

Web mining

Etienne G. Tajeuna

November 17, 2022



Plan

Introduction

Forage de structure web

Forage de contenu web

Introduction

- Le Web pourrait être vu comme un repertoire de données,

Introduction

- Le Web pourrait être vu comme un repertoire de données,
- Il est caractérisé par :

Introduction

- Le Web pourrait être vu comme un repertoire de données,
- Il est caractérisé par :
 - ▶ son grand volume de données : en perpétuelle augmentation, la quantité de données disponibles est immense.

Introduction

- Le Web pourrait être vu comme un repertoire de données,
- Il est caractérisé par :
 - ▶ son grand volume de données : en perpétuelle augmentation, la quantité de données disponibles est immense.
 - ▶ l'hétérogénéité des données : Les données présentes sont de types et natures diverses.

Introduction

- Le Web pourrait être vu comme un repertoire de données,
- Il est caractérisé par :
 - ▶ son grand volume de données : en perpétuelle augmentation, la quantité de données disponibles est immense.
 - ▶ l'hétérogénéité des données : Les données présentes sont de types et natures diverses.
 - ▶ l'interconnexion du contenu : Les liens d'une page vers d'autres pages créent une toile facilitant la navigation et la découverte des contenus.

Introduction

- Le Web pourrait être vu comme un repertoire de données,
- Il est caractérisé par :
 - ▶ son grand volume de données : en perpétuelle augmentation, la quantité de données disponibles est immense.
 - ▶ l'hétérogénéité des données : Les données présentes sont de types et natures diverses.
 - ▶ l'interconnexion du contenu : Les liens d'une page vers d'autres pages créent une toile facilitant la navigation et la découverte des contenus.
 - ▶ le bruit : Vu comme un média libre, toute personne est donc capable d'y poster une information qui pourrait être importante ou pas.

Introduction

- Le Web pourrait être vu comme un repertoire de données,
- Il est caractérisé par :
 - ▶ son grand volume de données : en perpétuelle augmentation, la quantité de données disponibles est immense.
 - ▶ l'hétérogénéité des données : Les données présentes sont de types et natures diverses.
 - ▶ l'interconnexion du contenu : Les liens d'une page vers d'autres pages créent une toile facilitant la navigation et la découverte des contenus.
 - ▶ le bruit : Vu comme un média libre, toute personne est donc capable d'y poster une information qui pourrait être importante ou pas.
 - ▶ sa vélocité : En perpétuelle croissance dû à l'ajout perpétuel de contenus par les utilisateurs.

Introduction

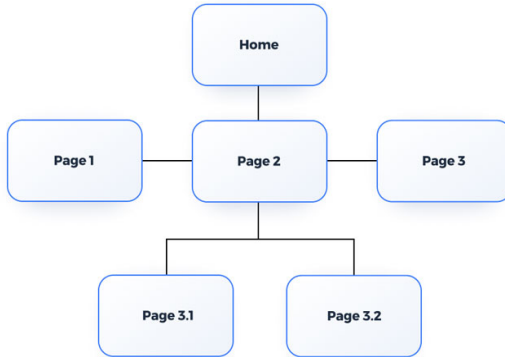
- Le Web pourrait être vu comme un repertoire de données,
- Il est caractérisé par :
 - ▶ son grand volume de données : en perpétuelle augmentation, la quantité de données disponibles est immense.
 - ▶ l'hétérogénéité des données : Les données présentes sont de types et natures diverses.
 - ▶ l'interconnexion du contenu : Les liens d'une page vers d'autres pages créent une toile facilitant la navigation et la découverte des contenus.
 - ▶ le bruit : Vu comme un média libre, toute personne est donc capable d'y poster une information qui pourrait être importante ou pas.
 - ▶ sa vélocité : En perpétuelle croissance dû à l'ajout perpétuel de contenus par les utilisateurs.
 - ▶ sa sociabilité : Les utilisateurs peuvent collaborer entre eux.

Introduction

- Bien qu'il soit vu comme un repertoire *fourre-tout*, où l'on peut incorporer tout type d'information, il n'en demeure pas moins que les liens entre les différentes pages de contenus constituent une structure topologique bien définie.

Introduction

- Bien qu'il soit vu comme un repertoire *fourre-tout*, où l'on peut incorporer tout type d'information, il n'en demeure pas moins que les liens entre les différentes pages de contenus constituent une structure topologique bien définie.



Introduction

- Ainsi présenté, le Web n'est plus seulement vu comme un repertoire *fourre-tout*, mais aussi une *toile* organisée où l'on pourrait naviguer.

Introduction

- Ainsi présenté, le Web n'est plus seulement vu comme un repertoire *fourre-tout*, mais aussi une *toile* organisée où l'on pourrait naviguer.
- La fouille de données a pour objectif l'extraction de connaissances à partir

Introduction

- Ainsi présenté, le Web n'est plus seulement vu comme un repertoire *fourre-tout*, mais aussi une *toile* organisée où l'on pourrait naviguer.
- La fouille de données a pour objectif l'extraction de connaissances à partir
 - ▶ de la structure formée par les liens entre pages/objets/communautés : Forage de structure.
Ce type de fouille vise essentiellement à analyser la structure du Web.

Introduction

- Ainsi présenté, le Web n'est plus seulement vu comme un repertoire *fourre-tout*, mais aussi une *toile* organisée où l'on pourrait naviguer.
- La fouille de données a pour objectif l'extraction de connaissances à partir
 - ▶ de la structure formée par les liens entre pages/objets/communautés : Forage de structure.
Ce type de fouille vise essentiellement à analyser la structure du Web.
 - ▶ du contenu des pages Web : Forage de contenu,
Ce type de fouille vise essentiellement à analyser le contenu des pages.

Introduction

- Ainsi présenté, le Web n'est plus seulement vu comme un repertoire *fourre-tout*, mais aussi une *toile* organisée où l'on pourrait naviguer.
- La fouille de données a pour objectif l'extraction de connaissances à partir
 - ▶ de la structure formée par les liens entre pages/objets/communautés : Forage de structure.
Ce type de fouille vise essentiellement à analyser la structure du Web.
 - ▶ du contenu des pages Web : Forage de contenu,
Ce type de fouille vise essentiellement à analyser le contenu des pages.
 - ▶ des données d'usages (typiquement des logs de serveurs web) : forage de données d'usage.

- Les premiers moteurs de recherche retrouvaient les documents pertinents seulement en fonction de la similarité du contenu avec la requête de l'utilisateur.

Forage de structure web

- Les premiers moteurs de recherche retrouvaient les documents pertinents seulement en fonction de la similarité du contenu avec la requête de l'utilisateur.
- l'augmentation rapide du nombre de pages Web a rendu très difficile l'ordonnancement des documents retrouvés.

Forage de structure web

- Les premiers moteurs de recherche retrouvaient les documents pertinents seulement en fonction de la similarité du contenu avec la requête de l'utilisateur.
- l'augmentation rapide du nombre de pages Web a rendu très difficile l'ordonnancement des documents retrouvés.
- il est très facile de *corrompre* le contenu d'une page de telle sorte à la rendre similaire à des très nombreuses requêtes.

Forage de structure web

- Les premiers moteurs de recherche retrouvaient les documents pertinents seulement en fonction de la similarité du contenu avec la requête de l'utilisateur.
- l'augmentation rapide du nombre de pages Web a rendu très difficile l'ordonnancement des documents retrouvés.
- il est très facile de *corrompre* le contenu d'une page de telle sorte à la rendre similaire à des très nombreuses requêtes.
- Pour contourner ce problème, il s'est donc posé la question d'exploiter les liens entre les pages. En d'autres termes, exploiter la structure topologique du web.

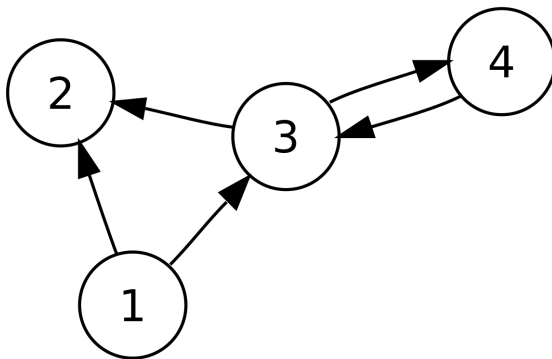
- La structure topologique du Web, peut-etre représentée par une structure de graphe

Forage de structure web

- La structure topologique du Web, peut-etre représentée par une structure de graphe
- où les noeuds représenteraient les différentes pages, tandis que les liens seraient les hyperliens.

- La structure topologique du Web, peut-etre représentée par une structure de graphe
- où les noeuds représenteraient les différentes pages, tandis que les liens seraient les hyperliens.
- De cette manière, on pourrait donc exploiter les propriétés statistiques du graphe et ainsi déduire de la pertinence de certaines pages Webs.

Forage de structure web



- Formellement, la structure topologique du Web ou d'un site Web pourrait être donnée par le graphe $\mathbf{G} = (\mathbf{V}, \mathbf{E})$, où

Forage de structure web

- Formellement, la structure topologique du Web ou d'un site Web pourrait être donnée par le graphe $\mathbf{G} = (\mathbf{V}, \mathbf{E})$, où
- \mathbf{V} serait l'ensemble des noeuds représentant les différentes pages webs et

Forage de structure web

- Formellement, la structure topologique du Web ou d'un site Web pourrait être donnée par le graphe $\mathbf{G} = (\mathbf{V}, \mathbf{E})$, où
- \mathbf{V} serait l'ensemble des noeuds représentant les différentes pages webs et
- $\mathbf{E} = \{(u, v) / u \in \mathbf{V}, v \in \mathbf{V}\}$ l'ensemble des liens représentant les hyperliens entre les pages webs.

- Formellement, la structure topologique du Web ou d'un site Web pourrait être donnée par le graphe $\mathbf{G} = (\mathbf{V}, \mathbf{E})$, où
- \mathbf{V} serait l'ensemble des noeuds représentant les différentes pages webs et
- $\mathbf{E} = \{(u, v) / u \in \mathbf{V}, v \in \mathbf{V}\}$ l'ensemble des liens representant les hyperliens entre les pages webs.
- Lorsque le graphe est non-dirigé on a, $(u, v) = (v, u)$ et lorsqu'il est dirigé on pourrait avoir $(u, v) \neq (v, u)$.

Forage de structure web – Importance d'une page web

- En fonction de sa position dans la structure topologique, on pourrait évaluer sa pertinence en utilisant les mesures de centralités

Forage de structure web – Importance d'une page web

- En fonction de sa position dans la structure topologique, on pourrait évaluer sa pertinence en utilisant les mesures de centralités
- Degré d'un noeud:

Forage de structure web – Importance d'une page web

- En fonction de sa position dans la structure topologique, on pourrait évaluer sa pertinence en utilisant les mesures de centralités
- Degré d'un noeud:

$$D(u) = |\{(u, v) \in \mathbf{E} / v \in \mathbf{V}\}|$$

Forage de structure web – Importance d'une page web

- En fonction de sa position dans la structure topologique, on pourrait évaluer sa pertinence en utilisant les mesures de centralités
- Degré d'un noeud:

$$D(u) = |\{(u, v) \in \mathbf{E} / v \in \mathbf{V}\}|$$

- Prestige d'un noeud (fonctionnel uniquement dans les graphes dirigés), Il représente encore le degré entrant.

Forage de structure web – Importance d'une page web

- En fonction de sa position dans la structure topologique, on pourrait évaluer sa pertinence en utilisant les mesures de centralités
- Degré d'un noeud:

$$D(u) = |\{(u, v) \in \mathbf{E} / v \in \mathbf{V}\}|$$

- Prestige d'un noeud (fonctionnel uniquement dans les graphes dirigés), Il représente encore le degré entrant.

$$P(u) = |\{(v, u) \in \mathbf{E} / v \in \mathbf{V}\}|$$

Forage de structure web – Importance d'une page web

- En fonction de sa position dans la structure topologique, on pourrait évaluer sa pertinence en utilisant les mesures de centralités
- Degré d'un noeud:

$$D(u) = |\{(u, v) \in \mathbf{E} / v \in \mathbf{V}\}|$$

- Prestige d'un noeud (fonctionnel uniquement dans les graphes dirigés), Il représente encore le degré entrant.

$$P(u) = |\{(v, u) \in \mathbf{E} / v \in \mathbf{V}\}|$$

- L'intermédierité:

Forage de structure web – Importance d'une page web

- En fonction de sa position dans la structure topologique, on pourrait évaluer sa pertinence en utilisant les mesures de centralités
- Degré d'un noeud:

$$D(u) = |\{(u, v) \in \mathbf{E} / v \in \mathbf{V}\}|$$

- Prestige d'un noeud (fonctionnel uniquement dans les graphes dirigés), Il représente encore le degré entrant.

$$P(u) = |\{(v, u) \in \mathbf{E} / v \in \mathbf{V}\}|$$

- L'intermédierité:

$$I(u) = \sum_{v \in \mathbf{V} \setminus \{u, w\}} \sum_{w \in \mathbf{V} \setminus \{u, v\}} SP_{v,w}(u)$$

avec $SP_{v,w}(u)$ le nombre de plus court chemins entre v et w en passant par u

- Se base sur le critère de prestige pour faire la classification des pages webs.

- Se base sur le critère de prestige pour faire la classification des pages webs.
- À partir d'une matrice stochastique décrivant la relation entre les différentes pages, on effectue le calcul de classement des pages.

- Se base sur le critère de prestige pour faire la classification des pages webs.
- À partir d'une matrice stochastique décrivant la relation entre les différentes pages, on effectue le calcul de classement des pages.
- Les étapes sont les suivantes:

- Se base sur le critère de prestige pour faire la classification des pages webs.
- À partir d'une matrice stochastique décrivant la relation entre les différentes pages, on effectue le calcul de classement des pages.
- Les étapes sont les suivantes:
 1. Construire la matrice d'adjacence A des pages webs liées

- Se base sur le critère de prestige pour faire la classification des pages webs.
- À partir d'une matrice stochastique décrivant la relation entre les différentes pages, on effectue le calcul de classement des pages.
- Les étapes sont les suivantes:
 1. Construire la matrice d'adjacence A des pages webs liées
 2. Transformer cette matrice sous forme stochastique A^*

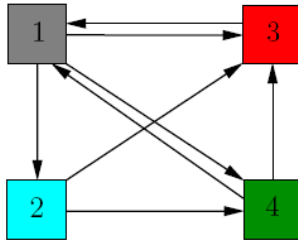
- Se base sur le critère de prestige pour faire la classification des pages webs.
- À partir d'une matrice stochastique décrivant la relation entre les différentes pages, on effectue le calcul de classement des pages.
- Les étapes sont les suivantes:
 1. Construire la matrice d'adjacence A des pages webs liées
 2. Transformer cette matrice sous forme stochastique A^*
 3. Résoudre l'équation $A^* \cdot X = X$ où $X = (x_1, x_2, \dots, x_n)$ est le vecteur de probabilité estimant la chance de parcourir les différentes pages webs.

- Il existe plusieurs méthodes pour résoudre l'équation,

- Il existe plusieurs méthodes pour résoudre l'équation,
- Méthode itérative:

Forage de structure web – PageRank

- Il existe plusieurs méthodes pour résoudre l'équation,
- Méthode itérative:



- On détermine la matrice d'adjacence A

- On détermine la matrice d'adjacence A
- Puis on la transforme en matrice stochastique

- On détermine la matrice d'adjacence A
- Puis on la transforme en matrice stochastique

$$A = \begin{bmatrix} 0 & 0 & 1 & \frac{1}{2} \\ \frac{1}{3} & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} & 0 & 0 \end{bmatrix}$$

- Marche aléatoire jusqu'à convergence:

- Marche aléatoire jusqu'à convergence:

$$\mathbf{v} = \begin{pmatrix} 0.25 \\ 0.25 \\ 0.25 \\ 0.25 \end{pmatrix}, \quad \mathbf{A}\mathbf{v} = \begin{pmatrix} 0.37 \\ 0.08 \\ 0.33 \\ 0.20 \end{pmatrix}, \quad \mathbf{A}^2\mathbf{v} = \mathbf{A}(\mathbf{A}\mathbf{v}) = \mathbf{A} \begin{pmatrix} 0.37 \\ 0.08 \\ 0.33 \\ 0.20 \end{pmatrix} = \begin{pmatrix} 0.43 \\ 0.12 \\ 0.27 \\ 0.16 \end{pmatrix}$$

$$\mathbf{A}^3\mathbf{v} = \begin{pmatrix} 0.35 \\ 0.14 \\ 0.29 \\ 0.20 \end{pmatrix}, \quad \mathbf{A}^4\mathbf{v} = \begin{pmatrix} 0.39 \\ 0.11 \\ 0.29 \\ 0.19 \end{pmatrix}, \quad \mathbf{A}^5\mathbf{v} = \begin{pmatrix} 0.39 \\ 0.13 \\ 0.28 \\ 0.19 \end{pmatrix}$$

$$\mathbf{A}^6\mathbf{v} = \begin{pmatrix} 0.38 \\ 0.13 \\ 0.29 \\ 0.19 \end{pmatrix}, \quad \mathbf{A}^7\mathbf{v} = \begin{pmatrix} 0.38 \\ 0.12 \\ 0.29 \\ 0.19 \end{pmatrix}, \quad \mathbf{A}^8\mathbf{v} = \begin{pmatrix} 0.38 \\ 0.12 \\ 0.29 \\ 0.19 \end{pmatrix}$$

- On pourrait aussi résoudre le problème en résolvant l'équation $A^* \cdot X = X$.

- On pourrait aussi résoudre le problème en résolvant l'équation $A^* \cdot X = X$.
- Dans notre exemple, on aurait le système d'équation

- On pourrait aussi résoudre le problème en résolvant l'équation $A^* \cdot X = X$.
- Dans notre exemple, on aurait le système d'équation

$$\begin{cases} x_1 = 1 \cdot x_3 + \frac{1}{2} \cdot x_4 \\ x_2 = \frac{1}{3} \cdot x_1 \\ x_3 = \frac{1}{3} \cdot x_1 + \frac{1}{2} \cdot x_2 + \frac{1}{2} \cdot x_4 \\ x_4 = \frac{1}{3} \cdot x_1 + \frac{1}{2} \cdot x_2 \end{cases}$$

- Problème avec la matrice stochastique

- Problème avec la matrice stochastique

$$\begin{bmatrix} 0 & 0 & 1 & 0 \\ \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}.$$

- Problème avec la matrice stochastique

$$\begin{bmatrix} 0 & 0 & 1 & 0 \\ \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}.$$

- Lorsqu'on a des valeurs nulles, on pourrait se retrouver dans une difficulté de classement

- Problème avec la matrice stochastique

$$\begin{bmatrix} 0 & 0 & 1 & 0 \\ \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}.$$

- Lorsqu'on a des valeurs nulles, on pourrait se retrouver dans une difficulté de classement

$$\begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 0 \\ \frac{1}{3} & 0 & \frac{1}{3} & \frac{1}{3} \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}^{10000} = \begin{bmatrix} 0.5 & 0 & 0.5 & 0 \end{bmatrix}$$

- Afin de limiter cet effet déplorable,

- Afin de limiter cet effet déplorable,
- Il faudrait s'assurer que la matrice de transition ait toujours des valeurs non-nulles,

- Afin de limiter cet effet déplorable,
- Il faudrait s'assurer que la matrice de transition ait toujours des valeurs non-nulles,
- Pour ce faire on va définir un coefficient de calibrage $q \in [0, 1]$ et ainsi calculer le facteur de *PageRank* d'une page u ($\mathcal{P}(u)$) comme suit:

- Afin de limiter cet effet déplorable,
- Il faudrait s'assurer que la matrice de transition ait toujours des valeurs non-nulles,
- Pour ce faire on va définir un coefficient de calibrage $q \in [0, 1]$ et ainsi calculer le facteur de *PageRank* d'une page u ($\mathcal{P}(u)$) comme suit:

$$\mathcal{P}(u) = \frac{1 - q}{N} + \sum_{v \in In_u} \frac{q\mathcal{P}(v)}{|Out_v|}$$

- Afin de limiter cet effet déplorable,
- Il faudrait s'assurer que la matrice de transition ait toujours des valeurs non-nulles,
- Pour ce faire on va définir un coefficient de calibrage $q \in [0, 1]$ et ainsi calculer le facteur de *PageRank* d'une page u ($\mathcal{P}(u)$) comme suit:

$$\mathcal{P}(u) = \frac{1 - q}{N} + \sum_{v \in In_u} \frac{qP(v)}{|Out_v|}$$

Avec N le nombre de pages webs, $In_u = \{v / (u, v) \in \mathbf{E}\}$,
 $Out_v = \{w / (v, w) \in \mathbf{E}\}$, $P()$ le prestige d'une page.

Forage de structure web – PageRank, Calibrage

$$\begin{bmatrix} 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 \\ \frac{1}{3} & 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 \\ \frac{1}{3} & 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} \frac{1-q}{6} & \frac{q}{3} + \frac{1-q}{6} & \frac{q}{3} + \frac{1-q}{6} & \frac{q}{3} + \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} \\ \frac{q}{3} + \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{q}{3} + \frac{1-q}{6} & \frac{q}{3} + \frac{1-q}{6} \\ q + \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} \\ q + \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} \\ q + \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} \\ q + \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} & \frac{1-q}{6} \end{bmatrix}$$

Forage de contenu web

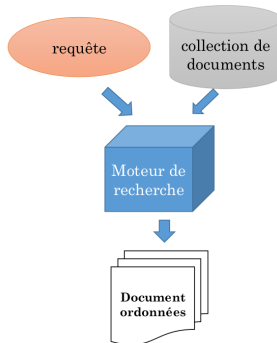
- Ici le contenu de la page web est prise en considération

Forage de contenu web

- Ici le contenu de la page web est prise en considération
- Un classement des pages web est effectué en fonction de leurs scores de pertinence par rapport à la requête.

Forage de contenu web

- Ici le contenu de la page web est prise en considération
- Un classement des pages web est effectué en fonction de leurs scores de pertinence par rapport à la requête.



- Étant donné D l'ensemble des contenus des pages webs.

- Étant donné D l'ensemble des contenus des pages webs.
- La sélection des documents peut se faire suivant une représentation booléenne:

- Étant donné D l'ensemble des contenus des pages webs.
- La sélection des documents peut se faire suivant une représentation booléenne:
 - ▶ Ici la requête de l'utilisateur est prise comme une suite logique.

- Étant donné D l'ensemble des contenus des pages webs.
- La sélection des documents peut se faire suivant une représentation booléenne:
 - ▶ Ici la requête de l'utilisateur est prise comme une suite logique.
 - ▶ Exemple: $((x \text{ AND } y) \text{ AND } (\text{NOT } z))$ signifie les documents contenant les termes x et y et non z

- Étant donné D l'ensemble des contenus des pages webs.
- La sélection des documents peut se faire suivant une représentation booléenne:
 - ▶ Ici la requête de l'utilisateur est prise comme une suite logique.
 - ▶ Exemple: $((x \text{AND} y) \text{AND} (\text{NOT} z))$ signifie les documents contenant les termes x et y et non z
 - ▶ le système récupère chaque document qui rend la requête logiquement vraie

- Étant donné D l'ensemble des contenus des pages webs.
- La sélection des documents peut se faire suivant une représentation booléenne:
 - ▶ Ici la requête de l'utilisateur est prise comme une suite logique.
 - ▶ Exemple: $((x \text{AND} y) \text{AND} (\text{NOT} z))$ signifie les documents contenant les termes x et y et non z
 - ▶ le système récupère chaque document qui rend la requête logiquement vraie
 - ▶ Ici l'appariement est exact, on ne tient pas compte de la pertinence d'un mot dans le document recherché.

- La sélection des documents peut se faire suivant une représentation vectorielle.

- La sélection des documents peut se faire suivant une représentation vectorielle.
 - ▶ Étant donnée la famille de termes $\mathbf{W} = \{W_1, W_2, \dots, W_n\}$ que l'on pourrait rencontrer sur une page web P ,

- La sélection des documents peut se faire suivant une représentation vectorielle.
 - ▶ Étant donnée la famille de termes $\mathbf{W} = \{W_1, W_2, \dots, W_n\}$ que l'on pourrait rencontrer sur une page web P ,
 - ▶ On pourrait représenter la page web P par le vecteur $P = (\omega_1, \omega_2, \dots, \omega_n)$ où ω_i , $1 \leq i \leq n$, est le nombre d'occurrence du mot W_i dans le contenu de la page P .

- La sélection des documents peut se faire suivant une représentation vectorielle.
 - ▶ Étant donnée la famille de termes $\mathbf{W} = \{W_1, W_2, \dots, W_n\}$ que l'on pourrait rencontrer sur une page web P ,
 - ▶ On pourrait représenter la page web P par le vecteur $P = (\omega_1, \omega_2, \dots, \omega_n)$ où ω_i , $1 \leq i \leq n$, est le nombre d'occurrence du mot W_i dans le contenu de la page P .
 - ▶ La requête de l'utilisateur est aussi représenté sous le même format

- La sélection des documents peut se faire suivant une représentation vectorielle.
 - ▶ Étant donnée la famille de termes $\mathbf{W} = \{W_1, W_2, \dots, W_n\}$ que l'on pourrait rencontrer sur une page web P ,
 - ▶ On pourrait représenter la page web P par le vecteur $P = (\omega_1, \omega_2, \dots, \omega_n)$ où ω_i , $1 \leq i \leq n$, est le nombre d'occurrence du mot W_i dans le contenu de la page P .
 - ▶ La requête de l'utilisateur est aussi représenté sous le même format
 - ▶ En fonction d'une mesure de similarité on peut donc faire le classement des pages en fonction de celles qui sont les plus similaires de la requête.