

spark-shell

Découvrir Apache Spark avec Scala.

Accéder à Spark à travers son shell, dans le répertoire de Spark : `cd bin/` puis `./spark-shell`

Exercice 1.

Écrire un programme Scala qui permet de lister tous les nombres impairs entre 10 et 100 (Fonction `filter()`)

Exercice 2.

Supposons un panier d'achat avec les articles suivants : « Lait », « Fromage », « Beignets », « Pommes », « Bananes ».

1. Représenter les articles ci-dessus dans une structure de données appropriée,
2. Définir une fonction de valeur qui supprime tous les articles de fruits du panier. En d'autres termes, le panier ne doit contenir que les articles « Lait », « Fromage » et « Beignets ».

Exercice 3.

1. Créer un vecteur avec les éléments numériques suivants : 0, 10, 20, 47, -2, 99, -98.
2. Écrire un programme Scala pour trouver le plus petit et le plus grand élément numérique dans le vecteur.

Exercice 4.

1. Créer une structure de données appropriée pour représenter les valeurs numériques suivantes : 2, 8, 19, 20, 25, 50, 100, 10.
2. Définir une fonction de valeur qui identifiera si une valeur numérique donnée est divisible par 2.
3. Définir la fonction `carre` qui retourne le carré d'un élément (ici le paramètre de la fonction)
4. Utiliser la fonction `map` pour appliquer un traitement, calculer le carré de chaque élément de la liste en utilisant la fonction `carre`
5. Utiliser une fonction anonyme pour faire le même traitement
6. Définir une fonction `Int→Boolean` qui retourne `true` si un nombre est pair
7. Utiliser la fonction `map` pour transformer la liste des entiers (carré des éléments) en une liste de booleans
Ex des resultat: `res4: List[Boolean] = List(false, true, ...)`
8. Appliquer la fonction `filter` pour restituer que les éléments pairs de la liste
9. Appliquer la fonction `foldLeft` pour calculer la somme des éléments de la liste

Exercice 5.

1. Créer le fichier texte « `input.txt` »:
« Spark est le projet open source, le plus actif dans le monde du Big Data.
Il a commencé comme un projet de recherche à l'UC Berkeley AMPLab en 2009 et était open source au début de 2010. »
2. Charger ce fichier texte en utilisant l'API `SparkContext`
3. Séparer les mots selon les caractères d'espacement (produire plusieurs RDD à partir d'un RDD : `flatMap`)
4. Produire un RDD à partir d'un autre RDD : `map`, `filter`, `reduceByKey`
 - 4.1. Appliquer un `map` sur les mots obtenus pour produire le couple (`<mot>,1`)
 - 4.2. Appliquer un `reduce` qui permet de faire la somme des 1 des mots identiques
5. Persistance : transférer les résultats vers un fichier « `file1out.count` »