

## Introduction to Machine Learning

Two-day workshop  
Duke University  
Adam Chekroud



### Background

- BA Experimental Psych, MSc Neuroscience (Oxford)
- PhD Psychology (Neuroscience)

### Interests

- Computational approaches to psychiatric illness
  - Optimising treatment selection in depression

## Objectives

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(\cdot)$   
Algorithm:

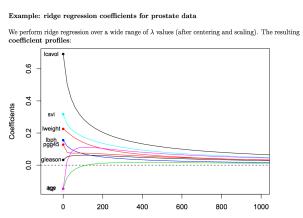
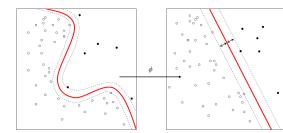
- Initialize model with a constant value:  

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma)$$
- For  $m = 1$  to  $M$ :
  - Compute so-called "pseudo-residuals":  

$$r_{im} = -\left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x_i)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n$$
  - Fit a base learner  $h_m(x)$  to pseudo-residuals, i.e. train
  - Compute multiplier  $\gamma_m$  by solving the following one-c $\ell$   

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i))$$
  - Update the model:  

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$
- Output  $F_M(x)$ .



# Objectives

## Thursday

- Introduction to ML
- Pipeline overview
- Feature selection (theory, code, practice)
- Algorithms (theory, code)
- Performance & **~Validation~**
- Review

## Friday

- Practical session and project brainstorming
- Research talk (@ DIBS)

# What is machine learning?

- *"The goal of machine learning is to program computers to use example data or past experience to solve a given problem"*
- *"Machine learning is a subfield of computer science that ... explores the study and construction of algorithms that ... can learn from and make predictions on data"*
- *"... algorithms operate by building a model from example inputs in order to make data-driven predictions or decisions"*
- Developing computer algorithms for doing — usually predicting — stuff.

# What is machine learning?

## *Doing [stuff]:*

- Systems that analyze past sales data to predict customer behavior
- Optimize robot behavior so that a task can be completed using minimum resources
- **Extract knowledge from bioinformatics data**
  - Diverse range of methods and approaches
  - Plays a central (often hidden) role in our life

# Many faces of machine learning

- Classification
- Clustering
- Dimensionality reduction
- Graphical models
- Genetic algorithms
- Deep learning
- Also many others that haven't yet filtered through to biological sciences: recommendation engines, bayesian networks, association learning, similarity matching
- *ML is not just classification!*

## Machine learning

Reinforcement Learning  
A machine learns to interact with its (dynamic) environment (e.g. playing a computer game, driving a car)

"Supervised" Learning

"Unsupervised" Learning

## Machine learning

### "Supervised" Learning

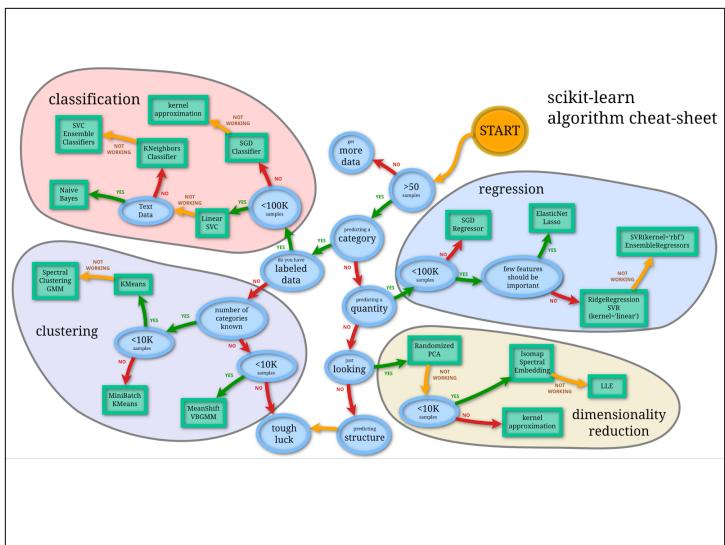
(a.k.a prediction, classification)

- We have labelled training data
- Use data to train algorithm
- Use algorithm to make predictions

### "Unsupervised" Learning

(inc. clustering, PCA)

- No labels on our data
- Use machines to *organise or find hidden structure* in our data
- Use this organisation to help our understanding



## The Perfect Milk Machine: How Big Data Transformed the Dairy Industry



ALEXIS C. MADRIGAL | MAY 1, 2012 | TECHNOLOGY

Dairy scientists are the Gregor Mendels of the genomics age, developing new methods for understanding the link between genes and living things, all while quadrupling the average cow's milk production since your parents were born.

## REVIEWS

### Machine learning applications in genetics and genomics

Maxwell W. Libbrecht<sup>1</sup> and William Stafford Noble<sup>1,2</sup>

**Abstract** | The field of machine learning, which aims to develop computer algorithms that improve with experience, holds promise to enable computers to assist humans in the analysis of large, complex data sets. Here, we provide an overview of machine learning applications for the analysis of genome sequencing data sets, including the annotation of sequence elements and epigenetic, proteomic or metabolomic data. We

Libbrecht & Noble, 2015, *Nat. Rev. Genetics*

## REPORTS

## Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring

T. R. Golub,<sup>1,2,\*†</sup> D. K. Slonim,<sup>1,†</sup> P. Tamayo,<sup>1</sup> C. Huard,<sup>1</sup>  
M. Gaasenbeek,<sup>1</sup> J. P. Mesirov,<sup>1</sup> H. Coller,<sup>1</sup> M. L. Loh,<sup>2</sup>  
J. R. Downing,<sup>3</sup> M. A. Caligiuri,<sup>4</sup> C. D. Bloomfield,<sup>4</sup>  
E. S. Lander<sup>1,5\*</sup>

Although cancer classification has improved over the past 30 years, there has been no general approach for identifying new cancer classes (class discovery) or for assigning tumors to known classes (class prediction). Here, a generic approach to cancer classification based on gene expression monitoring by DNA microarrays is described and applied to human acute leukemias as a test case.

Early example of ML in biology

- Discover new cancer classes
- Assign tumors to classes
- ~11k citations!

Golub et al. 1999, *Science*

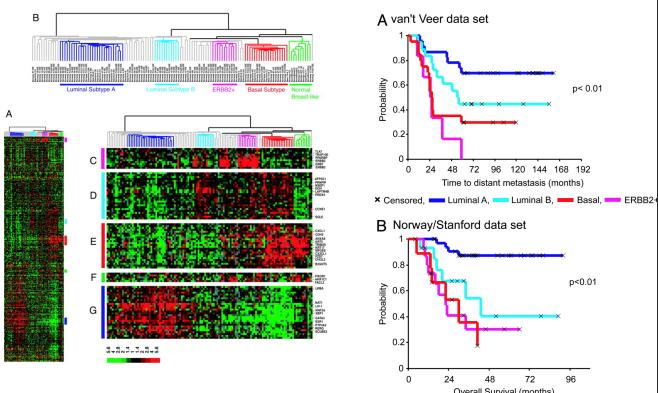
## Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications

Therese Sørlie<sup>a,b,c</sup>, Charles M. Perou<sup>a,d</sup>, Robert Tibshirani<sup>b</sup>, Turid Aas<sup>d</sup>, Stephanie Geisler<sup>b</sup>, Hilde Johnsen<sup>b</sup>, Trevor Hastie<sup>b</sup>, Michael B. Eisen<sup>b</sup>, Matt van de Rijn<sup>c</sup>, Stefanie S. Jeffrey<sup>b</sup>, Thor Thorsen<sup>b</sup>, Hanne Quist<sup>b</sup>, John C. Matese<sup>e</sup>, Patrick O. Brown<sup>e</sup>, David Botstein<sup>b</sup>, Per Eystein Lønning<sup>b</sup>, and Anne-Lise Børresen-Dale<sup>b,n</sup>

## Repeated observation of breast tumor subtypes in independent gene expression data sets

Therese Sørlie<sup>a</sup>, Robert Tibshirani<sup>b</sup>, Joel Parker<sup>b</sup>, Trevor Hastie<sup>b</sup>, J. S. Marron<sup>b</sup>, Andrew Nobel<sup>b</sup>, Shbing Deng<sup>b</sup>, Hilde Johnsen<sup>c</sup>, Robert Pesich<sup>c</sup>, Stephanie Geisler<sup>c</sup>, Janos Demeter<sup>c</sup>, Charles M. Perou<sup>c,d</sup>, Per E. Lønning<sup>c</sup>, Patrick O. Brown<sup>c</sup>, Anne-Lise Børresen-Dale<sup>c,e</sup>, and David Botstein<sup>b,f,g</sup>

Sorlie et al. 2001; 2003, *PNAS*

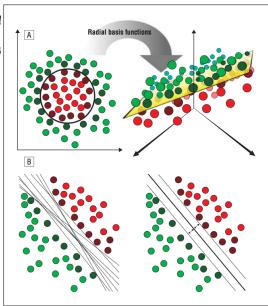


**ORIGINAL ARTICLE**

## Use of Neuroanatomical Pattern Classification to Identify Subjects in At-Risk Mental States of Psychosis and Predict Disease Transition

Nikolaos Koutsouleris, MD; Eva M. Meisenzahl, MD; Christos I. Thomas Frodl, MD; Johanna Schuecker, BA; Gisela Schmitt, Maximilian Reiser, MD; Hans-Jürgen Möller, MD; Christian G.

- Can we use a brain scan to predict conversion to psychosis?
  - Support vector machine classifier (86% accuracy)



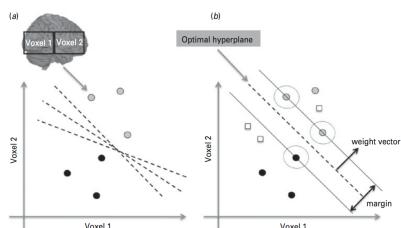
Koutsouleris et al, 2009. Archives

*Psychological Medicine* (2012), 42, 1037–1047. © Cambridge University Press 2011  
doi:10.1017/S0033291711002005

**ORIGINAL ARTICLE**

## Individualized prediction of illness course at the first psychotic episode: a support vector machine MRI study

J. Mourao-Miranda<sup>1,2</sup>, A. A. T. S. Reinders<sup>3,4</sup>, V. Rocha-Rego<sup>1</sup>, J. Lappin<sup>1</sup>, J. Rondina<sup>1</sup>, C. Morgan<sup>5</sup>, K. D. Morgan<sup>6</sup>, P. Fearon<sup>5</sup>, P. B. Jones<sup>5</sup>, G. A. Doody<sup>6</sup>, R. M. Murray<sup>5</sup>, S. Kapur<sup>3</sup> and P. Dazzan<sup>3,7\*</sup>



doi:10.1093/brain/aws084  
**BRAIN**  
A JOURNAL OF NEUROLOGY

Brain 2012; 135; 1508–1521 | 1508

## Multi-centre diagnostic classification of individual structural neuroimaging scans from patients with major depressive disorder

Benson Mwangi,<sup>1</sup> Klaus P. Ebmeier,<sup>2</sup> Keith Matthews<sup>4</sup> and J. Douglas Steele<sup>1</sup>

*Psychological Medicine* (2015), 45, 2805–2812. © Cambridge University Press 2015  
doi:10.1017/S0033291715000768

**ORIGINAL ARTICLE**

## Identifying neuroanatomical signatures of anorexia nervosa: a multivariate machine learning approach

L. Lavagnino<sup>1\*</sup>, F. Amianto<sup>2</sup>, B. Mwangi<sup>1</sup>, F. D'Agata<sup>2</sup>, A. Spalatro<sup>2</sup>, G. B. Zunta-Soares<sup>1</sup>, G. Abbate Daga<sup>2</sup>, P. Mortara<sup>2</sup>, S. Fassino<sup>2</sup> and J. C. Soares<sup>1</sup>

<sup>1</sup>UT Center of Excellence on Mood Disorders, Department of Psychiatry and Behavioral Sciences, UT Houston Medical School, Houston, TX, USA

<sup>2</sup>Department of Neuroscience, AOI San Giovanni Battista, Turin, Italy

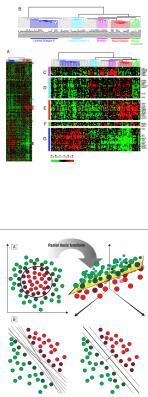
Koutsouleris et al, 2009. Archives

#### Clinical Questions

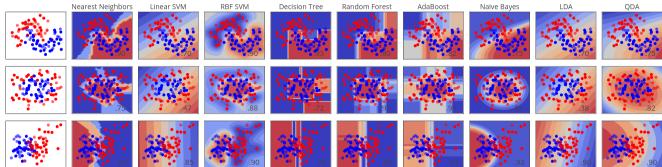
1. Prognosis:
  - A. Can we predict who will develop a disease based on a baseline scan (e.g. fMRI, sMRI)?
2. Treatment Response:
  - A. Can we predict treatment response based on things that we have right now?
    - i) Even better if they are cheap/available! (patient report, blood assay, biomedical test/examination, behavioural task, also genetics, brain scans)
3. Interpretation:
  - A. How to interpret classifiers/regression weights?
  - B. Can we find the most relevant brain regions for diagnosis/prognosis?
  - C. How likely is it that this model is actually real?**
4. (Diagnosis:
  - A. Can we classify groups of subjects (e.g. patients vs. controls) using structural sMRI/fMRI scans?
  - B. Can we combine information from different imaging modalities and/or clinical information?
  - C. Are patients outliers with respect to a "normal population"?)

## Interim summary

- Machine learning is a subfield of computer science concerned with getting computers to do useful stuff for us
- Helpful to consider supervised, unsupervised, and reinforcement learning as the three divisions
  - ML is not just classification! - Computational statistics
- Popular analysis tool in psychiatry

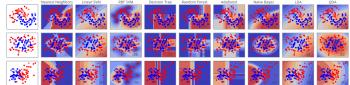


## What is all the fuss about?



- Mechanistic reasons
- Pragmatic reasons

## What is all the fuss about?



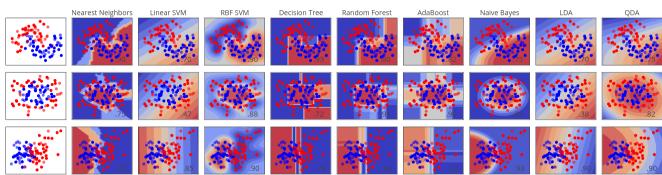
### Mechanistic reasons

- Develop systems for learning *complicated* relationships in data
  - Driven by: larger data sets, more computational resources
  - Allows us to discover new knowledge

### Pragmatic reasons

- Make better predictive tools
  - Often, without domain-specific knowledge (c.f. kaggle!)
- Performance usually outperforms traditional statistical procedures
- ML can replace boring tasks
  - Frees up humans for human things (like treating patients)

## What is all the fuss about?



## Summary

- Machine learning is a subfield of computer science concerned with getting computers to do useful stuff for us
  - ML is not just classification! - Computational statistics
- Popular analysis tool in psychiatry
  - Find complex relationships—more data, more resources
  - Build better tools to improve diagnosis and treatment

