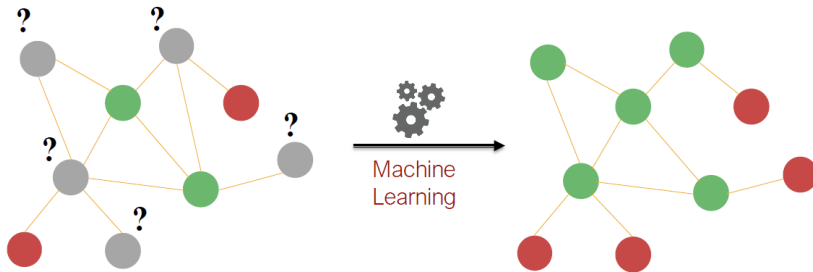


Apprentissage de représentation de graphes

Daniel Aloise <daniel.aloise@polymtl.ca>

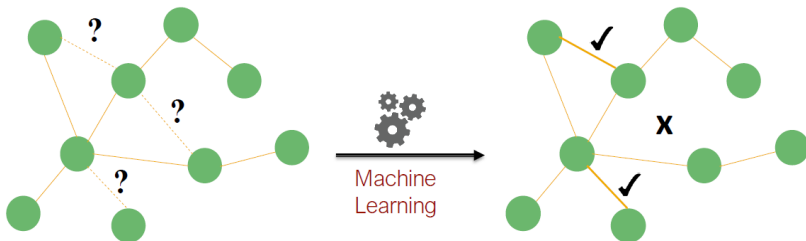
(basé sur les slides du professeur J. Leskovec, Stanford)

Example : classification collective



Source : Leskovec, 2021

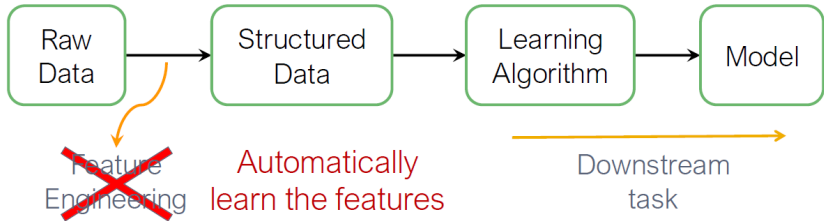
Exemple : prédiction des liens



Source : Leskovec, 2021

Génie d'attributs

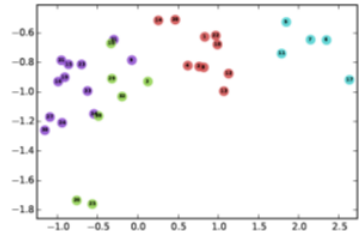
- Le pipeline de la fouille des données nécessite que les données soient préparées d'abord
 - Extraction d'attributs : **automatique vs. fait à la main.**



Source : Leskovec, 2021

Incorporation des sommets embedding

représentation vectorielle dont les valeurs à l'intérieur représentent l'information du mot



Intuition : recherche une incorporation des sommets dans un espace de dimension d de sorte que les sommets **similaires** dans le graphe soient proches les uns des autres.

Setup

V : ensemble de sommets

A : ensemble d'arrêtes

- Nous allons assumer un graphe $G = (V, A)$ où :
 - V est l'ensemble de sommets
 - A est une matrice d'adjacence (possiblement pondérée)
- Aucun attribut sur les sommets ou information supplémentaire n'est utilisée !
on n'a pas la topologie du graphe. ex: une personne on n'a pas son age, son poids, sa taille. On n'a pas considéré ces informations pour faire du embedding

Apprentissage de l'incorporation de sommets

Encodeur: Fonction qui prend un sommet U et qui va le transformer dans un vecteur z_u

- 1 Définition d'un **encodeur** (c'est-à-dire un *mapping* des sommets à l'espace d'incorporation)
- 2 Définition d'une **fonction de similarité** entre les sommets (c'est-à-dire une mesure de similarité dans le graphe d'origine).
Fonction de similarité: pouvoir comparer la similarité des sommets
- 3 Optimisation des paramètres de l'encodeur pour que :

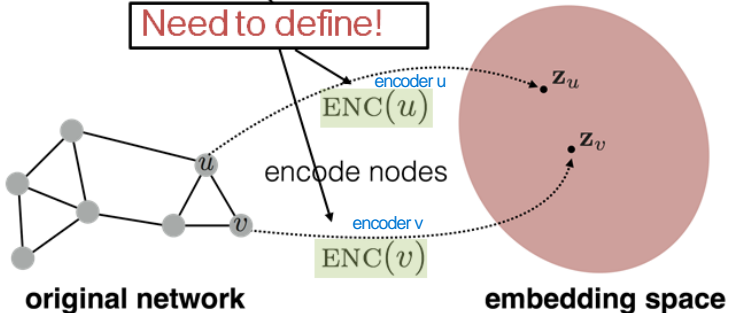
$$\text{similarite}(u, v) \approx z_v^T z_u$$

similarité des sommets u et v = produit scalaire des représentations est de 0.8 par exemple

Apprentissage de l'incorporation de sommets

Goal: $\text{similarity}(u, v) \approx \mathbf{z}_v^\top \mathbf{z}_u$

Need to define!



Source : Leskovec, 2021

ex: si on veut que chaque sommet soit représenté par 10 valeurs numériques, on aura une matrice qui aura 10 lignes et qui a autant de colonnes qu'il y a de noeuds dans notre graphe. Pour chaque colonne, la 1ere colonne est la représentation du sommet sommet 1, la 2e colonne représente le 2e sommet etc...

Codage simple

Z : matrice qui contient les valeurs réelles

$$ENC(v) = Zv$$

$$\underline{Z} \in \mathbb{R}^{d \times |V|}$$

d : nb de valeurs pour représenter les sommets

matrice où chaque colonne correspond à l'incorporation d'un sommet dans l'espace \mathbb{R}^d .

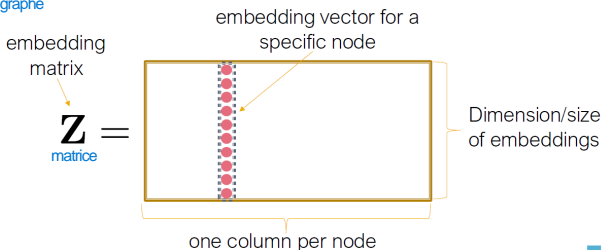
[ce qui nous apprenons]

$$\underline{v} \in \mathbb{I}^{|V|}$$

vecteur qui fait la taille de notre vocabulaire, donc le nb de sommets dans notre graphe

vecteur où la seule position non-nulle correspond à l'indice du sommet v .

v vecteur



Codage simple

- Bref, chaque sommet se voit attribuer un vecteur d'incorporation unique.
 - plusieurs méthodes : node2vec, DeepWalk, LINE, etc.
- La distinction clé entre les méthodes de codage simple est la façon dont elles définissent **la similarité** entre les sommets dans le graphe d'origine.
- Deux sommets devraient-ils avoir des incorporations similaires s'ils :
 - sont connectés ?
 - ont des voisins communs ?
 - ont des « rôles structurels » similaires ? .. si les sommets contrôlent des informations

Similarité basée sur la matrice d'adjacence

Problème: Approche locale, on considère uniquement les voisins directs. Si un noeud est lié, alors ils sont similaires. Si un noeud n'est pas lié ils sont différents, peu importe de la distance.

- La **fonction de similarité** est juste le poids de l'arête entre les sommets u et v dans le graphe d'origine.
- Intuition : les produits scalaires entre les incorporations de sommets se rapprochent de l'existence des arêtes.

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \| \mathbf{z}_u^\top \mathbf{z}_v - \mathbf{A}_{u,v} \|^2$$

Loss (l'erreur) \nearrow perte (ce qui on veut minimiser)

Similarité dans l'espace d'incorporation \nwarrow

Matrice d'adjacence du graphe \swarrow les arêtes (arc) dans notre graphe

si les arêtes sont reliés, les vecteurs sont similaires, sinon ils ne sont pas similaires

Similarité basée sur la matrice d'adjacence

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \|z_u^T z_v - A_{u,v}\|^2$$

- L'objectif est d'obtenir la matrice d'incorporation $Z \in \mathbb{R}^{d \times |V|}$ que minimize la perte \mathcal{L} .
 - Option 1 : Utiliser la méthode de descente de gradient stochastique.
 - Option 2 : Utiliser des méthodes de décomposition de matrices (e.g. SVD).

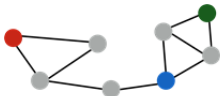
Problème: On considère uniquement les voisins directs!

Similarité basée sur la matrice d'adjacence

$$\mathcal{L} = \sum_{(u,v) \in V \times V} \|z_u^T z_v - A_{u,v}\|^2$$

- Inconvénients

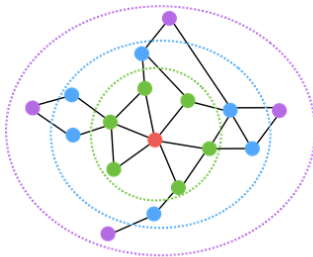
- Temps $O(|V|^2)$ (considère tous les pairs de sommets)
 - Nous pouvons optimiser à $O(|E|)$ - réf. sur Moodle
- L'optimisation considère uniquement des connections locales et directes.



e.g. le sommet **bleu** est plus similaire au sommet **vert** qu'au sommet **rouge**, bien qu'ils ne soient pas directement connectés.

Similarité *multi-hop*

- Le dernier inconvénient peut être contourné partiellement en utilisant de **similarités multi-hops**



- Rouge:** sommet cible
 - Vert:** voisinage 1-hop
 - \mathbf{A} (i.e., matrice d'adjacence)
 - Bleu:** voisinage 2-hop
 - \mathbf{A}^2
 - Mauve:** voisinage 3-hop
 - \mathbf{A}^3
- sommets situés à une puissance 2, peut être 3 ou 4

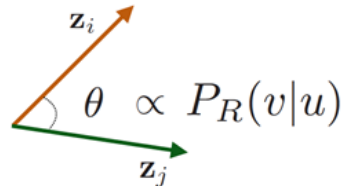
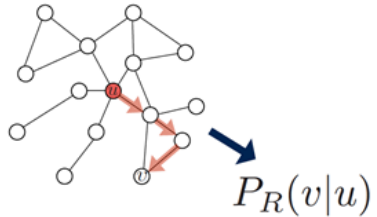
- Idée : $\min \mathcal{L} = \sum_{(u,v) \in V \times V} \|z_u^T z_v - A_{u,v}^k\|^2$, où k est le niveau de connectivité considéré.

Incorporation par marches aléatoires

$\mathbf{z}_u^\top \mathbf{z}_v \approx$
 Probabilité que u et v ^{apparaissent en mm temps} co-occurent dans une marche aléatoire sur le graphe.

Incorporation par marches aléatoires

- 1 Estimez la probabilité de visiter le sommet v lors d'une marche aléatoire à partir du sommet u en utilisant une stratégie de marche aléatoire R .
- 2 Optimisez l'incorporation pour encoder ces statistiques de marche aléatoire.



angle des 2 embedding doit être proportionnel à la proba de voir le sommet v sachant qu'on connaît le sommet u

Pourquoi des marches aléatoires ?

AVANTAGES

Expressivité : définition stochastique **flexible** de la similarité entre les sommets qui intègre à la fois des informations de voisinage locales et d'ordre supérieur.

Efficience : efficace efficace car pas besoin de prendre en compte toutes les paires de sommets lors de l'entraînement ; il suffit de considérer les paires qui coexistent lors de marches aléatoires.

Optimisation de marches aléatoires

- 1 Faite de courtes marches aléatoires à partir de chaque sommet du graphe en utilisant une stratégie R .
- 2 Pour chaque sommet u collectez $N_R(u)$, c.-à.d. l'ensemble* de sommets visités par des marches aléatoires débutant au sommet u .
- 3 Maximisez :

$$\sum_{\substack{\text{somme pr} \\ \text{chacun des} \\ \text{sommets de} \\ \text{départ}}} \sum_{v \in N_R(u)} \log(P(v|z_u)) \Rightarrow \text{objectif de maximum likelihood}$$

z_u : embedding du sommet u , le sommet de départ
 ressemblance maximale

où $P(v|z_u)$ est noté comme la probabilité (prédite) de visiter le sommet v sur des marches aléatoires à partir du sommet u .

* $N_R(u)$ peut contenir des sommets répétés

Optimisation de marches aléatoires

- De façon équivalente, on peut minimiser :

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} -\log(P(v|z_u))$$

- Ensuite, on paramètre $P(v|z_u)$ avec **softmax**.

$$P(v|z_u) = \frac{\exp(z_u^T z_v)}{\sum_{n \in V} \exp(z_u^T z_n)}$$

Proba d'obtenir le sommet v
à partir du sommet u

Pourquoi softmax ?

Nous voulons que le sommet v soit le plus similaire au sommet u parmi tous les sommets du graphe.

Intuition : $\sum_i \exp(x_i) \approx \max_i \exp(x_i)$

Optimisation des marches aléatoires

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} - \log \left(\frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)} \right)$$

somme pour tous les sommets u

Somme pour les sommets v visités à partir de u

Probabilité prédite pour que u et v co-occurrent ensemble dans une marche aléatoire

L'optimization de l'incorporation par marches aléatoires = trouver des incorporations \mathbf{z}_u qui minimisent \mathcal{L} .

Optimisation des marches aléatoires

Pourtant faire cela naïvement coûte trop cher!!

$$\mathcal{L} = \sum_{u \in V} \sum_{v \in N_R(u)} -\log \left(\frac{\exp(\mathbf{z}_u^\top \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^\top \mathbf{z}_n)} \right)$$

Les deux sommes
imbriquées entraîne une
complexité de $O(|V|^2)$

Le terme de normalisation du softmax est le coupable... peut-on l'approximer?

Optimisation des marches aléatoires

- Solution : *échantillonnage négatif*

$$\frac{\exp(\mathbf{z}_u^T \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^T \mathbf{z}_n)}$$

$$\approx \log(\sigma(\mathbf{z}_u^T \mathbf{z}_v)) - \sum_{i=1}^k \log(\sigma(\mathbf{z}_u^T \mathbf{z}_{\bar{n}_i}))$$

- $\sigma(\cdot)$ est la fonction sigmoïde - rend chaque terme entre 0 et 1.
- L'ensemble \bar{n} est sélectionné à partir des **échantillons négatifs** = sommets pas atteints par les marches aléatoires.
- Au lieu de normaliser par rapport à tous les sommets, on normalise juste par rapport à k échantillons négatifs \Rightarrow **plus rapide**

Optimisation des marches aléatoires

- Solution : *échantillonnage négatif*

$$\frac{\exp(\mathbf{z}_u^T \mathbf{z}_v)}{\sum_{n \in V} \exp(\mathbf{z}_u^T \mathbf{z}_n)}$$

$$\approx \log(\sigma(\mathbf{z}_u^T \mathbf{z}_v)) - \sum_{i=1}^k \log(\sigma(\mathbf{z}_u^T \mathbf{z}_{\bar{n}_i}))$$

- Les échantillons négatifs sont sélectionnés selon une probabilité proportionnelle à leur degrés.
- Deux considérations par rapport à k (# d'échantillons négatifs) :
 - Un k élevé donne des estimations plus robustes.
 - Un k élevé correspond aussi à un biais plus élevé sur les événements négatifs
 - En pratique, on utilise k entre 5 et 20 (très petite pour des graphes massives)



En résumé

- ① Faite de courtes marches aléatoires à partir de chaque sommet du graphe en utilisant une stratégie R .
- ② Pour chaque sommet u collectez $N_R(u)$, c.-à.d. l'ensemble de sommets visités par des marches aléatoires débutant au sommet u .
- ③ Maximisez :

$$\sum_{u \in V} \sum_{v \in N_R(u)} \log(P(v|z_u)) \Rightarrow \text{objectif de maximum likelihood}$$

Pour cela, vous pouvez utiliser une **descente de gradient stochastique**.

Descente de gradient stochastique

Au lieu d'évaluer les gradients sur tous les sommets, évaluez-les pour chaque sommet individuellement.

- Initialisez \mathbf{z}_u avec des valeurs aléatoires pour tous les sommets $u \in V$ chacune des colonnes initialisés aléatoirement
- Itérer jusqu'à la convergence : $\mathcal{L}^u = \sum_{v \in N_r(u)} -\log(P(v|\mathbf{z}_u))$
 - Échantillonnez un sommet u , ensuite pour tout v calculez le dérivé $\frac{\partial \mathcal{L}^u}{\partial \mathbf{z}_v}$
 - Pour tout v , calculez $\mathbf{z}_v \leftarrow \mathbf{z}_v - \eta \frac{\partial \mathcal{L}^u}{\partial \mathbf{z}_v}$

Comment devrait-on marcher aléatoirement ?

- Jusqu'à présent, nous avons décrit comment optimiser des incorporations des graphes étant donné une stratégie de marche aléatoire *R*.
- Quelles stratégies pouvons-nous utiliser pour exécuter ces marches aléatoires ?
 - Idée la plus simple : **Juste faire des marches aléatoires non-biaisées de taille fixe à partir de chaque sommet** (i.e., *DeepWalk* - Perozzi et al., 2013.)
 - *node2vec* marches aléatoires biaisées qui arbitrent entre une recherche en profondeur (vision globale) et une recherche en largeur (vision locale) - Grover and Leskovec, 2016.
 - D'autres possibilités existent encore.

Comment pouvons-nous utiliser les incorporations ?

z_i : représentation
vectorielle des points

Clustering : regroupement des sommets/points basés sur z_i

Classification des sommets : prédiction de l'étiquette du
sommets i basés sur z_i

Prédiction de liens : prédiction de l'arête (i, j)

- Nous pouvons dans ce cas enchaîner, sommer, multiplier ou prendre la différence entre les incorporations.

Fouille des graphes : classification des graphes

- Nous pouvons sommer ou calculer la moyenne des incorporations de tous les sommets du graphe.

Conclusion

- **Idée de base** : Incorporer des sommets de sorte que les distances dans l'espace d'incorporation reflètent les similarités entre les sommets dans le graphe d'origine.
- Différentes notions de similarités :
 - basées sur l'adjacence (c.-à.-d similaires si connectés).
 - définition multi-hop.
 - basées sur des marches aléatoires.

Conclusion

- Alors, quelle méthode devrait-on utiliser ?
- Aucune méthode ne l'emporte dans tous les cas
 - **node2vec** marche mieux pour la classification tandis que l'approche **multi-hop** marche mieux pour la prédiction de liens.
- Les approches par marches aléatoires sont normalement plus rapides.

- Problèmes en commun :

- Les méthodes de codage simple ne partagent aucun paramètre entre les sommets. peu importe la façon on représente un sommet, ça n'impacte pas la façon dont on représente les autres sommets, donc limitation

- Ils n'exploitent pas les attributs des sommets. Si 1 sommet avait des attributs, on a des caractéristiques sur cet individu mais on ne le prend pas en compte, donc limitation
- Possible solution : Graph Neural Networks (GNNs)

