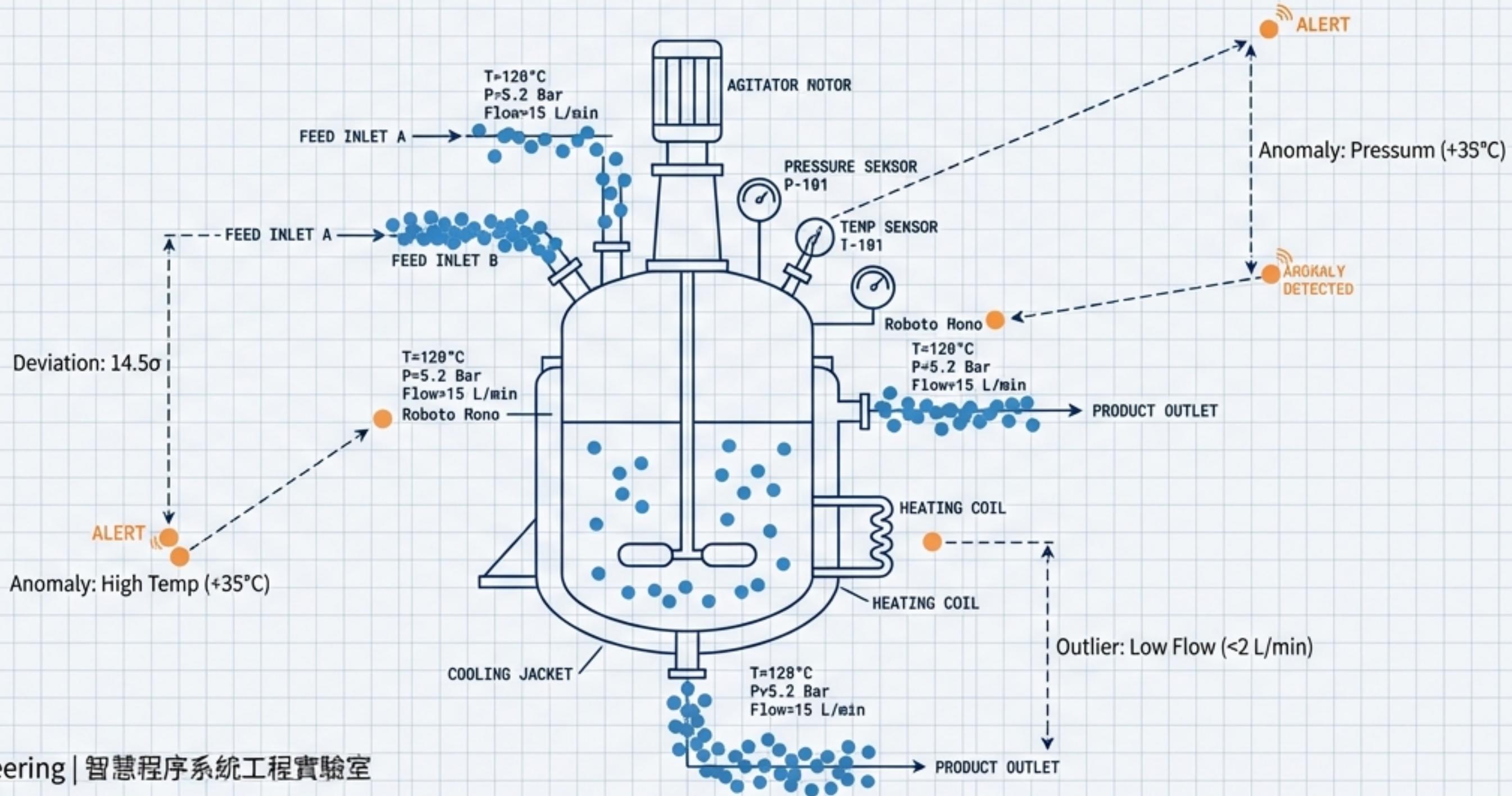


Unit 07 區域性離群因子 (LOF)

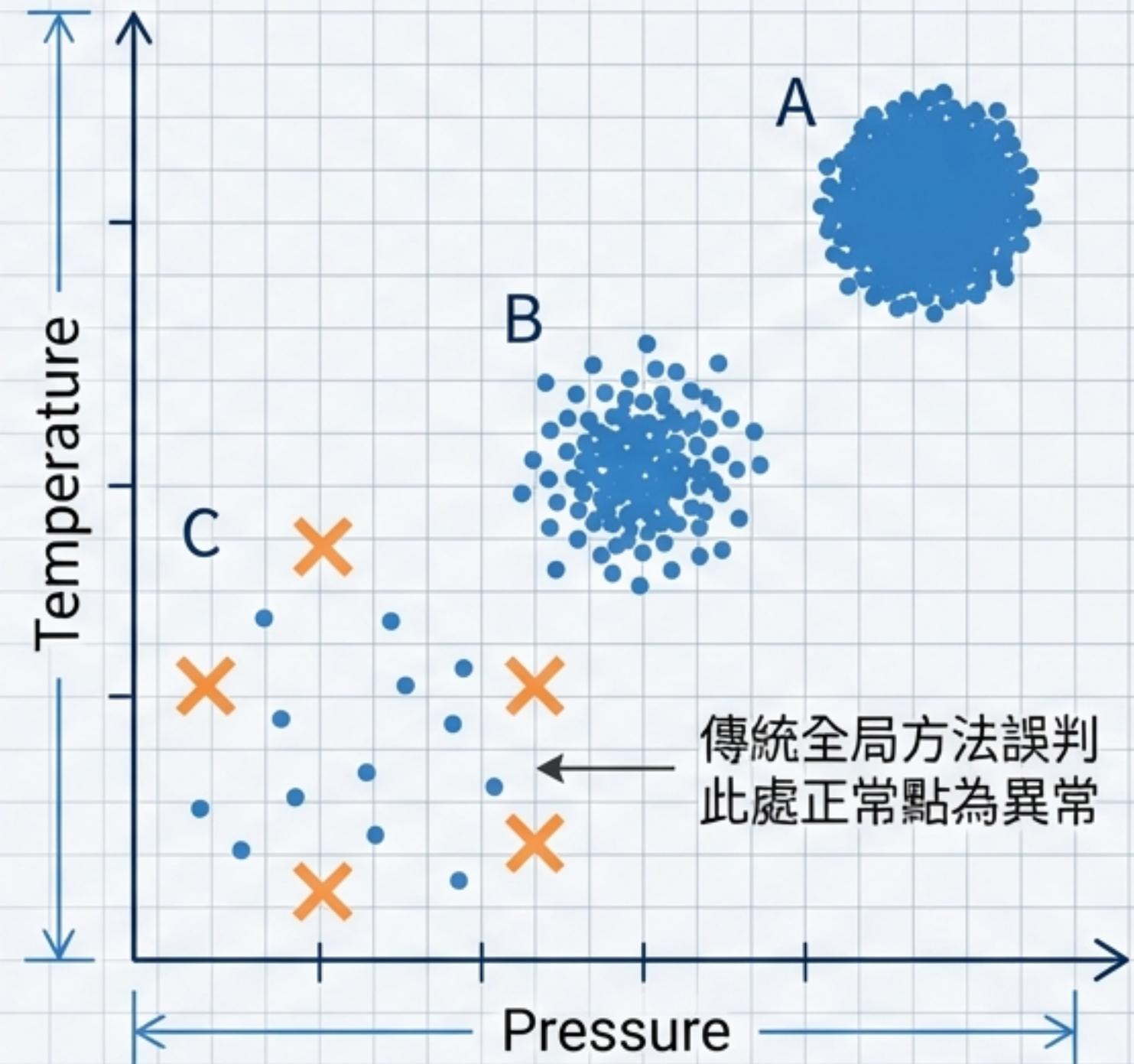
應對多密度化工數據的異常檢測方案



工業現場的挑戰：多模式製程數據 (Multi-Modal Processes)

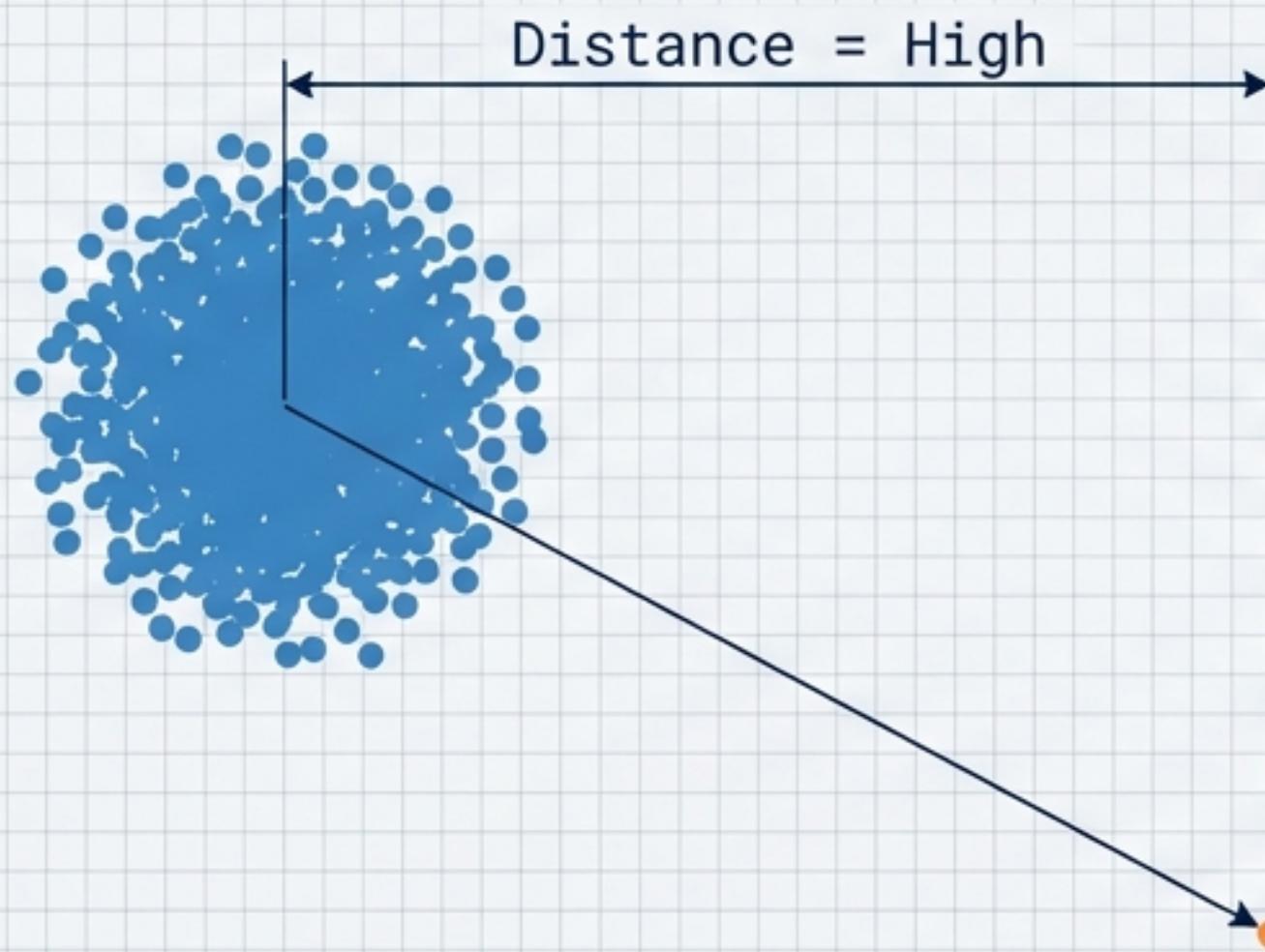
操作配方 (Operating Recipes)

- 1. Recipe A: 高溫高壓 (80-90°C)
 → [數據極度密集]
- 2. Recipe B: 中溫中壓 (60-70°C)
- 3. Recipe C: 低溫低壓 (40-50°C)
 → [數據相對稀疏]

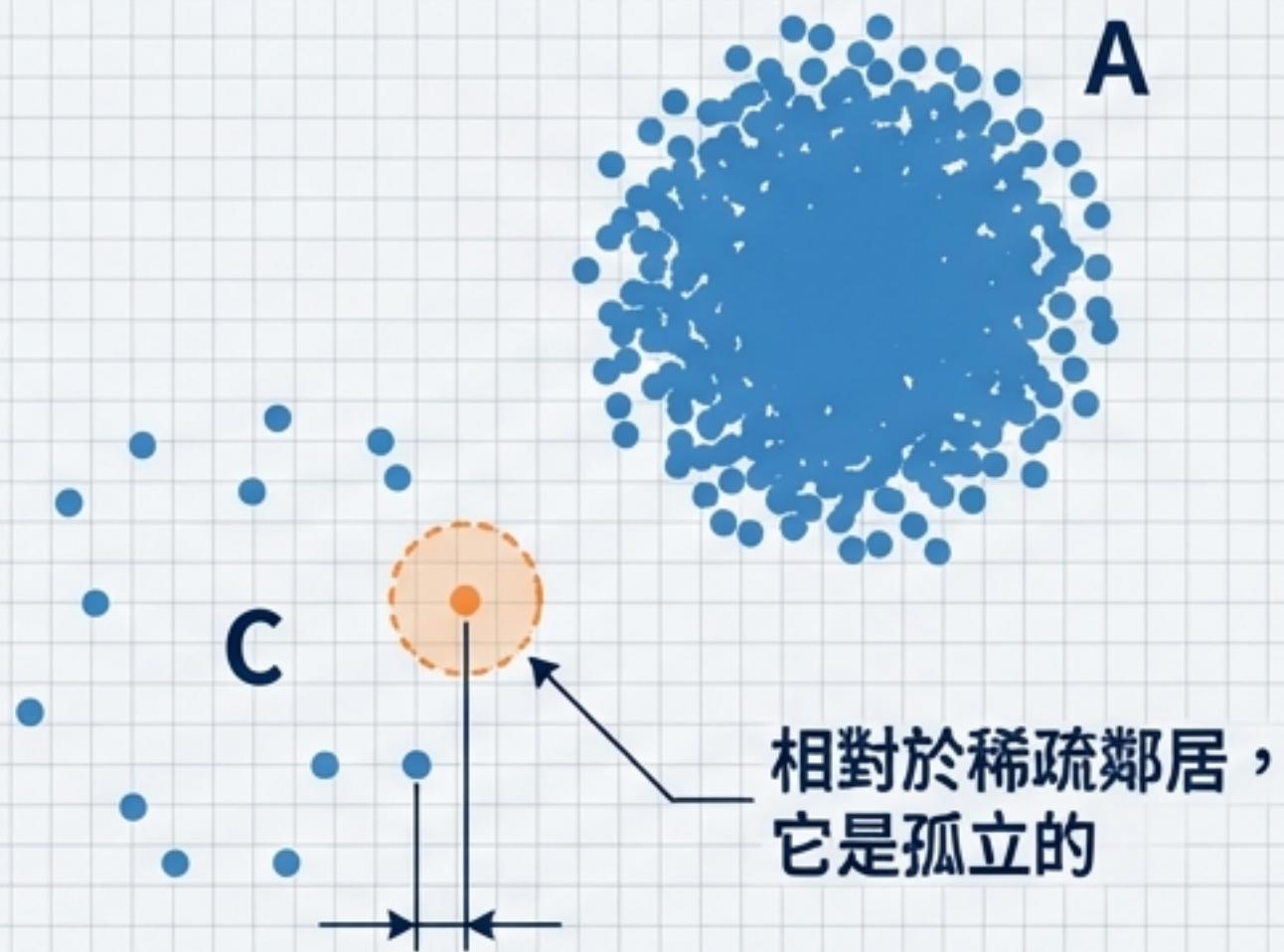


LOF 核心理念：異常是相對於鄰居而言

全局觀點 (Global View)

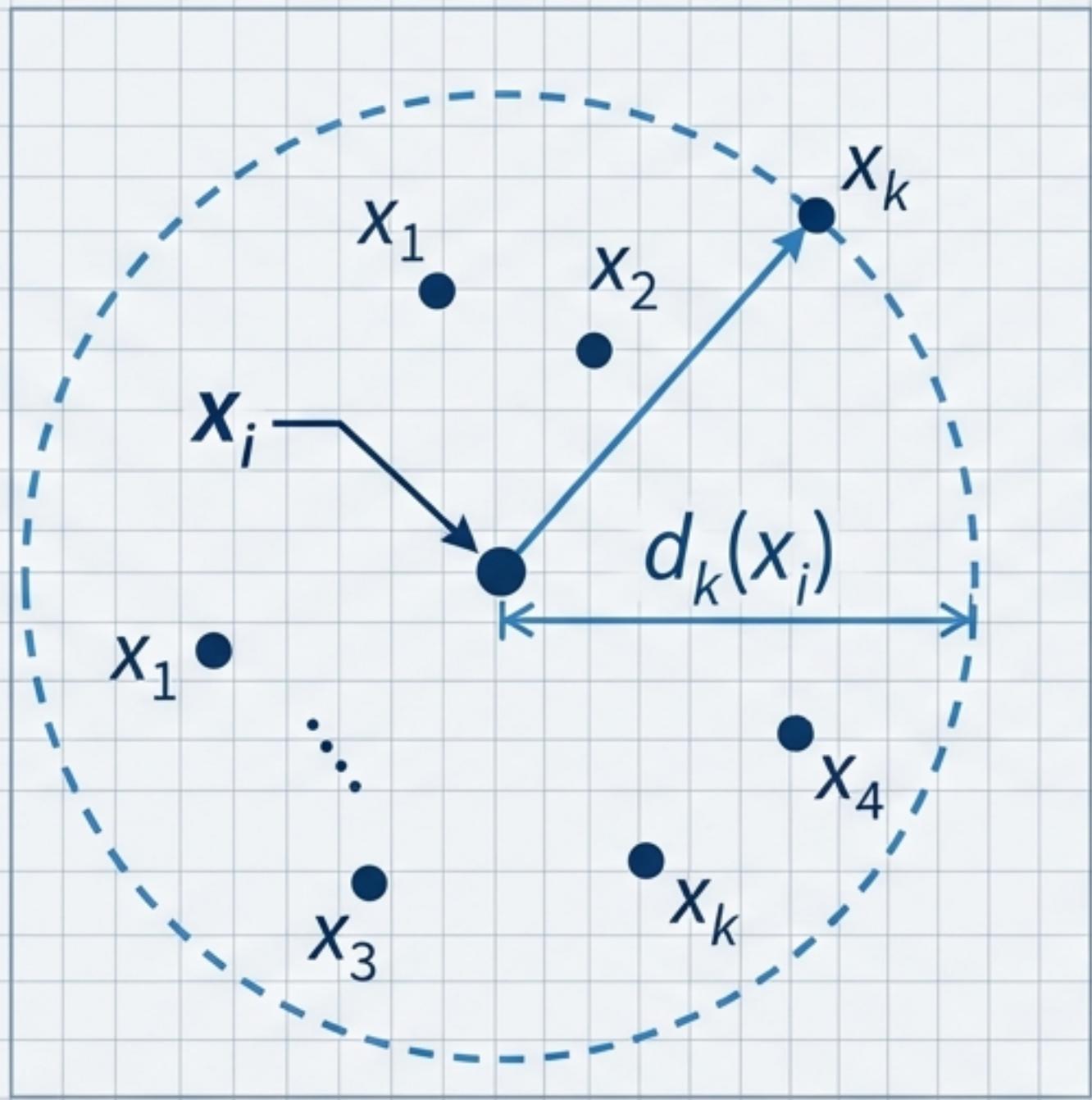


局部觀點 (Local View)



「異常點的局部密度明顯低於其鄰近點的局部密度。」

機制詳解 (1)：定義檢查範圍 (k-Distance)



定義 Definition

k -distance: 第 k 個最近鄰居的距離

k -neighborhood (N_k)：圓圈內的所有點

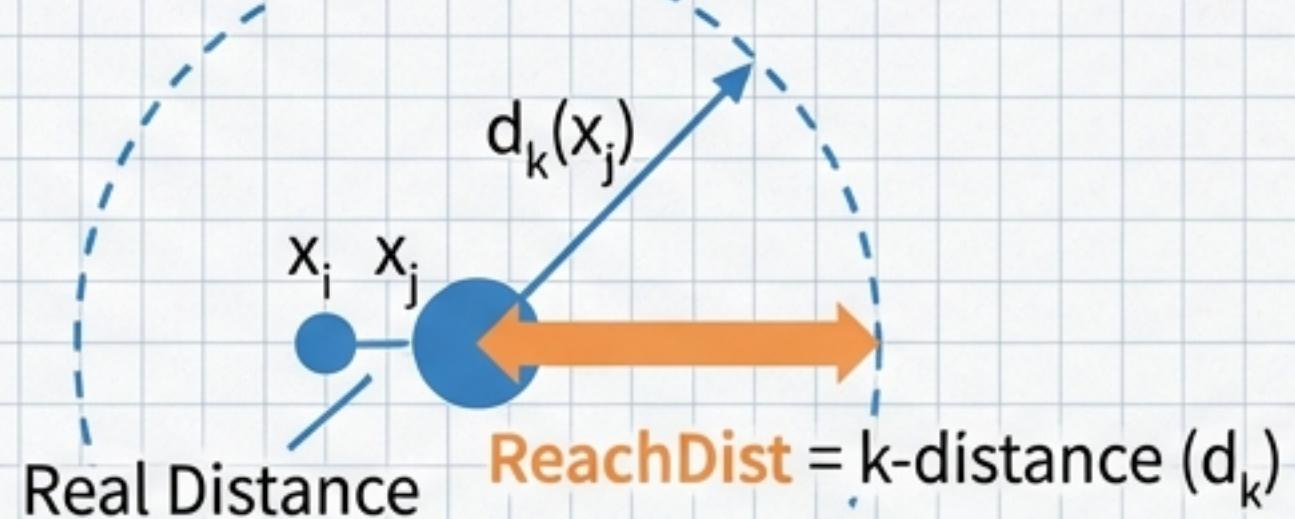
工程類比：為每個感測器讀值定義動態
的檢查半徑

$$d_k(x_i) = \text{distance}(x_i, x_i^{(k)})$$

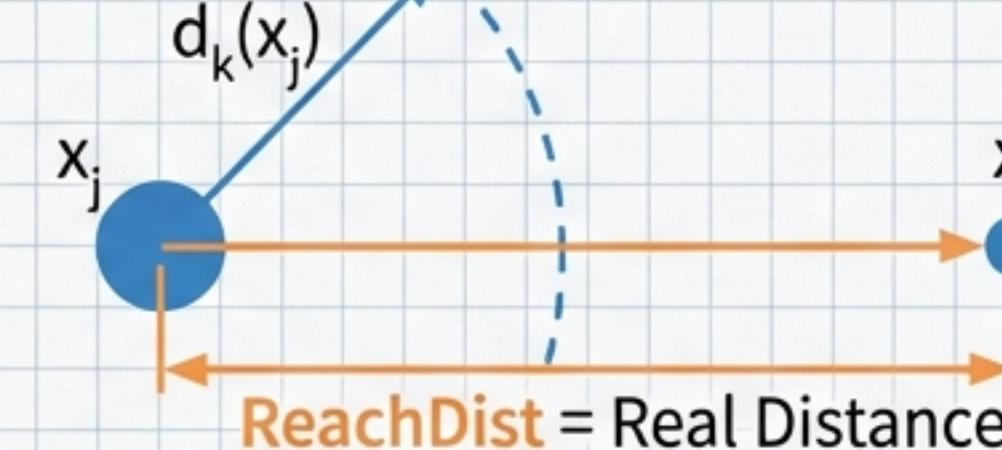
機制詳解 (2)：可達距離 (Reachability Distance)

目的：消除統計噪聲，確保密度估計的穩定性 (Stability).

Case A: Inside Core (Close)



Case B: Outside Core (Far)



$$\text{reach-dist}_k(x_i, x_j) = \max\{d_k(x_j), \text{distance}(x_i, x_j)\}$$

機制詳解 (3) : LRD 與 LOF 得分計算



Python 實作：Scikit-Learn 工具箱

```
from sklearn.neighbors import LocalOutlierFactor

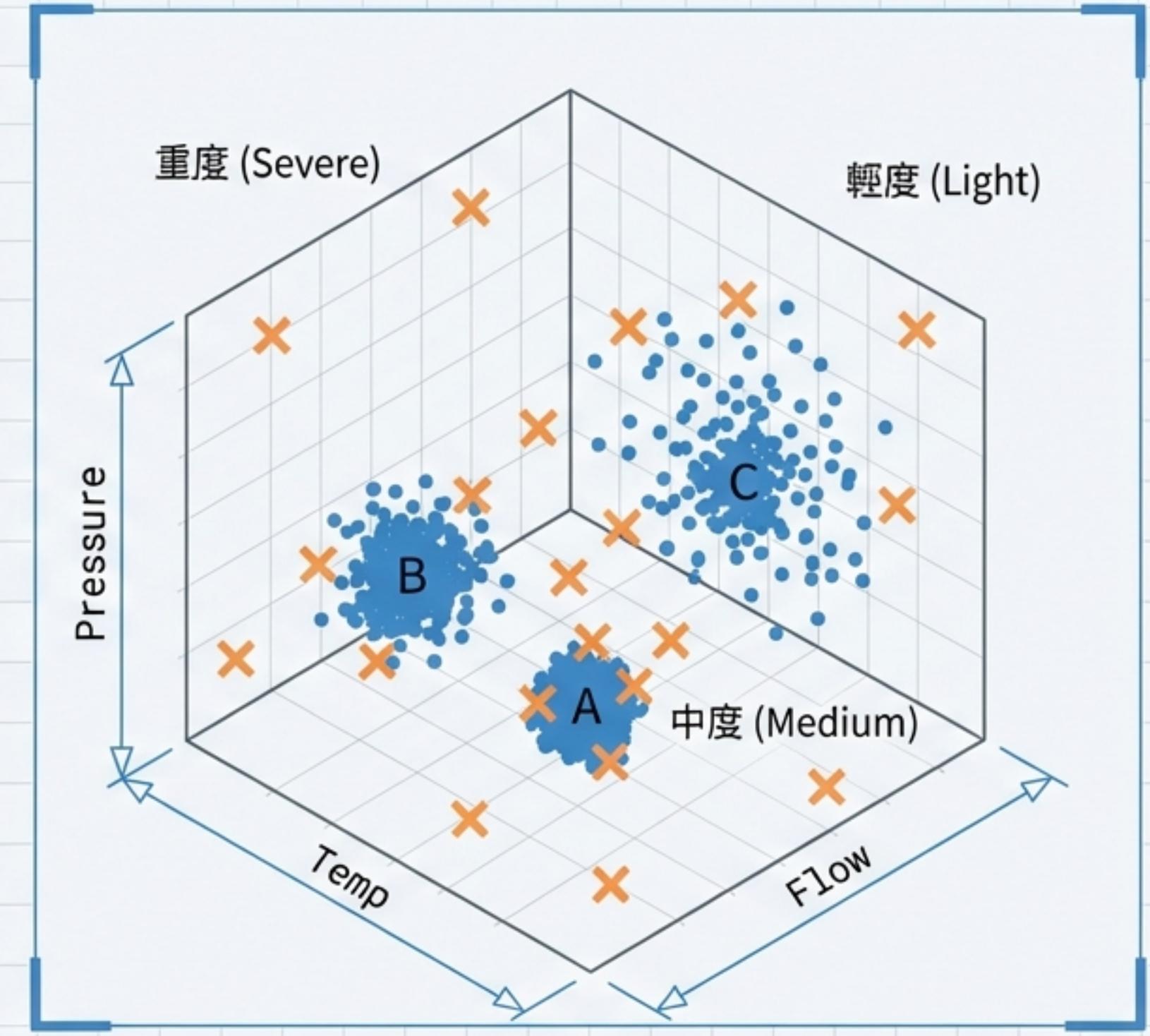
# 初始化模型：針對多密度化工數據
lof = LocalOutlierFactor(
    n_neighbors=15, # 關鍵參數 k (檢查範圍) → 決定敏感度的核心參數
    contamination='auto', # 自動決定異常比例
    novelty=False # False: 檢測歷史數據 (Outlier Detection)
    # True: 檢測新數據 (Novelty Detection)
)

# 訓練並取得結果 (-1 為異常，1 為正常)
y_pred = lof.fit_predict(X_train)
lof_scores = -lof.negative_outlier_factor_ # 轉為正數得分
```

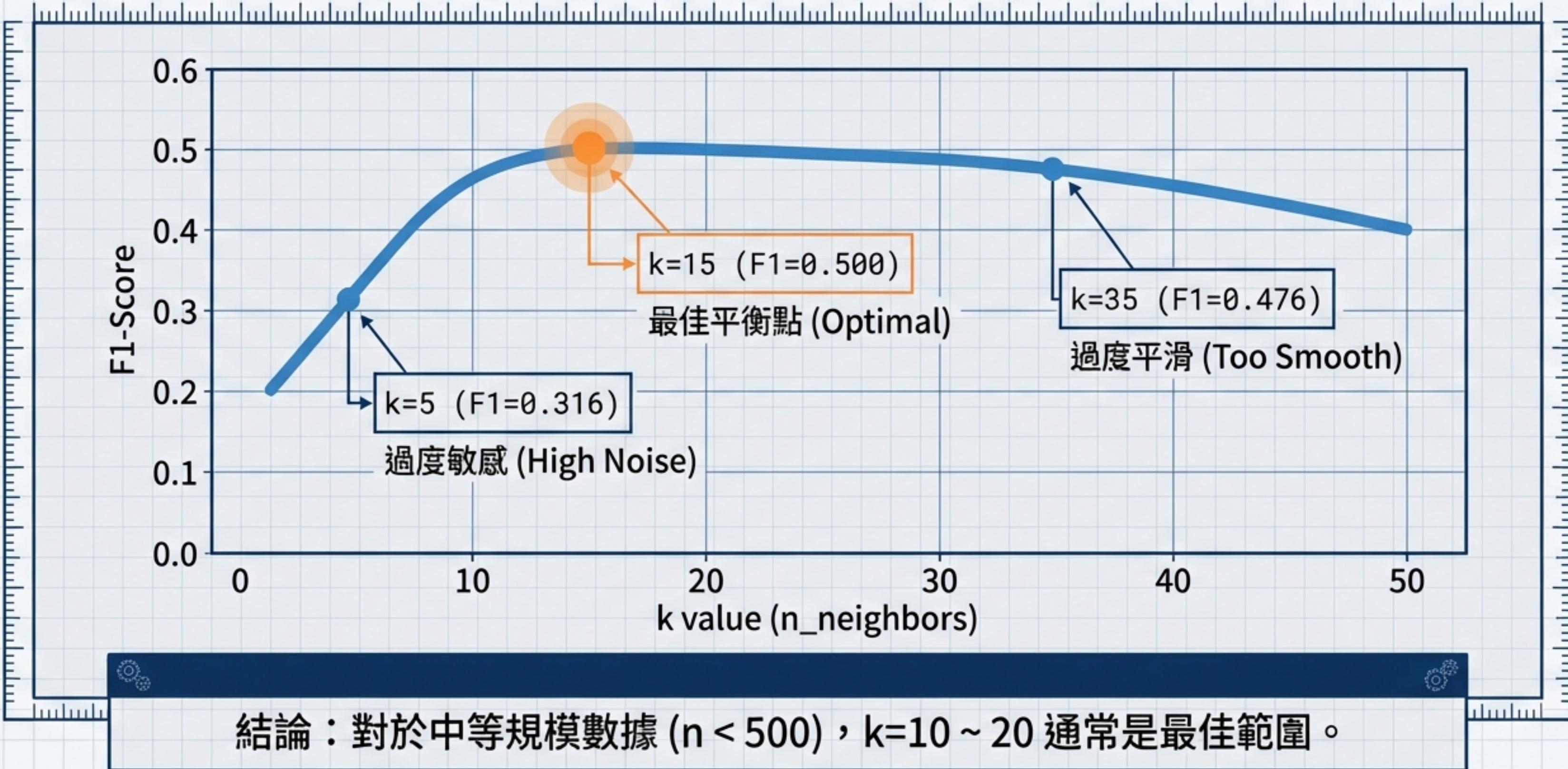
實驗設計：批次反應器模擬 (Batch Reactor Simulation)

數據生成參數

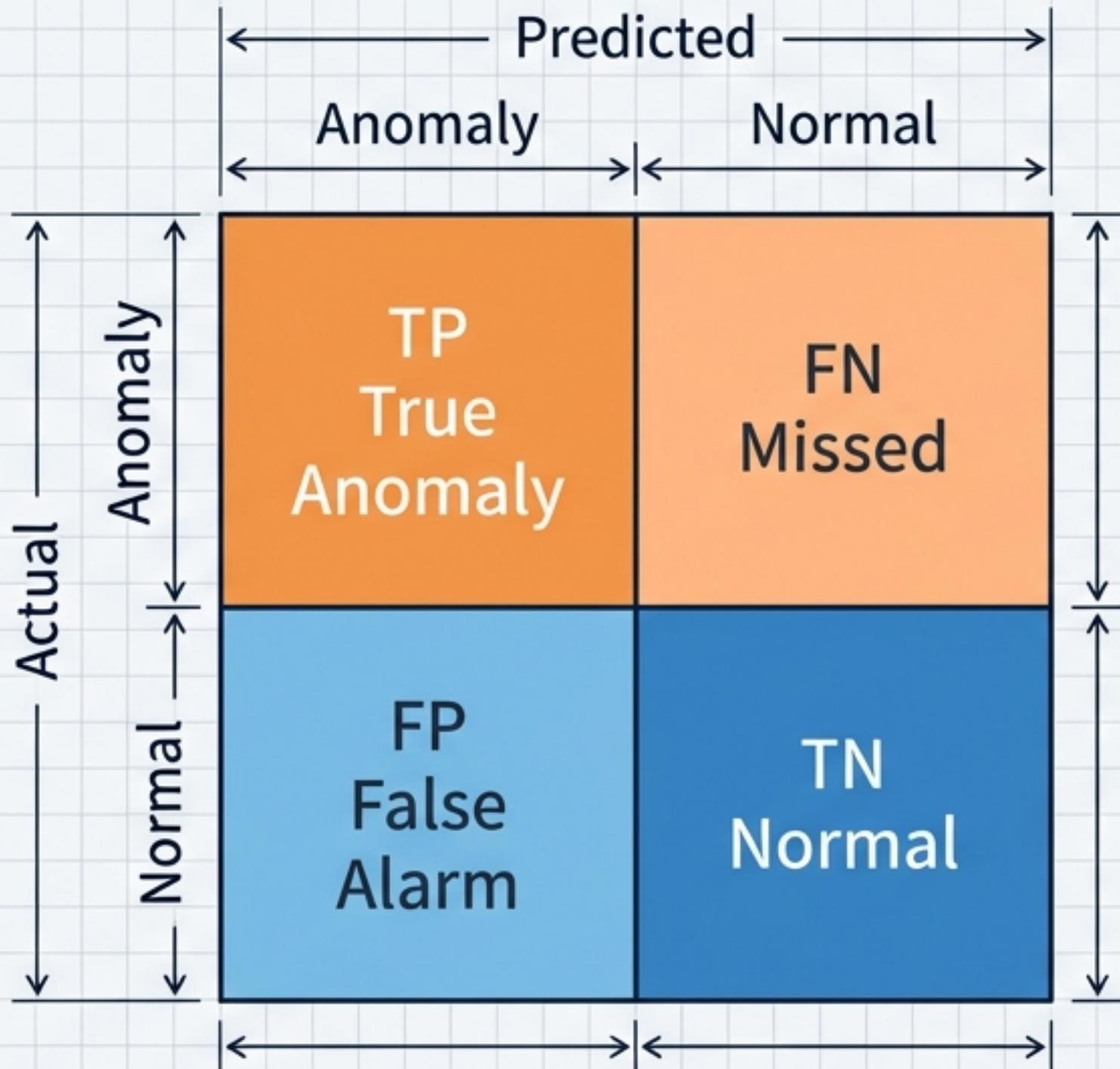
- 配方數量: 3 (Dense, Medium, Sparse)
- 正常數據變異性: +30-40% (模擬真實噪聲)
- 異常設計 (Hierarchical Anomalies):
 - 輕度 (Light): 40%
 - 中度 (Medium): 40%
 - 重度 (Severe): 20%



超參數分析：尋找最佳的 'k' (n_neighbors)



效能評估：真實數據的考驗



Model Metrics (k=15)

Precision: 0.500

50% Reliability

Recall: 0.500

Caught 50% anomalies

AUC: 0.835

Good Separability

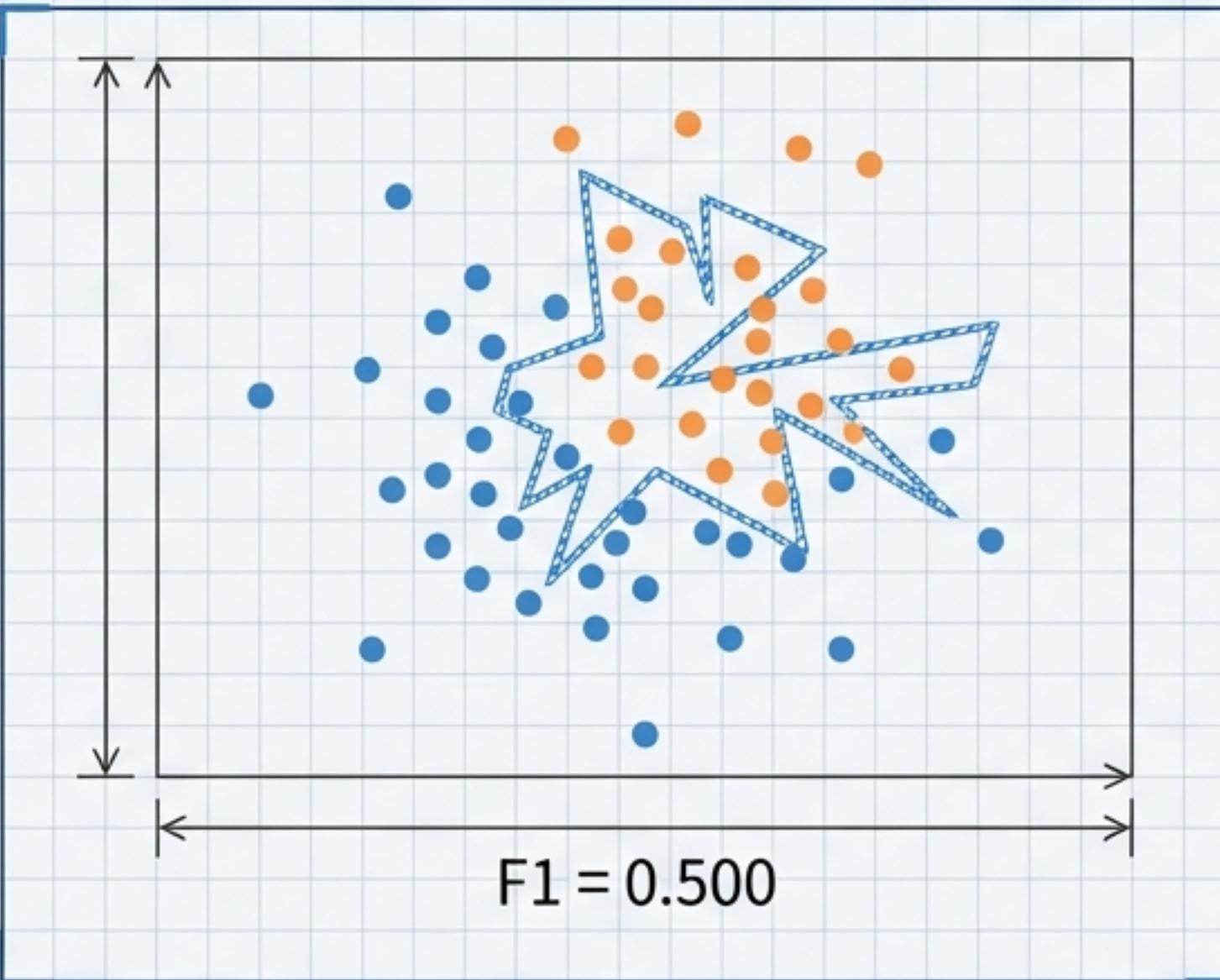
Why not AUC 1.0?

真實工業數據存在模糊邊界與輕微異常。

這是合理的現實結果 (Realistic Outcome).

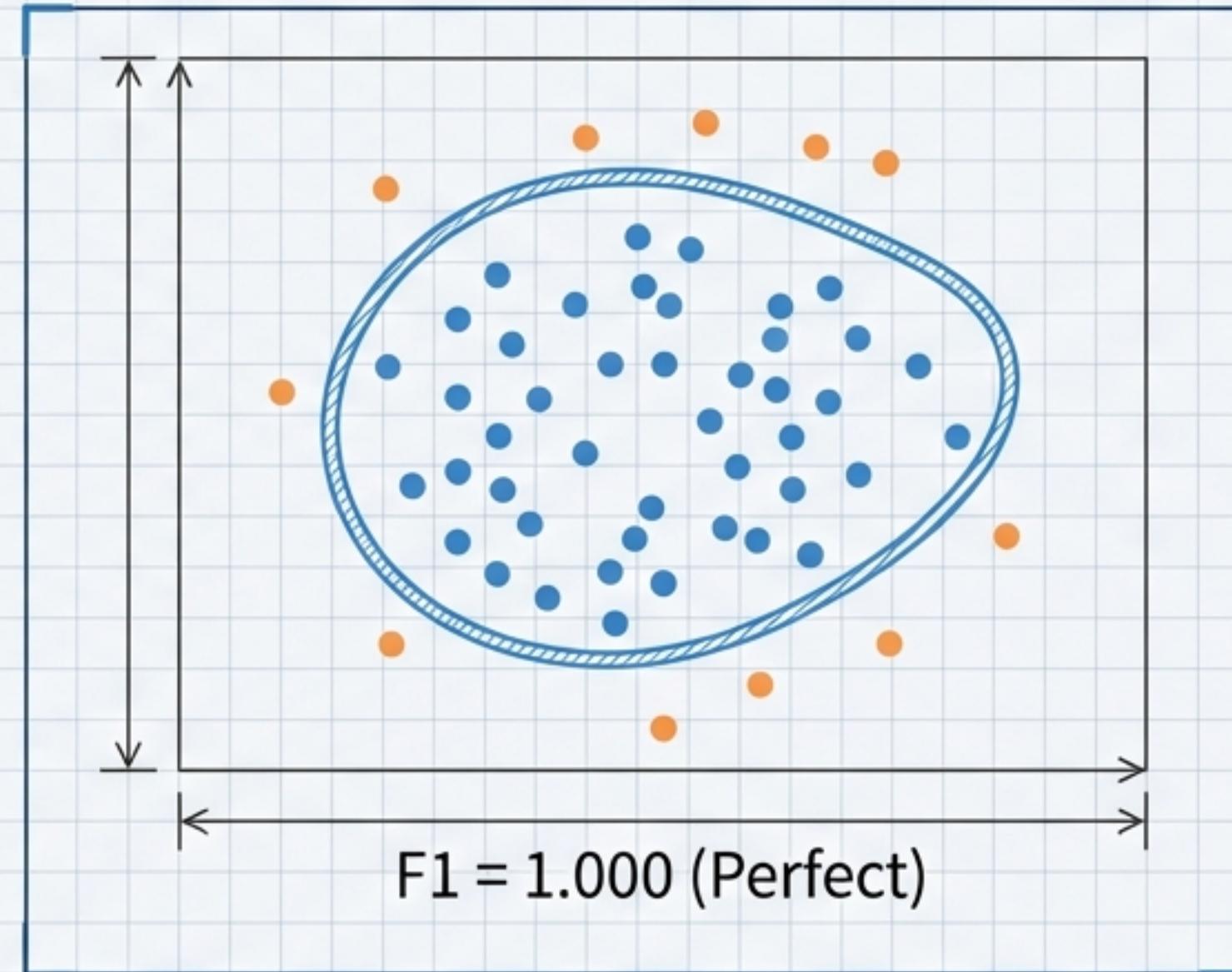
進階應用：新穎性檢測 (Novelty Detection)

Outlier Detection (Mixed Data)



Noto Sans TC

Novelty Detection (Clean History)



Roboto Mono

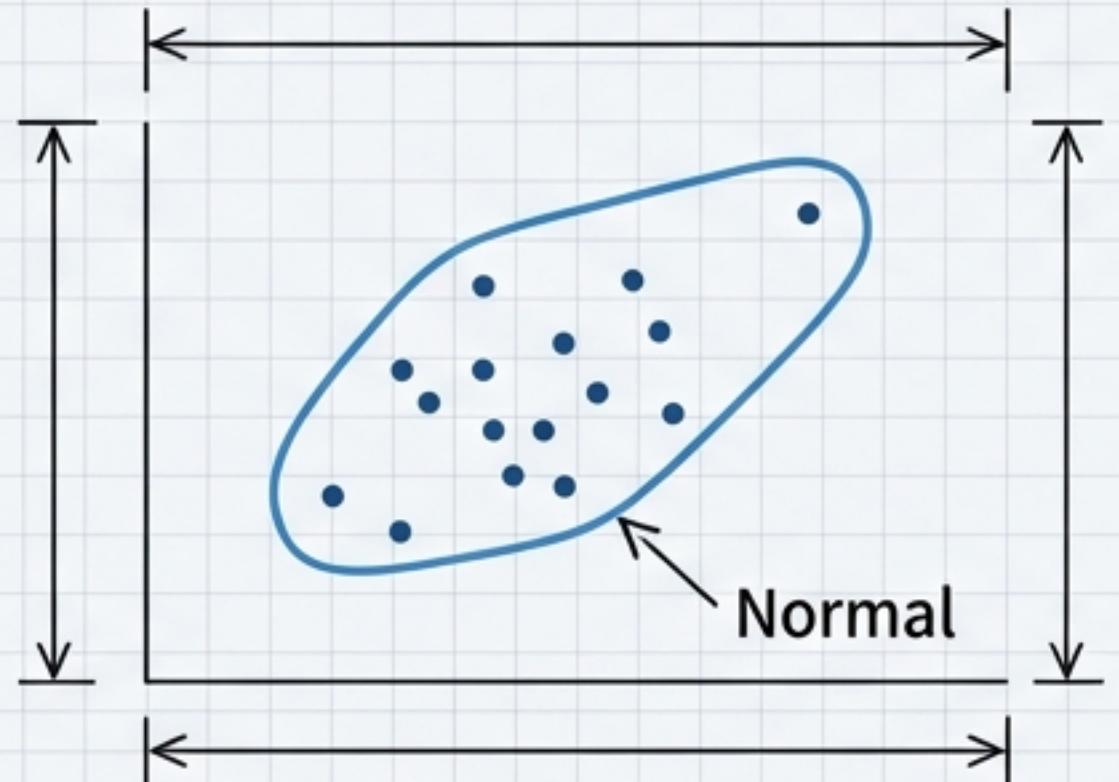
Training on Golden Batch only

應用場景：當擁有純淨的歷史數據時，切換 `novelty=True` 可實現最佳即時監控。

演算法對決：LOF vs. Isolation Forest

LOF (Local Outlier Factor)

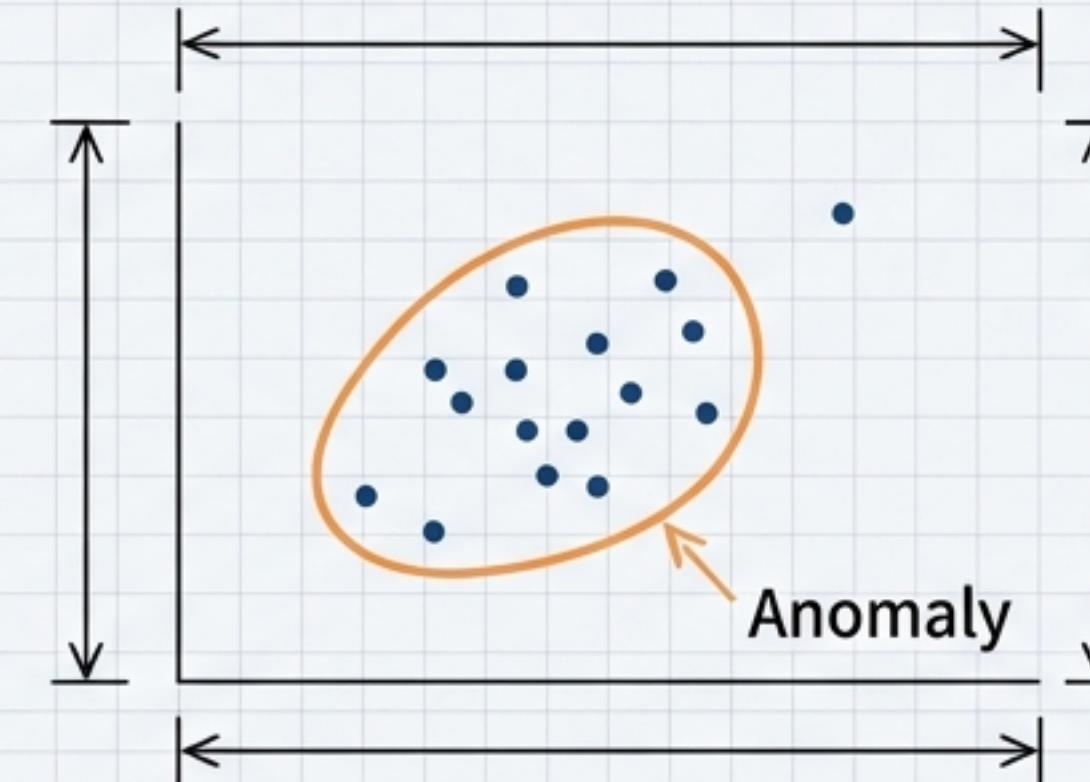
F1 Score: 0.500 ✓



適應局部密度：正確識別稀疏配方為正常。

Isolation Forest

F1 Score: 0.400 !



全局隔離：誤判稀疏配方為異常 (False Alarm)。

Isolation Forest 難以處理多密度 (Multi-Density) 數據集。



策略總結：LOF 的優勢與代價



- ✓ 處理多密度數據 (Multi-Density Handling)
- ✓ 無需假設分佈 (No Gaussian Assumption)
- ✓ 提供量化異常分數 (Quantifiable Score)

- ⚠ 計算成本較高 $O(n^2)$ (Computationally Expensive)
- ⚠ 對參數 k 敏感 (Sensitive to Hyperparameters)
- ⚠ 高維度表現下降 (Curse of Dimensionality)

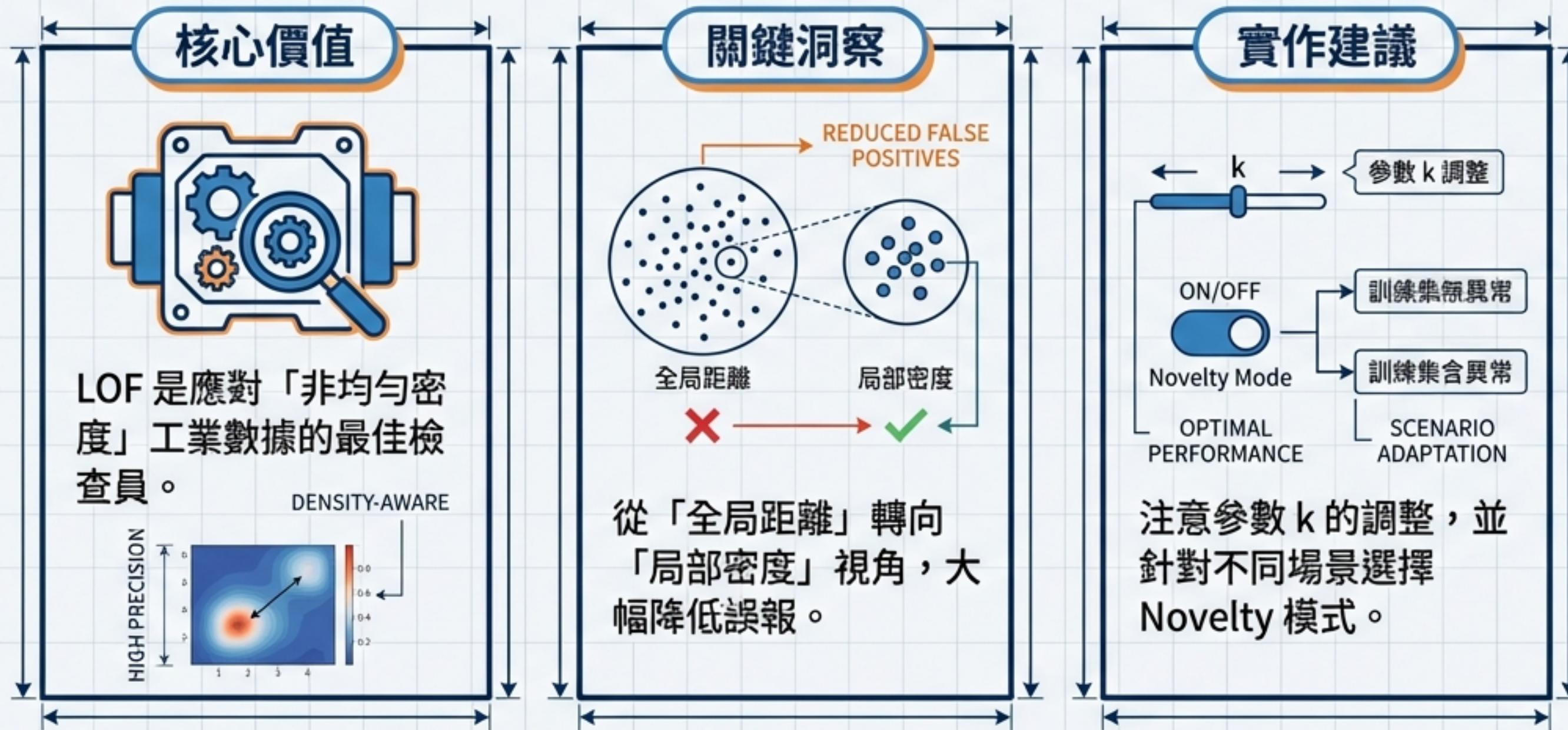
最佳適用：複雜、中等規模的製造業數據集
(Complex, Medium-sized Manufacturing Data)。

實務部署指南 (Deployment Guide)



Code is read much more often than it is written. — 請保持代碼與模型的可維護性。

結論與下一步



START CODING

