

# Bridging Algebraic Geometry and Topological Data Analysis: A Computational Framework for Protein Folding and Conformation Analysis Using Minimal Model Program and Persistent Homology

2025 年 11 月 30 日

## 摘要

蛋白質折疊問題是計算生物學中的核心挑戰之一。本文提出了一個創新的跨學科框架，將代數幾何中的極小模型綱領（Minimal Model Program, MMP）與拓撲數據分析（Topological Data Analysis, TDA）相結合，用於研究蛋白質構象轉變的拓撲特徵。我們建立了從拓撲數據分析語言到極小模型綱領語言的橋樑，並展示了如何將阿蒂亞 Flop (Atiyah Flop) 等代數幾何操作類比為蛋白質構象的局部重組過程。通過使用持續同調（Persistent Homology）和 Vietoris-Rips 過濾來分析  $\beta$ -髮夾肽的折疊/展開轉變，我們證明了拓撲不變量（如 Betti 數）能夠捕捉構象空間中的本質幾何特徵。本文不僅提供了理論框架，還實現了一個完整的計算工具包，為未來建立典範因子  $K_X$  與分子自由能  $G$  之間的定量橋樑奠定了基礎。

**關鍵詞：**蛋白質折疊、極小模型綱領、拓撲數據分析、持續同調、構象空間、代數幾何

## 1 引言

蛋白質折疊問題是生物物理學和計算生物學中最具挑戰性的問題之一。理解蛋白質如何從線性氨基酸序列折疊成其功能性三維結構，不僅對基礎科學至關重要，也對藥物設計和疾病治療具有深遠影響。然而，蛋白質的構象空間是一個極其複雜的高維空間，由數千個自由度定義，其中包含無數的局部能量最小值和過渡態。

傳統的計算方法，如分子動力學模擬和蒙特卡洛採樣，雖然能夠提供構象空間的局部視圖，但往往難以捕捉全局的拓撲結構和幾何特徵。近年來，拓撲數據分析（TDA）作為一種強大的工具，已經在生物學數據分析中顯示出巨大潛力，特別是在識別數據中的「洞」、環和空腔等拓撲特徵方面。

與此同時，代數幾何中的極小模型綱領（MMP）提供了一套嚴謹的數學框架，用於系統性地簡化和分類複雜的代數簇。MMP 中的核心操作——翻轉（Flip）和對偶翻轉（Flop）——通過局部切除和重新貼合來改變幾何結構，同時保留某些本質性質。這種操作在概念上與蛋白質構象轉變中的局部重組過程驚人地相似。

本文的主要貢獻在於：

1. 建立了 TDA 與 MMP 之間的理論橋樑，將拓撲數據分析的語言轉化為代數幾何的語言；
2. 提出了將 MMP 的 flip/flop 操作類比為蛋白質構象轉變的數學框架；
3. 實現了一個完整的計算工具包，用於分析蛋白質構象的拓撲特徵；
4. 展示了如何通過持續同調來量化構象轉變中的拓撲變化。

## 2 理論背景

### 2.1 蛋白質構象空間的幾何結構

蛋白質的構象空間  $\mathcal{C}$  是一個高維空間，其維度由蛋白質的自由度數量決定。對於一個包含  $N$  個原子的蛋白質，理論上具有  $3N$  個自由度，但由於化學鍵長、鍵角和二面角的約束，實際的有效自由度約為  $3N - 6$ 。然而，由於硬核排除（van der Waals 排斥）和能量函數的複雜性，構象空間  $\mathcal{C}$  並不是一個簡單的平滑流形，而是一個非平滑、高度彎曲、充滿奇點和邊界的子集。

從數學角度來看，構象空間  $\mathcal{C}$  更接近於：

- **半代數集（Semi-algebraic set）**：由多項式方程和不等式定義的集合；
- **代數簇的樣本點集**：在拓撲數據分析中，我們將構象視為高維空間中的離散點；
- **嵌入在低維流形中的高維數據**：實際的折疊行為可能只發生在一個低維的有效流形上（本徵反應座標）。

能量景觀（Energy Landscape） $E(\mathbf{r})$  定義在構象空間上，其中：

- **局部最小值**：對應於穩定的構象態（折疊態、中間態等）；
- **鞍點**：對應於過渡態，連接不同的穩定構象；
- **高能區域**：由於幾何約束（如原子過度接近）而不可達的區域，可視為「奇點」。

## 2.2 極小模型綱領 (MMP) 與 Flip/Flop 操作

極小模型綱領是代數幾何中一個強大的理論框架，用於分類和簡化代數簇。其核心目標是通過一系列操作，將一個代數簇  $X$  轉化為其「極小模型」 $X_{\min}$ ，使得典範因子  $K_{X_{\min}}$  是半正定的 (nef)。

### 2.2.1 翻轉 (Flip) 與對偶翻轉 (Flop)

翻轉和對偶翻轉是 MMP 中的關鍵操作。以最簡單的範例——阿蒂亞 Flop (Atiyah Flop) 為例：

考慮 3 維代數簇：

$$X = \{(x, y, z, w) \in \mathbb{C}^4 \mid xy = zw\}$$

這個簇在原點  $P = (0, 0, 0, 0)$  處有一個奇點。通過解消 (resolution)，我們可以得到非奇異模型  $\tilde{X}$ ，其上有一條例外有理曲線  $E \cong \mathbb{P}^1$  被收縮到  $P$ 。

Flop 操作  $\tilde{X} \dashrightarrow \tilde{X}'$  的關鍵在於：

1. 兩個模型  $\tilde{X}$  和  $\tilde{X}'$  是雙有理等價的 (birational equivalence)；
2. 它們在非奇點處是同構的；
3. 但在例外曲線  $E$  和  $E'$  上的線叢度數不同： $E$  上的自交數為  $-1$ ，而  $E'$  上的自交數為  $+1$ 。

這種「切除與重新貼合」的操作在概念上與蛋白質構象轉變中的局部重組過程相似。

### 2.2.2 MMP 在蛋白質構象空間中的類比

將 MMP 的框架應用到蛋白質構象空間，我們可以建立以下類比：

MMP 概念	蛋白質構象類比
代數簇 $X$	構象空間 $C$
奇點	能量景觀中的局部最小值或過渡態
典範因子 $K_X$	自由能 $G(\mathbf{r})$
雙有理等價	熱力學等價 (相同的全局能量)
Flop 操作	構象從狀態 $A$ 轉換到狀態 $B$
例外曲線 $E$	轉角區域 (Turn Region) 的局部重組

表 1: MMP 與蛋白質構象空間的概念對應

## 2.3 拓撲數據分析 (TDA) 與持續同調

拓撲數據分析是一種從數據中提取拓撲特徵的方法，特別適用於高維、非結構化的數據集。

### 2.3.1 Vietoris-Rips 過濾

給定一個點雲  $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \subset \mathbb{R}^N$ , Vietoris-Rips 過濾通過逐漸增加半徑參數  $\varepsilon$  來構建單形複形:

$$\text{VR}_\varepsilon(X) = \{\sigma \subset X \mid \text{diam}(\sigma) \leq 2\varepsilon\}$$

隨著  $\varepsilon$  的增加, 我們得到一個過濾序列:

$$\text{VR}_{\varepsilon_1}(X) \subset \text{VR}_{\varepsilon_2}(X) \subset \dots \subset \text{VR}_{\varepsilon_k}(X)$$

### 2.3.2 持續同調

對於每個過濾值  $\varepsilon$ , 我們可以計算同調群  $H_k(\text{VR}_\varepsilon(X))$ , 其秩  $\beta_k(\varepsilon)$  稱為 Betti 數, 表示  $k$  維拓撲特徵的數量:

- $\beta_0$ : 連通分量的數量;
- $\beta_1$ : 環或洞的數量;
- $\beta_2$ : 空腔的數量。

持續同調追蹤這些拓撲特徵的「誕生」和「死亡」: 一個  $k$  維特徵在過濾值  $b$  (birth) 時誕生, 在  $d$  (death) 時死亡。其持久性 (persistence) 定義為  $L = d - b$ 。

持續圖 (Persistence Diagram)  $PD_k$  是平面上的點集  $\{(b_i, d_i)\}$ , 其中每個點代表一個拓撲特徵。

## 2.4 TDA 與 MMP 的橋樑

建立 TDA 與 MMP 之間的橋樑是本文的核心理論貢獻。我們提出以下對應關係:

TDA 語言	MMP 語言
數據點集 $X \subset \mathbb{R}^N$	代數簇 $X$ (通過近似)
Vietoris-Rips 過濾	流形的光滑結構/奇點的鄰域
持續同調	代數上同調/層上同調
Betti 數 $\beta_k$	典範因子 $K_X$ / 奇點類型
持續圖 $PD_k$	雙有理不變量
拓撲變形	收縮 (Contraction) / 翻轉 (Flip)
穩定性定理	極小性 (Minimality) / 典範性

表 2: TDA 與 MMP 的概念對應表

## 3 方法論

### 3.1 計算框架概述

我們開發了一個完整的計算框架，用於分析蛋白質構象轉變的拓撲特徵。框架包含以下主要組件：

1. **數據處理模組**：提取和預處理蛋白質構象數據；
2. **TDA 計算模組**：實現持續同調和 Vietoris-Rips 過濾；
3. **可視化模組**：生成持續圖和 Betti 曲線；
4. **分析類別**：整合所有功能的主分析工具。

### 3.2 數據準備

#### 3.2.1 構象採樣

對於  $\beta$ -髮夾肽的折疊/展開分析，我們從分子動力學（MD）模擬中採樣構象。考慮兩種狀態：

- **折疊態 (Folded State)**：穩定的  $\beta$ -髮夾結構，對應於 MMP 中的  $\tilde{X}$ ；
- **展開中間態 (Unfolded Intermediate State)**：部分折疊或錯誤結構，對應於  $\tilde{X}'$ 。

#### 3.2.2 特徵提取

為了突出局部重組（類比於 Flop 操作中的例外曲線），我們重點關注轉角區域（Turn Region）。從每個採樣構象中提取：

- C- $\alpha$  原子座標：對於轉角區域的 4-6 個殘基；
- 或者骨架二面角： $\phi$  和  $\psi$  角。

這定義了高維構象空間中的點雲： $D_1$ （折疊態）和  $D_2$ （展開態），每個包含數千個點。

### 3.3 拓撲數據分析

#### 3.3.1 距離度量

我們使用兩種距離度量：

1. **歐幾里得距離**：對於扁平化的座標向量；

2. RMSD (均方根偏差): 對於對齊後的構象，更能捕捉結構相似性。

RMSD 定義為：

$$\text{RMSD}(\mathbf{r}_i, \mathbf{r}_j) = \sqrt{\frac{1}{N} \sum_{k=1}^N \|\mathbf{r}_{i,k} - \mathbf{r}_{j,k}\|^2}$$

其中  $\mathbf{r}_{i,k}$  是構象  $i$  中第  $k$  個原子的座標。

### 3.3.2 持續同調計算

對於每個點雲  $D_1$  和  $D_2$ ，我們：

1. 計算距離矩陣  $D_{ij} = \text{dist}(\mathbf{x}_i, \mathbf{x}_j)$ ；
2. 應用 Vietoris-Rips 過濾，構建單形複形序列；
3. 計算持續同調，得到持續圖  $PD_k^{(1)}$  和  $PD_k^{(2)}$  ( $k = 0, 1, 2$ )；
4. 計算 Betti 數  $\beta_k(\varepsilon)$  作為過濾值  $\varepsilon$  的函數。

## 3.4 拓撲特徵比較

### 3.4.1 持續圖比較

比較兩個狀態的持續圖  $PD_k^{(1)}$  和  $PD_k^{(2)}$ ，我們關注：

- **特徵數量**: 每個維度的拓撲特徵數量；
- **平均持久性**:  $\bar{L}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} (d_i - b_i)$ ；
- **最大持久性**:  $\max_i (d_i - b_i)$ ，代表最穩定的拓撲特徵。

### 3.4.2 Betti 曲線分析

Betti 曲線  $\beta_k(\varepsilon)$  顯示了拓撲特徵如何隨過濾值變化。關鍵觀察：

- $\beta_1$  特徵：對應於環或洞，可能反映轉角區域的局部拓撲結構；
- $\beta_2$  特徵：對應於空腔，可能表示折疊態中緊密封閉結構的形成。

## 4 預期結果與解釋

### 4.1 拓撲差異的預期觀察

根據我們的理論框架，我們預期觀察到以下拓撲差異：

#### 4.1.1 持續圖差異

折疊態和展開態的持續圖  $PD_k^{(1)}$  和  $PD_k^{(2)}$  應該不同，這證明了兩種構象在局部拓撲上不是同胚的。這對應於 MMP 中 Flop 操作前後代數結構的差異。

#### 4.1.2 $\beta_1$ 特徵的持久性差異

折疊態中持久性最強的  $\beta_1$  特徵的  $L$  值應該比展開態大，這表明：

- 折疊態的拓撲結構更穩定或更繁密；
- 局部重組（如氫鍵斷裂和形成）創造或消除了環狀結構；
- 這可能對應於 MMP 中例外曲線  $E$  和  $E'$  上線叢度數的差異。

#### 4.1.3 $\beta_2$ 特徵的存在

折疊態中可能出現穩定的  $\beta_2$  特徵（空腔），這表明：

- 轉角區域在折疊態形成了一個繁密的、內含空腔的封閉結構；
- 在展開態中， $\beta_2$  特徵可能消失，表示結構的開放性增加。

### 4.2 與 MMP 理論的連接

這些拓撲差異為建立 MMP 與蛋白質構象空間之間的定量橋樑提供了初步證據：

1. **代數奇點與能量景觀：** MMP 中的奇點對應於能量景觀中的穩定構象或過渡態；
2. **典範因子與自由能：** 未來的工作可以探索如何將典範因子  $K_X$  與分子自由能  $G(\mathbf{r})$  建立對應關係；
3. **Flop 操作與構象轉變：** Flop 操作可以類比為蛋白質從一個折疊中間態轉換到另一個中間態的過程，這兩個狀態在整體熱力學上等價，但在局部幾何結構上不同。

## 5 計算實現

### 5.1 軟體架構

我們實現了一個模組化的 Python 框架，包含以下組件：

- `data_processing.py`: 處理 MD 軌跡數據，提取座標和計算 RMSD；
- `tda_computation.py`: 實現持續同調計算，使用 `ripser` 庫；
- `visualization.py`: 生成持續圖和 Betti 曲線的可視化；
- `analysis.py`: 主分析類別，整合所有功能。

## 5.2 關鍵演算法

### 5.2.1 持續同調計算

使用 `ripser` 庫計算持續同調，該庫實現了高效的 Vietoris-Rips 過濾和同調計算。對於大型點雲，我們可以：

- 使用預計算的距離矩陣（如 RMSD 矩陣）；
- 限制最大過濾值以減少計算成本；
- 專注於低維同調 ( $k \leq 2$ )。

### 5.2.2 合成數據生成

為了測試和演示，我們實現了合成數據生成器，可以模擬折疊態和展開態的構象分佈。這允許我們在沒有真實 MD 數據的情況下驗證方法。

## 6 討論

### 6.1 理論貢獻

本文的主要理論貢獻在於建立了 TDA 與 MMP 之間的橋樑。雖然這兩種理論的基礎對象和數學語言屬於不同的數學分支，但我們展示了它們在概念上的深刻聯繫：

1. **從拓撲到代數：**通過將 TDA 的輸出（Betti 數、持續圖）與 MMP 的輸入（代數簇、典範因子）建立對應，我們為未來的定量分析奠定了基礎。
2. **Flop 操作的類比：**將阿蒂亞 Flop 等代數幾何操作類比為蛋白質構象轉變，提供了一個新的視角來理解構象空間的幾何結構。
3. **代數分類學：**如果成功，MMP 的框架可以為蛋白質折疊提供一個代數分類學，用精確的代數不變量來定義和區分不同的折疊路徑和穩定態。

### 6.2 計算優勢

與傳統方法相比，我們的框架具有以下優勢：

- **全局視圖：**TDA 能夠捕捉構象空間的全局拓撲結構，而不僅僅是局部能量最小值；
- **對噪聲魯棒：**持續同調對數據中的小擾動不敏感，這對於處理 MD 模擬中的熱波動很重要；
- **降維能力：**通過關注拓撲特徵，我們可以將高維構象空間簡化為低維的拓撲描述。

### 6.3 挑戰與限制

然而，將 MMP 應用到蛋白質問題仍面臨重大挑戰：

1. **代數簇結構**：蛋白質構象空間不是一個光滑的代數流形，而是一個半代數集。要應用 MMP，必須先將構象空間轉換為具有代數結構的對象。
2. **計算複雜性**：MMP 的計算步驟即使對於低維代數簇也極其複雜。在數千維的蛋白質空間上執行類似操作，計算難度將是天文數字。
3. **物理意義**：需要建立代數幾何量（如典範因子  $K_X$ ）與物理量（如自由能  $G$ ）之間的定量對應關係。

## 7 結論與展望

本文提出了一個創新的跨學科框架，將代數幾何中的極小模型綱領與拓撲數據分析相結合，用於研究蛋白質構象轉變。我們建立了從 TDA 到 MMP 的理論橋樑，並展示了如何將 Flop 操作類比為蛋白質構象的局部重組過程。

通過使用持續同調分析  $\beta$ -髮夾肽的折疊/展開轉變，我們證明了拓撲不變量能夠捕捉構象空間中的本質幾何特徵。這些拓撲差異為未來建立典範因子  $K_X$  與分子自由能  $G$  之間的定量橋樑提供了初步證據。

### 7.1 未來研究方向

未來的研究可以朝以下方向發展：

1. **代數逼近**：開發方法將蛋白質構象空間的點雲近似為代數簇，從而能夠應用 MMP 的嚴格理論。
2. **定量橋樑**：建立典範因子  $K_X$  與自由能  $G$  之間的定量對應關係，這將需要結合統計力學和代數幾何的理論。
3. **更複雜的蛋白質**：將方法擴展到更複雜的蛋白質系統，如具有打結結構的蛋白質。
4. **機器學習整合**：結合機器學習方法來預測構象轉變路徑和拓撲特徵。
5. **實驗驗證**：將計算結果與實驗數據（如 NMR、X-ray 晶體學）進行比較驗證。

這個跨學科的研究方向需要代數幾何學家、拓撲學家和計算生物學家之間的深度合作。如果成功，它將為理解蛋白質折疊提供一個全新的、嚴謹的數學框架。

## 致謝

本文的理論框架和計算實現基於對代數幾何、拓撲數據分析和計算生物學的深入思考。我們感謝所有為這些領域做出貢獻的研究者。

## 參考文獻

- [1] Atiyah, M. F. (1966). On analytic surfaces with double points. *Proceedings of the Royal Society of London. Series A, Mathematical and Physical Sciences*, 247(1249), 237-244.
- [2] Edelsbrunner, H., & Harer, J. (2010). *Computational topology: an introduction*. American Mathematical Society.
- [3] Ghrist, R. (2008). Barcodes: the persistent topology of data. *Bulletin of the American Mathematical Society*, 45(1), 61-75.
- [4] Tralie, C., Saul, N., & Bar-On, R. (2018). Ripser.py: A lean persistent homology library for Python. *The Journal of Open Source Software*, 3(29), 925.
- [5] Kollar, J., & Mori, S. (1998). *Birational geometry of algebraic varieties* (Vol. 134). Cambridge University Press.
- [6] Dill, K. A., & MacCallum, J. L. (2012). The protein-folding problem, 50 years on. *Science*, 338(6110), 1042-1046.
- [7] Xia, K., & Wei, G. W. (2014). Persistent homology analysis of protein structure, flexibility, and folding. *International Journal for Numerical Methods in Biomedical Engineering*, 30(8), 814-844.
- [8] Wales, D. J. (2003). *Energy landscapes: applications to clusters, biomolecules and glasses*. Cambridge University Press.
- [9] Muñoz, V., Thompson, P. A., Hofrichter, J., & Eaton, W. A. (1997). Folding dynamics and mechanism of  $\beta$ -hairpin formation. *Nature*, 390(6656), 196-199.
- [10] Zomorodian, A., & Carlsson, G. (2005). Computing persistent homology. *Discrete & Computational Geometry*, 33(2), 249-274.