Université Badji Mokhtar - Annaba

Faculté des Sciences Département d'Informatique

Travaux Pratiques : Data Mining

Sujet : Fusion et Extraction des N-Grams d'un fichier FASTA

Réalisé par : Wissem Mebarek

Djoudi islem

Encadré par : M. Boulemden

Année universitaire : 2025 / 2026

1. Introduction

Dans le domaine de la bio-informatique, l'analyse de séquences d'ADN permet de mieux comprendre la structure et le fonctionnement du matériel génétique. Ces séquences sont généralement représentées dans des fichiers au format FASTA, un format textuel standardisé qui contient des identifiants et des séquences d'acides nucléiques.

Le présent travail pratique (TP) s'inscrit dans le module de **Data Mining**, et vise à manipuler ce type de données à l'aide du langage Python. Plus précisément, il s'agit de :

- Fusionner plusieurs fichiers FASTA pour créer une base de séquences unique.
- Extraire et analyser des motifs de séquences appelés **n-grams**, permettant de détecter des répétitions ou similarités entre séquences.

Ce TP illustre ainsi les étapes de prétraitement, de manipulation et d'analyse de données textuelles appliquées au domaine biologique.

2. Objectifs du TP

Les objectifs du TP sont les suivants :

- Manipulation des fichiers FASTA : comprendre leur structure et les lire efficacement avec Python.
- **Fusion des données biologiques** : regrouper plusieurs fichiers pour créer un seul corpus cohérent.
- Extraction d'informations : utiliser les n-grams pour détecter des motifs de taille variable (1, 2, 3, ...).
- **Visualisation et interprétation** : présenter les résultats de manière claire et exploitable.

Ces étapes constituent une première approche du processus de *data preprocessing* en Data Mining.

3. Programme fusion.py

Ce premier programme est dédié à la fusion des fichiers FASTA. L'objectif est de combiner plusieurs jeux de séquences en un seul fichier merged_file.fasta, afin de simplifier les étapes d'analyse ultérieures.

3.1. Fonction parse_fasta(filepath)

Cette fonction lit un fichier FASTA ligne par ligne. Lorsqu'une ligne commence par le symbole ">", elle contient des informations d'en-tête (identifiant, nom de l'organisme,

etc.). Les lignes suivantes contiennent les séquences nucléotidiques. Chaque séquence est enregistrée dans un dictionnaire contenant quatre clés :

```
entry : identifiant de la séquence,
entry_name : nom d'entrée,
organism : organisme d'origine,
sequence : séquence d'ADN complète.
```

Cette fonction retourne une liste de dictionnaires simulant un tableau de données.

```
def parse_fasta(filepath):
   entry = entry_name = organism = ""
    sequence = "'
    with open(filepath, "r") as f:
        for line in f:
            line = line.strip()
            if line.startswith(">"):
                # Sauvegarde de la séquence précédente
                if entry:
                    data.append({
                        "entry": entry,
                        "entry_name": entry_name,
                        "organism": organism,
                        "sequence": sequence
                    sequence = ""
                parts = line.split()
                id_parts = parts[0].split("|")
                entry = id_parts[1] if len(id_parts) > 1 else ""
                entry_name = id_parts[2] if len(id_parts) > 2 else ""
                organism = ""
                for part in parts:
                    if part.startswith("OS="):
                        organism = part[3:]
                        break
            else:
                sequence += line
    if entry:
       data.append({
            "entry": entry,
            "entry_name": entry_name,
            "organism": organism,
            "sequence": sequence
    return data
```

FIGURE 1 – Lecture et parsing des séquences FASTA

3.2. Fonction merge_fasta_files(input_folder)

Cette fonction a pour but de fusionner plusieurs fichiers FASTA situés dans un même dossier. Elle parcourt chaque fichier du répertoire, lit son contenu à l'aide de la fonction précédente et ajoute les séquences dans une seule liste globale. Cela permet de regrouper toutes les séquences en un seul ensemble prêt à être sauvegardé.

```
def merge_fasta_files(file_list):
    """Fusionne plusieurs fichiers FASTA en un seul DataFrame (liste de dictionnaires)."""
    merged_data = []
    for file in file_list:
        print(f"Lecture de {file} ...")
        merged_data.extend(parse_fasta(file)) # ajoute les lignes du fichier
    return merged_data
```

FIGURE 2 – Fusion des fichiers FASTA en une seule structure

3.3. Fonction save_as_fasta(data, output_file)

Une fois la fusion terminée, cette fonction sauvegarde toutes les séquences dans un nouveau fichier FASTA. Chaque entrée est formatée selon la norme FASTA :

```
>IDENTIFIANT NOM_ENTRÉE OS=ORGANISME
SEQUENCE
```

Cette étape finalise la première partie du TP, en créant le fichier merged_file.fasta utilisé pour l'extraction des n-grams.

```
def save_as_fasta(df, output_path):
    """Sauvegarde la liste de dictionnaires dans un fichier FASTA."""
   with open(output_path, "w") as f:
           entry = row.get("entry", "")
           entry_name = row.get("entry_name", "")
           organism = row.get("organism", "")
sequence = row.get("sequence", "")
           header = f">sp|{entry}|{entry_name} OS={organism}\n"
           f.write(header)
           for i in range(0, len(sequence), 60):
               f.write(sequence[i:i+60] + "\n")
# Liste de tes 5 fichiers FASTA
files = [
    "bovine.fasta",
    "mouse.fasta",
    "zibrafish.fasta"
df_merged = merge_fasta_files(files)
# Nombre total d'entrées
print(f"\nNombre total de séquences fusionnées : {len(df_merged)}")
output_path = "merged_files.fasta"
save_as_fasta(df_merged, output_path)
```

FIGURE 3 – Enregistrement du fichier FASTA fusionné

4. Programme ngram.py

Ce second programme se concentre sur l'analyse du fichier fusionné. L'objectif est d'extraire des motifs (\mathbf{n} -grams) à partir des séquences ADN et d'observer leur fréquence d'apparition. Un \mathbf{n} -gramme est une sous-chaîne de longueur n:

```
    1-gram = A, T, G, C
    2-gram = AT, CG, TG, etc.
    3-gram = ATT, TGA, GCC, etc.
```

4.1. Fonction extract_ngrams(df, n)

Cette fonction par court chaque séquence du fichier fusionné. Pour une taille n donnée, elle génère tous les motifs consécutifs de cette longueur et calcule le ur fréquence. Le résultat est un dictionnaire associant chaque motif à son nombre d'occurrences pour chaque séquence.

```
def extract_ngrams(df, n=1):
   """Extrait les n-grammes de taille n pour chaque séquence."""
   results = []
   for row in df:
       seq = row["sequence"]
       entry_id = row["entry"]
       # Dictionnaire des fréquences
       grams_counts = {}
       for i in range(len(seq) - n + 1):
           gram = seq[i:i+n]
           grams_counts[gram] = grams_counts.get(gram, 0) + 1
       # Associe l'ID et les fréquences
       features = {"ID": entry id}
       features.update(grams counts)
       results.append(features)
   return results
```

FIGURE 4 - Code de la fonction extract_ngrams()

Cette méthode est essentielle car elle permet de transformer une séquence textuelle en un ensemble de caractéristiques numériques, étape cruciale dans toute démarche de Data Mining.

4.2. Fonction show_results(features, n)

Cette fonction s'occupe d'afficher les résultats des n-grams extraits. Elle affiche notamment :

- le nombre total de séquences traitées;
- pour chaque séquence, la liste des n-grams détectés et leur fréquence.

```
def show_results(features, n):
    """Affiche les 10 premières séquences avec leurs n-grammes."""
    if not features:
        print(f"\nAucun {n}-gramme extrait.")
        return

print(f"\n=== Extraction des {n}-grammes ===")
    print(f"Nombre total de séquences traitées : {len(features)}\n")

for i, f in enumerate(features[:10]): # 10 premières séquences
    print(f"--- Séquence {i+1} ---")
    print("ID :", f["ID"])
    grams_sorted = sorted((k, v) for k, v in f.items() if k != "ID")
    for gram, count in grams_sorted:
        print(f" {gram}: {count}", end=" ")
        print("\n" + "-" * 40)
```

FIGURE 5 - Code de la fonction show_results()

4.3. Résultats obtenus

Les résultats ci-dessous illustrent l'exécution du programme pour les valeurs de n=1, n=2 et n=3. Ces trois cas permettent de comparer la granularité des motifs extraits : plus n est grand, plus les motifs sont précis, mais aussi plus rares.

FIGURE 6 – Résultats des séquences pour les 1-grams

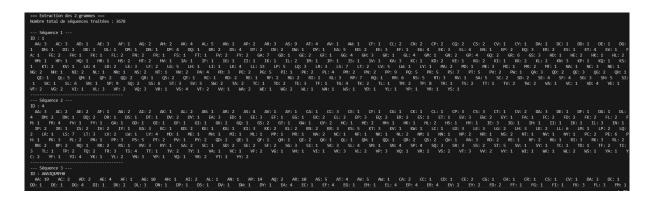


FIGURE 7 – Résultats des séquences pour les 2-grams

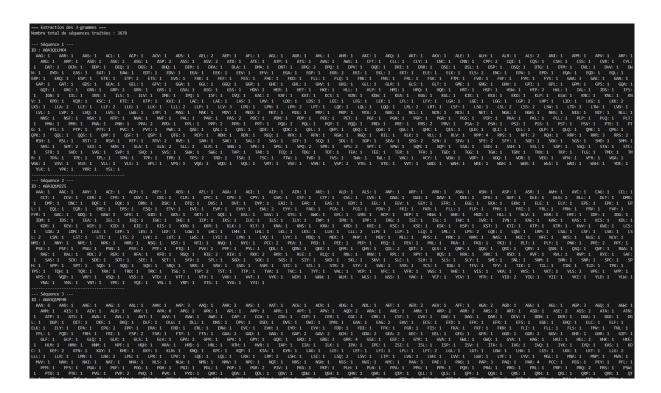


FIGURE 8 – Résultats des séquences pour les 3-grams

5. Conclusion

À travers ce TP, nous avons mis en œuvre deux aspects essentiels du Data Mining appliqué à la bio-informatique :

- 1. le **prétraitement des données biologiques**, via la lecture, la fusion et la normalisation des séquences FASTA;
- 2. l'analyse de motifs à travers l'extraction des n-grams pour identifier les combinaisons les plus fréquentes.

Ce travail met en évidence l'importance de l'automatisation et du traitement textuel dans l'analyse de données biologiques. Les concepts de n-grams sont également utilisés dans d'autres domaines, tels que le traitement automatique du langage naturel (NLP), la détection de similarité ou encore la classification de séquences.