

Programmation et projet encadré - L7TI005

Lancement du projet de groupe

Yoann Dupont yoann.dupont@sorbonne-nouvelle.fr

Pierre Magistry pierre.magistry@inalco.fr

2024-2025

Université Sorbonne-Nouvelle
INALCO
Université Paris-Nanterre

Ce qu'on attend aujourd'hui



Une feuille d'exercices avec tout ce qu'il y aura à faire vous sera donné. On va résumer ici ce qu'il y a dedans et ce qu'on va faire en plus.

Chaque exercice va demander d'ajouter une colonne au tableau existant.

Avant de traiter une page : vérifier qu'elle est valide

Dans le mini-projet, on a vu comment récupérer le code HTTP. Cependant, seul le code 200 indique une réussite.

Avant de traiter une page : vérifier qu'elle est valide

Dans le mini-projet, on a vu comment récupérer le code HTTP. Cependant, seul le code 200 indique une réussite.

On ne veut traiter que les pages qui fonctionnent. Il faut donc ajouter au script quelques vérifications :

1. si le code est 200, on traite l'URL normalement ;
2. tout autre code indique une erreur : on affiche donc un message sur l'erreur standard et on passe à l'URL suivante (tous les résultats, hors URL et code HTTP, devront être remplacés par des valeurs bidons).

Avant de traiter une page : vérifier qu'elle est valide

Dans le mini-projet, on a vu comment récupérer le code HTTP. Cependant, seul le code 200 indique une réussite.

On ne veut traiter que les pages qui fonctionnent. Il faut donc ajouter au script quelques vérifications :

1. si le code est 200, on traite l'URL normalement ;
2. tout autre code indique une erreur : on affiche donc un message sur l'erreur standard et on passe à l'URL suivante (tous les résultats, hors URL et code HTTP, devront être remplacés par des valeurs bidons).

Pour cet exercice, il faudra voir comment interrompre plus ou moins brutalement une boucle.

Préparer le texte 1 : récupérer les sources de données

Jusqu'à présent, on a juste récupéré des informations pour voir si on arrivait à communiquer avec les différents sites web.

Jusqu'à présent, on a juste récupéré des informations pour voir si on arrivait à communiquer avec les différents sites web.

À partir de là, on va commencer les vrais traitements. Mais il faut déjà avoir du matériau :

1. les pages HTML
2. les contenus textuels des pages HTML

Préparer le texte 1 : récupérer les sources de données

Jusqu'à présent, on a juste récupéré des informations pour voir si on arrivait à communiquer avec les différents sites web.

À partir de là, on va commencer les vrais traitements. Mais il faut déjà avoir du matériau :

1. les pages HTML
2. les contenus textuels des pages HTML

À partir des outils déjà vus, on va récupérer le tout et le stocker.

Préparer le texte 2 : gérer l'encodage

On a vu que les pages pouvaient avoir des encodages différents. Quand on veut faire des traitements, on normalise certains aspects du texte. Le premier aspect à normaliser est l'encodage : on voudra que tout soit en UTF-8.

On a vu que les pages pouvaient avoir des encodages différents. Quand on veut faire des traitements, on normalise certains aspects du texte. Le premier aspect à normaliser est l'encodage : on voudra que tout soit en UTF-8.

En partant de ce qu'on avait la semaine dernière :

1. vérifier l'encodage d'une page web
2. convertir le texte en UTF-8 si son encodage est différent (demandera peut-être un peu de recherche)

Préparer le texte 2 : gérer l'encodage

On a vu que les pages pouvaient avoir des encodages différents. Quand on veut faire des traitements, on normalise certains aspects du texte. Le premier aspect à normaliser est l'encodage : on voudra que tout soit en UTF-8.

En partant de ce qu'on avait la semaine dernière :

1. vérifier l'encodage d'une page web
2. convertir le texte en UTF-8 si son encodage est différent (demandera peut-être un peu de recherche)

À la fin, tous les textes qu'on va traiter seront en UTF-8.

Premiers traitements : compter (et compter encore)

Une fois le texte préparé, on va compter les occurrences de chaque mot dans une page.

On fera ensuite les comptages des bigrammes.

Premiers traitements : les contextes d'apparition

Tout simplement, récupérer un empan de texte qui contient notre mot d'intérêt.
On prendra comme contexte les lignes qui sont autour de notre mot.

Premiers traitements : un concordancier

Un concordancier est un moyen pratique de regarder précisément les usages d'un mot donné. Il s'agit d'une version plus lisible de l'exercice précédent.

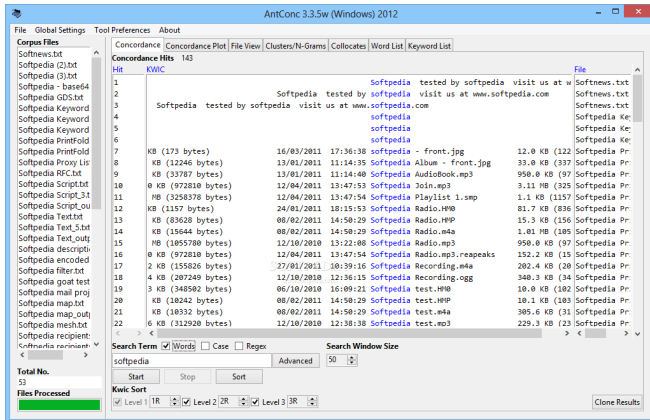


Figure 1: Un exemple de concordancier avec AntConc (source de l'image : <https://www.softpaz.com>)

Point projet



Discutons !

On va passer vous voir pour discuter un peu de vos projets, des vos choix de mots et de langue.

Le rendu

Le rendu sera en groupe :

Le rendu sera en groupe :

1. page d'accueil du projet liée aux tableaux (il faudra mettre en commun vos scripts, sorties et styles)

Le rendu sera en groupe :

1. page d'accueil du projet liée aux tableaux (il faudra mettre en commun vos scripts, sorties et styles)
2. vérifiez bien que les liens marchent aussi sur internet (pas uniquement votre ordinateur)

Le rendu sera en groupe :

1. page d'accueil du projet liée aux tableaux (il faudra mettre en commun vos scripts, sorties et styles)
2. vérifiez bien que les liens marchent aussi sur internet (pas uniquement votre ordinateur)
3. c'est en général à ce moment qu'on a souvent des problèmes de mise en commun. Toujours pull avant commit et push.

Le rendu sera en groupe :

1. page d'accueil du projet liée aux tableaux (il faudra mettre en commun vos scripts, sorties et styles)
2. vérifiez bien que les liens marchent aussi sur internet (pas uniquement votre ordinateur)
3. c'est en général à ce moment qu'on a souvent des problèmes de mise en commun. Toujours pull avant commit et push.

Un seul rendu final pour le projet le 12/01/2025. Vous gérez votre temps, mais attention, plusieurs feuilles à faire !

Le rendu sera en groupe :

1. page d'accueil du projet liée aux tableaux (il faudra mettre en commun vos scripts, sorties et styles)
2. vérifiez bien que les liens marchent aussi sur internet (pas uniquement votre ordinateur)
3. c'est en général à ce moment qu'on a souvent des problèmes de mise en commun. Toujours pull avant commit et push.

Un seul rendu final pour le projet le 12/01/2025. Vous gérez votre temps, mais attention, plusieurs feuilles à faire !

Il faut continuer à avoir un journal, mais on va passer à un journal de groupe. On indiquera qui écrit l'entrée du journal pour plus de clarté.