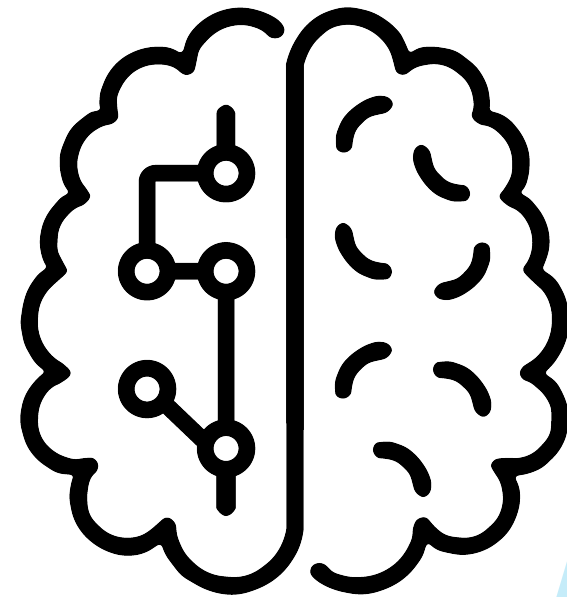
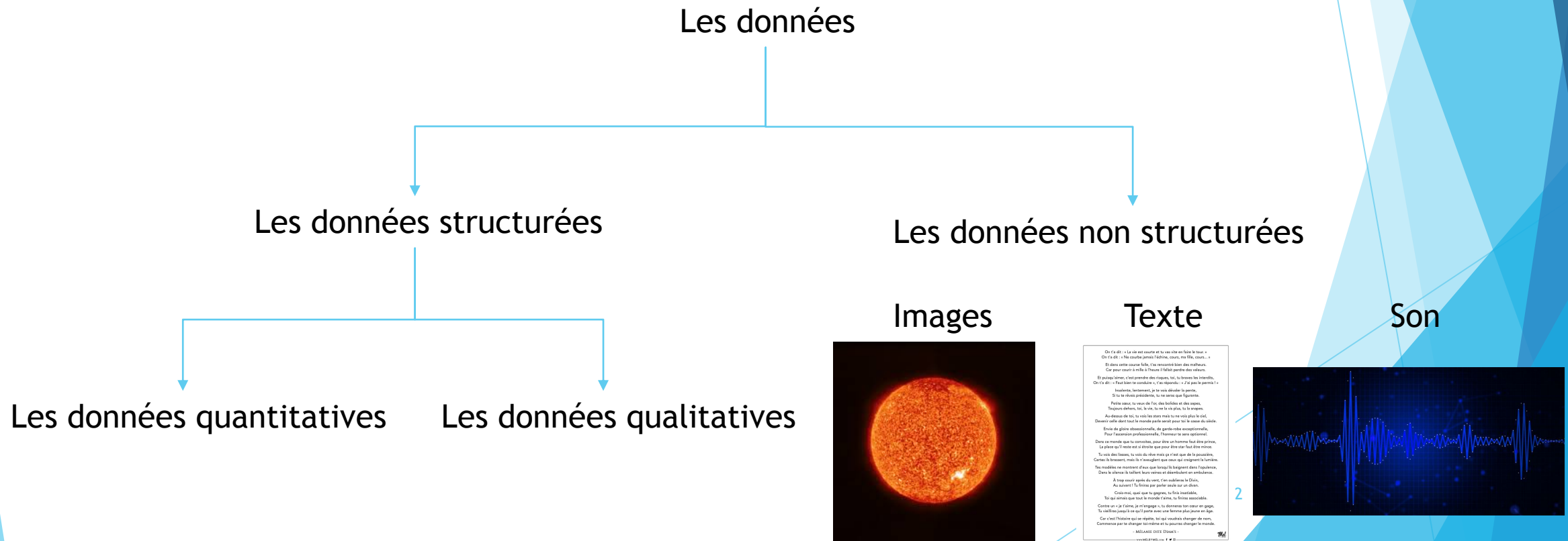


Bag of words

AIForYou - Morgan Gautherot



Les catégories de données



Texte

- MÉLANIE DITE DIAM'S -

[illegible]

Le vecteur de vocabulaire

$$V = \begin{bmatrix} arachide \\ \vdots \\ bain \\ \vdots \\ voiture \\ \vdots \\ zumba \end{bmatrix}$$

D'une phrase à un sac de mots

Doc 1

Le thé est bouillant



Traitements du texte

Doc 1

['thé', 'être', 'bouillant']



Transformation en bag of words

$V = [arachide, \dots, bain, bouillant, \dots, être, \dots, thé, \dots, voiture, \dots, zumba]$

$D_1 = [\quad 0, \dots, \quad 0, \quad 1, \dots, \quad 1, \dots, \quad 1, \dots, \quad 0, \dots, \quad 0]_5$

D'un corpus à un dataframe

Doc 1



Tu me donnes envie de faire table rase du passé, et de t'aider à faire table rase du tien. Tu me donnes envie d'être pour toi ce que je n'ai jamais su être pour quelqu'un d'autre. Tu me donnes envie de t'offrir mon cœur, mes ambitions, mon amour sur un plateau d'argent. Tu me donnes envie de t'aimer à en faire trembler le monde entier. Tu me donnes envie de parler de sentiments et de bonheur tout en conjuguant mes verbes au présent. Tu me donnes envie de devenir quelqu'un qui t'apportera fierté et douceur. Tu me donnes envie de décrocher la Lune, de croire en l'impossible. Tu me donnes envie de vivre, bordel. Parce que tu es l'étincelle au fond de mes yeux, mon sourire terriblement sincère. Parce que je ne suis plus seule désormais, tu es là. Et c'est si précieux.

$= D_1$

Doc 3

006

L'AMITIÉ

L'amitié est un repère, une ligne d'horizon. Une lumière et une porte secrète vers les émotions... Elle nous aide par temps durs, elle fait rire, elle rassure. Elle ne juge pas, elle accepte, elle est là... Tantôt près, tantôt loin, mais toujours à distance de coeurs, notre amitié c'est du bonheur. Notre amitié, c'est une plante que je cultive avec grands soins, le trésor de mon petit jardin. Une pierre précieuse qui m'éclaire dans le noir, et quelqu'un qui connaît toute mon histoire... J'avance sur mon chemin, tu avances sur le tien. L'horizon devant nous est infini, tout comme notre amitié, merci.

$= D_3$

Doc 2

Le supermarché

Sophie fait les courses au supermarché avec son neveu Clément. Après avoir garé sa voiture sur le parking, elle a pris un chariot dans lequel Clément s'est installé. Tous deux pénètrent maintenant dans le magasin. Les clients sont nombreux à cause des promotions et des soldes. Sophie sort sa liste des courses pour ne rien oublier. Elle regarde les marques et compare les prix. Elle choisit des produits frais, des conserves et du poisson surgelé qu'elle met dans son chariot. Clément est turbulent, il s'agite sur son siège, il ne cesse de réclamer des bonbons, des biscuits, du chocolat, des jouets. Sophie le gronde gentiment puis elle regarde sa montre, il est déjà seize heures. Elle doit se hâter de rentrer car sa sœur elle est invitée avec Jean à dîner chez la maman de Clément, sa sœur Lucie. Sophie termine les courses, elle se dirige vers la caisse où elle attend calmement son tour. A la caisse, elle place ses achats sur le tapis roulant puis elle règle sa note avec sa carte bancaire. A présent, elle regagne le parking, cherche quelques instants sa voiture. Elle la repère au fond de l'allée puis se dirige vers son véhicule. Elle ouvre le coffre et y met les courses puis elle installe Clément dans son siège pour enfant. Elle part retrouver Jean et se préparer pour la soirée.

$= D_2$

Doc 4

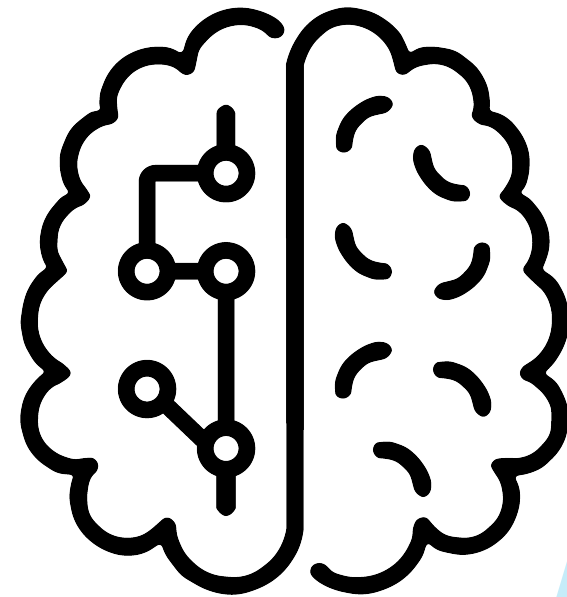
On l'a dit : « La vie est courte et tu vas vite en faire le tour. »
On l'a dit : « Ne courais jamais l'histoire, cours, moi fille, cours. »
Et dans cette course folle, t'es rencontrée bien des malheurs.
Car pour courir à contre l'heure t'as fait perdre des valeurs.
Et puis, hélas, c'est pendant des heures, toi, la brune les cheveux, On l'a dit : « Faut bien le conduire », t'es répondu : « J'ai pas le permis ! »
Insolente, lentement, je te vais dévoiler la gente.
Si tu te dirais présente, tu ne serais que l'égérie.
Petite sœur, tu veux de l'or, des bouillottes et des sapes.
Toujours effrontée, toi, la vie, tu ne te lais plus, tu la mènes.
Au-dessus de toi, tu vois les stars mais tu ne vois plus le ciel.
Devenir ça te doit tout le monde parle sans toi le caser du siècle.
Entre de glorie obsolescence, de grande maladeuse.
Pour l'accession professionnelle, l'honneur te sera optionnel.
Dans ce monde que tu convoites, pour être un homme faut être prince.
La reine qu'il t'en a dit se dit que pour être une faut être reine.
Tu vois des lasses, tu vois du rêve mais ça n'est que de la poussière.
Certes tu brèves, mais tu n'imagines que nous qui courons la vieillesse.
Tu mérites un moment d'être que l'orgueil du langage dans l'apathie.
Dans le silence laissent leurs veines et déboulent en ambulance.
À trop savoir après de vent, t'en maitrises le Divin.
Au secours ! Tu feras pas partie mais tu en disais.
Crois-moi, quoi que tu gagnes, tu fais insupportable.
Toi qui aimais que tout le monde t'aime, tu finis détestable.
Contre un je t'aime, je m'engage », tu donnes ton cœur en gage.
Tu vas-tu jusqu'à ce qu'il parte avec une femme plus jeune en âge.
Car c'est l'histoire qui se répète, toi qui voulais changer de route.
Commence par te changer toi-même et tu pourras changer le monde.

$= D_4$

V	D_1	D_2	D_3	D_4	...	D_n
mot_1	0	1	0	0	...	0
mot_2	0	0	0	1	...	0
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
mot_m	0	1	0	0	0	0

Le TF-IDF

AIForYou - Morgan Gautherot



Définition du TF-IDF

TF = Term Frequency

IDF = Inverse Document Frequency

Term Frequency

$$tf_{w,d} = \frac{n_{w,d}}{\sum_k n_{k,d}}$$

d est un document de notre jeu de données

w est un mot de notre document

$n_{w,d}$ est le nombre d'occurrences du mot w du document d

Inverse Document Frequency

$$idf_w = \log\left(\frac{|D|}{|\{d_i: t_w \in d_i\}|}\right)$$

D désigne tous les documents de notre jeu de données

w est un mot de notre document

$f(w, D)$ est le nombre de document D contenant le mot w

Calcul du TF-IDF

$$tfidf_{w,d} = tf_{w,d} * idf_w$$

$$tfidf_{w,d} = \frac{\text{fréquence du terme dans le document}}{\text{fréquence du terme dans le corpus}}$$

Vecteur de vocabulaire

$$V = \begin{bmatrix} arachide \\ \vdots \\ bain \\ \vdots \\ voiture \\ \vdots \\ zumba \end{bmatrix}$$

D'une phrase à un vecteur

Doc 1

Le thé est bouillant



Traitements du texte

Doc 1

['thé', 'être', 'bouillant']



Transformation en vecteur

$V = [arachide, \dots, bain, bouillant, \dots, être, \dots, thé, \dots, voiture, \dots, zumba]$

$D_1 = [\quad 0, \dots, \quad 0, tfidf_{bouillant}, \dots, tfidf_{être}, \dots, tfidf_{thé}, \dots, \quad 0, \dots, \quad 0]$

D'un corpus à une matrice



Doc 1

Tu me donnes envie de faire table rase du passé, et de t'aider à faire table rase du tien. Tu me donnes envie d'être pour toi ce que je n'ai jamais su être pour quelqu'un d'autre. Tu me donnes envie de t'offrir mon cœur, mes ambitions, mon amour sur un plateau d'argent. Tu me donnes envie de t'aimer à en faire trembler le monde entier. Tu me donnes envie de parler de sentiments et de bonheur tout en conjuguant mes verbes au présent. Tu me donnes envie de devenir quelqu'un qui t'apportera fierté et douceur. Tu me donnes envie de décrocher la Lune, de croire en l'impossible. Tu me donnes envie de vivre, bordel. Parce que tu es l'étincelle au fond de mes yeux, mon sourire terriblement sincère. Parce que je ne suis plus seule désormais, tu es là. Et c'est si précieux.

$= D_1$

Doc 3

006
L'AMITIÉ
L'amitié est un repère,
une ligne d'horizon. Une lumière
et une porte secrète vers les émotions...
Elle nous aide par temps durs, elle fait rire,
elle rassure. Elle ne juge pas, elle accepte,
elle est là... Tantôt près, tantôt loin,
mais toujours à distance de coeurs, notre
amitié c'est du bonheur. Notre amitié, c'est
une plante que je cultive avec grands soins,
le trésor de mon petit jardin. Une pierre
précieuse qui m'éclaire dans le noir, et
quelqu'un qui connaît toute mon histoire...
J'avance sur mon chemin, tu avances sur
le tien. L'horizon devant nous est infini,
tout comme notre amitié, merci.

$= D_3$

Doc 2

Sophie fait les courses au supermarché avec son neveu Clément. Après avoir garé sa voiture sur le parking, elle a pris un chariot dans lequel Clément s'est installé. Tous deux pénétrèrent maintenant dans le magasin. Les clients sont nombreux à cause des promotions et des soldes. Sophie sort sa liste des courses pour ne rien oublier. Elle regarde les marques et compare les prix. Elle choisit des produits frais, des conserves et du poisson surgelé qu'elle met dans son chariot. Clément est turbulent, il s'agite sur son siège, il ne cesse de réclamer des bonbons, des biscuits, du chocolat, des jouets. Sophie le gronde gentiment puis elle regarde sa montre, il est déjà seize heures. Elle doit se hâter de rentrer car sa sœur elle est invitée avec Jean à dîner chez la maman de Clément, sa sœur Lucie. Sophie termine les courses, elle se dirige vers la caisse où elle attend calmement son tour. A la caisse, elle place ses achats sur le tapis roulant puis elle règle sa note avec sa carte bancaire. A présent, elle regagne le parking, cherche quelques instants sa voiture. Elle la repère au fond de l'allée puis se dirige vers son véhicule. Elle ouvre le coffre et y met les courses puis elle installe Clément dans son siège pour enfant. Elle part retrouver Jean et se préparer pour la soirée.

$= D_2$

Doc 4

On l'a dit : « La vie est courte et tu vas vite en faire le tour. »
On l'a dit : « Ne courais jamais l'histoire, cours, moi fille, cours. »
Et dans cette course folle, t'es rencontré bien des malheurs.
Car pour courir à contre l'heure d'été j'ai perdu mes valeurs.
Et puis, hélas, c'est pendant des années, toi, la brune, les intellidits,
On l'a dit : « Faut bien le conduire », t'es répondu : « J'ai pas le permis ! »
Insolente, lentement, je te vais dévoiler la gente,
Si tu ne devais rien attendre, tu ne serais que l'apogée.
Petite sœur, tu veux de l'or, des bouillottes et des sapins,
Toujours effronté, toi, le vie, tu ne te lais plus le ciel.
Au-dessus de toi, tu vois les stars mais tu ne vois plus le ciel.
Devenir ce que d'être tout le monde parle pour toi le caser du siècle.
Entre de glorie obsolescence, de grande et exceptionnelle,
Pour l'accession professionnelle, l'honneur te sera optionnel.
Dans ce monde que tu convoites, pour être un homme faut être prince,
La reine qui t'aime ne se dit pas que pour être une fille reine.
Tu vois des lasses, tu vois du rêve mais ça n'est que de la possession,
Certes tu braves, mais tu n'imagines que ceux qui corrigent la leçon.
Tu mérites un moment d'être que lorsque tu braves dans l'apogée.
Dans le silence laissent leurs vagues et déambulent en ambivalence.
À trop savoir après de vent, t'en mérites la Diva,
Au moins ! Tu feras pas partie mais tu en auras.
Crois-moi, quoi que tu gagnes, tu feras inévitable.
Toi qui saines que tout le monde t'aime, tu feras exceptionelle.
Contre un je t'aime, je m'engage », tu donneras ton cœur en gage,
Tu verras jusqu'à ce qu'il parte avec une femme plus jeune en âge.
Car c'est l'histoire qui se répète, toi qui voudrais changer de route,
Commence par te changer toi-même et tu pourras changer le monde.

$= D_4$

V	D_1	D_2	D_3	D_4	...	D_n
mot_1	0	$fidf_{mot_1}$	0	0	...	0
mot_2	0	0	0	$fidf_{mot_2}$...	0
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
mot_m	0	$fidf_{mot_m}$	0	0	0	0

Notebook 2 - Implémentation du TF-IDF

Objectifs :

- ▶ Implémenter le bag of word pour un corpus de texte
- ▶ Implémenter le TF-IDF pour un corpus de texte