

PDB Autofill



Henry Lee, You-Hsin Chen, and Jenny Bennett

UNIVERSITY *of* WASHINGTON



Background

- > **Our goal:** Classify the reason for missing electron densities of protein crystals in PDB by using a random forest model
- > **Stretch:** Predict the missing density coordinates in proteins with a neural network model

						Sequence Number
ATOM	437	CG2	VAL	A	47	47
ATOM	438	H	VAL	A	47	
ATOM	439	N	GLY	A	52	
ATOM	440	CA	GLY	A	52	

```
REMARK 465 MISSING RESIDUES
REMARK 465 THE FOLLOWING RESIDUES WERE NOT LOCATED IN THE
REMARK 465 EXPERIMENT. (M=MODEL NUMBER; RES=RESIDUE NAME; C=CHAIN
REMARK 465 IDENTIFIER; SSSEQ=SEQUENCE NUMBER; I=INSERTION CODE.)
REMARK 465
REMARK 465 M RES C SSSEQI
REMARK 465 GLY A 48
REMARK 465 GLY A 49
REMARK 465 ILE A 50
REMARK 465 GLY A 51
REMARK 470
REMARK 470 MISSING ATOM
REMARK 470 THE FOLLOWING RESIDUES HAVE MISSING ATOMS (M=MODEL NUMBER;
REMARK 470 RES=RESIDUE NAME; C=CHAIN IDENTIFIER; SSEQ=SEQUENCE NUMBER;
REMARK 470 I=INSERTION CODE):
REMARK 470 M RES CSSEQI ATOMS
REMARK 470 PHE A 53 CG CD1 CD2 CE1 CE2 CZ
REMARK 470 LYS A 60 CG CD CE NZ
```

Use Cases

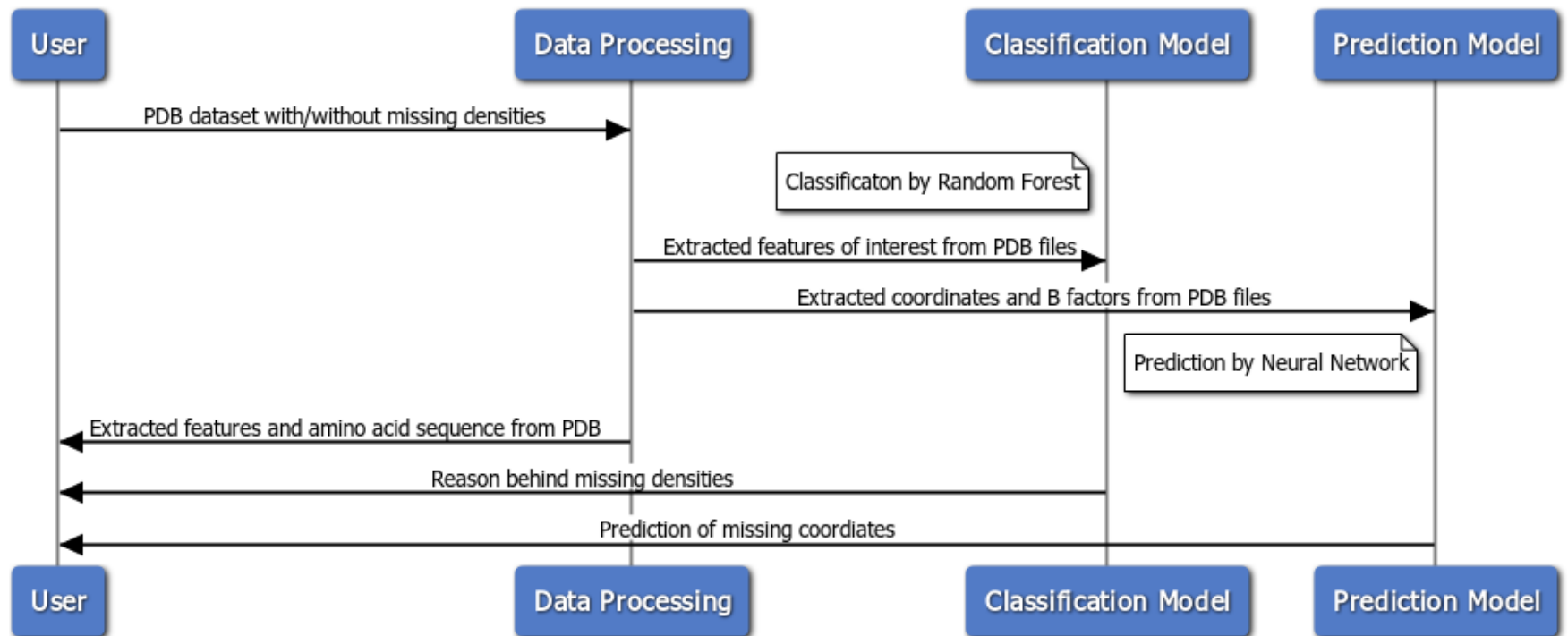
Alice is a computational scientist who will be using the tool to determine the expected structure of her proteins of interest. Alice will use this package to:

- > **Classify** why her proteins of interest have missing densities in the PDB.
- > **Predict the missing densities** in her proteins of interest.
- > **Extract sequences and features** in a DataFrame to perform further analysis on her own.

Joe is new to data science and is interested in learning more about PDB files. He will use the package to:

- > **Download a sample dataset** of PDB files.
- > **Extract interesting features** from the PDB files.
- > **Classify and predict** the missing densities that are not captured in the files.

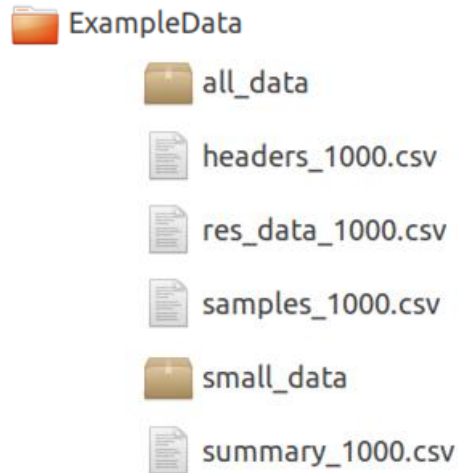
PDB Autofill



www.websequencediagrams.com

Downloading Data

- Application:
 - 1) Users can download data we used in example and do the following PDB analysis.
 - 2) The package will create a folder called ExampleData on the working directory and download two compressed .zip files which include 3000 .pdb and 20 .pdb files respectively and four .csv files which are our input files.



Data Processing

- Application:

- 1) Count how many .pdb files in the folder has missing residues.
- 2) Create new dataset which contains 50% proteins w/ missing residues and 50% proteins w/o missing residues for training data in order to avoid imbalanced class of follow-up classified model construction

1a0j	2020/2/28 下午 10:30	PDB 檔案	636 KB
1a3b	2020/2/28 下午 10:30	PDB 檔案	252 KB
1a7v	2020/2/28 下午 10:30	PDB 檔案	205 KB
1a09	2020/2/28 下午 10:30	PDB 檔案	262 KB
1a46	2020/2/28 下午 10:30	PDB 檔案	252 KB
1a52	2020/2/28 下午 10:30	PDB 檔案	363 KB
1a55	2020/2/28 下午 10:30	PDB 檔案	254 KB
1a95	2020/2/28 下午 10:30	PDB 檔案	429 KB
1ad4	2020/2/28 下午 10:30	PDB 檔案	381 KB

Name of
proteins



Biopython
glob

statistics

Total number of samples: 2970
Samples have missing residues: 2422
Samples without missing residues:: 548
There are 81.55% samples have missing residue



Random
pandas

Protein	
0	2y39
1	2o73
2	3d5m
3	1gey
4	4y79
...	...



Load pdb files we chose
and go extracting
information from them



UNIVERSITY of
WASHINGTON

Data Processing

- ### Dataset with some headers(resolution/statistics of b factor/has missing residue)

residues -> getting properties
(sequence length/ hydrophobic...)

	2y39	2o73	3d5m	1gey	4y79	3gem	2z91	2jfk	3ueo	3qun
0	GLY	ASP	SER	THR	ILE	SER	GLN	GLN	GLY	SER
1	ASP	ILE	MET	ILE	VAL	ALA	LEU	SER	LEU	GLY
2	LEU	ASN	SER	THR	GLY	PRO	LEU	MET	PHE	LEU
3	HIS	VAL	TYR	ASP	GLY	ILE	GLU	ARG	SER	VAL
4	GLU	VAL	THR	LEU	GLN	LEU	SER	LEU	GLN	PRO

UNIVERSITY of
WASHINGTON

Random Forest

• Application :

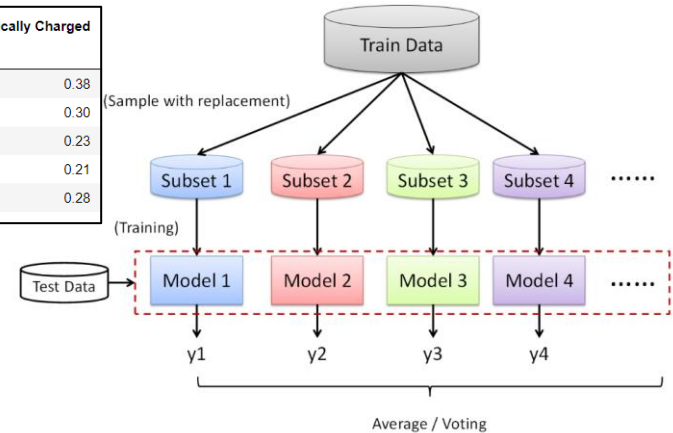
- 1) Users can build up own random forest model, and list importance of each feature in the model and mean accuracy under 5-fold cross-validation.

	has_missing_residues	Sequence Length	resolution	b_factor_gt50	b_factor_max	Electrically Charged
Protein						
2y39	True	110.0	1.41	0	46.46	0.38
2o73	True	992.0	1.80	4	65.24	0.30
3d5m	True	1116.0	2.20	3225	107.96	0.23
1gey	True	335.0	2.30	190	87.44	0.21
4y79	True	287.0	2.10	282	98.34	0.28



- 1) b_factor_gt50 0.181701
- 2) Sequence Length 0.153594
- 3) b_factor_max 0.146267
- 4) Electrically Charged 0.126673
- 5) resolution 0.118797
- 6) Hydrophobic 0.097705
- 7) Nonpolar Side Chains 0.088332
- 8) Special 0.086930

Accuracy of the prediction in 5-fold cross-validation = 71.27%



In our project, we used 1000 proteins with these 8 features to build up Random Forest model (500 decision trees). The average accuracy is 71.27 %. To enhance the performance, users could try adding more useful features, data or stacking with a boosting model.

Random Forest

- Application :
 - 2) Users who are going to do experiment can use RF model we have already built to test if it tends to get missing residues under this circumstance.
 - 3) Users can get most used features of nodes that each input data went through.

```
RF_model = joblib.load('RandomForest_model.pkl') #load model
```

input data

Prediction:
has missing residues
(yes/no)

In each decision tree

```
# The path of single input data goes through the random forest model
feature_count_accum = []
for j, tree in enumerate(RF_model.estimators_):
    # matrix of nodes that input data go through (boolean)
    dense_matrix = tree.decision_path(x[150].reshape(1, -1)).todense()
    # transform to array
    dense_sample = np.array(dense_matrix)[0]
    # extract number of nodes that input data goes through
    node_position = np.where(dense_sample == 1)[0]
    feature_count = []
```

```
for i in range(len(node_position)):
    number = node_position[i]
    feature_count.append(feature_name[tree.tree_.feature[number]])
feature_count accum.extend(feature_count)
```

Which feature that the node is ruled by

Matrix shows nodes that input data goes through (value = 1)

[illegible]

Transform to position (number) of nodes

W

UNIVERSITY of

Most used features of nodes that input data went through [('Nonpolar Side Chains', 1050), ('b_factor_max', 945), ('b_factor_gt50', 930), ('Sequence Length', 891), ('resolution', 833), ('Electrically Charged', 797), ('Hydrophobic', 645), ('Special', 597)]

Data Processing

- Application:

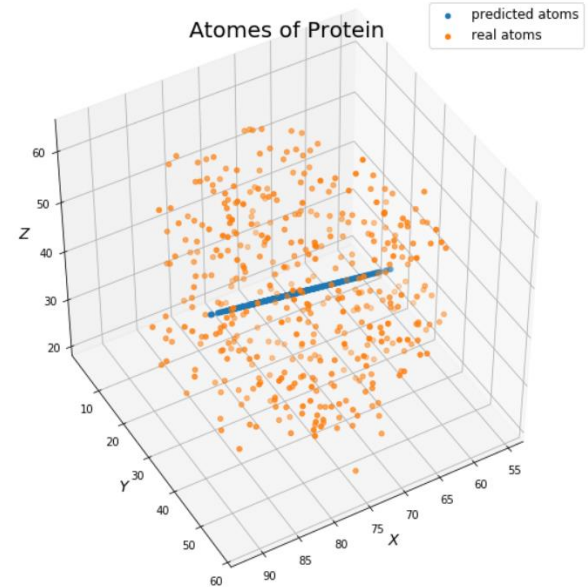
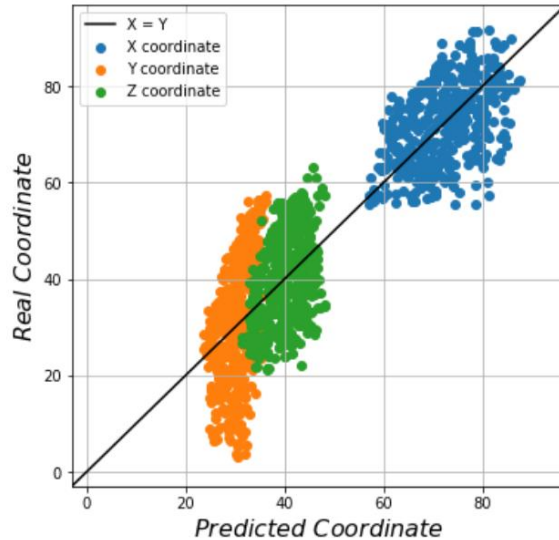
- 1) Creates a dataframe which contains the atom coordinates, the b values, the previous atom coordinates, the following atom coordinates.
- 2) Creates a dataframe which contains missing residues, the corresponding protein chain and sequences for each protein

	res name	chain	ssseq
2y39	[SER, HIS, ARG, ASN, GLU, ALA, GLY, HIS]	[A, A, A, A, A, A, A, A, A]	[31, 32, 33, 34, 35, 36, 37, 38]
2o73	[MET, LEU, SER, ASP, ILE, GLN, THR, LYS, LEU, ...]	[A, A, A, A, A, A, A, A, A, B, B, B, B, B, ...]	[1, 167, 168, 169, 170, 171, 172, 173, 174, 1, ...]
3qun	[GLY, HIS, GLY, ALA, ILE, ARG, ASP, HIS, ASP, ...]	[A, A, A, A, A, A, A, A, A, A, A, A, A, A, A, ...]	[57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 6, ...]
1gey	[MET, SER, THR, VAL, THR, PRO, TYR, GLN, SER, ...]	[A, A, A, A, A, A, A, A, A, A, A, A, A, A, A, ...]	[1, 2, 3, 4, 18, 19, 20, 21, 22, 25, 26, 27, 2, ...]
2jfk	[MET, HIS, HIS, HIS, HIS, HIS, SER, SER, ...]	[A, A, A, A, A, A, A, A, A, A, A, A, A, A, A, ...]	[399, 400, 401, 402, 403, 404, 405, 406, 407, ...]
3gem	[HIS, MSE, THR, LEU, SER, GLN, PRO, LYS, ASP, ...]	[A, A, A, A, A, A, A, A, A, A, A, A, A, A, A, ...]	[0, 1, 2, 3, 4, 183, 184, 185, 186, 187, 188, ...]
3d5m	[GLU, LYS, GLY, SER, LEU, SER, ARG, ALA, ARG, ...]	[A, A, A, A, A, A, A, A, A, A, A, A, A, A, A, ...]	[150, 151, 152, 563, 564, 565, 566, 567, 568, ...]
2z91	[SER, ALA, ALA, GLN, THR, ASN, ARG, ASP, CYS, ...]	[A, A, A, A, A, A, A, A, A, A, A, A, B, B, C, C, ...]	[128, 129, 130, 131, 132, 133, 213, 214, 215, ...]
4y79	[ARG, GLY, LEU, PRO, LYS, ALA, LYS, SER, HIS, ...]	[A, A, A, A, A, A, A, A, A, A, A, A, A, A, A, ...]	[245, 246, 247, 248, 249, 250, 251, 252, 253, ...]
3ueo	[GLY, PRO, LEU, GLY, SER, GLU, GLU, SER, LEU, ...]	[A, A, A, A, A, A, A, A, A, A, A, A, A, A, A, ...]	[544, 545, 546, 547, 548, 549, 550, 584, 585, ...]

Neural Networks

- Application:

- 1) Users can use their protein data as input and set the validation data size and seed to check if the multilayer perceptron model works well.





Thanks for listening