

خطة تنفيذ المشروع لكورس مقدمة في الجينومات الحاسوبية

بناءً على الملفات المرفقة، سأقدم خطة شاملة مقسمة إلى مراحل واضحة مع خطوات تنفيذ **Challenge_2025.pdf** لتنفيذ المشروع المطلوب في ملف كل مرحلة. كما سأشرح وظيفة كل ملف وما يحتويه وكيفية استخدامه، وأحدد المخرجات المطلوب تسليمها بدقة. الهدف هو تحقيق متطلبات المشروع باستخدام البرمجة وعلم الأحياء مع الاستفادة من جميع الملفات المقدمة.

وظيفة كل ملف ومحتواه وكيفية استخدامه

1. extract_sequences.py

- لاستخراج تسلسلات جينية من تسلسل **extract_flanked_region** تُسمى Python الوظيفة: يحتوي على دالة مع مراعاة اتجاه السلسلة (upstream و downstream) كروموسومي بناءً على مواقع البداية والنهاية مع إضافة مناطق مجاورة (strand).
- وعدد، (based- يأخذ معطيات مثل التسلسل الكروموسومي، الاتجاه (+ أو -)، مواقع البداية والنهاية (1 Python المحتوى: كود القواعد المجاورة، ويُرجع التسلسل المستخرج).
- بناءً على بيانات الطفرات في ملفات **genes.csv** كيفية الاستخدام: سأستخدم هذه الدالة لاستخراج التسلسلات الجينية من ملف **100 test_muts_data.csv** مع إضافة مناطق مجاورة (مثل 10 قواعد) لتحليل الطفرات، **train_muts_data.csv** و **test_muts_data.csv**.

2. test_feats.csv

- لعينات الاختبار، وهي عدد الطفرات لكل جين لكل عينة (features) الوظيفة: يحتوي على بيانات ميزات (...). وأعمدة لـ 100 جين مع قيم عددية تمثل عدد الطفرات (مثل 0، 1، 2 **case_id** يحتوي على أعمدة CSV المحتوى: جدول **train_feats.csv** باستخدام نموذج تعلم آلي مدرب على (labels) كيفية الاستخدام: سأستخدمه كمدخل للتنبؤ بالتسميات.

3. Challenge_2025.pdf

- الوظيفة: يحتوي على وصف المشروع والمتطلبات الأساسية.
- النص بالعبرية والعربية مختلط وغير مفهوم، لكن يمكن استنتاج أن OCR المحتوى: تعليمات غير واضحة تمامًا بسبب مشاكل الـ **Error = sum(df_known_test['Label'] != df_known_test['predict_label']) / len(df_known_test)**: المطلوب هو تحليل بيانات الطفرات وبناء نموذج للتنبؤ بالتسميات مع تقديم تقرير. يظهر معادلة لحساب الخطأ.
- كيفية الاستخدام: سأعتمد على هذا الملف لفهم الهدف العام (تصنيف العينات بناءً على الطفرات) وأستخدم المعادلة لتقييم النموذج.

4. train_feats.csv

- (labels) الوظيفة: بيانات تدريب تحتوي على ميزات العينات مع تسمياتها.
- قيم مثل 1.0 أو 2.0)، وأعمدة لـ 100 جين مع عدد الطفرات **case_id**، **Label** يحتوي على CSV المحتوى: جدول
- كيفية الاستخدام: سأستخدمه لتدريب نموذج تعلم آلي للتنبؤ بالتسميات بناءً على عدد الطفرات.

5. 100_genes.csv

- (الوظيفة: يحتوي على معلومات الجينات المرجعية (تسلسلات ومواقع التسلسل النووي)، Sequence، (- الاتجاه + أو strand، (اسم الجين) **gene** يحتوي على أعمدة CSV المحتوى: جدول (معلومات موقع الجين على الكروموسوم) **Info** و
- لتحليل التأثير الجيني **extract_sequences.py** كيفية الاستخدام: سأستخدم التسلسلات لاستخراج مناطق الطفرات باستخدام

6. test_muts_data.csv

- الوظيفة: يحتوي على بيانات الطفرات التفصيلية لعينات الاختبار.

- **case_id**، **Gene_name**، **Chromosome**، **Start_Position**، **End_Position**، **Mut_Strand**، **Variant_Classification** (مثل **Missense_Mutation**) ومعلومات الأليلات.
- وللتحقق من دقة **genes.csv** كـ **كيفية الاستخدام**: سأستخدمه لتحليل الطفرات واستخراج التسلسلات لمقارنتها مع **100 test_feats.csv**.

7. train_muts_data.csv

- **الوظيفة**: يحتوي على بيانات الطفرات التفصيلية لعينات التدريب مع التسميات.
- **Label** مع إضافة عمود **test_muts_data.csv** مشابه لـ **CSV المحتوى**: جدول.
- **train_feats.csv** كـ **كيفية الاستخدام**: سأستخدمه لتحليل الطفرات وتدريب النموذج وربط التسميات مع

الخطوة المقترحة لتنفيذ المشروع

المرحلة 1: تحليل البيانات وفهمها

- **المدة الزمنية**: 2-3 أيام
- **الخطوات**:
 1. **Python قراءة الملفات**: استيراد جميع الملفات باستخدام (CSV لـ pandas مكتبة).
 2. **فحص البيانات**:
 - **test_muts_data.csv** و **test_feats.csv**: التحقق من توافق عدد الطفرات في
 - **train_muts_data.csv** و **train_feats.csv**: التحقق من توافق
 - مع الجينات في ملفات الطفرات **genes.csv** فحص أسماء الجينات في **100**
 3. **معالجة البيانات الأولية**:
 - (تنظيف البيانات (إزالة القيم المفقودة إن وجدت).
 - (إلى قيم عددية (مثل 1 و 2 و **train_feats.csv** تحويل التسميات في
 4. وعدد الطفرات الكلي في ملفات الميزات (**Missense** فهم **العلاقات**: تحليل العلاقة بين الطفرات التفصيلية (مثل

المرحلة 2: استخراج التسلسلات الجينية

- **المدة الزمنية**: 3-4 أيام
- **الخطوات**:
 1. **إعداد البيانات**:
 - من (**Start_Position**, **End_Position**, **Mut_Strand**) استخراج مواقع الطفرات **test_muts_data.csv** و **train_muts_data.csv**.
 - باستخدام أسماء الجينات **genes.csv** ربطها مع التسلسلات في **100**
 2. **تعديل الدالة**:
 - مع إضافة معلمة لتحديد حجم المناطق **extract_sequences.py** من **extract_flanked_region** استخدام
 - **upstream=10**، **downstream=10** المجاورة (مثل
 3. **استخراج التسلسلات**:
 - لكل طفرة، استخراج التسلسل المرجعي مع المناطق المجاورة
 - كما في الدالة (**reverse complement**) إذا كان الاتجاه "-", حساب التكامل العكسي
 4. تحتوي على (**train_sequences.csv** و **test_sequences.csv** تخزين النتائج: إنشاء ملفات جديدة (مثل
 - التسلسلات المستخرجة مع تفاصيل الطفرات

المرحلة 3: بناء وتدريب نموذج التعلم الآلي

- المدة الزمنية: 4-5 أيام
- الخطوات:
 1. تحضير البيانات:
 - (Label): كمدخل (الميزات: عدد الطفرات لكل جين، الهدف **train_feats.csv** استخدام
 - scikit-learn من **train_test_split** تقسيم البيانات إلى تدريب (80%) واختبار داخلي (20%) باستخدام
 2. اختيار النموذج:
 - لأنها مناسبة لتصنيف متعدد الفئات Logistic Regression أو Random Forest تجربة نماذج بسيطة مثل
 3. تدريب النموذج:
 - scikit-learn تدريب النموذج على بيانات التدريب باستخدام مكتبة
 - GridSearchCV باستخدام Random Forest ضبط المعلمات (مثل عدد الأشجار في
 4. تقييم النموذج:
 - **Challenge_2025.pdf** استخدام معادلة الخطأ من
 - على بيانات $Error = \frac{\sum(predictions \neq true_labels)}{\len(true_labels)}$ الاختبار الداخلية
 - F1-score و Accuracy حساب مقاييس أخرى مثل
 5. **train_muts_data.csv** التحسين: إذا كان الأداء ضعيفاً، تجربة ميزات إضافية (مثل نوع الطفرة من

المرحلة 4: التنبؤ وتحليل النتائج

- المدة الزمنية: 2-3 أيام
- الخطوات:
 1. التنبؤ على بيانات الاختبار:
 - **test_feats.csv** استخدام النموذج المدرب للتنبؤ بتسميات
 2. مقارنة التسلسلات:
 - مع التسلسلات المرجعية لتحديد تأثير الطفرات (مثل تغيير **test_muts_data.csv** مقارنة التسلسلات المستخرجة من (الأحماض الأمينية
 3. إنشاء ملف النتائج:
 - **test_feats.csv** لكل عينة في **predict_label** و **case_id** يحتوي على CSV ملف

المرحلة 5: كتابة التقرير وإعداد التسليم

- المدة الزمنية: 2-3 أيام
- الخطوات:
 1. كتابة التقرير:
 - مقدمة عن المشروع وأهدافه
 - (وصف المنهجية (استخراج التسلسلات، بناء النموذج، التنبؤ
 - (النتائج (دقة النموذج، أمثلة لتسلسلات الطفرات
 - **Challenge_2025.pdf** في OCR الاستنتاجات والتحديات (مثل مشاكل
 2. إعداد الملفات:
 - (predictions.csv) ملف التنبؤات
 - ملفات التسلسلات المستخرجة
 - المستخدم مع تعليقات Python كود
 - 3.