$\tilde{\mathbb{P}} \times \mathbb{P} : (x', E_{\psi}(x))$ , avec  $x'$  pouvant être une partie d'un autre objet  $X'$  = exemples adverses

Plusieurs remarques :

- $M$ , la fonction qui calcule la MI, est approximée par une famille de fonctions paramétriques et l'occurrence un discriminateur  $D_{\omega} : \mathbb{P} \times \mathbb{P} \rightarrow \mathbb{R}$



a scheme for node classification problem is given below :



En mon sup, l'approche est utilisée pour avoir des rep. des graphes qui sont utilisés pour faire de la classification (entre 2 et 5 classes, jusqu'à 5000 graphes). La classif est faite avec LIBSVM en faisant varier C, mais pas d'info sur le kernel

## Deep Graph InfoMax (Velickovic et al., ICLR 2019)

Publié de manière concomitante.

DGI est aussi basé sur DIM pour maximiser la MI entre une représentation globale et une représentation locale au niveau du « patch » (init. pour les images).

Ça encourage donc une info globale présente un peu partout dans l'objet.

DGI est aussi une approche contrastive car on apprend à distinguer (classification) entre des paires (info locale, info globale) qui sont associées ou non au même objet. L'encodeur est basé sur un GCN pour apprendre des  $h_i \in \mathbb{R}^{F'}$ ,  $F' =$  taille espace latent

La représentation globale, appelée ici « summary vector », est calculée une fois encore avec une fonction appelée READOUT  $R$  :

$$\vec{s} = R(\mathcal{E}(X, A))$$

$\downarrow$   $\downarrow$   $\downarrow$   
 readout encodeur features

À noter que je n'ai pas trouvé comment est finalement calculé  $R$ ...

Comme proxy pour maximiser la MI, ils utilisent un discriminateur

$D: \mathbb{R}^F \times \mathbb{R}^F \rightarrow \mathbb{R}$  t.q.  $D(\vec{h}_i, \vec{s})$  est un score de proba :

( $\rightarrow$  élevé si  $i$  appartient au grapho résumé  
 $\rightarrow$  faible sinon (exemples négatifs))

Les exemples négatifs sont générés ici avec une procédure de corruption.

La fonction objectif est inspirée par DIM mais retranscrit par une cross entropie

$$\mathcal{L} = \frac{1}{N+M} \left( \sum_{i=1}^N \mathbb{E}_{(X_i, A)} [\log D(\vec{h}_i, \vec{s})] + \sum_{j=1}^M \mathbb{E}_{(\tilde{X}_j, \tilde{A})} [\log (1 - D(\vec{h}_j, \vec{s}))] \right)$$

$\uparrow$   $\uparrow$   $\uparrow$   $\uparrow$   
 pairs + pairs - pairs + pairs -

C'est très similaire à (1) + (2) utilisé dans InfoGraph.

Un lien est fait avec le GATN qui optimisent une fonction similaire.

Des dérivations théoriques montrent que optimiser (3) revient à optimiser la MI