

UNIVERSITY OF HAMBURG

MASTER THESIS

Predicting Protein Crystallization Conditions using Machine Learning

Author:

Michael HÜPPE

Supervisors:

Prof. Dr. Fabian Michael Kern

Prof. Dr. Aymelt Itzen

*A thesis submitted in fulfillment of the requirements
for the degree of Master of Science*

in the

Machine Learning in Bio Informatics
Department of Computer Science

Tuesday 18th November, 2025

Declaration of Authorship

I, Michael HÜPPE, declare that this thesis titled, “Predicting Protein Crystallization Conditions using Machine Learning” and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

“This one is for the boys with the booming system”

Nicki Minaj

UNIVERSITY OF HAMBURG

Abstract

Faculty Name
Department of Computer Science

Master of Science

Predicting Protein Crystallization Conditions using Machine Learning

by Michael HÜPPE

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centered vertically so can expand into the blank space above the title too...

Acknowledgements

The acknowledgments and the people to thank go here, don't forget to include your project advisor...

Contents

Declaration of Authorship	iii
Abstract	vii
Acknowledgements	ix
1 Decoding the Crystal Recipe: Predicting Protein Crystallization Conditions via Machine Learning	1
1.1 Introduction	1
1.2 Related Works	2
1.2.1 Challenge of X-ray Crystallography	2
1.2.2 Predicting Protein Crystallization	2
1.2.3 Predicting Crystallization Conditions from Sequence (and Structure)	3
2 Theory	5
2.1 Proteins	5
2.1.1 Protein Crystallization	5
2.2 Trees	5
3 Data	7
3.1 The Protein database	7
3.1.1 Acquisition	7
3.1.2 Format	7
3.1.3 Structure	8
3.1.4 Description	9
3.2 Data Normalization	11
3.2.1 Perturbation of missing values	11
3.2.2 Details Parsing	11
3.2.2.1 Pipeline	11
3.2.2.2 Parsing Quality	11
3.3 Data Analysis	11
3.3.1 Univariate Analysis	11
3.3.1.1 Input	11
3.3.1.2 Label	13
3.3.1.3 Label	13
3.3.2 Multivariate Analysis	13
3.3.3 Outlier and Anomaly Detection	13
3.3.4 Missing Data Analysis	13
4 Results	15
5 Discussion	17

A Data Appendix	19
Bibliography	21

List of Figures

3.1	Restructured database format.	9
3.2	Sequence Length and Atom Counts	11
3.3	Sequence Length and Atom Counts	12
3.4	Derived Sequence Features	12
3.5	Derived Surface Features	13

List of Tables

3.1 Percentage of Missing data per Attribute	10
--	----

List of Abbreviations

Physical Constants

Speed of Light $c_0 = 2.997\,924\,58 \times 10^8 \text{ m s}^{-1}$ (exact)

List of Symbols

a	distance	m
P	power	W (J s ⁻¹)
ω	angular frequency	rad

For/Dedicated to/To my...

Chapter 1

Decoding the Crystal Recipe: Predicting Protein Crystallization Conditions via Machine Learning

1.1 Introduction

Protein crystallization is the process of arranging purified protein molecules into a highly ordered, repeating lattice that forms a crystal. This crystalline state is essential for X-ray crystallography, the most widely used method for determining high-resolution protein structures. When X-rays are diffracted by a protein crystal, the resulting patterns allow reconstruction of the electron density and ultimately the atomic structure of the protein.

Crystallization is crucial because structural information provides fundamental insights into protein function and interactions which build the basis of structure-based drug design. Here binding sites derived from the structure enables rational development of small molecules that inhibit or modulate the protein's biological function.

Although AlphaFold predictions often align remarkably well with experimentally determined structures, they are not a substitute for them. Terwilliger *et al.* (2024) argue that AlphaFold accelerates, but cannot replace, experimental structure determination due to its varying local accuracies and occasional failures at the global structural level. This reinforces that, despite the breakthroughs brought by AlphaFold, further methodological innovation in structure determination and crystallography remains essential. Beyond crystallography, protein crystals are also used in neutron diffraction, cryo-electron microscopy benchmarking and biophysical studies of stability and folding.

Because most proteins do not readily form diffraction-quality crystals, the main bottleneck in structural biology is identifying the crystallization conditions under which a given protein will form suitable crystals (Mall *et al.*, 2025). Despite advances in structure prediction, identifying suitable crystallization conditions, such as buffer type, pH, salts, and precipitants remains largely empirical and often requires screening hundreds of combinations. This process is both time-consuming and costly, with a high failure rate.

At the same time, extensive data from the Protein Data Bank (PDB) and accurate structural predictions from models like AlphaFold2 have become widely available (Jumper *et al.*, 2021). These resources offer a unique opportunity to explore whether machine learning models can predict suitable crystallization conditions directly from protein sequence and/or structure.

1.2 Related Works

1.2.1 Challenge of X-ray Crystallography

Crystallizing a protein is often the bottleneck in X-ray crystallography. Only a small fraction (roughly 2–10%) of proteins produce diffraction-quality crystals, meaning over 90% of crystallization trials fail (Mall *et al.*, 2025). This trial-and-error process is costly, where >70% of the total expense in structure determination is spent on attempts not producing crystals of diffraction quality (Mall *et al.*, 2025). Achieving crystals requires finding the right crystallization conditions (e.g. precipitant chemicals, salts, pH, temperature), but currently these must be determined empirically by screening hundreds or thousands of conditions, which demands a lot of protein (McPherson & Gavira, 2013). This reality motivates computational approaches to predict either whether a protein is likely to crystallize (crystallization propensity) or even which specific conditions might lead to crystals. The former is a widely researched topic with models such as CrystalP2 (Kurgan *et al.*, 2009) or PPCPred (Mizianty & Kurgan, 2011), the latter however has received little to no attention (Jin *et al.*, 2022). While crystallization robots have eased the burden of manual search, discovering optimal conditions still means testing thousands of solutions, wasting protein in the process (Wilson & DeLucas, 2014). Any such predictive model could significantly reduce experimental screening, saving time and cost.

1.2.2 Predicting Protein Crystallization

As mentioned predicting crystallization propensity is a well established research area. Early **classical ML** relied on hand-crafted features such as amino-acid composition, the proteins isoelectric point, hydrophobicity, disorder, predicted structure and classifiers (SVM, RF, LR, GB). Matinyan *et al.* (2024) summarizes multiple tools such as:

- *XtalPred*/*XtalPred*-RF: feature-distribution scoring.
- *TargetCrys*: two-layer SVM ensemble.
- *Crysalis* (2016): integrated predictions + mutation suggestions.
- *BCrystal* (2020): XGBoost + SHAP-based feature selection.
- *DCFCrystal* (2021): cascaded RF stages (expression → purification → crystallization), with a membrane-protein branch.

These achieved moderate accuracy (60–75%, MCC 0.4–0.6) but demanded expert feature engineering. Deep learning automatically extracts sequence patterns as shown in superior performances presented by:

- *DeepCrystal* (2019): multi-scale CNN on one-hot sequences; ~ 83% accuracy, MCC 0.66.
- *CLPred* (2020): CNN+BLSTM; accuracy 85%, MCC 0.70.
- *ATTCrys* (2021): adds multi-head attention; MCC 0.72.
- Structure-infused: *SADeepCry* (2022) uses autoencoder+self-attention on sequence+predicted structure; *GMapCrys* (2023) employs GNNs on AlphaFold2 contact maps.

1.2.3 Predicting Crystallization Conditions from Sequence (and Structure)

While predicting “will it crystallize?” is useful, a more ambitious goal is to predict the actual crystallization conditions that would make a given protein form crystals. This is a multi-output prediction (the combination of reagents, concentrations, pH, etc. that will work) and more challenging. Thousands of successful crystallization recipes are known (recorded in databases like the Protein Data Bank), but each protein is unique and may crystallize in different conditions, often unpredictably with no known patterns to predict crystallization conditions (Zhang *et al.*, 2022). However, this does not stem from a lack of attention, Kirkwood *et al.* (2015) for example studied correlations between a protein’s isoelectric point and the pH of its crystallization buffer. However, only weak or inconsistent trends were found. Ultimately, Kirkwood *et al.* (2015) conclude that these trends are not sufficiently robust to guide initial crystallization-pH selection. Interestingly, even proteins with high sequence similarity did **not** necessarily crystallize under similar conditions, highlighting that small sequence/structure differences can lead to different optimal crystallization cocktails indicating a multidimensional problem.

Liao and Sun (2025) emphasize the importance of crystal packing and inter-molecular packing interfaces in determining crystallization conditions. In their work, they present **Molecular Assembly Simulation in Crystal Lattice (MASCL)**, a framework for simulating crystal packing using AlphaFold combined with symmetrical docking. Crystallization conditions are predicted using a patch-based method that quantifies molecular interface similarity between proteins. For a given target protein, proteins with the most similar physicochemical interface descriptors are used as reference points from which test crystallization conditions are chosen. Meaning no de-novo crystallization conditions are constructed.

This pipeline of constructing a “crystal fingerprint” for a given protein and comparing it to previously crystallized proteins was evaluated on lysozyme, a common model protein in crystallization studies. In this test case, the proposed AAI-PatchBag approach successfully identified conditions yielding crystals with the desired packing characteristics. However, the use of lysozyme as a model system has been repeatedly criticized (Chayen & Saridakis, 2001). Particularly in the context of assessing prediction accuracy lysozyme benchmarks should be assessed with caution, because it is well known and used for its unusually high crystallizability (Ghosh, 2023). It crystallizes across a broad range of pH values without loss of crystal quality (Iwai *et al.*, 2008), and also tolerates wide variations in temperature and salt concentration (Ataka & Asai, 1988).

Nevertheless, for lysozyme specifically, Liao and Sun (2025) show that similarity in crystal packing information has a stronger influence on predicting successful crystallization conditions than sequence or structural homology alone.

In contrast Lee *et al.* (2019) introduced a proof-of-concept deep learning model to map protein sequence to de-novo crystallization conditions. They parsed crystallization records from PDB entries and framed the task as a multi-label classification: for a given protein sequence, predict which “crystallization terms” (e.g. specific buffers, salts, precipitants like PEG, etc.) appeared in successful recipes for that protein. Essentially, the model learns associations between sequence features and the types of reagents or techniques that tend to be used. Remarkably, this sequence-to-condition model did show predictive power. A simple 1-layer CNN could achieve a weighted F1-score around 0.46 on held-out proteins, substantially better than random guessing in this high-dimensional space. The CNN outperformed a fully-connected network, suggesting that local sequence motifs (captured by convolutional filters) were

informative for certain crystallization agents. For example, the authors noted that hydrophilic or charged residues in the sequence strongly influenced buffer and salt predictions (likely because they affect the protein's isoelectric point and solubility). This indicates real biochemical signal: proteins rich in acidic/basic residues might require certain pH buffers or salt conditions, etc., whereas hydrophobic patches might correlate with needing precipitating agents like PEGs or additives.

It's important to emphasize that this research is still early-stage. An F1 of 0.45 means the model is far from perfectly pinpointing the exact crystallization recipe, but it is better than trial-and-error alone and demonstrates that sequence patterns can inform what conditions are likely to work. As a future direction, the authors suggested incorporating more variables (e.g. predicting optimal pH and temperature as continuous values, not just class labels) and using the approach to focus screening on a smaller set of candidate conditions. Nonetheless, this is a promising frontier: even a modest predictor that suggests, say, the top 10 most likely crystallization cocktails for a new protein (instead of blindly testing 1000) would be hugely valuable.

Chapter 2

Theory

2.1 Proteins

2.1.1 Protein Crystallization

2.2 Trees

Chapter 3

Data

3.1 The Protein database

The **Protein Data Bank (PDB)** is the central repository for experimentally determined three-dimensional structures of proteins, nucleic acids, and complexes. As of 2025, it contains over 220,000 entries, the majority of which are proteins solved via X-ray crystallography. Each entry includes not only the atomic coordinates of the protein structure but also extensive metadata about the experimental conditions used during crystallization. The following outlines the data **acquisition**, **format**, and **structure** to enable data preprocessing and analysis.

3.1.1 Acquisition

Using the RCSB Search API, all entries solved by X-ray diffraction (about 80%) were queried. Entries that did not contain any information about the crystallization at all or the entry was otherwise incomplete were not downloaded. A Python script then converted these identifiers into download links for the mmCIF files (.cif.gz), wrote them to text files, and split the list into several parts to enable parallel downloads via wget. The complete set of structures was downloaded in compressed mmCIF format and stored locally (after decompression to .cif where needed). Due to the large number of entries, the download process took approximately two days. However, the total download time naturally depends on the available internet bandwidth.

3.1.2 Format

The **PDB** provides structural data in two formats: PDB and CIF. The newer CIF format was chosen because, unlike the fixed-width PDB format, it does not impose size limitations on large structures. In addition, it can represent complex features such as branched carbohydrates and offers greater detail and flexibility than its predecessor. The data items are in the format of `'_' + category name + '.' + attribute name`. Data categories can be either saved in either key-value or in tabular format. These can be easily parsed into a column vector as is in CSV files. For example the crystallization conditions are saved in key-value format and look like the following for 3P4V:

<code>_exptl_crystal_grow.crystal_id</code>	1
<code>_exptl_crystal_grow.method</code>	'VAPOR DIFFUSION, SITTING DROP'
<code>_exptl_crystal_grow.apparatus</code>	None
<code>_exptl_crystal_grow.atmosphere</code>	None
<code>_exptl_crystal_grow.pH</code>	9.5
<code>_exptl_crystal_grow.temp</code>	298.0
<code>_exptl_crystal_grow.pdbx_details</code>	'3.2M (NH4)2SO4, 0.1M Glycine, ...'
<code>_exptl_crystal_grow.time</code>	None

A category is stored in tabular format when a token defines multiple values. In this case, `loop_` is followed by rows of data-item names, with data values separated by whitespace. Notably, the protein’s structural information i.e., the `atom_sites` category is represented in this format, as shown in the following example:

```

loop_
_atom_site.group_PDB
_atom_site.id
_atom_site.type_symbol
_atom_site.label_atom_id
_atom_site.label_alt_id
_atom_site.label_comp_id
_atom_site.label_asym_id
_atom_site.label_entity_id
_atom_site.label_seq_id
_atom_site.pdbx_PDB_ins_code
_atom_site.Cartn_x
_atom_site.Cartn_y
_atom_site.Cartn_z
ATOM 1 N N . VAL A 1 1 ? 6.204 16.869 4.854 1.00 49.05 ...
ATOM 2 C CA . VAL A 1 1 ? 6.913 17.759 4.607 1.00 43.14 ...
ATOM 3 C C . VAL A 1 1 ? 8.504 17.378 4.797 1.00 24.80 ...
ATOM 4 O O . VAL A 1 1 ? 8.805 17.011 5.943 1.00 37.68 ...

```

3.1.3 Structure

In structural biology, one rarely needs to work with all proteins deposited in the **PDB**. Consequently, the **PDB** is designed as an entry-oriented resource: users typically retrieve and parse the complete dataset for a single protein of interest. This design is well aligned with the standard workflow in structural biology, where researchers focus on a small number of proteins but require all available structural, experimental, and metadata associated with those entries.

In this project, however, the goal is fundamentally different: to identify patterns across many proteins, using all deposited entries in the database. Parsing a single CIF file takes ~100 ms using the fastest available library, **Gemmi** which amounts to approximately 9 hours for the entire PDB. Repeating such operations would be not feasible for data analysis where data typically has to be read multiple times across sessions.

To address this, the first step after downloading the database was to restructure it into a feature-based data model, in which information is grouped by category rather than by entry. Each category is stored as a separate Parquet file (also convertible to CSV), where each attribute becomes a column spanning all proteins. This enables highly targeted access: if an analysis requires only a single attribute, the corresponding column can be read directly without loading irrelevant data.

Although the initial restructuring (reading, grouping, and writing) takes ~11 hours, it drastically improves downstream performance. For example, extracting all crystallization-condition information from the entire database now takes ~10s, compared to the 9 hours required to parse every CIF individually.

Storage efficiency is also greatly improved. The full PDB occupies ~46 GB and allocates ~110 GB, making it impossible to load into memory. In contrast, the derived

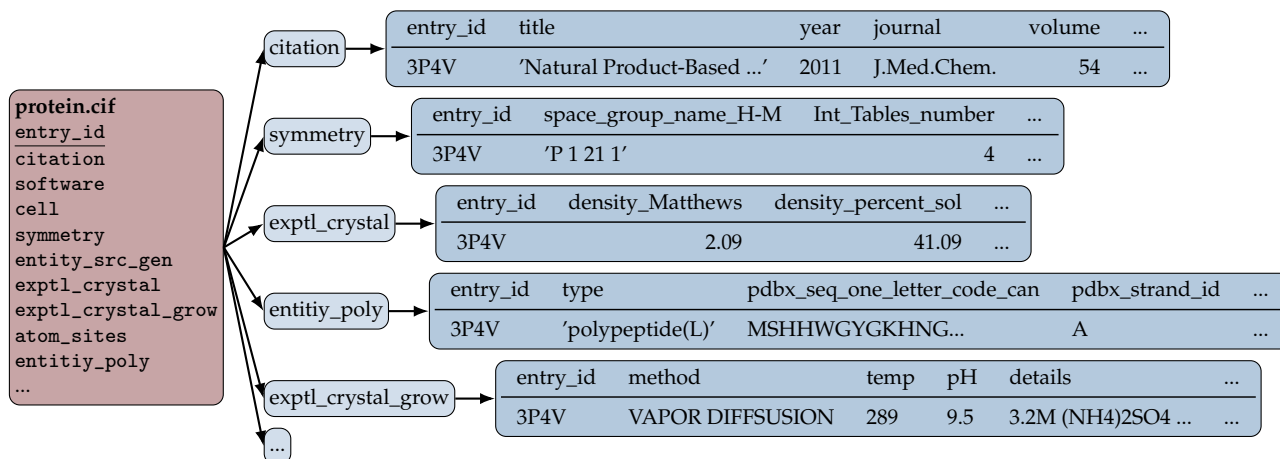


FIGURE 3.1: Restructured database format.

dataset containing only crystallization conditions is about 20 MB, small enough to be memory-resident and repeatedly queried without overhead.

This dramatic reduction in both time and storage requirements arises because most information in CIF files is irrelevant for crystallization-condition prediction. Metadata such as software provenance or publication information is not needed for this project and can be safely excluded from the feature-based representation.

However, to enable analyses that combine multiple categories, the restructured database must be joinable in a way that allows the original information to be reconstructed. This was accomplished using a star schema, in which each category table is indexed by the unique `entry_id` of the protein. This design ensures that information from different categories (such as crystallization conditions and sequence data) can be linked by performing joins on the shared entry identifier. The restructured database format is depicted in [Figure 3.1](#).

3.1.4 Description

After filtering for majorly incomplete entries the total number of deposited proteins was 196743 as of 01.09.2025. Each of these entries has around 700 attributes across 68 categories. However, as mentioned before the minority of these attributes/categories is of importance when predicting crystallization conditions. More specifically, the important categories and their attributes are the following:

Input Categories

The `entity_poly` category of the mmCIF format describes each polymeric entity in a macromolecular structure. A polymer represents a sequence of linked monomers (e.g., amino acids or nucleotides), and this category provides essential metadata about its type, sequence, and representation in the structural model. `entity_id` Unique identifier linking the polymer to the corresponding entry in the `entity` category.

`type` Specifies the polymer type, such as polypeptide(L), DNA, RNA, or polysaccharide.

`nstd_linkage` Indicates whether the polymer contains non-standard chemical linkages (e.g., cross-links or modified connectivity).

`nstd_monomer` Flags the presence of non-standard or modified monomers within the sequence.

`pdbx_seq_one_letter_code` The polymer sequence in one-letter code, including symbols for modified residues.

`pdbx_seq_one_letter_code_can` Canonicalized one-letter sequence where modified residues are mapped to their closest standard equivalents.

`pdbx_strand_id` Lists the chain identifiers (e.g., A, B) representing this polymer in the structure.

`pdbx_target_identifier` Optional external identifier used mainly in structural genomics pipelines.

The `atom_site` category contains the atomic coordinates and related information that define the three-dimensional structure of the macromolecule. Each row corresponds to a single atom and records properties such as atom name, element, residue identifier, chain, Cartesian coordinates, occupancy, and atomic displacement parameters.

Label Categories

The category `exptl_crystal_grow` defines the crystallization conditions. It contains 17 attributes that describe the experimental setup and methodology used to grow the protein crystal. However, as seen in [Table 3.1](#) for 13 of the attributes the majority of entries do not contain any information.

Attribute	Missing percentage
<code>method</code>	10.78
<code>temp</code>	10.79
<code>pH</code>	22.30
<code>pdbx_pH_range</code>	75.24
<code>pdbx_details</code>	0.04
<code>temp_details</code>	98.56
<code>apparatus</code>	99.84
<code>atmosphere</code>	99.84
<code>details</code>	99.83
<code>method_ref</code>	99.84
<code>pressure</code>	99.83
<code>pressure_esd</code>	99.84
<code>seeding</code>	99.83
<code>seeding_ref</code>	99.84
<code>temp_esd</code>	99.84
<code>time</code>	99.84

TABLE 3.1: Percentage of Missing data per Attribute

Thus, the attributes of relevance are the following:

`method`: Describes the crystallization technique employed, such as vapor diffusion, batch crystallization, or microbatch methods (in free text).

`pH`: Specifies the pH of the crystallization solution, which strongly influences protein stability and crystal formation (numerical).

`temp`: Records the temperature at which the crystallization experiment was performed, typically given in Kelvin or Celsius (in free text).

`pdbx_details`: Provides free-text experimental details, such as buffer components,

precipitants, additives, or other conditions important for reproducing the crystallization setup (in free text).

Moreover, it might be of importance how well the protein has crystallized. Crystal quality measures are primarily provided in the `exptl_crystal` category, which describes the physical properties of the crystal, and in the `diffn` and `reflns` categories, which contain diffraction statistics such as resolution, completeness, and R-factors that reflect the overall quality of the crystal.

3.2 Data Normalization

3.2.1 Perturbation of missing values

Typically,

3.2.2 Details Parsing

How the details were parsed to retrieve a uniform label representation

3.2.2.1 Pipeline

Detailed explanation of the pipeline

3.2.2.2 Parsing Quality

How the parsing quality was determined.

3.3 Data Analysis

3.3.1 Univariate Analysis

Histograms, distributions, summary statistics

3.3.1.1 Input

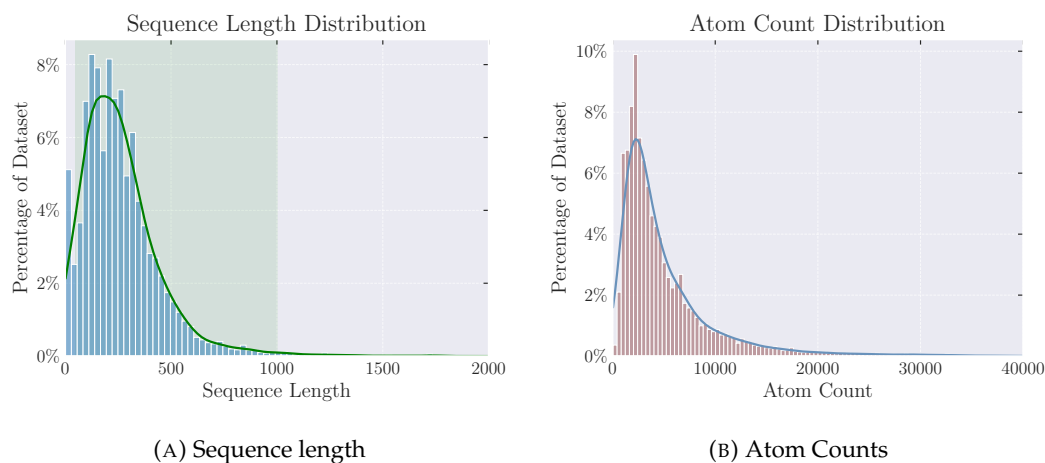


FIGURE 3.2: Sequence Length and Atom Counts

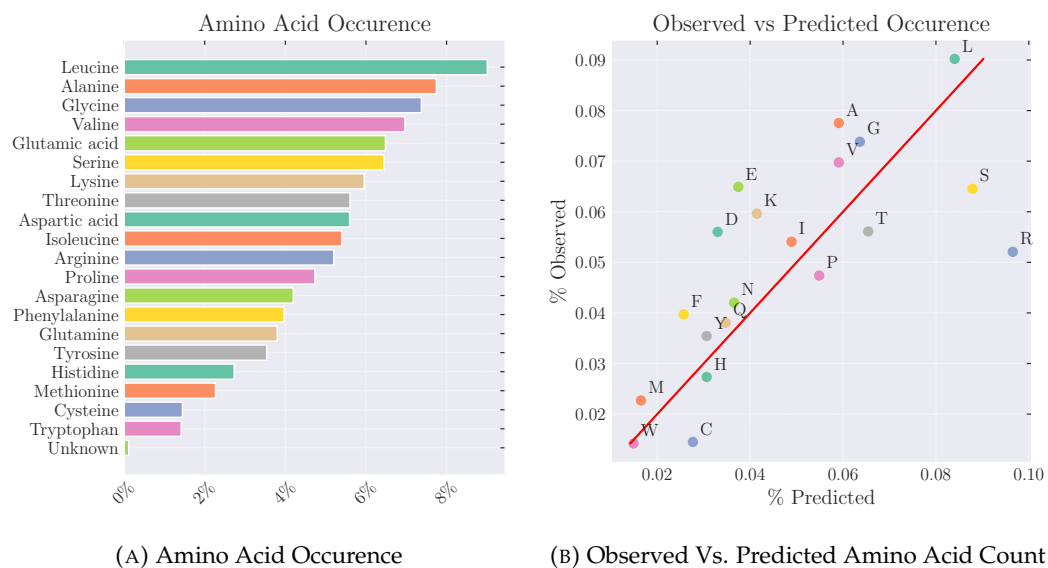


FIGURE 3.3: Sequence Length and Atom Counts

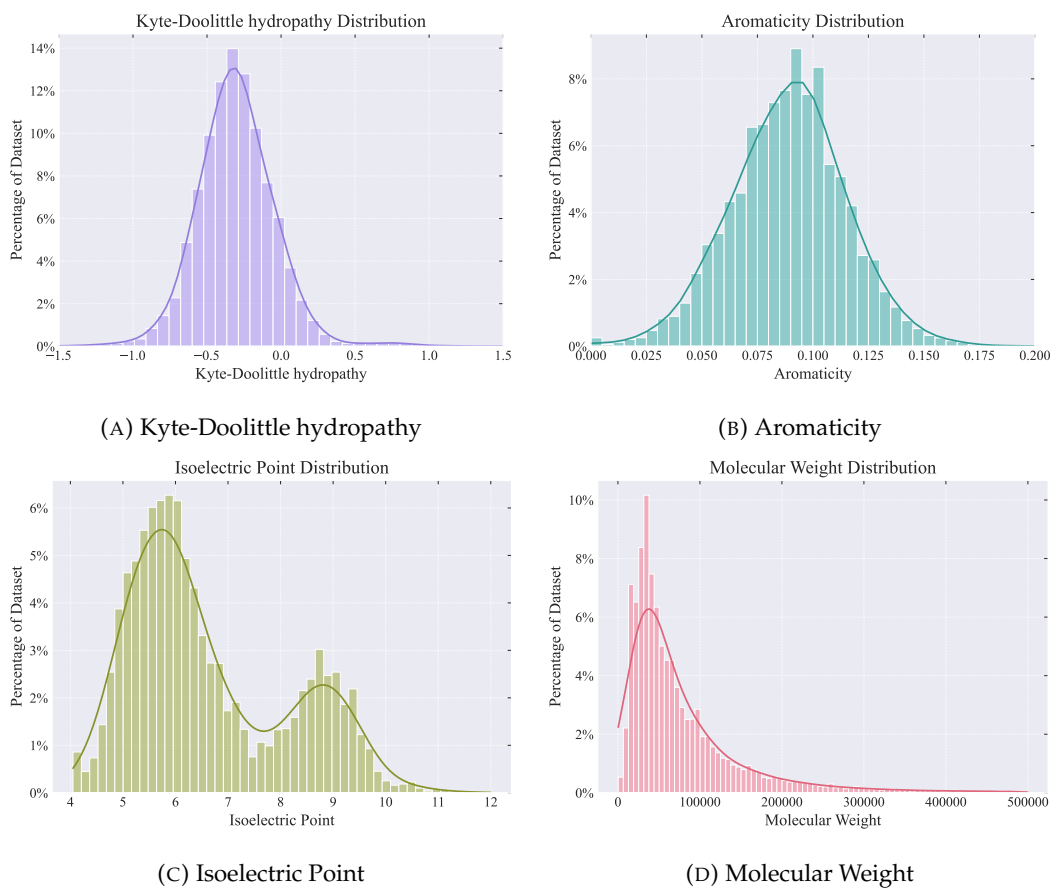


FIGURE 3.4: Derived Sequence Features

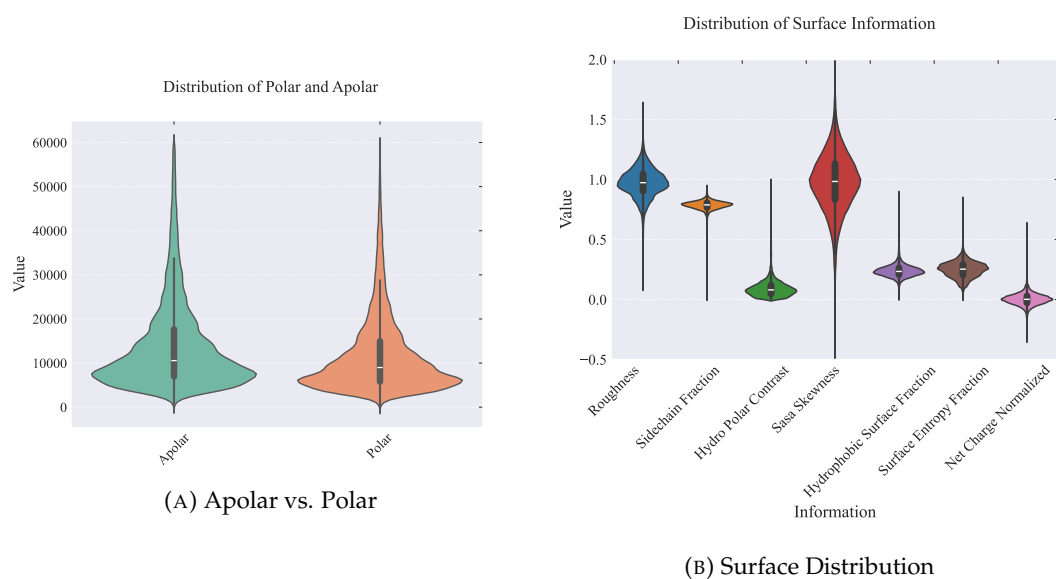
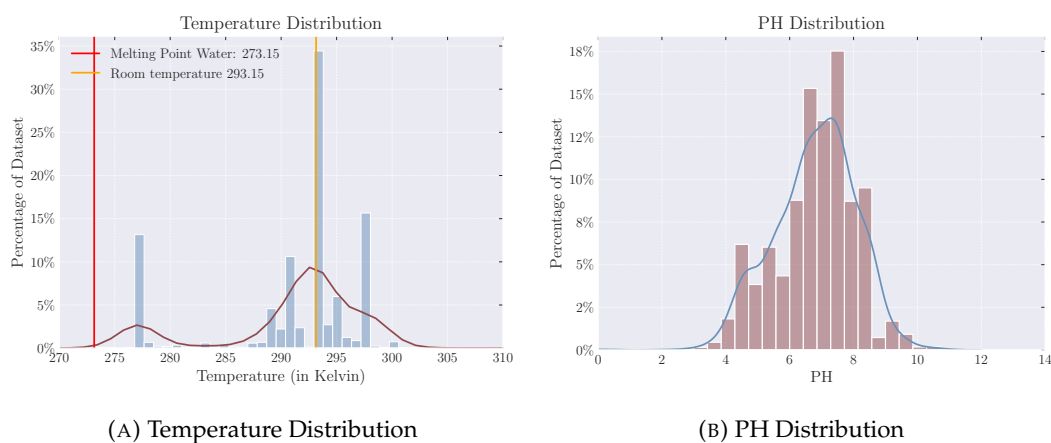


FIGURE 3.5: Derived Surface Features



3.3.1.2 Label

3.3.1.3 Label

3.3.2 Multivariate Analysis

Feature-feature relationships, correlation matrices, scatter plots

3.3.3 Outlier and Anomaly Detection

Identifying unusual values or patterns

3.3.4 Missing Data Analysis

Assessing missingness, patterns, potential mechanism

Chapter 4

Results

Chapter 5

Discussion

Appendix A

Data Appendix

Bibliography

- Ataka, M., & Asai, M. (1988). Systematic studies on the crystallization of lysozyme: Determination and use of phase diagrams. *Journal of Crystal Growth*, 90(1), 86–93. [https://doi.org/10.1016/0022-0248\(88\)90302-8](https://doi.org/10.1016/0022-0248(88)90302-8)
- Chayen, N. E., & Saridakis, E. (2001). Is lysozyme really the ideal model protein? *Journal of Crystal Growth*, 232(1), 262–264. [https://doi.org/10.1016/S0022-0248\(01\)01203-9](https://doi.org/10.1016/S0022-0248(01)01203-9)
- Ghosh, R. (2023). Membrane-Based Micro-Volume Dialysis Method for Rapid and High-Throughput Protein Crystallization. *Processes*, 11(7), 2148. <https://doi.org/10.3390/pr11072148>
- Iwai, W., Yagi, D., Ishikawa, T., Ohnishi, Y., Tanaka, I., & Niimura, N. (2008). Crystallization and evaluation of hen egg-white lysozyme crystals for protein pH titration in the crystalline state. *Journal of Synchrotron Radiation*, 15(Pt 3), 312–315. <https://doi.org/10.1107/S0909049507059559>
- Jin, C., Shi, Z., Kang, C., Lin, K., & Zhang, H. (2022). TLCrys: Transfer Learning Based Method for Protein Crystallization Prediction. *International Journal of Molecular Sciences*, 23(2), 972. <https://doi.org/10.3390/ijms23020972>
- Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Židek, A., Potapenko, A., Bridgland, A., Meyer, C., Kohl, S. A. A., Ballard, A. J., Cowie, A., Romera-Paredes, B., Nikolov, S., Jain, R., Adler, J., ... Hassabis, D. (2021). Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873), 583–589. <https://doi.org/10.1038/s41586-021-03819-2>
- Kirkwood, J., Hargreaves, D., O’Keefe, S., & Wilson, J. (2015). Using isoelectric point to determine the pH for initial protein crystallization trials. *Bioinformatics*, 31(9), 1444–1451. <https://doi.org/10.1093/bioinformatics/btv011>
- Kurgan, L., Razib, A. A., Aghakhani, S., Dick, S., Mizianty, M., & Jahandideh, S. (2009). CRYSTALP2: Sequence-based protein crystallization propensity prediction. *BMC Structural Biology*, 9(1), 50. <https://doi.org/10.1186/1472-6807-9-50>
- Lee, H., Wu, Z. H., Corbi-Verge, C., Mok, M., Kang, S., Liao, S., Zhang, Z., & Garton, M. (2019). De novo crystallization condition prediction with deep learning [Accessed: 2025-07-31]. *Advances in Neural Information Processing Systems*, 32. https://mlcb.github.io/mlcb2019_proceedings/papers/paper_3.pdf
- Liao, K.-J., & Sun, Y.-J. (2025). Using AlphaFold and Symmetrical Docking to Predict Protein–Protein Interactions for Exploring Potential Crystallization Conditions. *Proteins*, 93(10), 1747–1766. <https://doi.org/10.1002/prot.26844>
- Mall, R., Kaushik, R., Martinez, Z. A., Thomson, M. W., & Castiglione, F. (2025). Benchmarking protein language models for protein crystallization. *Scientific Reports*, 15(1), 2381. <https://doi.org/10.1038/s41598-025-86519-5>
- Matinyan, S., Filipcik, P., & Abrahams, J. P. (2024). Deep learning applications in protein crystallography. *Acta Crystallographica. Section A, Foundations and Advances*, 80(Pt 1), 1–17. <https://doi.org/10.1107/S2053273323009300>

- McPherson, A., & Gavira, J. A. (2013). Introduction to protein crystallization. *Acta Crystallographica. Section F, Structural Biology Communications*, 70(Pt 1), 2–20. <https://doi.org/10.1107/S2053230X13033141>
- Mizianty, M. J., & Kurgan, L. (2011). Sequence-based prediction of protein crystallization, purification and production propensity. *Bioinformatics*, 27(13), i24–i33.
- Terwilliger, T. C., Liebschner, D., Croll, T. I., Williams, C. J., McCoy, A. J., Poon, B. K., Afonine, P. V., Oeffner, R. D., Richardson, J. S., Read, R. J., & Adams, P. D. (2024). AlphaFold predictions are valuable hypotheses and accelerate but do not replace experimental structure determination. *Nature Methods*, 21(1), 110–116. <https://doi.org/10.1038/s41592-023-02087-4>
- Wilson, W. W., & DeLucas, L. J. (2014). Applications of the second virial coefficient: Protein crystallization and solubility. *Acta Crystallographica Section F: Structural Biology Communications*, 70(5), 543–554. <https://doi.org/10.1107/S2053230X1400867X>
- Zhang, X., Xu, Z., Zhou, J., Xing, X., & Li, L. (2022). Enhancement of Protein Crystallization Using Nano-Sized Metal–Organic Framework. *Crystals*, 12(5), 578. <https://doi.org/10.3390/cryst12050578>