



UiO • Universitetet i Oslo

# Machine Learning in Computational Biology: Overview

IN-BIOS5000/IN-BIOS9000

Milena Pavlović  
Department of Informatics

[milenpa@ui.no](mailto:milenpa@ui.no)

# Disclaimer

I am a machine learning researcher, not a biologist:  
you are the experts there!

# Learning aims

- ❑ Key points should be the intuition and high-level understanding of what machine learning is, types of problems it can help solving
- ❑ Machine learning is not a black box: every choice we make has a meaning
- ❑ Overall understanding data representation and a machine learning algorithm
- ❑ High-level understanding of machine learning workflow, comparison and uncertainty related to it

What is your previous experience with machine learning?



# AI in Nobel Prizes

Illustrations: Niklas Elmehed

## THE NOBEL PRIZE IN CHEMISTRY 2024



David  
Baker

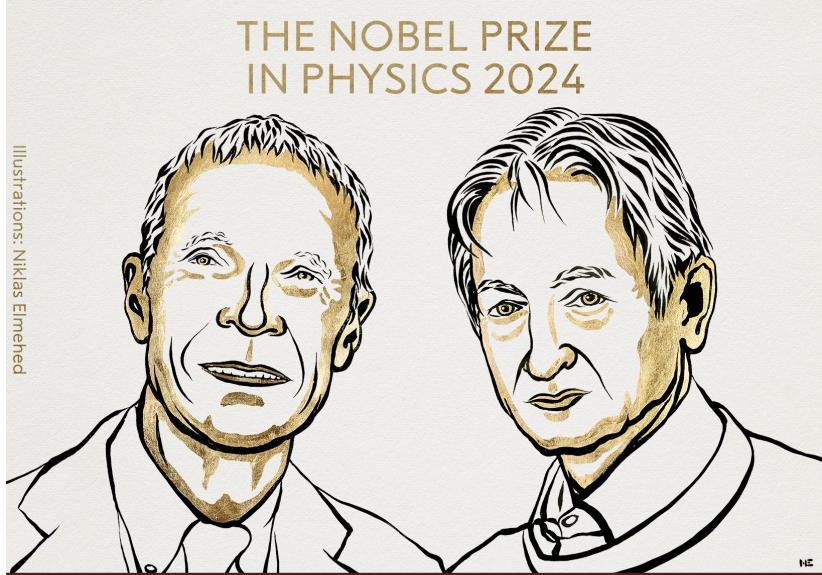
"for computational  
protein design"

Demis  
Hassabis

"for protein structure prediction"

John M.  
Jumper

THE ROYAL SWEDISH ACADEMY OF SCIENCES



John J. Hopfield

"for foundational discoveries and inventions  
that enable machine learning  
with artificial neural networks"

Geoffrey E. Hinton

THE ROYAL SWEDISH ACADEMY OF SCIENCES

# AI for Science

Review | Published: 02 August 2023

## Scientific discovery in the age of artificial intelligence

Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak, Shengchao Liu, Peter Van Katwyk, Andreea Deac, Anima Anandkumar, Karianne Bergen, Carla P. Gomes, Shirley Ho, Pushmeet Kohli, Joan Lasenby, Jure Leskovec, Tie-Yan Liu, Arjun Manrai, Debora Marks, Bharath Ramsundar, Le Song, Jimeng Sun, Jian Tang, ... Marinka Zitnik  + Show authors

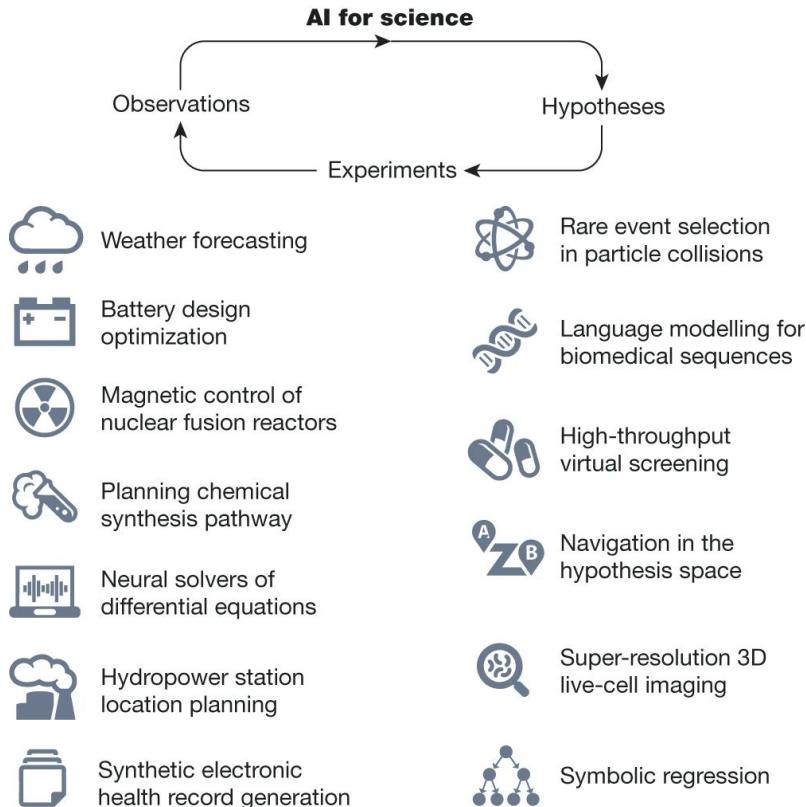
Nature 620, 47–60 (2023) | [Cite this article](#)

Perspective | Published: 06 March 2024

## Artificial intelligence and illusions of understanding in scientific research

Lisa Messeri  & M. J. Crockett 

Nature 627, 49–58 (2024) | [Cite this article](#)



# AI for Science

---

Review | Published: 02 August 2023

## Scientific discovery in the age of artificial intelligence

Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak,  
Shengchao Liu, Peter Van Katwyk, Andreea Deac, Anima Anandkumar, Karianne Bergen, Carla P.  
Gomes, Shirley Ho, Pushmeet Kohli, Joan Lasenby, Jure Leskovec, Tie-Yan Liu, Arjun Manrai, Debora  
Marks, Bharath Ramsundar, Le Song, Jimeng Sun, Jian Tang, ... Marinka Zitnik  + Show authors

[Nature](#) 620, 47–60 (2023) | [Cite this article](#)

Although scientific practices and procedures vary across stages of scientific research, the development of AI algorithms cuts across traditionally isolated disciplines (Fig. 1). Such algorithms can enhance the design and execution of scientific studies. They are becoming indispensable tools for researchers by optimizing parameters and functions<sup>4</sup>, automating procedures to collect, visualize, and process data<sup>5</sup>, exploring vast spaces of candidate hypotheses to form theories<sup>6</sup>, and generating hypotheses and estimating their uncertainty to suggest relevant experiments<sup>7</sup>.

Let's start with a very simple example...

# Sequencing technologies provide data which can be examined for biological properties

CDR3	V gene	J gene	Species
CAAAERNTGELFF	TRBV28*01	TRBJ2-2*01	HomoSapiens
CAAGVENTGELFF	TRBV5-6*01	TRBJ2-2*01	HomoSapiens
CAAQATNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAAQDSNTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAAQMTNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAAQNLTNTGELFF	TRBV15*01	TRBJ2-2*01	HomoSapiens
CAARDQRDLNTGELFF	TRBV2*01	TRBJ2-2*01	HomoSapiens
CAASDPNTGELFF	TRBV12-3*01	TRBJ2-2*01	HomoSapiens
CAASEMNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CACQELNTGELFF	TRBV30*01	TRBJ2-2*01	HomoSapiens
CAEGELNTGELFF	TRBV7-2*01	TRBJ2-2*01	HomoSapiens
CAGADSNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGDYLNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGGDPNTGELFF	TRBV7-9*01	TRBJ2-2*01	HomoSapiens
CAGGDSNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGGRGNTGELFF	TRBV12-3*01	TRBJ2-2*01	HomoSapiens
CAGGVNPNTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAGQDLNTGELFF	TRBV7-2*01	TRBJ2-2*01	HomoSapiens
CAGQNLNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAGQRANTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAIADANTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAIGDENTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAIGDRNTGELFF	TRBV5-5*01	TRBJ2-2*01	HomoSapiens
CAIGDRSSGEQYF	TRBV5-4*01	TRBJ2-7*01	HomoSapiens
CAIQDLNTGELFF	TRBV13*01	TRBJ2-2*01	HomoSapiens
CAIQESNTGELFF	TRBV10-3*01	TRBJ2-2*01	HomoSapiens
CAIQYANTGELFF	TRBV15*01	TRBJ2-2*01	HomoSapiens
CAIRTSGMLNTGELFF	TRBV2*01	TRBJ2-2*01	HomoSapiens

Motifs and data from VDJdb (Bagaev et al. 2020)

# Sequencing technologies provide data which can be examined for biological properties

CDR3	V gene	J gene	Species
CAAAERNTGELFF	TRBV28*01	TRBJ2-2*01	HomoSapiens
CAAGVENTGELFF	TRBV5-6*01	TRBJ2-2*01	HomoSapiens
CAAQATNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAAQDSNTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAAQMTNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAAQNLTNTGELFF	TRBV15*01	TRBJ2-2*01	HomoSapiens
CAARDQRDLNTGELFF	TRBV2*01	TRBJ2-2*01	HomoSapiens
CAASDPNTGELFF	TRBV12-3*01	TRBJ2-2*01	HomoSapiens
CAASEMNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CACQELENNTGELFF	TRBV30*01	TRBJ2-2*01	HomoSapiens
CAEGELENNTGELFF	TRBV7-2*01	TRBJ2-2*01	HomoSapiens
CAGADSNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGDYLNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGGDPNTGELFF	TRBV7-9*01	TRBJ2-2*01	HomoSapiens
CAGGDSNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGGRGNTGELFF	TRBV12-3*01	TRBJ2-2*01	HomoSapiens
CAGGVNPNTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAGQDLENNTGELFF	TRBV7-2*01	TRBJ2-2*01	HomoSapiens
CAGQNLNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAGQRANTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAIADANTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAIGDENTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAIGDRNTGELFF	TRBV5-5*01	TRBJ2-2*01	HomoSapiens
CAIGDRSSGEQYF	TRBV5-4*01	TRBJ2-7*01	HomoSapiens
CAIQDLNTGELFF	TRBV13*01	TRBJ2-2*01	HomoSapiens
CAIQESNTGELFF	TRBV10-3*01	TRBJ2-2*01	HomoSapiens
CAIQYANTGELFF	TRBV15*01	TRBJ2-2*01	HomoSapiens
CAIRTSGMLNTGELFF	TRBV2*01	TRBJ2-2*01	HomoSapiens

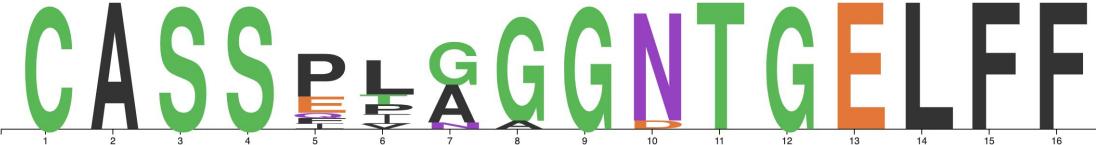
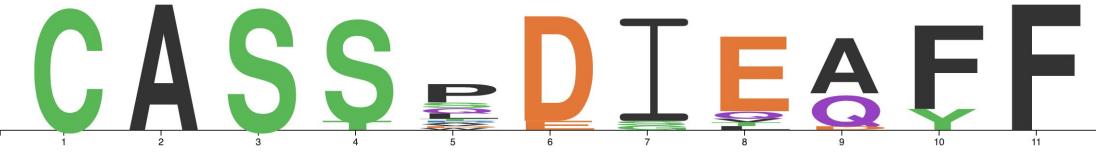
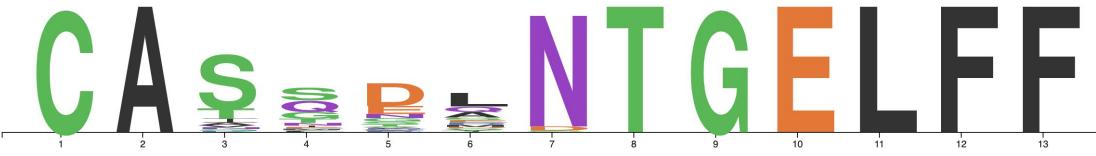
discover  
motifs in  
the data



# Sequencing technologies provide data which can be examined for biological properties

CDR3	V gene	J gene	Species
CAAAERNTGELFF	TRBV28*01	TRBJ2-2*01	HomoSapiens
CAAGVENTGELFF	TRBV5-6*01	TRBJ2-2*01	HomoSapiens
CAAQATNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAAQDSNTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAAQMTNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAAQNLTNTGELFF	TRBV15*01	TRBJ2-2*01	HomoSapiens
CAARDQRDLNTGELFF	TRBV2*01	TRBJ2-2*01	HomoSapiens
CAASDPNTGELFF	TRBV12-3*01	TRBJ2-2*01	HomoSapiens
CAASEMNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CACQELENNTGELFF	TRBV30*01	TRBJ2-2*01	HomoSapiens
CAEGELENNTGELFF	TRBV7-2*01	TRBJ2-2*01	HomoSapiens
CAGADSNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGDYLNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGGDPNTGELFF	TRBV7-9*01	TRBJ2-2*01	HomoSapiens
CAGGDSNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGGRGNTGELFF	TRBV12-3*01	TRBJ2-2*01	HomoSapiens
CAGGVNPNTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAGQDLNTGELFF	TRBV7-2*01	TRBJ2-2*01	HomoSapiens
CAGQNLNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAGQRANTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAIADANTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAIGDENTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAIGDRNTGELFF	TRBV5-5*01	TRBJ2-2*01	HomoSapiens
CAIGDRSSGEQYF	TRBV5-4*01	TRBJ2-7*01	HomoSapiens
CAIQDLNTGELFF	TRBV13*01	TRBJ2-2*01	HomoSapiens
CAIQESNTGELFF	TRBV10-3*01	TRBJ2-2*01	HomoSapiens
CAIQYANTGELFF	TRBV15*01	TRBJ2-2*01	HomoSapiens
CAIRTSGMLNTGELFF	TRBV2*01	TRBJ2-2*01	HomoSapiens

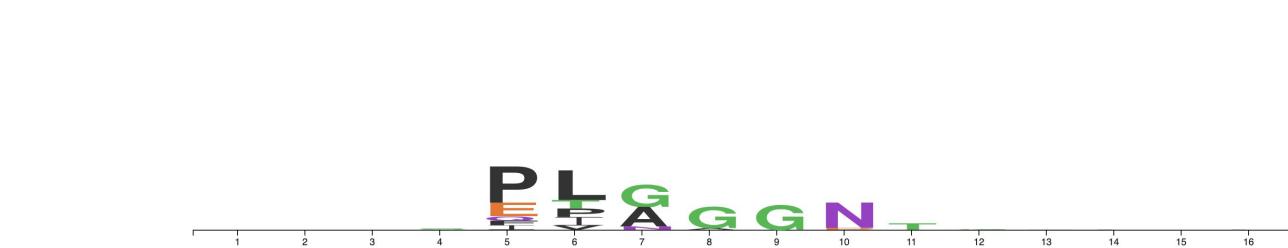
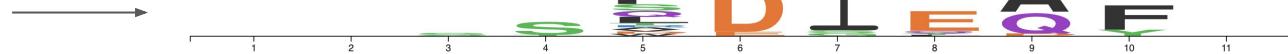
discover  
motifs in  
the data



# Sequencing technologies provide data which can be examined for biological properties

CDR3	V gene	J gene	Species
CAAAERTNTGELFF	TRBV28*01	TRBJ2-2*01	HomoSapiens
CAAGVENTGELFF	TRBV5-6*01	TRBJ2-2*01	HomoSapiens
CAAQATNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAAQDSNTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAAQMTNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAAQNLTNTGELFF	TRBV15*01	TRBJ2-2*01	HomoSapiens
CAARDQRDLNTGELFF	TRBV2*01	TRBJ2-2*01	HomoSapiens
CAASDPNTGELFF	TRBV12-3*01	TRBJ2-2*01	HomoSapiens
CAASEMNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CACQELNTGELFF	TRBV30*01	TRBJ2-2*01	HomoSapiens
CAEGELNTGELFF	TRBV7-2*01	TRBJ2-2*01	HomoSapiens
CAGADSNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGDYLNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGGDPNTGELFF	TRBV7-9*01	TRBJ2-2*01	HomoSapiens
CAGGDSNTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAGGRGNTGELFF	TRBV12-3*01	TRBJ2-2*01	HomoSapiens
CAGGVNPNTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAGQDLNTGELFF	TRBV7-2*01	TRBJ2-2*01	HomoSapiens
CAGQNLNTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAGQRANTGELFF	TRBV19*01	TRBJ2-2*01	HomoSapiens
CAIADANTGELFF	TRBV5-1*01	TRBJ2-2*01	HomoSapiens
CAIGDENTGELFF	TRBV7-8*01	TRBJ2-2*01	HomoSapiens
CAIGDRNTGELFF	TRBV5-5*01	TRBJ2-2*01	HomoSapiens
CAIGDRSSGEQYF	TRBV5-4*01	TRBJ2-7*01	HomoSapiens
CAIQDLNTGELFF	TRBV13*01	TRBJ2-2*01	HomoSapiens
CAIQESNTGELFF	TRBV10-3*01	TRBJ2-2*01	HomoSapiens
CAIQYANTGELFF	TRBV15*01	TRBJ2-2*01	HomoSapiens
CAIRTSGMLNTGELFF	TRBV2*01	TRBJ2-2*01	HomoSapiens

discover  
motifs and  
remove  
genetic  
background



Motifs and data from VDJdb (Bagaev et al. 2020)

# Sequencing technologies provide data which can be examined for biological properties

- ❑ One way to approach an analysis: make a position weight matrix showing product multinomial distribution of amino acids
- ❑ But what if we want to predict if a sequence is specific to a virus or not?



# Machine learning is a powerful approach to discovering patterns in (biological) data

- ❑ A set of methods that allow for making inferences about the data
- ❑ Example: will the receptor bind to the virus or not? - we can fit a logistic regression model on receptor data and then predict binding for new receptors

CAAAERNNTGELFF	+
CAAGVENTGELFF	-
CAAQATNTGELFF	+
CAAQDSNTGELFF	-
CASSADIEQFF	-
CASSADVEAFF	+
CASSASYYEQYF	+
.....	
raw labeled data	

# Machine learning is a powerful approach to discovering patterns in (biological) data

- ❑ A set of methods that allow for making inferences about the data
- ❑ Example: will the receptor bind to the virus or not? - we can fit a logistic regression model on receptor data and then predict binding for new receptors

CAAAERNNTGELFF	+
CAAGVENTGELFF	-
CAAQATNTGELFF	+
CAAQDSNTGELFF	-
CASSADIEQFF	-
CASSADVEAFF	+
CASSASYYEQYF	+
.....	

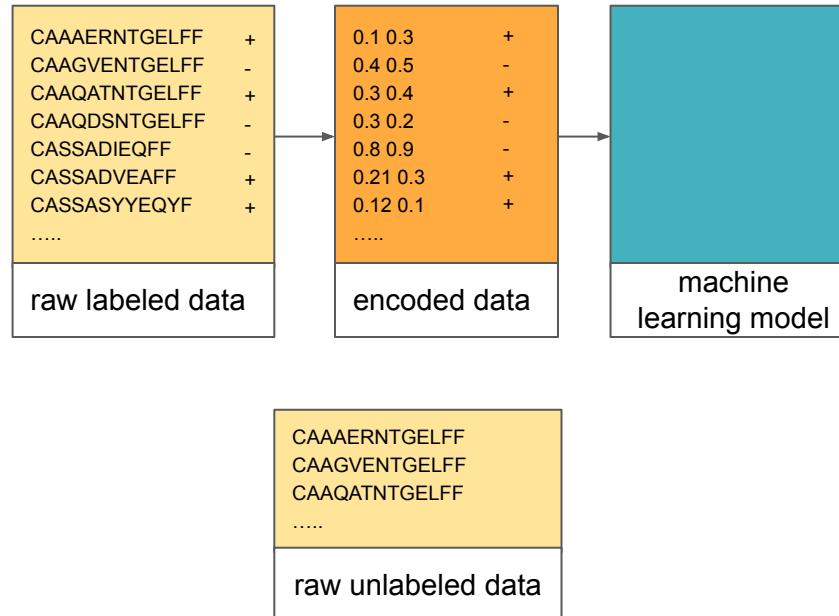
raw labeled data

CAAAERNNTGELFF
CAAGVENTGELFF
CAAQATNTGELFF
.....
.....

raw unlabeled data

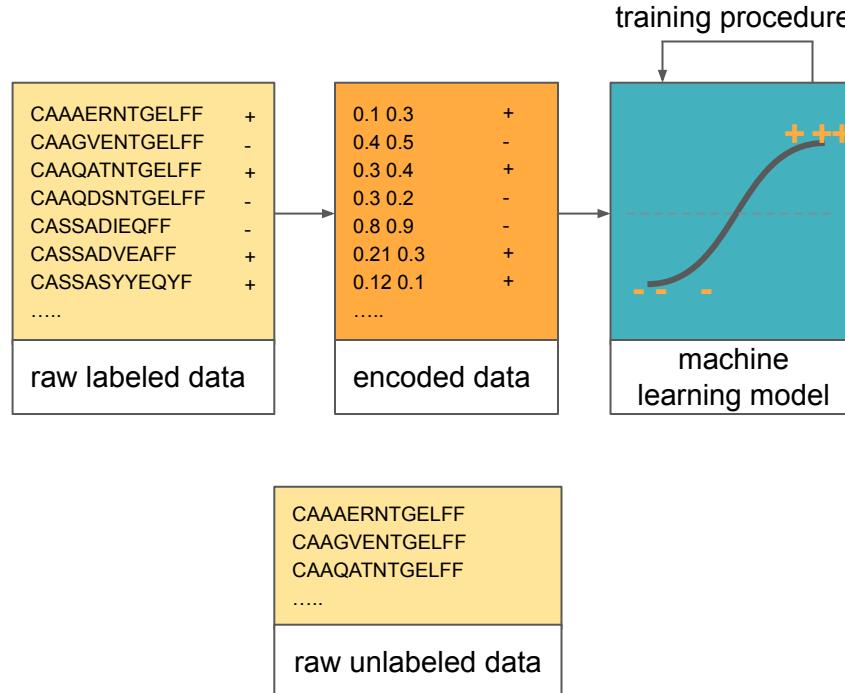
# Machine learning is a powerful approach to discovering patterns in (biological) data

- ❑ A set of methods that allow for making inferences about the data
- ❑ Example: will the receptor bind to the virus or not? - we can fit a logistic regression model on receptor data and then predict binding for new receptors



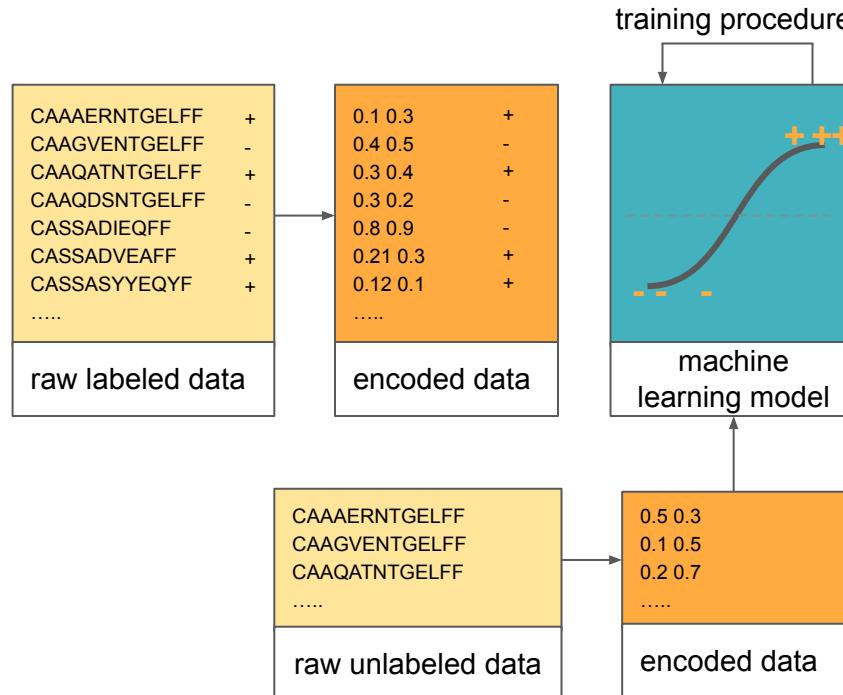
# Machine learning is a powerful approach to discovering patterns in (biological) data

- ❑ A set of methods that allow for making inferences about the data
- ❑ Example: will the receptor bind to the virus or not? - we can fit a logistic regression model on receptor data and then predict binding for new receptors



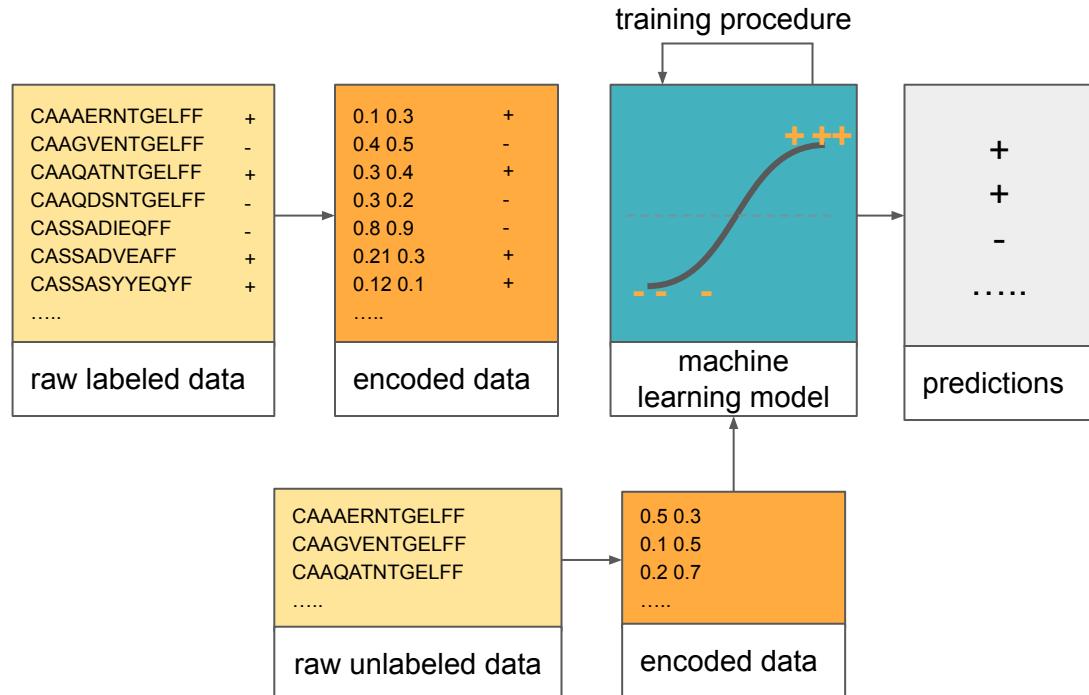
# Machine learning is a powerful approach to discovering patterns in (biological) data

- ❑ A set of methods that allow for making inferences about the data
- ❑ Example: will the receptor bind to the virus or not? - we can fit a logistic regression model on receptor data and then predict binding for new receptors



# Machine learning is a powerful approach to discovering patterns in (biological) data

- ❑ A set of methods that allow for making inferences about the data
- ❑ Example: will the receptor bind to the virus or not? - we can fit a logistic regression model on receptor data and then predict binding for new receptors



# Some concrete ML applications

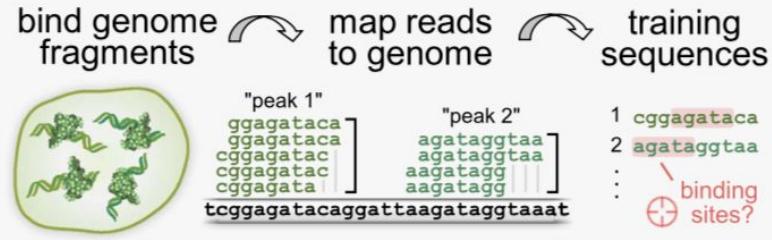
# Transcription factor binding prediction

Transcription factors are proteins which bind to certain sites in DNA and regulate transcription of genes

Given a set of DNA sequences for which we know if they will bind or not, how can we predict if a transcription factor will bind to a new DNA sequence?

Classification problem!

## measuring specificity with sequencing



Leung et al. 2016

Published: 27 July 2015

## Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning

Babak Alipanahi, Andrew Delong, Matthew T Weirauch & Brendan J Frey [✉](#)

*Nature Biotechnology* 33, 831–838(2015) | Cite this article

# Identification of new cell types

Single-cell RNA-seq clustering for identification of new cell types:

A dimensionality reduction technique is applied to normalized count data

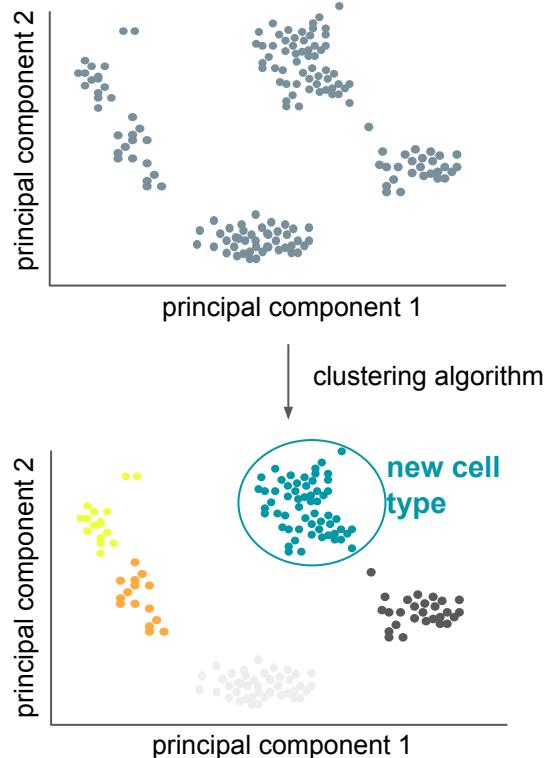
Clustering the data (using e.g., k-means algorithm) can reveal new cell types

Review Article | Published: 07 January 2019

## Challenges in unsupervised clustering of single-cell RNA-seq data

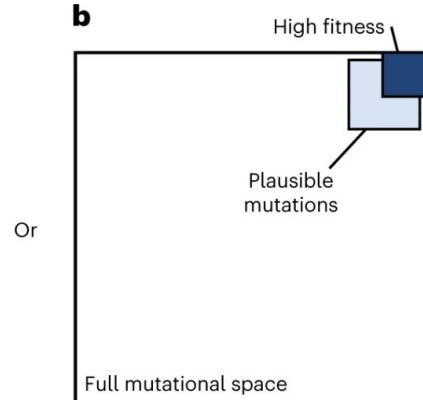
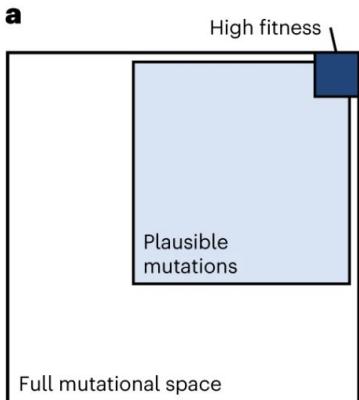
Vladimir Yu Kiselev, Tallulah S. Andrews & Martin Hemberg 

*Nature Reviews Genetics* **20**, 273–282(2019) | Cite this article



# Antibody design

An ensemble of protein language models predicts amino acid substitutions to improve antibody characteristics

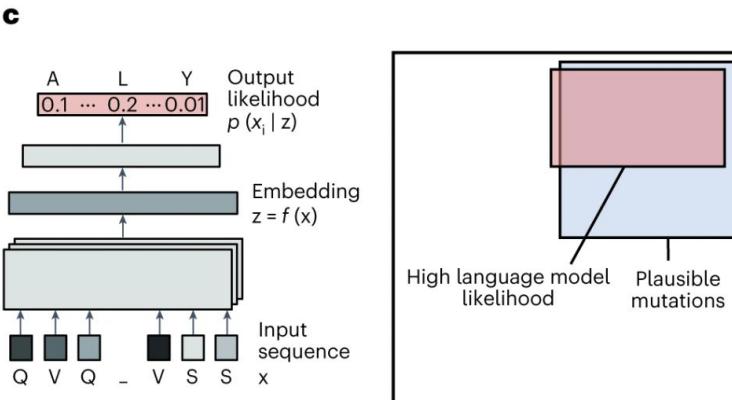


Article | [Open access](#) | Published: 24 April 2023

## Efficient evolution of human antibodies from general protein language models

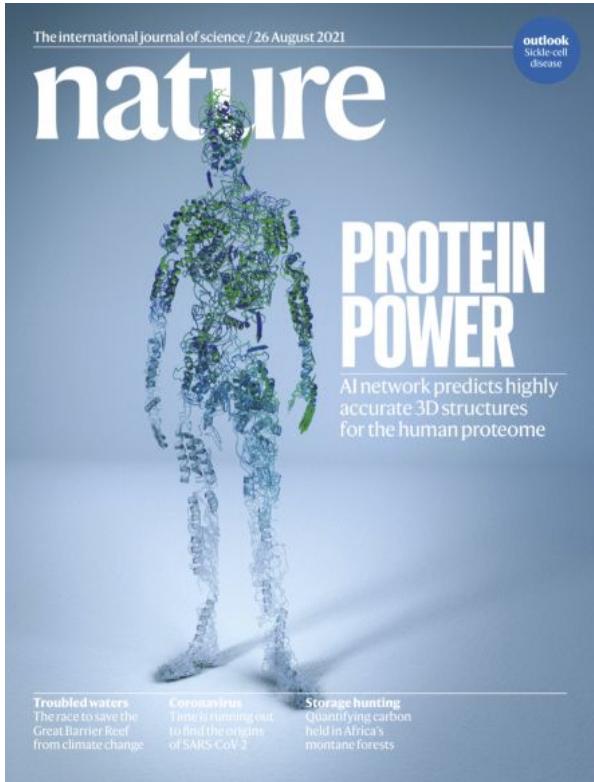
Brian L. Hie , Varun R. Shanker, Duo Xu, Theodora U. J. Bruun, Payton A. Weidenbacher, Shaogeng Tang, Wesley Wu, John E. Pak & Peter S. Kim

*Nature Biotechnology* 42, 275–283 (2024) | [Cite this article](#)



Experiment: use high language model likelihood to sample evolutionarily plausible mutations, measure fitness

# Protein Structure Prediction



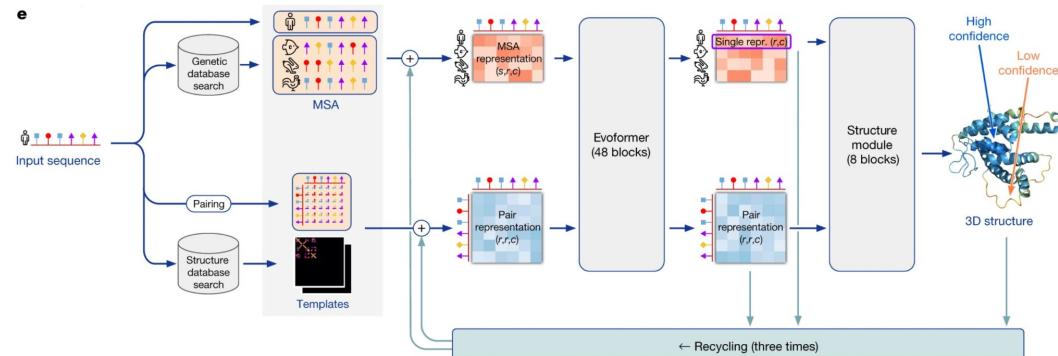
Article | Open Access | Published: 15 July 2021

## Highly accurate protein structure prediction with AlphaFold

John Jumper , Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishabh Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michał Zieliński, Martin Steinegger, Michałina Pacholska, Tamas Berghammer, Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli & Demis Hassabis

— Show fewer authors

*Nature* 596, 583–589 (2021) | [Cite this article](#)



# Design of cis-regulatory elements

Article | [Open access](#) | Published: 23 October 2024

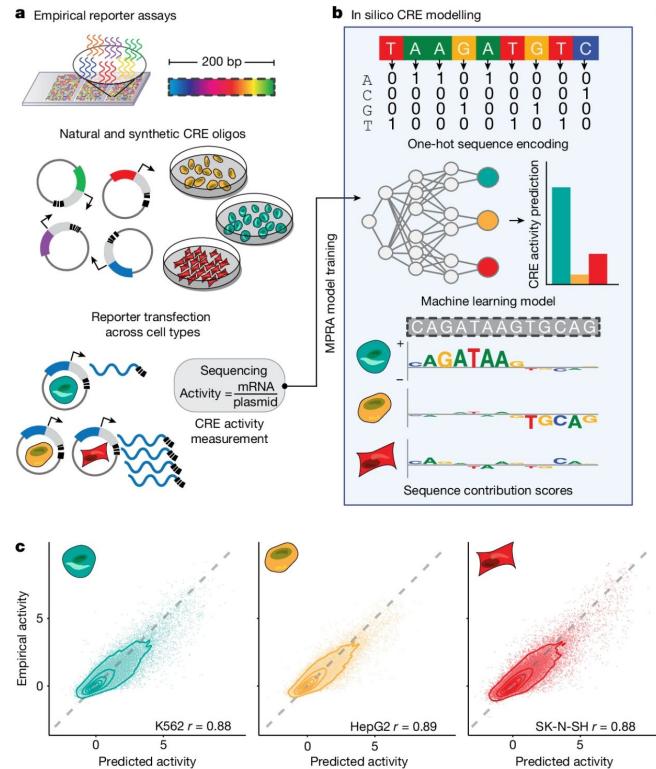
## Machine-guided design of cell-type-targeting *cis*-regulatory elements

Sager J. Gosai , Rodrigo I. Castro , Natalia Fuentes, John C. Butts, Kousuke Mouri, Michael Alasoadura, Susan Kales, Thanh Thanh L. Nguyen, Ramil R. Noche, Arya S. Rao, Mary T. Joy, Pardis C. Sabeti, Steven K. Reilly  & Ryan Tewhey 

*Nature* **634**, 1211–1220 (2024) | [Cite this article](#)

A convolutional neural network predicts cell-type-specific CRE effects directly from the nucleotide sequence across cell lines

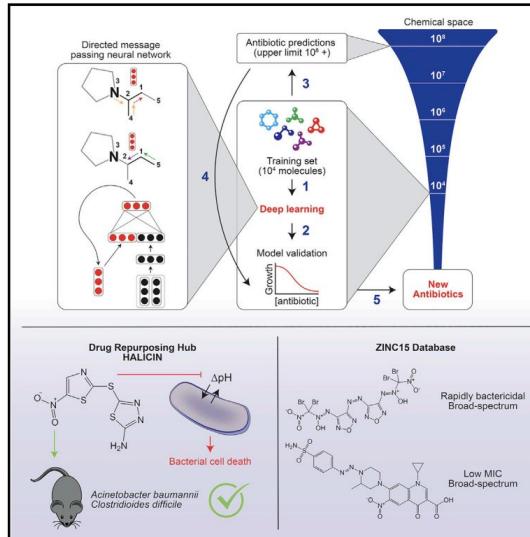
Design of novel CREs with target functionality through sequence activity prediction and iterative optimization algorithms



# Antibiotic discovery

## A Deep Learning Approach to Antibiotic Discovery

### Graphical Abstract



### Authors

Jonathan M. Stokes, Kevin Yang,  
Kyle Swanson, ..., Tommi S. Jaakkola,  
Regina Barzilay, James J. Collins

### Correspondence

regina@csail.mit.edu (R.B.),  
jimjc@mit.edu (J.J.C.)

### In Brief

A trained deep neural network predicts antibiotic activity in molecules that are structurally different from known antibiotics, among which Halicin exhibits efficacy against broad-spectrum bacterial infections in mice.

## Highlights

- A deep learning model is trained to predict antibiotics based on structure
- Halicin is predicted as an antibacterial molecule from the Drug Repurposing Hub
- Halicin shows broad-spectrum antibiotic activities in mice
- More antibiotics with distinct structures are predicted from the ZINC15 database

# Machine learning is applied to a variety of problems in computational biology

Article | [Open access](#) | Published: 22 October 2024

## METASPACE-ML: Context-specific metabolite annotation for imaging mass spectrometry using machine learning

Bishoy Wadie, Lachlan Stuart, Christopher M. Rath, Bernhard Drotleff, Sergii Mamedov & Theodore Alexandrov 

*Nature Communications* 15, Article number: 9110 (2024) | [Cite this article](#)

Article | Published: 03 April 2023

## Deep learning supported discovery of biomarkers for clinical prognosis of liver cancer

Junhao Liang, Weisheng Zhang, Jianghui Yang, Meilong Wu, Qionghai Dai , Hongfang Yin , Ying Xiao  & Lingjie Kong 

*Nature Machine Intelligence* 5, 408–420 (2023) | [Cite this article](#)

Article | [Open access](#) | Published: 15 March 2024

## Foundation model for cancer imaging biomarkers

Suraj Pai, Dennis Bontempi, Ibrahim Hadzic, Vasco Prudente, Mateo Sokač, Tafadzwa L. Chaunzwa, Simon Bernatz, Ahmed Hosny, Raymond H. Mak, Nicolai J. Birkbak & Hugo J. W. L. Aerts 

*Nature Machine Intelligence* 6, 354–367 (2024) | [Cite this article](#)

Article | [Open access](#) | Published: 10 September 2021

## Unified AI framework to uncover deep interrelationships between gene expression and Alzheimer's disease neuropathologies

Nicasia Beebe-Wang, Safiye Celik, Ethan Weinberger, Pascal Sturmfelz, Philip L. De Jager, Sara Mostafavi  & Su-In Lee 

*Nature Communications* 12, Article number: 5369 (2021) | [Cite this article](#)

Article | Published: 01 May 2023

## Uncovering expression signatures of synergistic drug responses via ensembles of explainable machine-learning models

Joseph D. Janizek, Ayse B. Dincer, Safiye Celik, Hugh Chen, William Chen, Kamila Naxerova  & Su-In Lee 

*Nature Biomedical Engineering* 7, 811–829 (2023) | [Cite this article](#)

Article | Published: 26 September 2022

## scBERT as a large-scale pretrained deep language model for cell type annotation of single-cell RNA-seq data

Fan Yang, Wenchuan Wang, Fang Wang, Yuan Fang, Duyu Tang, Junzhou Huang, Hui Lu  & Jianhua Yao 

*Nature Machine Intelligence* 4, 852–866 (2022) | [Cite this article](#)

# Computational biology poses unique challenges for machine learning

- ❑ Dimensionality & dataset size
- ❑ Signal and noise in the data
- ❑ Unknown ground truth and weakly labeled datasets
- ❑ Selection bias

Keep in the data generation process in mind!

# Machine learning in computational biology - outline

- Introduction to machine learning:
  - What is machine learning, types of problems, assumptions, workflow, generalization
- Machine learning models and algorithms:
  - Discriminative vs generative models, supervised models (logistic and linear regression, kNN, neural networks), unsupervised models (dimensionality reduction, clustering)
- Data representation:
  - Considerations and examples, one-hot encoding, feature engineering, representation learning
- Model comparison and uncertainty:
  - Model assessment, model selection, uncertainty, cross-validation
- Transparency and reproducibility

# Machine learning in computational biology - outline

- **Introduction to machine learning:**
  - What is machine learning, types of problems, assumptions, workflow, generalization
- Machine learning models and algorithms:
  - Discriminative vs generative models, supervised models (logistic and linear regression, kNN, neural networks), unsupervised models (dimensionality reduction, clustering)
- Data representation:
  - Considerations and examples, one-hot encoding, feature engineering, representation learning
- Model comparison and uncertainty:
  - Model assessment, model selection, uncertainty, cross-validation
- Transparency and reproducibility

# What is machine learning?

“Machine learning refers to extracting patterns from raw data.”

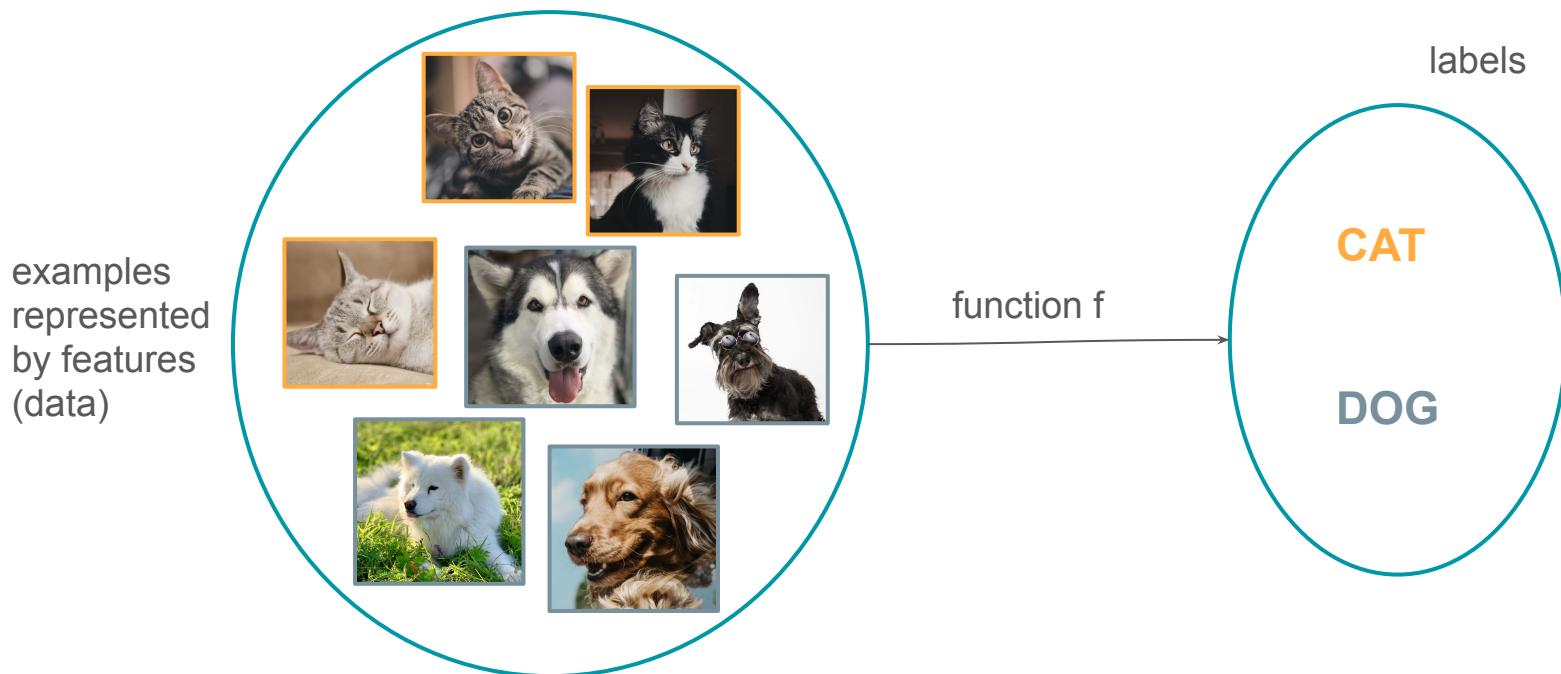
# What is machine learning?

“Machine learning refers to extracting patterns from raw data.”

“A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E.”

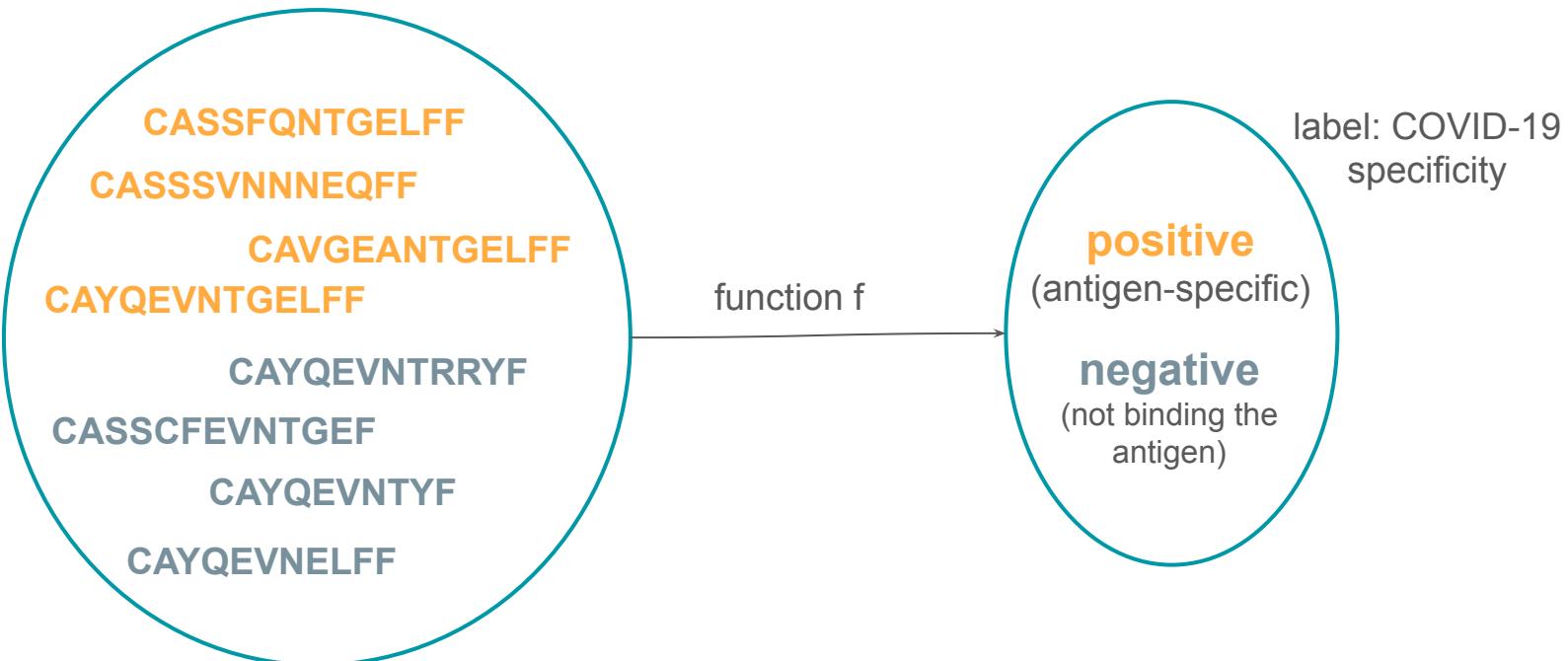
Mitchell 1997

# Machine learning as a function approximation task

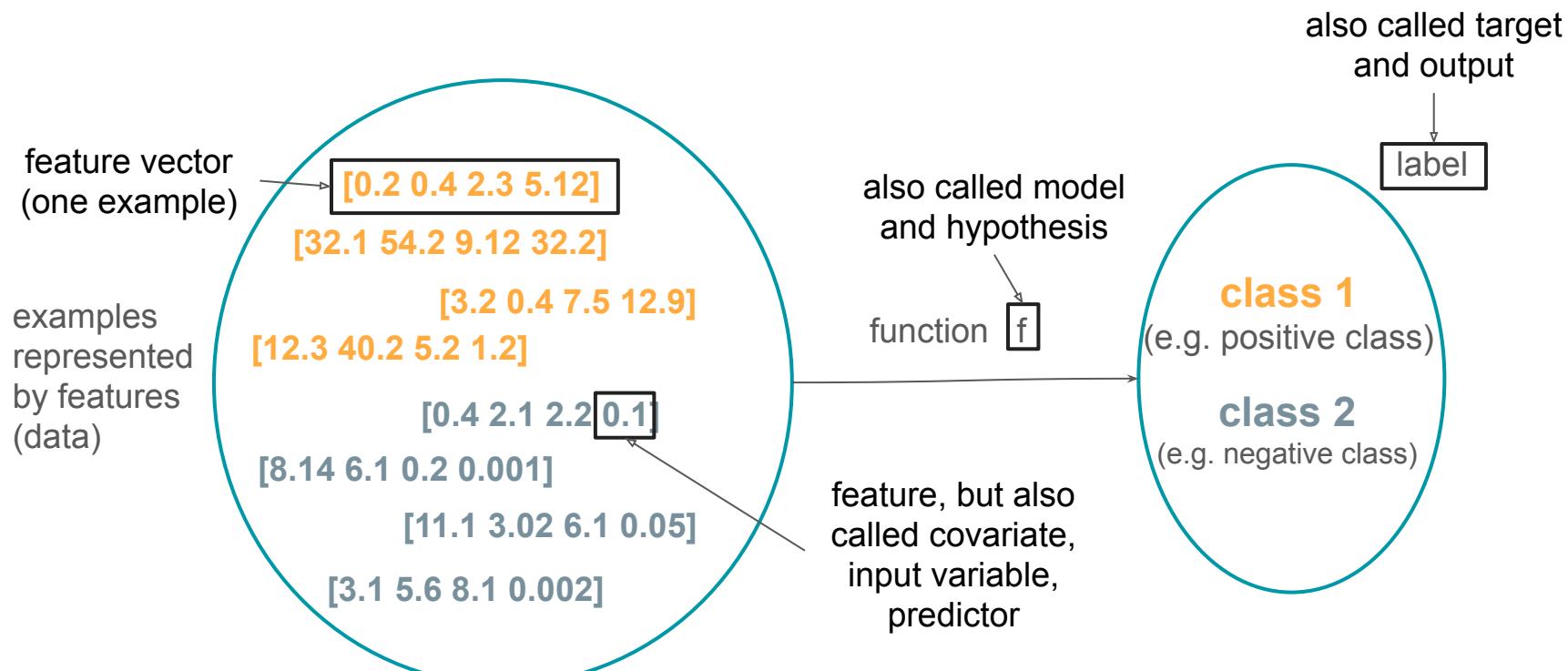


# Machine learning as a function approximation task

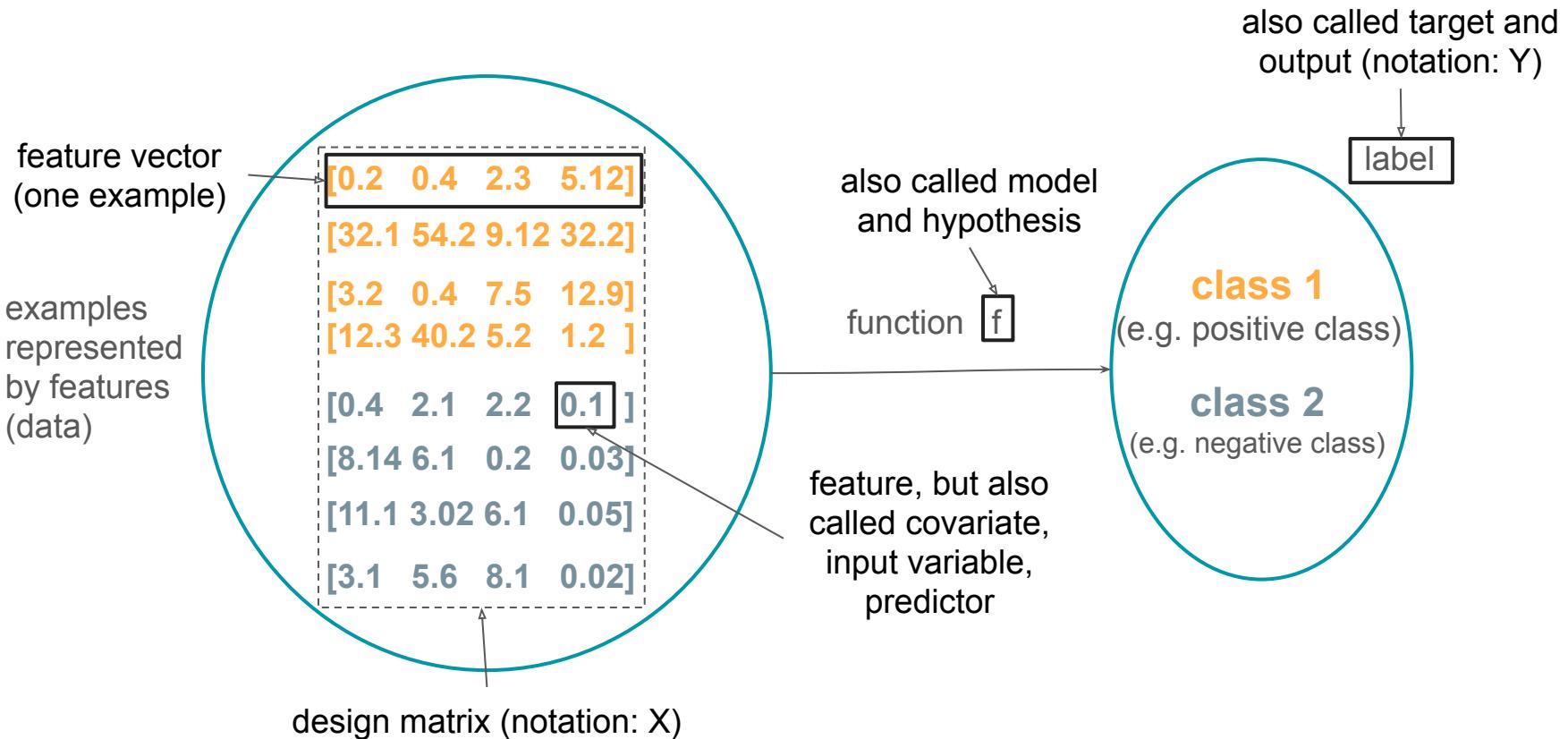
examples  
represented  
by features:  
immune  
receptor  
sequences



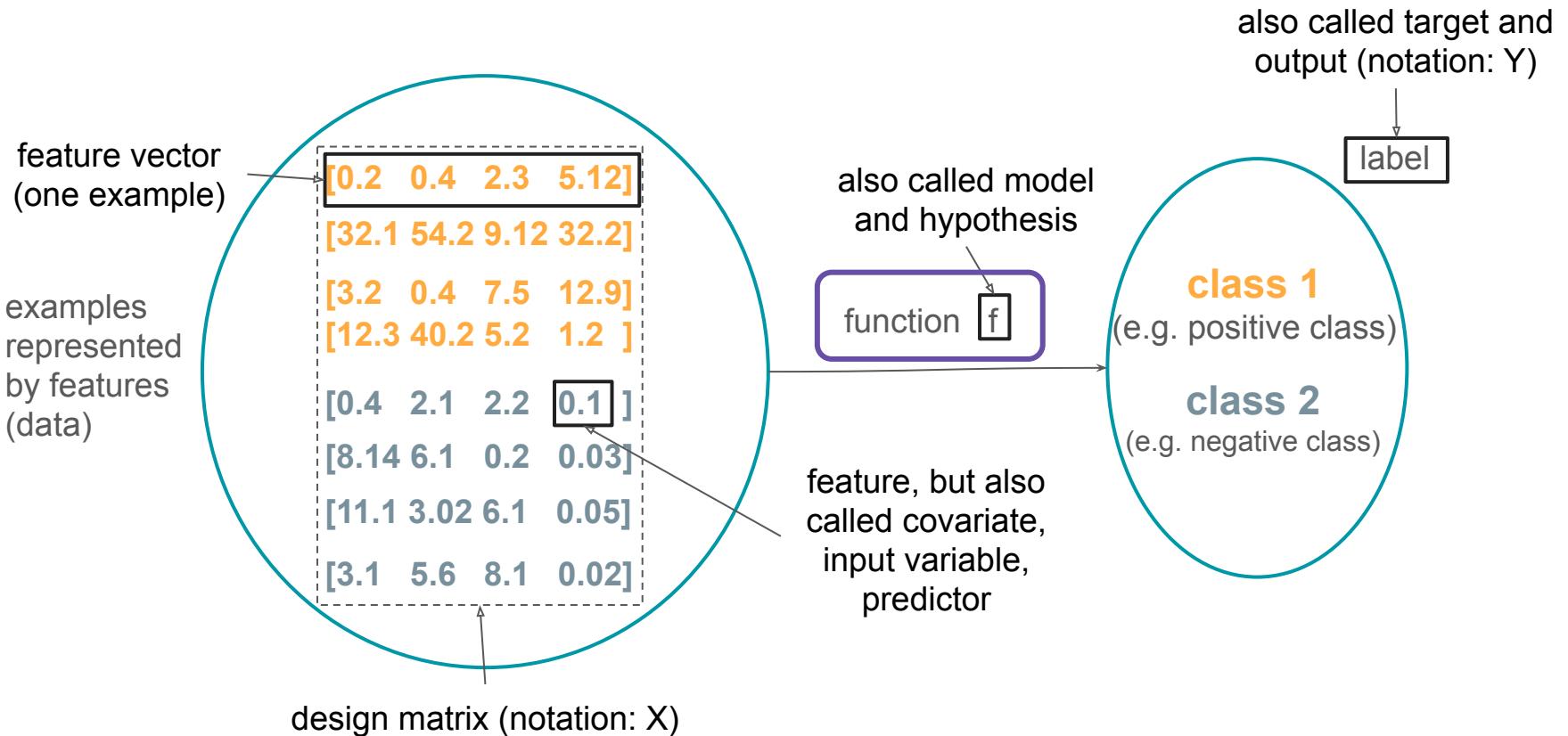
# Machine learning as a function approximation task



# Machine learning as a function approximation task



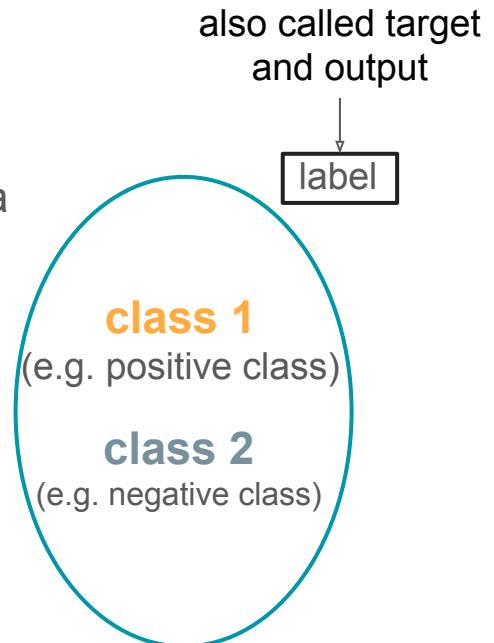
# Machine learning as a function approximation task



# Types of problems in machine learning: what does the function f do

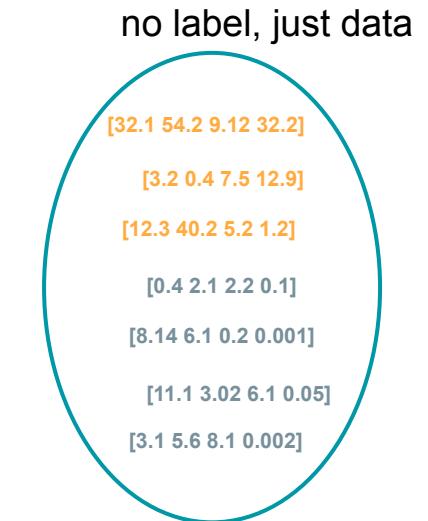
1. **Supervised**: for each example we know the label

- a. Label can be a discrete value - a class (e.g., a receptor is antigen-specific or not, the picture contains a dog or a cat): **classification** or
- b. a continuous value (e.g., binding affinity, house price): **regression**



# Types of problems in machine learning: what does the function f do

1. **Supervised: classification and regression**
2. **Unsupervised:**
  - a. the data we have has a lot of features, and we want to see if there is a structure in the data
  - b. there is no explicit label
  - c. example: there is a set of cells and we want to see if we can group them and see if there are new groups which could indicate new cell types
  - d. clustering, semi-supervised learning (proxy supervised tasks from unlabeled data), density estimation



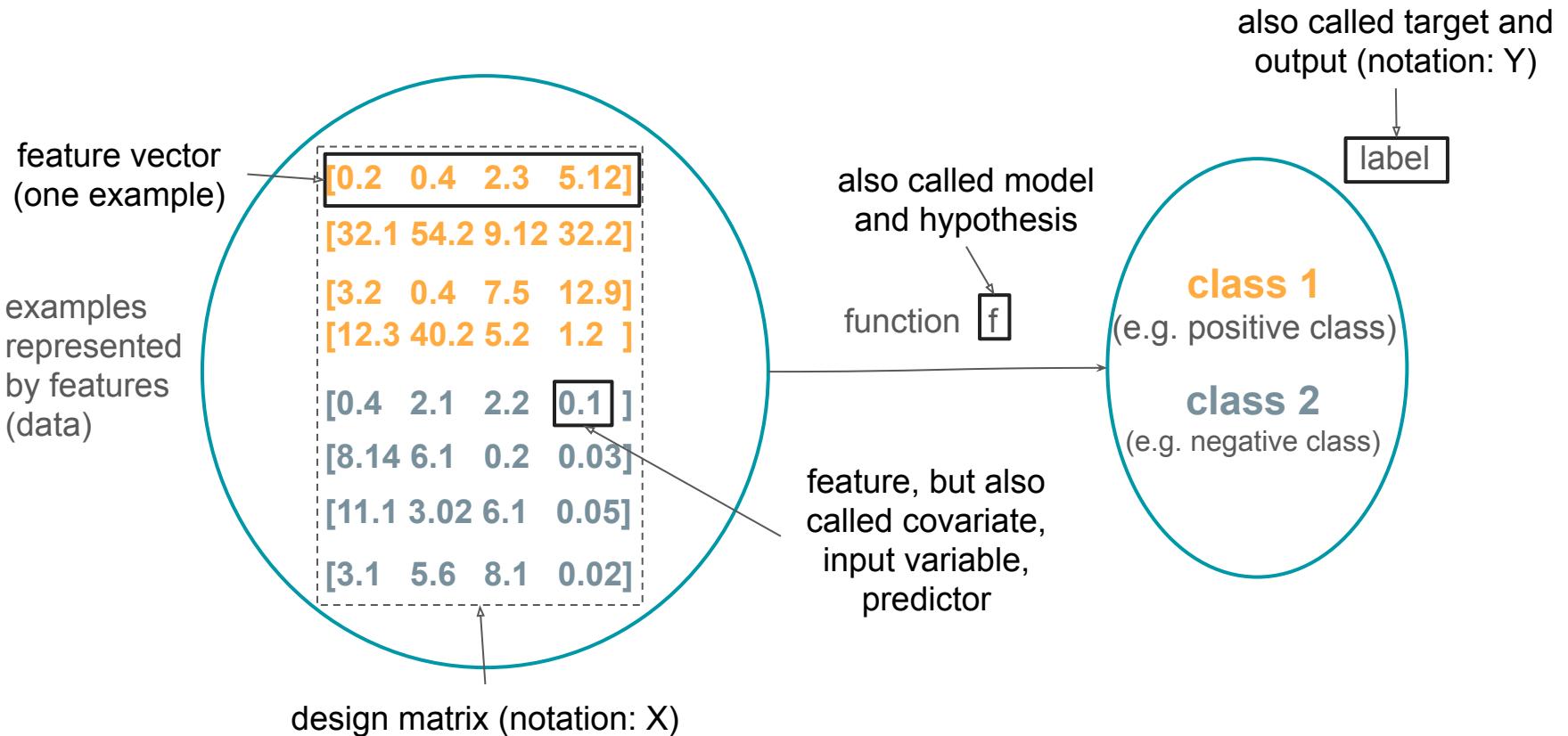
# Types of problems in machine learning: what does the function f do

1. **Supervised: classification and regression**
2. **Unsupervised**
3. **Reinforcement learning:**
  - a. dataset is not fixed, the program interacts with environment
  - b. used when choosing a sequence of actions: we don't know the label - don't know the optimal sequence of actions, but we know how good an action is
  - c. example: discover optimal dosing policy for a medication

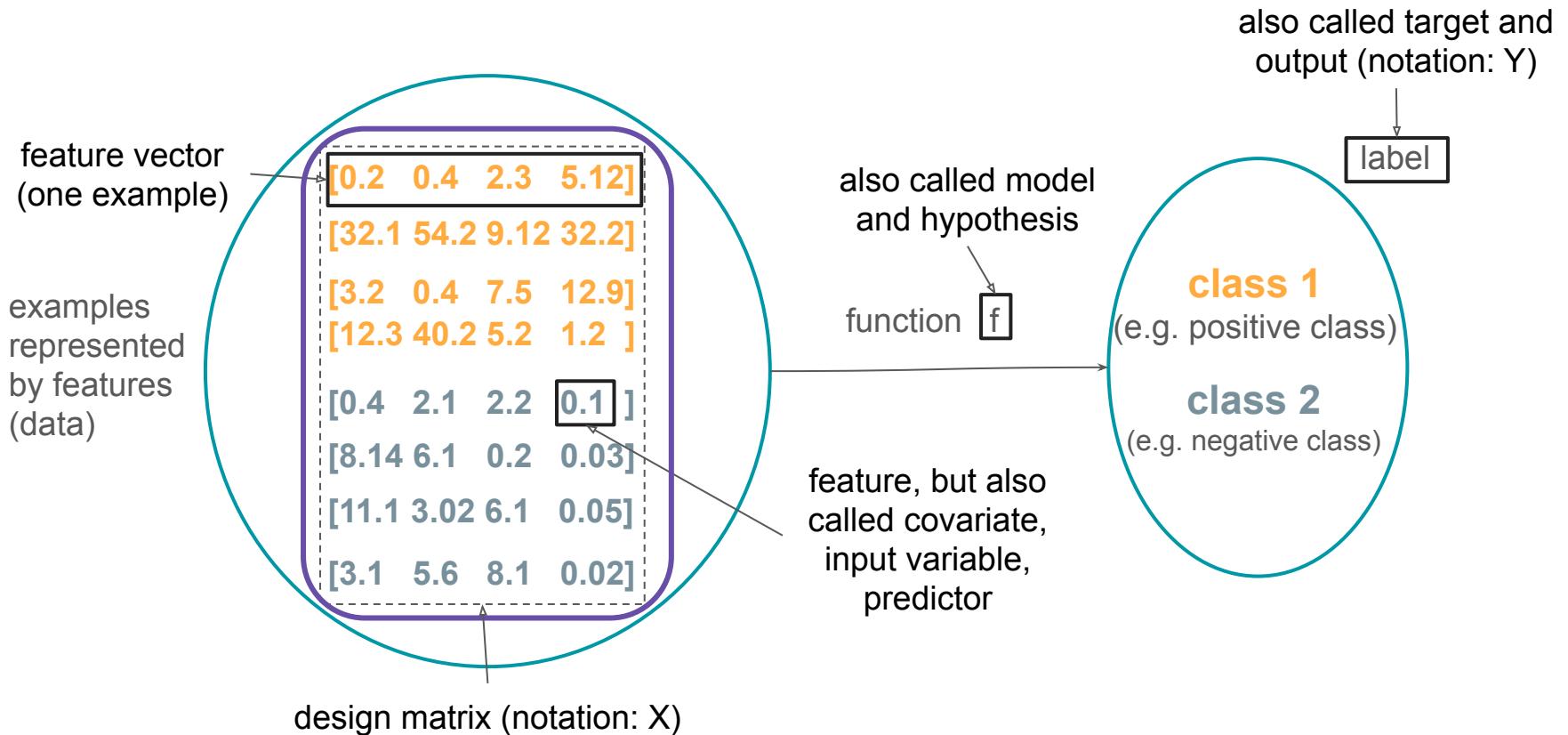
# Types of problems in machine learning: what does the function f do

1. **Supervised: classification and regression**
2. **Unsupervised**
3. **Reinforcement learning**

# Machine learning as a function approximation task



# Machine learning as a function approximation task



# What do we assume about the data?

- ❑ Data generation process produces the data  
(data generation process results in a probability distribution  $p_{data}$ )
- ❑ We assume:
  - ❑ Examples in each dataset are independent of each other
  - ❑ When we want to use the machine learning model on some new data to predict a label: these new data come from the same data generation process (same probability distribution)

# What do we assume about the data?

- ❑ Data generation process produces the data  
(data generation process results in a probability distribution  $p_{data}$ )
- ❑ We assume:
  - ❑ Examples in each dataset are independent of each other
  - ❑ When we want to use the machine learning model on some new data to predict a label: these new data come from the same data generation process (same probability distribution)

# What do we assume about the data?

- ❑ Data generation process produces the data  
(data generation process results in a probability distribution  $p_{data}$ )
- ❑ We assume:
  - ❑ Examples in each dataset are independent of each other
  - ❑ When we want to use the machine learning model on some new data to predict a label: these new data come from the same data generation process (same probability distribution)

i.i.d.  
assumption



examples are  
independent  
and identically  
distributed

# What do we assume about the data?

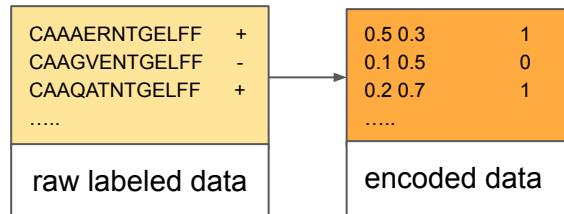
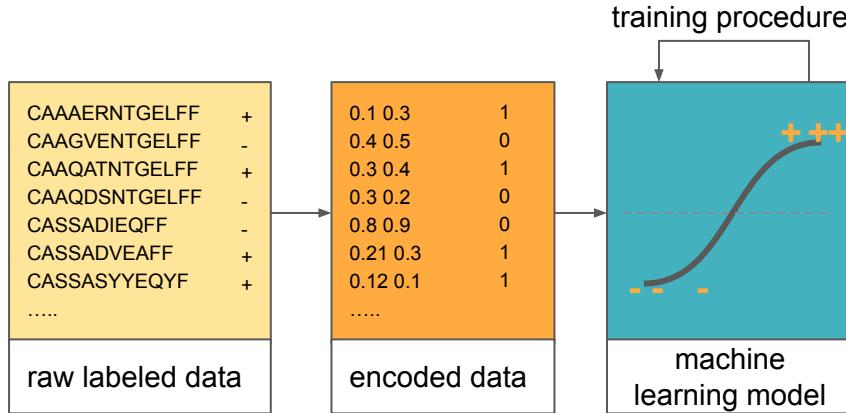
- ❑ Data generation process produces the data  
(data generation process results in a probability distribution  $p_{data}$ )
- ❑ We assume:
  - ❑ Examples in each dataset are independent of each other
  - ❑ When we want to use the machine learning model on some new data to predict a label: these new data come from the same data generation process (same probability distribution)
- ❑ With these assumptions satisfied (or approximately satisfied), we choose the data representation and estimate the function

i.i.d.  
assumption

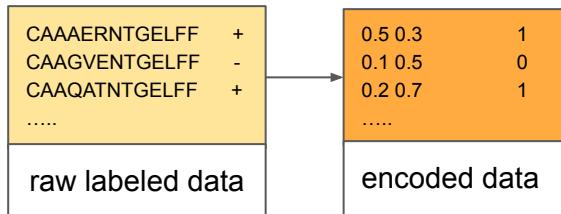
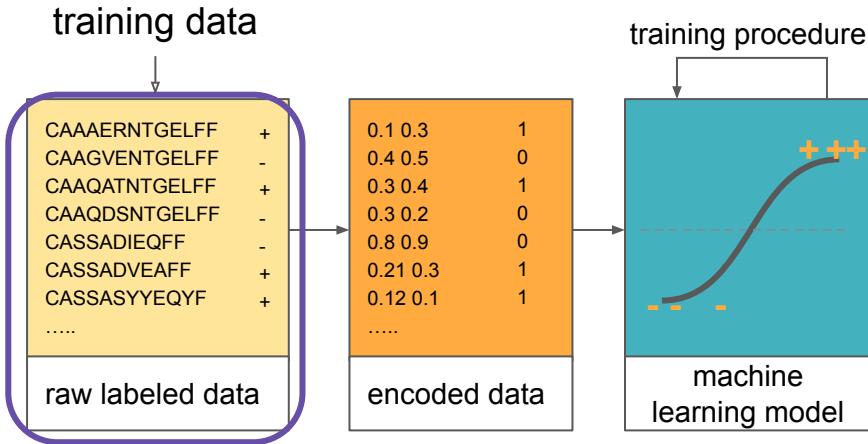
↓

examples are  
independent  
and identically  
distributed

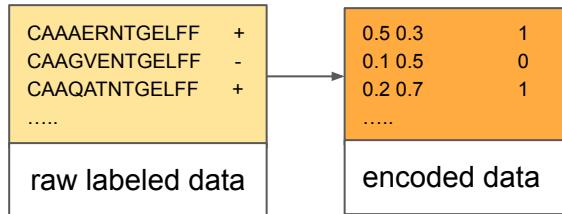
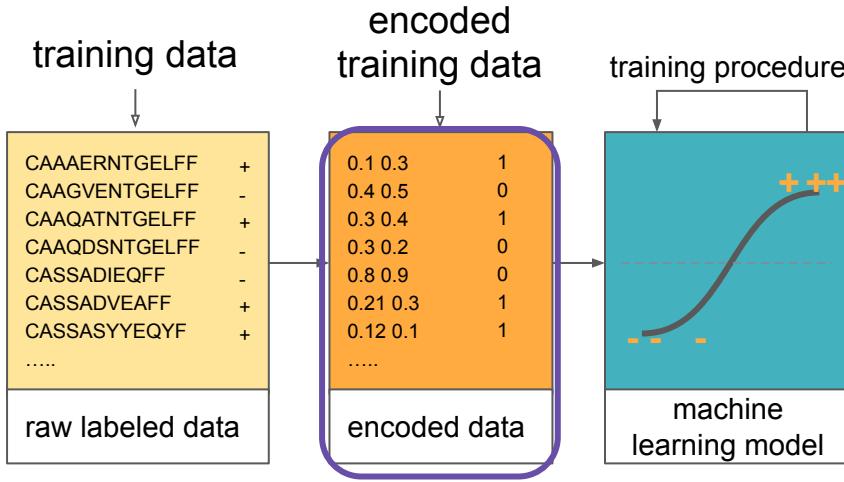
# Estimating the function (training procedure)



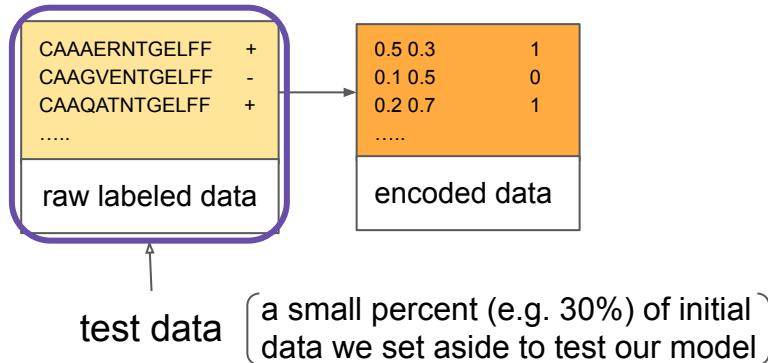
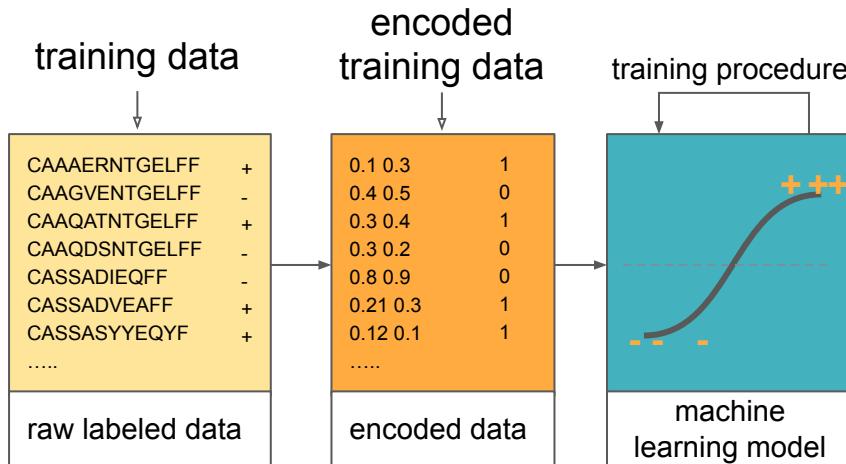
# Estimating the function (training procedure)



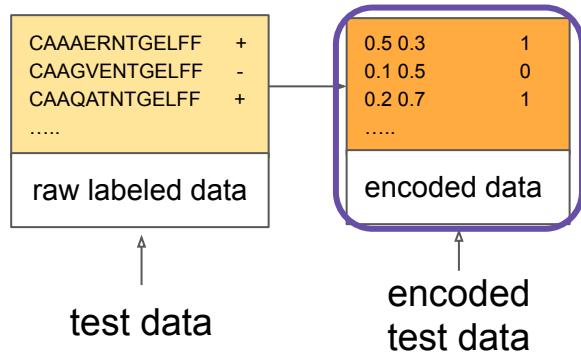
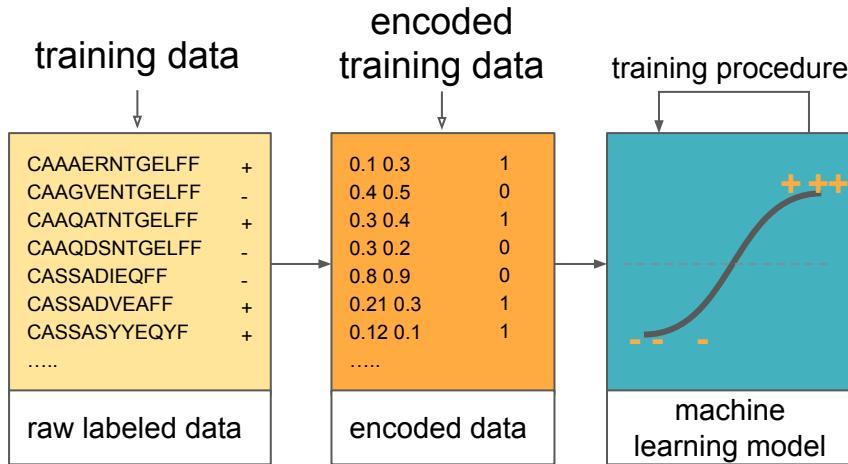
# Estimating the function (training procedure)



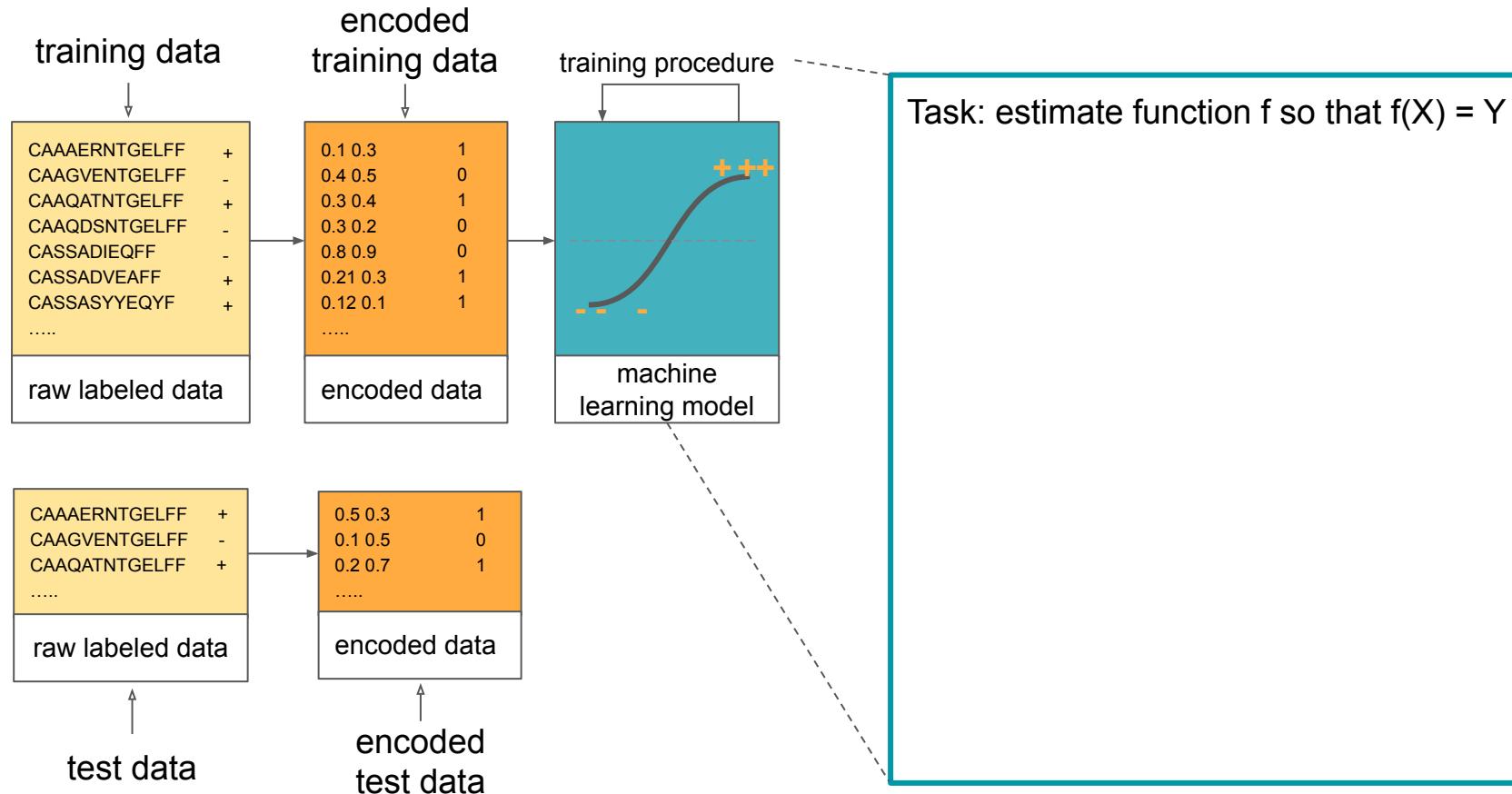
# Estimating the function (training procedure)



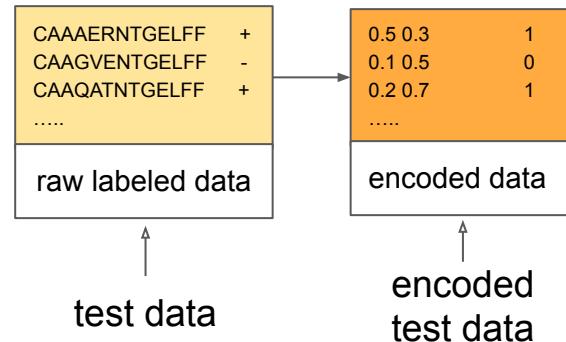
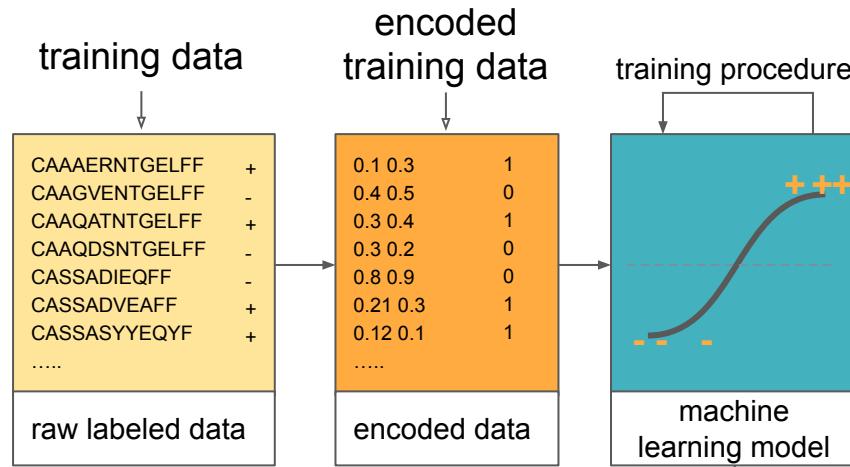
# Estimating the function (training procedure)



# Estimating the function (training procedure)



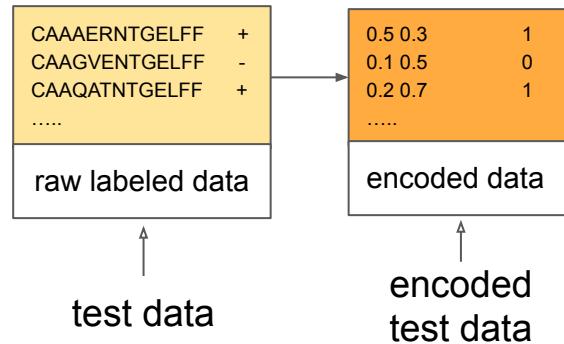
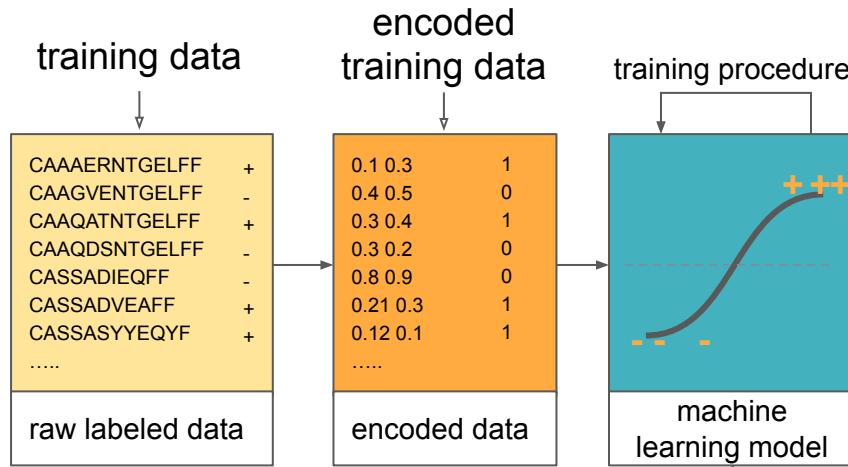
# Estimating the function (training procedure)



Task: estimate function  $f$  so that  $f(X) = Y$

Training procedure:

# Estimating the function (training procedure)

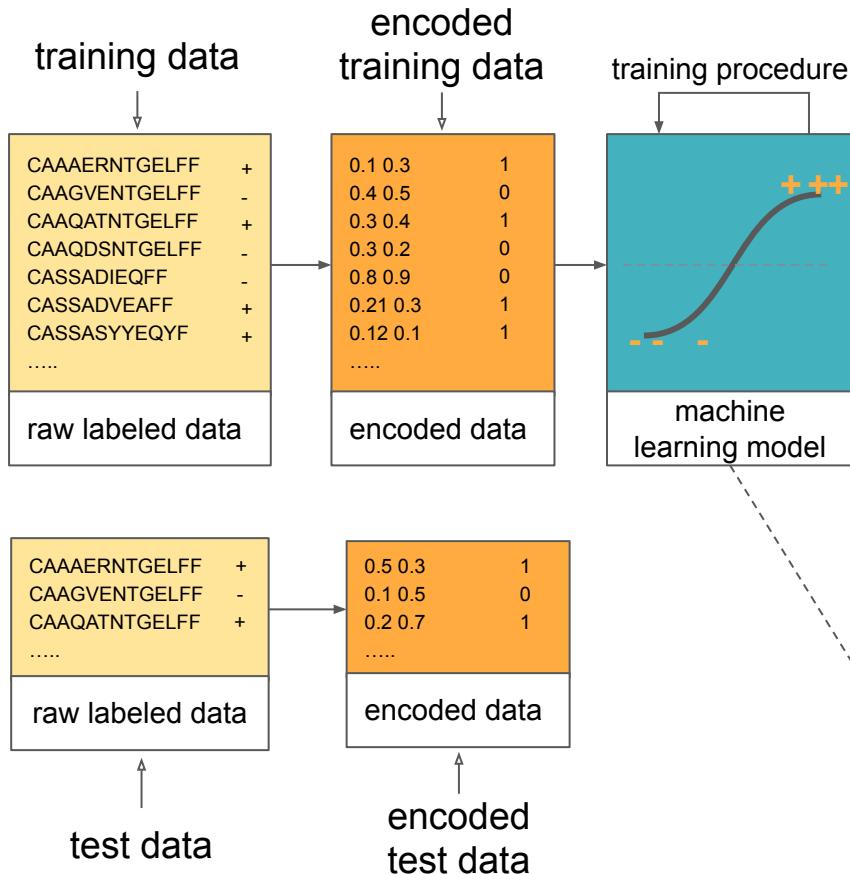


Task: estimate function  $f$  so that  $f(X) = Y$

Training procedure:

1. Start with some function  $f$  with some parameters

# Estimating the function (training procedure)



Task: estimate function  $f$  so that  $f(X) = Y$

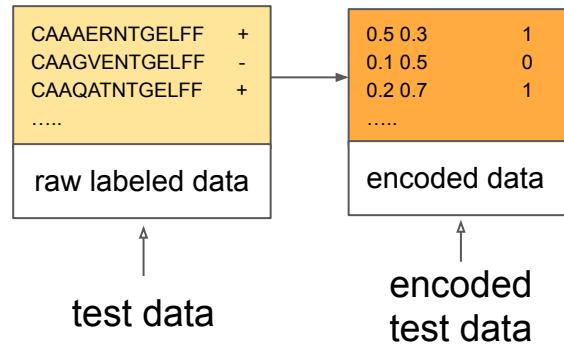
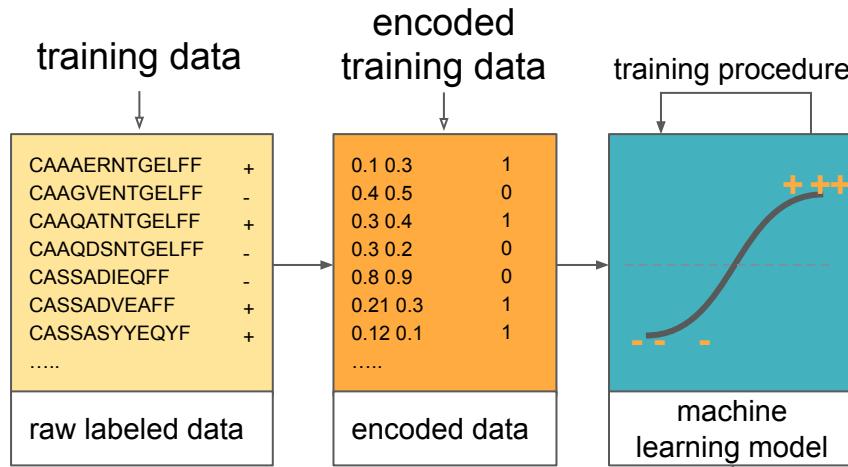
Training procedure:

1. Start with some function  $f$  with some parameters  
for example, logistic regression:

$$g(\omega x + b) = (1 + e^{-(\omega x + b)})^{-1}$$

$$f(x) = \begin{cases} 1, & g(\omega x + b) \geq 0.5 \\ 0, & g(\omega x + b) < 0.5 \end{cases}$$

# Estimating the function (training procedure)

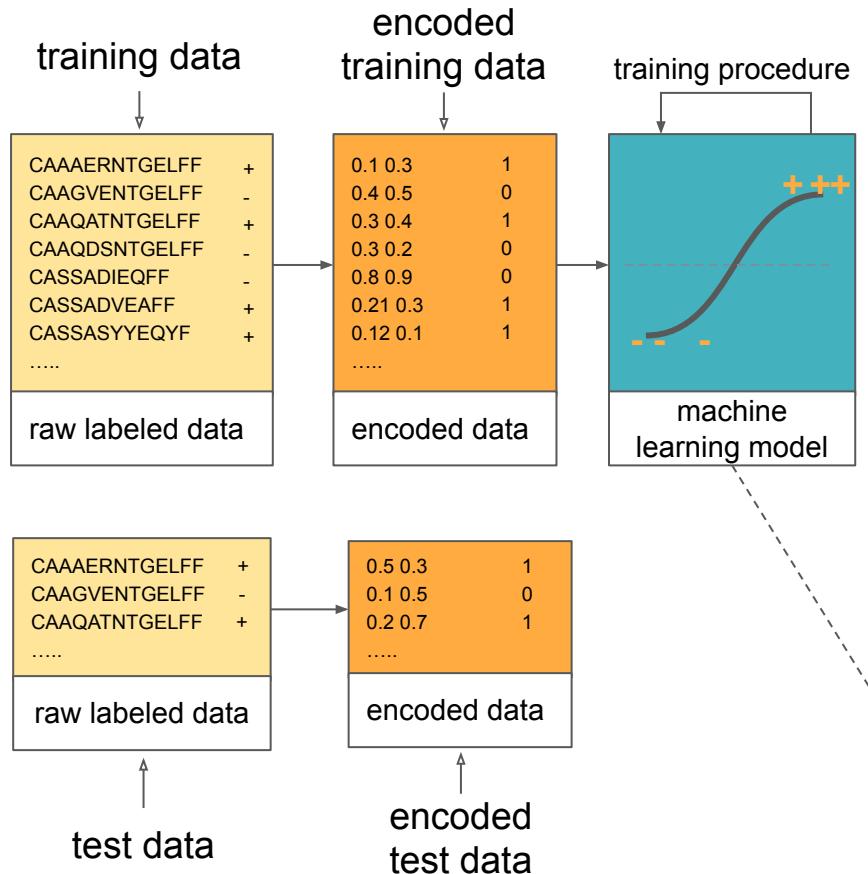


Task: estimate function  $f$  so that  $f(X) = Y$

Training procedure:

1. Start with some function  $f$  with some parameters (e.g., logistic regression)

# Estimating the function (training procedure)

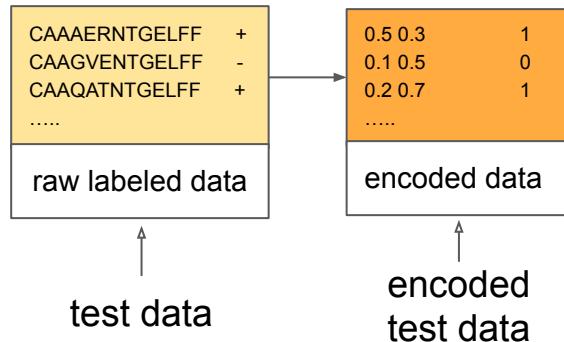
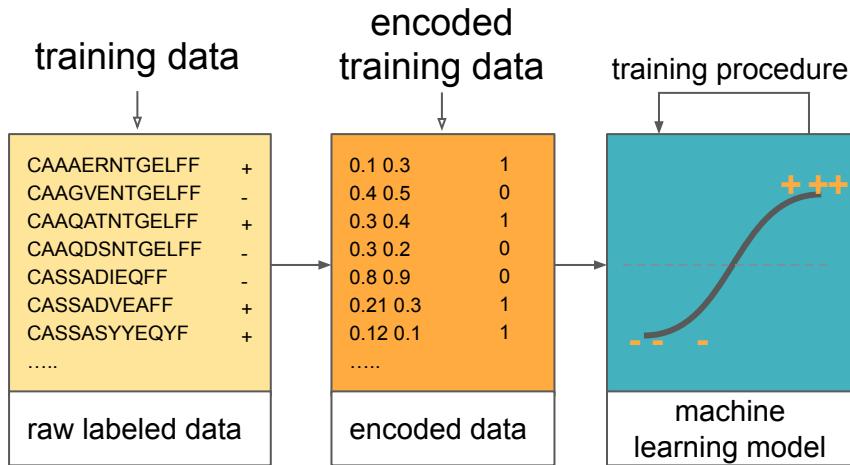


Task: estimate function  $f$  so that  $f(X) = Y$

Training procedure:

1. Start with some function  $f$  with some parameters (e.g., logistic regression)
2. While training:

# Estimating the function (training procedure)



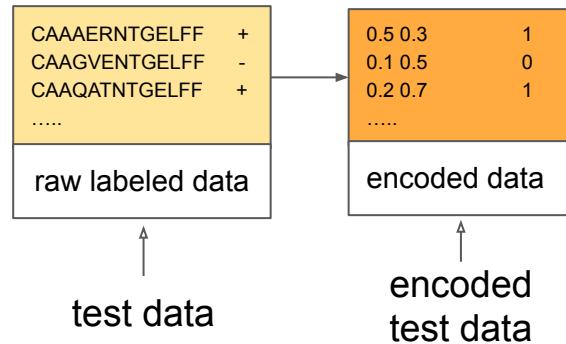
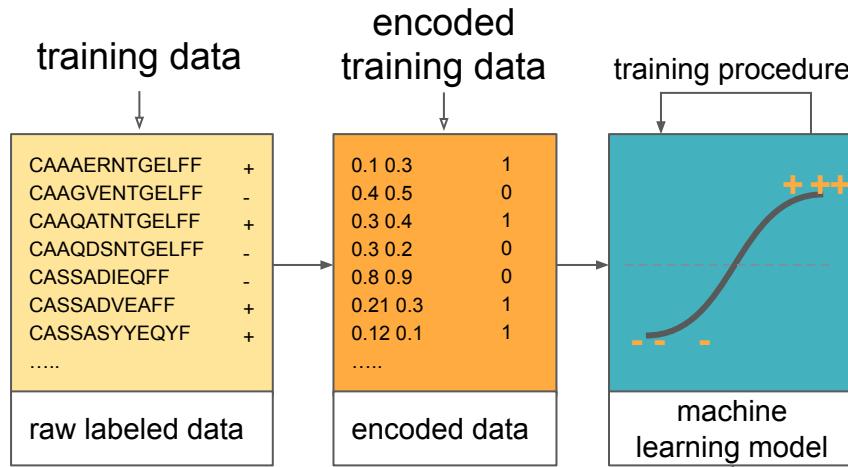
Task: estimate function  $f$  so that  $f(X) = Y$

Training procedure:

1. Start with some function  $f$  with some parameters (e.g., logistic regression)
2. While training:

max number of iterations was not reached and predictions are not good enough

# Estimating the function (training procedure)

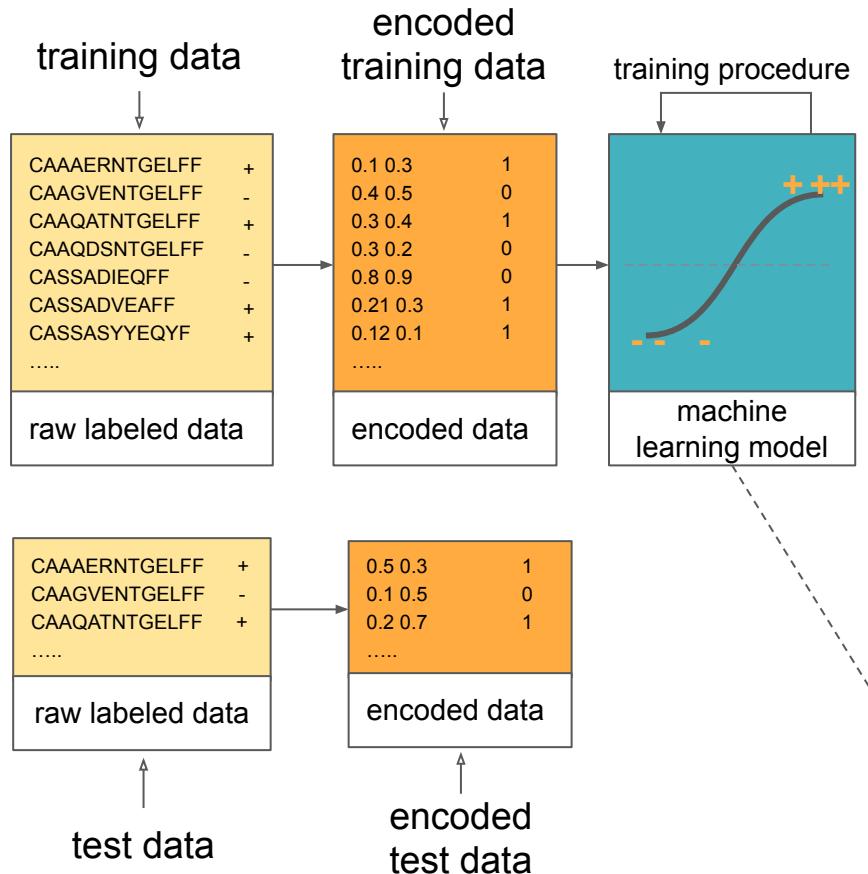


Task: estimate function  $f$  so that  $f(X) = Y$

Training procedure:

1. Start with some function  $f$  with some parameters (e.g., logistic regression)
2. While training:

# Estimating the function (training procedure)

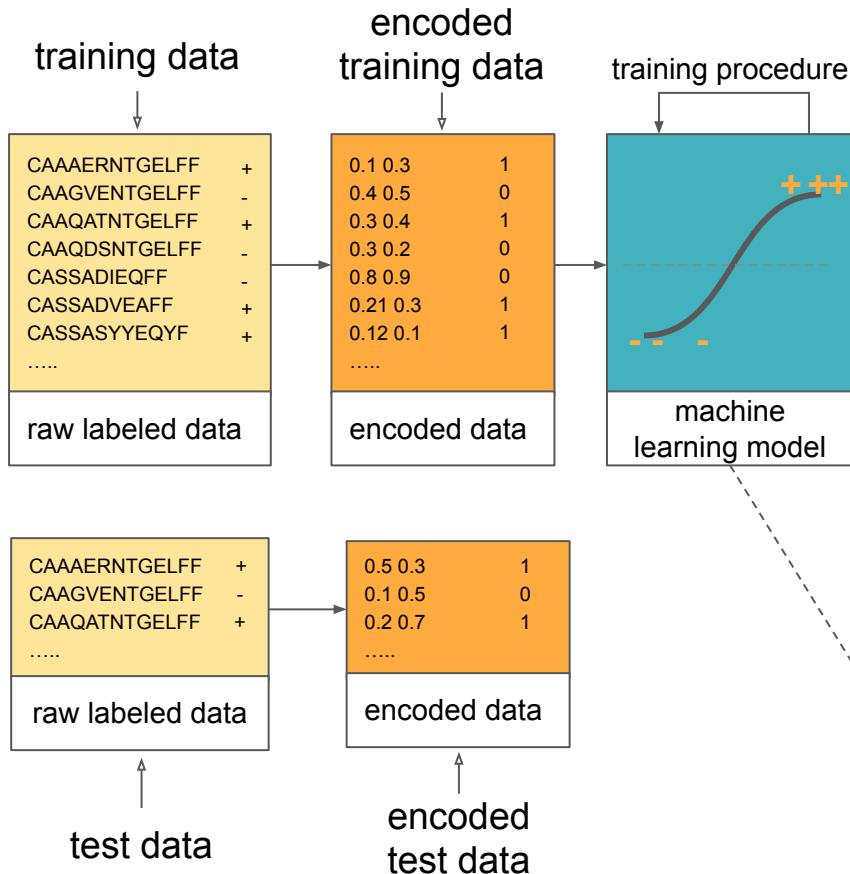


Task: estimate function  $f$  so that  $f(X) = Y$

Training procedure:

1. Start with some function  $f$  with some parameters (e.g., logistic regression)
2. While training:
  - a. Predict the label  $Y$  from the encoded training data  $X$

# Estimating the function (training procedure)

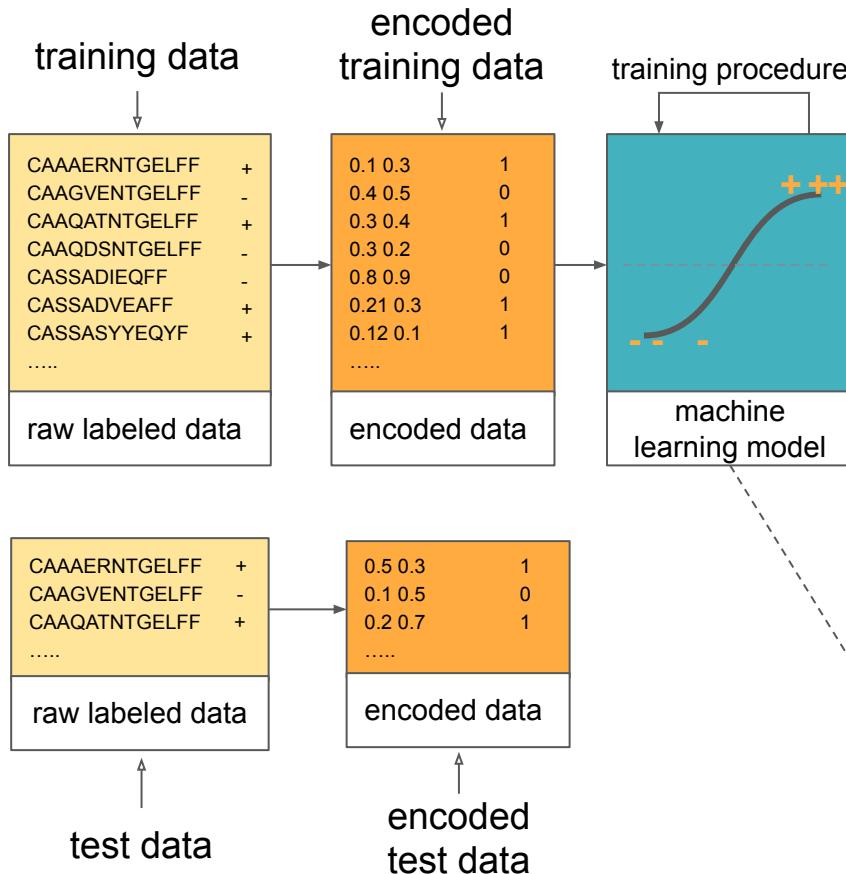


Task: estimate function  $f$  so that  $f(X) = Y$

Training procedure:

1. Start with some function  $f$  with some parameters (e.g., logistic regression)
2. While training:
  - a. Predict the label  $Y$  from the encoded training data  $X$
  - b. Compute the cost function: how much predictions deviate from the label  $Y$

# Estimating the function (training procedure)

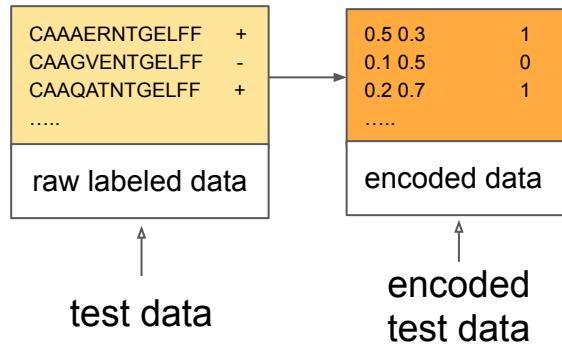
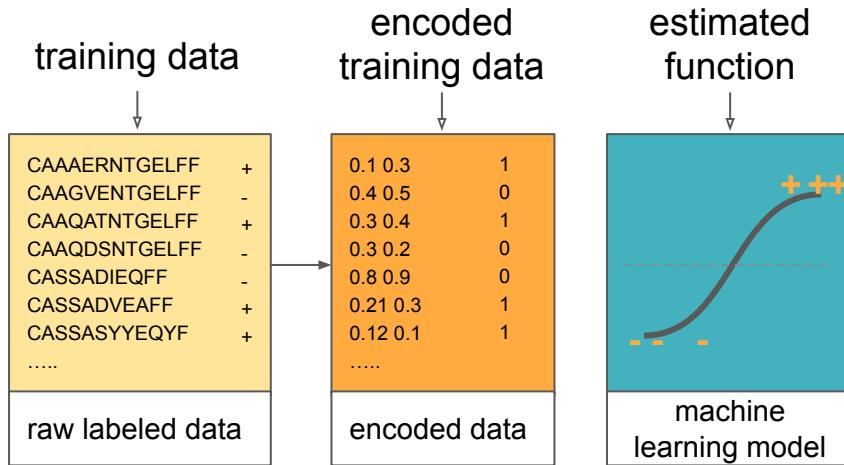


Task: estimate function  $f$  so that  $f(X) = Y$

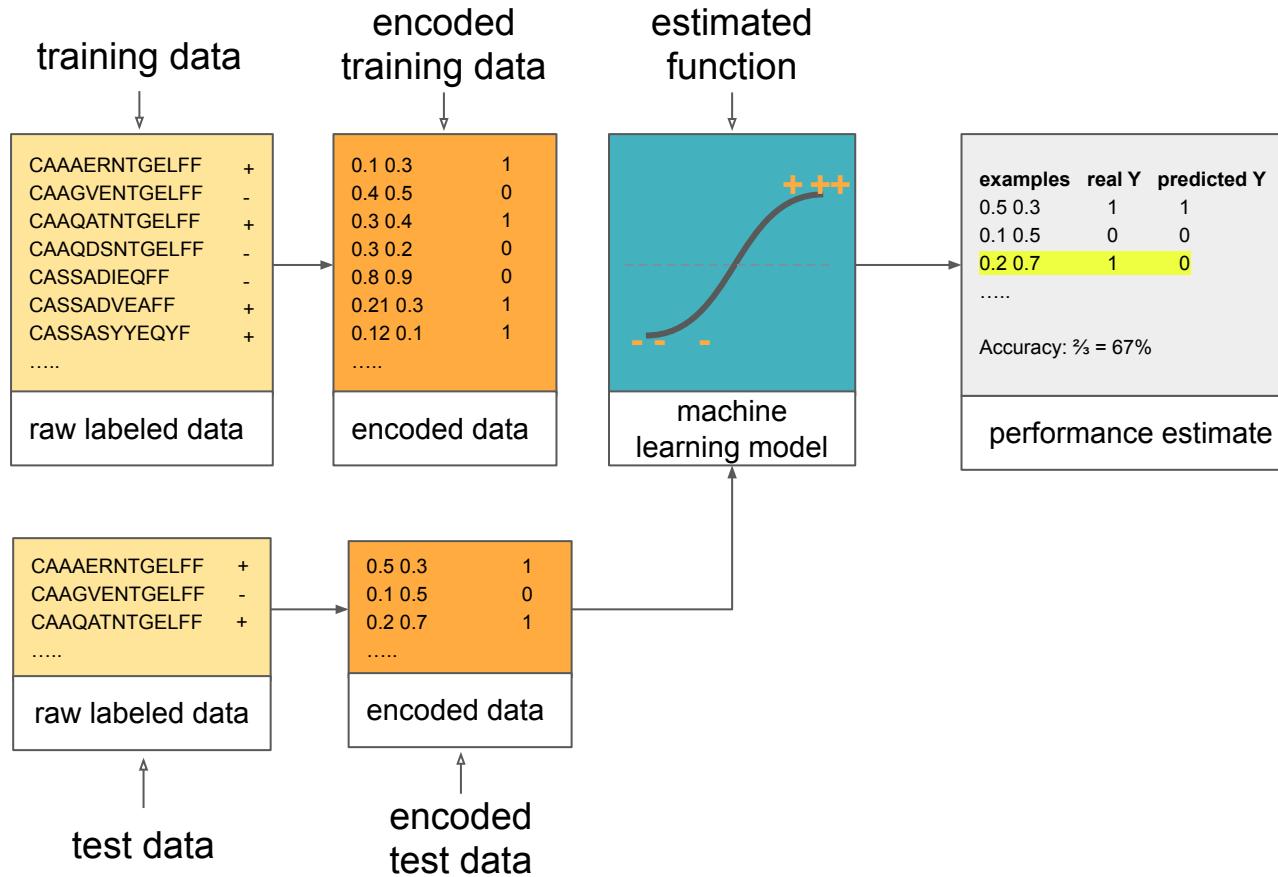
Training procedure:

1. Start with some function  $f$  with some parameters (e.g., logistic regression)
2. While training:
  - a. Predict the label  $Y$  from the encoded training data  $X$
  - b. Compute the cost function: how much predictions deviate from the label  $Y$
  - c. Update the parameters of the function  $f$  to reduce the cost function so that we get better predictions

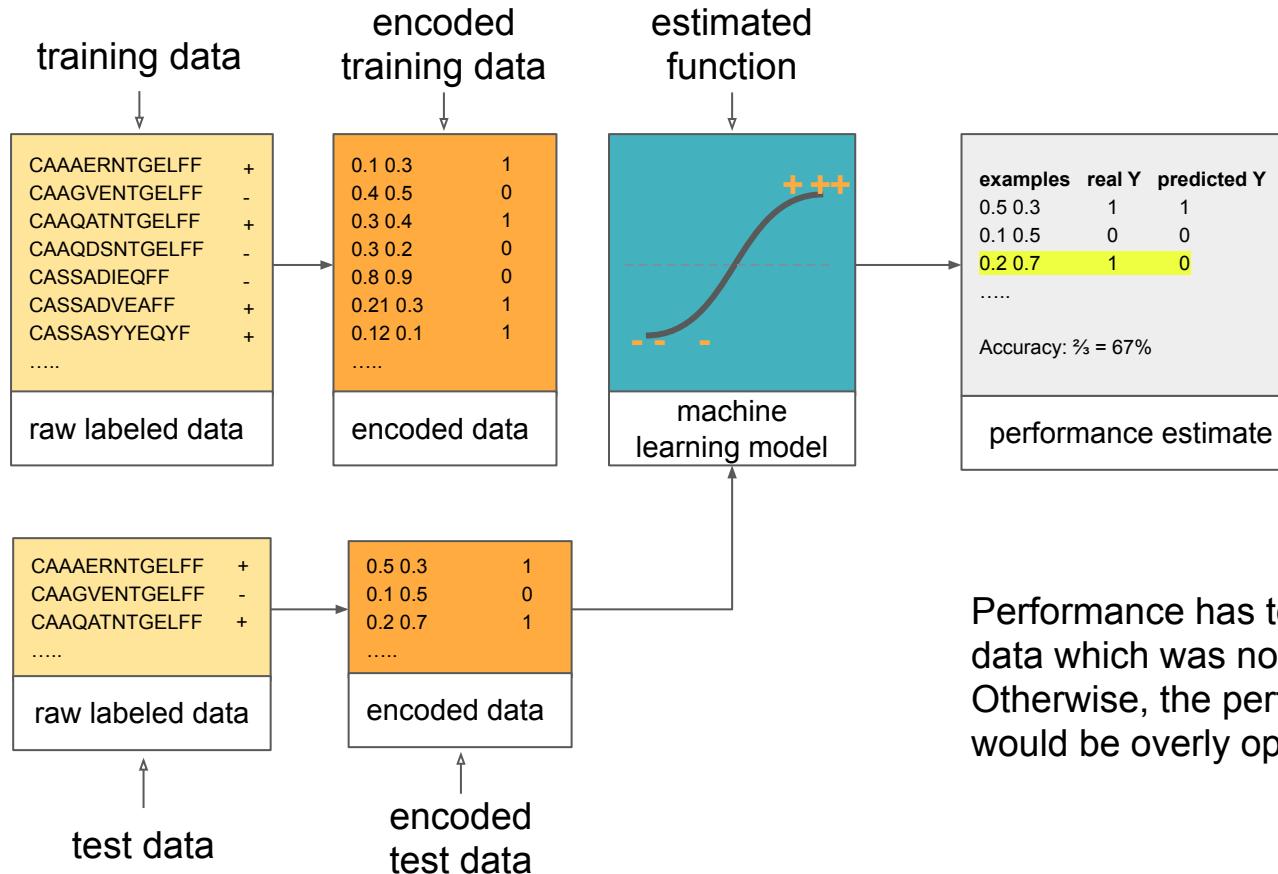
# Estimating the function (training procedure)



# Estimating the function (training procedure)



# Estimating the function (training procedure)



Performance has to be estimated on data which was not used during training. Otherwise, the performance estimate would be overly optimistic.

# Performance metrics - classification

- ❑ Depends on the problem and the data
  - ❑ Classification (label values come from a discrete set):

for binary classification, this equality holds:

$$\text{accuracy} = \frac{\text{number of correct predictions}}{\text{total number of predictions}} \downarrow = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

true positives                    false positives                    false negatives  
true negatives

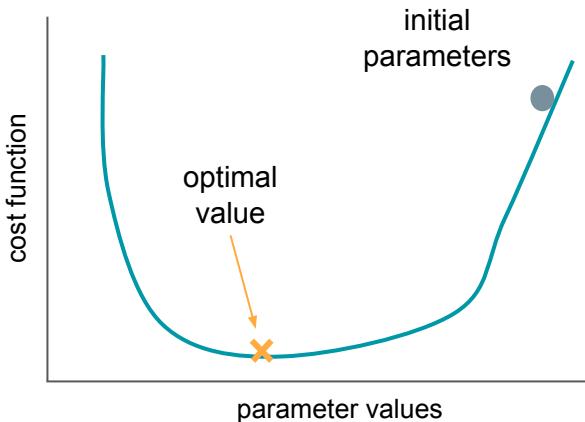
Other metrics: balanced accuracy, precision, recall, sensitivity, specificity, ROC curve, AUC, cross-entropy

# Training the machine learning model

- We want to minimize the cost function
- For instance, we can use optimization algorithm called gradient descent:

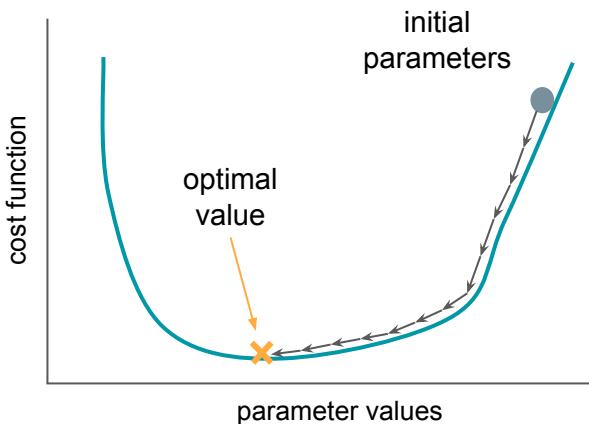
Repeat until optimal solution / max number of iterations:

1. Find derivative of the cost function w.r.t. each of the parameters of the model
2. Update each parameter incrementally using the cost function as a starting point for the update computation



# Training the machine learning model

- We want to minimize the cost function
- For instance, we can use optimization algorithm called gradient descent

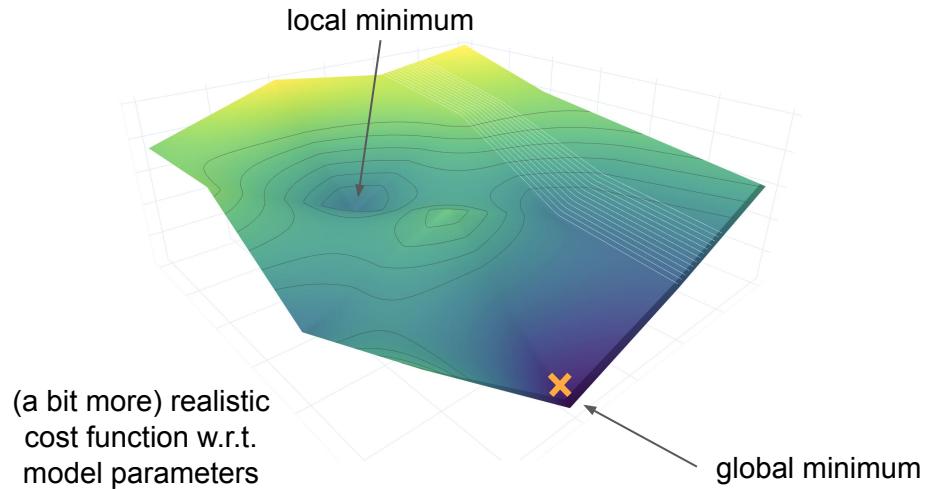
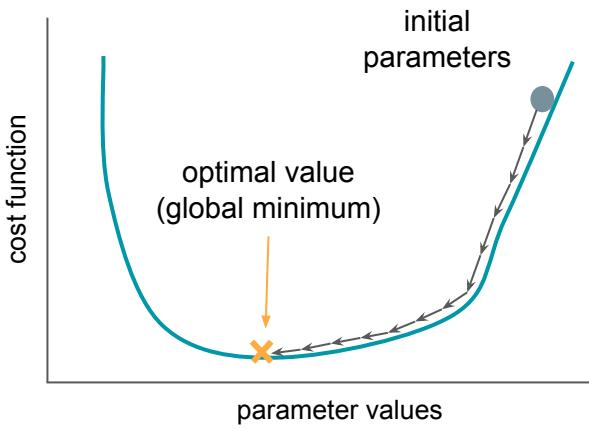


Repeat until optimal solution / max number of iterations:

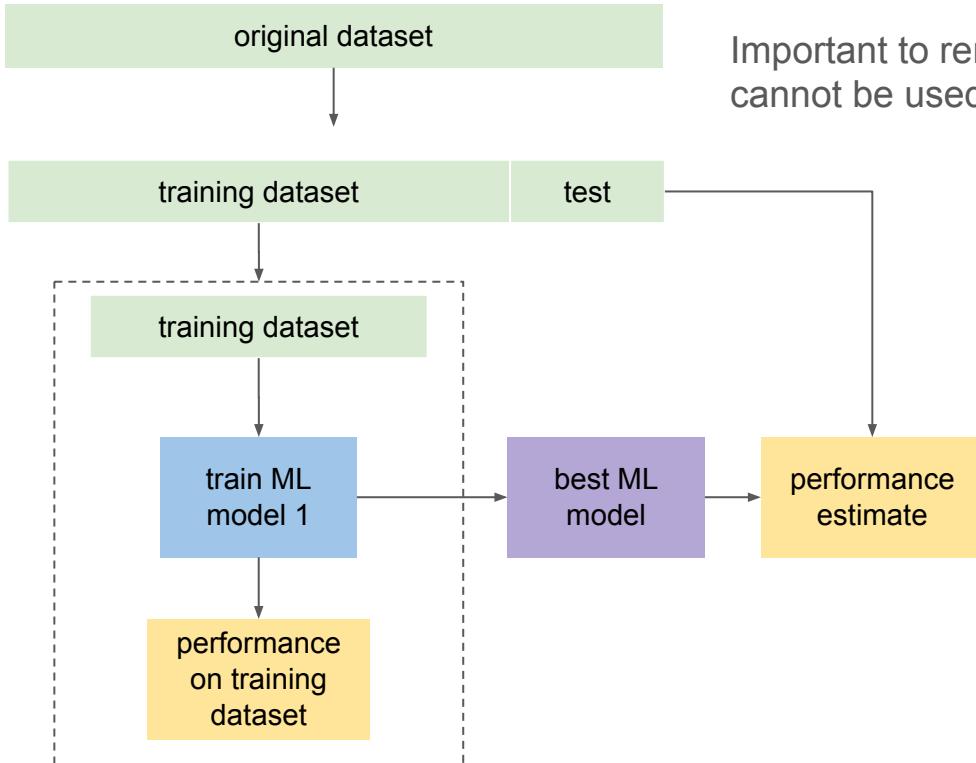
1. Find derivative of the cost function w.r.t. each of the parameters of the model
2. Update each parameter incrementally using the cost function as a starting point for the update computation

# Training the machine learning model

- We want to minimize the cost function
- For instance, we can use optimization algorithm called gradient descent



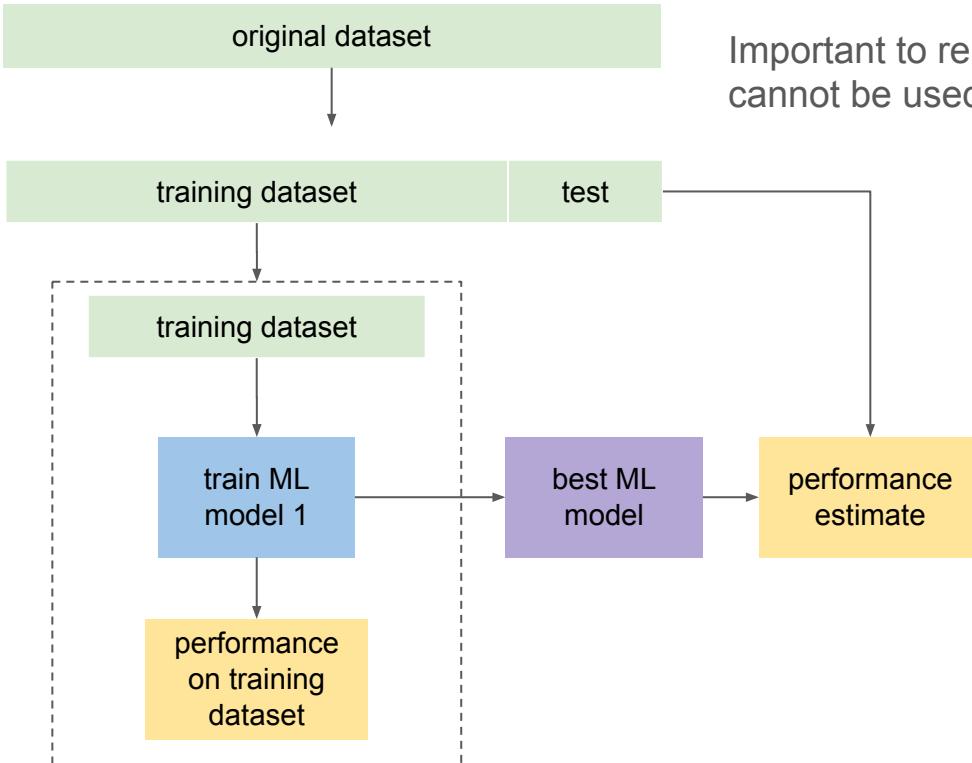
# Machine learning workflow



One way to set up a machine learning workflow

Important to remember: data used to assess the performance cannot be used during training

# Machine learning workflow



One way to set up a machine learning workflow

Important to remember: data used to assess the performance cannot be used during training

Performance on the test data (not seen during training) will typically be worse than performance on validation

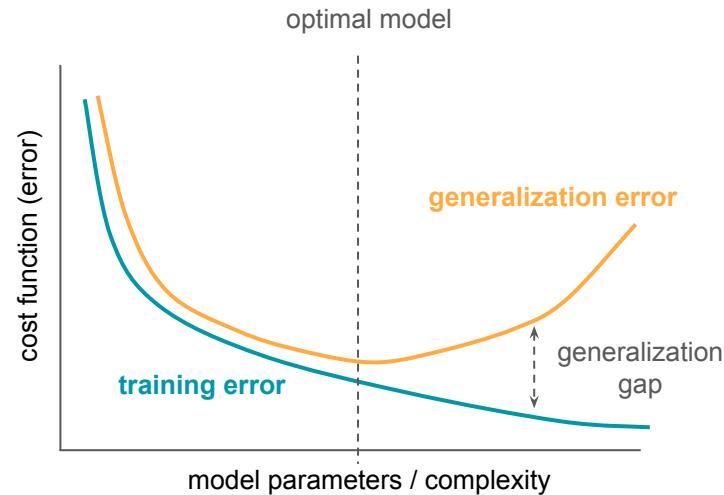
And we will come back to this later...

# Generalization in ML

- ❑ Generalization is the ability of an ML model to perform well on previously unseen data
- ❑ We use error on the test set as an estimate of generalization error
- ❑ Generalization error is the expected error on new data

We want a model which will have:

- ❑ Small error on the training set
- ❑ Small gap between training set error and test set error



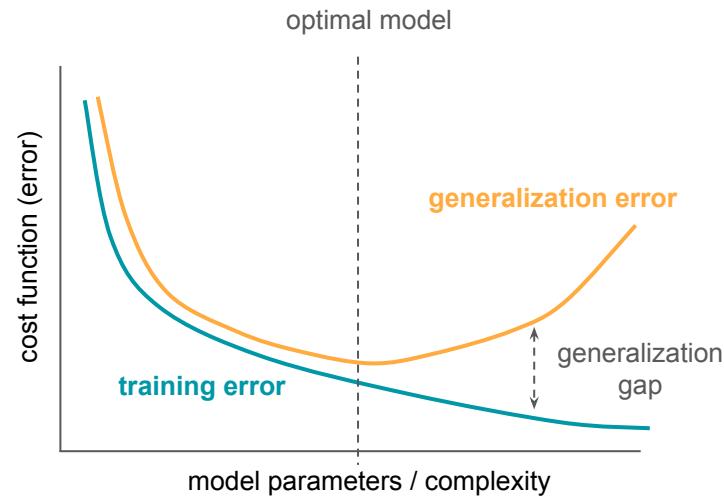
# Generalization in ML

- ❑ Generalization is the ability of an ML model to perform well on previously unseen data
- ❑ We use error on the test set as an estimate of generalization error
- ❑ Generalization error is the expected error on new data

We want a model which will have:

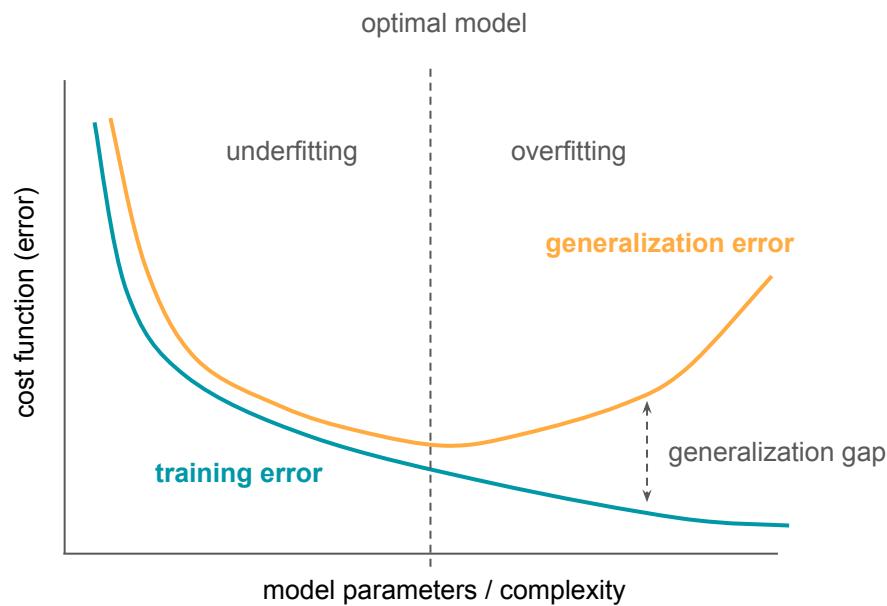
- ❑ Small error on the training set
- ❑ Small gap between training set error and test set error

Remember that we can talk about generalization like this only if the i.i.d. assumption at least approximately holds.



# Overfitting and underfitting

- ❑ **Underfitting:** the model was not able to learn from the training data - it had high training error
- ❑ **Overfitting:** the generalization gap is too large because the model fit the training data too well but failed to extract patterns which would enable good performance on the new (test) data



# References

Goodfellow IJ, Bengio Y, Courville A. *Deep Learning*. MIT Press; 2016. <https://mitpress.mit.edu/books/deep-learning>

Mitchell T. *Machine Learning*. McGraw Hill; 1997. <http://www.cs.cmu.edu/~tom/mlbook.html>

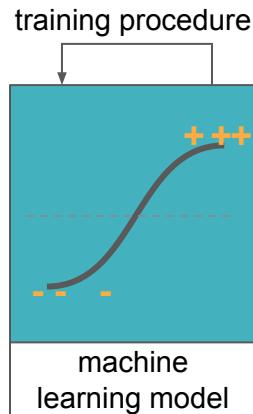
**Murphy K.** *Probabilistic Machine Learning: An Introduction*. MIT Press; 2022. <https://probml.github.io/pml-book/book1.html>  
[Chapter 1]

# Machine learning in computational biology - outline

- Introduction to machine learning:
  - What is machine learning, types of problems, assumptions, workflow, generalization
- **Machine learning models and algorithms:**
  - Discriminative vs generative models, supervised models (logistic and linear regression, kNN, neural networks), unsupervised models (dimensionality reduction, clustering)
- Data representation:
  - Considerations and examples, one-hot encoding, feature engineering, representation learning
- Model comparison and uncertainty:
  - Model assessment, model selection, uncertainty, cross-validation
- Transparency and reproducibility

# ML models

- We mentioned logistic regression before - a simple model for binary classification



Task: estimate function  $f$  so that  $f(X) = Y$

Training procedure:

- Start with some function  $f$  with some parameters  
for example, logistic regression:

$$g(\omega x + b) = (1 + e^{-(\omega x + b)})^{-1}$$

$$f(x) = \begin{cases} 1, & g(\omega x + b) \geq 0.5 \\ 0, & g(\omega x + b) < 0.5 \end{cases}$$

# Some terminology regarding ML models and algorithms

- ❑ **Learning algorithm:** a function that, given a set of examples and their labels, constructs a model, e.g., logistic regression
- ❑ **Model:** a function which was fit to the data using the learning algorithm, e.g., logistic regression with specific coefficients

Dietterich 1998

Usually model and learning algorithm are used interchangeably but they mean slightly different things

# Capacity of the model

- ❑ A model's capacity is its ability to fit a wide variety of functions, for instance:  
linear regression:

$$\hat{y} = b + \omega x$$

polynomial regression:

$$\hat{y} = b + \omega_1 x + \omega_2 x^2$$

# Capacity of the model

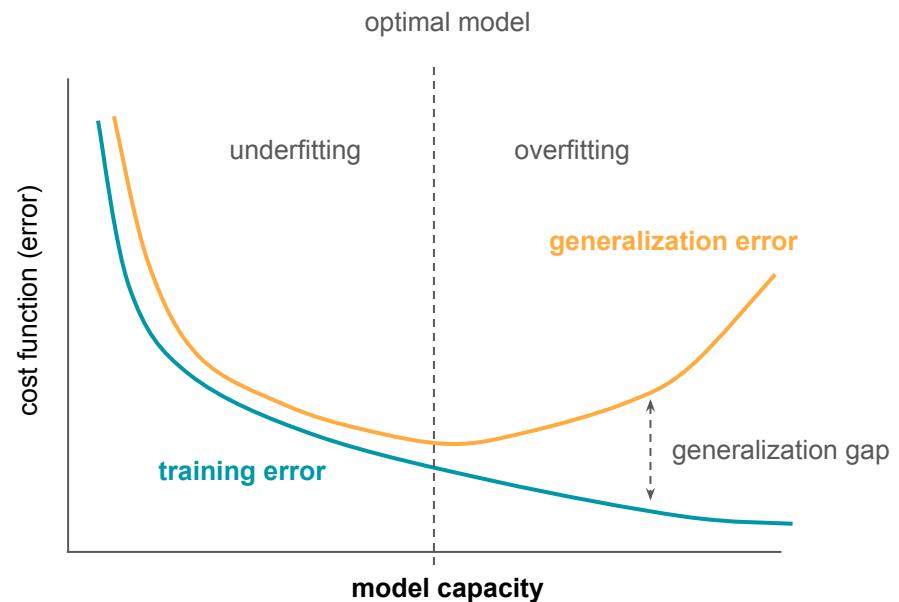
- ❑ A model's capacity is its ability to fit a wide variety of functions, for instance:

linear regression:

$$\hat{y} = b + \omega x$$

polynomial regression:

$$\hat{y} = b + \omega_1 x + \omega_2 x^2$$



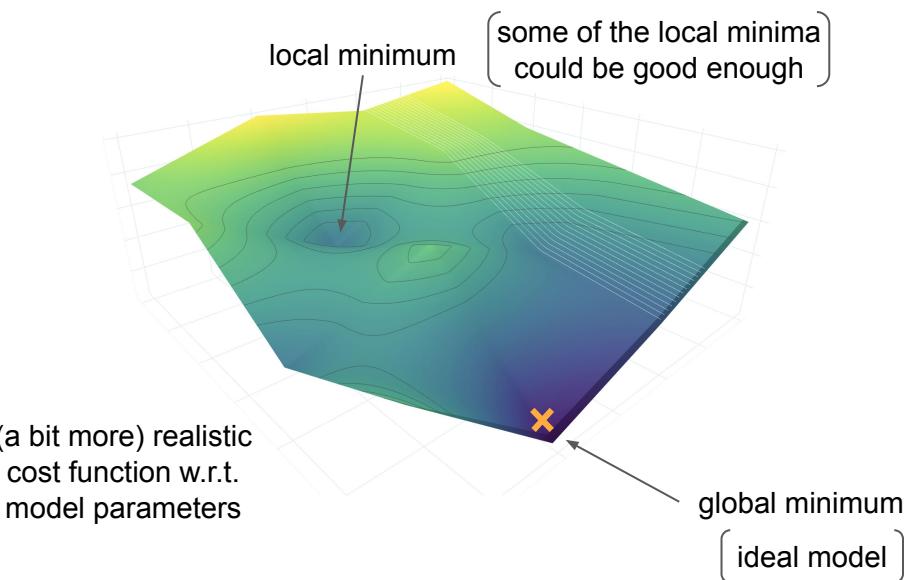
# Searching through a hypothesis space using optimization algorithms

Fitting the parameters of the model is an optimization problem (there are different optimization algorithms, but one example is **gradient descent**)

We can get stuck in local minima

If the performance is good enough, we can accept that “suboptimal” model

Performance is measured by cost function [a lot of engineering can happen there]



# Regularization restricts what models can be learned

Regularization – restricting the values of parameters of the model that can be learned (e.g., based on our domain knowledge) so that the model performs better on new data

Typical forms of regularization (also called penalty):

L1 (lasso):

$$\text{regularization (parameters)} = \sum_i |\text{parameter}_i|$$

L2 (ridge):

$$\text{regularization (parameters)} = \sum_i \text{parameter}_i^2$$

# Logistic Regression

## □ Binary classification:

function computing log-odds for the positive class

$$g(\omega x + b) = (1 + e^{-(\omega x + b)})^{-1}$$

model parameters (coefficients)

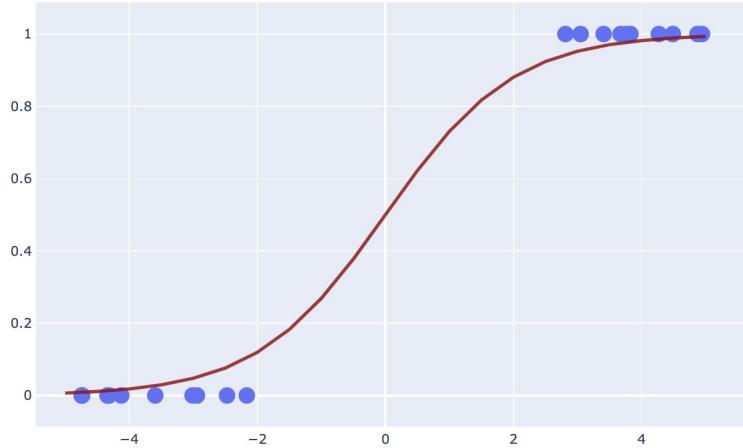
model parameter (bias or intercept)

linear combination of features

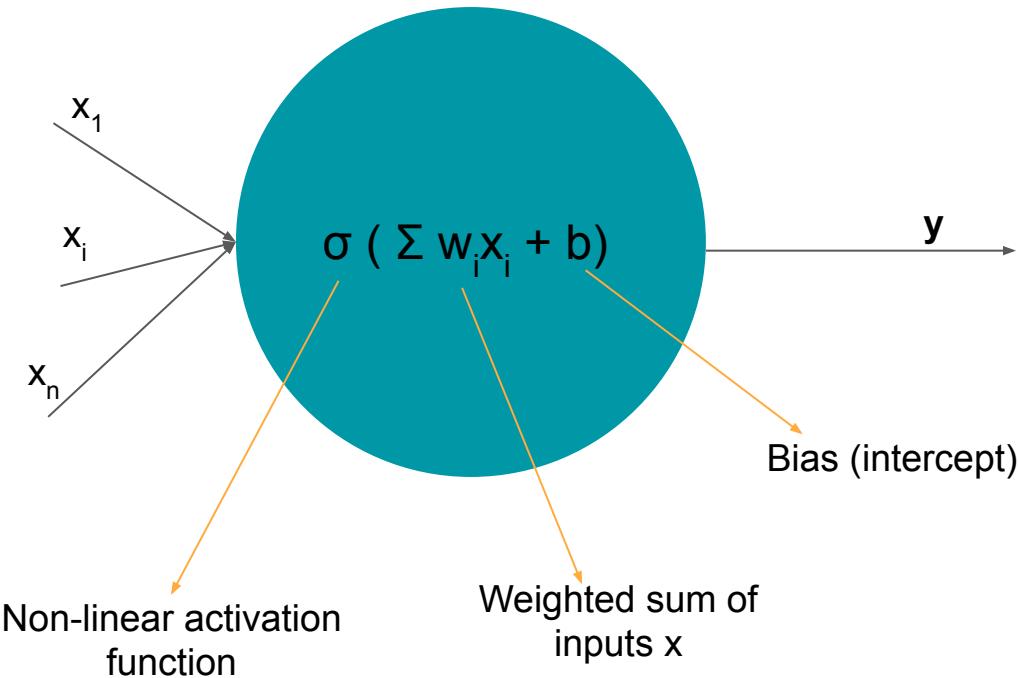
data (design matrix or feature vector)

function making the class prediction based on the threshold from log-odds value

$$f(x) = \begin{cases} 1, & g(\omega x + b) \geq 0.5 \\ 0, & g(\omega x + b) < 0.5 \end{cases}$$

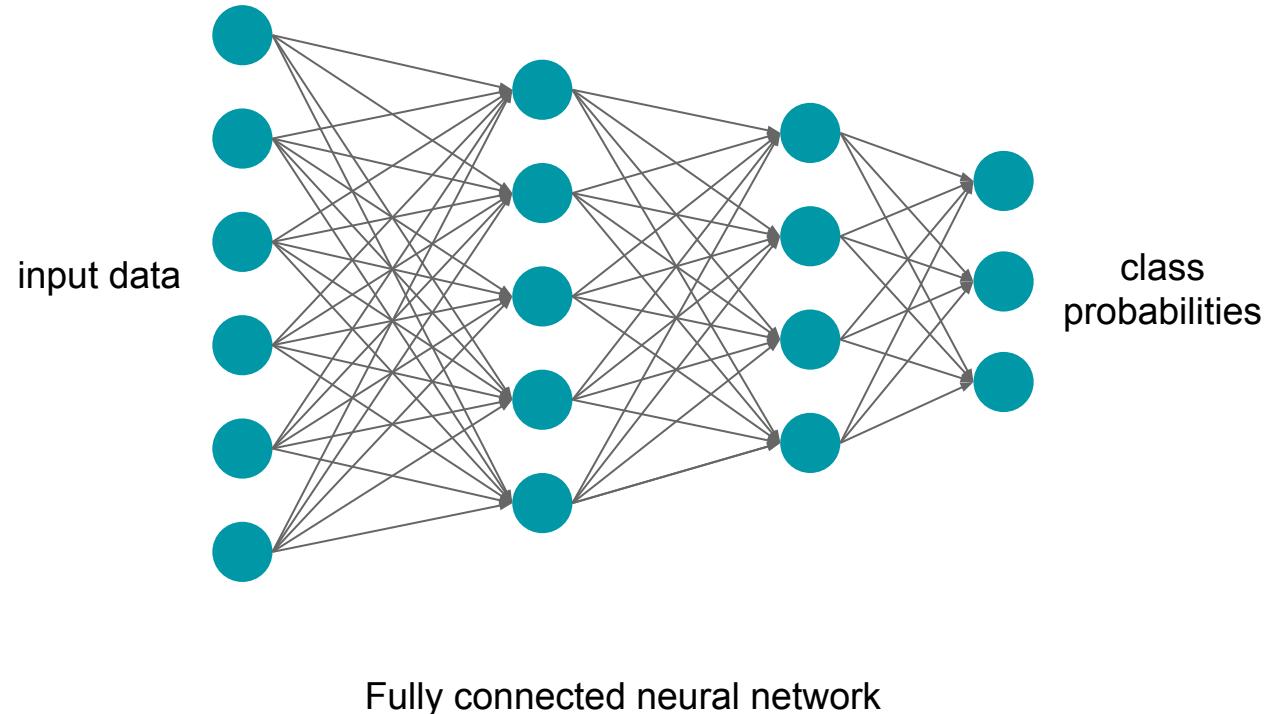


# Single nodes in the neural network do something similar



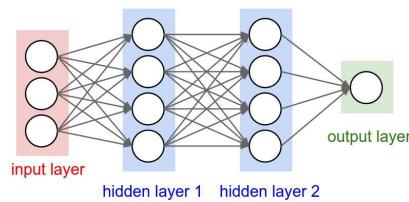
# Neural Networks

- ❑ Nodes in neural networks are organized into layers
- ❑ Number of nodes in the layer and number of layers are hyperparameters (not optimized during training, but instead set manually)
- ❑ Hierarchical structure makes them very powerful



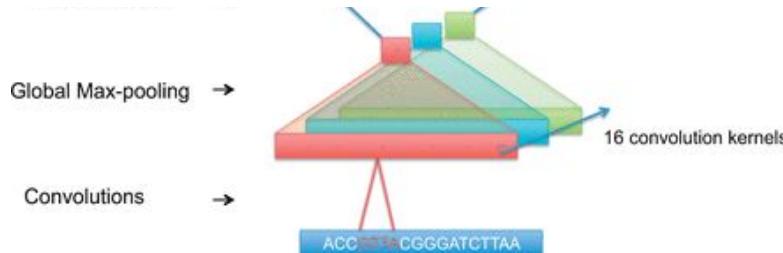
# Types of layers in neural networks

- ❑ **Fully connected networks:**  
can approximate almost any function



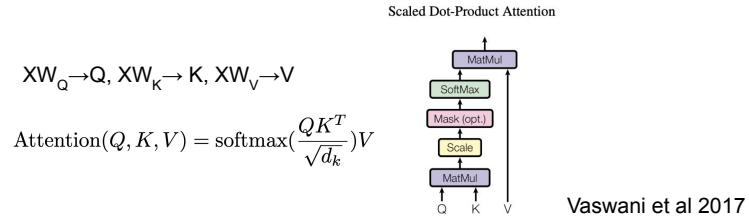
<https://cs231n.github.io/neural-networks-1/>

- ❑ **Convolutional neural networks:**  
detect position-invariant local patterns



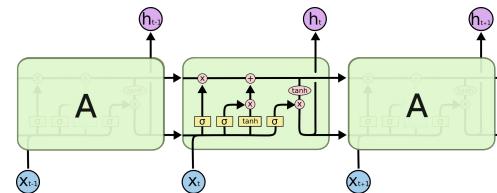
Zeng et al 2016

- ❑ **Attention layers:**  
emphasizes certain parts of the layer before based on its content



Vaswani et al 2017

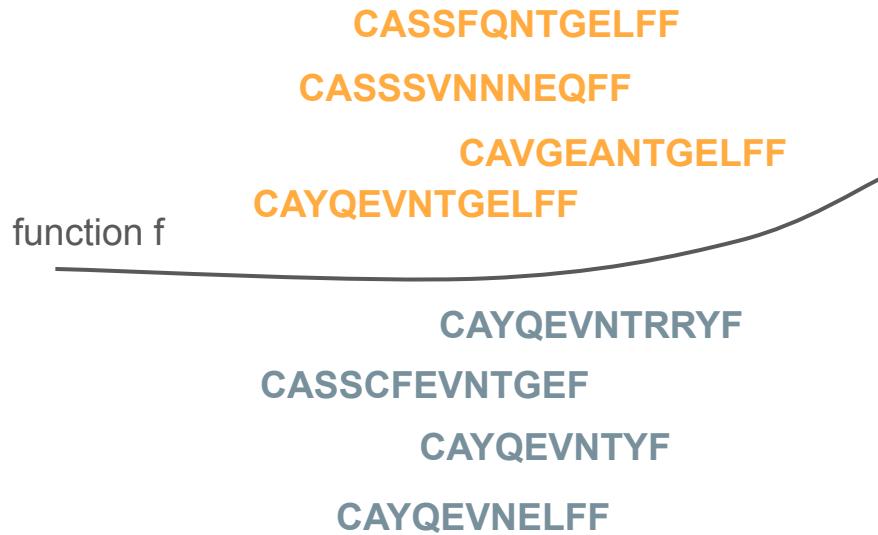
- ❑ **Recurrent neural networks:**  
can be Turing-complete, often used for long(er)-term dependencies in e.g., sequence data



<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

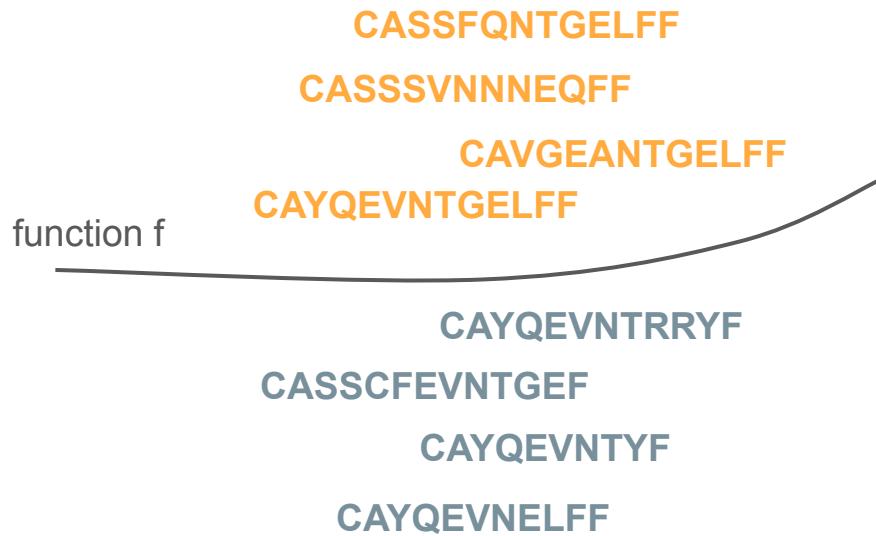
# Generative models

# Generative models

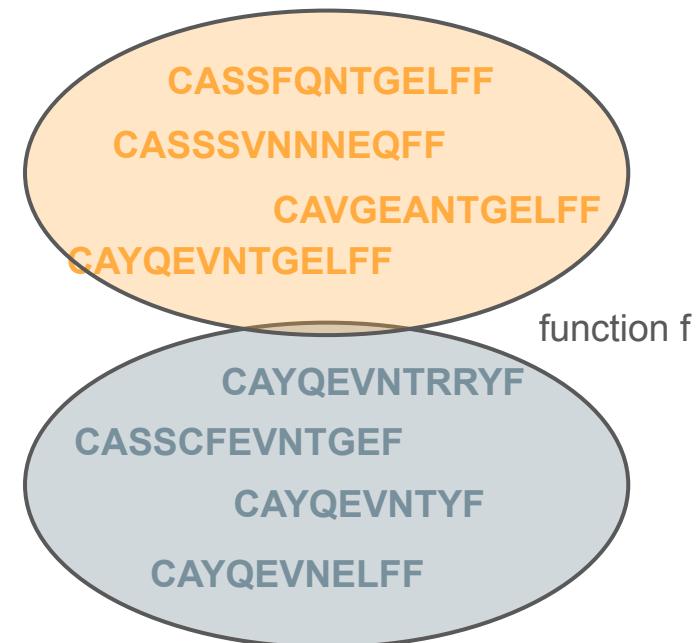


discriminative model learns the  
boundary between classes

# Generative models

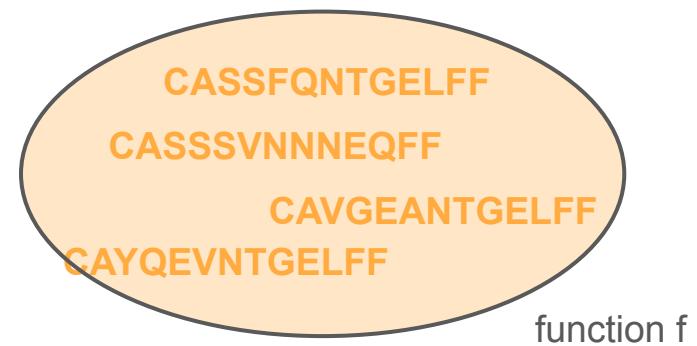


discriminative model learns the  
boundary between classes



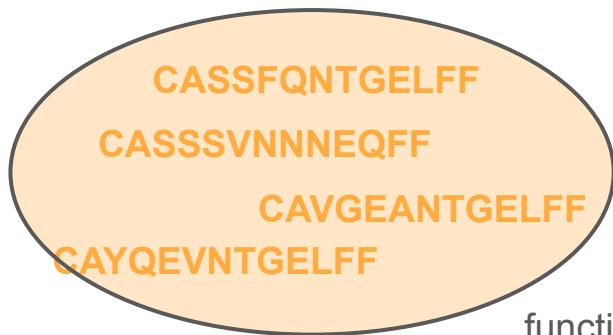
generative model learns the  
probability distribution of the data

# Generative models



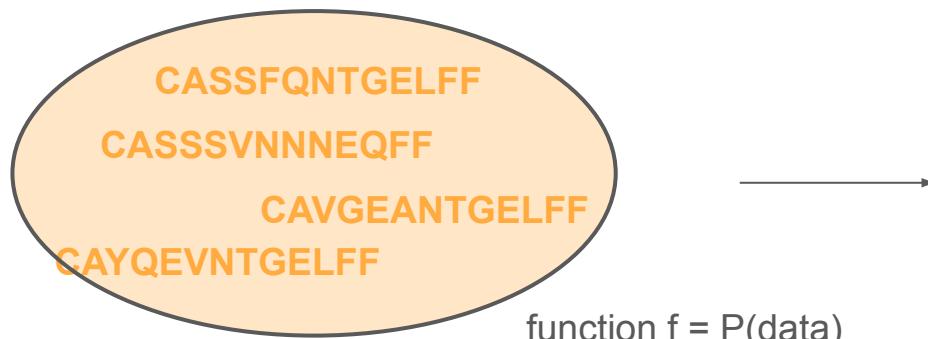
generative model learns the  
probability distribution of the data

# Generative models



function  $f = P(\text{data})$

# Generative models

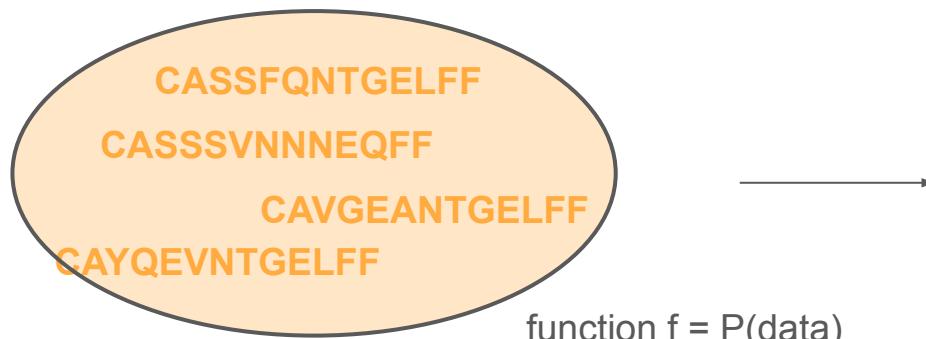


Probability distribution of  
antigen-specific sequences

## Goals:

- Generate new data [e.g., sequences with same specificity]
- Density estimation
- Imputation of missing values

# Generative models



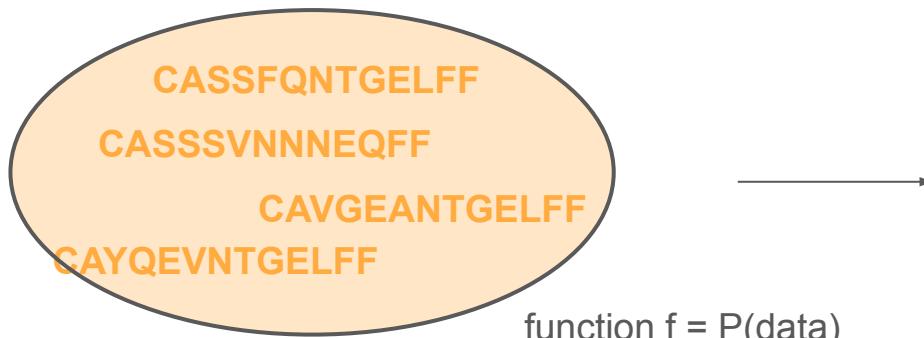
Probability distribution of  
antigen-specific sequences

## Goals:

- Generate new data [e.g., sequences with same specificity]
- Density estimation
- Imputation of missing values

Examples: variational autoencoders,  
generative adversarial networks, diffusion  
models, mixture models...

# Generative models



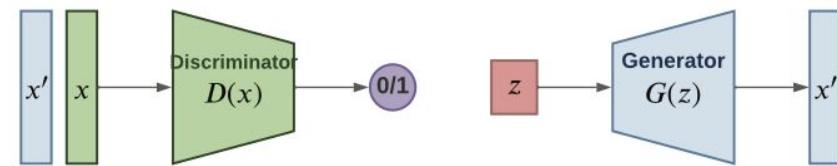
Probability distribution of  
antigen-specific sequences

Goals:

- Generate new data [e.g., sequences with same specificity]
- Density estimation
- Imputation of missing values

Example: generative adversarial network

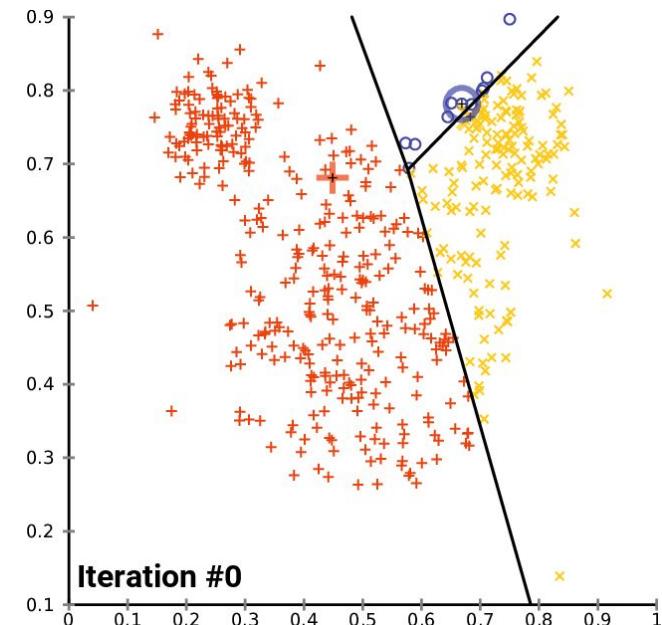
GAN:  
Adversarial  
training



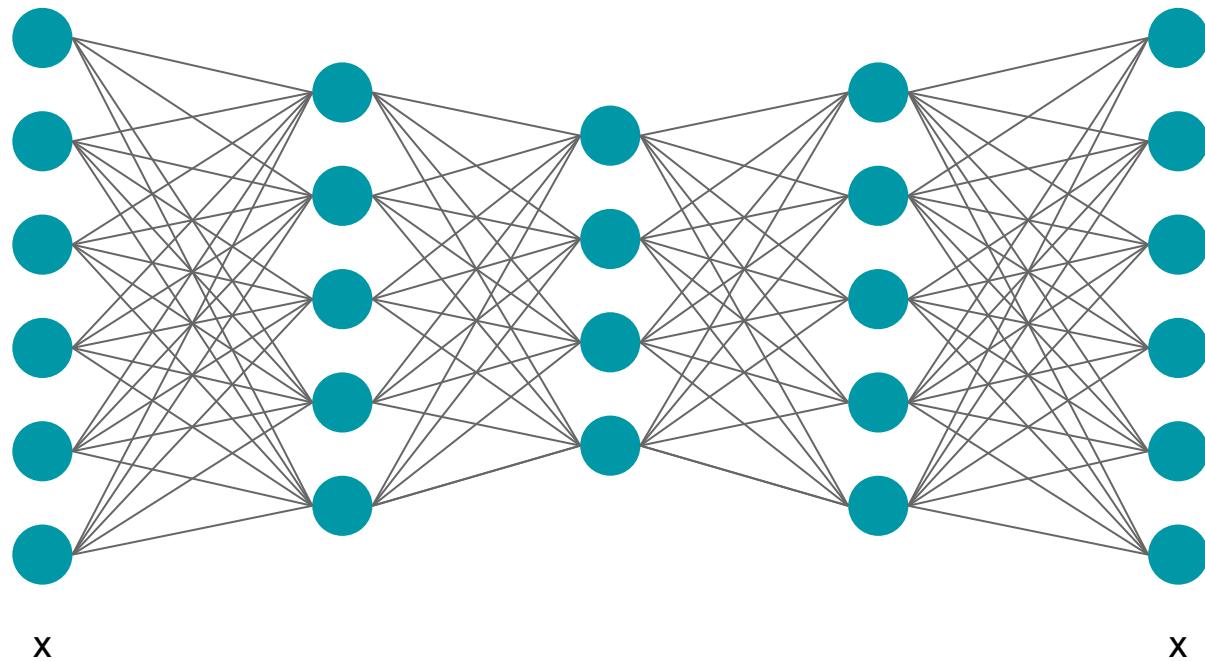
GAN (Goodfellow et al. 2014); image from: Murphy 2023

# Unsupervised learning - k-means clustering

- ❑ Given a set of data points, split them into  $k$  clusters so that the distances between points in the cluster are minimal
- ❑ Algorithm:
  - ❑ Pick  $k$  random points as cluster centers
  - ❑ Repeat until there are no more changes:
    - ❑ For each point, compute distances to each cluster center and assign it to the nearest cluster
    - ❑ Recompute cluster centroids as means of points in the cluster

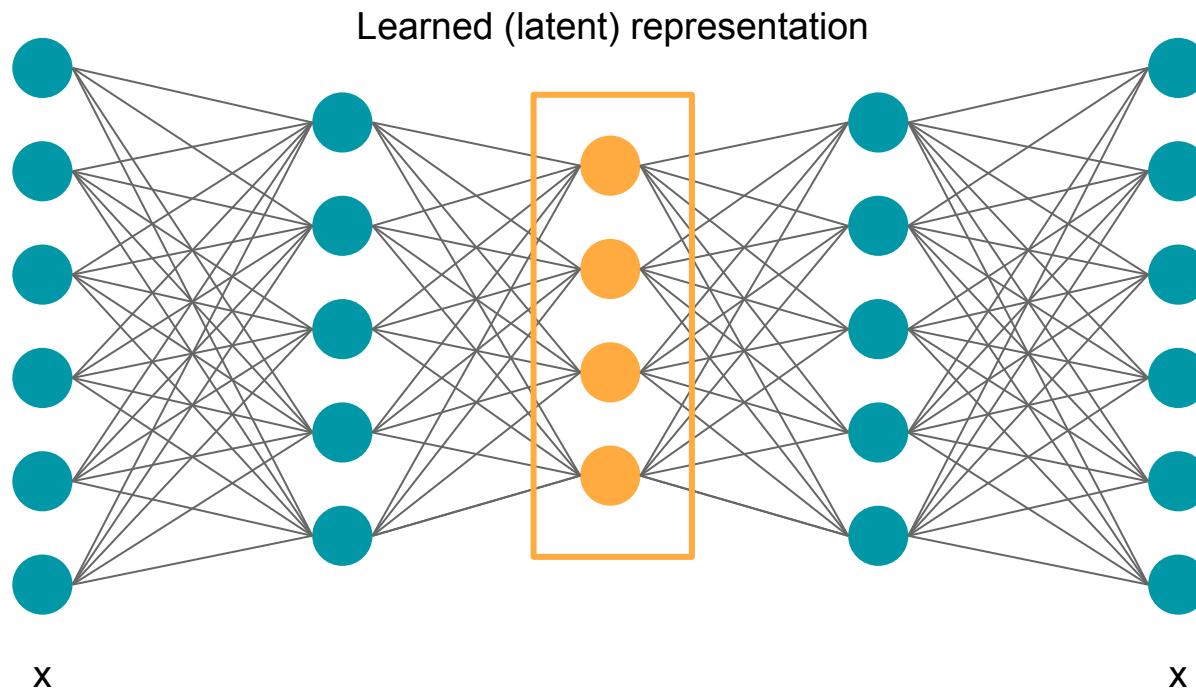


# Unsupervised learning - autoencoder



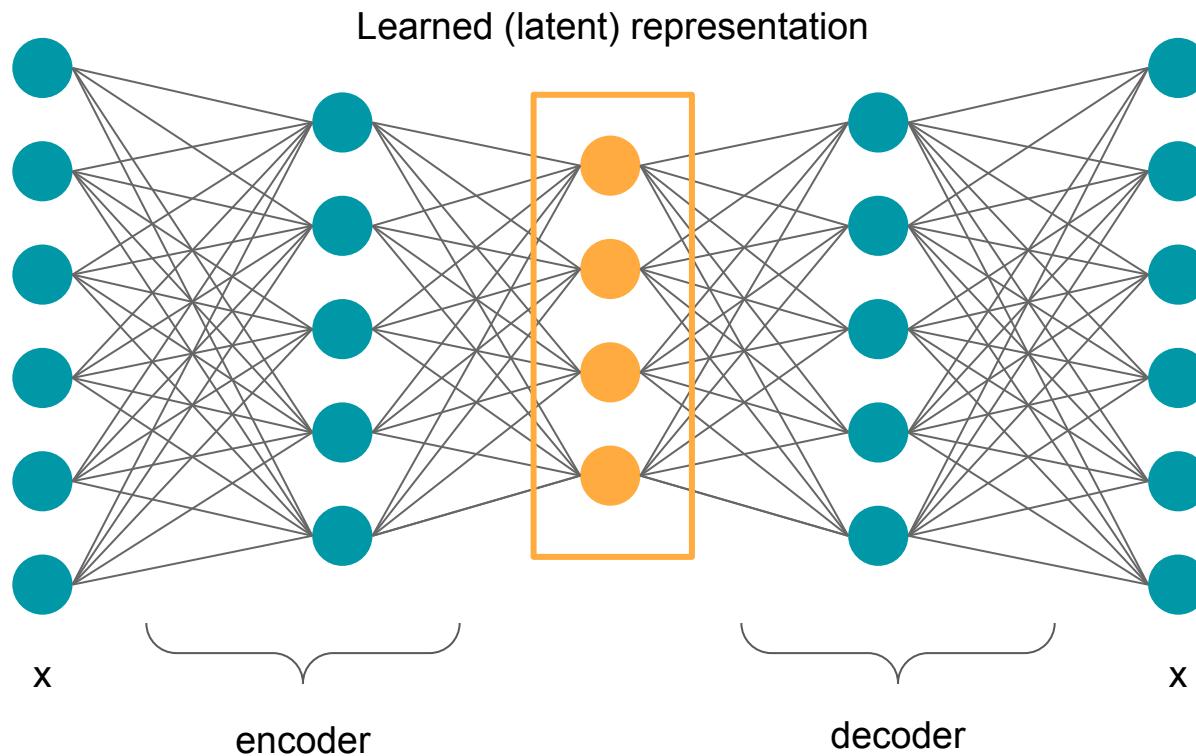
- ❑ Autoencoder is a neural network trained to attempt to copy its input to its output

# Unsupervised learning - autoencoder



- ❑ Autoencoder is a neural network trained to attempt to copy its input to its output while passing through a latent representation

# Unsupervised learning - autoencoder



- ❑ Autoencoder is a neural network trained to attempt to copy its input to its output while passing through a latent representation
- ❑ Learned representation can have useful properties: reduced dimensionality, easy to visualize, but there are other tasks as well

# References

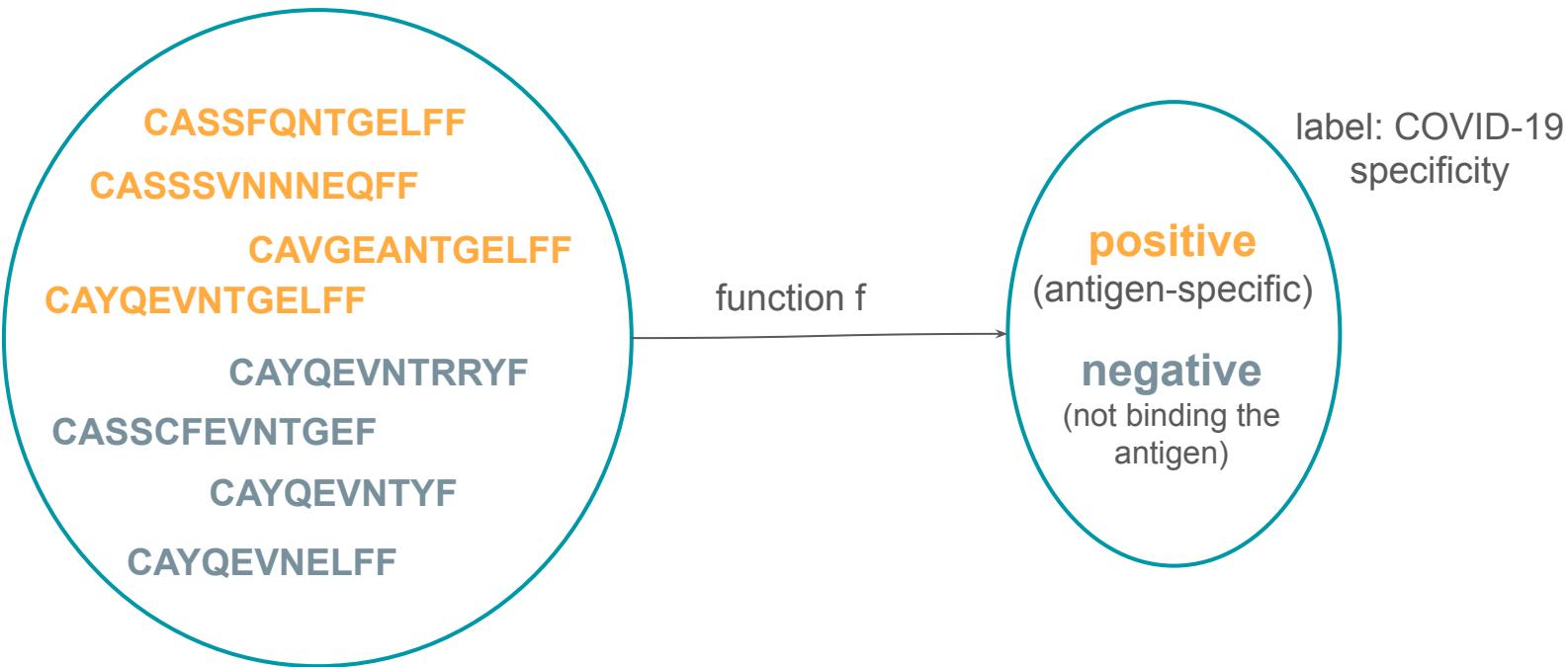
- Dietterich TG. Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms. *Neural Comput.* 1998;10(7):1895–1923. doi:[10.1162/089976698300017197](https://doi.org/10.1162/089976698300017197)
- Goodfellow IJ, Bengio Y, Courville A. *Deep Learning*. MIT Press; 2016. <https://mitpress.mit.edu/books/deep-learning> (especially chapter 5 - Machine Learning Basics)
- Killoran, Nathan, Leo J. Lee, Andrew Delong, David Duvenaud, and Brendan J. Frey. ‘Generating and Designing DNA with Deep Generative Models’. *ArXiv:1712.06148 [Cs, q-Bio, Stat]*, 17 December 2017. <http://arxiv.org/abs/1712.06148>.
- Zeng, Haoyang, Matthew D. Edwards, Ge Liu, and David K. Gifford. ‘Convolutional Neural Network Architectures for Predicting DNA–Protein Binding’. *Bioinformatics* 32, no. 12 (15 June 2016): i121–27. <https://doi.org/10.1093/bioinformatics/btw255>.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., & Bengio, Y. (2014). Generative Adversarial Networks. *arXiv:1406.2661 [Cs, Stat]*. <http://arxiv.org/abs/1406.2661>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention Is All You Need. *arXiv:1706.03762 [Cs]*. <http://arxiv.org/abs/1706.03762>
- Murphy, K. P. (2023). *Probabilistic Machine Learning: Advanced Topics*. MIT Press. (Chapter 20: advanced text covering generative models)

# Machine learning in computational biology - outline

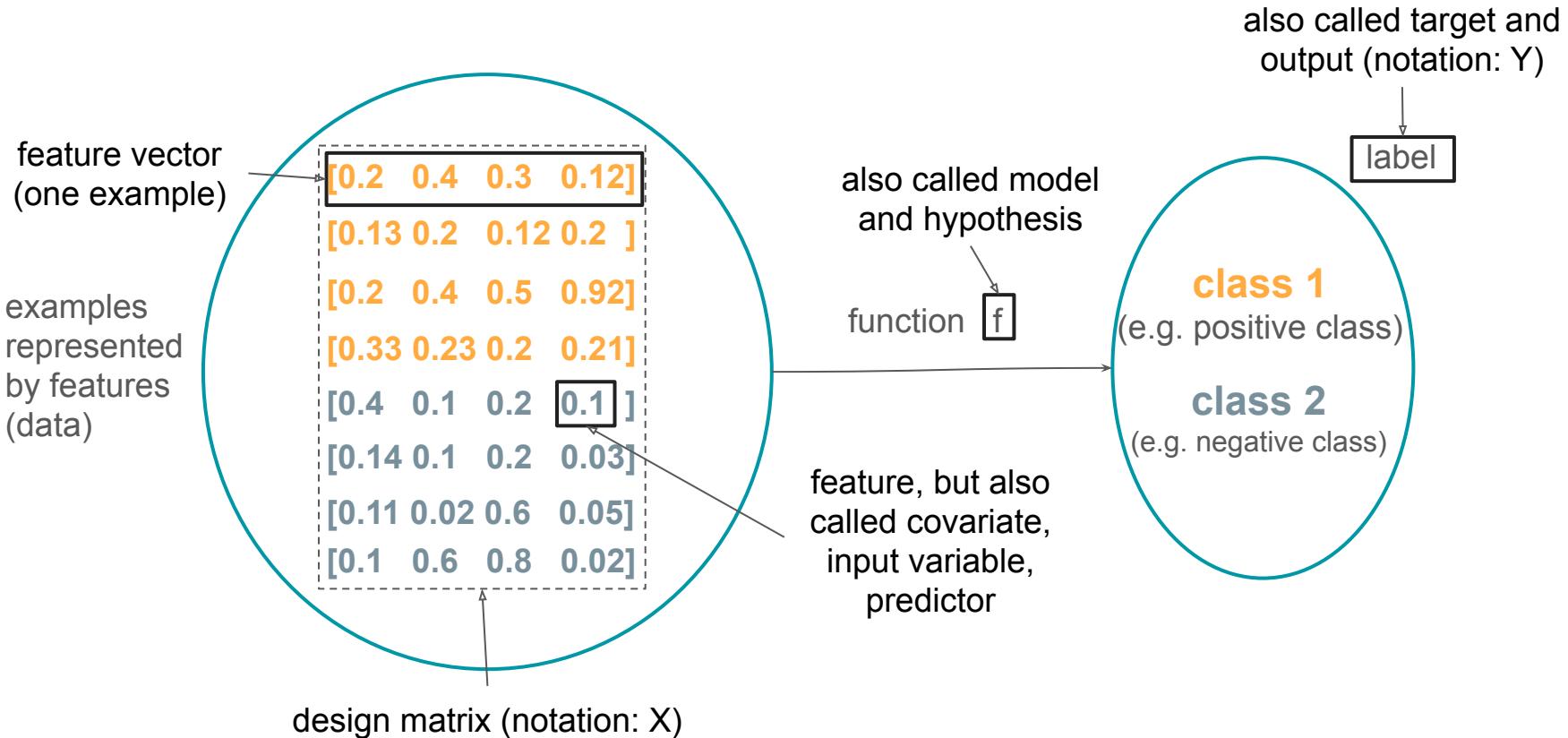
- Introduction to machine learning:
  - What is machine learning, types of problems, assumptions, workflow, generalization
- Machine learning models and algorithms:
  - Discriminative vs generative models, supervised models (logistic and linear regression, kNN, neural networks), unsupervised models (dimensionality reduction, clustering)
- **Data representation:**
  - Considerations and examples, one-hot encoding, feature engineering, representation learning
- Model comparison and uncertainty:
  - Model assessment, model selection, uncertainty, cross-validation
- Transparency and reproducibility

# We are given a set of sequences... but algorithms only understand numbers!

examples represented by features:  
immune receptor sequences



# We can represent sequences by their physicochemical properties, for instance



# Some examples of data representation (encoding)

- ❑ One-hot encoding
- ❑ K-mer frequencies

Data representation heavily depends on data, so in different domains, there will be different representations:

When classifying images with classical approaches: number of edges, objects

When predicting the length of the trip: number of traffic lights, time of day

When predicting if an email is a spam or not: certain words, presence/absence of personal name

# Some examples of data representation (encoding)

- ❑ One-hot encoding
  - ❑ K-mer frequencies
- We have to be careful how we choose features - we must not introduce information that should not be there!

Data representation heavily depends on data, so in different domains, there will be different representations:

When classifying images with classical approaches: number of edges, objects

When predicting the length of the trip: number of traffic lights, time of day

When predicting if an email is a spam or not: certain words, presence/absence of personal name

# One-hot encoding

- ❑ A common way to represent categorical data where only one value can be chosen: rows represent the possible values
- ❑ Also called *dummy variables* in statistics

nucleotide sequence: AATGC



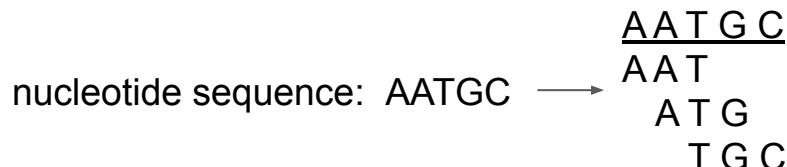
A A T G C

is it A	1	1	0	0	0
is it C	0	0	0	0	1
is it G	0	0	0	1	0
is it T	0	0	1	0	0

one-hot encoding

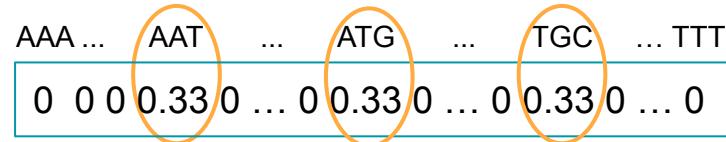
# K-mer frequency

- ❑ Often used for sequence representation
- ❑ k-mers are (optionally overlapping) subsequences of length k

nucleotide sequence: AATGC →  
  
A A T G C  
A A T  
A T G  
T G C

present 3-mers: AAT, ATG, TGC

all possible 3-mers: AAA, AAC, AAG, AAT, ACA, ..., TTT  
( $4^3=64$  combinations)



k-mer frequency encoding (k=3)

# ML algorithm performance heavily depends on data representation

- Data representation refers to choosing and constructing features
- We don't always know in advance which features are the best for the problem: we have to know the domain:

CASSFQNTGELYF  
CASSSVNNNEYFF  
CAVGEANTGELFF  
CAYQEVTGELFF  
  
CAYQEVNTRRYF  
CASSCFEVNTGEF  
CAYQEVTNTYF  
CAYQEVLNELFF

raw data

represent sequence as  
1 if it contains Y  
(tyrosine) and 0 if not:  
 $[1, 1, 0, 1, 1, 0, 1, 1]^T$

data representation:  
option 1

represent the sequence  
by k-mer frequency:  
 $[0, \dots, 0.02, 0.03, \dots, 0]$

data representation:  
option 2

Which one is better?

# Feature engineering & feature selection

- ❑ Feature engineering: together with domain experts, ML researchers would discuss and derive features which they believe could be useful for the model

Example: for biological sequences, there are a few popular alternatives like k-mer frequencies and physicochemical properties
- ❑ This way a lot of features could be constructed and the best ones would be selected as a part of fitting the model (feature selection)

# Representation learning

- ❑ Most often in context of neural networks: the many layers of the network learn a hierarchical, alternative representation of the (raw) data that was provided as input – usually much better than manually derived features
- ❑ Inductive biases

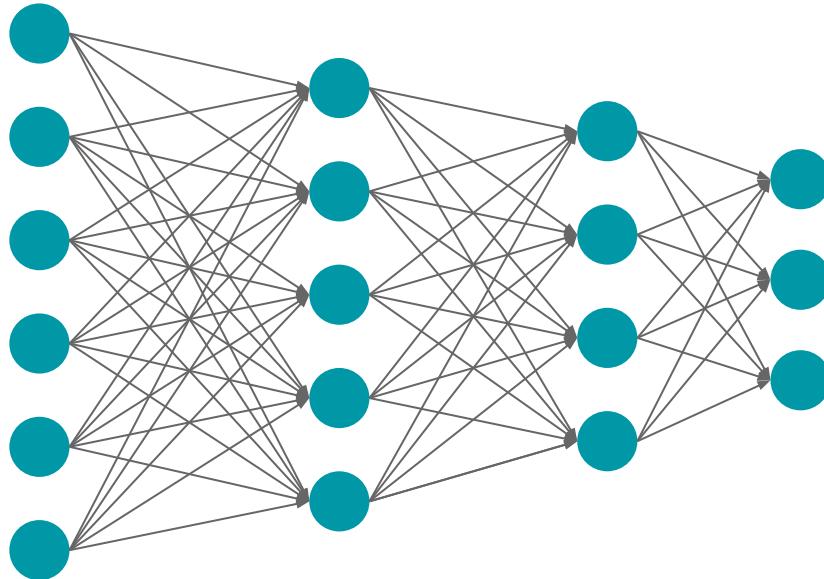
IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, VOL. 35, NO. 8, AUGUST 2013

## Representation Learning: A Review and New Perspectives

Yoshua Bengio, Aaron Courville, and Pascal Vincent

**Abstract**—The success of machine learning algorithms generally depends on data representation, and we hypothesize that this is because different representations can entangle and hide more or less the different explanatory factors of variation behind the data. Although specific domain knowledge can be used to help design representations, learning with generic priors can also be used, and the quest for AI is motivating the design of more powerful representation-learning algorithms implementing such priors. This paper reviews recent work in the area of unsupervised feature learning and deep learning, covering advances in probabilistic models, autoencoders, manifold learning, and deep networks. This motivates longer term unanswered questions about the appropriate objectives for learning good representations, for computing representations (i.e., inference), and the geometrical connections between representation learning, density estimation, and manifold learning.

# Representation learning - hidden layers in neural networks can be seen as different representations

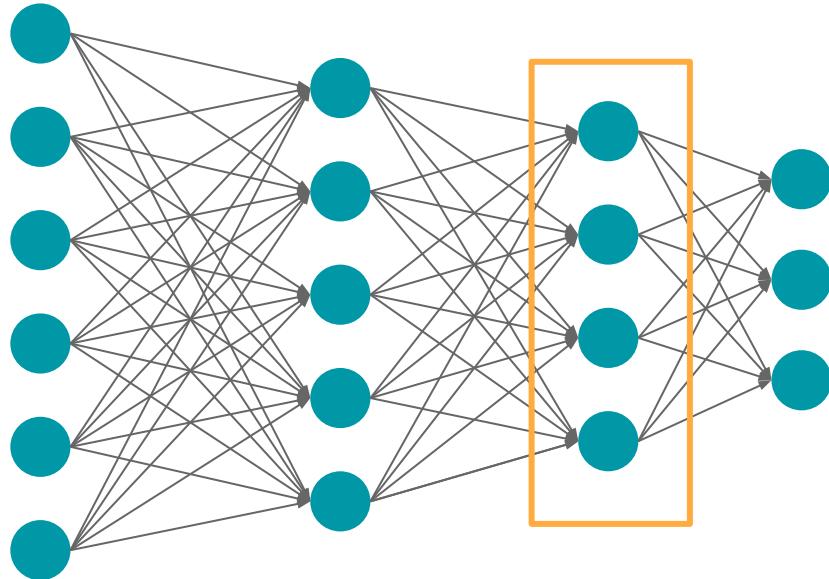


Deep neural network

“A good representation is the one that makes the learning task easier.”

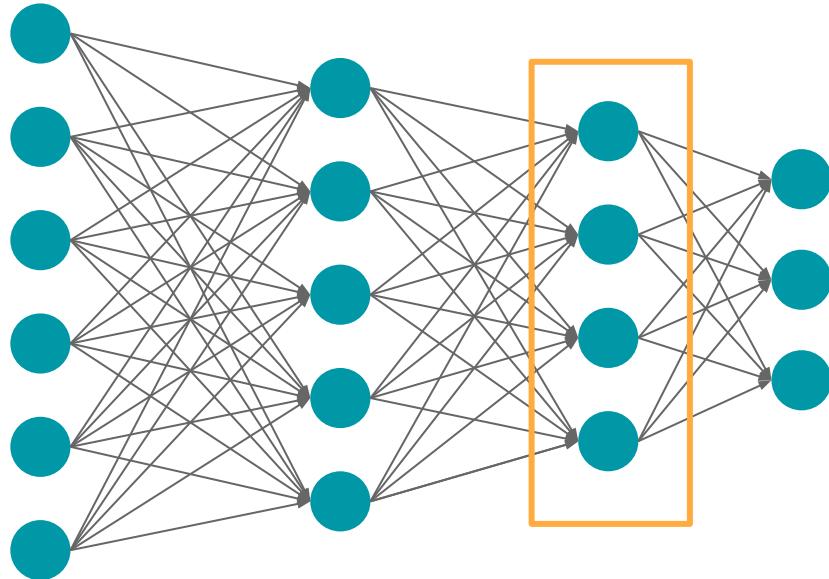
Goodfellow et al. 2016

# Representation learning - hidden layers in neural networks can be seen as different representations

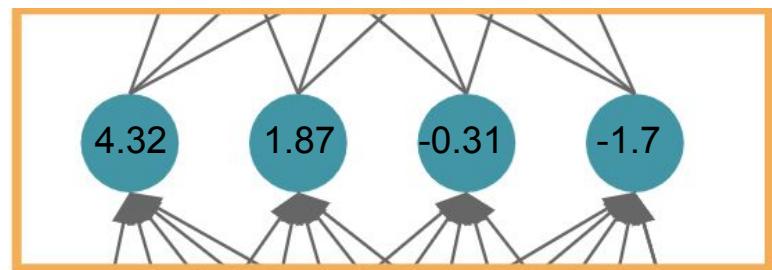


Deep neural network

# Representation learning - hidden layers in neural networks can be seen as different representations

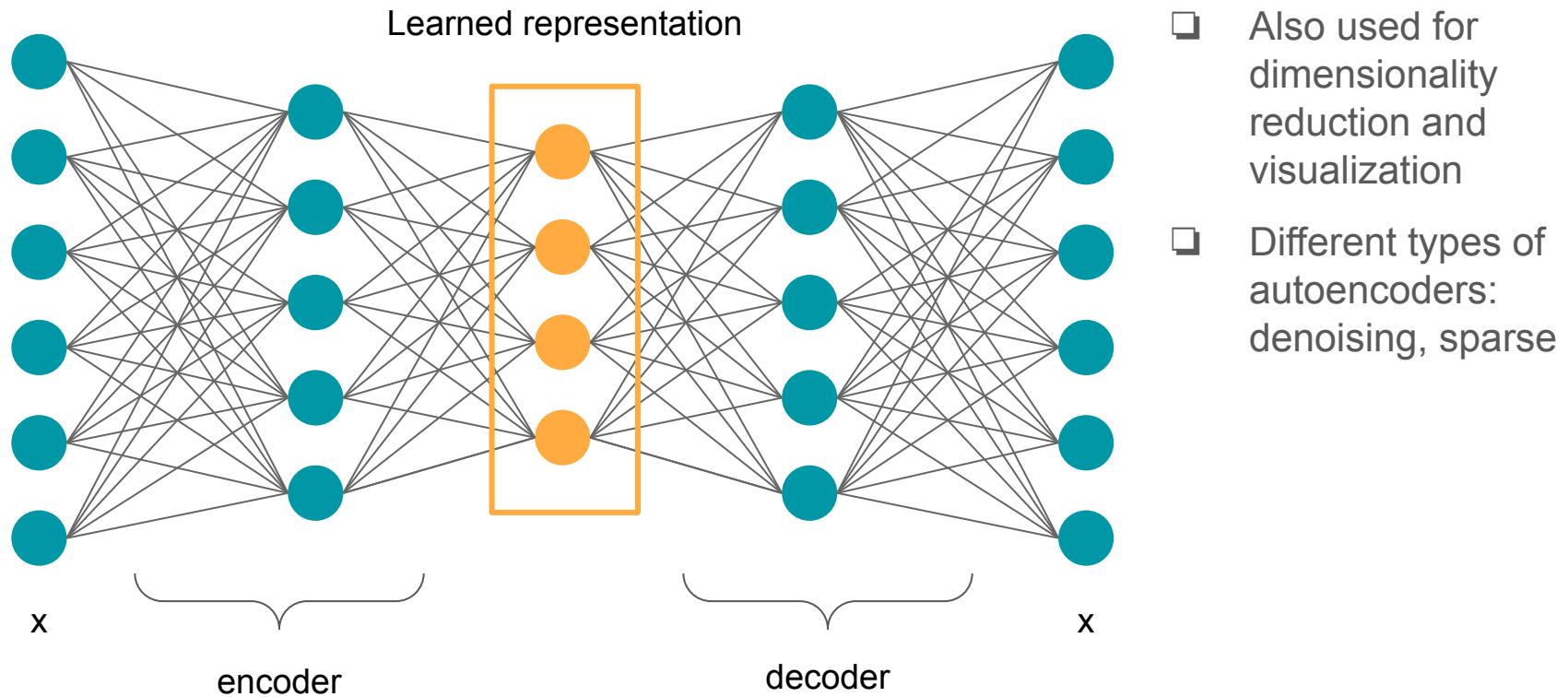


Deep neural network



New data representation

# Representation learning with autoencoders



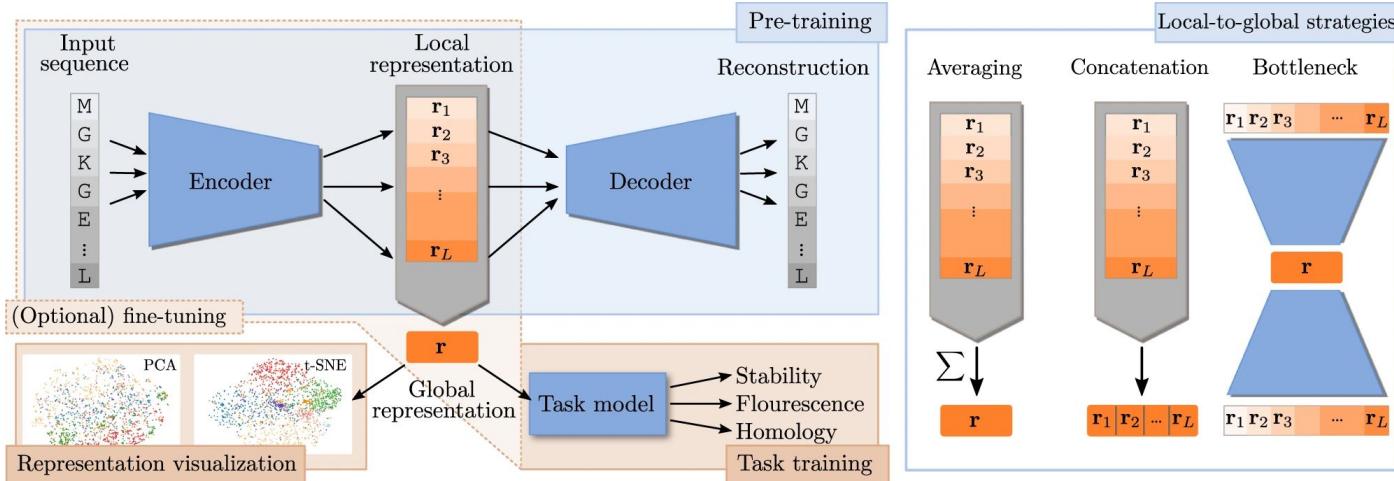
# Representation learning: protein sequence example

Article | [Open access](#) | Published: 08 April 2022

## Learning meaningful representations of protein sequences

[Nicki Skafte Detlefsen, Søren Hauberg & Wouter Boomsma](#) 

[Nature Communications](#) 13, Article number: 1914 (2022) | [Cite this article](#)



# Machine learning in computational biology - outline

- Introduction to machine learning:
  - What is machine learning, types of problems, assumptions, workflow, generalization
- Machine learning models and algorithms:
  - Discriminative vs generative models, supervised models (logistic and linear regression, kNN, neural networks), unsupervised models (dimensionality reduction, clustering)
- Data representation:
  - Considerations and examples, one-hot encoding, feature engineering, representation learning
- **Model comparison and uncertainty:**
  - **Model assessment, model selection, uncertainty, cross-validation**
- Transparency and reproducibility

# Model selection and model assessment

In a typical workflow, we perform two tasks:

**Model selection:** estimating the performance of different models in order to choose the best one

**Model assessment:** having chosen the final model, estimating its prediction error (generalization error) on new data

# Model selection and model assessment

In a typical workflow, we perform two tasks:

**Model selection:** estimating the performance of **different models** in order to choose the best one



different learning algorithms or same learning algorithms, but different hyperparameters (e.g., different number of layers in a neural network)

**Model assessment:** having chosen the final model, estimating its prediction error (generalization error) on new data

# Evaluation of an ML algorithm

- ❑ To do model selection (and model assessment) we need to know how to do model evaluation
- ❑ A suitable cost function has to be chosen

# Evaluation of a ML algorithm

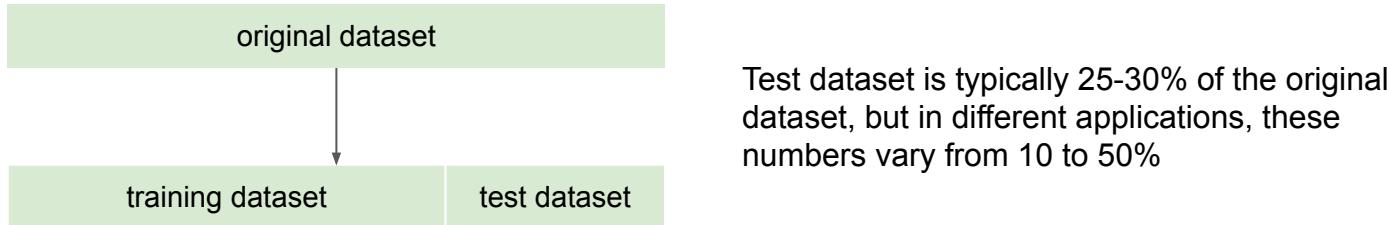
- ❑ Two ML algorithms are different if and only if their percentage of correct classifications would be different on average when trained on a training set of a given fixed size and tested on all data points in the population

Dietterich 1998

Of course we cannot test on the entire population, but we should be aware that the difference we see in one particular result might not be the true difference between algorithms

# Evaluation of an ML algorithm: random holdout test set

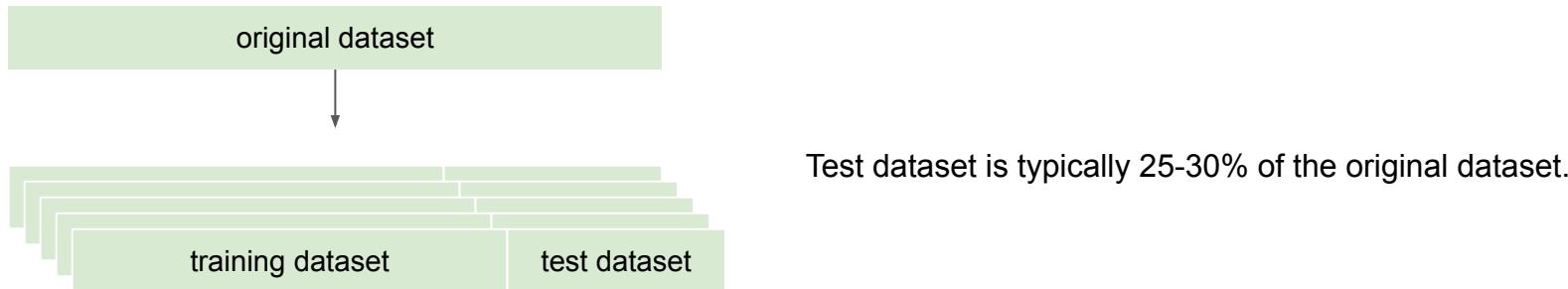
- ❑ The simplest scenario: split the dataset to train/test dataset



- ❑ Estimate performance of the algorithm on the test set

# Evaluation of a ML algorithm: random holdout test sets

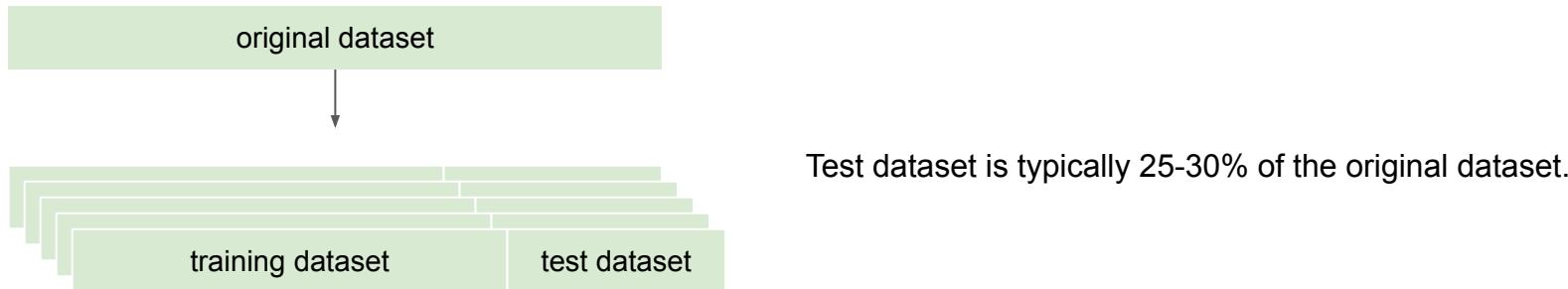
- ❑ A bit more robust scenario: split the dataset randomly into multiple train/test datasets



- ❑ Estimate performance of the algorithm as the average of the performances on test sets

# Evaluation of a ML algorithm: random holdout test sets

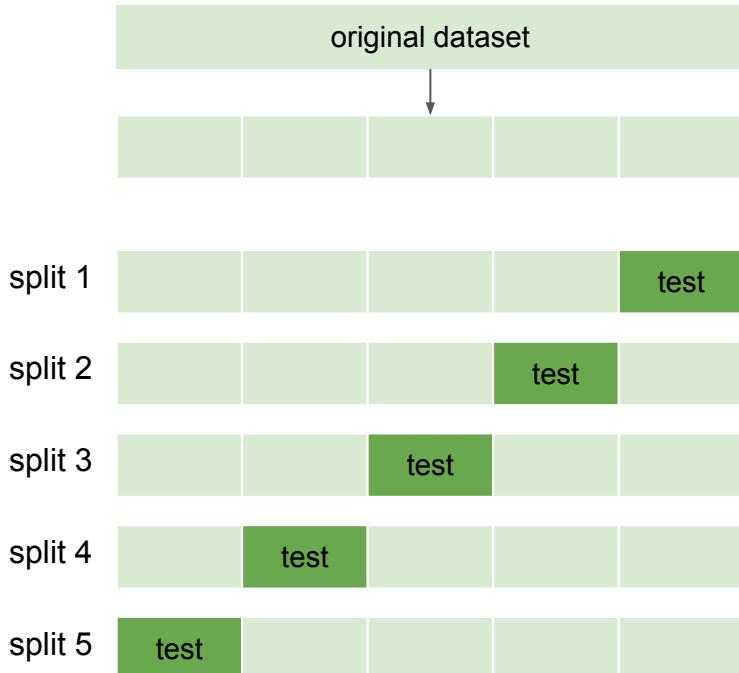
- ❑ A bit more robust scenario: split the dataset randomly into multiple train/test datasets



- ❑ Estimate performance of the algorithm as the average of the performances on test sets
- ❑ The same examples could be in multiple test sets: biased estimate of the performance!

# K-fold cross validation

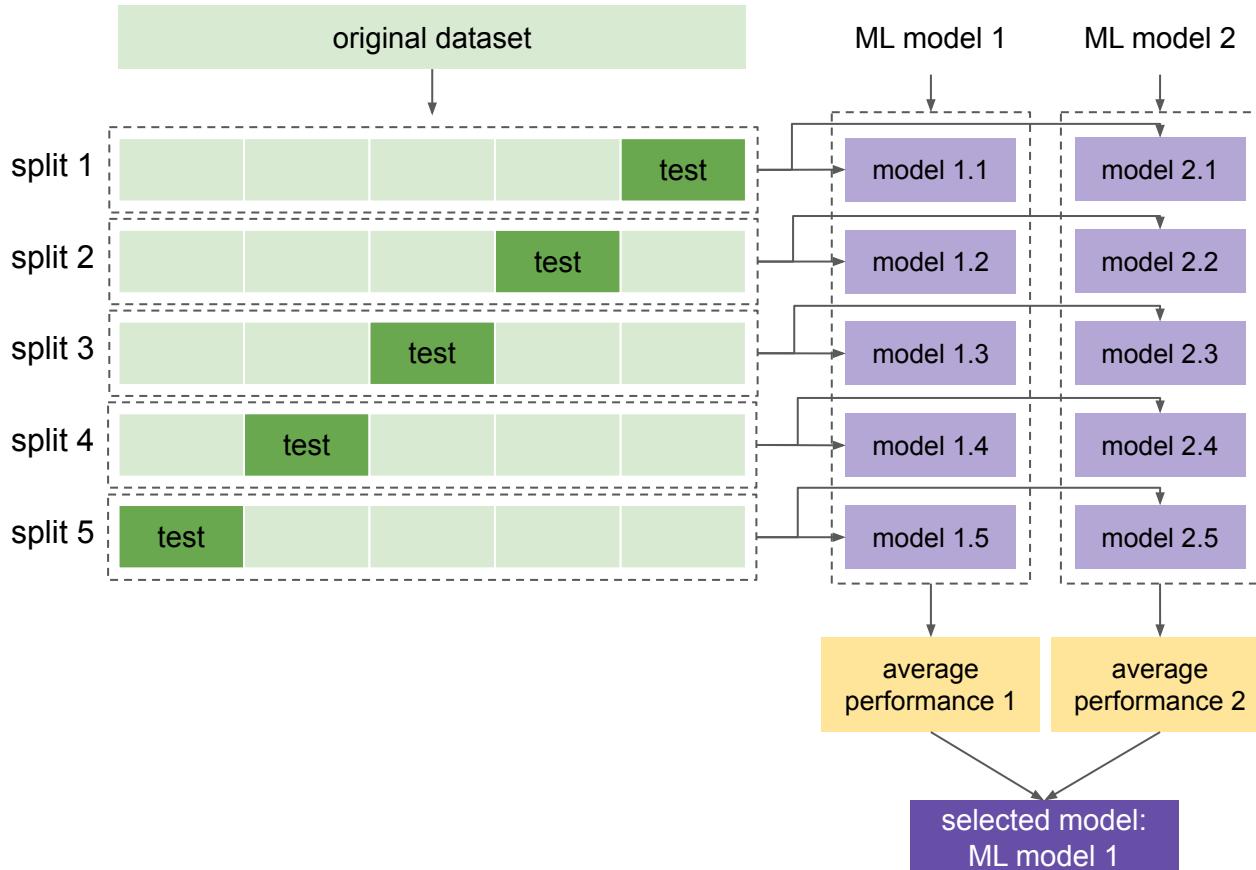
- ❑ K-fold cross validation (CV) refers to splitting the data into k training and test set pairs so that each example is once a part of the test set



- ❑ The model is trained and evaluated on each train/test set pair
- ❑ The estimated performance is the obtained as average over all test set performances
- ❑ Typical values of k: 5, 10
- ❑ Leave-one-out CV

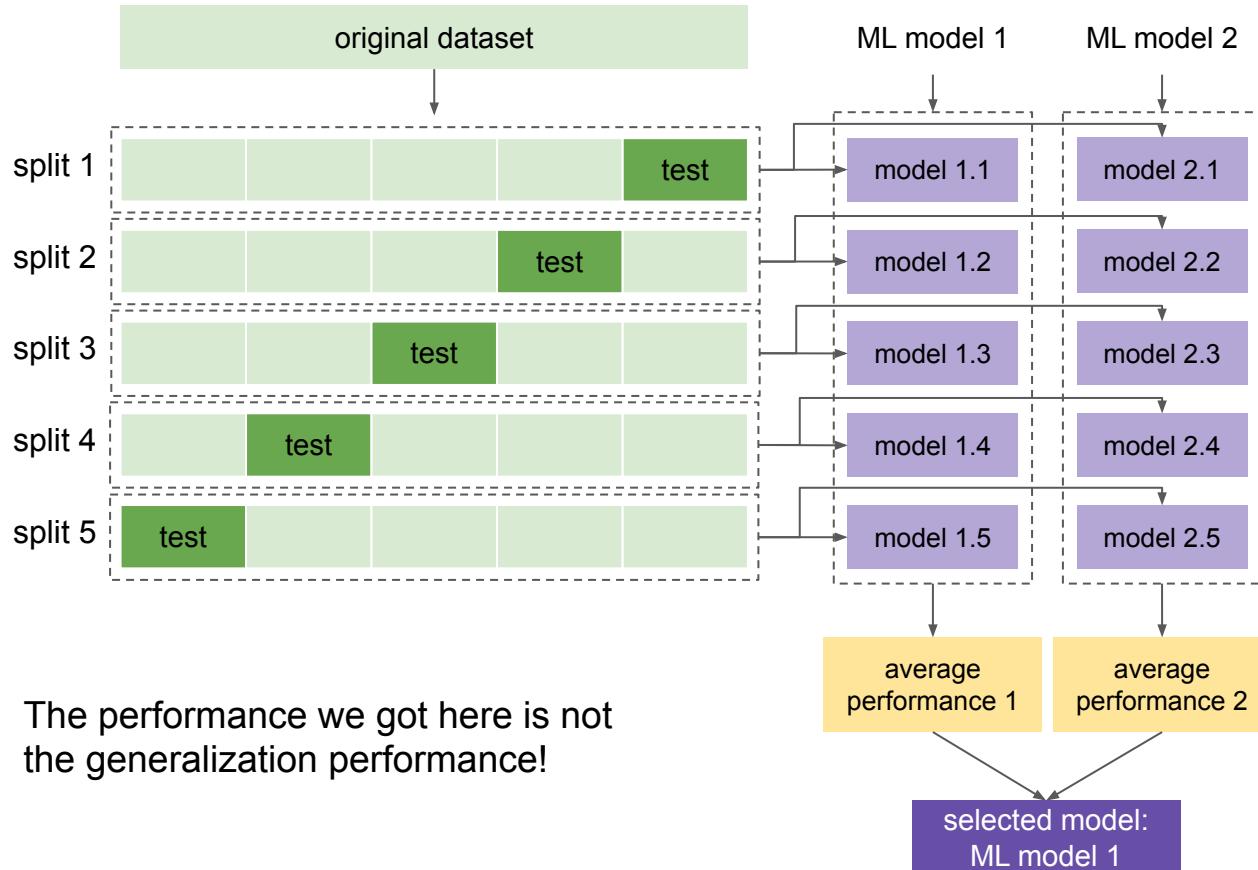
# How to do model selection

1. Split the original dataset to train/test dataset per one of the previous strategies (e.g., 5-fold CV)
2. For each model compute the performance as the average of performances on test sets (as described in the previous slide)
3. Select the model with best average performance

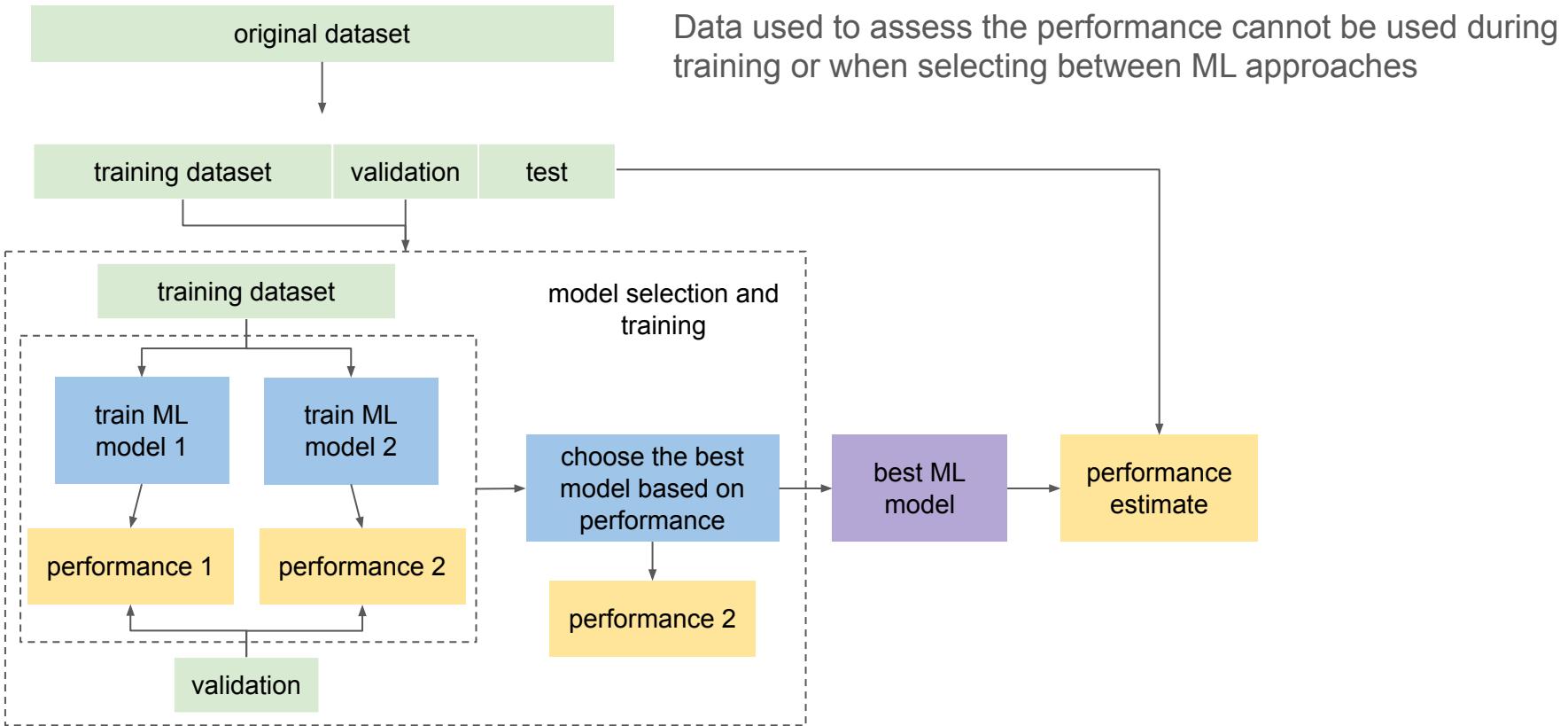


# How to do model selection

1. Split the original dataset to train/test dataset per one of the previous strategies (e.g., 5-fold CV)
2. For each model compute the performance as the average of performances on test sets (as described in the previous slide)
3. Select the model with best average performance

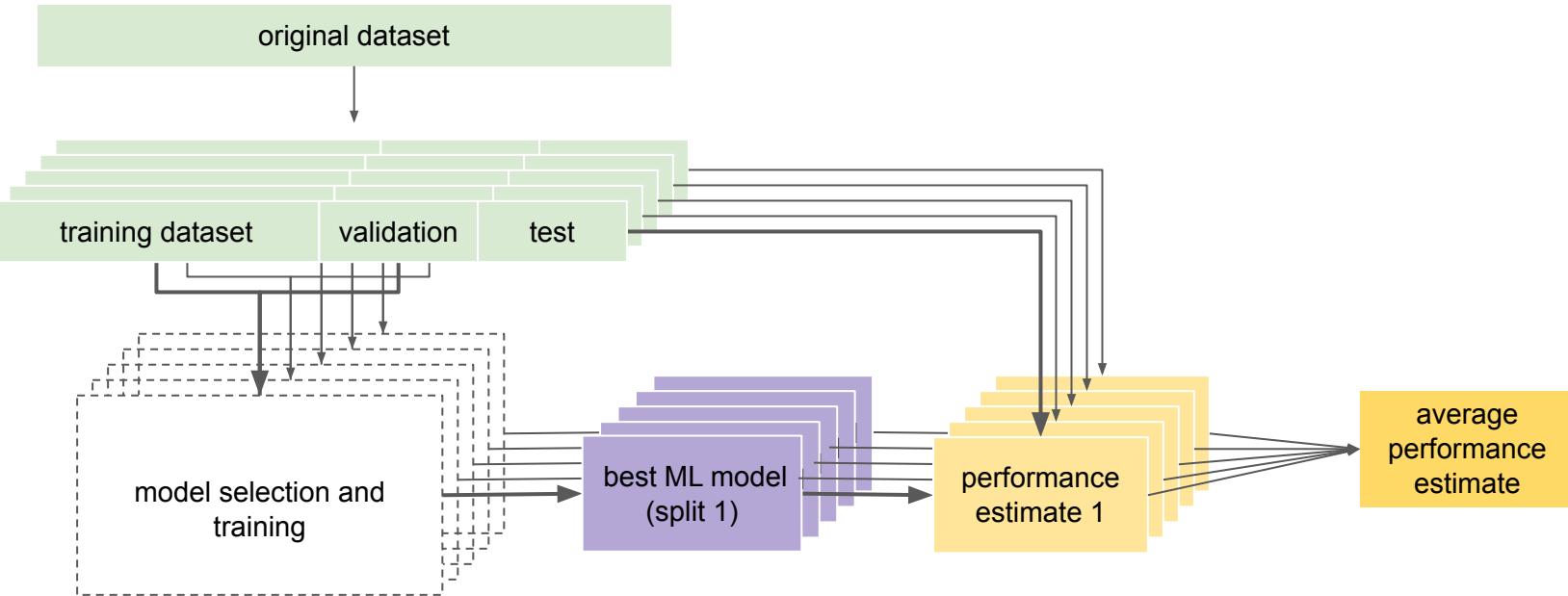


# Model assessment uses a separate test set not used during model selection



# Model assessment with multiple test sets: a more robust estimate

- We can employ all the same techniques (e.g., CV) to split the data, do model selection as described before and use separate test sets to estimate generalization performance



# Recommendations for ML in biology

## DOME: recommendations for supervised machine learning validation in biology

Ian Walsh, Dmytro Fishman, Dario Garcia-Gasulla, Tiina Titma, Gianluca Pollastri, ELIXIR Machine Learning Focus Group, Jennifer Harrow & Fotis E. Psomopoulos & Silvio C. E. Tosatto

*Nature Methods* 18, 1122–1127 (2021) | [Cite this article](#)

**Table 1 | Supervised ML in biology: concerns, the consequences they impart and recommendations**

Broad topic	Be on the lookout for	Consequences	Recommendation(s)
Data	<ul style="list-style-type: none"><li>Inadequate data size &amp; quality</li><li>Inappropriate partitioning, dependence between train and test data</li><li>Class imbalance</li><li>No access to data</li></ul>	<ul style="list-style-type: none"><li>Data not representative of domain application</li><li>Unreliable or biased performance evaluation</li><li>Cannot check data credibility</li></ul>	<ul style="list-style-type: none"><li><b>Use independent optimization (training) and evaluation (testing) sets.</b> This is especially important for meta algorithms, where independence of multiple training sets must be shown to be independent of the evaluation (testing) sets.</li><li><b>Release data, preferably using appropriate long-term repositories, and include exact splits.</b></li><li>Offer sufficient evidence of data size &amp; distribution being representative of the domain.</li></ul>
Optimization	<ul style="list-style-type: none"><li>Overfitting, underfitting and illegal parameter tuning</li><li>Imprecise parameters and protocols given</li></ul>	<ul style="list-style-type: none"><li>Reported performance is too optimistic or too pessimistic</li><li>The model models noise or misses relevant relationships</li><li>Results are not reproducible</li></ul>	<ul style="list-style-type: none"><li><b>Clarify that evaluation sets were not used for feature selection, preprocessing steps or parameter tuning.</b></li><li><b>Report indicators on training and testing data that can aid in assessing the possibility of under- or overfitting; for example, train vs. test error.</b></li><li><b>Release definitions of all algorithmic hyperparameters, regularization protocols, parameters and optimization protocol.</b></li><li>For neural networks, release definitions of training and learning curves.</li><li>Include explicit model validation techniques, such as N-fold cross-validation.</li></ul>
Model	<ul style="list-style-type: none"><li>Unclear if black box or interpretable model</li><li>No access to resulting source code, trained models &amp; data</li><li>Execution time impractical</li></ul>	<ul style="list-style-type: none"><li>An interpretable model shows no explainable behavior</li><li>Cannot cross compare methods &amp; reproducibility, or check data credibility</li><li>Model takes too much time to produce results</li></ul>	<ul style="list-style-type: none"><li><b>Describe the choice of black box or interpretable model. If interpretable, show examples of interpretable output.</b></li><li>Release documented source code + models + executable + user interface/webserver + software containers.</li><li>Report execution time averaged across many repeats. If computationally tough, compare to similar methods.</li></ul>
Evaluation	<ul style="list-style-type: none"><li>Performance measures inadequate</li><li>No comparisons to baselines or other methods</li><li>Highly variable performance</li></ul>	<ul style="list-style-type: none"><li>Biased performance measures reported</li><li>The method is falsely claimed as state-of-the-art</li><li>Unpredictable performance in production</li></ul>	<ul style="list-style-type: none"><li><b>Compare with public methods &amp; simple models (baselines).</b></li><li><b>Adopt community-validated measures and benchmark datasets for evaluation.</b></li><li>Compare related methods and alternatives on the same dataset.</li><li>Evaluate performance on a final independent held-out set.</li><li><b>Use confidence intervals/error intervals and statistical tests to gauge prediction robustness.</b></li></ul>

Key recommendations are bolded.

# Model assessment and selection for unsupervised learning

- ❑ Not as straightforward: if there are no labels - how to assess the quality of the models?
- ❑ One example for clustering validation: there is still a need for having validation/test data

Received: 1 March 2021 | Revised: 16 November 2021 | Accepted: 27 November 2021  
DOI: 10.1002/widm.1444

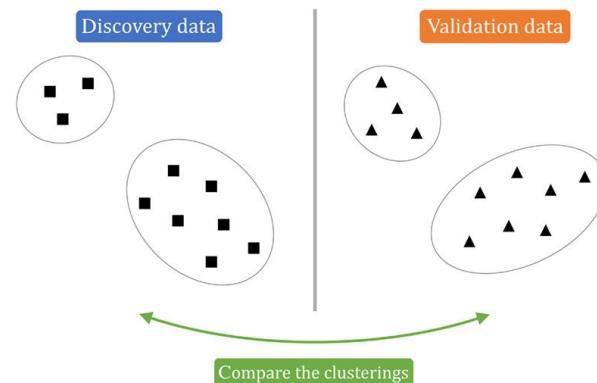
ADVANCED REVIEW

 WIREs  
DATA MINING AND KNOWLEDGE DISCOVERY

WILEY

## Validation of cluster analysis results on validation data: A systematic framework

Theresa Ullmann<sup>1</sup>  | Christian Hennig<sup>2</sup>  | Anne-Laure Boulesteix<sup>1</sup> 



# References

- Hastie T, Tibshirani R, Friedman J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition*. 2nd ed. Springer-Verlag; 2009. doi:[10.1007/978-0-387-84858-7](https://doi.org/10.1007/978-0-387-84858-7)
- Dietterich TG. Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms. *Neural Comput.* 1998;10(7):1895–1923. doi:[10.1162/089976698300017197](https://doi.org/10.1162/089976698300017197)
- Goodfellow IJ, Bengio Y, Courville A. *Deep Learning*. MIT Press; 2016. <https://mitpress.mit.edu/books/deep-learning>
- Ullmann, T., Hennig, C., & Boulesteix, A.-L. (2022). Validation of cluster analysis results on validation data: A systematic framework. *WIREs Data Mining and Knowledge Discovery*, 12(3), e1444. <https://doi.org/10.1002/widm.1444>

# Machine learning in computational biology - outline

- Introduction to machine learning:
  - What is machine learning, types of problems, assumptions, workflow, generalization
- Machine learning models and algorithms:
  - Discriminative vs generative models, supervised models (logistic and linear regression, kNN, neural networks), unsupervised models (dimensionality reduction, clustering)
- Data representation:
  - Considerations and examples, one-hot encoding, feature engineering, representation learning
- Model comparison and uncertainty:
  - Model assessment, model selection, uncertainty, cross-validation
- Transparency and reproducibility

# Transparency & interpretability are crucial for ML

- ❑ “**Transparency** refers to how easily a stakeholder can examine the model and understand, or explain, how the model operates when combining the inputs to produce output, regardless how accurate the model is”.

Wainberg et al. 2018

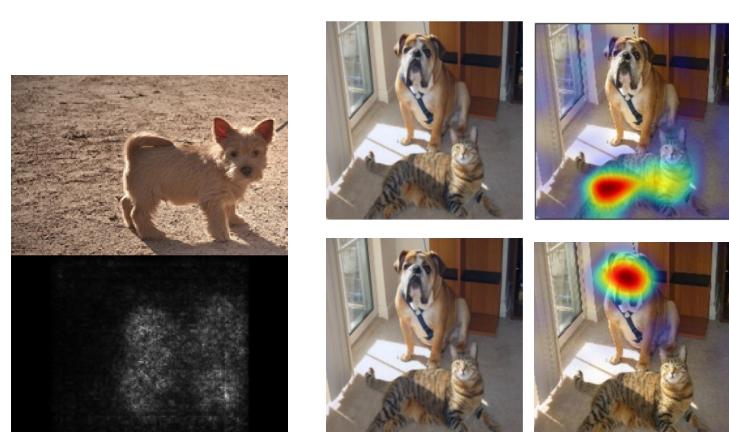
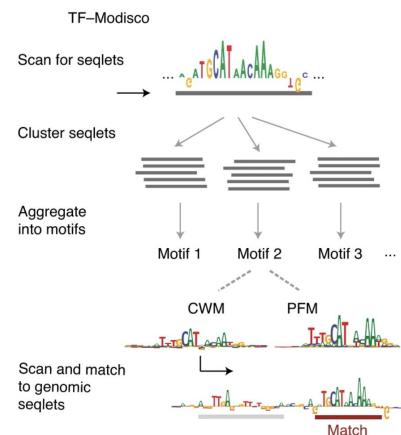
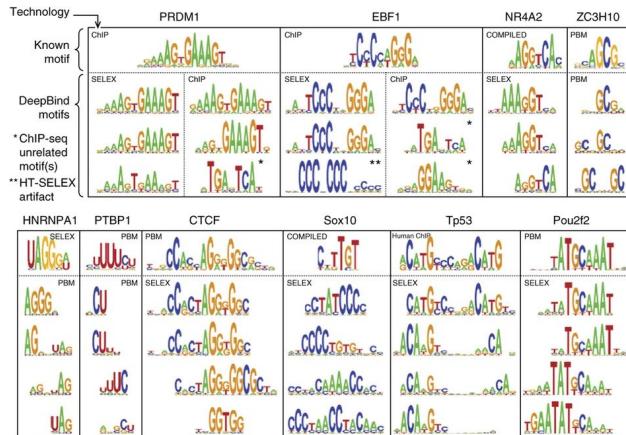
- ❑ Transparency would enable experts to check the model prediction and use it or check it more
- ❑ **Interpretability:** “understanding the patterns in the data may be just as important as fitting the data”

Ching et al. 2018

# Interpretability of neural networks

## Some techniques:

- Visualizing filters in (the first layers of) CNNs
- Clustering high contribution regions to derive motifs
- Backpropagating w.r.t. input data (saliency maps)
- Visualizing class activation maps



# Interpretable ML for computational biology

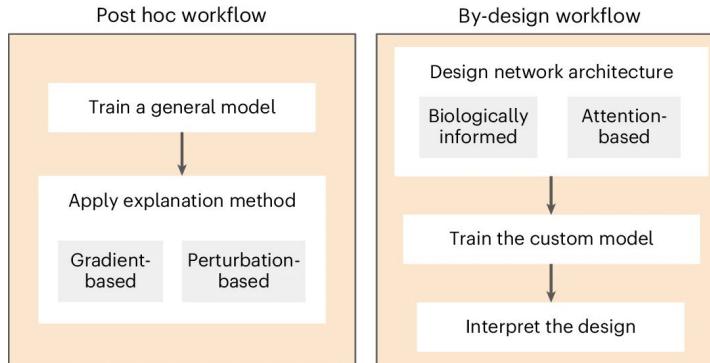
Perspective | Published: 09 August 2024

## Applying interpretable machine learning in computational biology—pitfalls, recommendations and opportunities for new developments

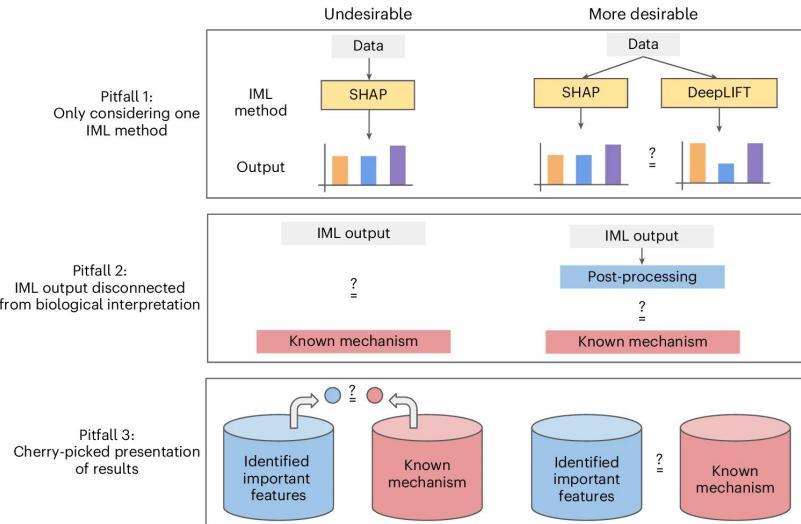
Valerie Chen, Muyu Yang, Wenbo Cui, Joon Sik Kim, Ameet Talwalkar & Jian Ma

Nature Methods 21, 1454–1461 (2024) | [Cite this article](#)

### Approaches:



### Common pitfalls:



# Reproducibility in ML

Reproducibility: “the provision of enough detail about study procedures and data so the same procedures could be exactly repeated”

Goodman et al. 2016

Achieving reproducibility:

- ❑ Keep track how each result is produced
- ❑ Record all intermediate results [in standardized formats]
- ❑ Store raw data behind all plots
- ❑ Provide public access to scripts, runs and results
- ❑ Archive the exact versions of all external programs used
- ❑ Use version control

Sandve et al. 2013

## Reproducibility standards for machine learning in the life sciences

[Benjamin J. Heil](#), [Michael M. Hoffman](#), [Florian Markowetz](#), [Su-In Lee](#), [Casey S. Greene](#)  & [Stephanie C. Hicks](#) 

*Nature Methods* **18**, 1132–1135 (2021) | [Cite this article](#)

	Bronze	Silver	Gold
Data published and downloadable	x	x	x
Models published and downloadable	x	x	x
Source code published and downloadable	x	x	x
Dependencies set up in a single command		x	x
Key analysis details recorded		x	x
Analysis components set to deterministic		x	x
Entire analysis reproducible with a single command			x

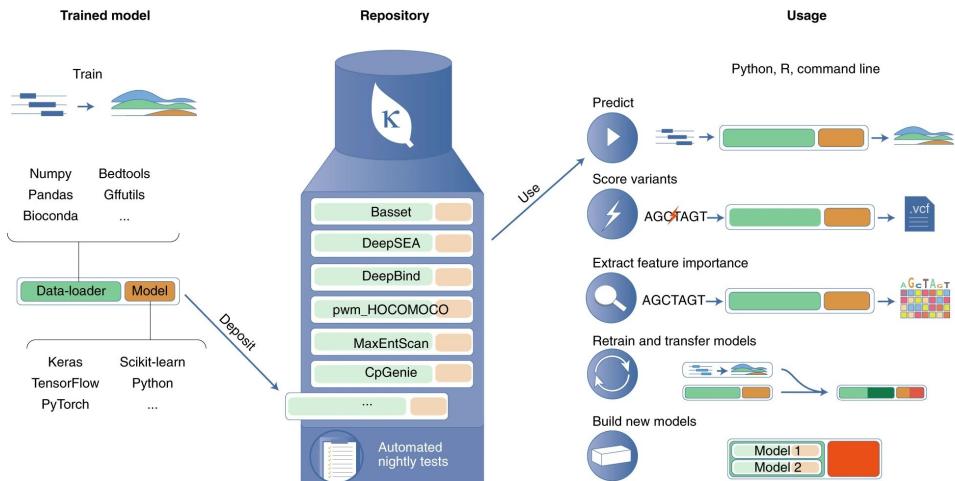
# Reproducibility in ML

- ❑ The code and models should be publicly available (GitHub, Zenodo, or other services)
- ❑ Possible to inspect, validate, (re)use, and improve models

**The Kipoi repository accelerates community exchange and reuse of predictive models for genomics**

Žiga Avsec Roman Kreuzhuber, Johnny Israeli, Nancy Xu, Jun Cheng, Avanti Shrikumar, Abhimanyu Banerjee, Daniel S. Kim, Thorsten Beier, Lara Urban, Anshul Kundaje , Oliver Stegle & Julien Gagneur

*Nature Biotechnology* 37, 592–600 (2019) | [Cite this article](#)



# References

Lipton ZC. The Mythos of Model Interpretability. *arXiv:160603490 [cs, stat]*. Published online March 6, 2017. <http://arxiv.org/abs/1606.03490>

Alipanahi, B., Delong, A., Weirauch, M. et al. Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning. *Nat Biotechnol* 33, 831–838 (2015). <https://doi.org/10.1038/nbt.3300>

Avsec, Ž., Weilert, M., Shrikumar, A. et al. Base-resolution models of transcription-factor binding reveal soft motif syntax. *Nat Genet* 53, 354–366 (2021). <https://doi.org/10.1038/s41588-021-00782-6>

Simonyan, Karen, Andrea Vedaldi, and Andrew Zisserman. ‘Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps’. *ArXiv:1312.6034 [Cs]*, 19 April 2014. <http://arxiv.org/abs/1312.6034>.

Selvaraju, Ramprasaath R., Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. ‘Grad-CAM:Visual Explanations from Deep Networks via Gradient-Based Localization’. In *2017 IEEE International Conference on Computer Vision (ICCV)*, 618–26, 2017. <https://doi.org/10.1109/ICCV.2017.74>.

Goodman SN, Fanelli D, Ioannidis JPA. What does research reproducibility mean? *Science Translational Medicine*. 2016;8(341):341ps12-341ps12. doi:[10.1126/scitranslmed.aaf5027](https://doi.org/10.1126/scitranslmed.aaf5027)

Sandve GK, Nekrutenko A, Taylor J, Hovig E. Ten Simple Rules for Reproducible Computational Research. *PLOS Computational Biology*. 2013;9(10):e1003285. doi:[10.1371/journal.pcbi.1003285](https://doi.org/10.1371/journal.pcbi.1003285)

Heil, B.J., Hoffman, M.M., Markowetz, F. et al. Reproducibility standards for machine learning in the life sciences. *Nat Methods* 18, 1132–1135 (2021). <https://doi.org/10.1038/s41592-021-01256-7>

Avsec, Ž., Kreuzhuber, R., Israeli, J. et al. The Kipoi repository accelerates community exchange and reuse of predictive models for genomics. *Nat Biotechnol* 37, 592–600 (2019). <https://doi.org/10.1038/s41587-019-0140-0>