

AI 的未來與機器學習的基石

廖晨鈞 112652011

December 4, 2025

1 AI 的未來能力

我認為 20 年後 AI 能做到、且對人類有重大意義的事是：針對患者的特定疾病狀態，從零開始設計出全新的藥物分子，並在該患者的人體模型中完整模擬其療效及分析藥物機理，精準治療患者的症狀。

這項能力將徹底改變現有的藥物開發模式。透過整合基因組學、蛋白質組學與臨床數據，AI 能突破傳統試錯法的限制，針對個人獨特的生理特徵，設計出副作用最小且療效最佳的專屬療法。此外，結合數位孿生技術進行虛擬試驗，能大幅縮短藥物研發週期，降低臨床風險，真正實現精準醫療的終極目標。

2 所需的成分與資源

資料 (Data)：

- 醫學數據：患者的基因體 (Genomics)、轉錄體 (Transcriptomics)、蛋白質結構。
- 高品質標註與反應數據：需要大量藥物與蛋白質交互作用的真實濕實驗數據 (Wet-lab data) 作為 Ground Truth。
- 化學空間數據：包含大量已知可合成的分子的化學屬性資料。

工具 (Tools)：

- 幾何深度學習 (Geometric Deep Learning)：針對分子與蛋白質的三維結構進行建模，精準預測藥物與標靶在空間中的結合構象與交互作用。
- 生成式模型 (Diffusion Models)：生成符合物理法則的原子級 3D 構象，確保藥物在真實環境中的結構穩定性與生物有效性。

硬體與環境 (Hardware / Environment)：

- 數位孿生模擬環境 (Digital Twin Environment)：一個高度仿真的虛擬生理環境，允許 AI 進行數百萬次的虛擬臨床試驗。

學習架構 (Learning Setup)：

- 強化學習 (RLHF)：AI 生成藥物後，由環境（數位孿生）給予反饋（有效性/毒性），並結合人類專家的審核來優化生成策略。

3 涉及的機器學習類型

此任務涉及非監督式學習與強化學習的組合。

- 理由 (Why):
 - 非監督式學習 (Unsupervised/Generative): 用於學習化學分子的語法與分佈 (如 VAE 或 Diffusion), 以此來創造新的藥物分子。
 - 強化學習 (RL): 因為藥物設計的目標空間接近無限大, 且需滿足高療效、低毒性、可合成等多重目標。我們無法標註所有分子的好壞, 必須讓 AI 在化學空間中探索, 根據模擬結果來優化其設計策略。
- 資料來源: 初始是現有的藥物數據庫 (如 ChEMBL, PubChem), 後期則利用生成模型創造的新分子結構進行自我訓練。
- 目標訊號 (Reward Signal): 來自數位人體模型的評估分數。

$$\text{Reward} = (+\text{療效分數}) - (\text{毒性分數}) - (\text{合成難度係數})$$

- 環境互動: AI 提出分子 \rightarrow 投入數位模型模擬 \rightarrow 獲得生理反應反饋 \rightarrow AI 修正設計。

4 簡化模型問題

為了實現 De Novo Design (從頭藥物設計) 的雛形, 我們使用 Variational Autoencoder (VAE, 變分自編碼器) 作為核心架構。

與傳統機器學習不同, VAE 是分子生成領域的經典模型, 其核心能力在於學習數據的潛在空間 (Latent Space), 進而具備創造新分子的能力。

4.1 問題設計 (Problem Design)

說明與對應目標

為了實現 20 年後精準設計藥物的願景, 首要任務是讓 AI 學會分子的化學語言。我們的目標是讓模型從現有藥物數據中歸納出分子組成的語法與規律。

定義

- 資料形式 (Data Format): 來自 MOSES 資料集的藥物分子, 經過 SMILES \rightarrow SELFIES \rightarrow Token ID 的預處理, 並進行 Padding 至固定長度 (64)。
- 輸入 (Input): 經 Token 化處理的 SELFIES 整數序列, 形狀為 (Batch_Size, 64)。
- 輸出 (Output): 重建後的 SELFIES 分子序列機率分佈, 形狀為 (Batch_Size, 64, Vocab_Size), 代表每個位置出現各個化學符號的機率。
- 任務目標 (Task): 訓練模型將離散分子壓縮為連續潛在向量, 並能從中解碼出新分子。

4.2 模型與方法 (Model & Method)

我們選擇 LSTM-based Variational Autoencoder (VAE) 作為 Toy Model 的模型。理由如下：

1. **LSTM:** 化學分子的字串表示本質上是序列數據 (Sequential Data)，類似於自然語言。LSTM (長短期記憶網路) 擅長捕捉序列中的長距離依賴關係 (例如分子開頭的環結構需要在結尾閉合)，非常適合處理化學語法。
2. **VAE:** 傳統的 Autoencoder 只能壓縮數據，其潛在空間不連續，無法進行生成。VAE 透過引入 KL 散度，強迫潛在空間呈現常態分佈 (Gaussian)。這創造了一個連續且平滑的空間，讓我們可以透過隨機採樣 (Random Sampling) 或插值 (Interpolation) 來探索未知的化學空間，進而發現新藥物。

- 模型架構：

- **Encoder:** 雙層 LSTM，將輸入序列編碼為隱藏狀態，再映射到潛在空間的均值 (μ) 和對數變異數 ($\log \sigma^2$)。
- **Latent Space:** 128 維的向量空間，使用 Reparameterization Trick 進行採樣。
- **Decoder:** 雙層 LSTM，接收潛在向量 z (重複擴展至序列長度) 作為輸入，逐步還原分子序列。

- Loss Function：

$$Loss = \text{Cross Entropy} + \beta \cdot \text{KL Divergence}$$

- Cross Entropy (重建損失)：衡量生成的分子與原始分子是否相同。
 - KL Divergence：衡量潛在分佈與標準常態分佈 $N(0, 1)$ 的差異， β 為 KL 權重。
- 訓練策略：採用 **KL Annealing** (退火) 策略。前 2 個 Epoch KL 權重為 0，隨後線性增加至 0.05。這防止了模型在初期忽略潛在變數造成後驗失效 (Posterior Collapse)，確保 Latent Space 既具有編碼能力又符合常態分佈。

4.3 實作與結果 (Results)

- 訓練過程：模型在大約 Epoch 9 後收斂，訓練與測試損失皆穩定下降至約 66.4。



Figure 1: 訓練過程與 Loss 曲線

- 生成與篩選機制：訓練完成後，模型通過以下步驟生成新分子：
 1. 隨機採樣：從標準常態分佈 $N(0, 1)$ 中隨機抽取潛在向量 z 。
 2. 解碼：Decoder 將 z 轉換成機率分佈，並使用 溫度採樣 (Temperature Sampling) (設為 1.0) 來調整生成的多樣性。
 3. 格式轉換：Index List \rightarrow SELFIES \rightarrow SMILES \rightarrow RDKit Mol 物件。
 4. 規則篩選：
 - 原子數限制：過濾掉重原子數小於 12 的分子 (太小的分子通常不是好藥物)。
 - 環的限制：
 - * 剔除含有大環 (> 7 元環) 的分子。
 - * 剔除只含有不穩定小環 (< 5 元環) 的分子 (保留至少有一個 5, 6, 7 元環的分子)。
 5. QED 評分：計算 QED (Quantitative Estimate of Drug-likeness) 分數，這是一個衡量分子是否適合作為藥物的綜合指標 (範圍 0~1)。
 6. 結果展示：在 500 個生成的分子中，篩選出了 70 個符合標準的分子，再從中挑選 QED 分數最高的 5 個分子畫出來。

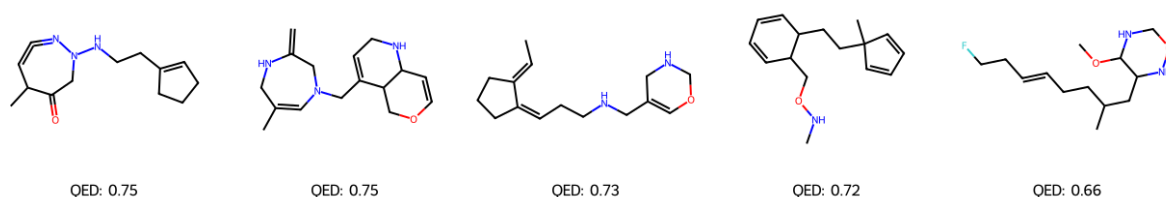


Figure 2: 生成的分子結構示意圖

4.4 討論 (Discussion)

這個實驗展示了 AI 從數據中學習化學語法的能力。透過 VAE，我們確實能創造出資料集中不存在的新分子。

雖然我們可以生成符合化學規則的分子，但模型完全不知道這些分子是否真的具備療效。目前的篩選僅基於簡單的物理規則 (如環的大小、藥物相似度)，並未將結合親和力、代謝率、脂溶性、毒性等因素納入考慮。

在 20 年後的願景中，我們不能只依賴隨機採樣，而是要基於患者的情況、基因組、蛋白質組，設計一個特製的藥物分子。此外，我們也需要將數位人體模型的反饋訊號 (如結合親和力、代謝率) 回傳給生成模型，這需要引入強化學習 (RL) 或條件生成 (Conditional Generation)。

使用 AI(Gemini) 協助生成 Python 程式碼以及潤飾報告文句。