



LOSCHMIDT
LABORATORIES

12. Artificial Intelligence in Life Sciences

B17430 Molecular Biotechnology

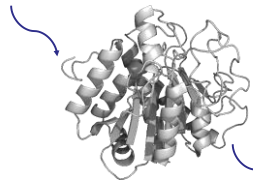
Outline

- Motivation
- Introduction to AI and ML
- Modern challenges in Bioengineering
- Basics of ML
- Recent applications

Motivation

Motivation

```
MSLGAKPFGEKKFIEIKGRRMAYIDEGTGDPILFQHGNTSSYLWRNIMPHC  
AGLRRLIACDLIGMGDSKLDPSGPERYAYAEHRDYLALWEALDLDGRVV  
LVVHDWGSALGFDWARRHREVRVQGIAYMEAIAMPIEWADFPEDRDLFQ  
AFRSQAGEELVQD
```



Function



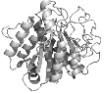
Motivation

Sequence

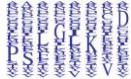
```

MSLGAKPFGEKFKIEIKGRRMAYIDETG
DPILFQHGNTPTSSYLWRNIMPHCAGLGR
LLACCLDGHGSSKLDPSGPIRYATYENR
DYLDALWEALDLDRVLYVYHDWGSAL
GFDWARRHRERVQGIAYMEAIAMPLEW
ADFFQQRDLFQAFRSQAGEELVLD
    
```

Structure



Evolution



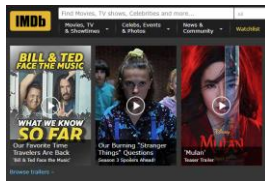
Main approaches:

- Experimental
- Rule-based
- **Machine learning**

Introduction to AI and ML

Introduction to AI and ML

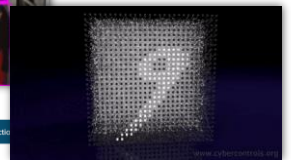
- Recommendation engines
- Gaming
- Image & speech recognition
- Anomaly detection
- Natural language processing
- Data mining
- ...



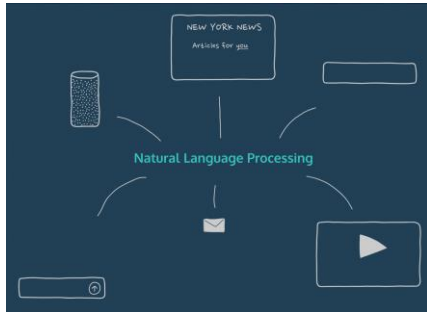
Introduction to AI and ML



Login	Non-Transactional Activities	Transactional
<ul style="list-style-type: none"> • Challenges • Device • Cookie • IP Address • Time of day • Network 	<ul style="list-style-type: none"> • View balance • View history • Updated address • Update email • Update password 	<ul style="list-style-type: none"> • Add new user • Change limits • Set up batch • Set up template • Add payers • ACH • Wire • Bill Pay • Loan Draw



Introduction to AI and ML



From towardsdatascience.com

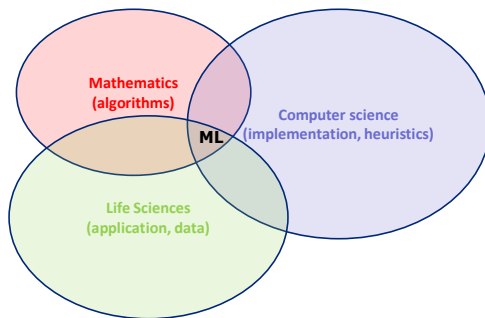
Introduction to AI and ML

The screenshot shows a website header with "DataRobot" and navigation links: PLATFORM, SOLUTIONS, SUCCESS, RESOURCES, PARTNERS, METEOR APP, and CONTACT US. The main content area has a title "Automated Machine Learning: The Competitive Edge You Need" and a sub-headline "Professional sports is a cutthroat industry, on and off the playing field. The smallest of competitive advantages can...". Below this is a soccer field diagram with player movement arrows. To the right is a list of benefits:

- Optimize Player Performance
- Predict — and Prevent — Injuries
- Project Prospect Improvement
- Consolidate Valuable Data
- Increase Profit Potential
- Improve Operational Efficiency

From www.analyticsindiamag.com/data-mining-strategy-development-football/

Introduction to AI and ML



Modern challenges in Bioengineering

Modern Challenges

1. A large number of **relevant features**;
 - From bio-physico-chemical to textual;
 - Heterogeneous data (e.g. clinical, imaging, and genomic data);
2. **Complex tasks** and large parameter space;
 - A single 300-amino-acid-long protein will have 300·19=5700 single-point variants!
3. **Large datasets** available and new data are collected.



Ideal ML setup = complex tasks + relevant features + abundant data

Basics of ML

Basics of ML: features

- Feature vector: $x = (x_1, x_2, \dots, x_n)$
- All features must be converted to numbers:

yes = 1, no = 0



1	1	1	0
1	0	1	0
1	1	0	0
1	0	0	0

ATG TACTGA

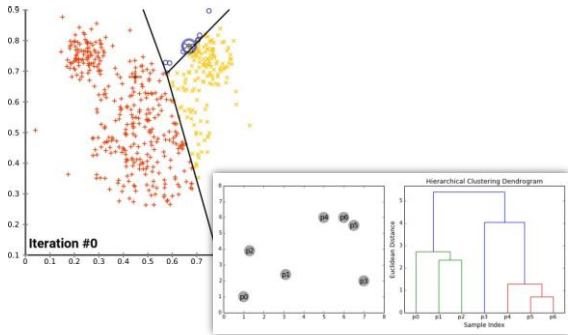
One-hot
encoding

	A	C	G	T
1	1	0	0	0
0	0	0	0	1
0	0	1	0	0
0	0	0	1	0
1	0	0	0	1
0	1	0	0	0
0	0	1	0	0
0	0	0	1	0
1	0	0	0	1

Basics of ML: types

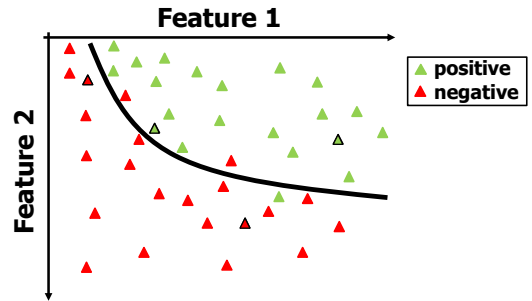
- **Unsupervised learning**
 - only features are available;
 - goal: cluster the data or reduce their dimensionality;
- **Supervised learning**
 - features and labels are available;
 - goal: learn to predict the label based on the features;
- Reinforcement learning, feature learning, anomaly detection, etc...

Basics of ML: unsupervised learning



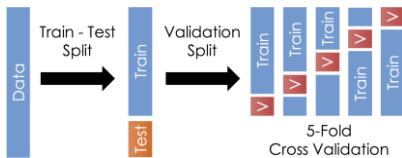
From en.wikipedia.org/wiki/K-means_clustering

Basics of ML: supervised learning



Basics of ML: supervised learning

- We want ML models to be generalizable = good at predicting labels for **previously unseen data**;
- It is essential to split the data into **training set** and **test sets** and use the latter for final evaluation only!
- To fine-tune an algorithm, **K-fold cross validation** is implemented:



Basics of ML: supervised learning

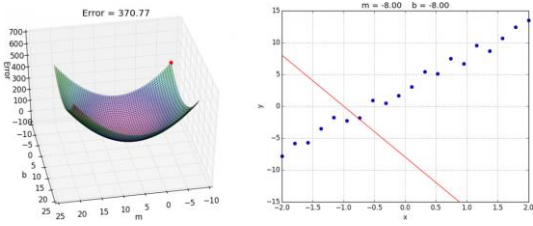
- Let $(x^{(i)}, y^{(i)})$ be our data set, where **x** are **feature values**, and **y** are **labels**;
- Any ML predictor is fundamentally a function **$f(x)$** :

$$f(\text{feature values}) = \text{label}$$

- Usually, a generic group of functions **$f(x, \beta)$** is chosen, where **β** is a set of parameters;
- Then we "train" the ML predictor: pick such **β^*** that **$f(x^{(i)}, \beta^*)$** is as close to **$y^{(i)}$** as possible.

Basics of ML: supervised learning

$$f(x) = m \cdot x + b$$



Recent applications

Recent applications: decision tree

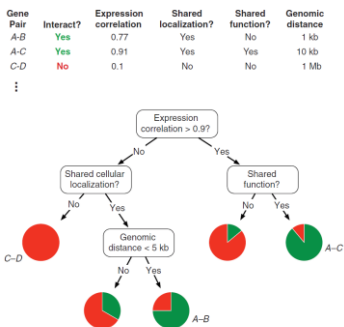
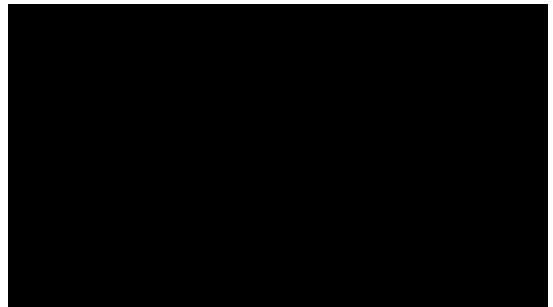


Figure 1 A hypothetical example of how a decision tree might predict protein-protein interactions.

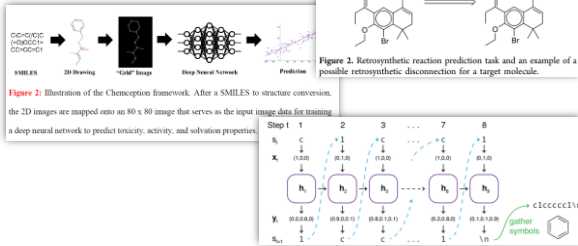
Kingsford, Carl, and Steven L. Salzberg. "What are decision trees?" *Nature biotechnology* 26.9 (2008): 1011.

Recent applications: neural networks



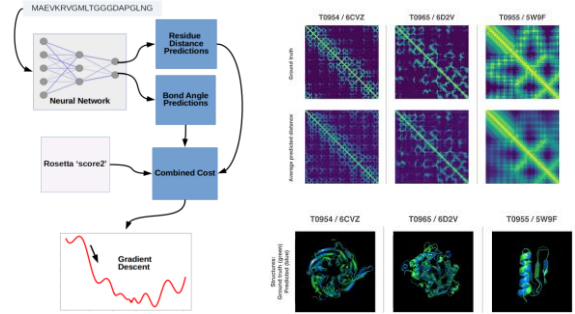
Recent applications: ANNs for drug design

- QSAR modeling for ligand binding activity;
- Structure generation;
- Synthetic pathway generation;



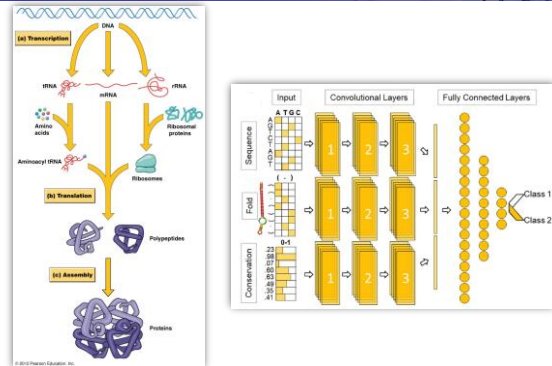
Goh et al. 2017 Chemception: a deep neural network with minimal chemistry knowledge matches the performance of expert-developed QSAR/QSPR models
 Segler et al. 2018 Planning chemical syntheses with deep neural networks and symbolic AI. Nature

Recent applications: AlphaFold for protein folding



<https://deeptmind.com/blog/article/alphafold>
<https://doornik.com/articles/a-summary-of-deepminds-protein-folding-upset-at-ca>

Recent applications: MuStARD for genome annotation



<http://www.mun.ca/biology/desmid/brian/BIOL2060-21/CB21.html>
 Georgakitis, GK et al. "MuStARD: a Deep Learning method for intra- and inter-species scanning identification of small RNA molecules." bioRxiv (2019)

Reading

- Eraslan G et al. 2019 "Deep learning: new computational modelling techniques for genomics." Nature Reviews Genetics. (especially pages 2-6)

