Algèbre Linéaire

Cours du jeudi 19 décembre

Jérôme Scherer



LA MATRICE DE GOOGLE

Soient P_i les (25 milliards de) pages du web. On pose

- **1** I_j le nombre de liens d'une page P_j .
- ② B_i l'ensemble des pages ayant un lien vers P_i .

LA MATRICE DE GOOGLE

Soient P_i les (25 milliards de) pages du web. On pose

- \bullet I_i le nombre de liens d'une page P_i .
- \bigcirc B_i l'ensemble des pages ayant un lien vers P_i .

DÉFINITION

Le rang de la page
$$P_i$$
 vaut $r(P_i) = \sum_{B_i} \frac{r(P_j)}{l_j}$.

LA MATRICE DE GOOGLE

Soient P_i les (25 milliards de) pages du web. On pose

- I_i le nombre de liens d'une page P_i .
- \bigcirc B_i l'ensemble des pages ayant un lien vers P_i .

DÉFINITION

Le rang de la page
$$P_i$$
 vaut $r(P_i) = \sum_{B_i} \frac{r(P_j)}{l_j}$.

Il faut connaître tous les $r(P_i)$ pour calculer $r(P_i)$?

SERGEY BRIN ET LARRY PAGE

C'est le point de départ de l'idée géniale de deux doctorants de



UN SYSTÈME D'ÉQUATIONS

Soit H la matrice dont les coefficients

$$H_{ij} = \begin{cases} \frac{1}{l_j} & \text{si } P_j \in B_i \\ 0 & \text{sinon} \end{cases}$$

Un système d'équations

Soit H la matrice dont les coefficients

$$H_{ij} = \begin{cases} \frac{1}{l_j} & \text{si } P_j \in B_i \\ 0 & \text{sinon} \end{cases}$$

Soit \overrightarrow{r} le vecteur "PageRank" dont les coefficients sont $r(P_i)$. Alors

$$\overrightarrow{r} = H\overrightarrow{r}$$

Un système d'équations

Soit H la matrice dont les coefficients

$$H_{ij} = \begin{cases} \frac{1}{l_j} & \text{si } P_j \in B_i \\ 0 & \text{sinon} \end{cases}$$

Soit \overrightarrow{r} le vecteur "PageRank" dont les coefficients sont $r(P_i)$. Alors

$$\overrightarrow{r} = H\overrightarrow{r}$$

OBSERVATION

Le vecteur PageRank est un vecteur propre pour la valeur propre 1.

REPRÉSENTATION



PROBLÈMES

Pages "dead end" et sous-réseaux.

CORRECTIONS

Pour régler leur compte aux pages sans issues, on remplace les colonnes de zéros par des colonnes de 1/n de sorte que la somme des coefficients de chaque colonne vaut 1. Cette matrice est appelée S.

CORRECTIONS

Pour régler leur compte aux pages sans issues, on remplace les colonnes de zéros par des colonnes de 1/n de sorte que la somme des coefficients de chaque colonne vaut 1. Cette matrice est appelée S.

Conséquence

Le nombre 1 est valeur propre de la matrice S.

Corrections

Pour régler leur compte aux pages sans issues, on remplace les colonnes de zéros par des colonnes de 1/n de sorte que la somme des coefficients de chaque colonne vaut 1. Cette matrice est appelée S.

Conséquence

Le nombre 1 est valeur propre de la matrice S.

Pour régler le problème des sous-réseaux, on choisit un coefficient c (probablement ≈ 0.85). Soit E la matrice dont tous les coefficients valent 1.

Corrections

Pour régler leur compte aux pages sans issues, on remplace les colonnes de zéros par des colonnes de 1/n de sorte que la somme des coefficients de chaque colonne vaut 1. Cette matrice est appelée S.

Conséquence

Le nombre 1 est valeur propre de la matrice S.

Pour régler le problème des sous-réseaux, on choisit un coefficient c (probablement ≈ 0.85). Soit E la matrice dont tous les coefficients valent 1.

$$G = cS + \frac{1-c}{n}E$$

LE SURFEUR DU WEB



Un surfeur du web sur la page P_i a une probabilité c de cliquer sur un lien de cette page, et une probabilité 1-c de recommencer sa recherche d'informations sur une page choisie aléatoirement.

La matrice Google est stochastique et strictement positive.

1 est valeur propre.

- 1 est valeur propre.
- 1 est la plus grande valeur propre parmi toutes les valeurs propres, réelles et complexes.

- 1 est valeur propre.
- 1 est la plus grande valeur propre parmi toutes les valeurs propres, réelles et complexes.
- **3** La dimension de E_1 vaut 1.

- 1 est valeur propre.
- 1 est la plus grande valeur propre parmi toutes les valeurs propres, réelles et complexes.
- **3** La dimension de E_1 vaut 1.
- **1** La deuxième valeur propre de G est de grandeur environ c.

- 1 est valeur propre.
- 1 est la plus grande valeur propre parmi toutes les valeurs propres, réelles et complexes.
- \odot La dimension de E_1 vaut 1.
- **1** La deuxième valeur propre de G est de grandeur environ c.
- La matrice G est diagonalisable.

La matrice Google est stochastique et strictement positive.

- 1 est valeur propre.
- 1 est la plus grande valeur propre parmi toutes les valeurs propres, réelles et complexes.
- \odot La dimension de E_1 vaut 1.
- **1** La deuxième valeur propre de G est de grandeur environ c.
- La matrice G est diagonalisable.

Conséquence

Le vecteur PageRank est bien défini! La somme des $r(P_i)$ vaut 1.

Il est impossible de résoudre un système de 25 milliards d'équations à 25 milliards d'inconnues.

Il est impossible de résoudre un système de 25 milliards d'équations à 25 milliards d'inconnues.

On calcule une approximation par une méthode inductive. On choisit un vecteur \overrightarrow{r}_0 et on calcule

$$\overrightarrow{r}_1 = G\overrightarrow{r}_0$$

Il est impossible de résoudre un système de 25 milliards d'équations à 25 milliards d'inconnues.

On calcule une approximation par une méthode inductive. On choisit un vecteur \overrightarrow{r}_0 et on calcule

$$\overrightarrow{r}_1 = G\overrightarrow{r}_0$$

et on itère. Comme toutes les valeurs propres autres que 1 sont plus petites que 1 strictement, la partie de \overrightarrow{r}_n relative à celles-ci tend vers zéro.

Il est impossible de résoudre un système de 25 milliards d'équations à 25 milliards d'inconnues.

On calcule une approximation par une méthode inductive. On choisit un vecteur \overrightarrow{r}_0 et on calcule

$$\overrightarrow{r}_1 = G\overrightarrow{r}_0$$

et on itère. Comme toutes les valeurs propres autres que 1 sont plus petites que 1 strictement, la partie de \overrightarrow{r}_n relative à celles-ci tend vers zéro.

Conséquence

Le vecteur PageRank $\overrightarrow{r} = \lim \overrightarrow{r}_n$.

Joyeux Noël!

