# Machine learning – introduction

### Till Bärnighausen

DS-I Africa Short course UKZN, Durban 23 January 2023



HEIDELBERG H UNIVERSITY FA HOSPITAL M









# **DS-LAfrica** Data Science for Health Discovery and Innovation in Africa

# Machine learning is a type of AI



AI = artificial intelligence, ML = machine learning, DL = deep learning

### **MSNBC**

With layoffs, tech companies are quickening the robot revolution

# With layoffs, tech companies are quickening the robot revolution

Google, Microsoft and others in Big Tech have announced massive job cuts, as the companies pivot to artificial intelligence projects. What does it mean?



Jan. 21, 2023, 2:36 AM EAT **By Ja'han Jones** 

The robot revolution appears to be in full swing.

Google is just the latest company to announce major job cuts said to align with its prioritization of artificial intelligence.

Industries far and wide are pivoting to AI, which – shameless plug – you may have anticipated if you read my end-of-year ReidOut Blog post on 2023 being an important year for the development of artificial intelligence technology.

our course located? **A FRAMEWORK OF MACHINE** LEARNING **METHODS** 



# **Clustering analyses can identify types**

#### **OVERVIEW**

- Centroid k means, k medoids, k prototypes ...
- *Distribution* EM, GMM, BMM, ...
- Connectivity hierarchical, ...
- **Density** DBSCAN, ADBSCAN, HDBSCAN, ...

EM = expectation-maximization, GMM = Gaussian mixture model, BMM = Bernoulli mixture model, DBSCAN = density-based spatial clustering of applications with noise, ADBSCAN = adaptive DBSCAN, HDBSCAN = hierarchical DBSCAN

# Dimension reduction can reduce complexity and identify latent constructs and

**OVERVIEW** 

- Linear PCA, SVM, LDA, ...
- Non-linear kernel PCA, FAMD, t-SNE, ...

PCA = principal component analysis, SVM = support vector machine, LDA = linear discrimnant analysis, FAMD = factor analysis for mixed data, t-SNE = t-distributed stochastic neighbor embedding

# Feature selection can reduce complexity and cost

**OVERVIEW** 

- Supervised
  - Filter information gain, correlation, chi-squared, ...
  - Wrapper forward, backward, stepwise, exhaustive
  - *Embedded* RF, regularization, ...
- **Unsupervised** variance, multicollinearity, incompleteness, ...

Causal analyses is the traditional mainstay of				
epidemiology overview		Control of unobserved confounding	Assumptions	
Experiments		Complete	Weak	
Quasi-experiments				
<ul> <li>Instrumental variable approaches</li> <li>Regressions discontinuity</li> </ul>	}	Complete	Less weak	
<ul> <li>Difference-in-differences approaches</li> <li>Fixed effects approaches</li> </ul>	}	Partial	Less weak	
<ul> <li>Non-experiments</li> </ul>				
<ul> <li>Regression</li> <li>Matching</li> <li>Stratification</li> </ul>	}	None	Strong	

Bärnighausen et al. Journal of Clinical Epidemiology 2017

## **Our course focuses on prediction**

#### **OVERVIEW**

- Regression
  - SVM
  - Decision trees
  - RF
  - kNN
  - Neural networks
  - ..

- Classification
  - kNN
  - Logistic regression
  - LDA
  - QDA

. . .

- Naïve Bayesian
- Neural networks

SVM = support vector machine, RF = random forest, kNN = k nearest neighbors, LDA = linear discriminant analysis, QDA = quadratic discriminant analysis

**Causal and** predictive analysis require very different approaches

COMPARISON

PSM = propensity score matching, IV = instrumental variable analysis, RDD = regression discontinuity design, kNN = k nearest neighbors

	Causation	Prediction	
Disciplines	<ul><li>Epidemiology</li><li>Economics</li></ul>	<ul><li>Machine learning</li><li>Computer science</li></ul>	
Foundation	<ul><li>Theory-based</li><li>Hypothesis testing</li></ul>	<ul> <li>Data-driven</li> </ul>	
Purposes	<ul><li>Understanding</li><li>Policy guidance</li><li>Regulatory approval</li></ul>	<ul><li>Intervention targeting</li><li>Intervention tailoring</li><li>Now- and forecasting</li></ul>	
Goal of approach	Minimize/eliminate bias	<ul> <li>Optimize bias-variance trade-off</li> </ul>	
Approach	Estimation	<ul> <li>Training-(validation)- testing</li> <li>Complexity reduction</li> </ul>	
Example	<ul> <li>Ordinary multiple regression</li> <li>PSM</li> <li>IV and RDD</li> </ul>	<ul> <li>Regularized multiple regression</li> <li>kNN</li> <li>Neural networks</li> </ul>	

For best prediction, we trade-off bias and variance

Source: https://www.geeksfor geeks.org/ml-biasvariance-trade-off/



# How can we reduce overfitting in ML?

#### APPROACHES

#### • Data

- Adding training data
- Data augmentation
- Adding noise

#### Features

- Feature selection
- Dimension reduction
- Regularization (Ridge, Lasso)

#### Algorithm-specific

- Pruning (trees)
- Bagging of weakly correlated trees (forests)
- Dropping out layers (neural networks)





High training error High test error Low training error Low test error Overfit (high variance)



Low training error High test error

Source:

https://www.linkedin.com/pulse/overfittingunderfitting-machine-learning-ml-concepts-com



### There are excellent introductory textbooks ...

#### **EXAMPLES**



#### **Basic to intermediary**

Springer Series in Statistics Trevor Hastie Robert Tibshirani Jerome Friedman **The Elements of Statistical Learning** Data Mining, Inference, and Prediction

Second Edition

🖄 Springer

... and online resources for machine learning







# In this short course, we will use the R programming language

**PROGRAMMING LANGUAGES** 

- **Procedural** C, Java
- Functional R, Scala
- **Object-oriented** Python, Java
- *Scripting* Python, Ruby
- *Logic* Prolog

# In the coming two weeks, we will cover ML concepts, concrete methods, and applications

**COURSE OVERVIEW** 

- General approaches
  - Feature selection
  - Validation / cross-validation
  - Parameter tuning to reduce overfitting
  - Performance metrics

#### Concrete methods

- Regression / logistic regression
- Ridge and lasso regression
- kNN
- Naïve Bayes
- LDA / QDA
- -SVM
- Trees and forests
- Boosting

## In the coming two weeks, we will cover both ML general approaches and concrete predictive methods COURSE OVERVIEW

- General approaches
  - Feature selection
  - Validation / cross-validation
  - Parameter tuning to reduce overfitting
  - Performance metrics

#### Concrete predictive methods

- Regression / logistic regression
- Ridge and lasso regression
- kNN
- Naïve Bayes
- LDA/QDA
- SVM
- Trees and forests
- Boosting