Machine Learning and Artificial Intelligence in MX

Melanie Vollmar Diamond Light Source Ltd UK



Outline

- Aim and idea for the work
- Why machine learning and artificial intelligence
- Training data and METRIX database for exploratory work
- Experimental phasing success
- Molecular replacement success
- Map traceability
- Future plans





Aim

Is there predictive power in crystallographic metrics (multiplicity, completeness, different R values, high and low resolution limit)

If yes, are they useful for anything

If yes then:

Create a set of tools using machine learning to help users/crystallographers to solve their protein structures.

During data collection, e.g. giving recommendations

During data analysis, e.g. data reduction and phasing



Why machine learning and artificial intelligence

close to 158,000 structures



Yea

Problems:

- Inconsistent crystallographic metrics across the PDB
- Most diffraction data not public
- Connection between diffraction data and deposited structure



Training data and METRIX database



Training data and METRIX database

METRIX content and schema



Pre-assessment and classifiers tried

Pre-assessment tried

•Linear Pearson's correlation coefficients

•Recursive feature elimination

Classifiers tried

•Support vector machine with linear kernel

•Support vector machine with RBF kernel

Decision tree

•Decision tree with Bagging

Decision tree with AdaBoost

Random forest

Extreme random forest

Combined results to identify most important decision making features;

Then retrain all classifiers and assess their performance; Python 3.x



703 samples; stratified test-train split (20/80) 3-fold cross-validation (20/80 split)



diamond

Important decision making features

d_{max} → low resolution cut-off d_{min} → high resolution cut-off https://github.com/ccp4/metrix_m



	tree with AdaBoost	classifier
ACC (%)	95	100
Class Error (%)	5	0
Sensitivity (%)	96	100
Specificity (%)	94	100
FPR (%)	6	0
Precision (%)	97	100
F1 score (%)	96	100
ROC AUC (%)	99	100
TP test set	90	94
TN test set	44	47
FP test set	3	0
FN test set	4	0

Decision

Porfoct







Probability cut-off for class 1: 80% https://github.com/ccp4/metrix_ml





Molecular replacement success

Pre-assessment and classifiers tried

Pre-assessment tried

•Linear Pearson's correlation coefficients

•Recursive feature elimination

Classifiers tried

•Support vector machine with linear kernel

•Support vector machine with RBF kernel

Decision tree

•Decision tree with Bagging

Decision tree with AdaBoost

Random forest

Extreme random forest

Combined results to identify most important decision making features:

Then retrain all classifiers and assess their performance; Python 3.x

1020 samples; stratified test-train split (20/80) 3-fold cross-validation (20/80 split)



Molecular replacement success



Important decision making features

Molecular replacement success



	Decision tree with AdaBoost	Perfect classifier
ACC (%)	96	100
Class Error (%)	4	0
Sensitivity (%)	93	100
Specificity (%)	97	100
FPR (%)	3	0
Precision (%)	94	100
F1 score (%)	93	100
ROC AUC (%)	99	100
TP test set	64	69
TN test set	131	135
FP test set	4	0
FN test set	5	0





Map traceability

"Good" vs "bad" map PDB entry: 4DNK

solvent flattening backbone tracing

solvent flattening no backbone tracing no solvent flattening no backbone tracing



Map traceability

cNN



Map traceability

cNN assessment and performance

	Train: traced Test: traced	Train: traced Test: solvent flattened, no build
ACC (%)	96	63
Class Error (%)	4	37
Sensitivity (%)	94	67
Specificity (%)	98	58
FPR (%)	2	42
Precision (%)	98	62
TP test set	621	442
TN test set	648	385
FP test set	12	275
FN test set	39	218

solvent flattening no backbone tracing



no solvent flattening no backbone tracing



3 CCP4 🐶 diamond

https://github.com/DiamondLightSource/python-topaz3

Future plans

- Molecular replacement or experimental phasing success
- -Feedback through Synchweb/ISPyB
- -Include other software
- Map traceability

Applying filters such as Gaussian, mean and median; data augmentation
deep cNN for image denoising (Deep Image Prior, Ulyanov, D., Vedaldi, A., Lempitsky, V., https://arxiv.org/abs/1711.10925)

- -ResNet and others
- -From 2D to 3D
- On the side

-General integration into and querying from Synchweb/ISPyB

- -Integration into CCP4 or some of its individual programs
- -Expanding crystallographic data analysis framework
- -Expanding METRIX, including public access
- -Point/space group classifier
- -Add other bioinformatics prediction tools





Acknowledgements:

Diamond Light Source: James Parkhurst Jenna Elliott (summer student 2018) Tim Guite Dominic Jaques (summer student 2016) Gwyndaf Evans Irakli Sikharulidze CCP4: David Waterman Eugene Krissinel

MRC-LMB: Garib Murshudov

University of Newcastle: Arnaud Baslé

Vollmar, M., Parkhurst, J. M., Jaques, D., Baslé, A., Murshudov, G. N., Waterman, D. and Evans, G. (2019) IUCrJ, The predictive power of data processing statistics, submitted



UK Research and Innovation

