

Using supervised machine learning algorithms for the building of predictive biomonitoring models.

9th Metabarcoding school
Villa de Leyva – 1-6th November 2019

Tristan Cordier

EFFECTS OF BIODIVERSITY ON ECOSYSTEM FUNCTIONING: A CONSENSUS OF CURRENT KNOWLEDGE

D. U. HOOPER,^{1,16} F. S. CHAPIN, III,² J. J. EWEL,³ A. HECTOR,⁴ P. INCHAUSTI,⁵ S. LAVOREL,⁶ J. H. LAWTON,⁷
D. M. LODGE,⁸ M. LOREAU,⁹ S. NAEEM,¹⁰ B. SCHMID,⁴ H. SETÄLÄ,¹¹ A. J. SYMSTAD,¹²
J. VANDERMEER,¹³ AND D. A. WARDLE^{14,15}

Biodiversity and ecosystem stability in a decade-long grassland experiment

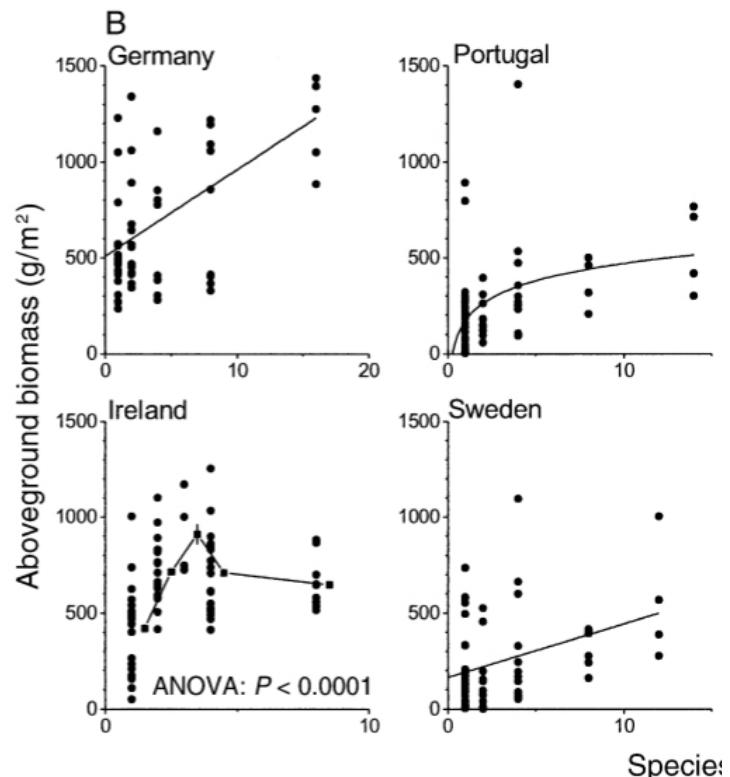
David Tilman¹, Peter B. Reich² & Johannes M. H. Knops³

REVIEW

doi:10.1038/nature11148

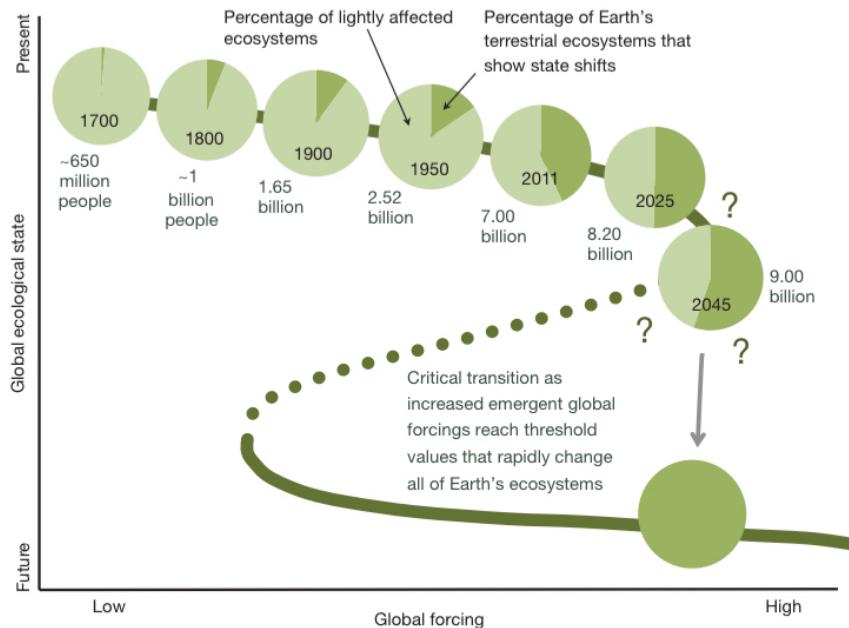
Biodiversity loss and its impact on humanity

Bradley J. Cardinale¹, J. Emmett Duffy², Andrew Gonzalez³, David U. Hooper⁴, Charles Perrings⁵, Patrick Venail¹, Anita Narwani¹, Georgina M. Mace⁶, David Tilman⁷, David A. Wardle⁸, Ann P. Kinzig⁵, Gretchen C. Daily⁹, Michel Loreau¹⁰, James B. Grace¹¹, Anne Larigauderie¹², Diane S. Srivastava¹³ & Shahid Naeem¹⁴

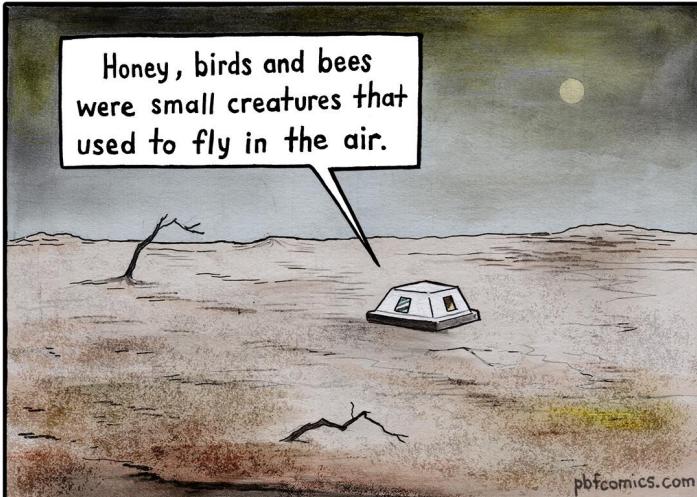
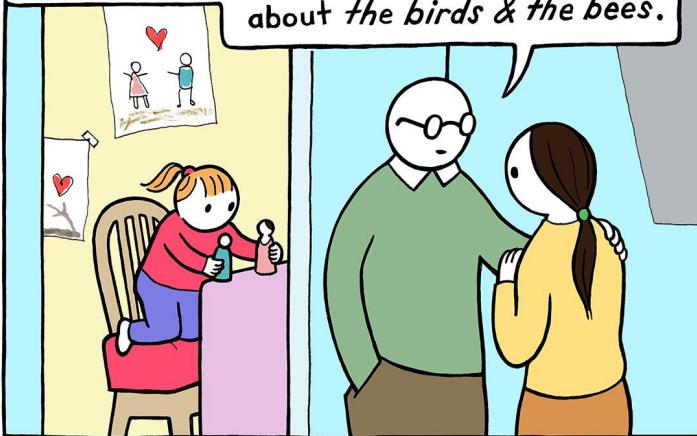


Approaching a state shift in Earth's biosphere

Anthony D. Barnosky^{1,2,3}, Elizabeth A. Hadly⁴, Jordi Bascompte⁵, Eric L. Berlow⁶, James H. Brown⁷, Mikael Fortelius⁸, Wayne M. Getz⁹, John Harte^{9,10}, Alan Hastings¹¹, Pablo A. Marquet^{12,13,14,15}, Neo D. Martinez¹⁶, Arne Mooers¹⁷, Peter Roopnarine¹⁸, Geerat Vermeij¹⁹, John W. Williams²⁰, Rosemary Gillespie⁹, Justin Kitzes⁹, Charles Marshall^{1,2}, Nicholas Matzke¹, David P. Mindell²¹, Eloy Revilla²² & Adam B. Smith²³



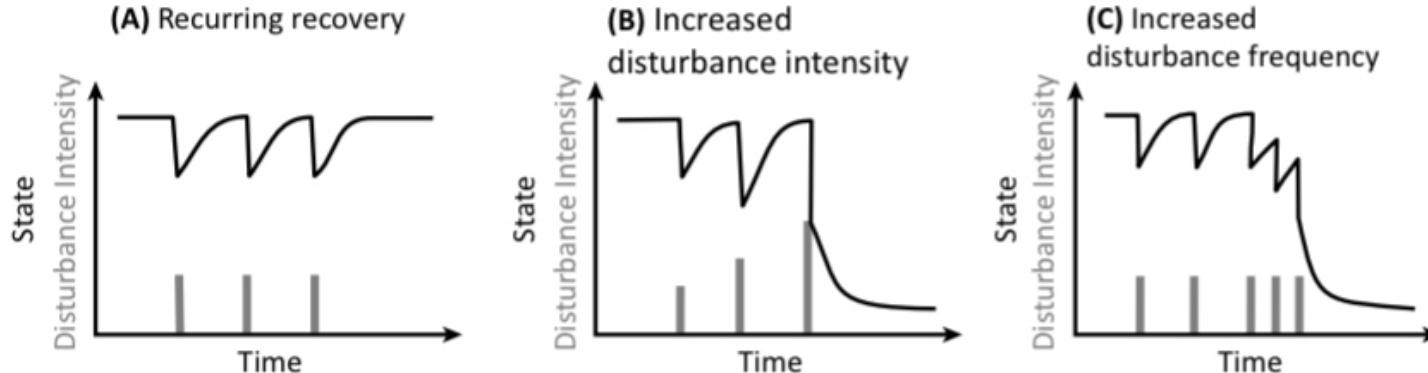
It's about time we had a talk with her
about the birds & the bees.



Review

Abrupt Change in Ecological Systems: Inference and Diagnosis

Zak Ratajczak,^{1,*} Stephen R. Carpenter,² Anthony R. Ives,¹ Christopher J. Kucharik,³
Tanjona Ramiadantsoa,¹ M. Allison Stegner,¹ John W. Williams,⁴ Jien Zhang,¹ and
Monica G. Turner^{1,*}



LETTERS

Biodiversity and ecosystem multifunctionality

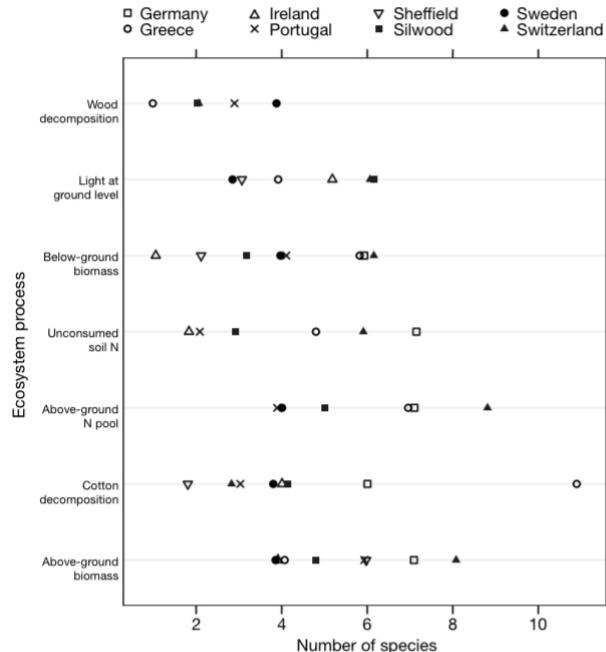
Andy Hector¹ & Robert Bagchi¹

Figure 1 | Number of species with desirable effects on the suite of ecosystem processes measured in the different BIODEPTH project experiments. The number of species was identified by the AIC-based multiple regression (and species with effects with undesirable signs were then excluded).

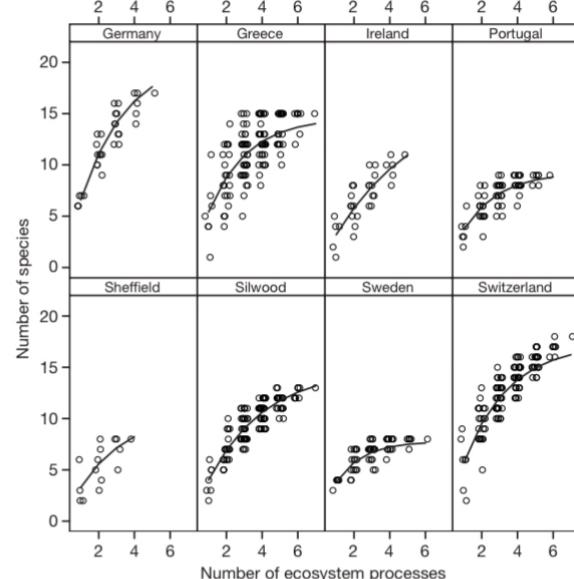
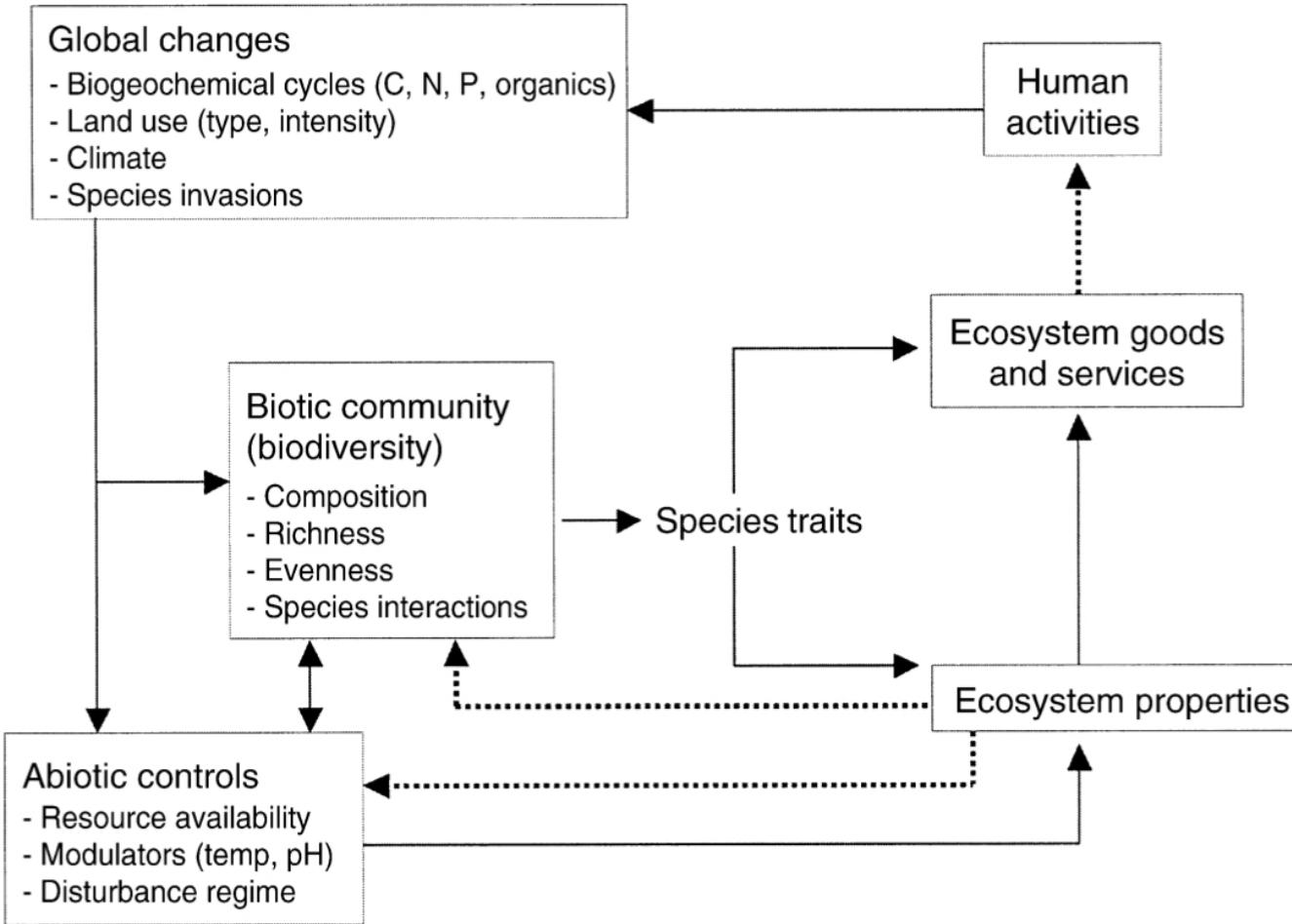
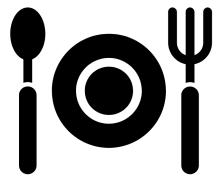
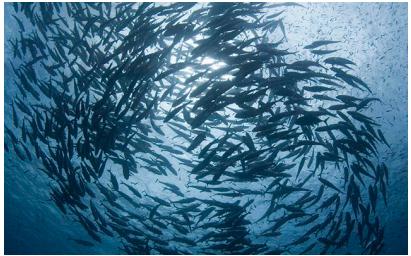


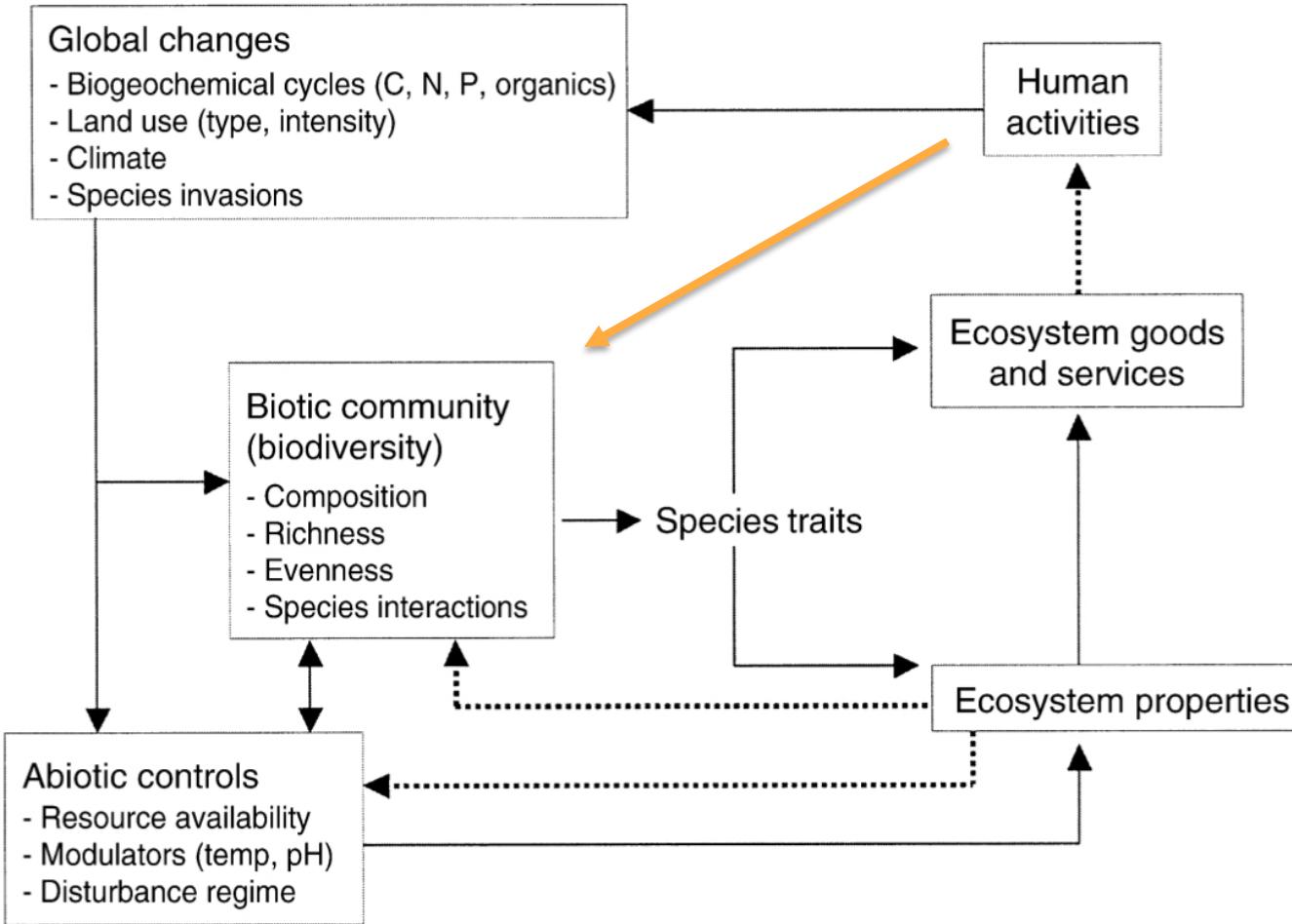
Figure 2 | Positive relationship between the range of ecosystem processes considered and the number of species that affect one or more aspect of ecosystem functioning. The points (jittered for clarity) show numbers of species required for all possible combinations of ecosystem processes. Lines are theoretical predictions from the model based on the average number of species required for a single process, \bar{x} , and the average overlap in the sets of species required for each pair of processes, \bar{o} , using equation (2).



(Marine) ecosystems services

- Sustainable development !





WHAT CONSTITUTES ECOSYSTEM HEALTH?

DAVID J. RAPPORT* *Perspectives in Biology and Medicine*, 33, 1 · Autumn 1989 | 121

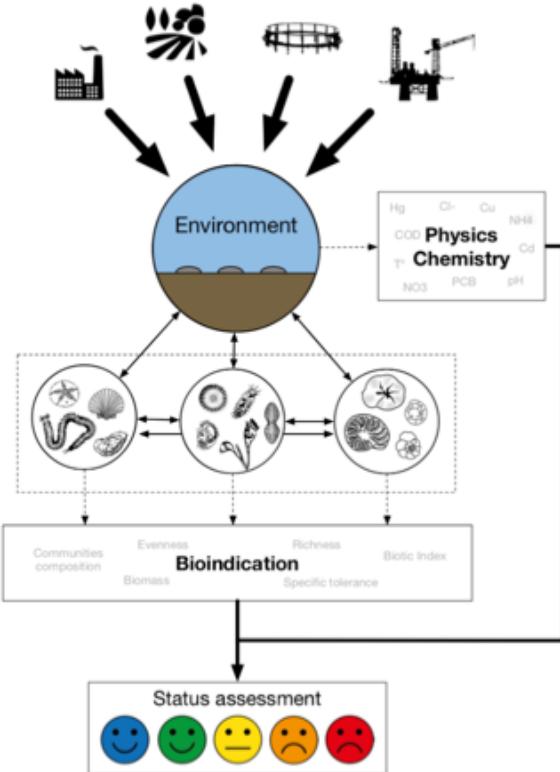
The Use of Metaphor in Science

In addressing the question of what constitutes ecosystem health, it is tempting to make use of metaphors drawn from human medicine. The use of metaphor, it has been argued, has not only a role in poetry, but also a legitimate place in science [1]. Its function in both areas is to stimulate associations, bringing into juxtaposition phenomena that might at first appear to have little connection but can be seen to be in some way related. In poetry the association might be quite fanciful, while in science the value of metaphor lies in pointing to phenomena in apparently different spheres that bear some structural identity.

What constitutes the health of Nature—that is, what are the suitable concepts and conventions to assess the condition of environment—is a question now being raised for the earth's major ecosystems (forests, lakes, seas, etc.) and indeed the entire biosphere. Ultimately, a healthy environment is essential for a healthy human population. It does not follow, however, that the appropriate standards for the health of nature need be based solely on criteria for human health. Ecosystems have a life of their own with or without human components, and it is this life that is receiving ever more attention as situations come to light in which ecosystems have become severely damaged [2].

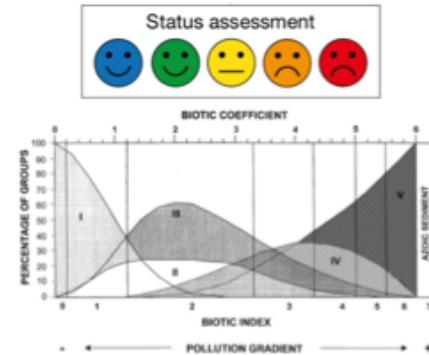
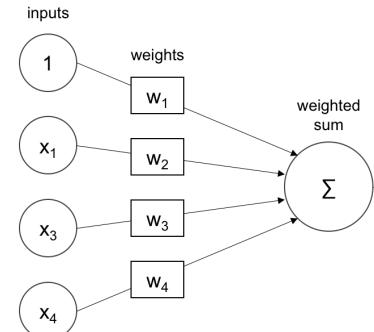


Impact assessment relies on biotic indices



AMBI index groups (Borja 2000)

- I – very sensitive species
- II – indifferent species
- III – tolerant species
- IV – second order opportunistic species
- V – first order opportunistic species



Bioindicators convey a cumulative measure of disturbances

Impact assessment relies on biotic indices

Ecological Applications, 9(2), 1999, pp. 699–713
© 1999 by the Ecological Society of America

THE INFANAL TROPHIC INDEX (ITI): ITS SUITABILITY FOR MARINE ENVIRONMENTAL MONITORING

DON MAURER,¹ HAI NGUYEN,² GEORGE ROBERTSON,² AND TOM GERLINGER²

¹Department of Biological Sciences, California State University, 1250 Bellflower Boulevard, Long Beach, California 90840 USA

²County Sanitation Districts of Orange County California, 10844 Ellis Avenue, Fountain Valley, California 92708-7018 USA

Biotic indices based on soil nematode communities for assessing soil quality in terrestrial ecosystems

Arantzazu Urzelai^a, Ana Jesús Hernández^b, Jesús Pastor^{a,*}

^aCentre of Environmental Sciences, CSIC, Serrano 115 dpto., E-28006 Madrid, Spain

^bArea of Ecology, University of Alcalá de Henares, Ctra. Madrid-Barcelona km 38, E-28771 Alcalá de Henares, Spain



Contents lists available at ScienceDirect

Marine Pollution Bulletin

journal homepage: www.elsevier.com/locate/marpolbul



A bacterial community-based index to assess the ecological status of estuarine and coastal environments

Eva Aylagas^{a,*}, Ángel Borja^{a,*}, Michael Tangherlini^b, Antonio Dell'Anno^b, Cinzia Corinaldesi^b, Craig T. Michell^c, Xabier Irigoien^{a,c,d}, Roberto Danovaro^b, Naaira Rodríguez-Ezpeleta^a

^a AZTI - Marine Research, Herrera Kalz, Portugalete 216 – 20110 Pasaia, Gipuzkoa, Spain

^b Department of Life and Environmental Sciences, Polytechnic University of Marche, Via Brecce Bianche, 60131 Ancona, Italy

^c Red Sea Research Center, King Abdullah University of Science and Technology (KAUST), Thuwal 23955-6900, Saudi Arabia

^d IKERBASQUE, Basque Foundation for Science, Bilbao, Spain



Pergamon

Marine Pollution Bulletin Vol. 40, No. 12, pp. 1100–1114, 2000

© 2000 Elsevier Science Ltd. All rights reserved

Printed in Great Britain

0025-326X/00 \$ - see front matter

PII: S0025-326X(00)00061-8

A Marine Biotic Index to Establish the Ecological Quality of Soft-Bottom Benthos Within European Estuarine and Coastal Environments

A. BORJA*, J. FRANCO and V. PÉREZ

Department of Oceanography and Marine Environment, Technological Institute for Fisheries and Food (AZTI), Av. Satrústegui 8, 20008 San Sebastián, Spain

Impact assessment relies on biotic indices

Ecological Applications, 9(2), 1999 pp. 699–713
© 1999 by the Ecological Society



Contents lists available at ScienceDirect



Ecological Indicators 18 (2012) 31–41

THE INFAUNAL

DON MA

¹Department

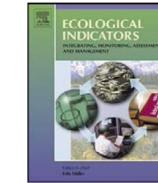
²Coun



Contents lists available at SciVerse ScienceDirect

Ecological Indicators

journal homepage: www.elsevier.com/locate/ecolind



I status of



zia Corinaldesi ^b, Craig T. Michell ^c,

Three hundred ways to assess Europe's surface waters: An almost complete overview of biological methods to implement the Water Framework Directive

Sebastian Birk ^{a,*}, Wendy Bonne ^b, Angel Borja ^c, Sandra Brucet ^b, Anne Courrat ^d, Sandra Poikane ^b, Angelo Solimini ^e, Wouter van de Bund ^b, Nikolaos Zampoukas ^b, Daniel Hering ^a

^a University of Duisburg-Essen, Faculty of Biology, Applied Zoology/Hydrobiology, Universitätsstraße 5, 45141 Essen, Germany

^b European Commission, Joint Research Centre, Institute for Environment and Sustainability, 21027 Ispra, Italy

^c AZTI-Tecnalia, Marine Research Division, Herrera Kaia, Portualdea s/n, 20110 Pasaia, Spain

^d CEMAGREF - UR EPBX, 50 avenue de Verdun, 33612 Cestas, France

^e Department of Public Health and Infectious Diseases, Sapienza University of Rome, Roma, Italy

Pollution Bulletin Vol. 40, No. 12, pp. 1100–1114, 2000
© 2000 Elsevier Science Ltd. All rights reserved
Printed in Great Britain
0025-326X/00 \$ - see front matter

Biotic indice
assassin

Arantzaz

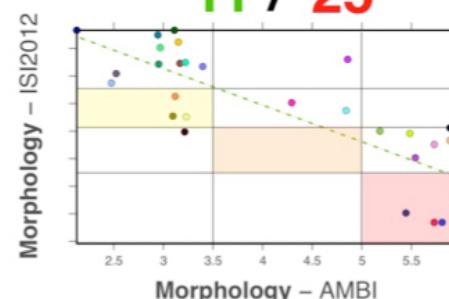
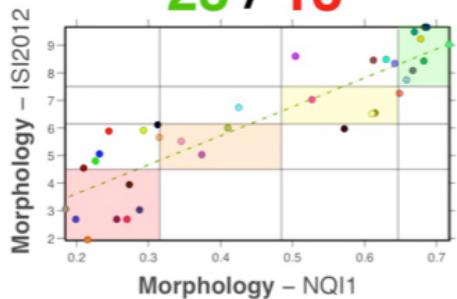
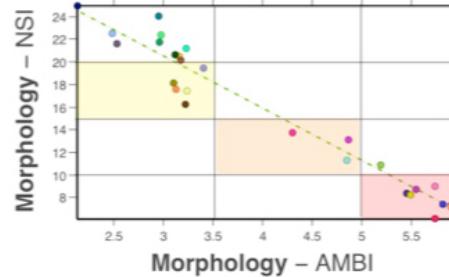
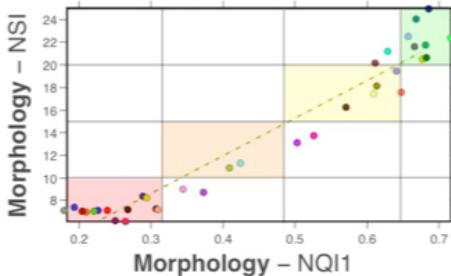
Establish the -Bottom Estuarine and

^a Centre of Environmental Sciences, CSIC, Serrano 115 dpto., E-28006 Madrid, Spain

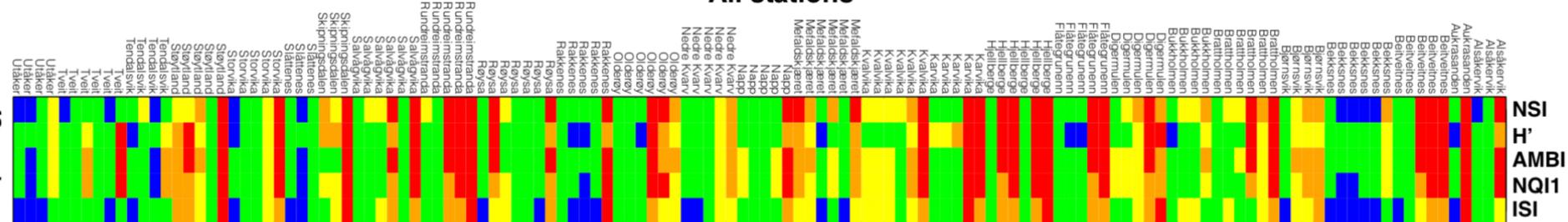
^b Area of Ecology, University of Alcalá de Henares, Ctra. Madrid-Barcelona km 38, E-28771 Alcalá de Henares, Spain

A. BORJA*, J. FRANCO and V. PÉREZ

Department of Oceanography and Marine Environment, Technological Institute for Fisheries and Food (AZTI), Av. Satrústegui 8, 20008 San Sebastián, Spain



All stations



EU Water Framework Directive (WFD) and Marine Strategy Framework Directive (MSFD)



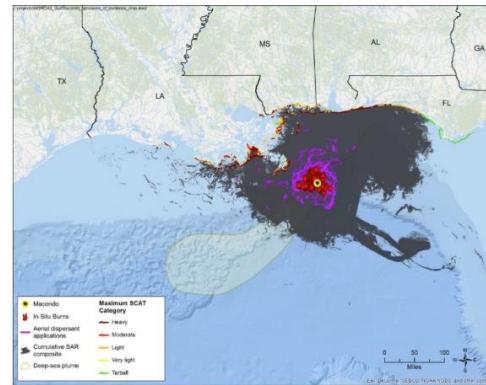
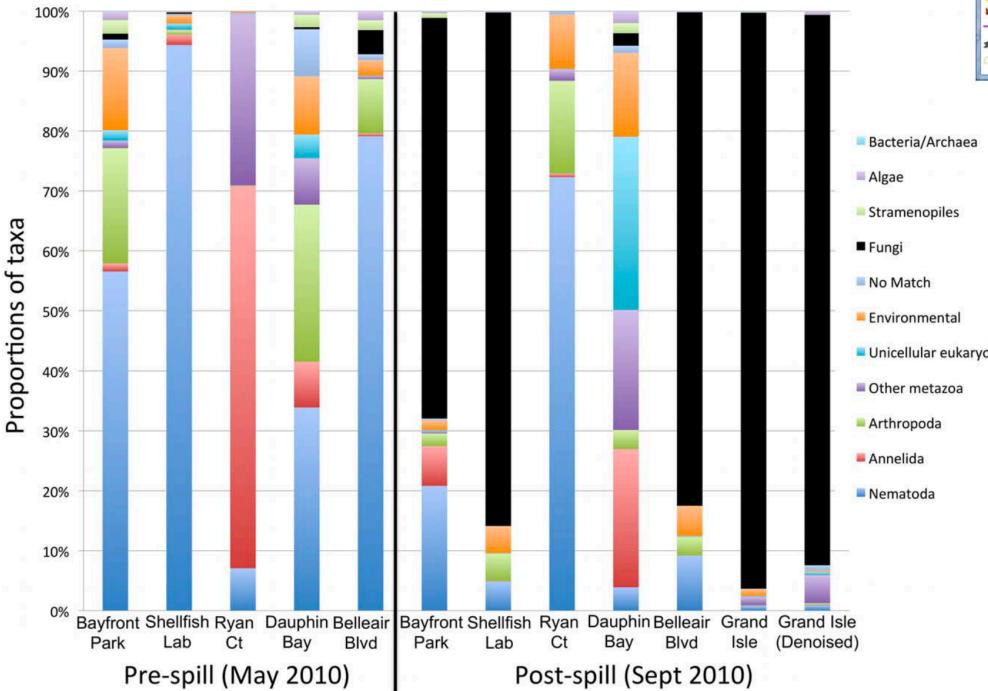
- Time consuming
- Only morphosp.
- Taxonomic skills



- Faster
- Full community
- Standardized
- Automatized

Dramatic Shifts in Benthic Microbial Eukaryote Communities following the Deepwater Horizon Oil Spill

Holly M. Bik^{1,2*}, Kenneth M. Halanych³, Jyotsna Sharma⁴, W. Kelley Thomas¹



OPEN

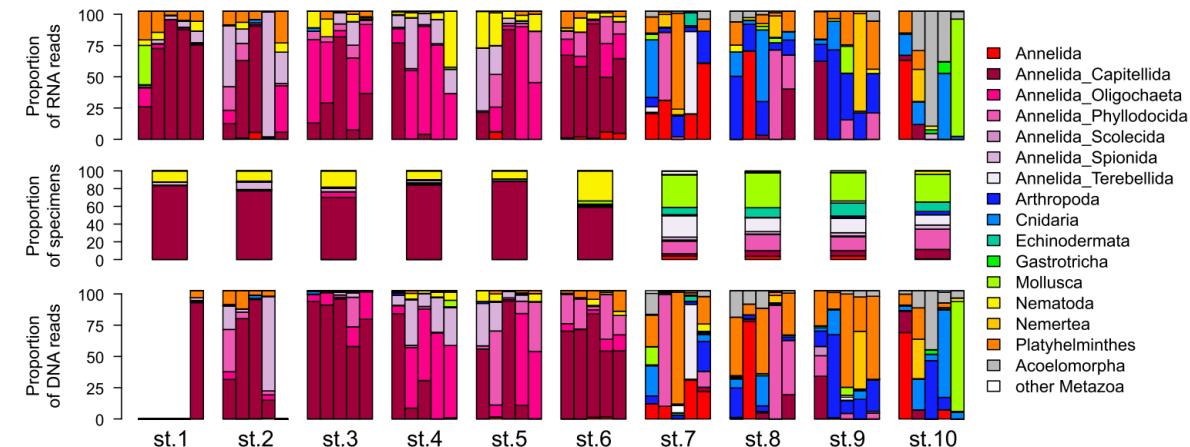
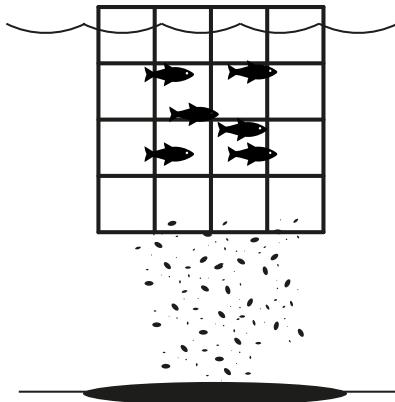
High-throughput sequencing and morphology perform equally well for benthic monitoring of marine ecosystems

Received: 01 April 2015

Accepted: 12 August 2015

Published: 10 September 2015

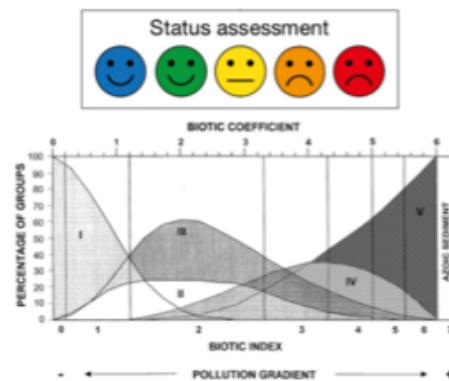
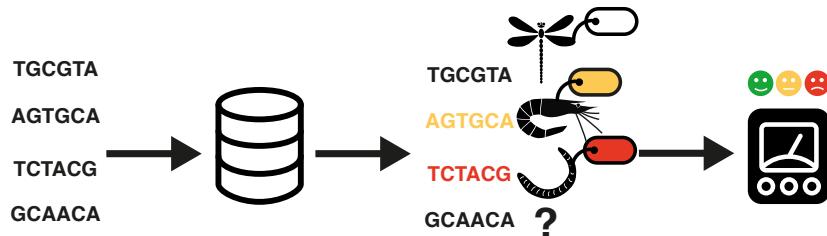
Franck Lejzerowicz^{1,†}, Philippe Esling^{1,2}, Loïc Pillet^{1,3}, Thomas A. Wilding⁴, Kenneth D. Black⁴ & Jan Pawłowski¹



A

Taxonomy-based approaches

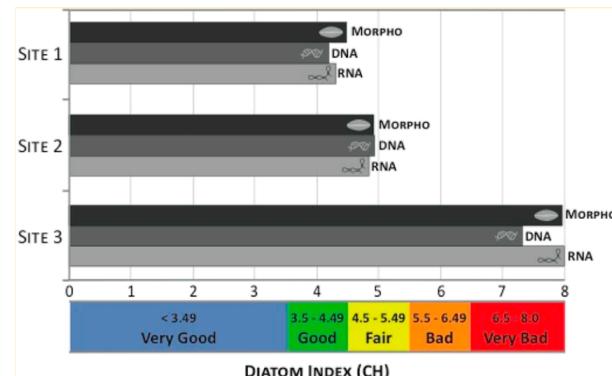
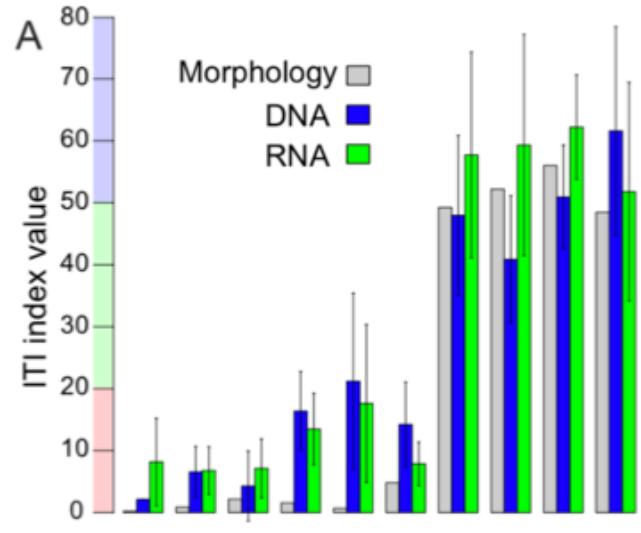
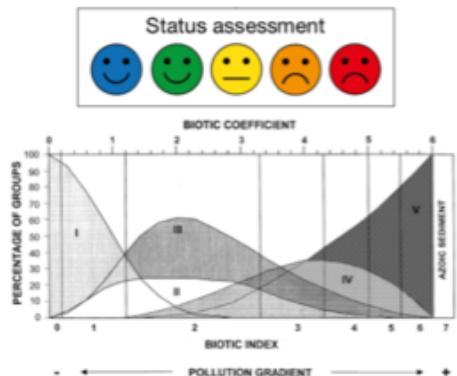
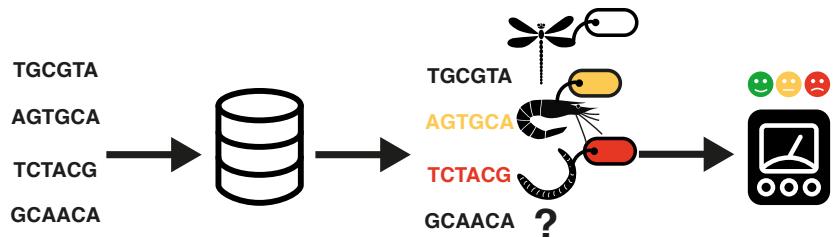
- Screening species / screening bioindicators



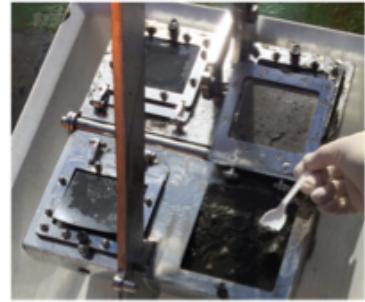
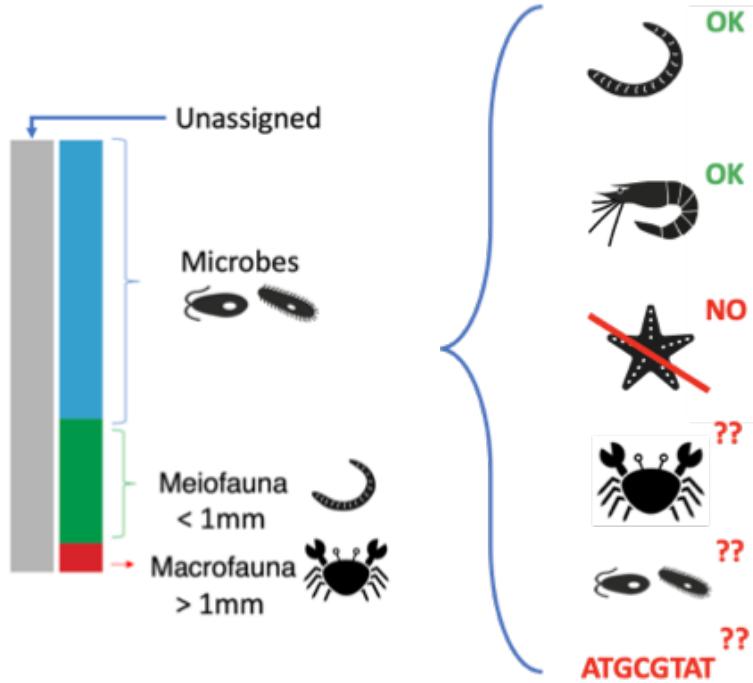
A

Taxonomy-based approaches

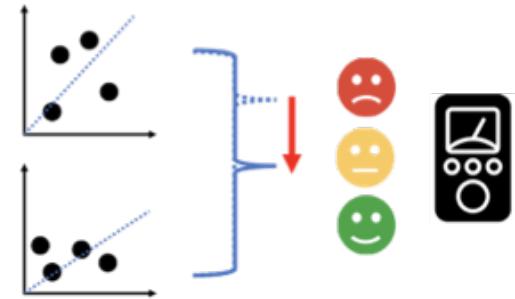
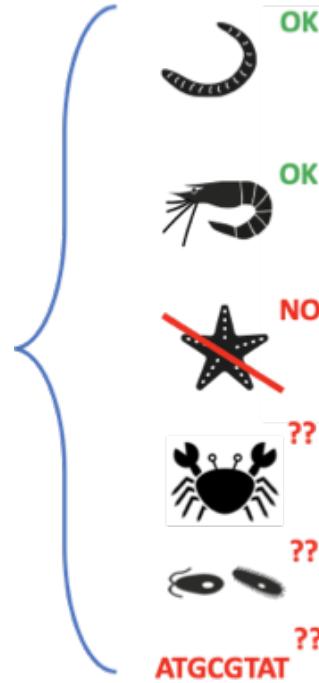
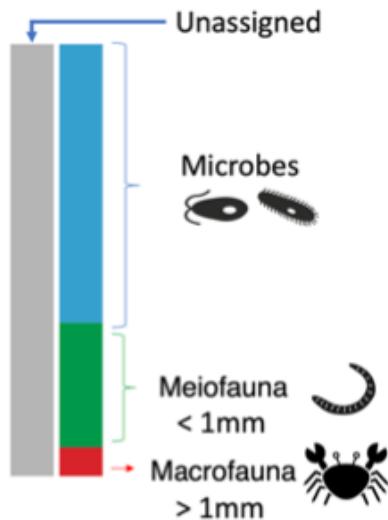
- Screening species / screening bioindicators



Sampling eDNA : What you actually get



Sampling eDNA : What you actually get

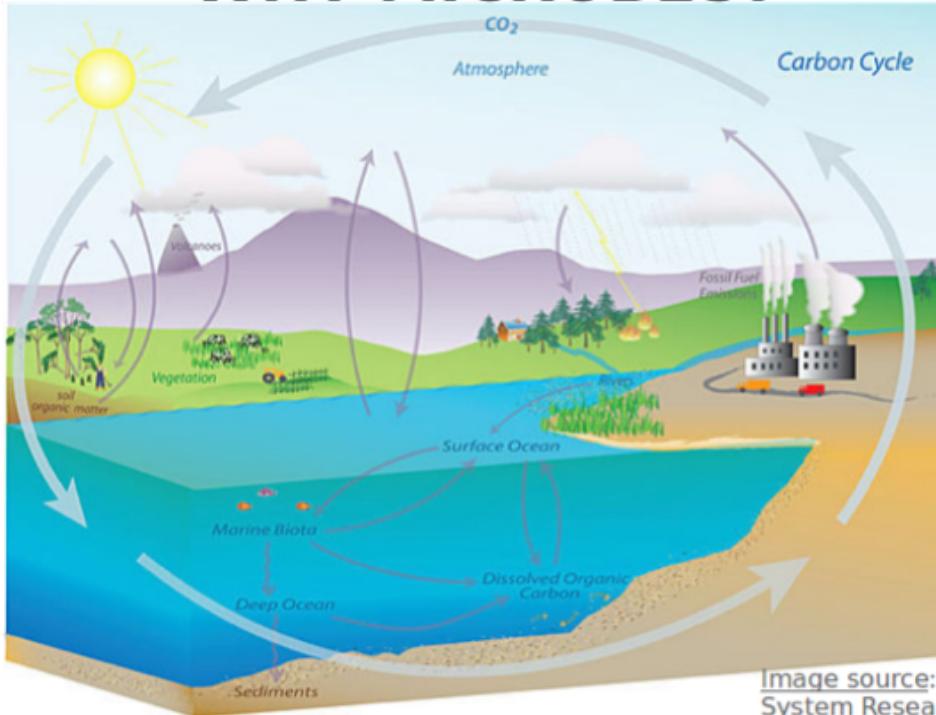


Reference DB? → Filling it!

Quantification? → Correction factor

Quest for the perfect PCR primers..

WHY MICROBES?



- Microbial life dominates
- Drivers of **global biogeochemical cycles**
- Sensitive and responsive to changing conditions

OPEN

Scientists' warning to humanity: microorganisms and climate change

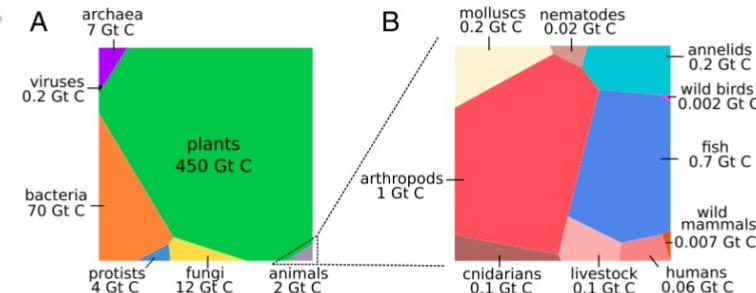
Ricardo Cavicchioli^{1*}, William J. Ripple², Kenneth N. Timmis³, Farooq Azam⁴, Lars R. Bakken⁵, Matthew Baylis¹⁶, Michael J. Behrenfeld⁷, Antje Boetius^{10, 8, 9}, Philip W. Boyd¹⁰, Aimée T. Classen¹¹, Thomas W. Crowther^{1, 2}, Roberto Danovaro^{1, 3, 14}, Christine M. Foreman¹⁵, Jef Huisman¹⁶, David A. Hutchins¹⁷, Janet K. Jansson^{10, 18}, David M. Karl¹⁹, Britt Koskella^{10, 20}, David B. Mark Welch^{10, 21}, Jennifer B. H. Martiny²², Mary Ann Moran^{10, 23}, Victoria J. Orphan²⁴, David S. Reay²⁵, Justin V. Remais^{10, 26}, Virginia I. Rich^{10, 27}, Brajesh K. Singh^{10, 28}, Lisa Y. Stein^{10, 29}, Frank J. Stewart³⁰, Matthew B. Sullivan^{10, 31}, Madeleine J. H. van Oppen^{10, 32, 33}, Scott C. Weaver³⁴, Eric A. Webb¹⁷ and Nicole S. Webster^{10, 33, 35}

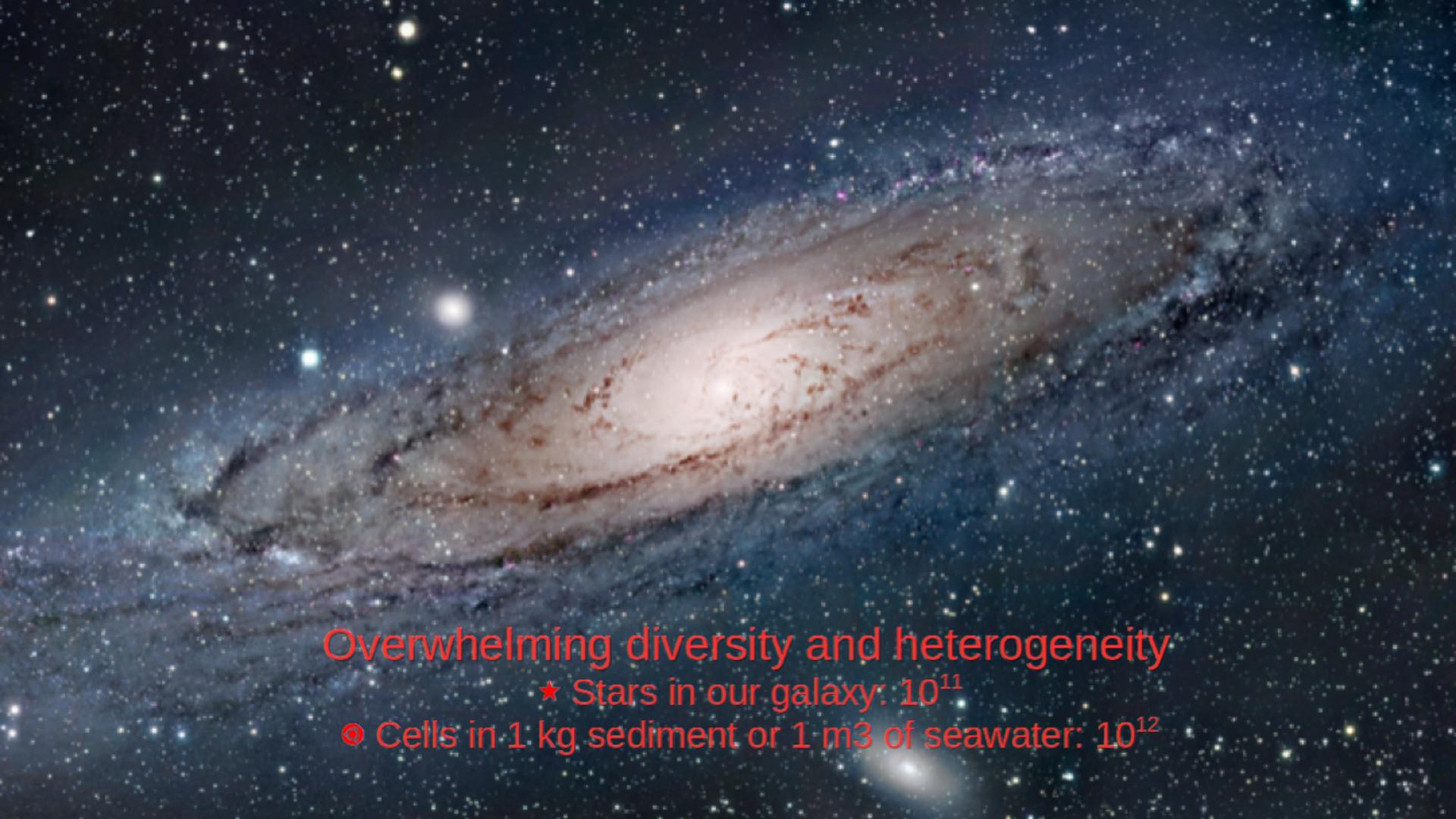
PNAS

The biomass distribution on Earth

Yinon M. Bar-On^a, Rob Phillips^{b,c}, and Ron Milo^{a,1}

Image source: NOAA Earth
System Research Laboratory,
Dept. of Commerce





Overwhelming diversity and heterogeneity

- ★ Stars in our galaxy: 10^{11}
- ◎ Cells in 1 kg sediment or 1 m³ of seawater: 10^{12}



Occurrence

... and the exploration ...
just started

Detection

Evolution

Time

Genes

Function

Molecule

Genomics

Activity

Composition

Long-term
analysis

Population
dynamics

Ecology

Structure

Autecology

Taxonomy

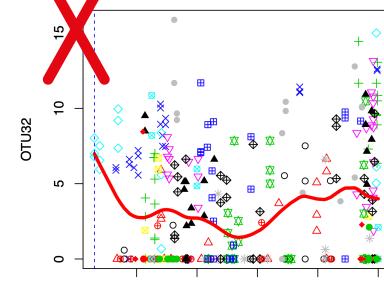
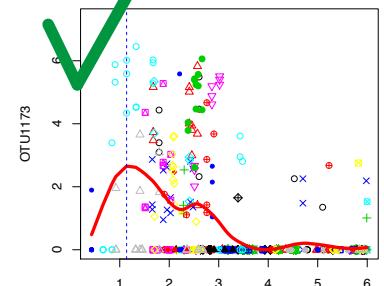
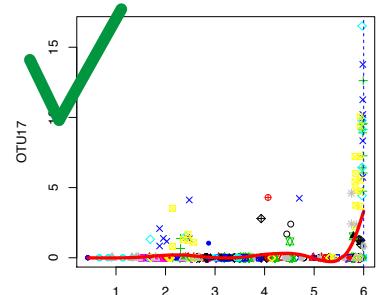
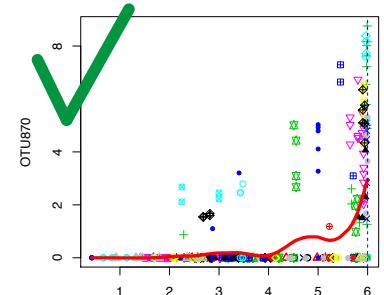
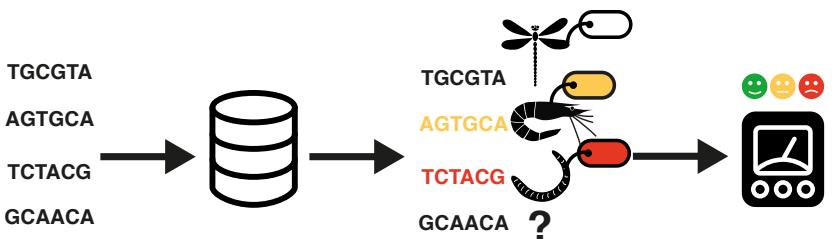
Abundance



A

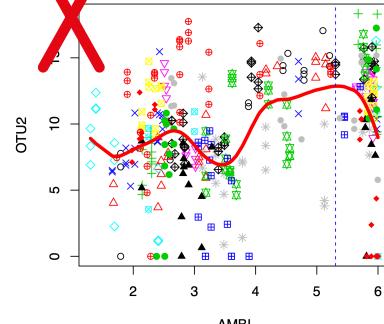
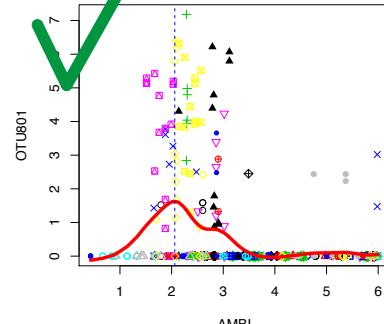
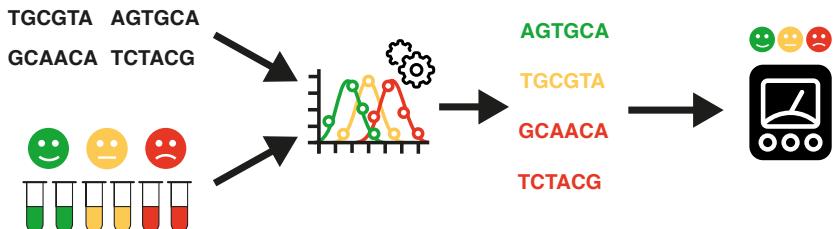
Taxonomy-based approaches

- Screening species / screening bioindicators

**B**

De novo approaches

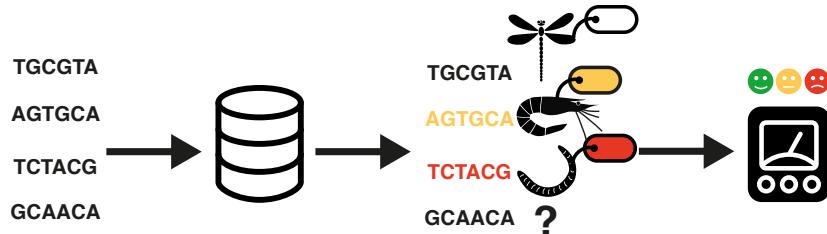
- Biondicators discovery / supervised learning



A

Taxonomy-based approaches

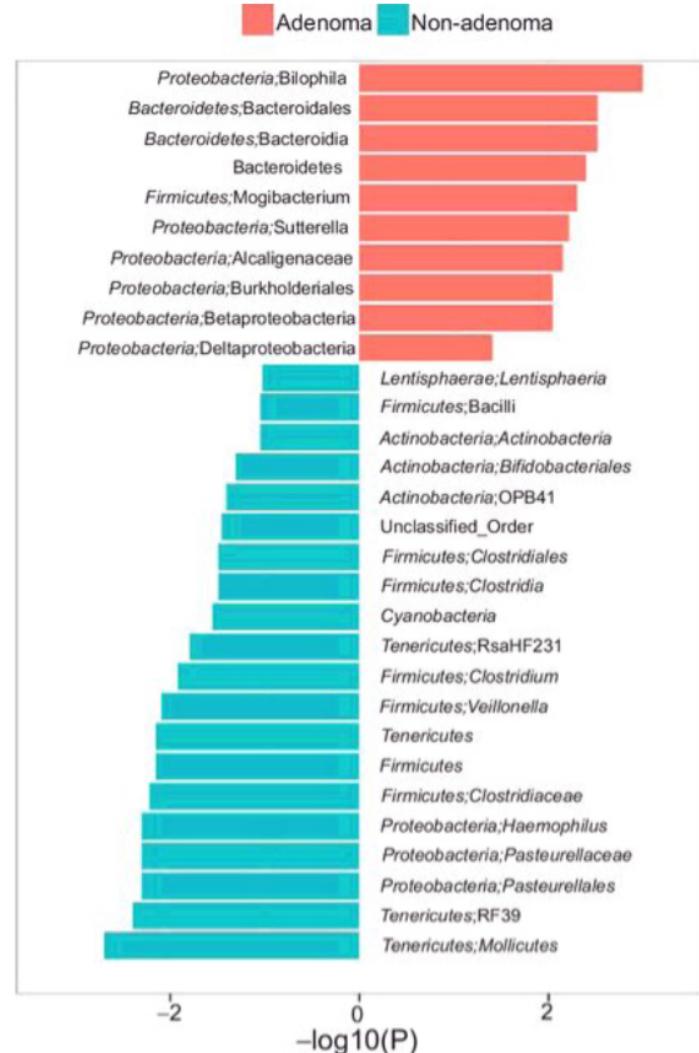
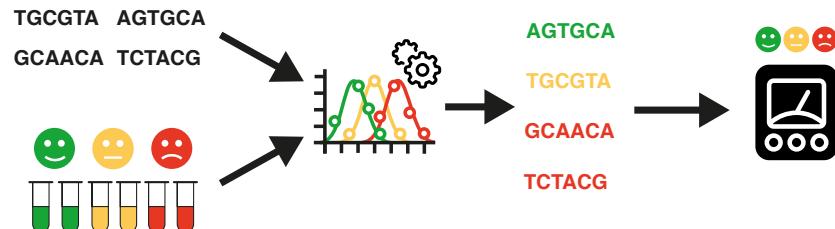
- Screening species / screening bioindicators



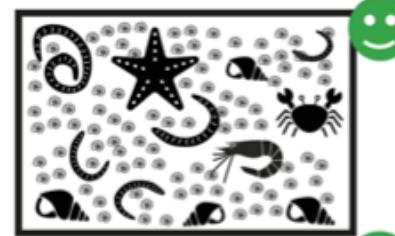
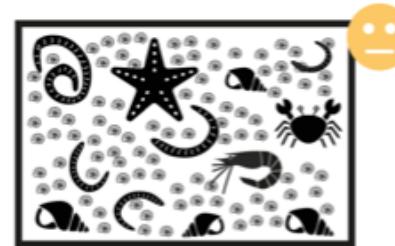
B

De novo approaches

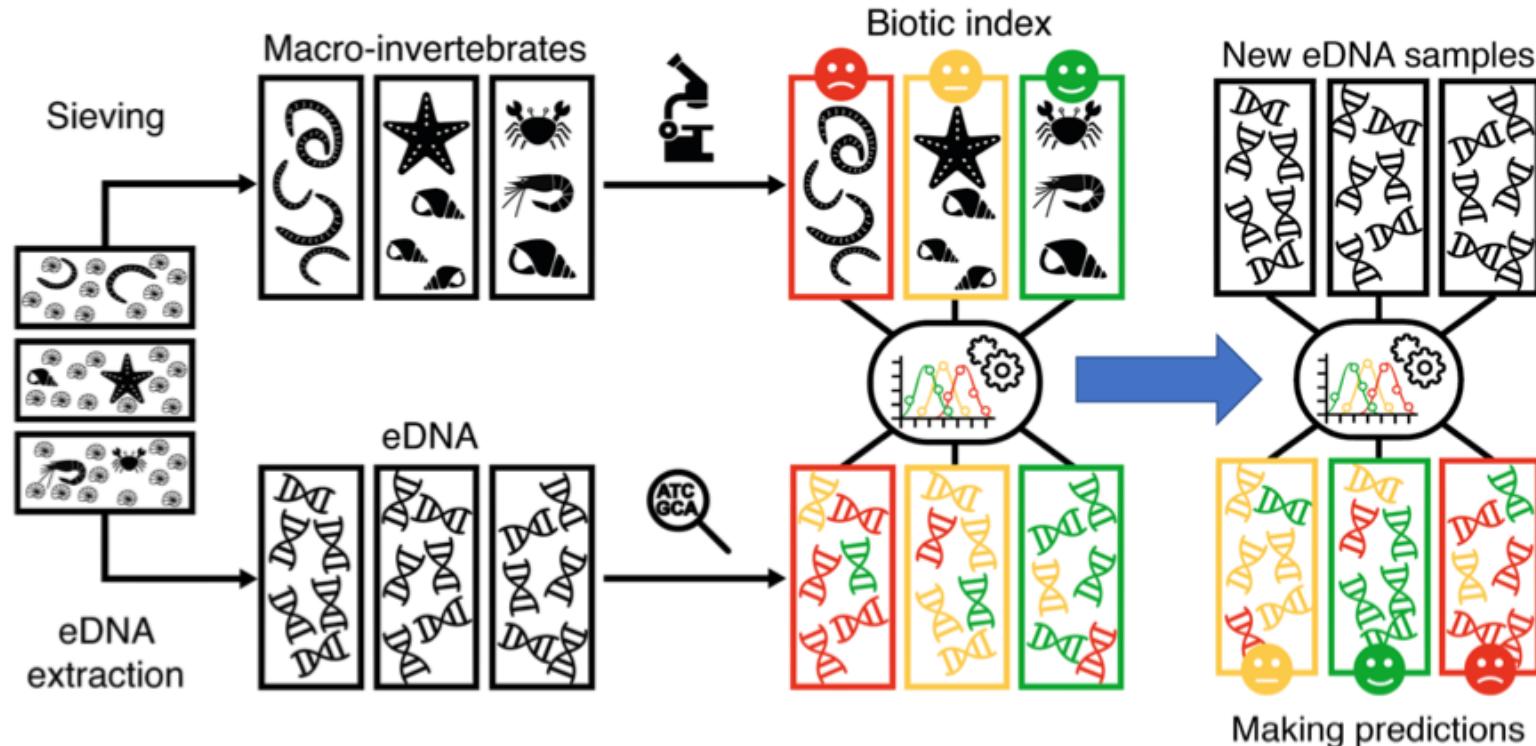
- Biondicators discovery / supervised learning



Biomonitoring: another classification problem → Machine Learning !



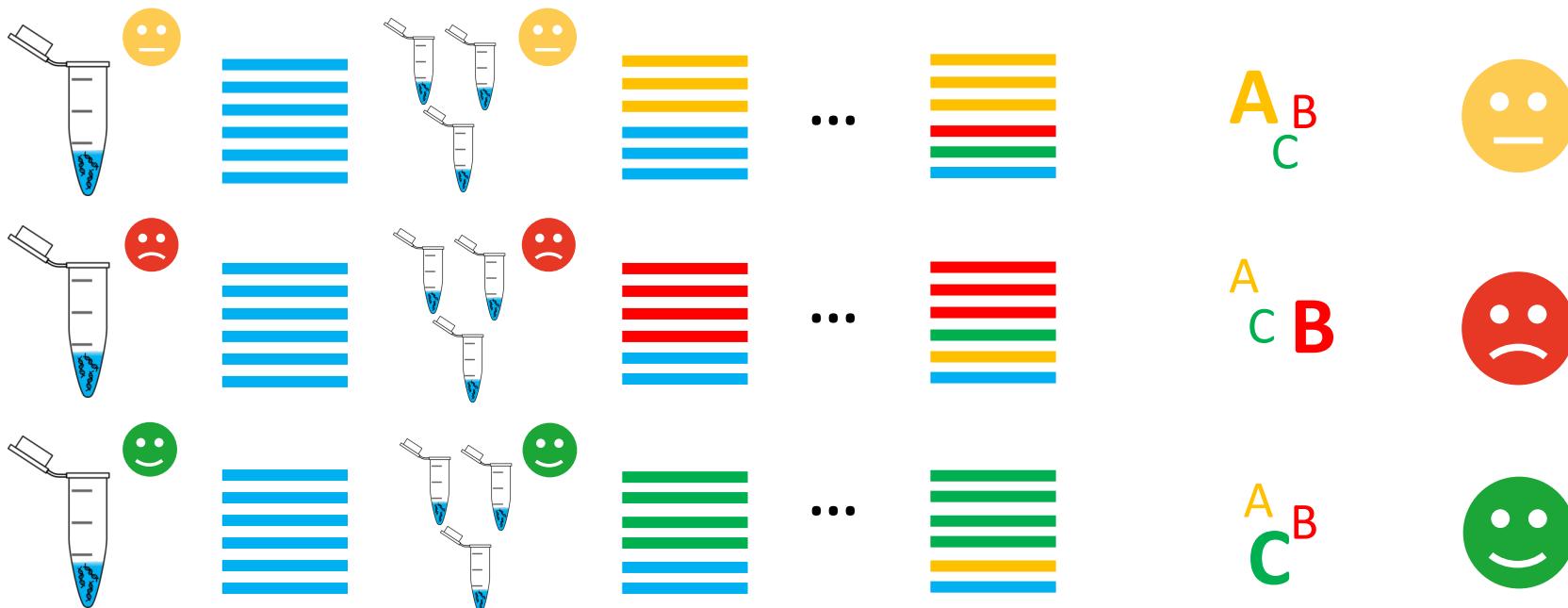
Biomonitoring: another classification problem → Machine Learning !



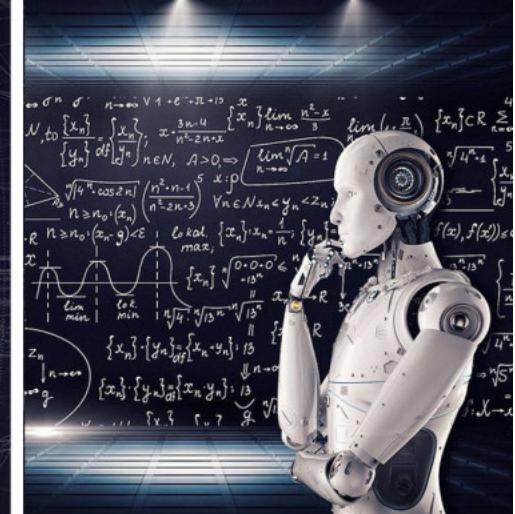
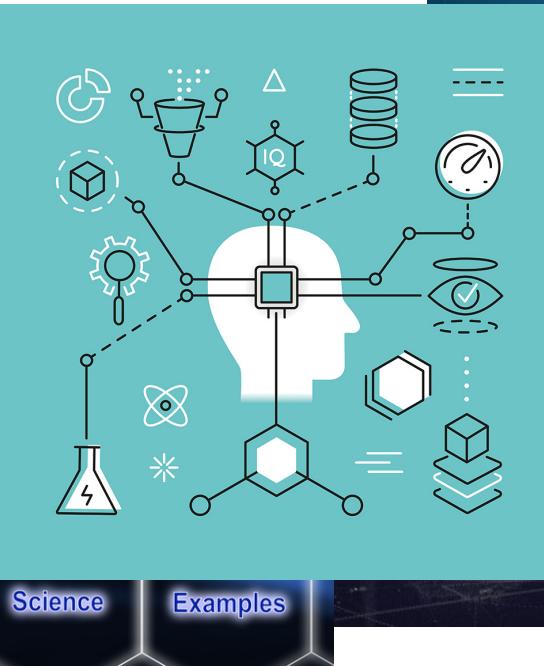
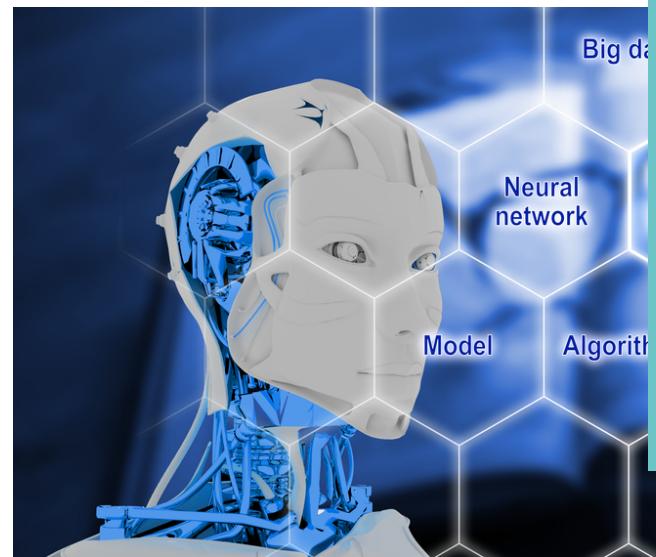
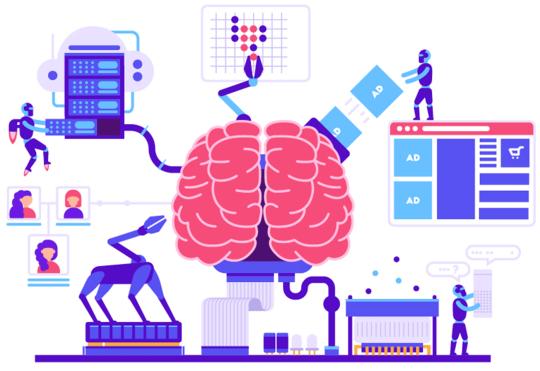
Biomonitoring: another classification problem → Machine Learning !

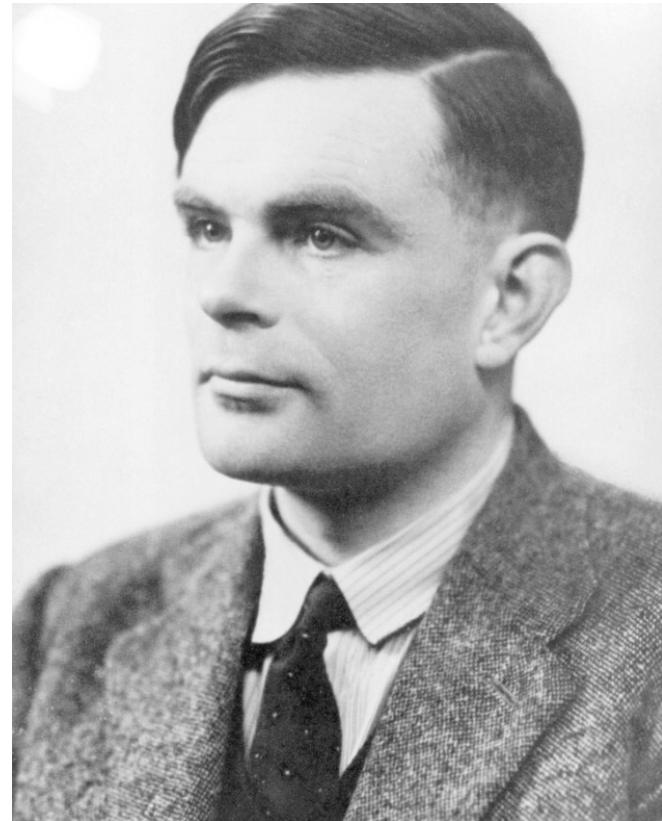
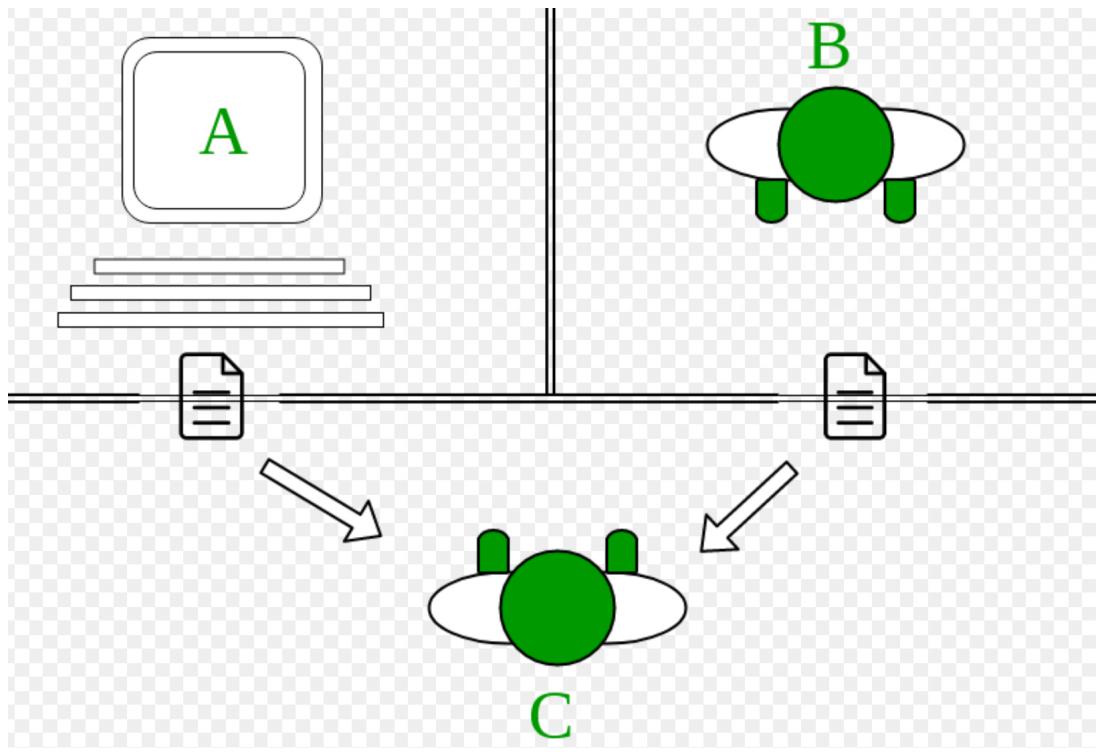


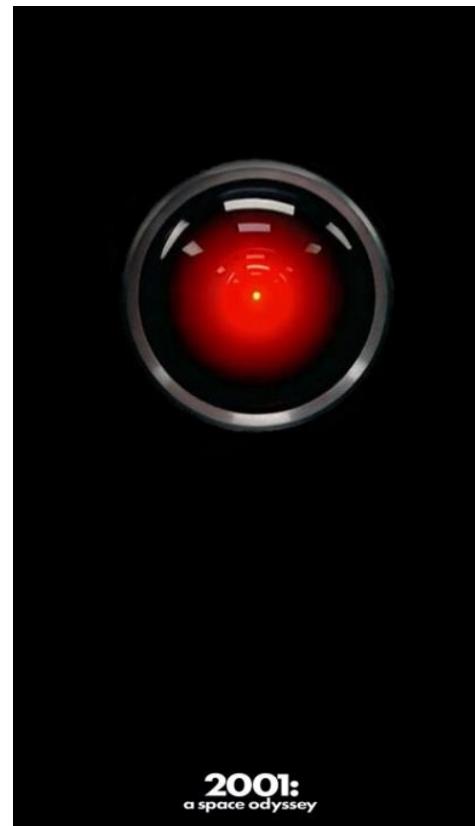
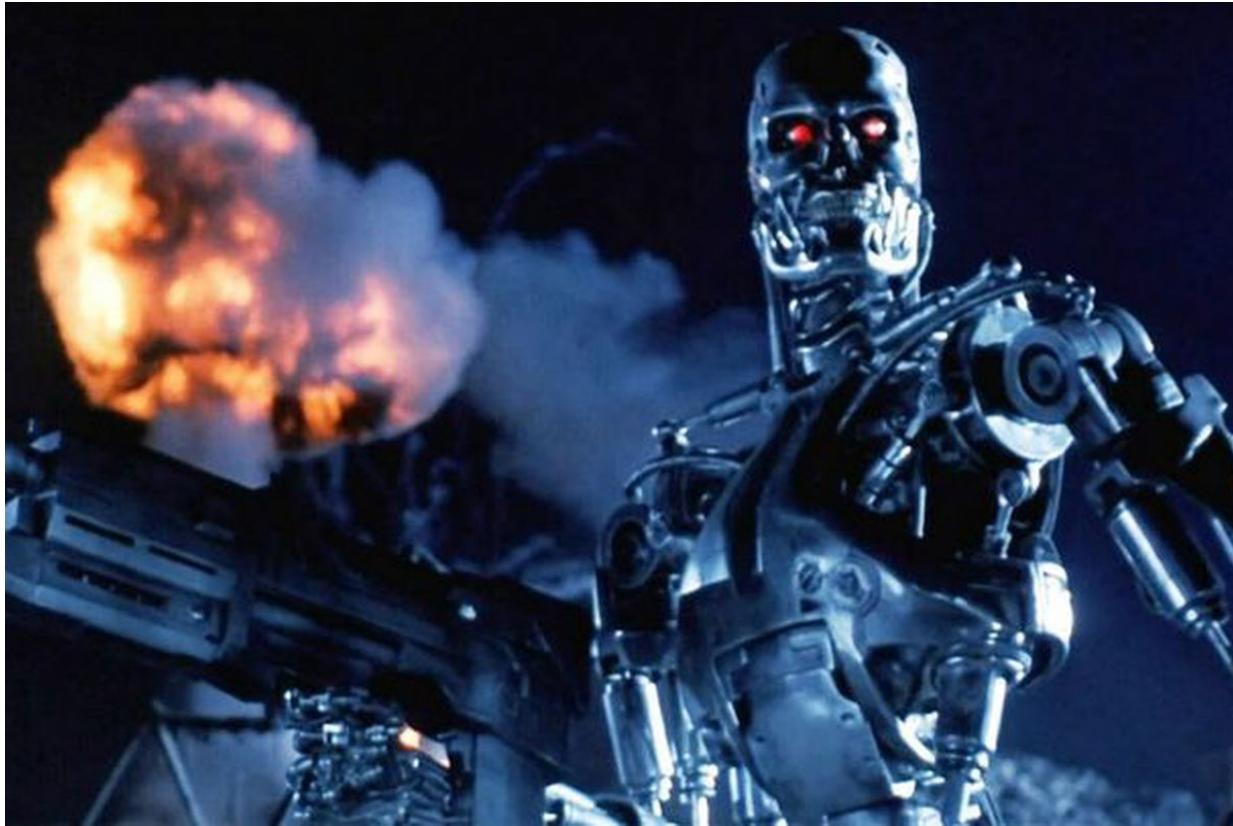
Biomonitoring: another classification problem → Machine Learning !



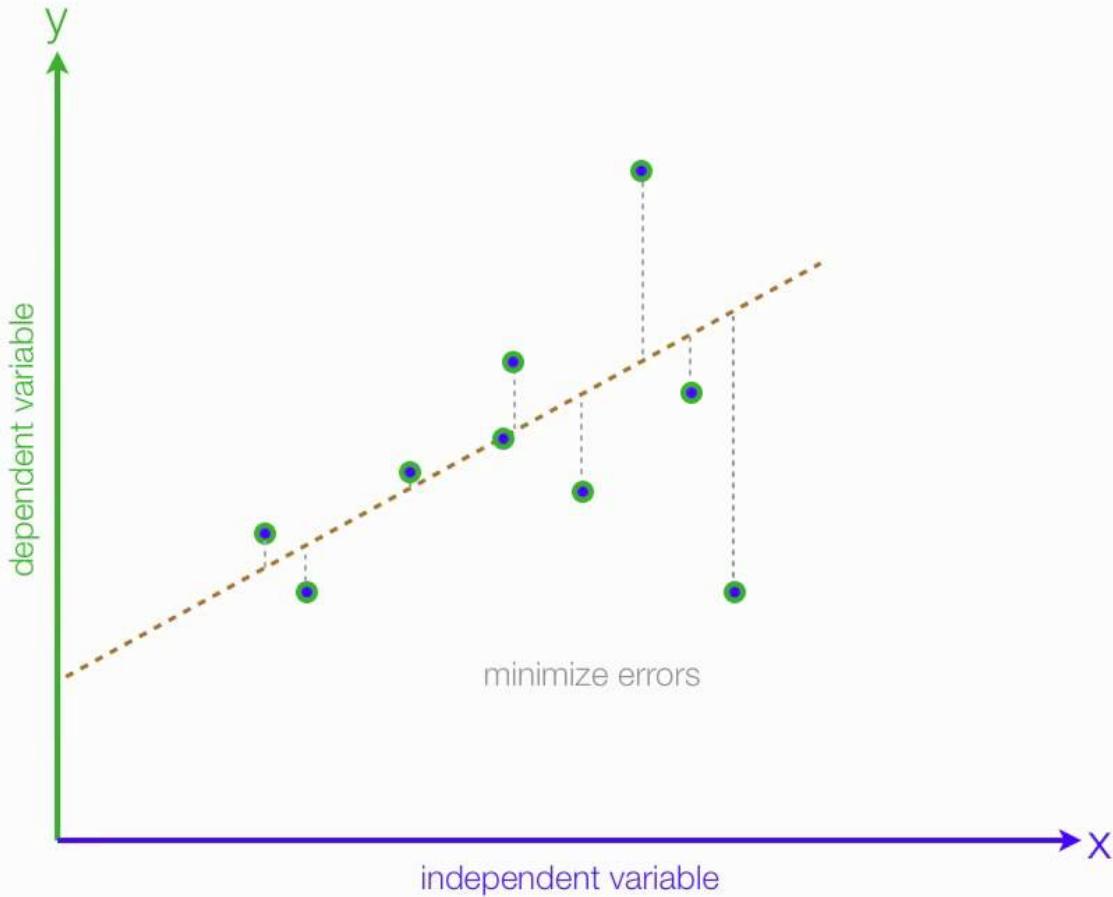
What is Machine Learning?







Association between things



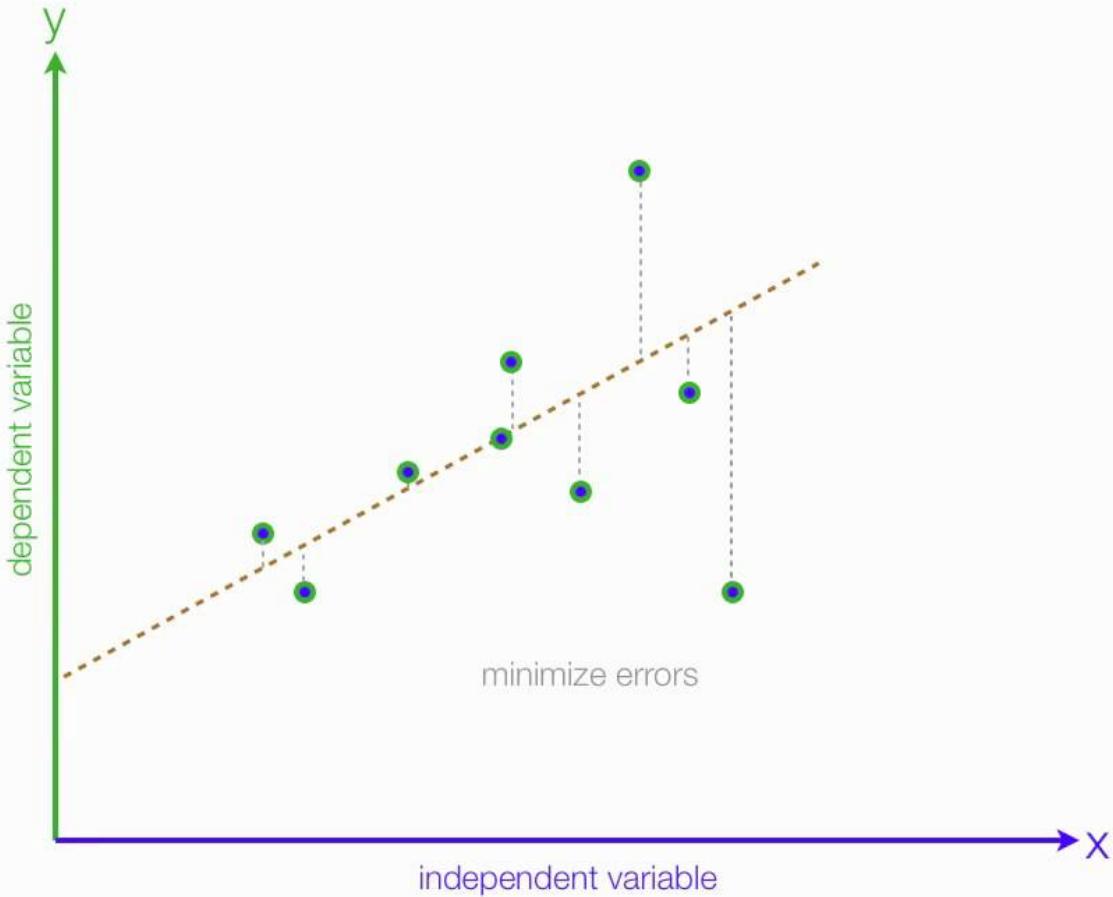
Statistical Modeling: The Two Cultures

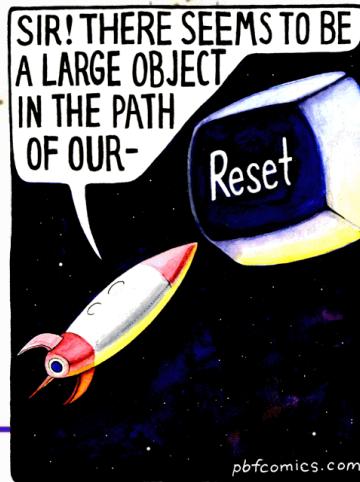
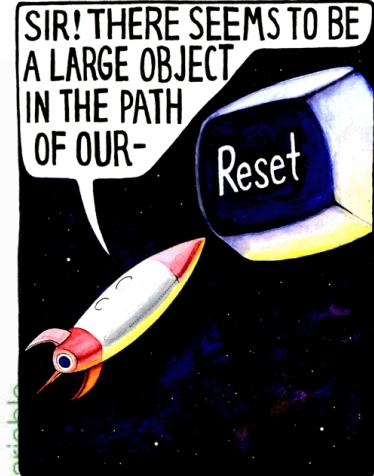
Leo Breiman



Abstract. There are two cultures in the use of statistical modeling to reach conclusions from data. One assumes that the data are generated by a given stochastic data model. The other uses algorithmic models and treats the data mechanism as unknown. The statistical community has been committed to the almost exclusive use of data models. This commitment has led to irrelevant theory, questionable conclusions, and has kept statisticians from working on a large range of interesting current problems. Algorithmic modeling, both in theory and practice, has developed rapidly in fields outside statistics. It can be used both on large complex data sets and as a more accurate and informative alternative to data modeling on smaller data sets. If our goal as a field is to use data to solve problems, then we need to move away from exclusive dependence on data models and adopt a more diverse set of tools.

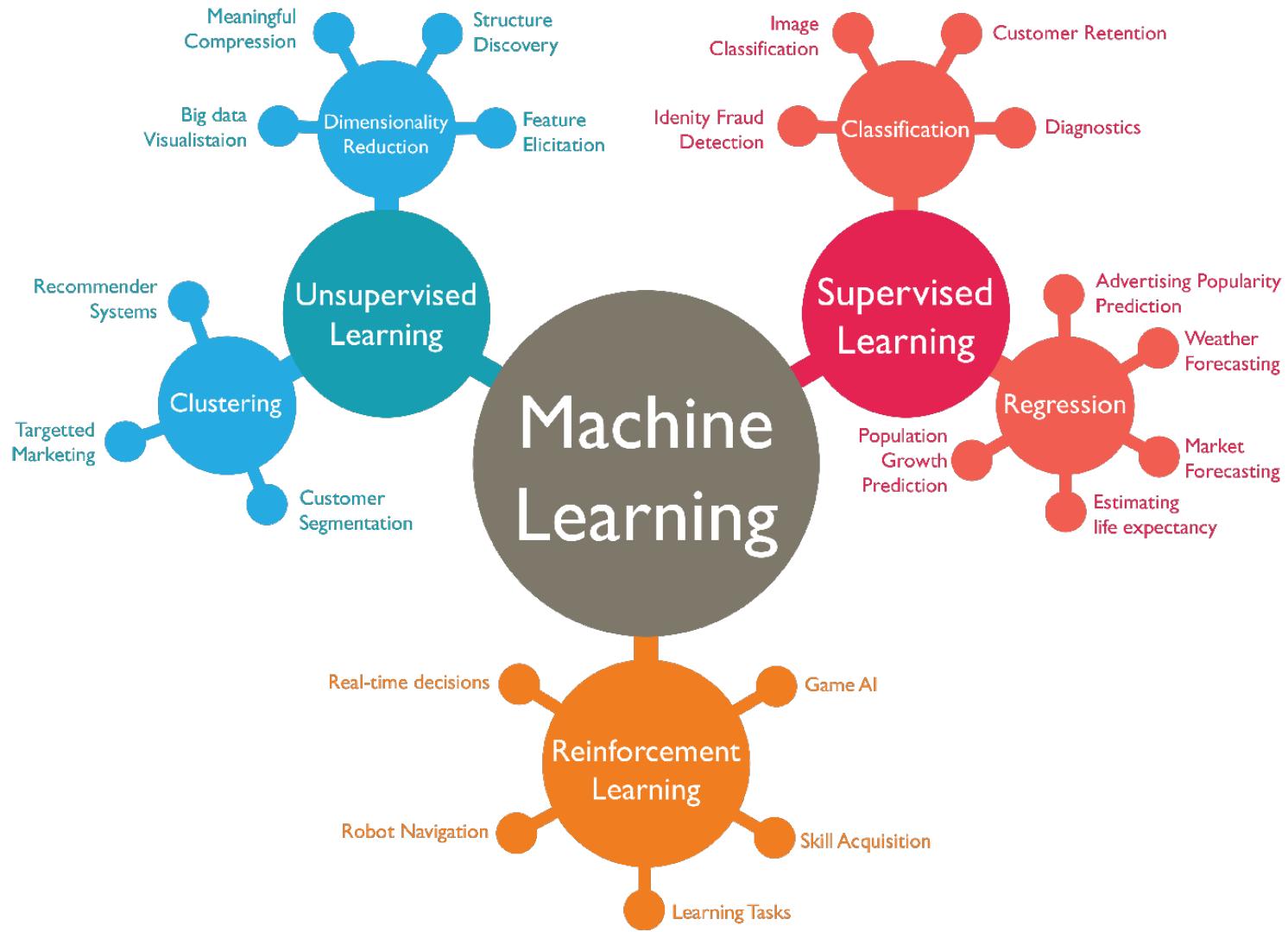
Association between things



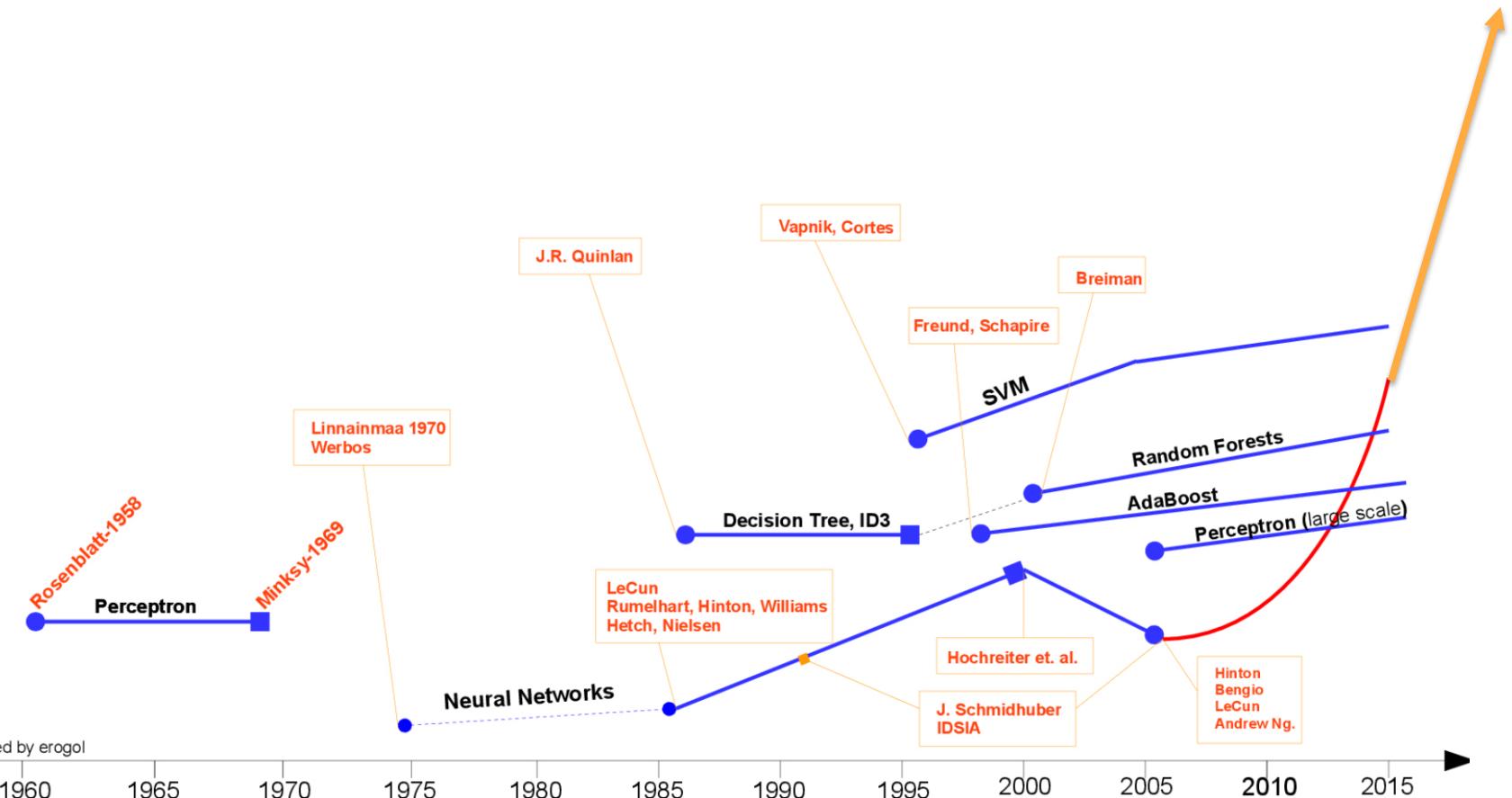


Association between things



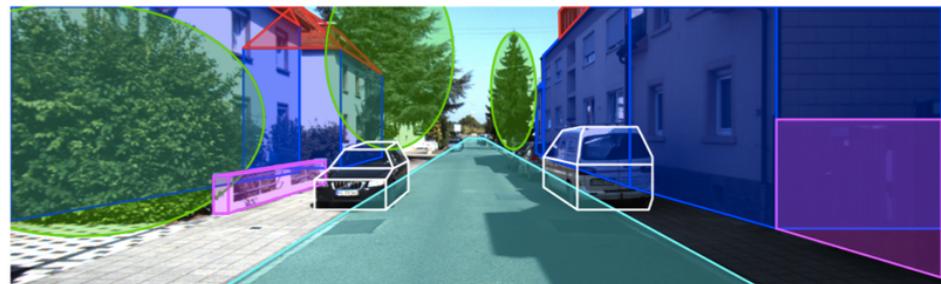
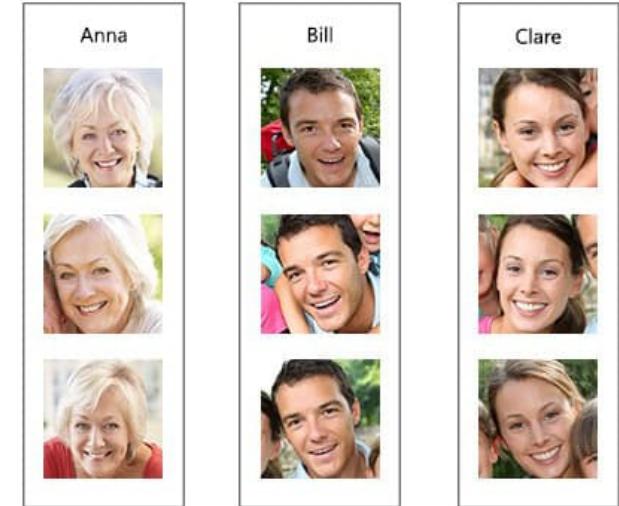
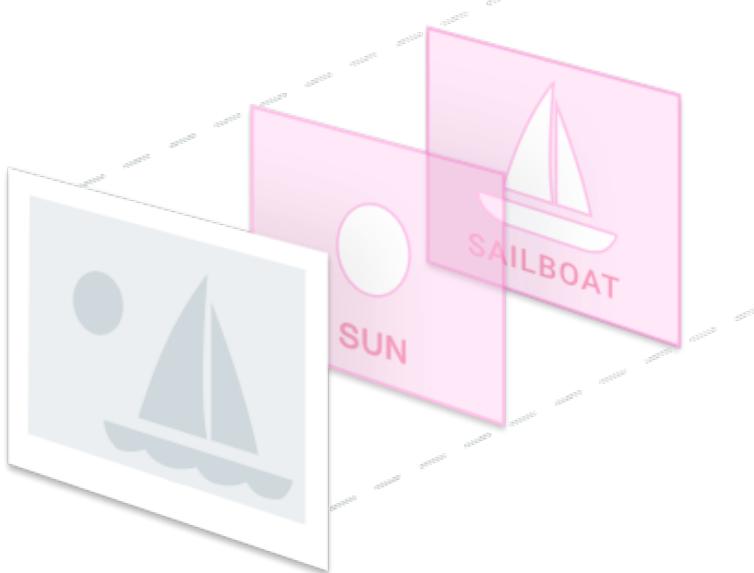


Subjective Popularity



Machine learning applications

- Widely used in classification problems
 - Computer vision tasks



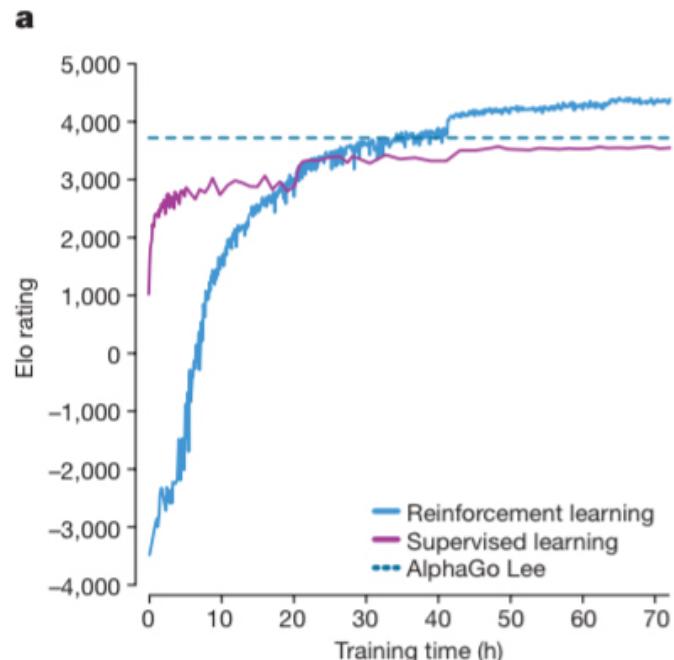
Machine learning applications

ARTICLE

doi:10.1038/nature24270

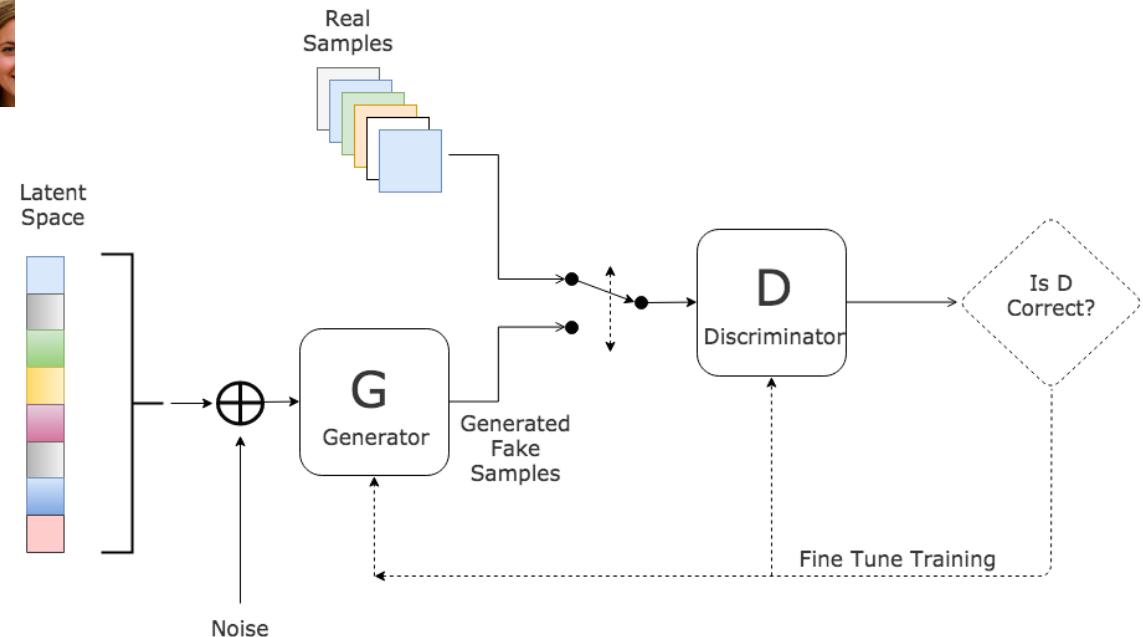
Mastering the game of Go without human knowledge

David Silver^{1*}, Julian Schrittwieser^{1*}, Karen Simonyan^{1*}, Ioannis Antonoglou¹, Aja Huang¹, Arthur Guez¹, Thomas Hubert¹, Lucas Baker¹, Matthew Lai¹, Adrian Bolton¹, Yutian Chen¹, Timothy Lillicrap¹, Fan Hui¹, Laurent Sifre¹, George van den Driessche¹, Thore Graepel¹ & Demis Hassabis¹



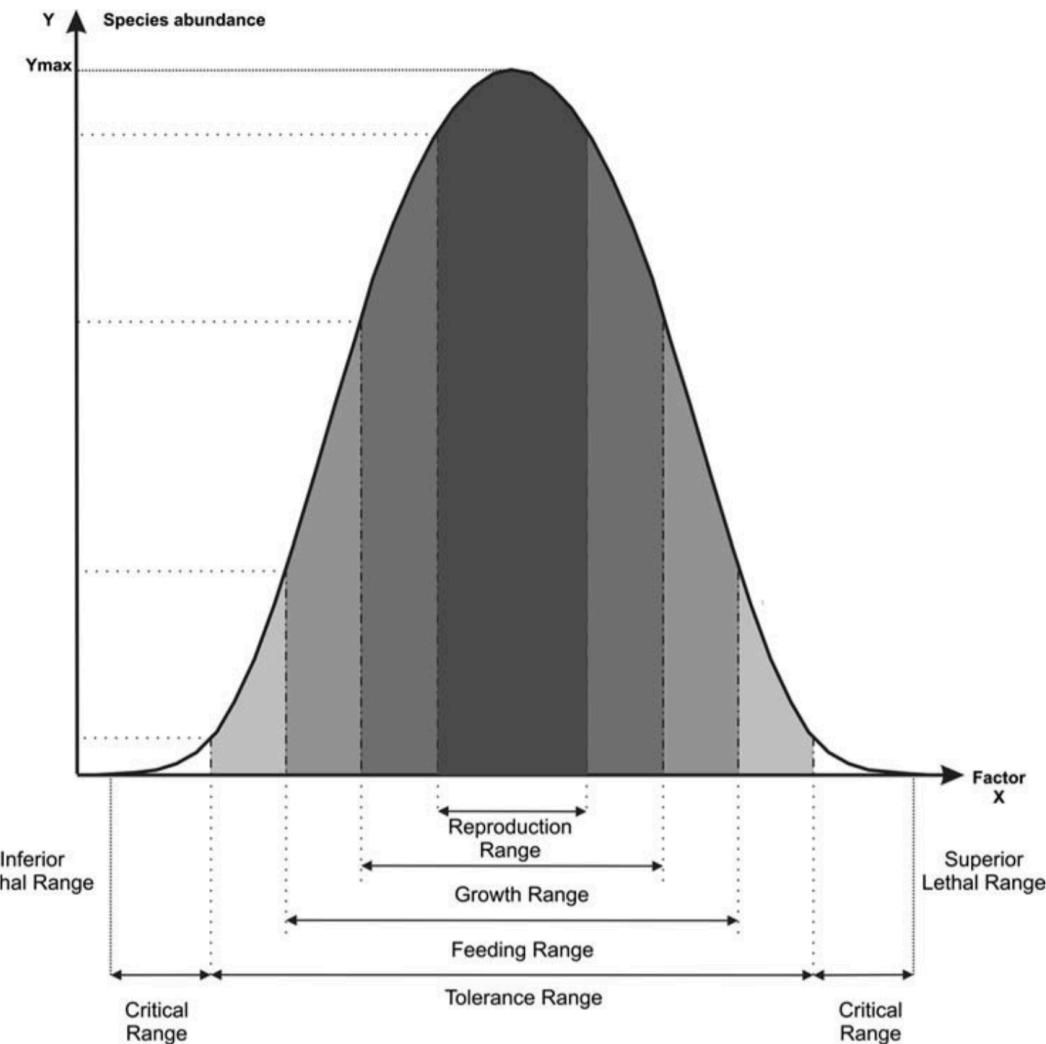


Generative Adversarial Network

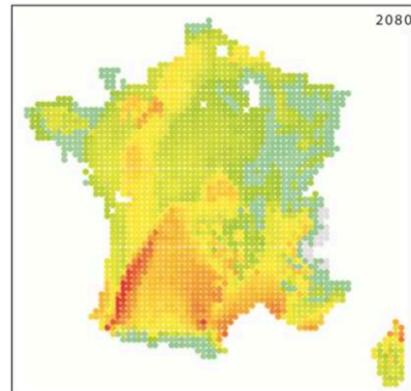
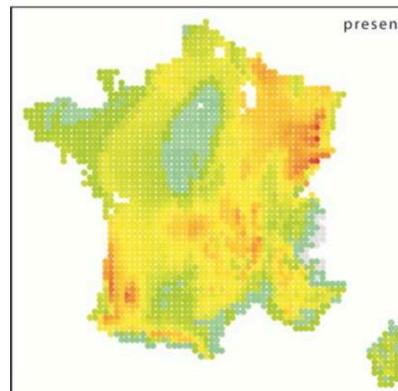
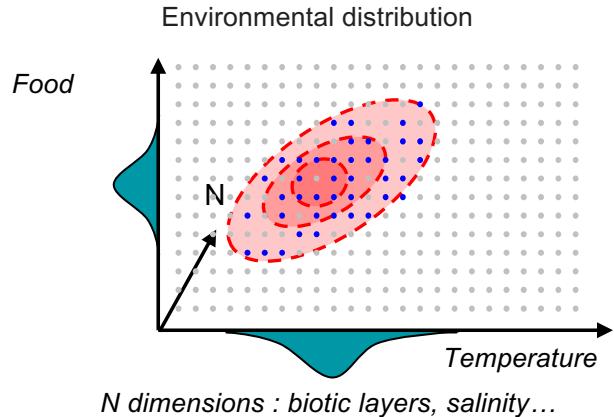
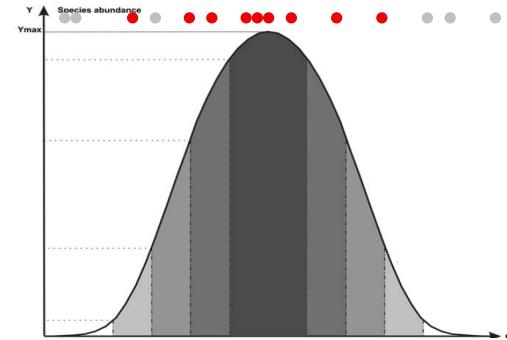
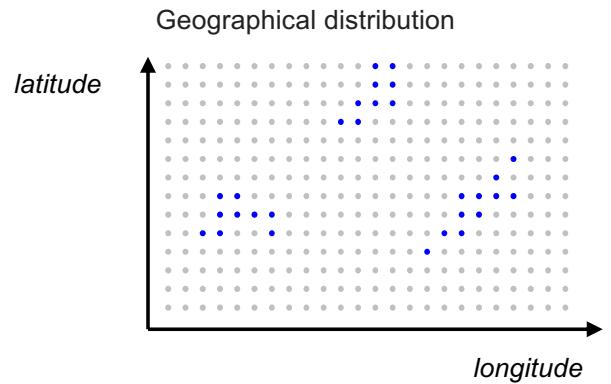


Ecological niche?

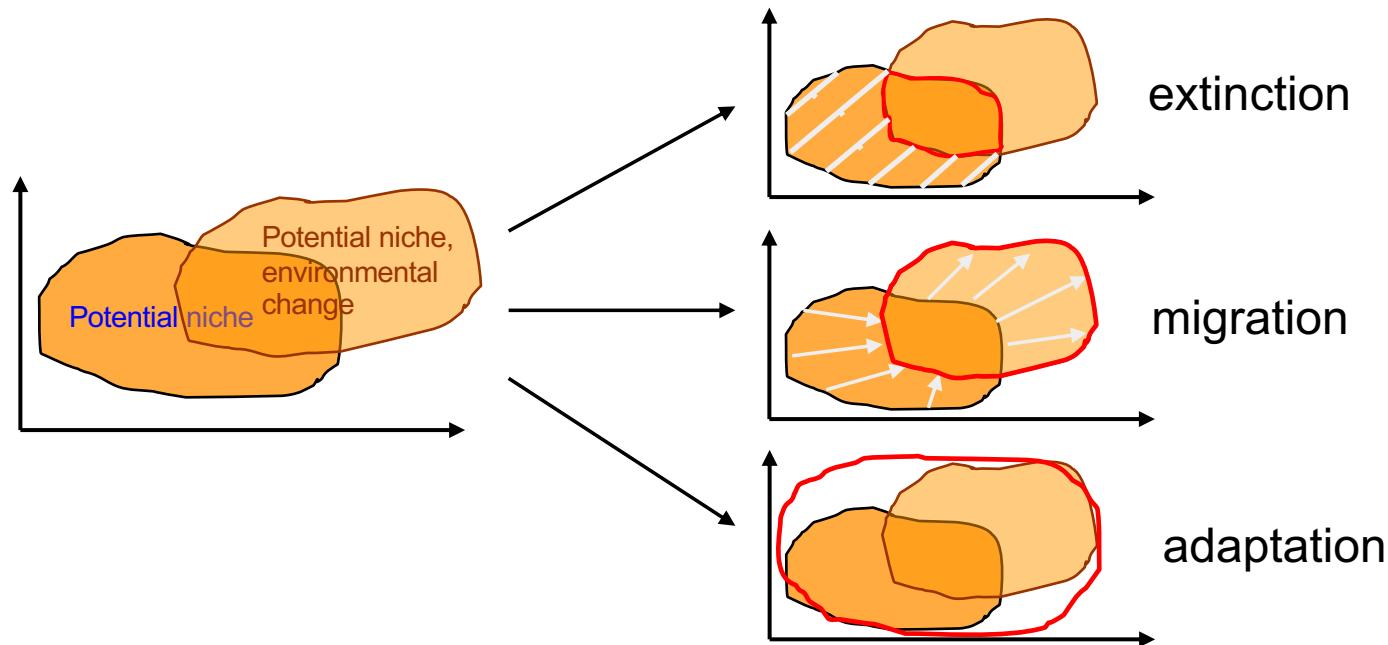
- Define the full set of environmental parameters (abiotic and biotic) that allow a species to live and reproduce.
(Hutchinson 1957)



Niche modeling

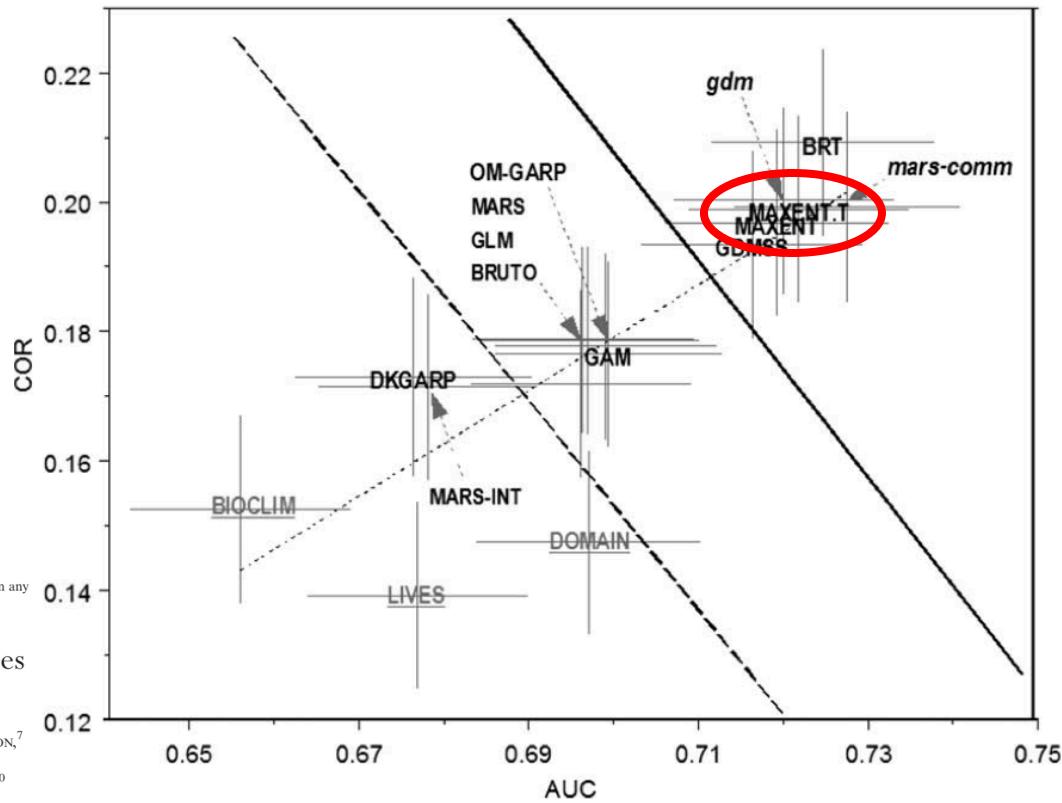


Which response to changes?



Tools for modelling

- BIOCLIM
- GLM
- GAM
- BRT
- **MAXENT** (Phillips, 2004)
- BIOMOD2 R package (Thuillier, 2009)



Ecological Monographs, 89(3), 2019, e01370
© 2019 The Authors. *Ecological Monographs* published by Wiley Periodicals, Inc. on behalf of Ecological Society of America
This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

A comprehensive evaluation of predictive performance of 33 species distribution models at species and community levels

ANNA NORBERG ^{1,34} NEREA ABREGO,^{2,3} F. GUILLAUME BLANCHET,⁴ FREDERICK R. ADLER,^{5,6} BARBARA J. ANDERSON,⁷ JANI ANTILA,¹ MIGUEL B. ARAÚJO,^{8,9,10} TAD DALLAS,¹ DAVID DUNSON,¹¹ JANE ELITH,¹² SCOTT D. FOSTER,¹³ RICHARD FOX,¹⁴ JANET FRANKLIN,¹⁵ WILLIAM GODSOE,¹⁶ ANTOINE GUISAN,^{17,18} BOB O'HARA,¹⁹ NICOLE A. HILL,²⁰ ROBERT D. HOLT,²¹ FRANCIS K. C. HUI,²² MAGNE HUSBY,^{23,24} JOHN ATLE KALAS,²⁵ ALEKSI LEHIKOINEN,²⁶ MIKA LUOTO,²⁷ HEIDI K. MOD,¹⁸ GRAEME NEWELL,²⁸ IAN RENNER,²⁹ TOMAS ROSLIN ,^{3,30} JANNE SOININEN ,²⁷ WILFRIED THUILLIER,³¹ JARNO VANHATALO,¹ DAVID WARTON,³² MATT WHITE,²⁸ NIKLAUS E. ZIMMERMANN,³³ DOMINIQUE GRAVEL,⁴ AND OTSO OVASKAINEN .^{1,2}

Elith et al., 2006

ML in genomics → Funct. annotation tax.classification

Predictive functional profiling of microbial communities using 16S rRNA marker gene sequences

Morgan G I Langille^{1,14}, Jesse Zaneveld^{2,14}, J Gregory Caporaso^{3,4}, Daniel McDonald^{5,6}, Dan Knights^{7,8}, Joshua A Reyes⁹, Jose C Clemente¹⁰, Deron E Burkepile¹¹, Rebecca L Vega Thurber², Rob Knight^{10,12}, Robert G Beiko¹ & Curtis Huttenhower^{9,13}

From Genomes to Phenotypes: Traitar, the Microbial Trait Analyzer

Aaron Weimann,^{a,b,c} Kyra Mooren,^{a,c} Jeremy Frank,^d Phillip B. Pope,^d

Andreas Bremges,^{a,b} Alice C. McHardy^{a,b,c}

Computational Biology of Infection Research, Helmholtz Center for Infection Research, Braunschweig, Germany^a; German Center for Infection Research (DZIF), Partner Site Hannover-Braunschweig, Braunschweig, Germany^b; Department for Algorithmic Bioinformatics, Heinrich Heine University, Düsseldorf, Germany^c; Department of Chemistry, Biotechnology and Food Science, Norwegian University of Life Sciences, Ås, Norway^d

Opinion

Next-Generation Global Biomonitoring: Large-scale, Automated Reconstruction of Ecological Networks

David A. Bohan,^{1,*} Corinne Vacher,² Alireza Tamaddoni-Nezhad,³ Alan Raybould,⁴ Alex J. Dumbrell,⁵ and Guy Woodward⁶

Naïve Bayesian Classifier for Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy^{▽†}

Qiong Wang,¹ George M. Garrity,^{1,2} James M. Tiedje,^{1,2} and James R. Cole^{1,*}

Center for Microbial Ecology¹ and Department of Microbiology and Molecular Genetics,² Michigan State University, East Lansing, Michigan 48824

IDTAXA: a novel approach for accurate taxonomic classification of microbiome sequences

Adithya Murali¹, Aniruddha Bhargava² and Erik S. Wright^{3*} 

ML in genomics → Medical research

MEDICAL RESEARCH

Machine learning classifies cancer

Brain tumours are often classified by visual assessment of tumour cells, yet such diagnoses can vary depending on the observer. Machine-learning methods to spot molecular patterns could improve cancer diagnosis. [SEE ARTICLE P.469](#)

DEREK WONG & STEPHEN YIP

Accurate diagnosis is essential for appropriate disease treatment. A core technique used to diagnose brain cancer today is the microscope-based analysis of tumour samples on glass slides, termed histology. However, this requires the appraisal of subtle cellular alterations, which in some cases may lead to different classifications for a given sample by different individuals. Nowadays, technological developments enable vast amounts of molecular data to be obtained and assessed for a tumour without the need for such subjective diagnostics. Machine-based-learning approaches are being developed to aid the diagnosis of clinical samples, and on page 469, Capper *et al.*¹ report such a method for classifying brain tumours on the basis of molecular patterns.

In 1926, a publication entitled *A Classification of the Tumors of the Glioma Group on a Histo-Genetic Basis with a Correlated Study of Prognosis*² by neurosurgeons Percival Bailey and Harvey Cushing provided early insight into the development, cellular characteristics and clinical consequences of glioma, a type of

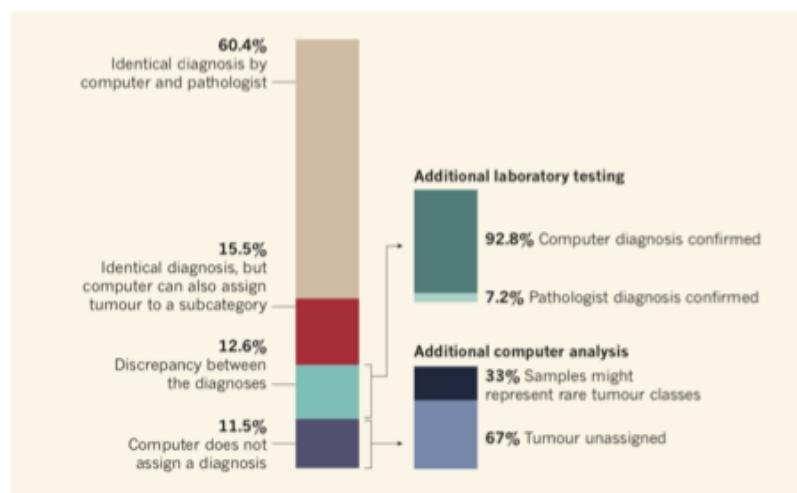
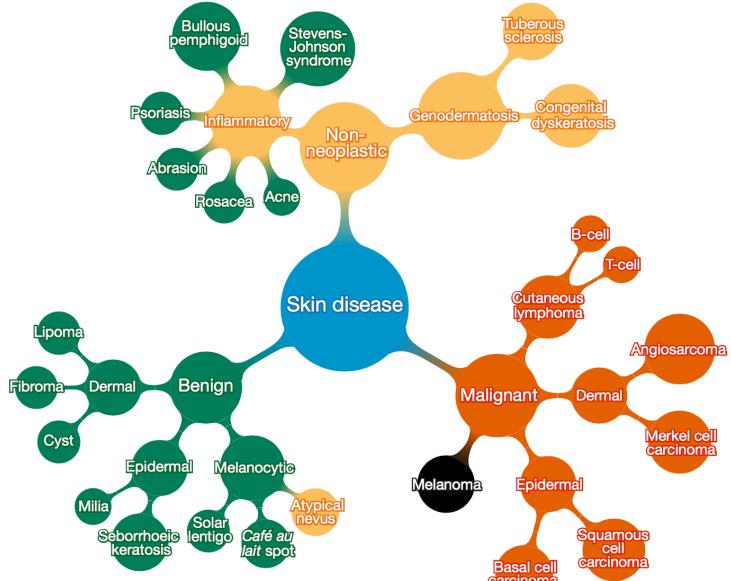


Figure 1 | Tumour classification using a machine-learning approach. Capper *et al.*¹ used a machine-learning approach to classify brain tumours on the basis of genome-wide patterns of a type of DNA alteration called methylation. The computer was trained using methylation data for tumour samples that had been diagnosed by pathologists using standard microscopy-based analysis or analysis of selected genes. After training, the computer was given 1,104 test cases. The authors compared the diagnoses made by the computer and by the pathologists. Although the machine was unable to diagnose all specimens, of the specimens that it classified, the machine-based diagnosis was more accurate or could assign tumours to more-specific subcategories than the classifications made by the pathologists.

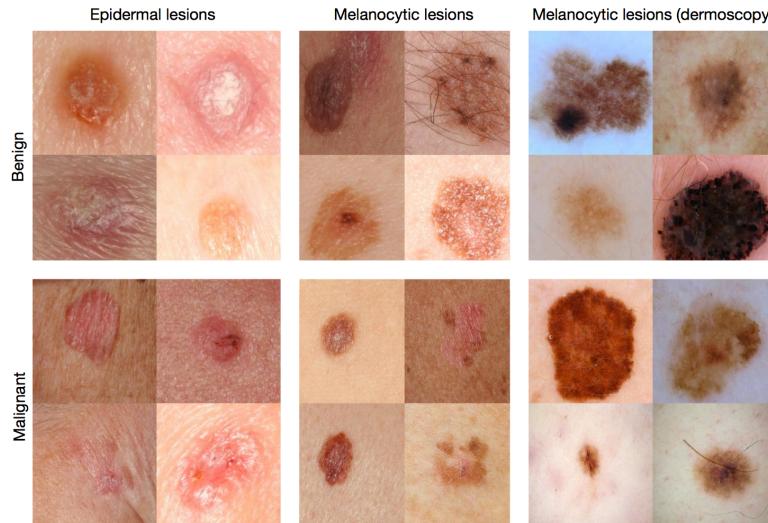
ML in genomics → Medical research

- Widely used in classification problems

- Cancer detection



Classifier	Three-way accuracy
Dermatologist 1	65.6%
Dermatologist 2	66.0%
CNN	69.4 ± 0.8%
CNN - PA	72.1 ± 0.9%



Esteva et al., 2017

ML in genomics → Marine research

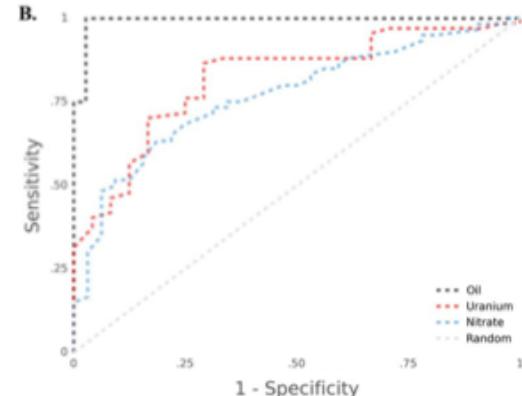
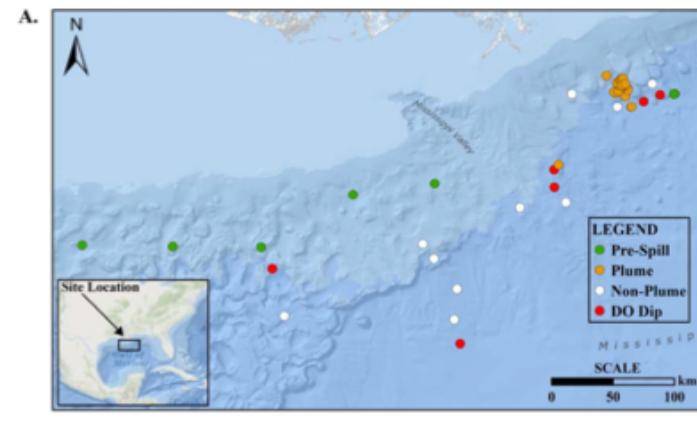
→ Impact assessment - biomonitoring

Natural Bacterial Communities Serve as Quantitative Geochemical Biosensors

Mark B. Smith,^a Andrea M. Rocha,^b Chris S. Smillie,^c Scott W. Olesen,^d Charles Paradis,^e Liyou Wu,^f James H. Campbell,^{b,m} Julian L. Fortney,^g Tonia L. Mehlhorn,^h Kenneth A. Lowe,^h Jennifer E. Earles,^h Jana Phillips,^h Steve M. Techtmann,^g Dominique C. Joyner,^g Dwayne A. Elias,^b Kathryn L. Bailey,^b Richard A. Hurt, Jr.,^b Sarah P. Preheim,^d Matthew C. Sanders,^d Joy Yang,^c Marcella A. Mueller,^h Scott Brooks,^h David B. Watson,^h Ping Zhang,^f Zhili He,^f Eric A. Dubinsky,ⁱ Paul D. Adams,^{j,l} Adam P. Arkin,^{j,l} Matthew W. Fields,^j Jizhong Zhou,^f Eric J. Alm,^{a,c,d} Terry C. Hazen^{b,e,g,k}

Here we show that statistical analysis of **DNA from natural microbial communities can be used to accurately identify environmental contaminants**, including uranium and nitrate at a nuclear waste site. In addition to contamination, **sequence data from the 16S rRNA gene alone can quantitatively predict a rich catalogue of 26 geochemical features** collected from 93 wells with highly differing geochemistry characteristics.

Smith et al., mBio, 2015



What tool do I need?

Continuous

- Clustering - dimensionality reduction
 - PCA
 - NMDS
 - K-means

Categorical

- Association analysis
 - Correspondence analysis
 - Factorial analysis

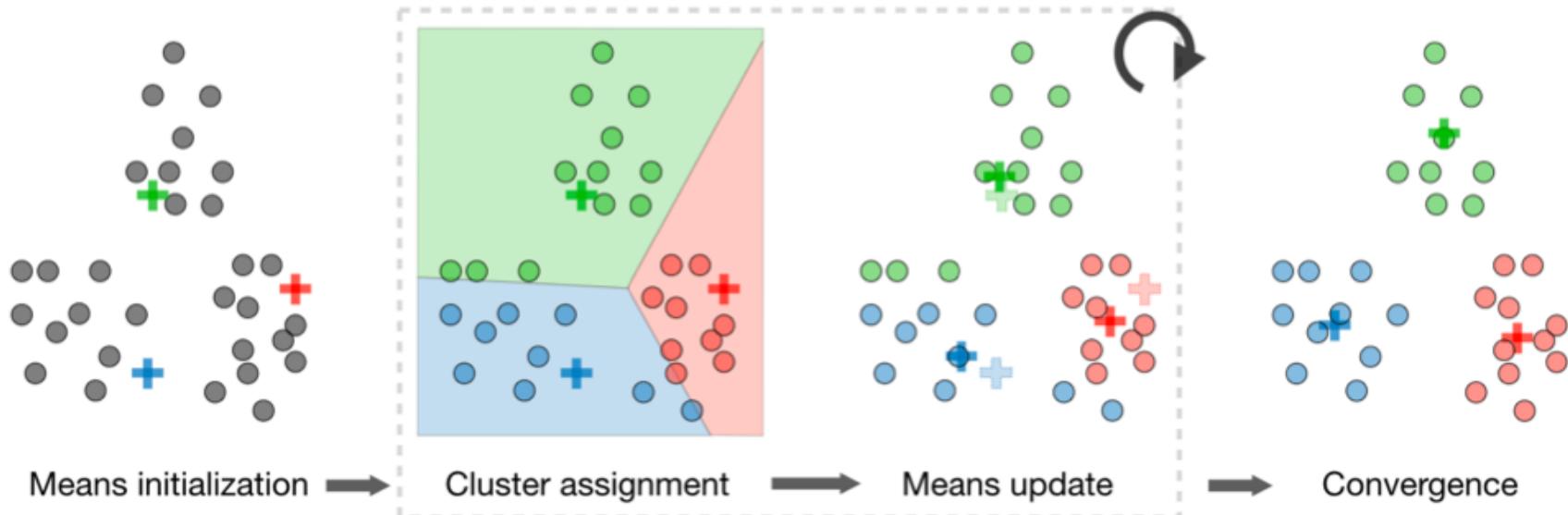
Unsupervised

Supervised

- Regression (linear or not)
 - Linear or polynomial models
 - Regression trees (random forest)
 - Neural networks
- Classification
 - Random forest
 - Neural networks
 - Logistic regression
 - SVMs
 - ...

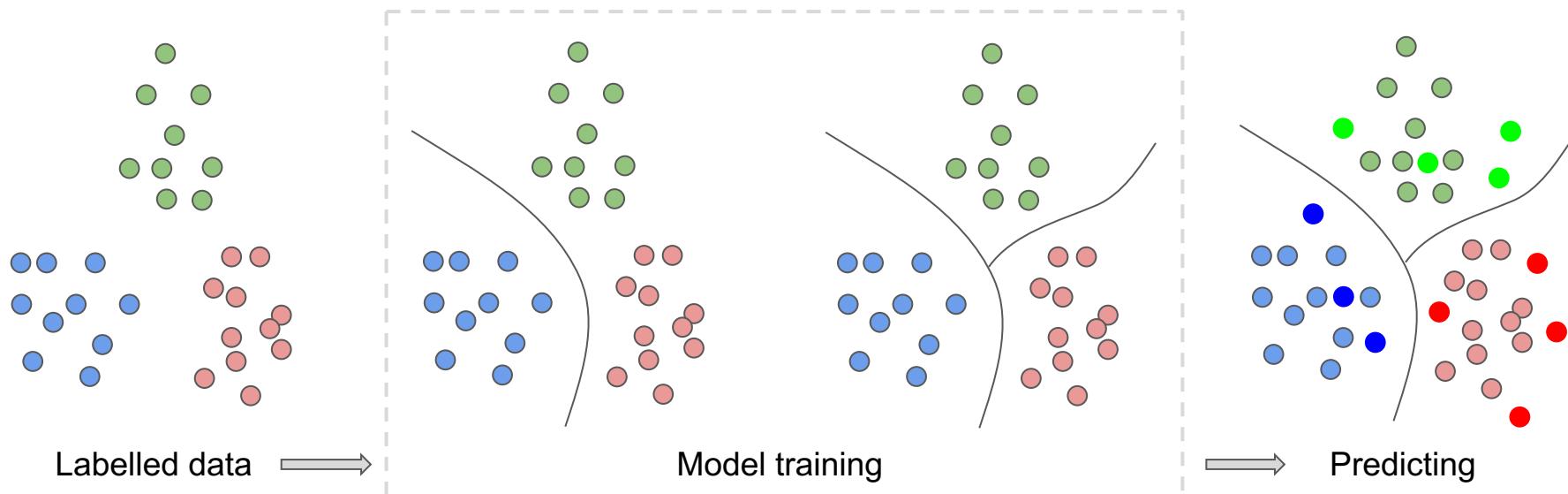
Type of problems in machine learning

Unsupervised: the solution is not known. — **Goal:** Find hidden patterns



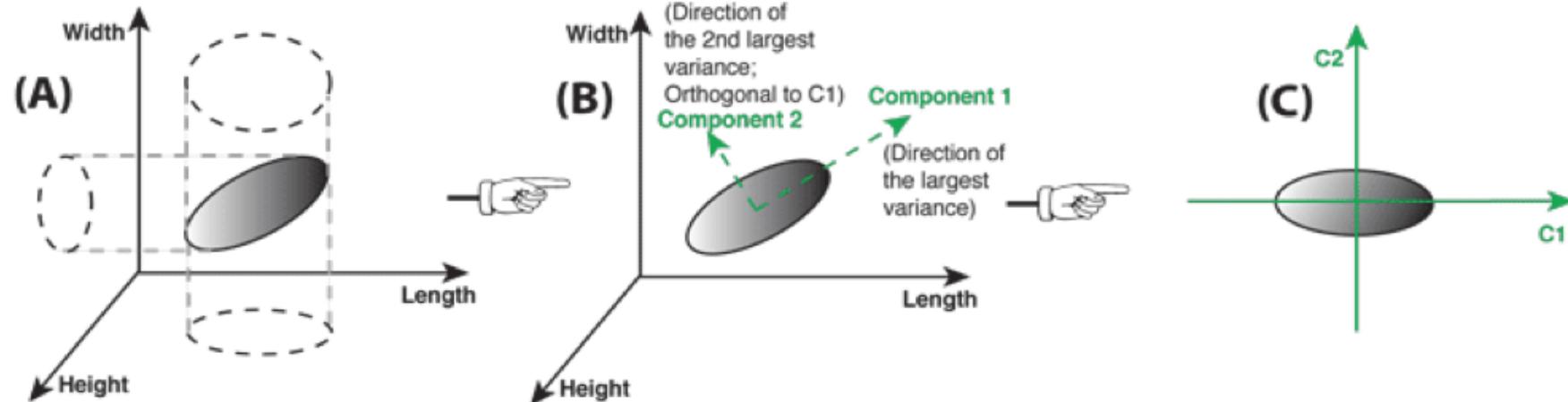
Type of problems in machine learning

Supervised: the solution is known. — **Goal:** Make predictions on new data



Unsupervised - Dimensionality reduction - PCA

Principal component analysis is all about how to choose a good coordinate system



$$\text{Data}(i) = f(L(i), W(i), H(i))$$

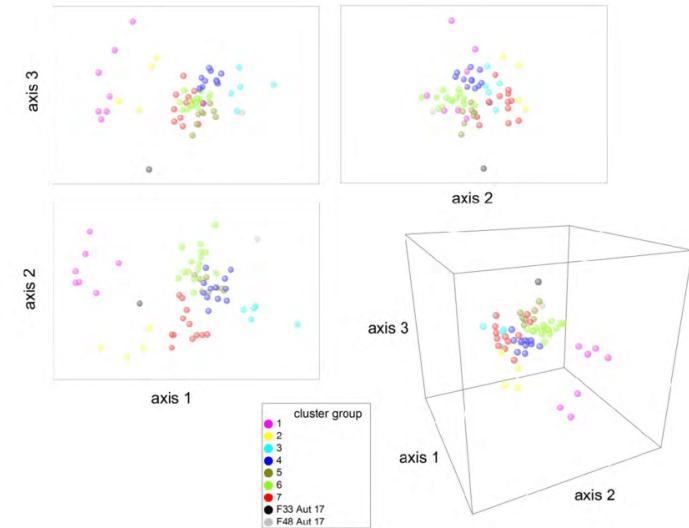
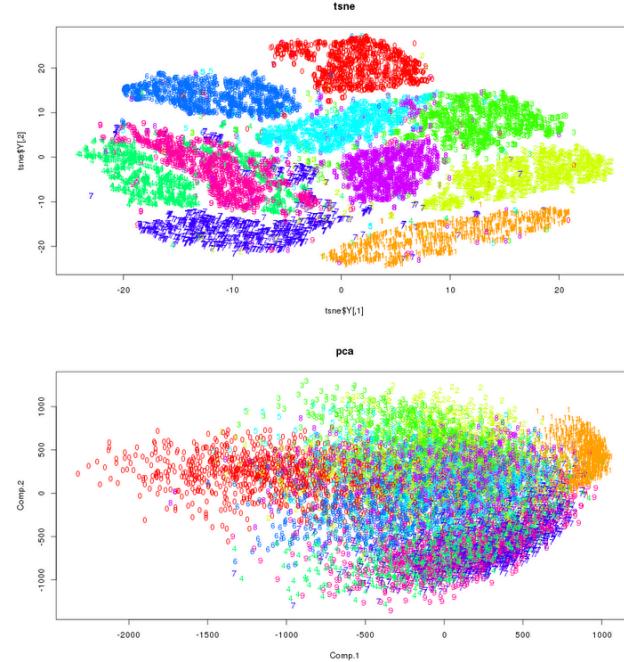
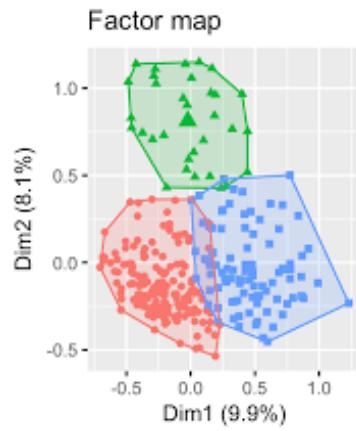
[3-dimensional function]

$$\text{Data}(j) = g(C1(j), C2(j))$$

[2-dimensional function]

- Projection of n-dimensional data on new coordinate axis using eigen-decomposition
- Uses actual OTU abundances with only linear transformation

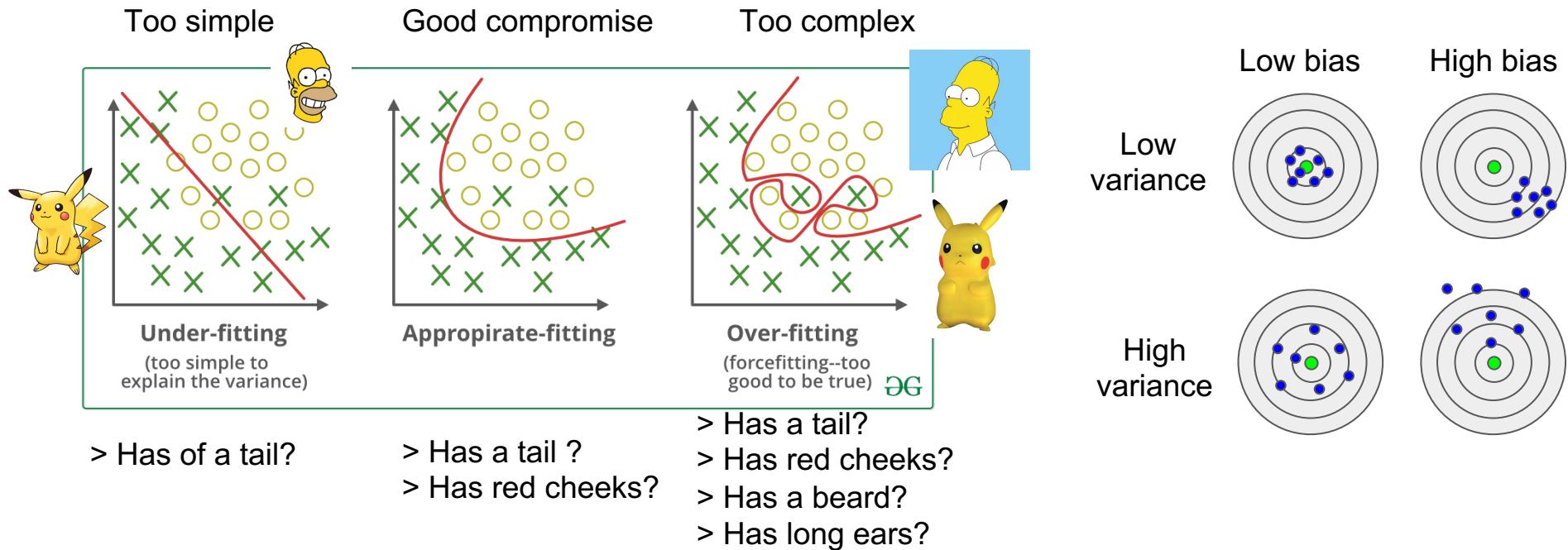
Unsupervised - Dimensionality reduction - PCA



- **Projection of n-dimensional data on new coordinate axis using eigen-decomposition**
- Uses actual OTU abundances with only linear transformation

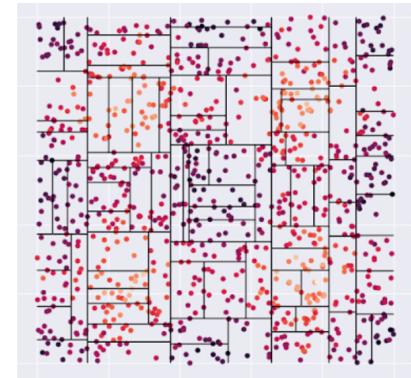
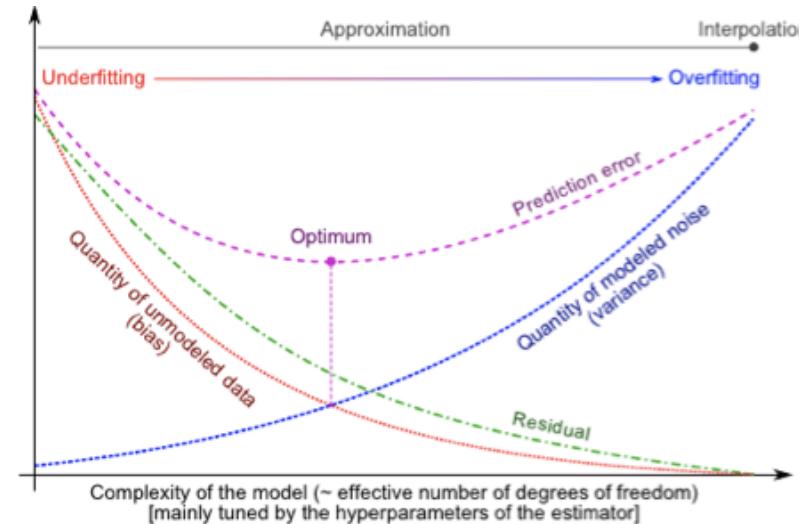
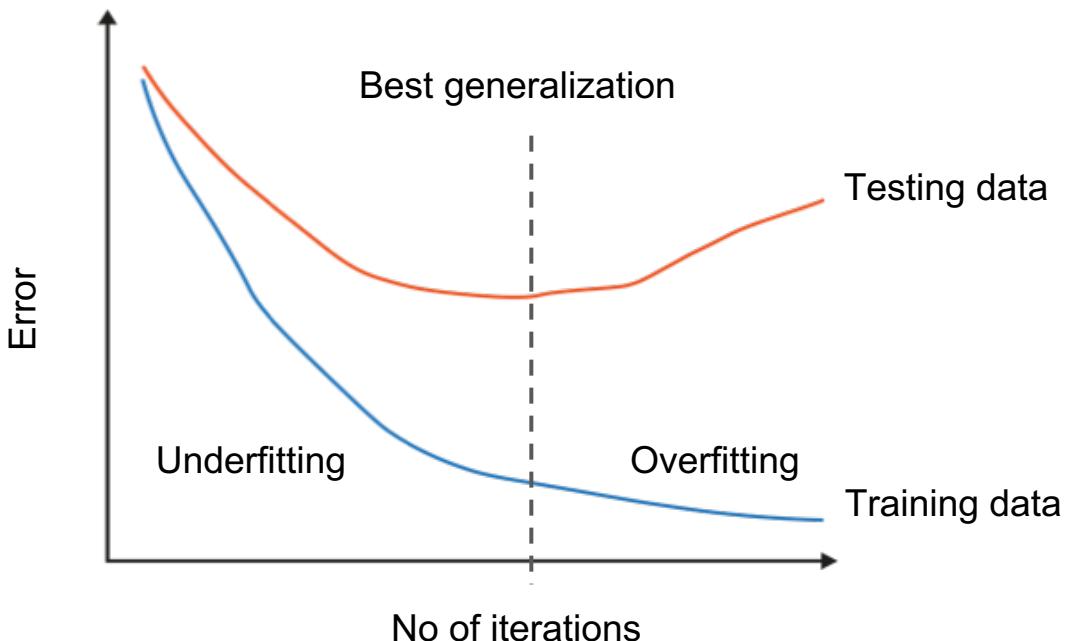
Supervised - predictive model training

- Inferring a **function** from **labelled training data** (subset of the real world!)
- The aim is to **minimize** the **error** of predictions in the real world.



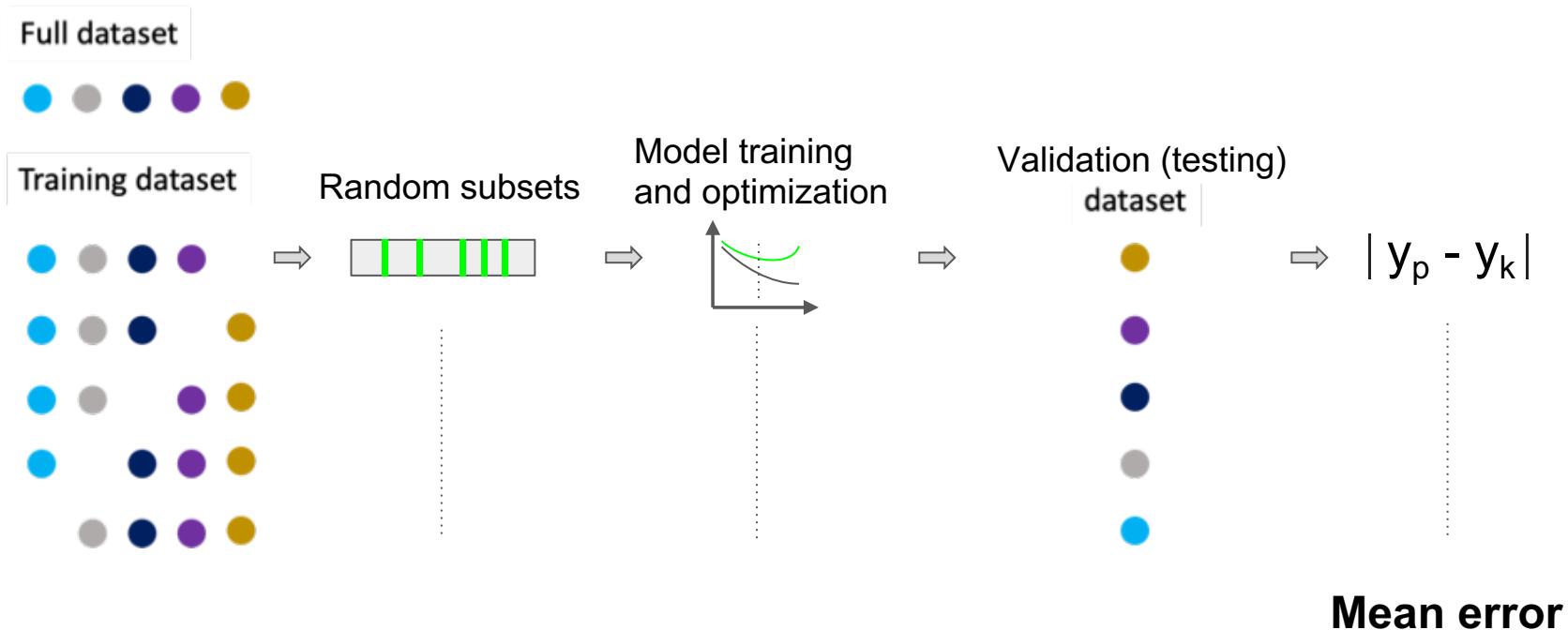
Supervised - predictive model training

- The bias-variance trade-off



Supervised - predictive model training

- Using cross-validation to measure *generalization error*

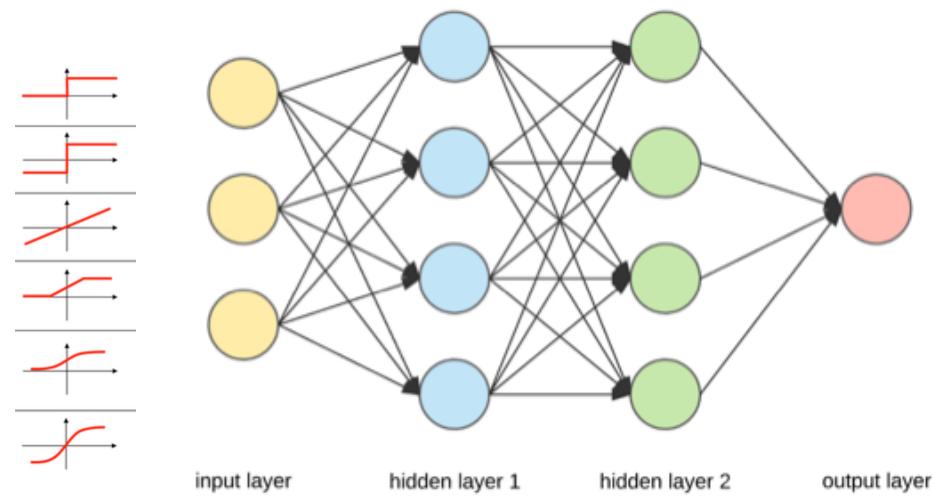


neural networks

Randomly initiate weights for several nodes in several layers of classifiers → “predict” training data and compare to ground truth (using a *loss function*) → backwards pass to try to correct (complex...)

Multiple parameters to play with :

- Number of layers / neurons
- Learning rate
- Activation function (sigmoid, linear...)



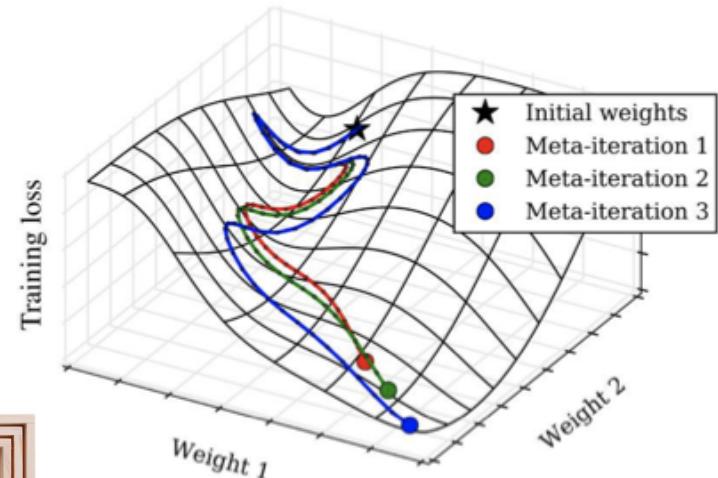
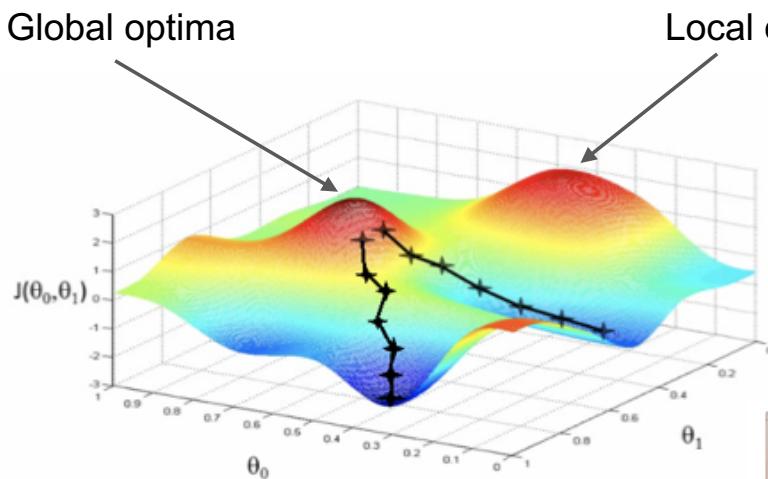
Supervised - predictive model training

- Hill-climbing vs gradient descent for model optimization



Global optima

Local optima

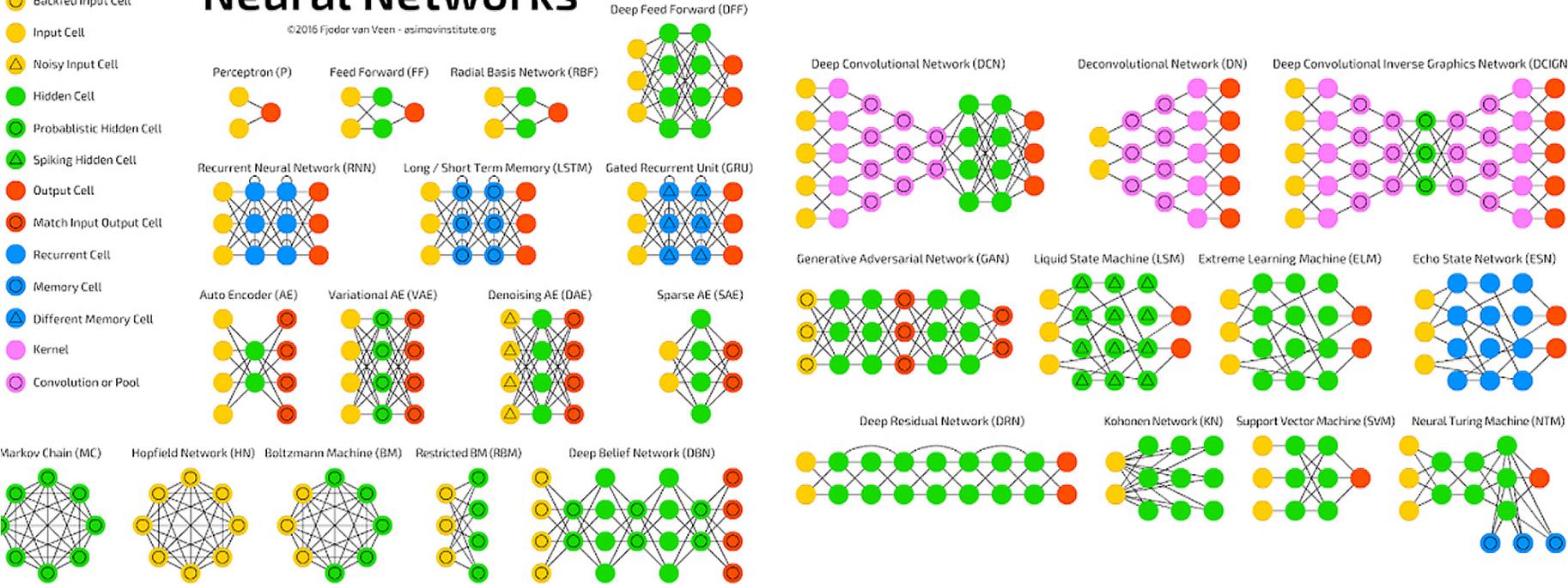


A mostly complete chart of

Neural Networks

©2016 Fjodor van Veen - osimovinstitute.org

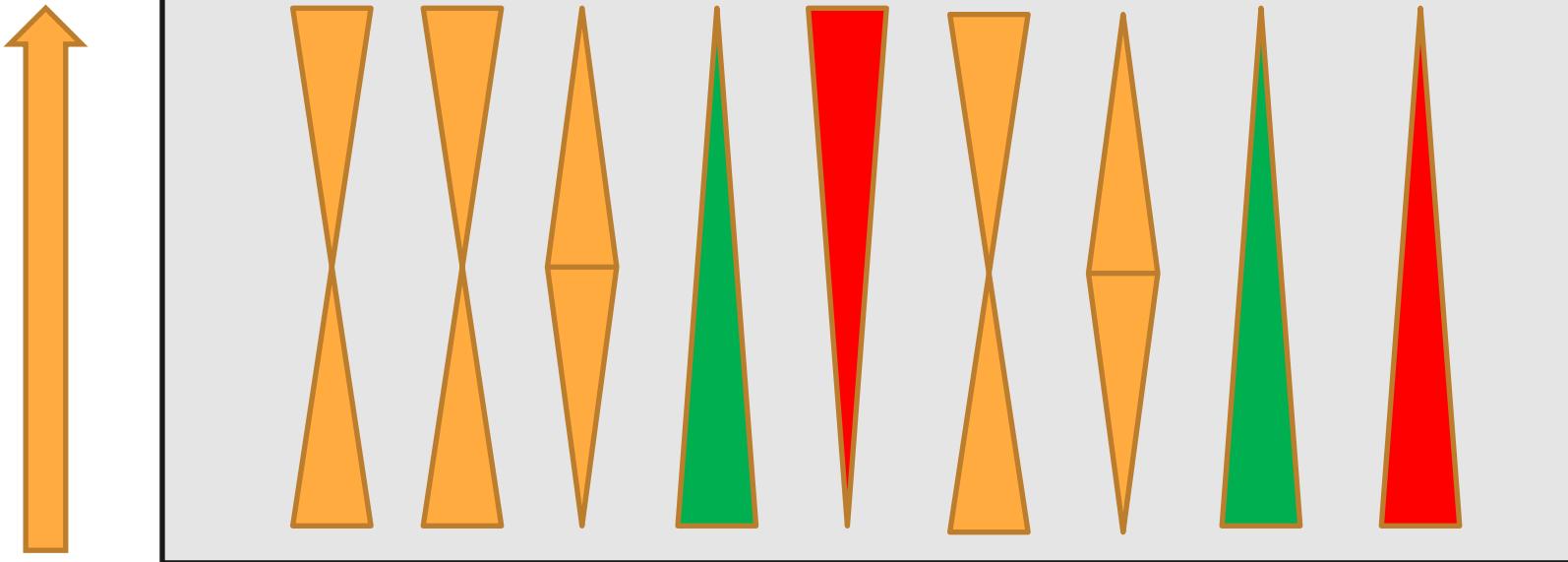
- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Different Memory Cell
- Kernel
- Convolution or Pool



Millions to billions of samples with thousands of variables

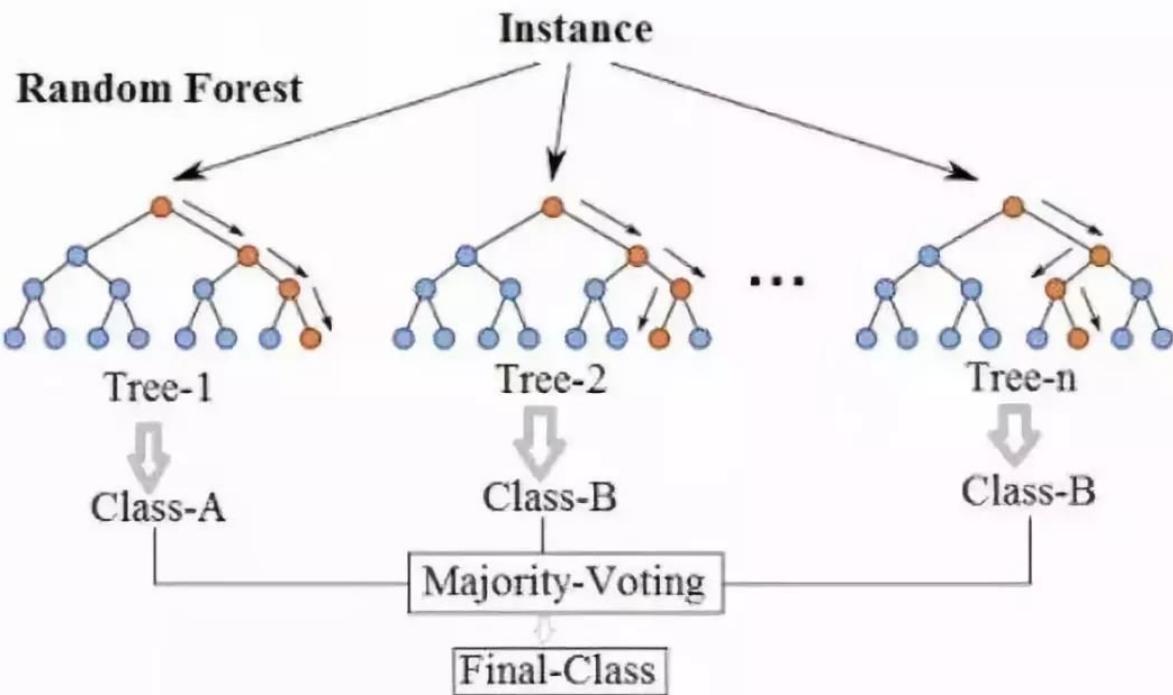
eDNA : thousand of samples with millions of variables (OTUs)

eDNA : thousand of samples with millions of variables (OTUs)
—> curse of dimentionality



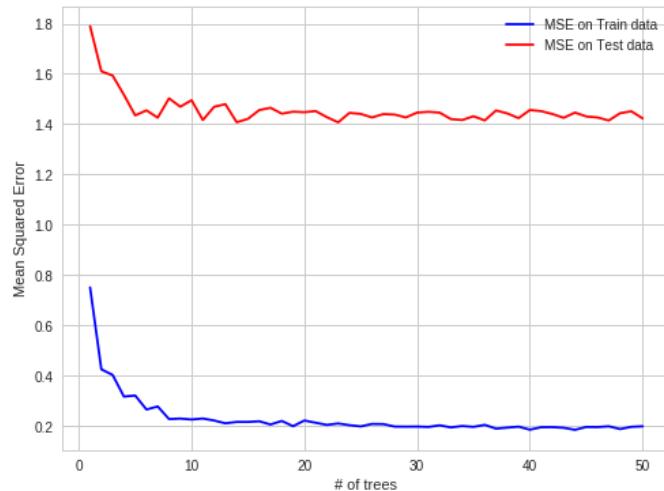


Random Forest Simplified



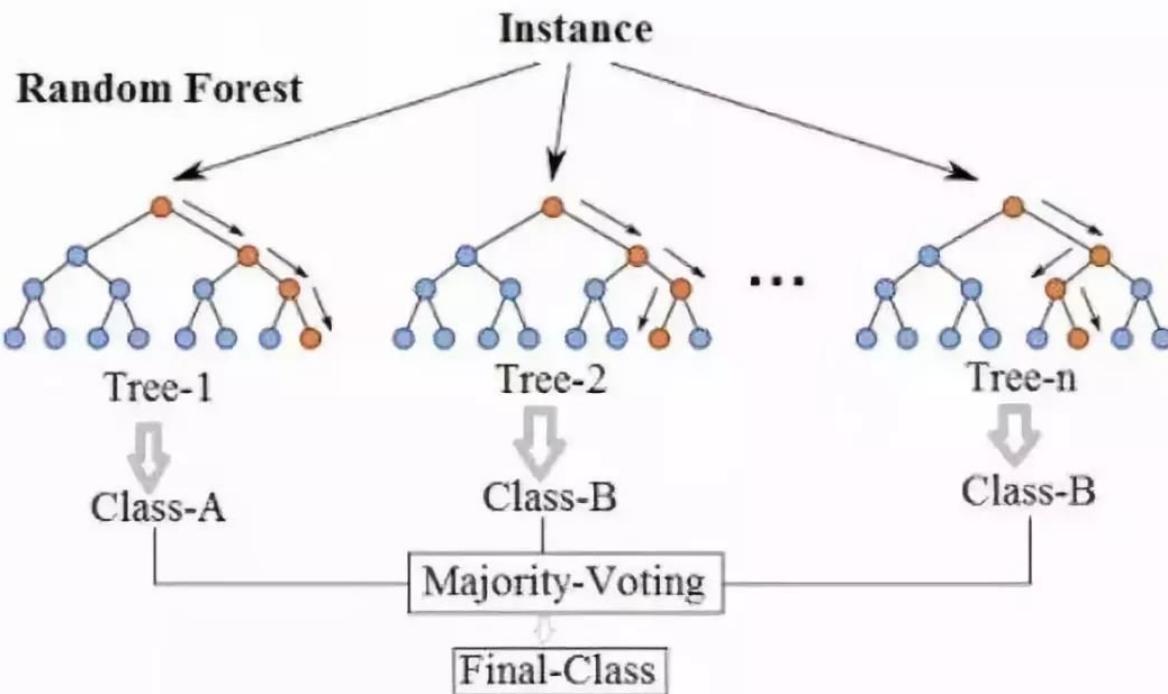
Two parameters:

- Number of trees



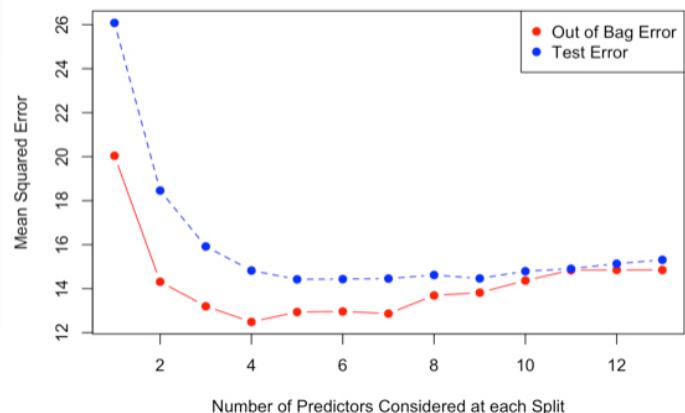


Random Forest Simplified

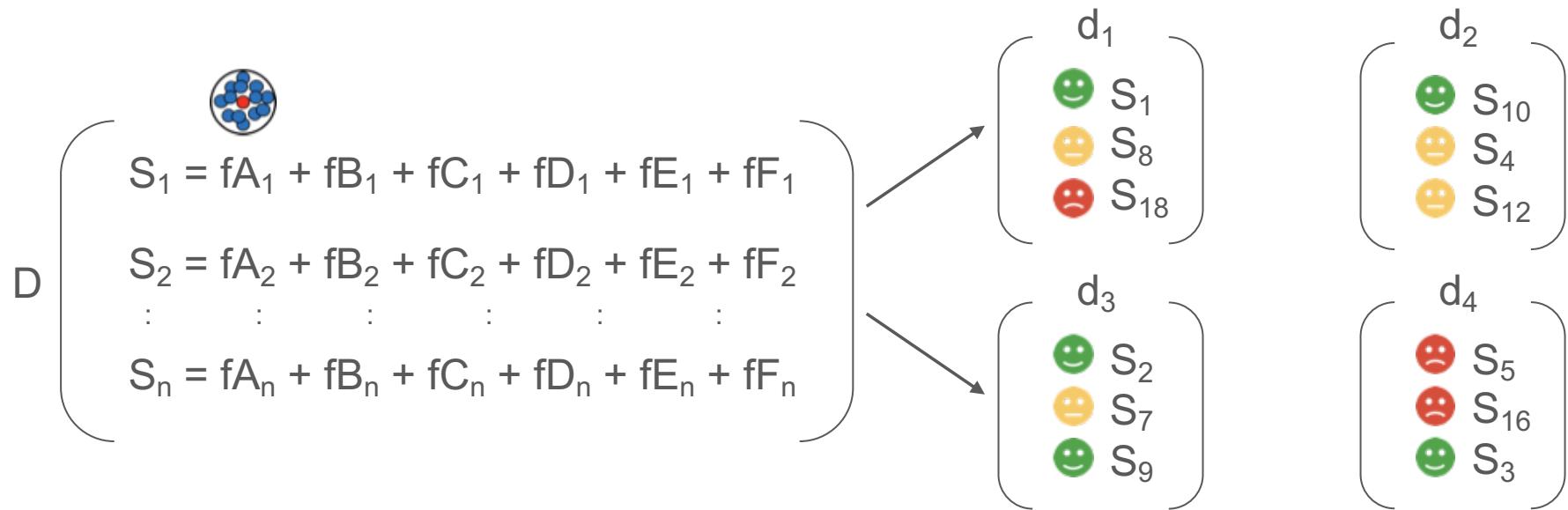


Two parameters:

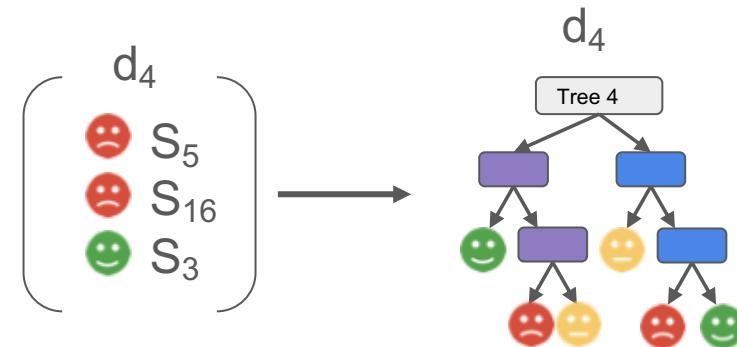
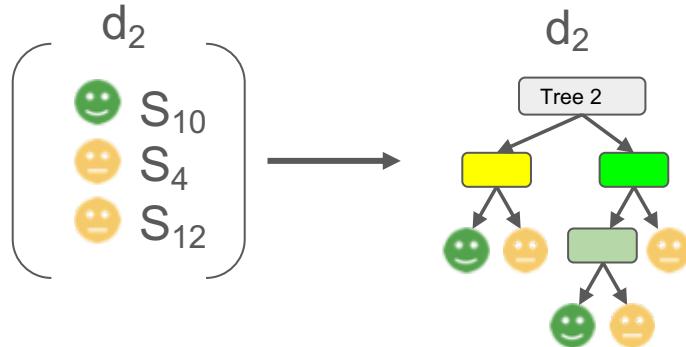
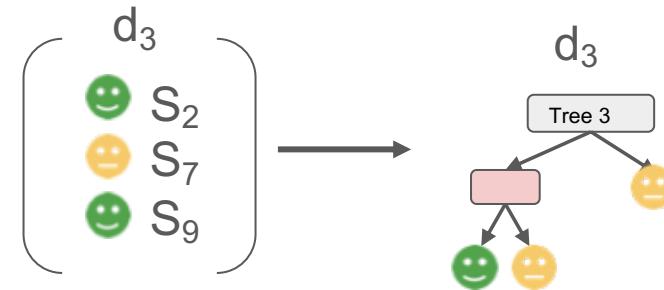
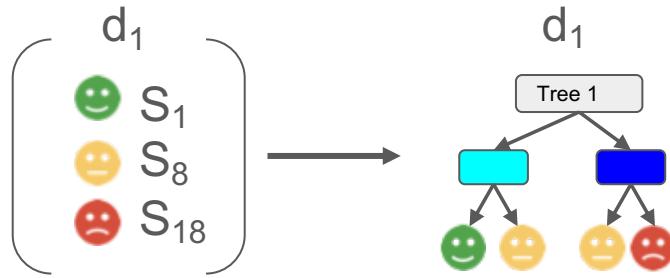
- Number of trees
- Number of features at each split

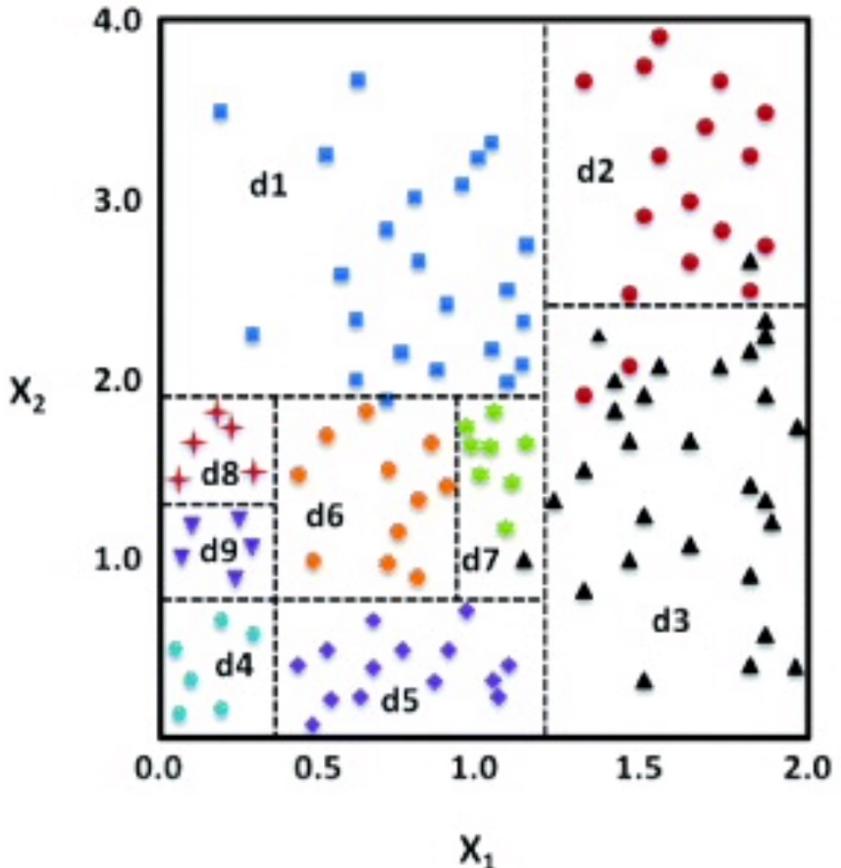


Random forest

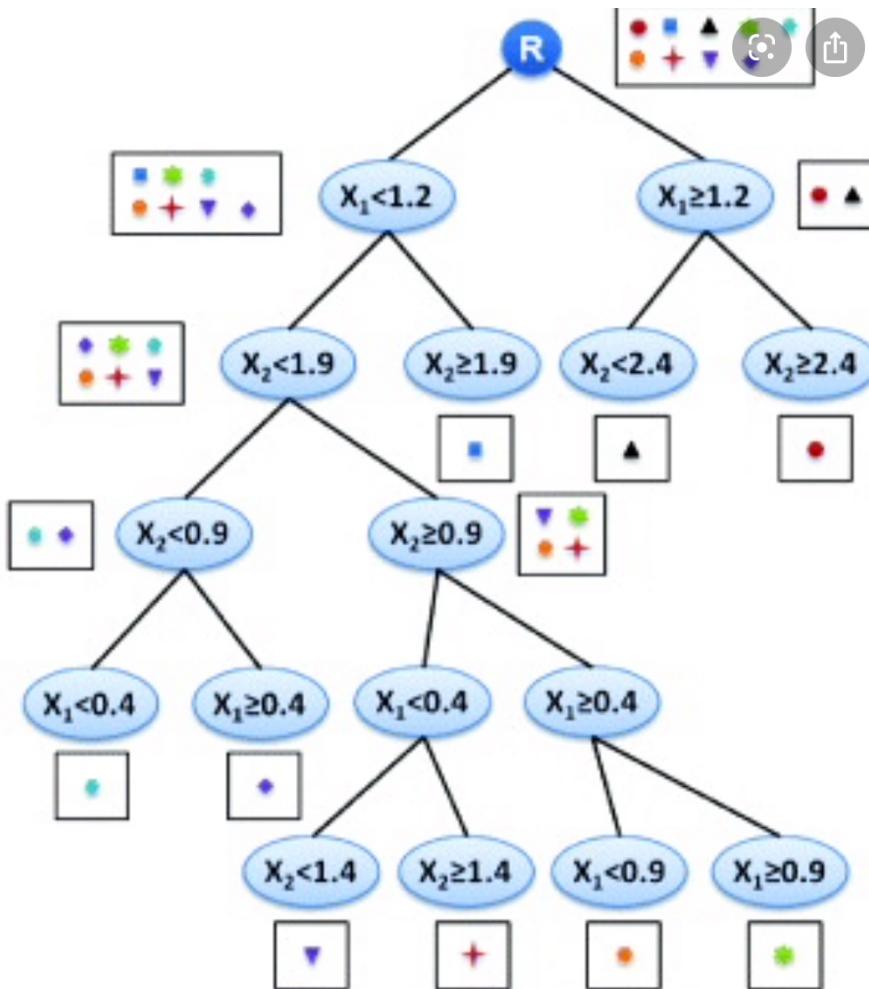


Random forest - training





(a)

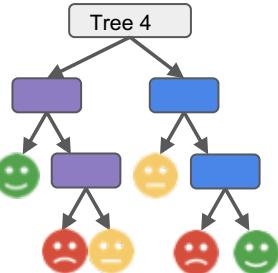
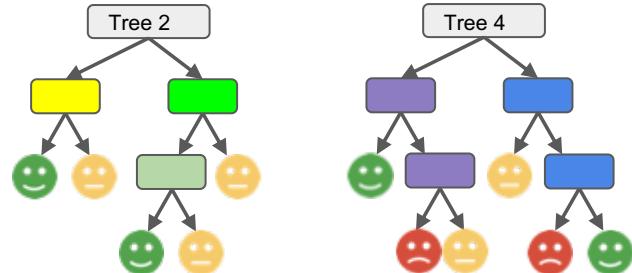
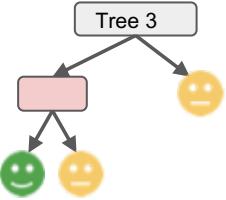
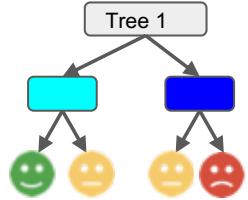


(b)

Random forest - predicting



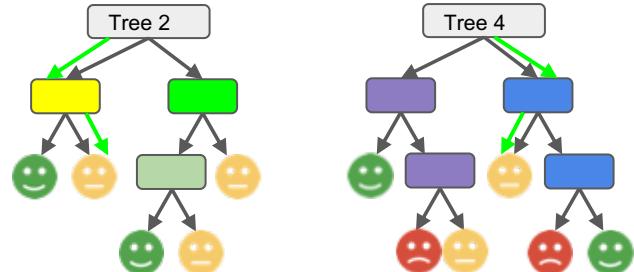
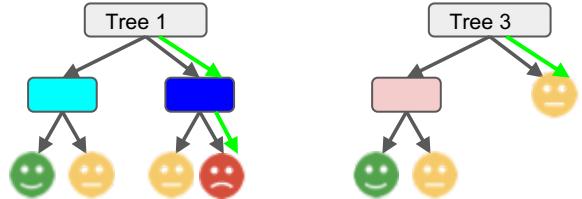
$$fA_p + fB_p + fC_p + fD_p + fE_p + fF_p$$



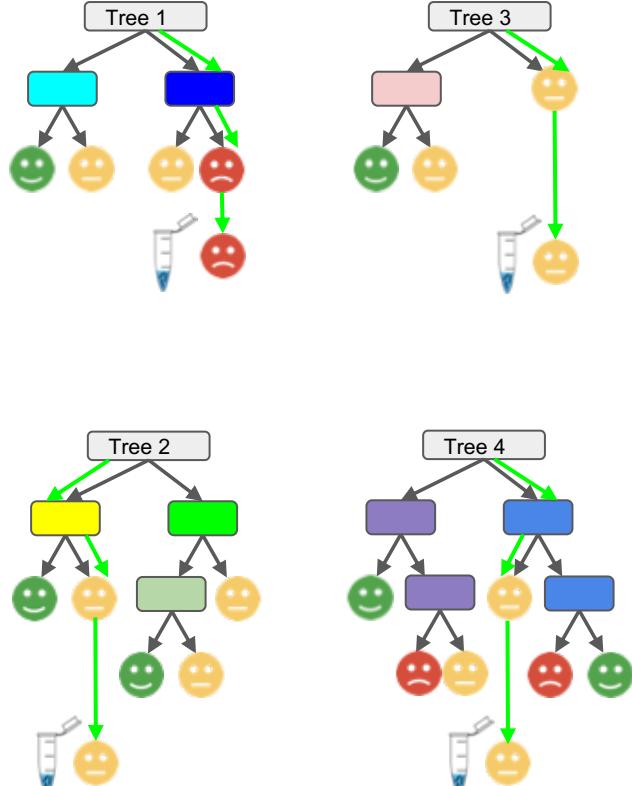
Random forest - predicting



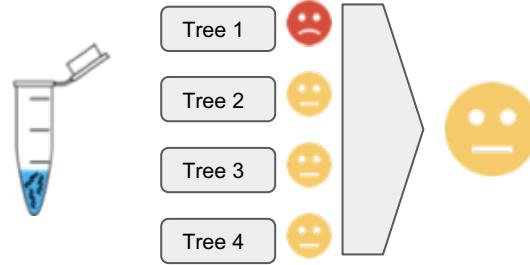
$$fA_p + fB_p + fC_p + fD_p + fE_p + fF_p$$



Random forest - predicting



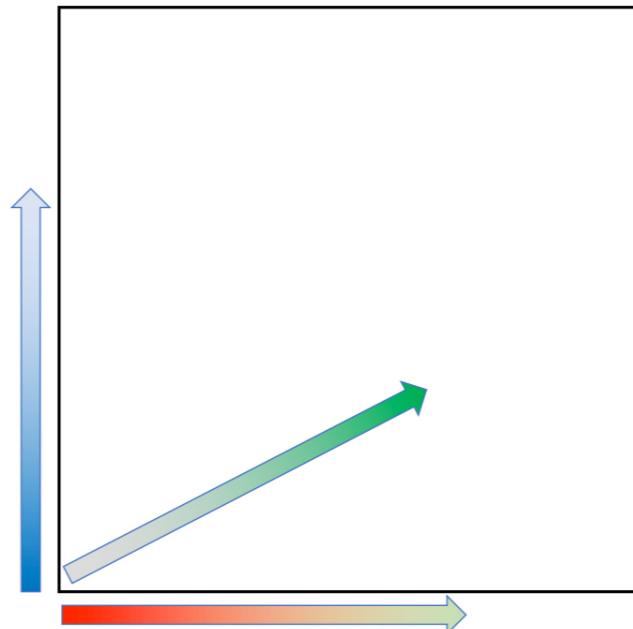
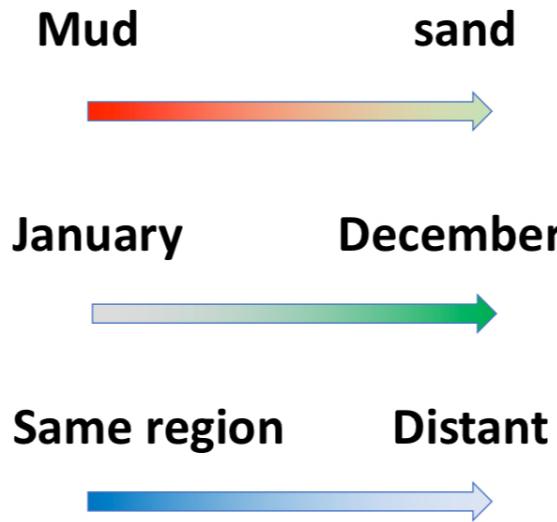
$$fA_p + fB_p + fC_p + fD_p + fE_p + fF_p$$



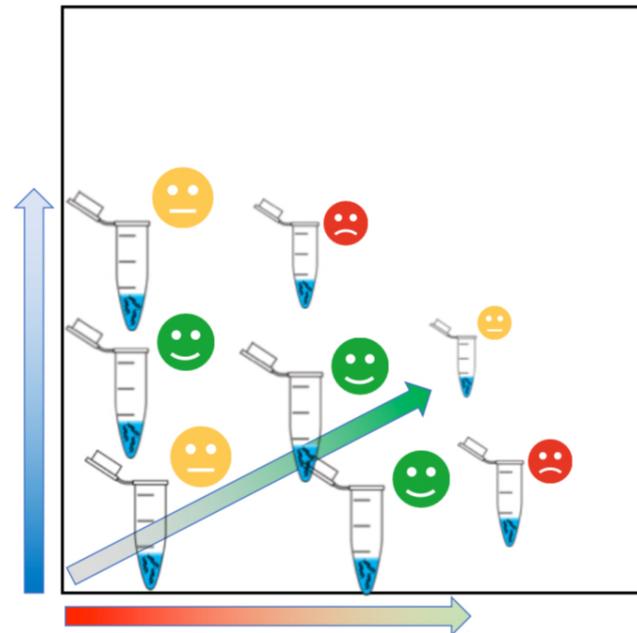
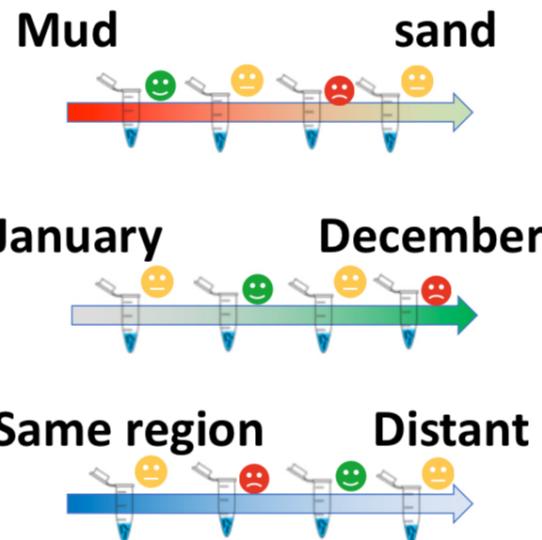
Collecting a *good* training set

- Know your system
 - what is my reference? concentration of pollutant? Natural / anthropised? An index?
 - Does it make sense to make so much effort? Sometimes easier to measure the pollutant
- Cover the gradient in a balanced effort and account for confounding factors (season, landscape attributes, geography...)
 - Your sampling, molecular and metadata collection protocols cannot change afterwards!
 - Have backups! Samples, eDNA and data...

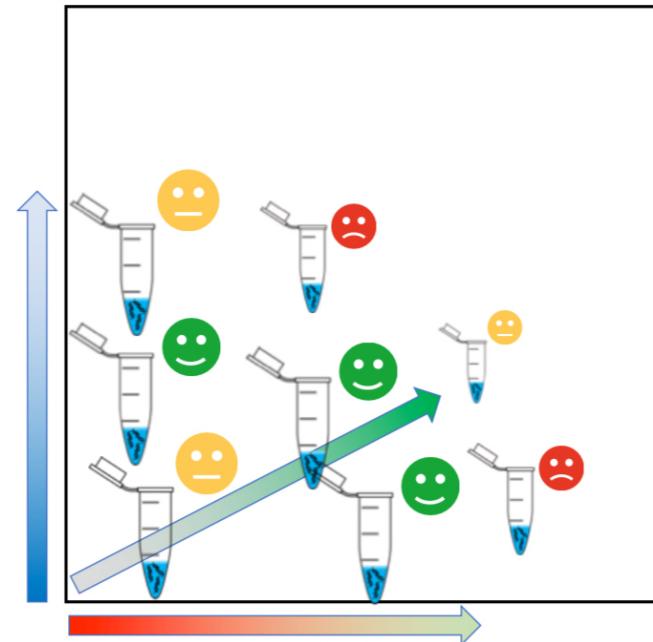
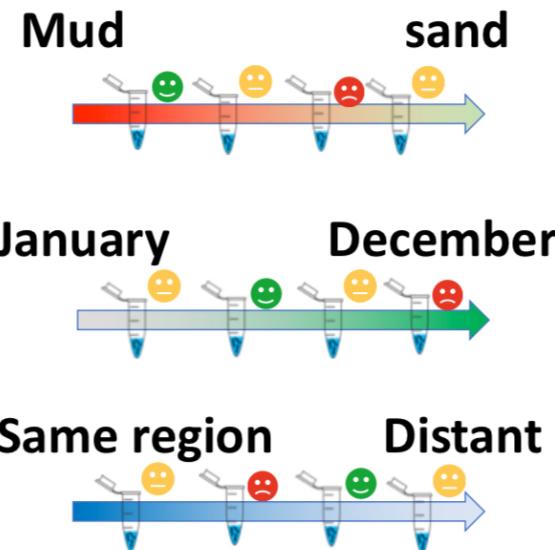
The SML compensates for the variations of benthic communities due to the biogeography, seasonality, or bottom features



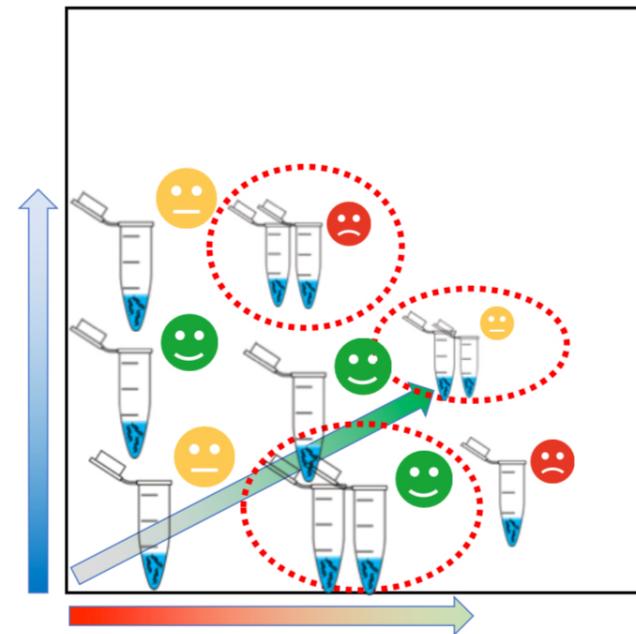
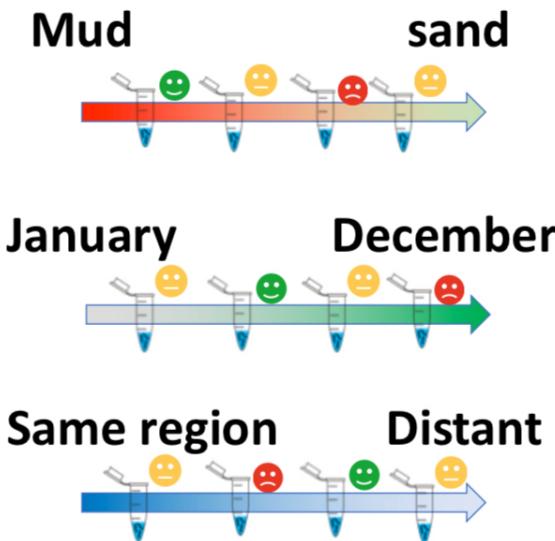
The SML compensates for the variations of benthic communities due to the biogeography, seasonality, or bottom features



The SML compensates for the variations of benthic communities due to the biogeography, seasonality, or bottom features



However, the SML requires a good range of sampling covering different regions, bottom conditions (and seasons)



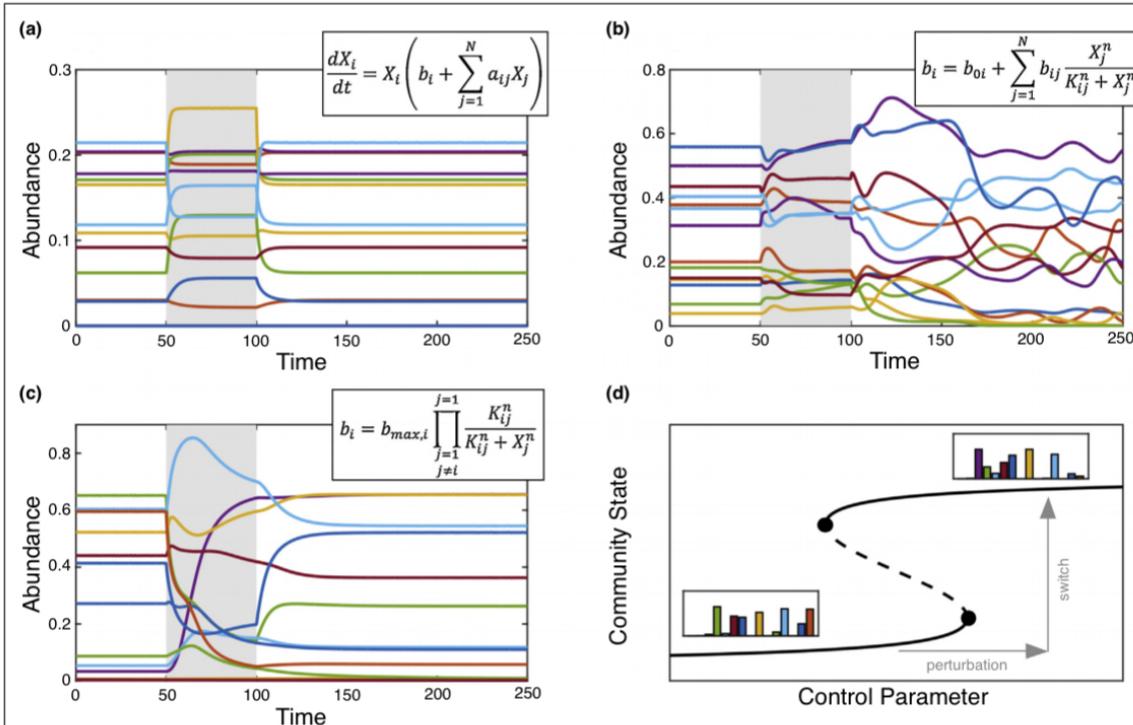


Microbial communities as dynamical systems

Didier Gonze^{1,2}, Katharine Z Coyte^{3,4}, Leo Lahti^{5,6,7} and Karoline Faust⁵



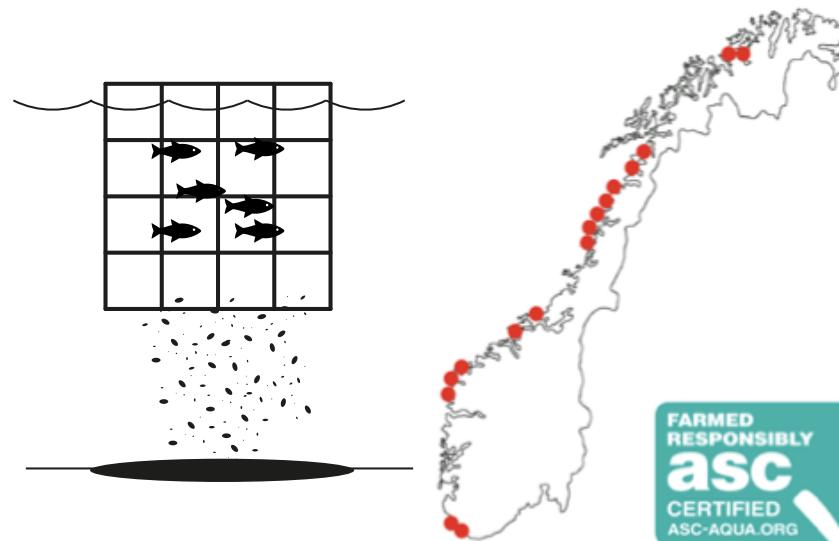
Figure 3



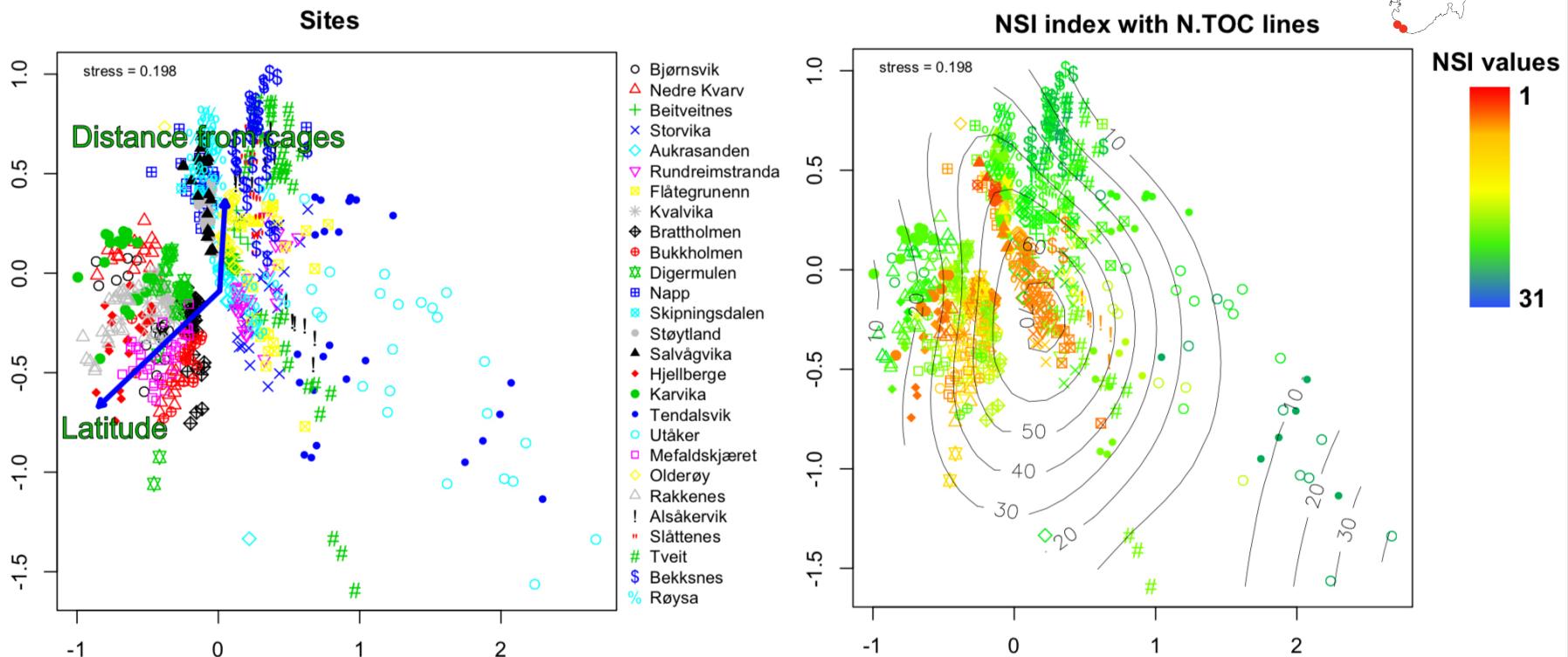
Testing the approach in a marine system



FISKERI- OG HAVBRUKNÆRINGENS
FORSKNINGSFOND



Controlling for other factors

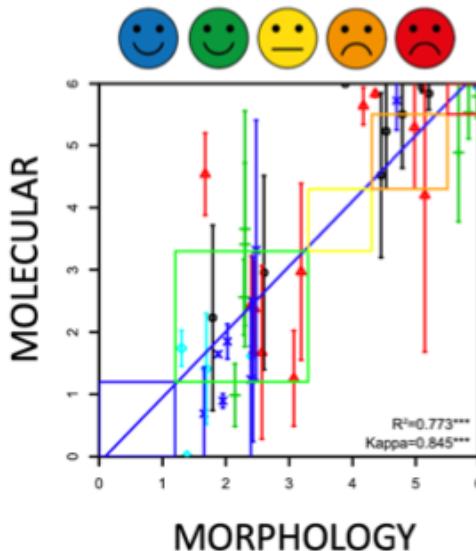


Taxonomy-based approach

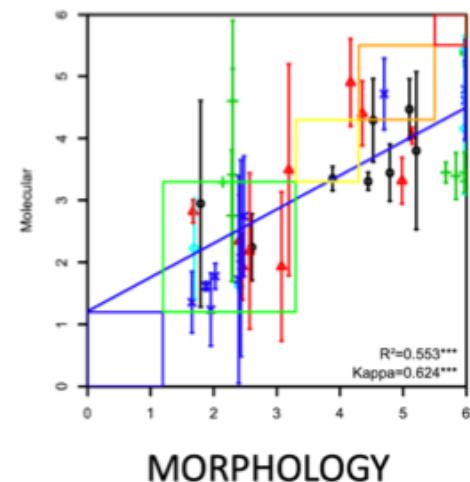
eDNA-taxonomy



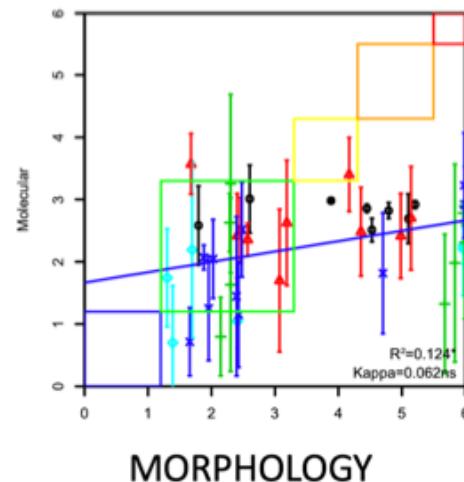
Eukaryotes V4



Eukaryotes V1V2



Eukaryotes V9



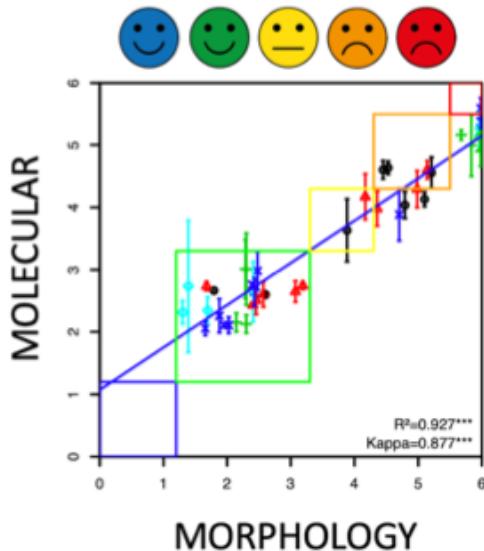
10-20% of
the data

De novo approach

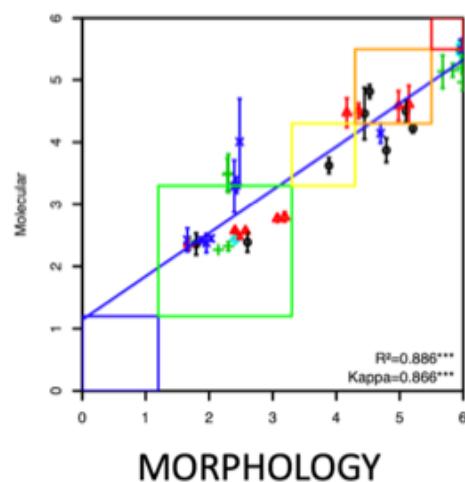
eDNA-ML



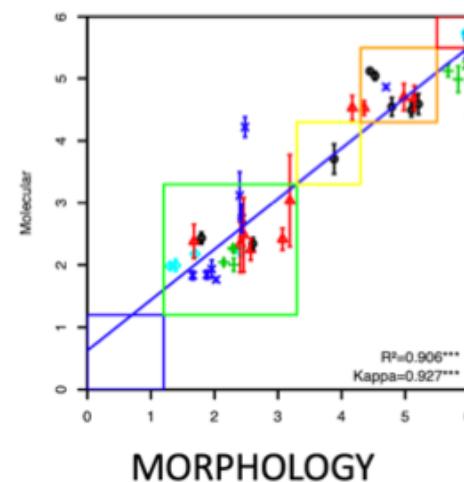
Eukaryotes V4



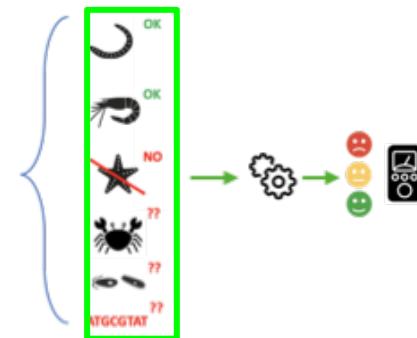
Eukaryotes V1V2



Eukaryotes V9



100% of
the data

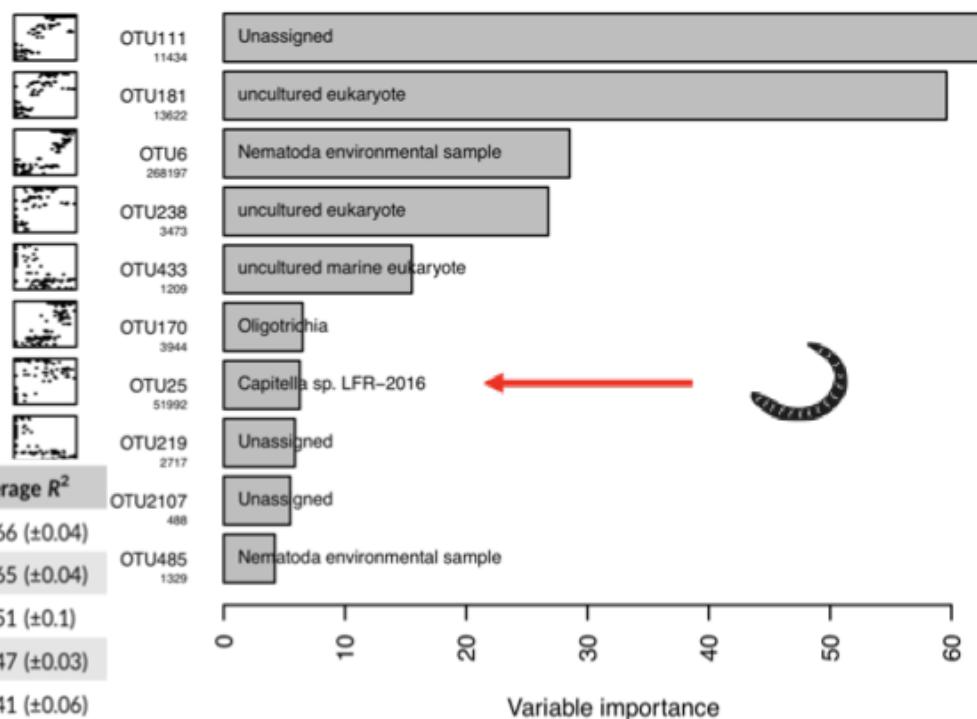


Identifying new bioindicators

Best ones are microbial!

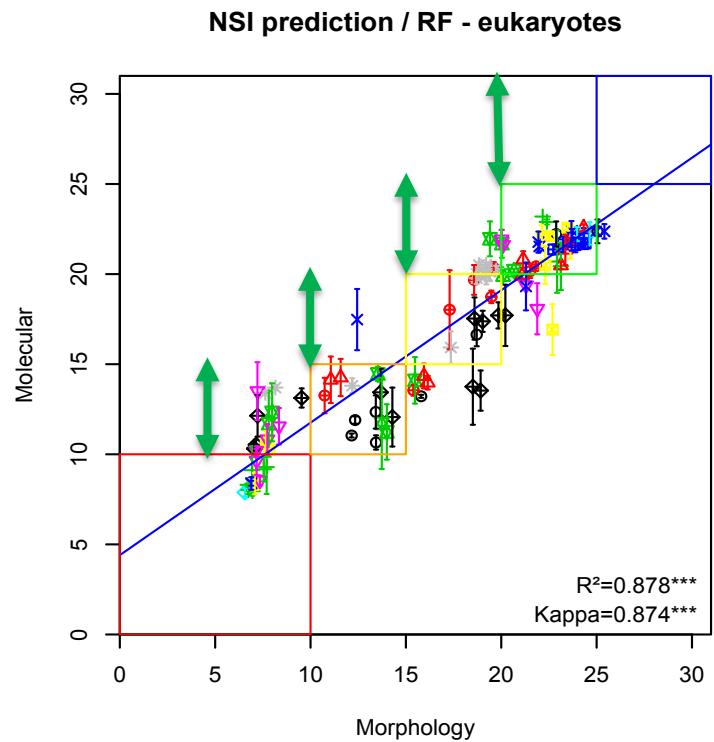
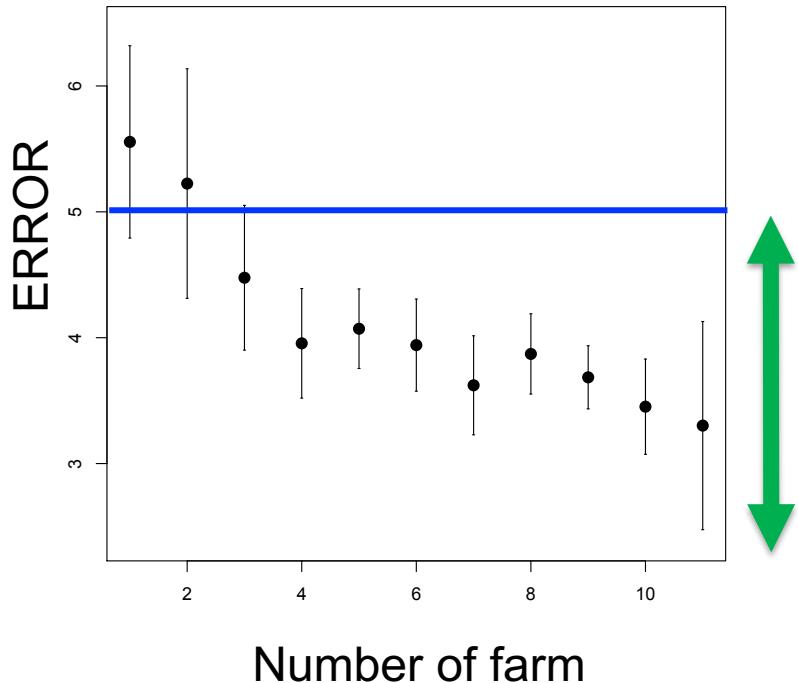


Marker	Taxonomic subgroups	OTUs	Percentage	Average R ²
Bacteria	Gemmatimonadetes	29	0.59	0.866 (± 0.04)
Bacteria	Proteobacteria	1,439	40.82	0.865 (± 0.04)
Bacteria	Planctomycetes	253	4.18	0.851 (± 0.1)
Bacteria	Actinobacteria	170	8.86	0.847 (± 0.03)
Bacteria	Verrucomicrobia	108	2.03	0.841 (± 0.06)
Bacteria	Acidobacteria	180	4.8	0.822 (± 0.07)
Eukaryotes V4	Apicomplexa	20	0.36	0.818 (± 0.09)
Eukaryotes V1V2	Apicomplexa	28	0.27	0.78 (± 0.11)
Metazoa V1V2	Nematoda	157	17.28	0.766 (± 0.06)
Bacteria	Chlorobi	7	0.23	0.766 (± 0.09)
Eukaryotes V9	Charophyta	13	0.15	0.718 (± 0.09)
Bacteria	Nitrospirae	9	0.62	0.717 (± 0.11)
Metazoa V9	Nematoda	18	6.16	0.694 (± 0.12)

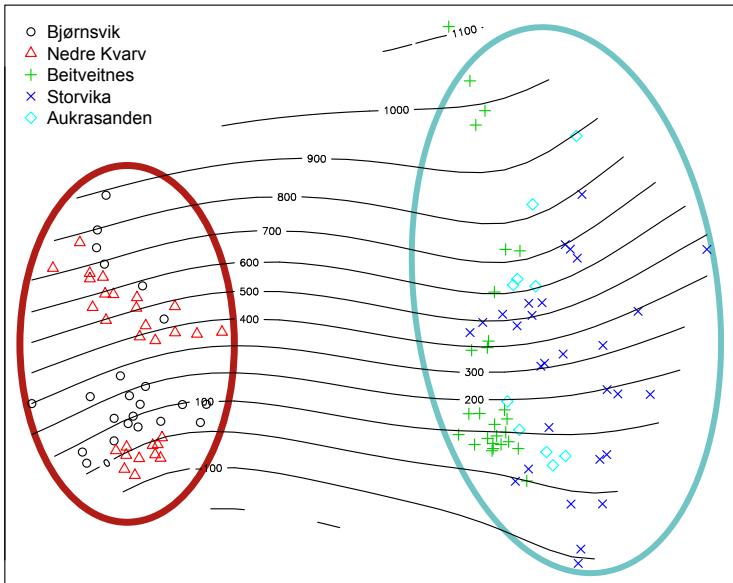
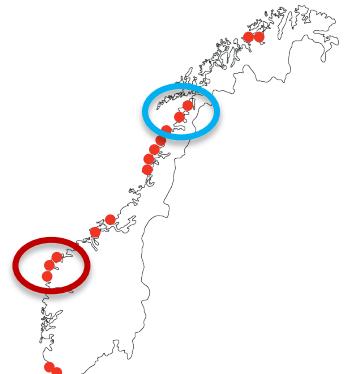


eDNA-based biomonitoring
should focus on microbial
communities

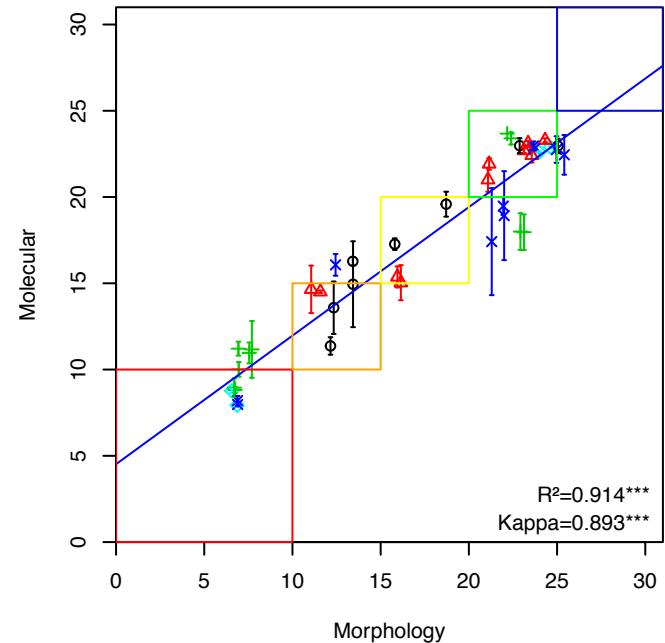
How many farms do we need ?



ML can compensates for geography...

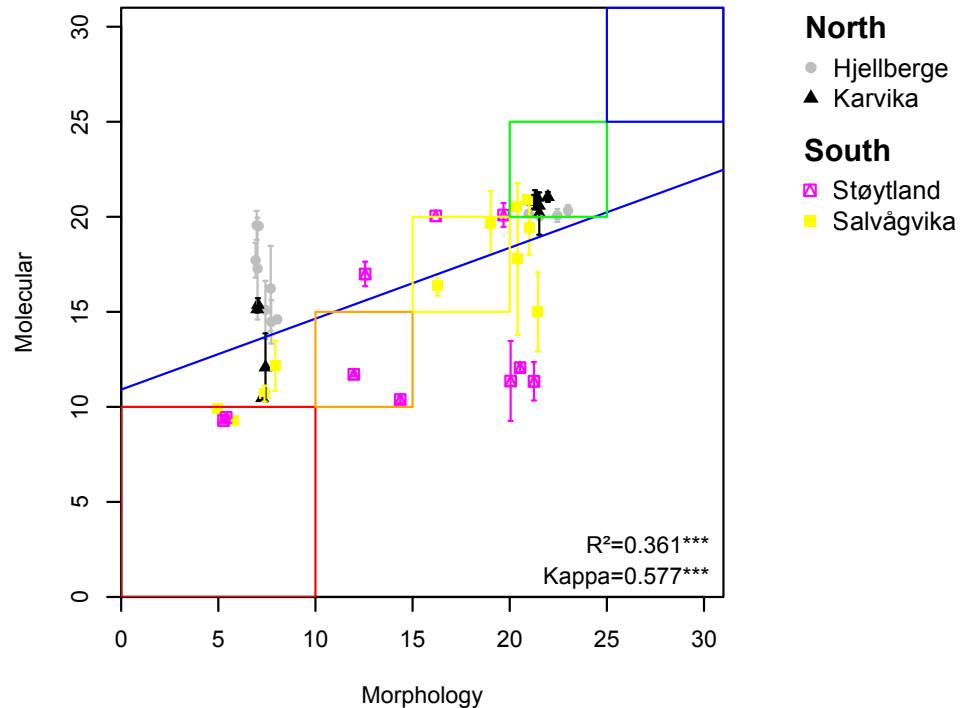
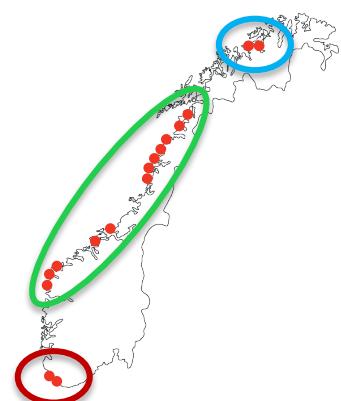


NSI prediction / RF – V1V2 eukaryotes All



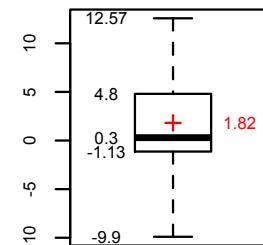
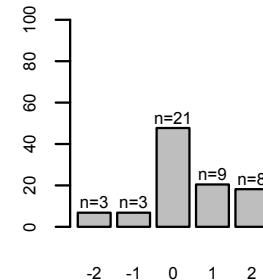
... to some extent

NSI prediction / RF - eukaryotes

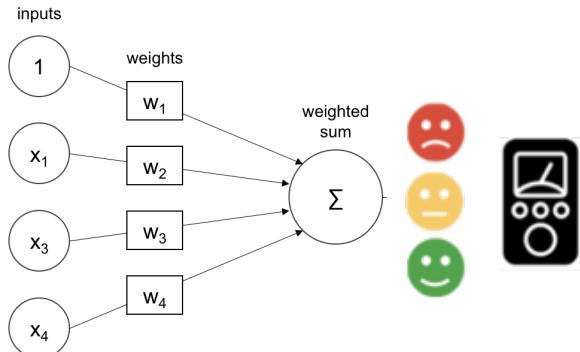
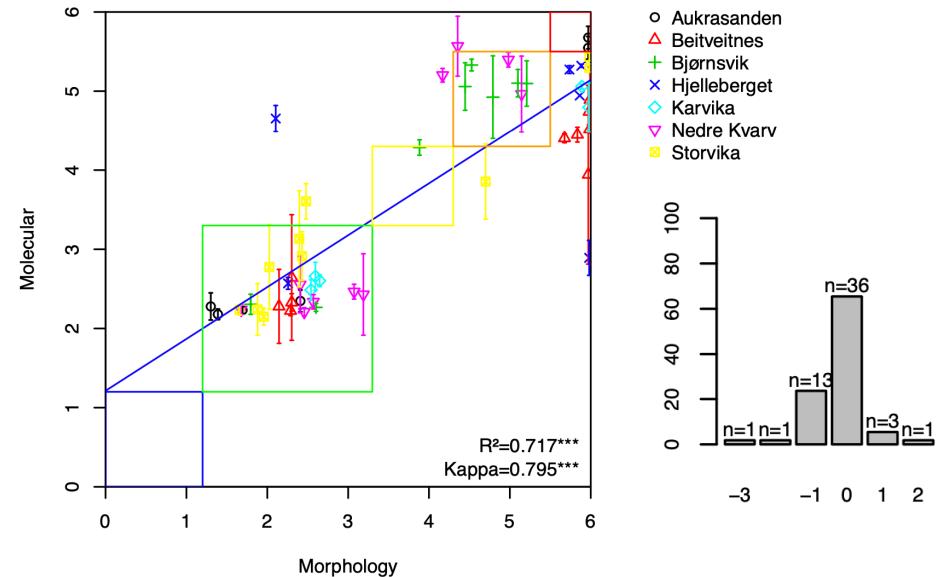
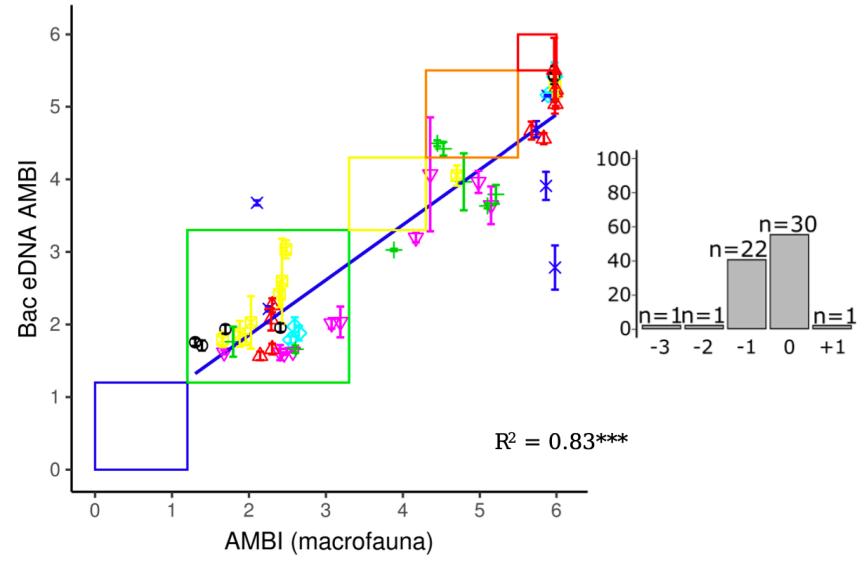


North
● Hjellberge
▲ Karvika

South
■ Støyland
■ Salvågvika



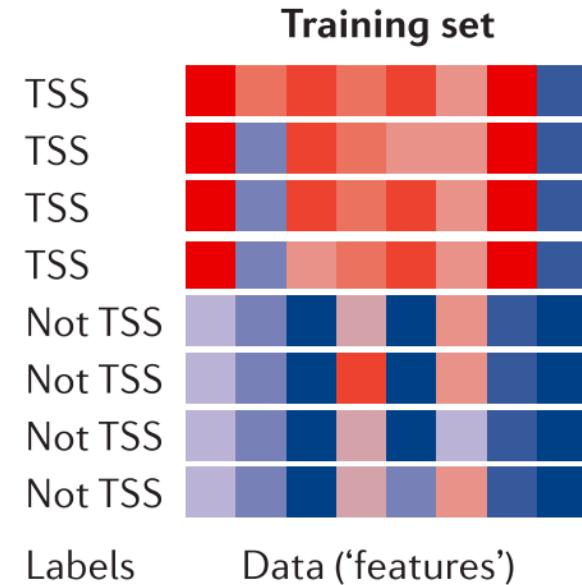
→ Increasing the spatial coverage of the training set might resolve this



Frühe et al., submit.

OTU-based indices VS Supervised Machine Learning

OTU_ID	Weight	Taxa
OTU252	5	Gastrotricha 1
OTU343	2	Nematoda 1
OTU665	4	????
OTU292	3	Bacteria 1
OTU305	1	Nematoda 2
OTU152	2	????
OTU259	4	Bacteria 2
OTU149	2	Fungi 1
OTU1621	3	Fungi 2
OTU1382	2	Bacteria 3
OTU1648	5	???



- + Easily interpretable
- + Fit the current standards and regulations
- Less accurate
- One-by-one approach

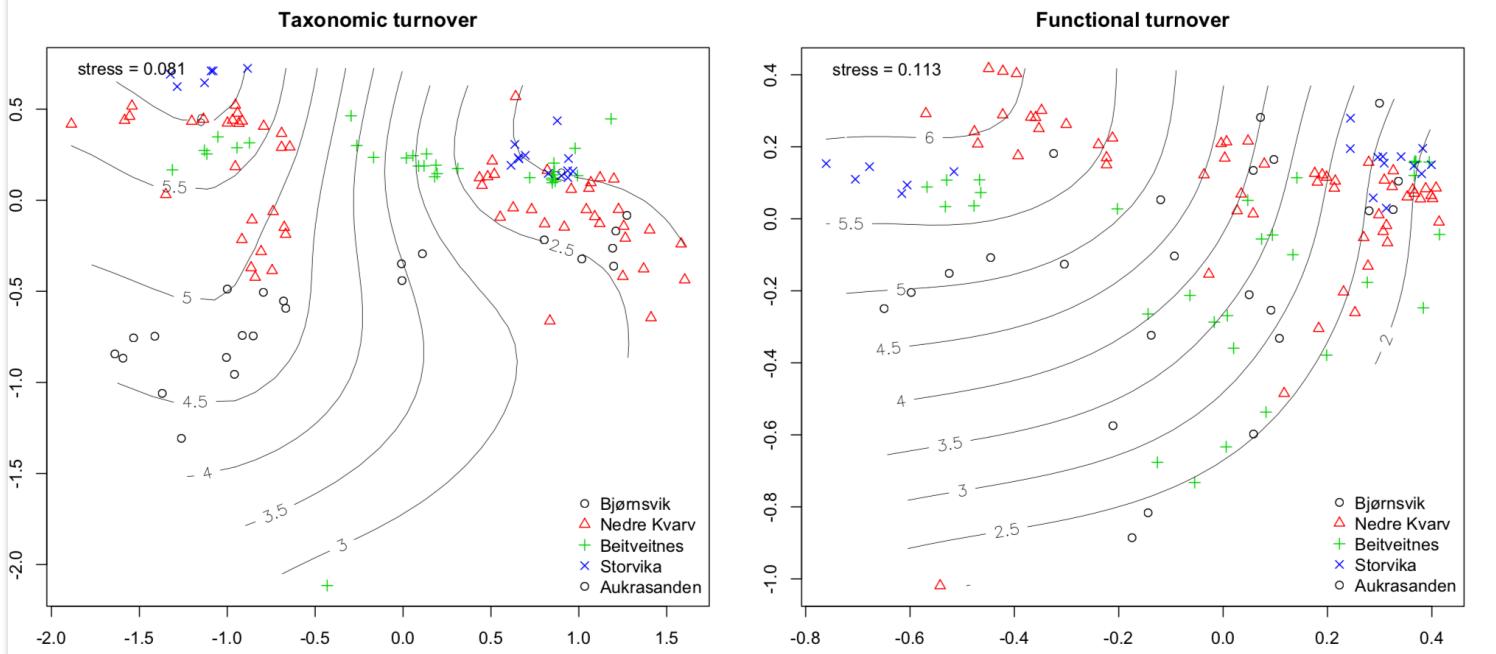
- + More accurate
- + Can identify meaningful co-occurrences
- Less easily interpretable
- Does not (yet) fit current standards

Who is there vs what they do?

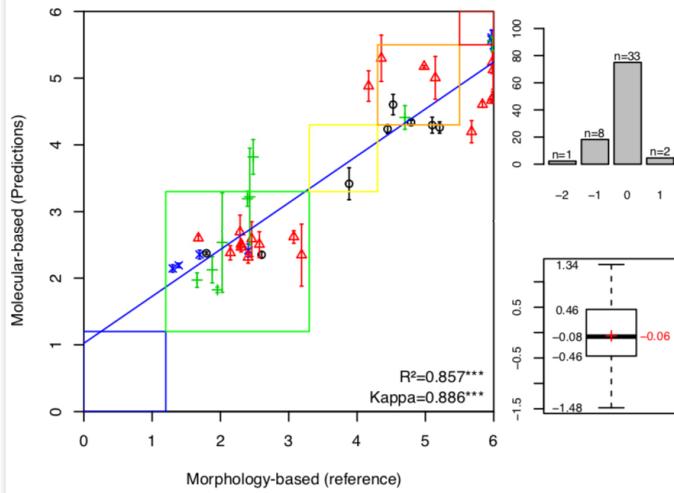
Sequence analysis

Tax4Fun: predicting functional profiles from metagenomic 16S rRNA data

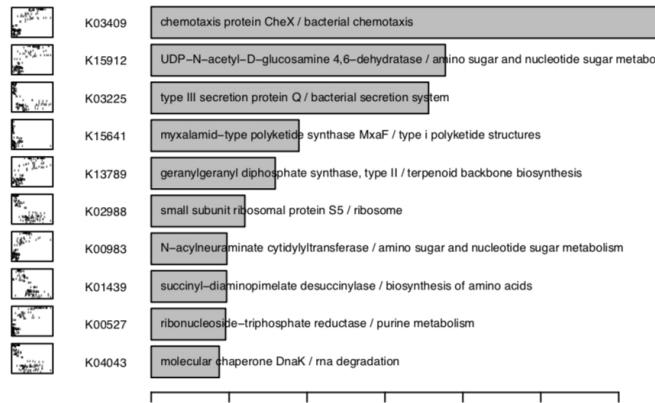
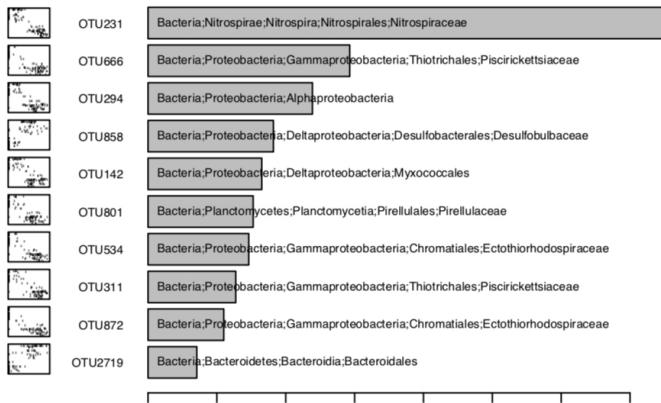
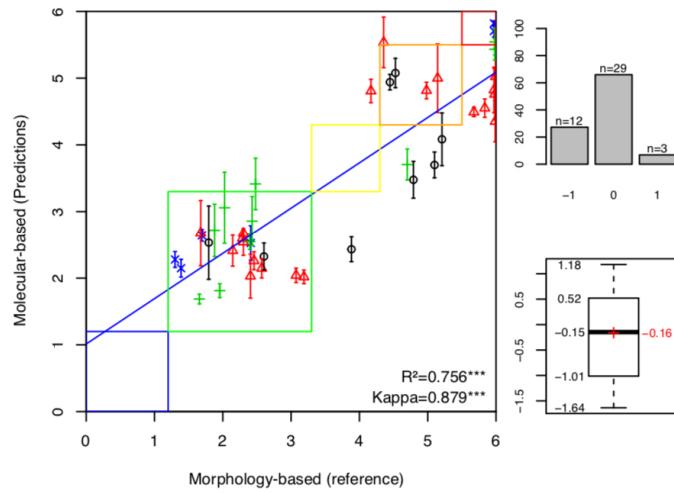
Kathrin P. Aßhauer^{1,*}, Bernd Wemheuer², Rolf Daniel² and Peter Meinicke¹



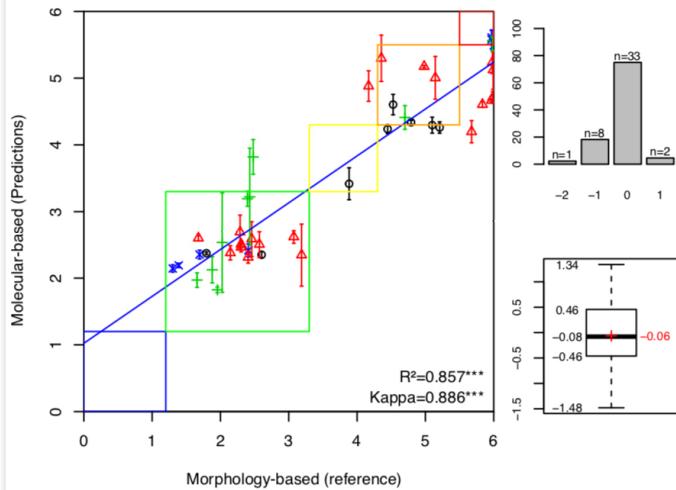
Taxonomic turnover



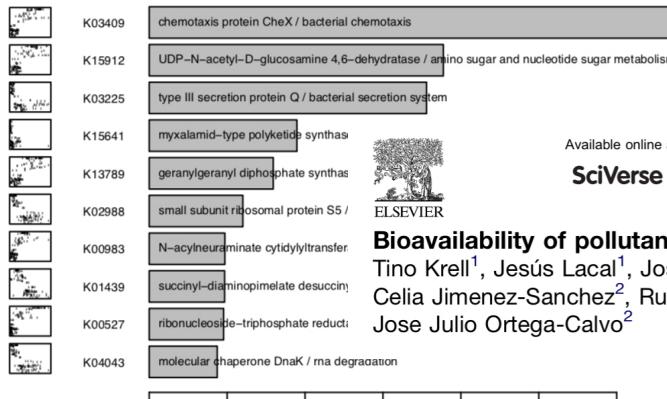
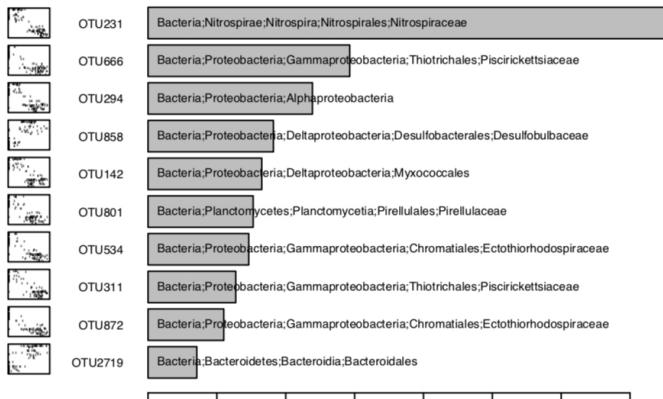
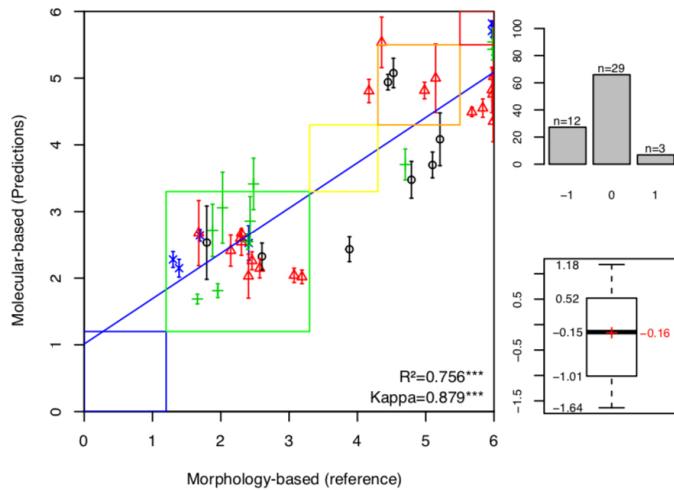
Functional turnover



Taxonomic turnover



Functional turnover



Available online at www.sciencedirect.com

SciVerse ScienceDirect

Current Opinion in
Biotechnology

Bioavailability of pollutants and chemotaxis

Tino Krell¹, Jesús Lacal¹, Jose Antonio Reyes-Dariás¹, Celia Jimenez-Sánchez², Rungroch Sungthong² and Jose Julio Ortega-Calvo²



Untitled.R * Untitled6 * Untitled4 * Untitled5 * dada2_table_FHF.R * metaBTutorial.R*

Source on Save Run Source

```
1 #####
2 #####
3 ### 9th Metabarcoding school - Colombia
4 ### AI/ML analysis of metabarcoding data for biomonitoring - 1st-6th November 2019
5 # Tristan Cordier
6 #####
7 ##### Setting the stage #####
8
9 # setting the working directory
10 #setwd("/path/to/the/folder/metaB_ml")
11 setwd("~/Desktop/POGO_ML_workshop_genomics/Metabarcoding school/")
12
13 # libraries we need
14 library('vegan')
15 library('matrixStats')
16 library('ranger')
17 library('irr')
18 ## for parallel CPU computation
19 library('doMC')
20 registerDoMC(cores = detectCores())
21
22 ##### Playing and manipulating the data #####
68:10 Playing and manipulating the data
```

Console Jobs

```
<--> # and do they match exactly?
> table(gsub("-", ".", met[,s]) == rownames(otu))

TRUE
629
>
> ### keep the dataset at this point for coming back to it
> otu_raw <- otu
>
> ### sum of samples
> row_S <- rowSums(otu_raw)
> plot(row_S, main="Sequencing depth - limit at 10,000 reads", ylab = "Reads counts", xlab = "sample ID",
+ pch=16, cex=.4)
> abline(h=10000, col="blue", lwd=2)
> plot(log(row_S+1), main="Sequencing depth - limit at 10,000 reads", ylab = "log(Reads counts + 1)", xlab =
+ "sample ID", yaxt="n", pch=16, cex=.4)
> axis(2, at=log(c(10,100,1000,10000,100000,1000000)),labels=c(10,100,1000,10000,100000,1000000),cex.axis=
+ 0.7, las=2, ylim=log(c(10,1000000)))
> abline(h=log(10000+1), col="blue", lwd=2)
> |
```

Environment History Connections

Global Environment

Data

- met 629 obs. of 15 variables
- otu 629 obs. of 4903 variables
- otu_raw 629 obs. of 4903 variables
- taxo Large matrix (9806 elements, 662.9 Kb)

Values

- row_S Named num [1:629] 140042 150430 74820 42 6...
- s "samples_names"
- seq_depth_cuto... 10000

Functions

- bar_plot function (otu_table, comp_, aggreg = F, ...)
- plot_ml Large function (989.1 Kb)
- sml_compo Large function (785.3 Kb)

Files Plots Packages Help Viewer

Zoom Export

Sequencing depth - limit at 10,000 reads

log(Reads counts + 1)

sample ID

Department of Genetics and Evolution

Molecular Systematics and Environmental Genomics

- Jan Pawlowski
- Maria Holzmann
- Laure Apothéloz Perret Gentil
- Emanuella Reo
- Florian Mauffrey
- Kristina Cermacova
- Sofia Caetano Wyler
- Léo Charvoz
- Jade Boutten
- Brian Bourrat

University of Kaiserslautern, Germany

Ecology group

- Thorsten Stoeck
- Dominik Forster
- Larissa Frühe

AZTI-tecnalia, Spain

Coastal monitoring group

- Anders Lanzén
- Angel Borja
- Naiara Rodríguez-Ezpeleta

