

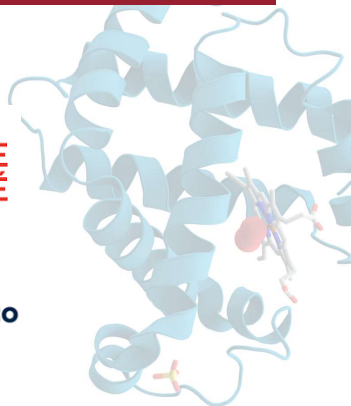
A machine learning playground

Unsupervised and supervised analysis of protein sequences

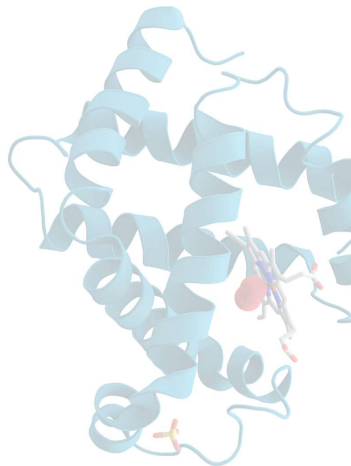
Matteo Allione



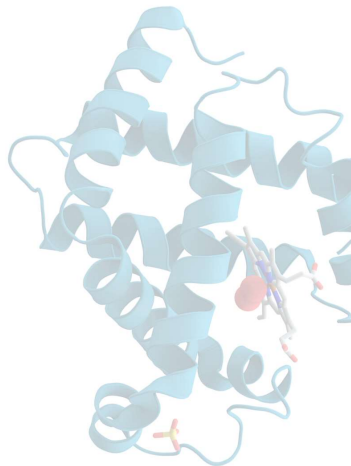
**Politecnico
di Torino**



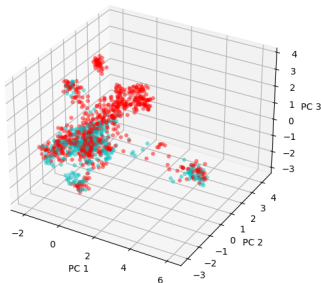
- 1 Data exploration
- 2 Clustering
- 3 Detecting functionality
- 4 Generative models



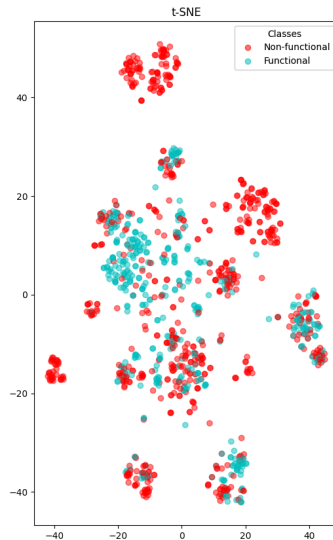
- 1 Data exploration
- 2 Clustering
- 3 Detecting functionality
- 4 Generative models



PCA and t-SNE



PCA with $\sim 10\%$ of the variance (left); t-SNE with
perplexity 40 (right)



Edit distance and ISOMAP

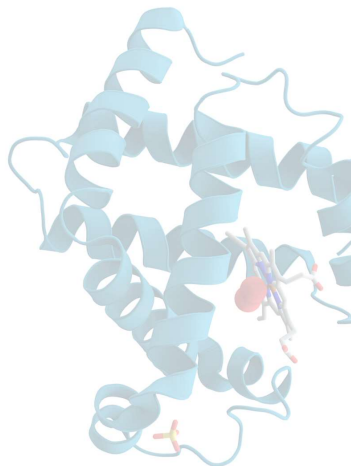
One-hot-encoding paradigm

...	I	L	A	D	L	-	...
...	0	0	1	0	0	0	...
...	0	0	0	0	0	0	...
...	0	0	0	0	0	0	...
...	0	0	0	1	0	0	...
...

$$E : \vec{p} \rightarrow \vec{s} \in \{0, 1\}^{20 \cdot L}$$



- preserves equality between letters
- allows efficient computation of Hamming-like distances



Edit distance and ISOMAP

One-hot-encoding paradigm

...	I	L	A	D	L	-	...
...	0	0	1	0	0	0	...
...	0	0	0	0	0	0	...
...	0	0	0	0	0	0	...
...	0	0	0	1	0	0	...
...

$$E : \vec{p} \rightarrow \vec{s} \in \{0, 1\}^{20 \cdot L}$$



- preserves equality between letters
- allows efficient computation of Hamming-like distances

But in biology we care about the number of mutations that can lead from a protein to another.

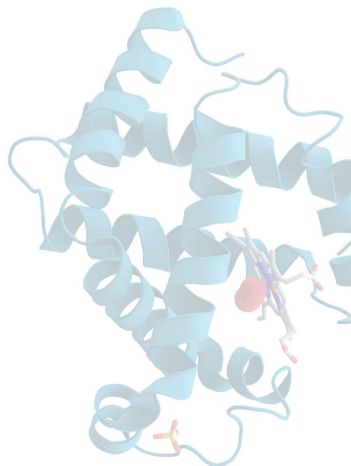
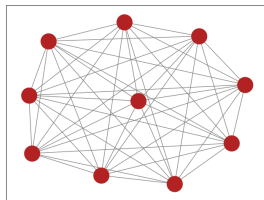
This is encoded in the **edit distance**: allowing deletion, insertion and substitution.

We can compute it efficiently recursively.

Edit distance and ISOMAP

$$\Delta^{ij} = \|z_i - z_j\|^2$$

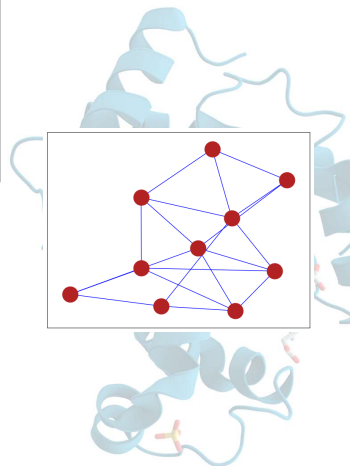
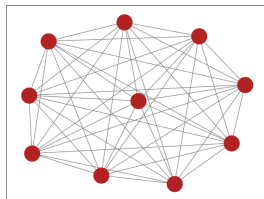
$$\hat{G} = -\frac{1}{2}(\mathbb{I} - \frac{1}{N}\mathbf{U})\Delta(\mathbb{I} - \frac{1}{N}\mathbf{U})$$



Edit distance and ISOMAP

$$\Delta^{ij} = \|z_i - z_j\|^2$$

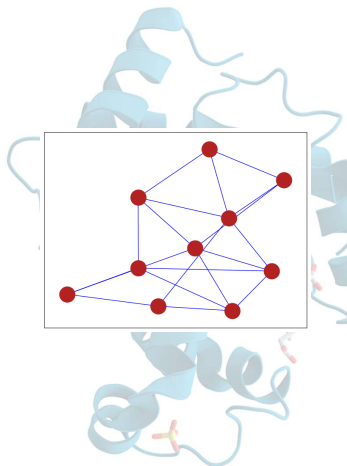
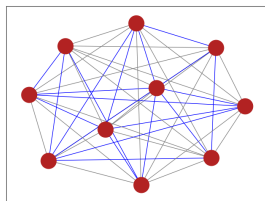
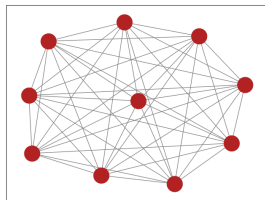
$$\hat{G} = -\frac{1}{2}(\mathbb{I} - \frac{1}{N}\mathbf{U})\Delta(\mathbb{I} - \frac{1}{N}\mathbf{U})$$

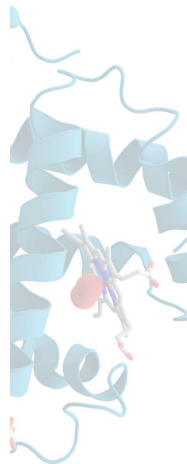
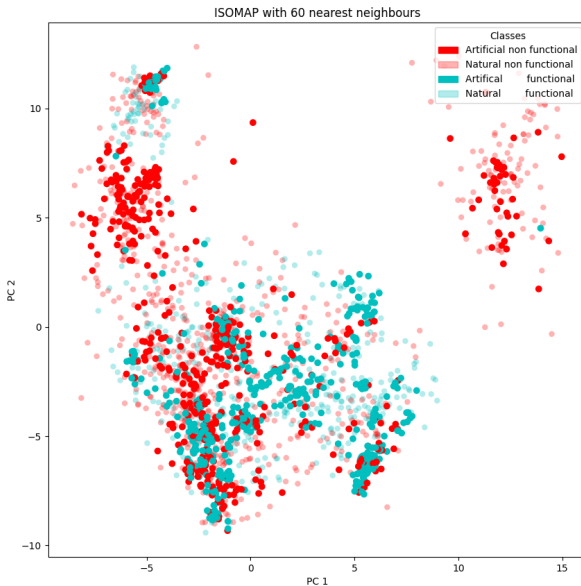


Edit distance and ISOMAP

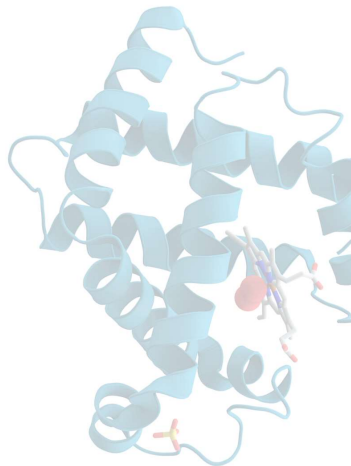
$$\Delta^{ij} = \|z_i - z_j\|^2$$

$$\hat{G} = -\frac{1}{2}(\mathbb{I} - \frac{1}{N}\mathbf{U})\Delta(\mathbb{I} - \frac{1}{N}\mathbf{U})$$

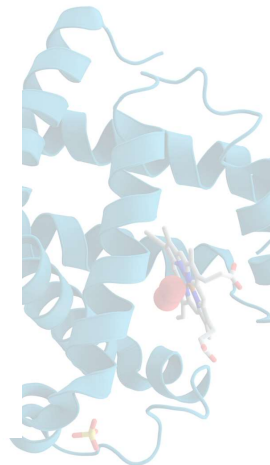
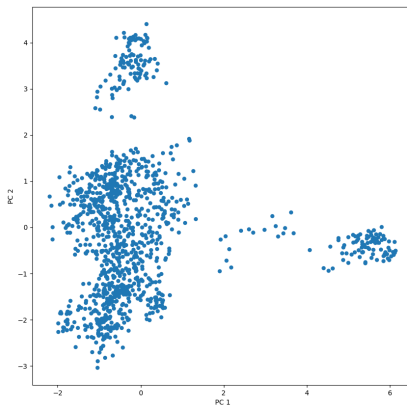




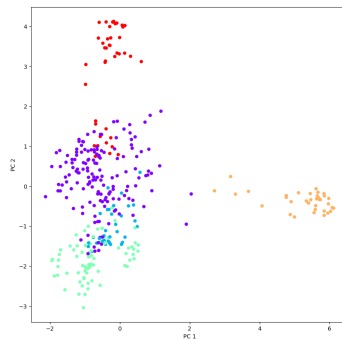
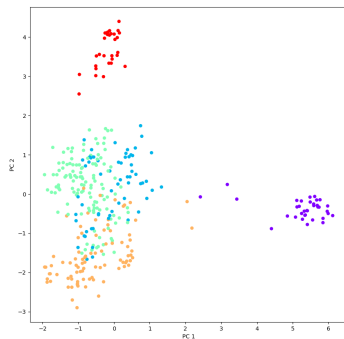
- 1 Data exploration
- 2 Clustering
- 3 Detecting functionality
- 4 Generative models



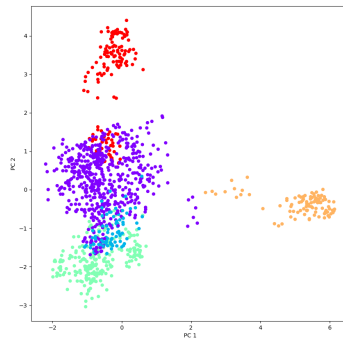
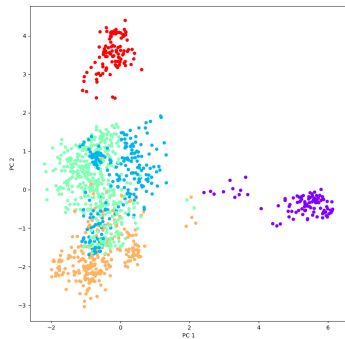
k-means and its consistency



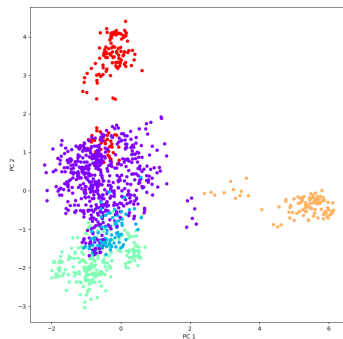
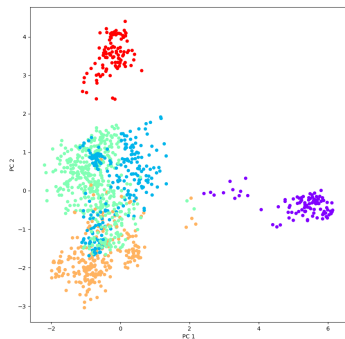
k-means and its consistency



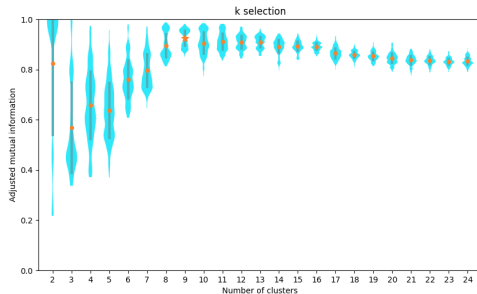
k-means and its consistency



k-means and its consistency

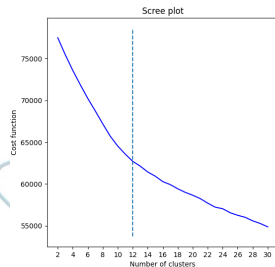
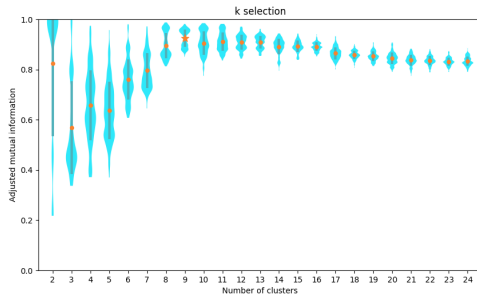


$$I(C_1, C_2) = \sum_{C_1, C_2} p(C_1, C_2) \log \frac{p(C_1, C_2)}{p(C_1)p(C_2)}$$



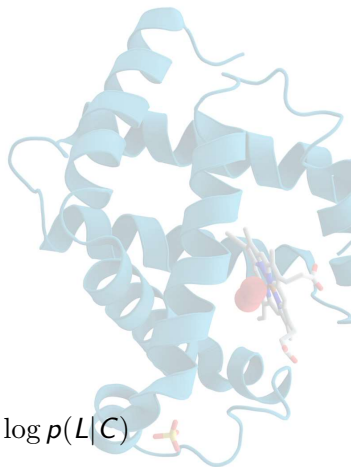
Caveat

To get meaningful information many k-means++ initializations are run for every crop. The mutual information must be 'adjusted'.

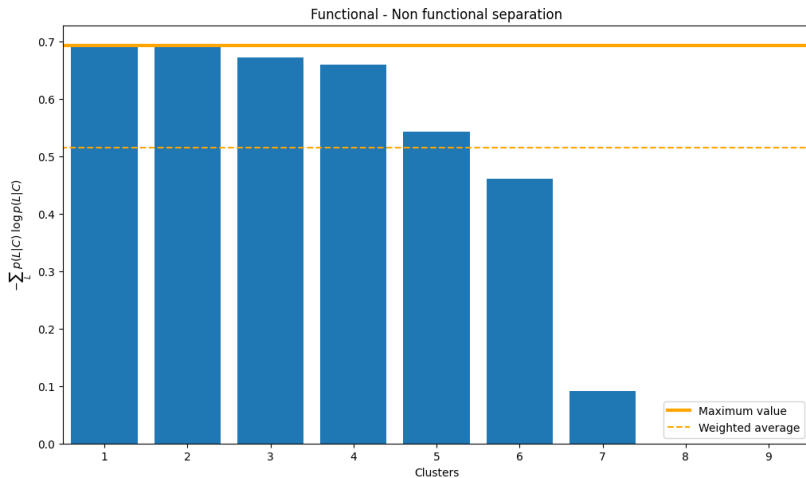


Caveat

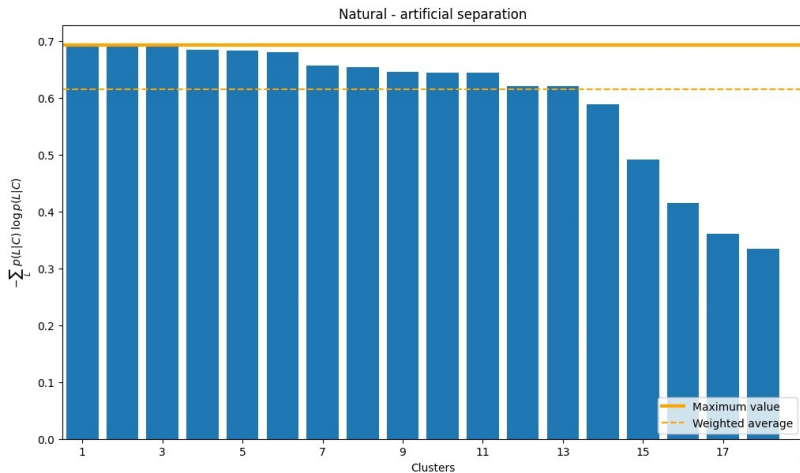
To get meaningful information many k-means++ initializations are run for every crop. The mutual information must be 'adjusted'.



$$\mathbf{H}[L|C] = - \sum_C p(C) \sum_L p(L|C) \log p(L|C)$$

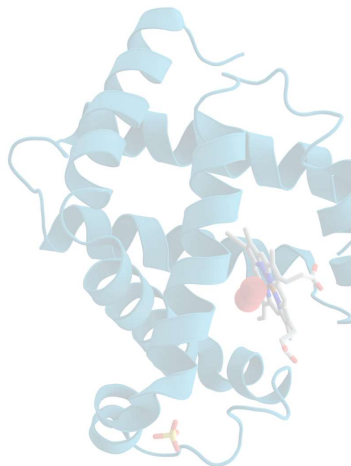


$$\mathbf{H}[L|C] = - \sum_C p(C) \sum_L p(L|C) \log p(L|C)$$

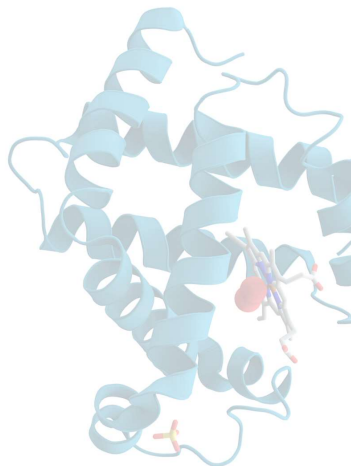
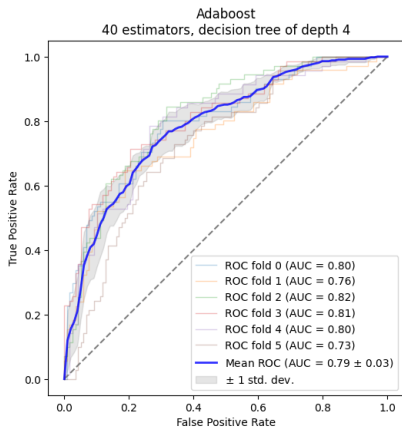


$$\mathbf{H}[L|C] = - \sum_C p(C) \sum_L p(L|C) \log p(L|C)$$

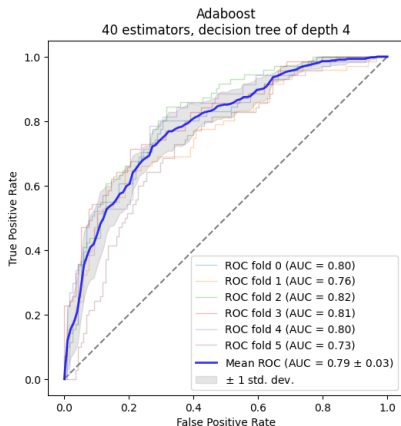
- 1 Data exploration
- 2 Clustering
- 3 Detecting functionality**
- 4 Generative models



Selecting the best model



Selecting the best model



	ROC auc
DT (4)	(0.75 ± 0.03)
DT (7)	(0.76 ± 0.02)
RF (40, 4)	(0.84 ± 0.02)
RF (40, 15)	(0.87 ± 0.02)
RF (300, 15)	(0.88 ± 0.01)
AL (50)	(0.88 ± 0.01)

DT: Decision Tree (depth)

RF: Random Forest

(# of estimators, depth)

AL: Adaboost with logistic regression
(# of estimators)

Networks

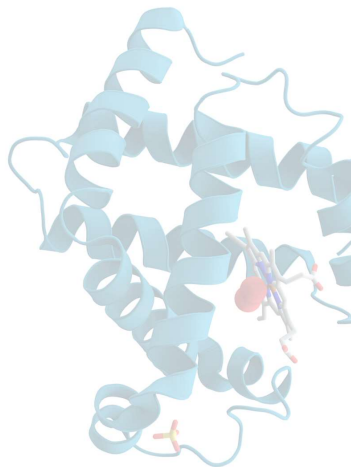
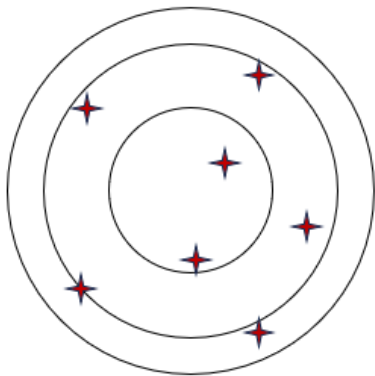
We tried many different architectures and used several tricks to solve some of the problems encountered:

- easily get 100% accuracy on the training set and **overfitting**:
 - inserting **dropout**
 - performing **early stopping**
- problems dealing with an **unbalanced** dataset:
 - using **weighted cost** functions
- **vanishing gradients** with CNN and some sigmoid functions:
 - **batch normalisation**

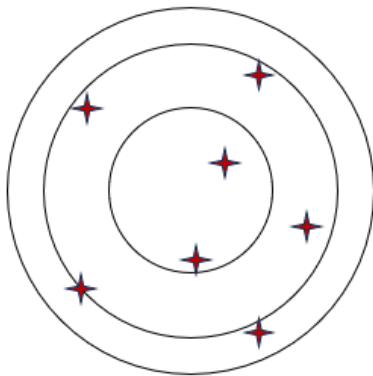
```
class EarlyStopping:
    def __init__(self, patience=10):
        self.counter = 0
        self.best_val = float("inf")
        self.val_loss = None
        self.alpha = 0.01
        self.training_loss = None
        self.validation_loss = None
        self.patience = patience

    def early_stop(self, validation_loss):
        self.val_loss = min(self.val_loss, validation_loss)
        self.validation_loss = validation_loss
        self.training_loss = min(self.training_loss, self.val_loss)
        self.counter += 1
        self.alpha = 1000 / (self.counter + 1)
        if (self.validation_loss > self.patience * self.alpha):
            return True
        return False
```

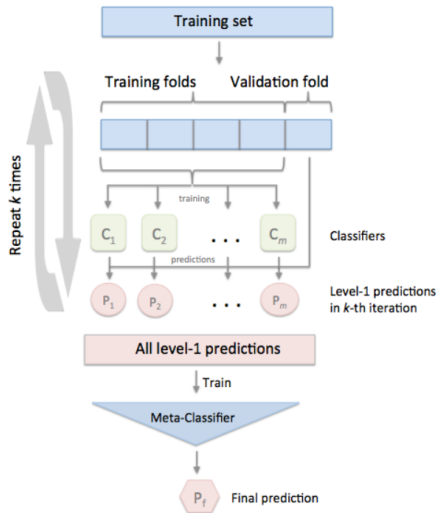

Metaclassifier



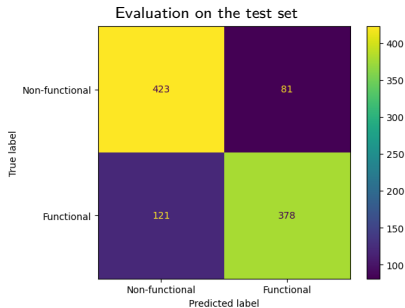
Metaclassifier



Right image from https://rasbt.github.io/mlxtend/user_guide/classifier/StackingCVClassifier/

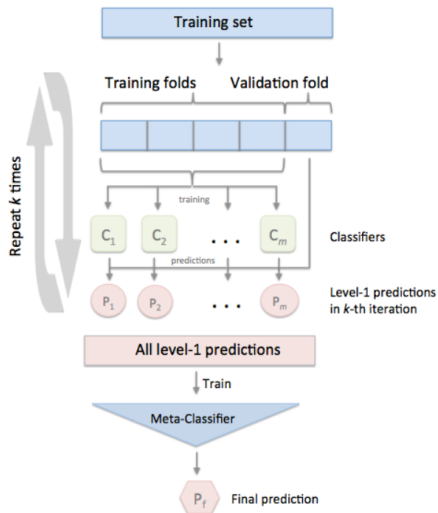


Metaclassifier

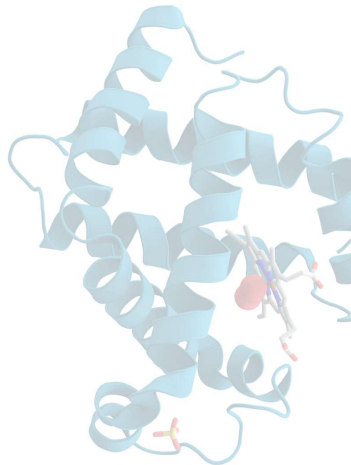


Notice

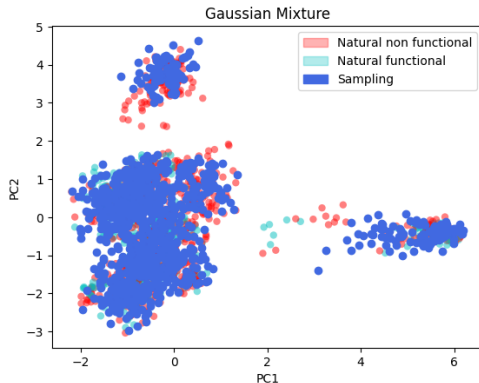
The results shown are obtained through the use of an Adaboost classifier as metaclassifier. With other classifiers it is possible to rank the feature importance and it can be seen that features from classifier that use Isomap embedding can be relevant. In the following parts, for fastness of the computations, we will remove them.



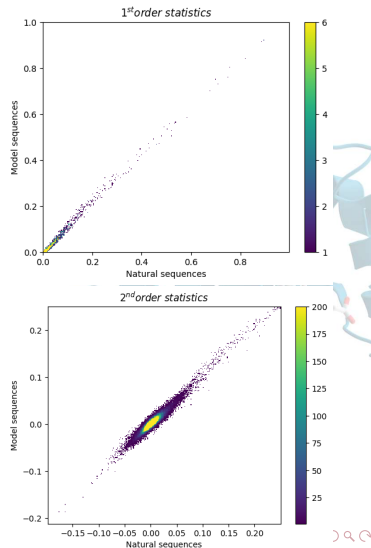
- 1 Data exploration
- 2 Clustering
- 3 Detecting functionality
- 4 **Generative models**



GMM



Using the previously defined classifier you get that 42.6% of the sampled sequences are functional. (Natural functional frequency $\sim 37.4\%$)



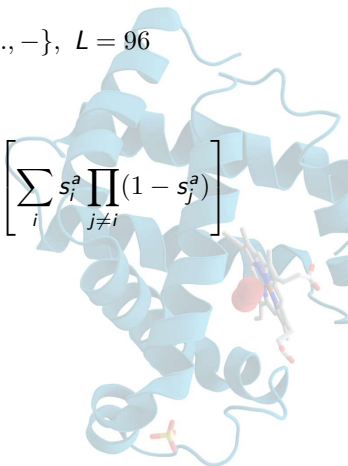
Towards an energy model

Field only pots model:

$$H(\underline{s}) = - \sum_{i=a}^L h_a(s_a) \quad , \quad s_a \in \{A, B, \dots, -\}, \quad L = 96$$

Ising-like model for one-hot-encoded proteins:

$$H(\underline{s}) = - \sum_{i,j,a,b} J_{ij}^{ab} s_i^a s_j^b - \sum_{i,a} h_i^a s_i^a + \gamma \sum_a \left[\sum_i s_i^a \prod_{j \neq i} (1 - s_j^a) \right]$$
$$s_i^a \in \{0, 1\}$$



Towards an energy model

Field only potts model:

$$H(\underline{s}) = - \sum_{i=a}^L h_a(s_a) \quad , \quad s_a \in \{A, B, \dots, -\}, \quad L = 96$$

Ising-like model for one-hot-encoded proteins:

$$H(\underline{s}) = - \sum_{i,j,a,b} J_{ij}^{ab} s_i^a s_j^b - \sum_{i,a} h_i^a s_i^a + \gamma \sum_a \left[\sum_i s_i^a \prod_{j \neq i} (1 - s_j^a) \right]$$

So far: pairwise interactions are sufficient for creating an effective generative model



Our question: is a minimally connected model enough to explain different features of the dataset?

W.P. Russ et al., 'An evolution-based model for designing chorismate mutase enzymes,' Nature, Jul 2020)

Towards an energy model

Field only pots model:

$$H(\underline{s}) = - \sum_{i=a}^L h_a(s_a) \quad , \quad s_a \in \{A, B, \dots, -\}, \quad L = 96$$

Ising-like model for one-hot-encoded proteins:

$$H(\underline{s}) = - \sum_{i,j,a,b} J_{ij}^{ab} s_i^a s_j^b - \sum_{i,a} h_i^a s_i^a + \gamma \sum_a \left[\sum_i s_i^a \prod_{j \neq i} (1 - s_j^a) \right]$$

So far: pairwise interactions are sufficient for creating an effective generative model



Our question: is a minimally connected model enough to explain different features of the dataset?

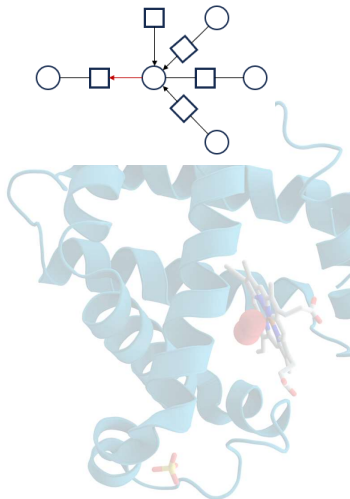
W.P. Russ et al., 'An evolution-based model for designing chorismate mutase enzymes,' Nature, Jul 2020)

Objective:

$$\{J_{ij}^*, h_i^*\} = \operatorname{argmax}_{\{J_{ij}, h_i\}} \mathcal{L}$$

$$\frac{\partial \mathcal{L}}{\partial h_i} = \beta(\hat{m}_i - \langle s_i \rangle_p) = 0$$

$$\frac{\partial \mathcal{L}}{\partial J_{ij}} = \beta(\hat{c}_{ij} - \langle s_i s_j \rangle_p) = 0$$

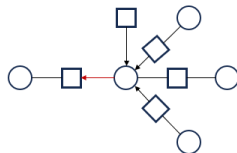


Objective:

$$\{J_{ij}^*, h_i^*\} = \operatorname{argmax}_{\{J_{ij}, h_i\}} \mathcal{L}$$

$$\frac{\partial \mathcal{L}}{\partial h_i} = \beta(\hat{m}_i - \langle s_i \rangle_p) = 0$$

$$\frac{\partial \mathcal{L}}{\partial J_{ij}} = \beta(\hat{c}_{ij} - \langle s_i s_j \rangle_p) = 0$$

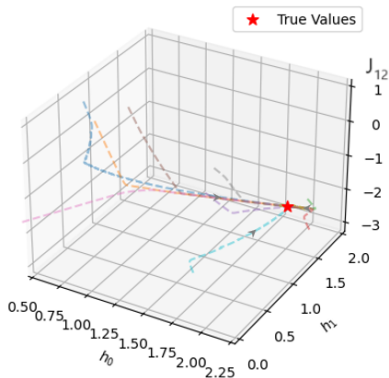


Chow Liu Theorem

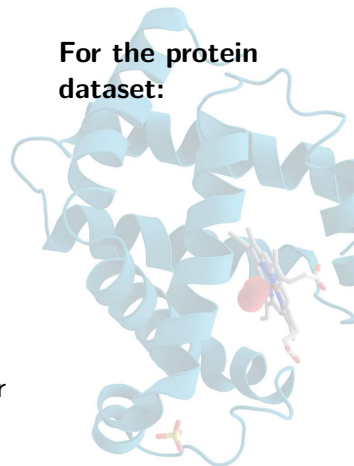
The best topology is the Maximum Spanning Tree of the complete graph with $w_{ij} = M_{ij}$

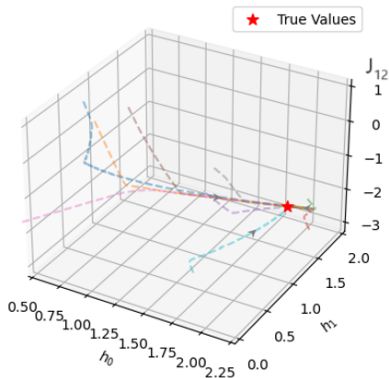
Belief propagation equations

$$b_i(s_i) = \frac{1}{z_i} \prod_{b \in \partial i} m_{bi}(s_i) \quad b_a(\underline{s}_a) = \frac{1}{z_a} \prod_{i \in \partial a} m_{ia}(x_i)$$



Dynamic of the procedure and convergence for
a 2 spin example.





Dynamic of the procedure and convergence for
a 2 spin example.

For the protein
dataset:

**WORK IN
PROGRESS...**