

EDA 流程濃縮摘要

交給夥伴時，只要照「目的 → 具體操作 → 成果 / 後續要做的事」這個順序講解，他就能快速複製流程。

步驟	目的 & 核心概念	具體做法	之後前處理要注意
1. 載入並合併三張表 <code>train+features+patient_notes</code>	把診斷詞、病歷全文、標註整合到同一筆資料	<code>pd.merge()</code> 兩次 (先用 <code>feature_num</code> ，再用 <code>pn_num</code>)	保持主鍵唯一；後面切分資料時用 <code>case_num</code> 避免洩漏
2. 解析 <code>annotation / location</code>	原欄位是「字串化 list」，需還原成真正的 Python list	<code>ast.literal_eval()</code> → <code>annotation_list</code> 、 <code>location_list</code>	任何轉換失敗回傳空 list，避免後續報錯
3. 建立 <code>has_annotation</code> 標記	判斷該筆病歷有無標註診斷詞	<code>has_annotation = len(annotation_list) > 0</code>	之後可用來做正負樣本抽樣 / stratified split
4. 病歷長度檢查	確認輸入長度是否超過 BERT 512 tokens	<code>AutoTokenizer</code> 算 <code>tok_len</code> ，畫直方圖	若 >512 比例高，考慮長序列模型或滑動視窗截斷
5. 類別分佈分析	總樣本數 & 標註率 是否嚴重失衡	- <code>value_counts()</code> 畫 Zipf plot - 排序後畫 Lorenz curve & Gini	Gini 越大→越不均；可決定 loss 加權 / 重新抽樣策略
6. 單診斷詞標註率	哪些 <code>feature_text</code> 容易被標註？	<code>pos_rate = pos / total</code> 排序觀看	幫助挑重點類別做錯誤分析或數據增強
7. 一致性檢查	標註數 ≠ 位置數	比對 <code>len(annotation_list)</code> vs <code>len(location_list)</code>	若不相等須人工排查 (目前 0 筆)
8. 重複病歷偵測	避免同一病歷被重複訓練多次	去除空白後做 MD5 → <code>pn_hash</code> ，統計出現次數	針對 <code>cnt > 1</code> 的病歷：- 只保留一筆或- 交叉驗證時分層抽樣

傳達給夥伴的重點

- 1. 先跑步驟 1-3，把欄位轉乾淨，`has_annotation` 做完才能正確分群。
- 2. 步驟 4 決定輸入長度策略 (>512 token 就要截斷或換模型)。
- 3. 步驟 5-6 看完不平衡情形，決定 **class weight / sampler**。
- 4. 步驟 7-8 做資料品質把關：標註對齊、刪重複。
- 5. 後續 資料前處理 (tokenize→label encoding) 時，一定以 `annotation_list` 和 `location_list` 為準，別用原始字串欄位。