# Machine Learning
# ML. 1. Intro

David Ríos Insua and  Roi Naveiro

# Objectives and schedule

A  broad overview of Machine Learning

- Supervised learning
- Unsupervised learning
- Reinforcement learning

# Machine learning

From Wikipedia

ML: the study of computer algos that improve automatically through experience. It is seen as a part of AI. ML algos build a model based on sample data, known as "training data", in order to make predictions or decisions without being explicitly programmed to do so….

A subset of ML is closely related to computational statistics, which focuses on making predictions using computers; but not all ML is statistical learning.

The study of mathematical optimization delivers methods, theory and application domains to the field of machine learning.

Data mining is a related field of study, focusing on exploratory data analysis through unsupervised learning.

In its application across business problems, machine learning is also referred to as predictive analytics.

# Some ML examples. Red matters!!!

Uncertainty is almost ubiquitous in ML:

- Given a certain transaction, is it fraudulent or not? If fraudulent should I stop it?

- Given the monitoring trace of an Inet device, are we facing an attack? Should I stop operations?

- Does this medical image correspond to a person with a certain illness? Should I make further tests?

- A person with these FB likes will buy this type of beer? Should I send him my brand add?

- A person with these tweets is conservative? Should I send her Brexit propaganda?

- Robots (or ADS): If robot performs this, How will the user react? And the environment? Consequently, what should the robot do?

# In applications, we'll need to go beyond

- Beyond a model with good fit…

- Beyond a model that predicts well…

- Fraud detection. Classification problem
  - Few false positives. FPR
  - Few false negatives. FNR
  - But what really matters are minimising monetary losses!!!

# In applications, we'll need to go beyond

- Beyond a model with good fit…
- Beyond a model that predicts well…
- Fraud detection. Classification problem
  - Few false positives. FPR
  - Few false negatives. FNR
  - But what really matters are minimising monetary losses!!!

- Reservoir system management. Forecasting model for inputs and demands

Feeds decision model e.g to minimize energy deficit, wasted water, given constraints…..

- Aviation safety risk management. Forecasting models for accidents and incidents, as well as their multple impacts

Feed a risk managment model: optimal safety resource allocation gven constraints…

# ML, Stats in modern times (Big Data) Computer Age Statistical Inference

**V**olume. Space scalability

**V**ariety. Text, images, sound, video,… video,…..

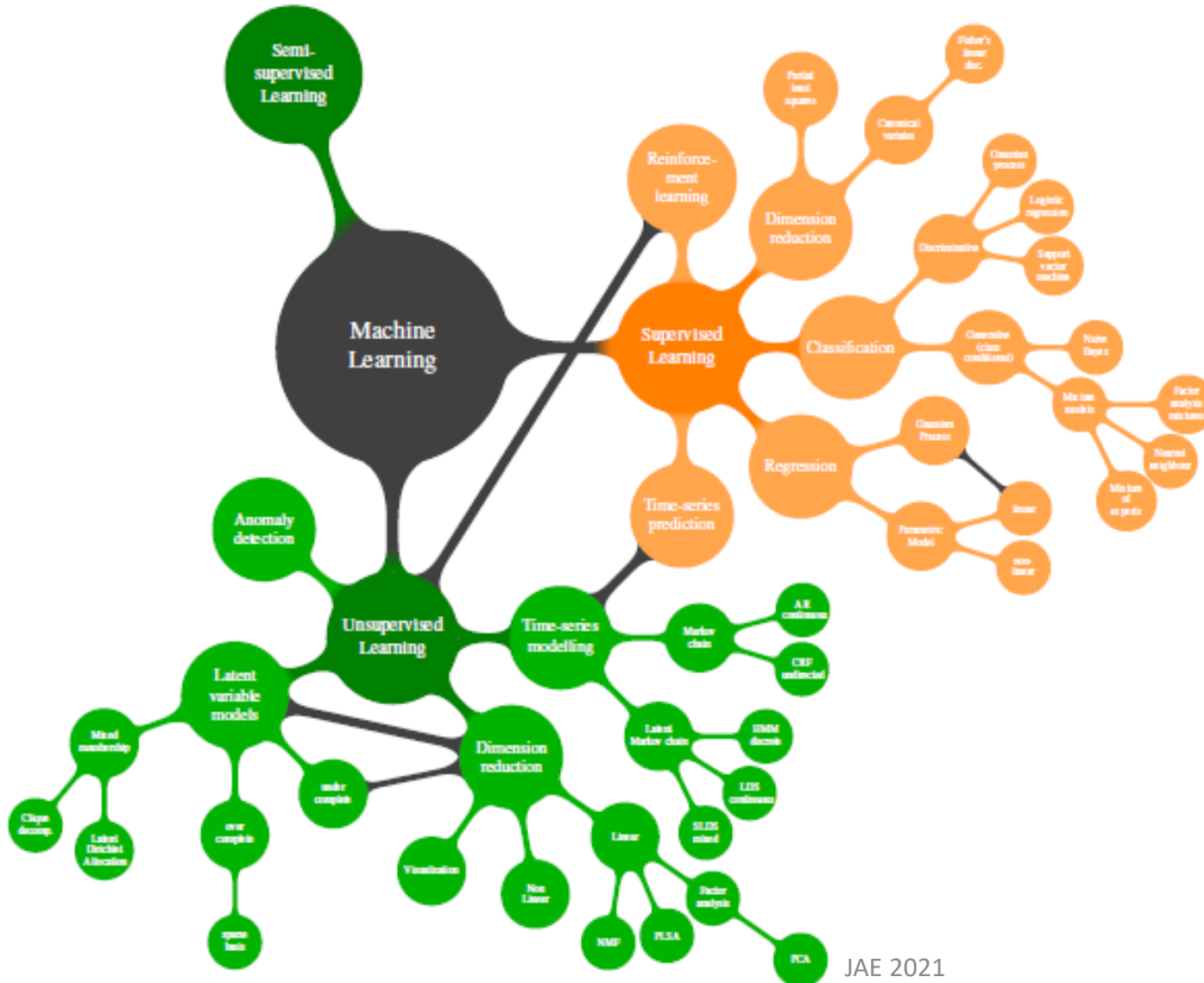**V**elocity. High frequency for data and decisions, time series, dynamic models. Time scalability

Create **v**alue by actually supporting decisions!!!

# Core themes

- Supervised learning: Pairs input-output available
  Regression, Classification

  Case: Cannabinoid detection

- Unsupervised learning: Outputs not available (or the inputs are the outsputs)
  Density estimation, clustering, outlier detection, Visualisation

  Case: Singular community detection

- Reinforcement learning: Decisions impacting outputs on-the-fly
  Markov decision processes

  Case: Autonomous driving systems

# With somewhat blurry borders….



Case: Lookalike modeling

JAE 2021

# A brief KitKat



VS

# A basic example

Consider a cyberattack recovery protocol for SMEs. Introduce a process. Want to assess it. E.g to compare it with another

Test it against 12 attacks, effective 9 times (e.g., system up in less than 1 hour)

# A basic example. Model

- Number X of successes in n trials (ii)
- Success probability in one trial
- Distribution of number of successes
- For X=9,

$$\theta_1$$

$$X \mid \theta_1 \sim Bin(12, \theta_1)$$

$$Pr(X = 9 \mid \theta_1) \propto \theta_1^9 (1 - \theta_1)^3, \quad \theta_1 \in [0, 1]$$

# A basic example. MLE

Likelihood

Log-likelihood

Maximise likelihood or maximize log likelihood . MLE

In this case, MLE is

Defects?

For future observations (e.g. 4 succeses in next 7 trials)

$$l(\theta_1) \propto \theta_1^9 (1-\theta_1)^3$$

$$h(\theta_1) = \log(l(\theta_1)) = 9 \log \theta_1 + 3 \log(1-\theta_1)$$

$$h'(\theta_1) = 0 \implies \hat{\theta}_1 = \frac{9}{12} = .75$$

$$Pr(Y=4 \mid \hat{\theta}_1, 7) = \binom{7}{4} .75^4 .25^3$$

# A basic example. Bayes

Prior, e.g.

Posterior

    Posterior mean

    Posterior mode (MAP)

Predictive

$$\pi(\theta_1) = 1$$

$$\pi(\theta_1|9) \propto 1 \times \theta_1^9 (1-\theta_1)^3 \sim \mathcal{B}e(10,4)$$

$$\frac{10}{14}$$

$$\frac{9}{12}$$

$$Pr(Y=4|9) = \int \binom{7}{4} \theta_1^4 (1-\theta_1)^3 \pi(\theta_1|9) \, d\theta_1$$

$$= \frac{\binom{7}{4}\binom{13}{3}}{\binom{20}{12}}$$

# And so…

We end up using this to make decisions. Which protocol to implement?

How would you choose between two protocols?

# Intro to Supervised Learning

# SL: ingredients

- Data available: examples, samples, instances,…

- Several observed variables: predictors, attributes, features, covariates, explanatory variables, independent variables,…

- Some of special interest: response(s), dependent variable(s), target(s), output(s), label(s),…

# SL: types of problems

1. Regression, response variable is continuous

2. Classification, response variable is discrete

3. Other:
   – Mixed (some continuous, some discrete)
   – Discrete but ordered
   – …

Predictors                                  Dependent variable

$$(x_1, \ldots, x_p)$$                                      $Y$

Some relation

$$Y = f(x) + \epsilon$$

Systematic info        Random term. Zero mean, Indep of x

Inference vs Prediction

# How do we estimate f?

Training data

$$\{(x_1, y_1), \ldots, (x_n, y_n)\}$$

$$\hat{f}(x_i) \approx y_i$$

For any observable

$$\hat{f}(x_0) \approx y_0$$

Parametric, e.g.

$$\hat{f}(x) = \beta_0 + \beta_1 x_1 + \ldots + \beta_p x_p$$

$$\hat{f} \text{ vs } (\beta_1, \ldots, \beta_p)$$

$$y \approx \beta_0 + \ldots + \beta_p x_p$$

Flexibility and overfitting

Non-parametric

Wider range, much larger #observations

# Flexibility vs Interpretability

(somewhat old figure from ISLR)



Rudin's paper!!!!

Deep Models!!!!

# Assessing accuracy

No free lunches

Quality of fit, e.g.

$$MSE = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - \hat{f}(x_i) \right)^2$$

    Training MSE    vs Test MSE

      Not      $\hat{f}(x_i) \approx y_i$        but      $\hat{f}(x_0) \approx y_0$
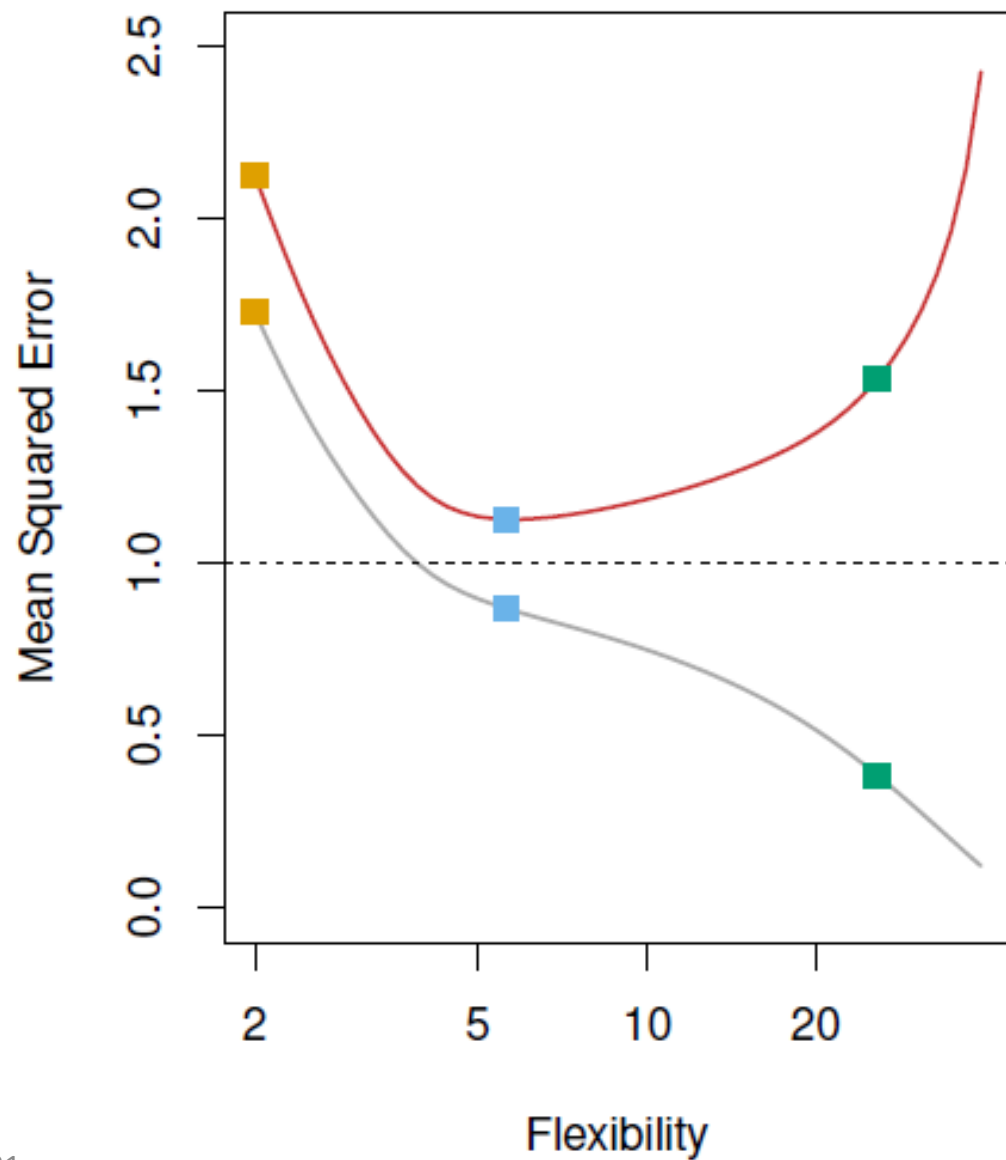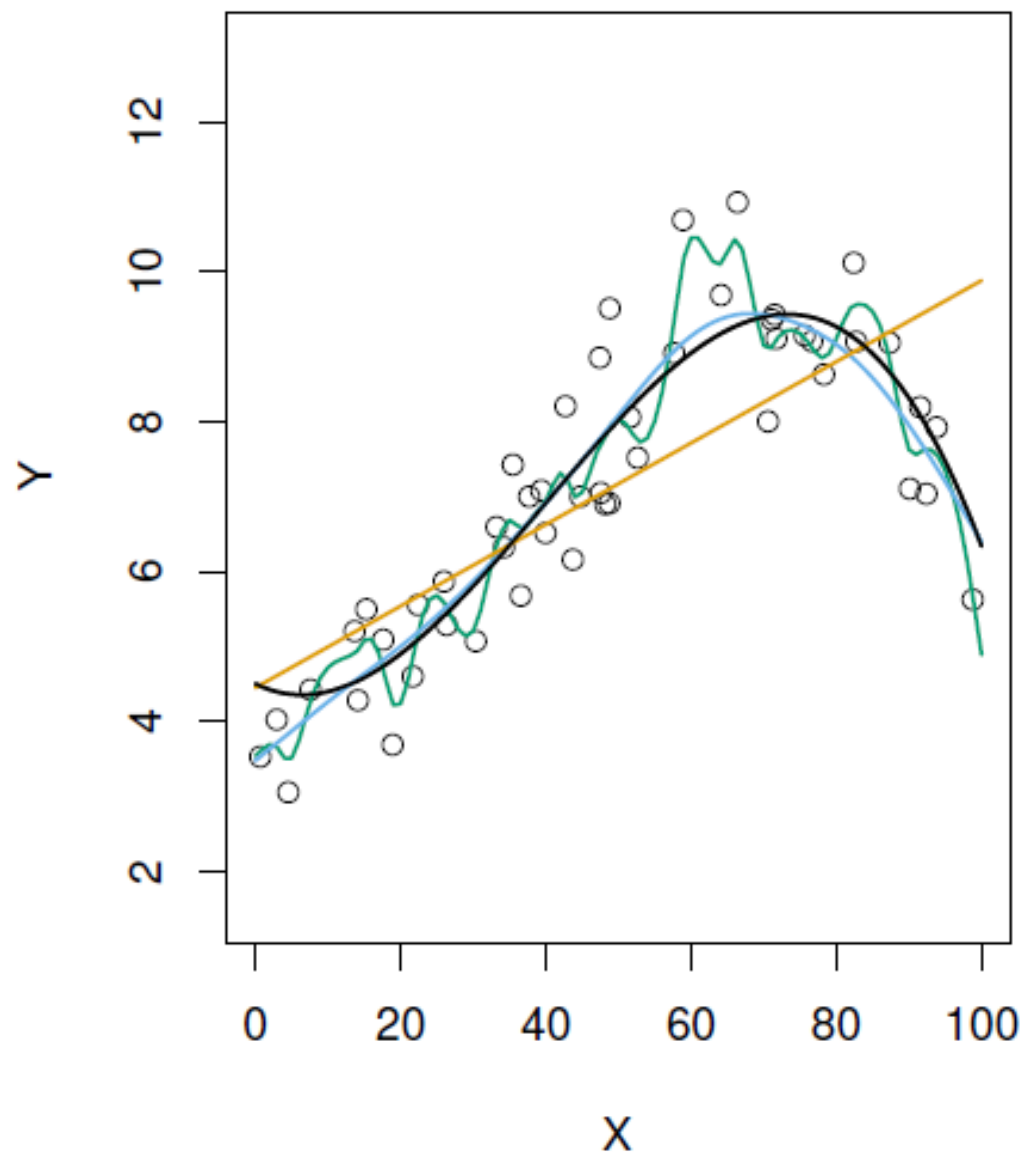
         Small       $AVE \left( y_0 - \hat{f}(x_0) \right)^2$

# Model selection

- Compute empirical risk over training set
- May be reduced almost arbitrarily increasing model complexity

e.g. based on polynomials

- Generalisation error over observations not used to train the model  (cannabinoid project)
- If no test set available, split data in two sets:
  - Training set
  - Test set

# Cross validation

- Hyperparameter choice to control model complexity
- Choose a third set for validation to select and compare models
- If data not plentiful, divide set in k partitions
- Use k-1 to train and the other to test: k models
- Cross validation error



- If k=n  leave-one-out cross validation

# Bias-variance tradeoff

- Assume model

$$Y = f(x) + \varepsilon$$

$$E(\varepsilon) = 0$$
$$Var(\varepsilon) = \sigma^2$$

- Exp. Pred. error (under quad. Loss)

$$EPE = E\left(Y - \hat{f}(x)\right)^2$$

- Decomposed as

Variance. How approximation changes if a different training set used
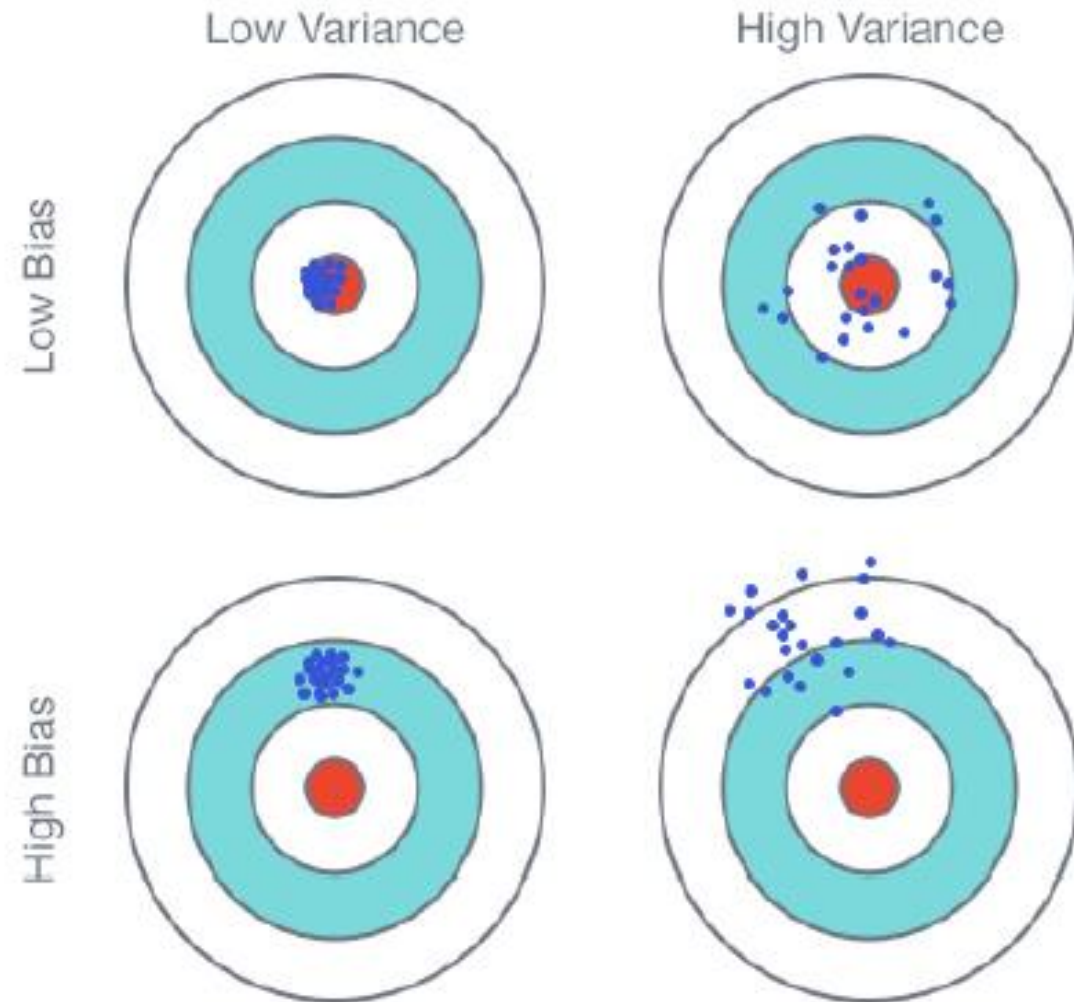
Bias. Error due to using much simpler model

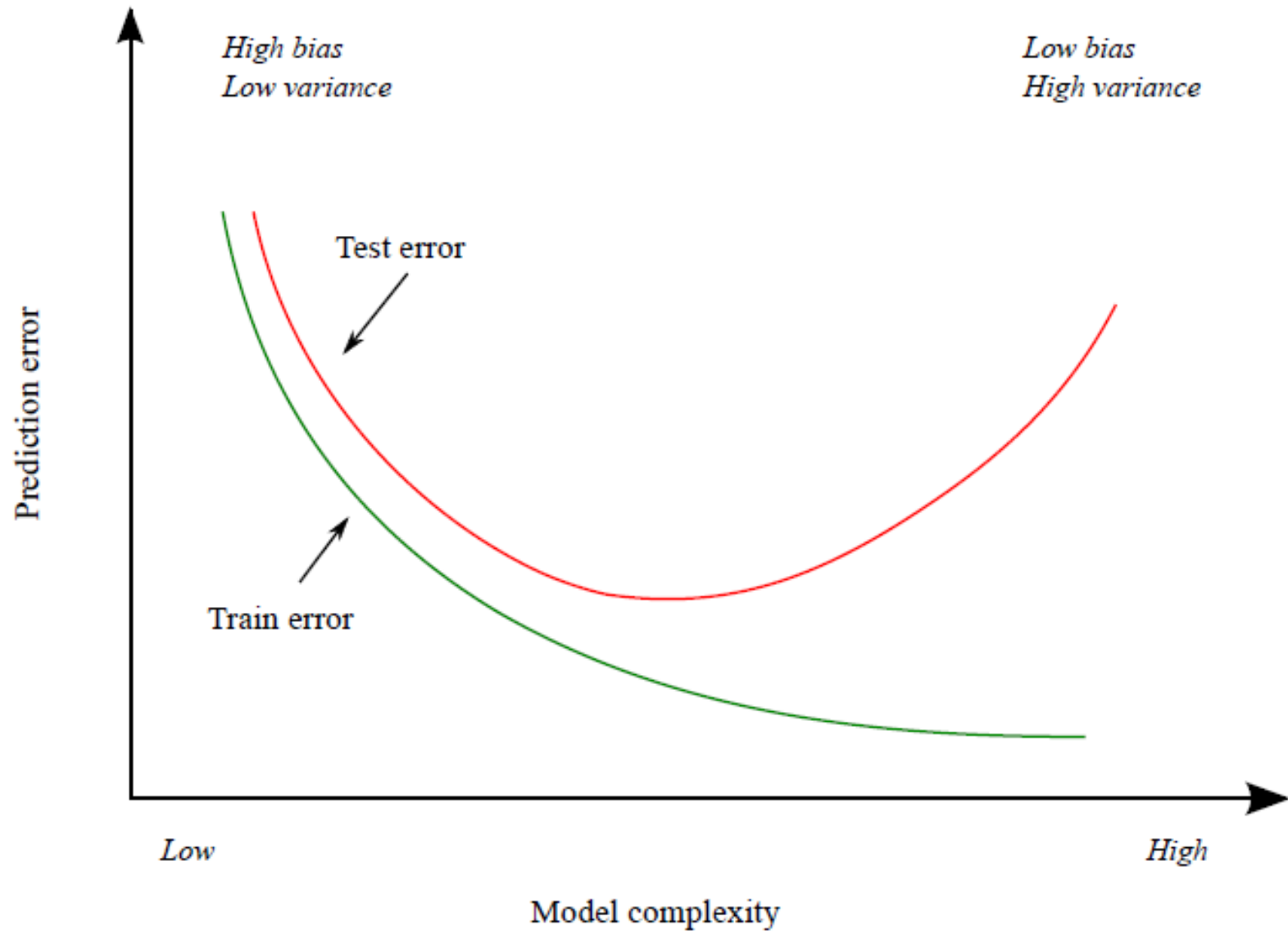Generally, more flexible method: variance increases, bias decreases

$$EPE = E\left(\hat{f}(x) - f(x)\right)^2 \quad \text{BIAS}^2$$
$$+$$
$$E\left(\hat{f}(x) - E(\hat{f}(x))\right)^2 \quad \text{VAR}$$
$$+$$
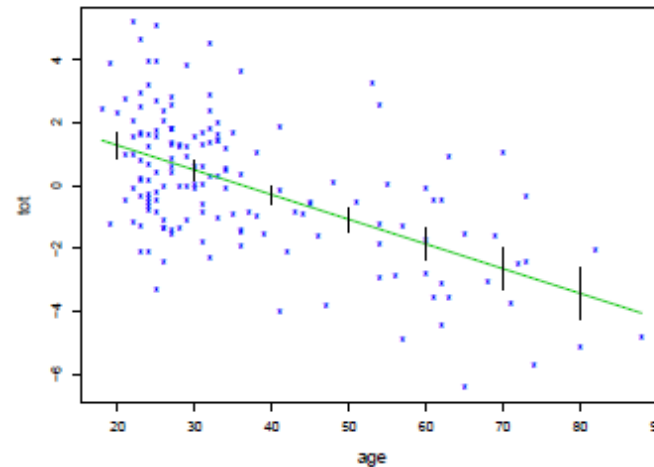$$\sigma^2 \quad \text{NOISE}$$

# Bias-variance tradeoff

# Regularisation

- Aim: reduce variance in exchange of a small bias

- Introduce sparsity

- Limit model complexity by adding a regularisation term

$$\min \sum_{i=1}^{n} (y_i - \beta x_i)^2 + \boxed{\lambda \sum_{j=1}^{l} \beta_j^l}$$

# Linear regression model. A typical example

Consider a study of kidney function. Data represent (x=age of person, y=tot, a composite measure of the overall function). Kidney function declines with age. We need to provide additional information concerning decline rate. This is important in managing kidney transplants.



Check
https://en.wikipedia.org/wiki/Simple_linear_regression

# Linear regression

Data structure. Response
Explanatory variables

Model

Likelihood

Log-likelihood

MLE



$$Y$$
$$(x_1, \ldots, x_n)$$

$$y_i = \beta_0 + \beta_1 x_{1i} + \cdots + \beta_p x_{pi}, \quad i = 1, \ldots, n$$
$$\varepsilon_i \sim N(0, \sigma^2) \text{ IND.}$$

$$\theta = (\beta_0, \ldots, \beta_p, \sigma)$$

$$p(\theta \mid \underline{x}) = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \prod_{i=1}^{n} \exp\left(-\frac{1}{2}\left(\frac{y_i - \beta x_i}{\sigma}\right)^2\right)$$

$$\max \quad -\frac{1}{2}\sum_{i=1}^{n}\left(\frac{y_i - \beta x_i}{\sigma}\right)^2 \quad \cdots$$

$$\hat{\beta} = (X^t X)^{-1} X^t y$$

$$s^2 = \frac{1}{(n-p)}(Y - X\beta)^t(Y - X\beta)$$

# Linear regression

If n or n and p  large



・ COMPUTE $\quad X = QR \qquad Q_{n \cdot p}$ ORTH. COLUMNS, $\quad R_{p \cdot p}$ UPPER TRIANG.

$$\left[ (X^T X)^{-1} = (R^t Q^t Q R)^{-1} = (R^T R)^{-1} = R^{-1} (R^{-1})^T \right]$$

・ COMPUTE $\quad R^{-1}$

・ SOLVE $\quad R \hat{\beta} = Q^T y$

$$\left[ \hat{\beta} = (X^T X)^{-1} X^T y = (R^T Q^T Q R)^{-1} R^T Q^T y \right.$$
$$\left. = (R^T R)^{-1} R^T Q^T y = R^{-1} Q^T y \right]$$

# Linear regression with regulariser

If p large (much larger than n)

$$\min \sum_{i=1}^{n} (y_i - x_i\beta)^2 + \lambda \sum \beta_i^2$$

$$\min \sum_{i=1}^{n} (y_i - x_i\beta)^2 + \lambda \cdot \left| \sum \beta_i \right|$$
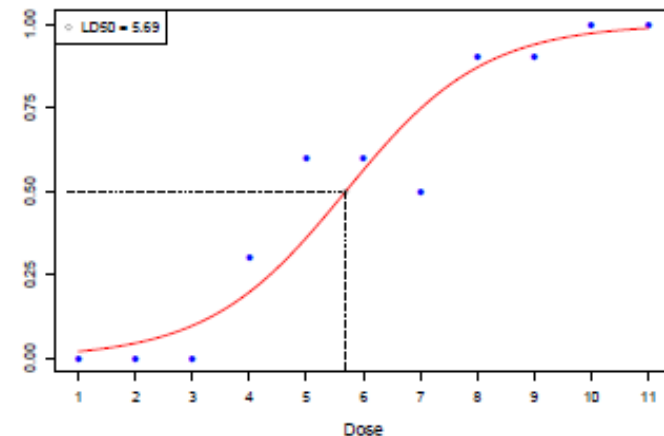
# Bayesian inference with linear regression  model

Model

Standard  noninformative prior

Posterior

$$Y_1 = x_1^T \beta + \varepsilon_1$$
$$\vdots$$
$$Y_n = x_n^T \beta + \varepsilon_n$$
$$\varepsilon_i \sim N(0, \sigma^2)$$

$$y \mid \beta, \sigma, X \sim N(X\beta, \sigma^2 I)$$

$$p(\beta, \sigma^2) = p(\beta \mid \sigma^2) \, p(\sigma^2) \propto \sigma^{-2}$$

$$p(\beta, \sigma^2 \mid y) \propto \underline{p(y \mid \beta, \sigma^2) \, p(\beta, \sigma^2)}$$

$$\beta \mid \sigma, y \sim N(\hat{\beta}, V_\beta \sigma^2)$$

$$V_\beta = (X^T X)^{-1}$$
$$\hat{\beta} = V_\beta X^T y$$

$$p(\sigma^2 \mid y) = \frac{p(\beta, \sigma^2 \mid y)}{p(\beta \mid \sigma^2, y)} \sim \text{Inv-}\chi^2 (n - p, s^2)$$

$$s^2 = \frac{1}{n-p} (y - X\hat{\beta})^T (y - X\hat{\beta})$$

# Logistic regression. A typical example

A new anti-cancer drug is being developed. Before human testing can begin, animal studies are needed to determine safe dosages. A bioassay or dose-response experiment is carried out: 11 groups of 10 mice are treated with an increasing dose of drug and the proportion of deaths are observed.



Check
https://en.wikipedia.org/wiki/Logistic_regression

# Intro to Unsupervised Learning

# Elements of unsupervised learning

Given

- Input space

- Training set

Objective

- Learn model

- Infer some property

- Sample

$$x \in \mathcal{X}$$

$$S = \{x_i\}_{i=1}^{N}$$

$$p(x)$$

# Taxonomy of unsupervised learning algos

- Density estimation
- Manifold learning: PCA, non-linear PCA, …
- Finding modes and groups: cluster analysis, mixture models,…
- Sampling: GANs, Autoencoders, Variational autoencoders,…

# Challenges in unsupervised learning

- High dimension of feature space

- Properties of interest more complex than parameter estimation

- No direct error quantification measure

# Paradigm: Principal component analysis (PCA)

Two views

- Orthogonal projection to lower dimension space to maximize variance
- Linear projection minimizing average projection cost= average quadratic distance between data and projections

Applications

- Dimension reduction
- Compression
- Visualization
- Extraction of predictor. PC Regression
- ….

# PCA: Maximum variance

Given

$$x_i \in \mathbb{R}^D \;,\; i = 1, \ldots, u$$

Find linear projection to space of smaller dimension maximizing variance of projected data

$$\pi : \mathbb{R}^D \rightarrow \mathbb{R}^M \;,\; M < D$$

# PCA: Maximum variance

- 1 dimensional projection
- Projection defined by
- Projection is
- Mean of projected data


- Variance of projected data

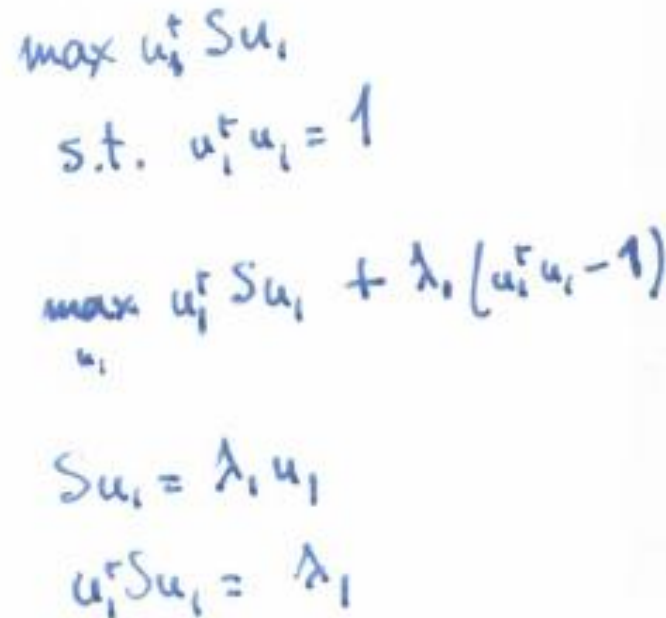$$M = 1, \mathbb{R}$$

$$u_1 \in \mathbb{R}^D \qquad (u_1^t u_1 = 1)$$

$$u_1^t x$$

$$\frac{1}{n} \sum_{i=1}^{n} u_i^t x = u_i^t \bar{x}$$

$$\frac{1}{n} \sum_{i=1}^{n} \left( u_i^t x_i - u_i^t \bar{x} \right)^2 = u_1^t S u_1$$

# PCA: Maximum variance

- Problem to be solved

- Lagrangian formulation

- Solution

$$\max u_i^t S u_i$$
$$\text{s.t.} \quad u_i^t u_i = 1$$

$$\max_{u_i} u_i^t S u_i + \lambda_i (u_i^t u_i - 1)$$

$$S u_i = \lambda_i u_1$$
$$u_i^t S u_i = \lambda_1$$

- Projection is eigenvector associated with first eigenvalue!!!

(and so on)

# Data compression

Projecting each D-dimension point to M

$$\hat{x}_i = \bar{x} + \sum_{j=1}^{M} (x_i^t - \bar{x} u_j) u_j$$

- M = 1

- M = 3

# Data compression



- M = 10



- M = 20

# Data compression

- M = 50

- M = 200

# Implementation challenges

- High dimensionality. What if D>>n
- n points in space of dimension D
- Computational complexity of computing eigenvectors

# Reinforcement learning

# RL: features

- Learning by interaction with environment
  - 'Cause-Effect' relations
  - Consequences of actions
  - What to do to achieve goals
- Goal directed learning: what to do to maximize a reward
  - Discover actions that yield most reward by trying them (trial and error search)
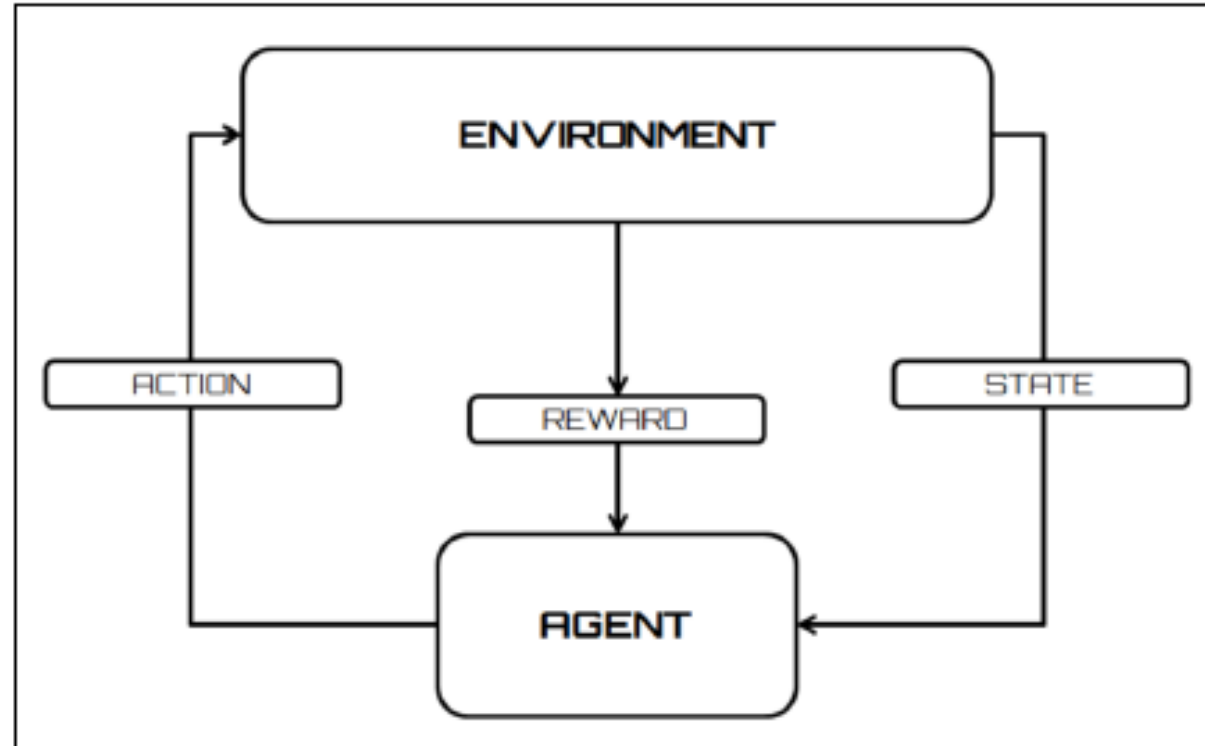  - Actions affect not only immediate reward but also affect environment (delayed reward)

# RL features

- Optimal control of incompletely known Markov decision processes

  - Schemes for sense-act-respond

  - Exploration (collect more info)-exploitation (best action)

  - Uncertainty about evolution of environment and rewards achieved

  - Sequential learning

# RL elements

- Agent
- Environment with states
- Policy
- Reward signal
- Value function
- Model of environment

- Model based methods
- Model free methods

# RL elements

# RL Elements: MDPs

- States
- Actions
- Transition
- Reward
- History
- LT Expected discounted utility
- Policy

$$s \in S$$

$$a \in A$$

$$T: S \times A \longrightarrow \Delta(S)$$

$$R: S \times A \longrightarrow \Delta(R)$$

$$\tau = (s_0, a_0, s_1, a_1, \ldots)$$

$$E_\tau \left( \sum_{t=0}^{\infty} \gamma^t R(a_t, s_t) \right)$$

$$\pi: S \longrightarrow \Delta(A)$$

# RL elements: Q-learning

$$Q(s,a) := (1-\alpha)\, Q(s,a) + \alpha\left(r(s,a) + \gamma \max_{a'} Q(s',a')\right)$$

# Conceptual Recap

# Recap: Classical vs Bayesian

Most approaches in ML (but not all, recall SVMs, RL…)

Once model fixed, we want to learn about it (its parameters)

| Classical | Bayesian |
| --- | --- |
| Parameters fixed | Parameters uncertain, prior |
| Given data, formulate likelihood | Given data, formulate likelihood |
| Maximize likelihood to find MLE (mimimum least squares, cross entropy,…) | Aggregate likelihood and prior to get posterior |
| Plug in MLE to make predictions | Use predictive distribution to make predictions |

Regularisers as bridges

And then used them for decision support !!!

# Inference in ML

Probabilistic model of observed variables x and latent variables z (includes parameters)

$$p(\mathbf{z}, \mathbf{x})$$

ML  e

$$\mathbf{z}^\star = \arg \max_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})$$

MAP e

$$\mathbf{z}^\star = \arg \max_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}) = \arg \max_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$$

Bayes e

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})} = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{\int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}}$$

Incorporates prior info
Estimates full distribution
Denominator (evidence) frequently intractable

# Recap: ML2

Likelihood

$$\ell(\theta|\underline{x}) = \prod_{i=1}^{n} f(x_i|\theta)$$

$$h(\theta) = \log\left(\ell(\theta|\underline{x})\right)$$

MLE

$$\max_{\theta} h(\theta) \longrightarrow \hat{\theta}$$

Predictions

$$f(x|\hat{\theta})$$

# Recap: BML

Prior

Likelihood

Posterior

Predictive

$$f(\theta)$$

$$\ell(\theta \mid x)$$

$$f(\theta \mid x) = \frac{f(x \mid \theta)\, f(\theta)}{f(x)} \propto f(x \mid \theta)\, f(\theta)$$

$$f(y \mid x) = \int f(y \mid \theta)\, f(\theta \mid x)\, d\theta$$

# Recap: RegML

$$\max\ h(\theta) + \lambda g(\theta)$$

$$g(\theta) = \sum \theta_i^2$$

$$g(\theta) = \sum |\theta_i|$$

$$l(\theta) \propto K \rightarrow MAP$$
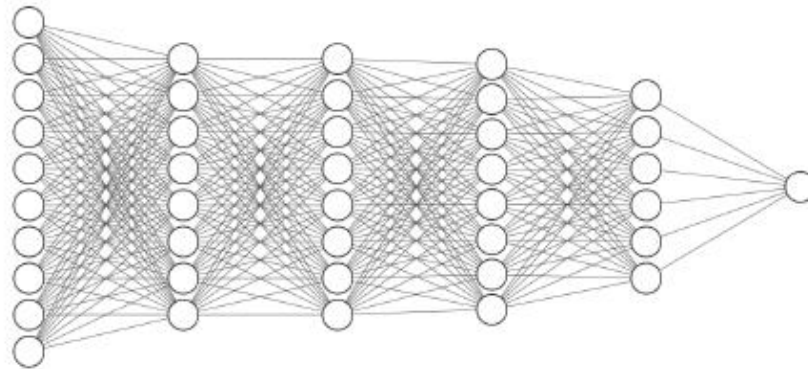
# Optimisation. Recall

Check e.g. Goodfellow et al Ch. 4 (+8)

# Optimization: Fitting neural nets (least squares, maximum likelihood)

$$y_j = \sum_{i=1}^{m} \beta_i \psi(x_k \omega_i) + \varepsilon_j$$

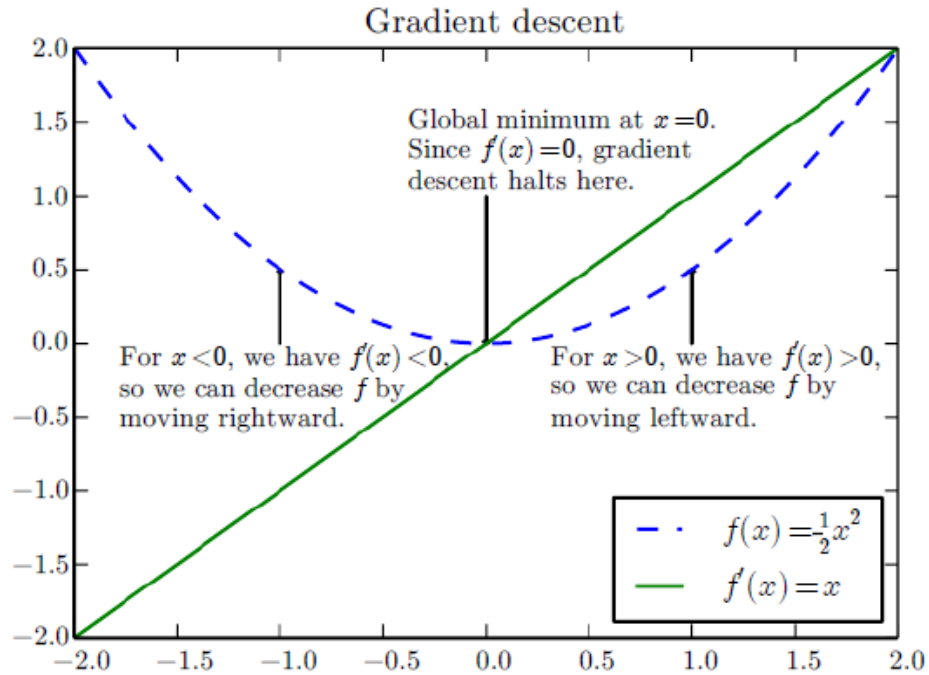$$\min_{\beta, w} \sum_{k=1}^{n} \left( y_k - \sum_{i=1}^{m} \beta_i \psi(x_k \omega_i) \right)^2$$



Input Layer ∈ ℝ¹⁰    Hidden Layer ∈ ℝ⁴  Output Layer ∈ ℝ¹



Input Layer ∈ ℝ¹⁰    Hidden Layer ∈ ℝ⁸   Hidden Layer ∈ ℝ⁸   Hidden Layer ∈ ℝ⁸   Hidden Layer ∈ ℝ⁸  Output Layer ∈ ℝ¹

IAF 2021

# Optimization: Using gradient info



Gradient descent

Global minimum at $x=0$.
Since $f'(x)=0$, gradient
descent halts here.

For $x<0$, we have $f'(x)<0$,
so we can decrease $f$ by
moving rightward.

For $x>0$, we have $f'(x)>0$,
so we can decrease $f$ by
moving leftward.

$- \cdot - \quad f(x)=\frac{1}{2}x^2$

$\underline{\quad\quad} \quad f'(x)=x$

$$f(x+\epsilon) \approx f(x) + \epsilon f'(x)$$

$$f(x - \epsilon\, \mathrm{sign}(f'(x))) \quad < \quad f(x)$$

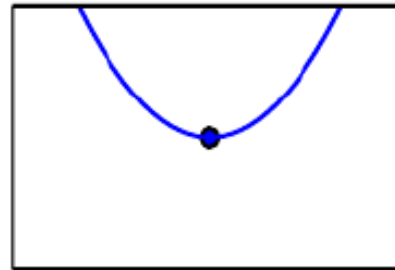$$x' = x - \epsilon \nabla_x f(x)$$

Learning rate

Until stopping condition
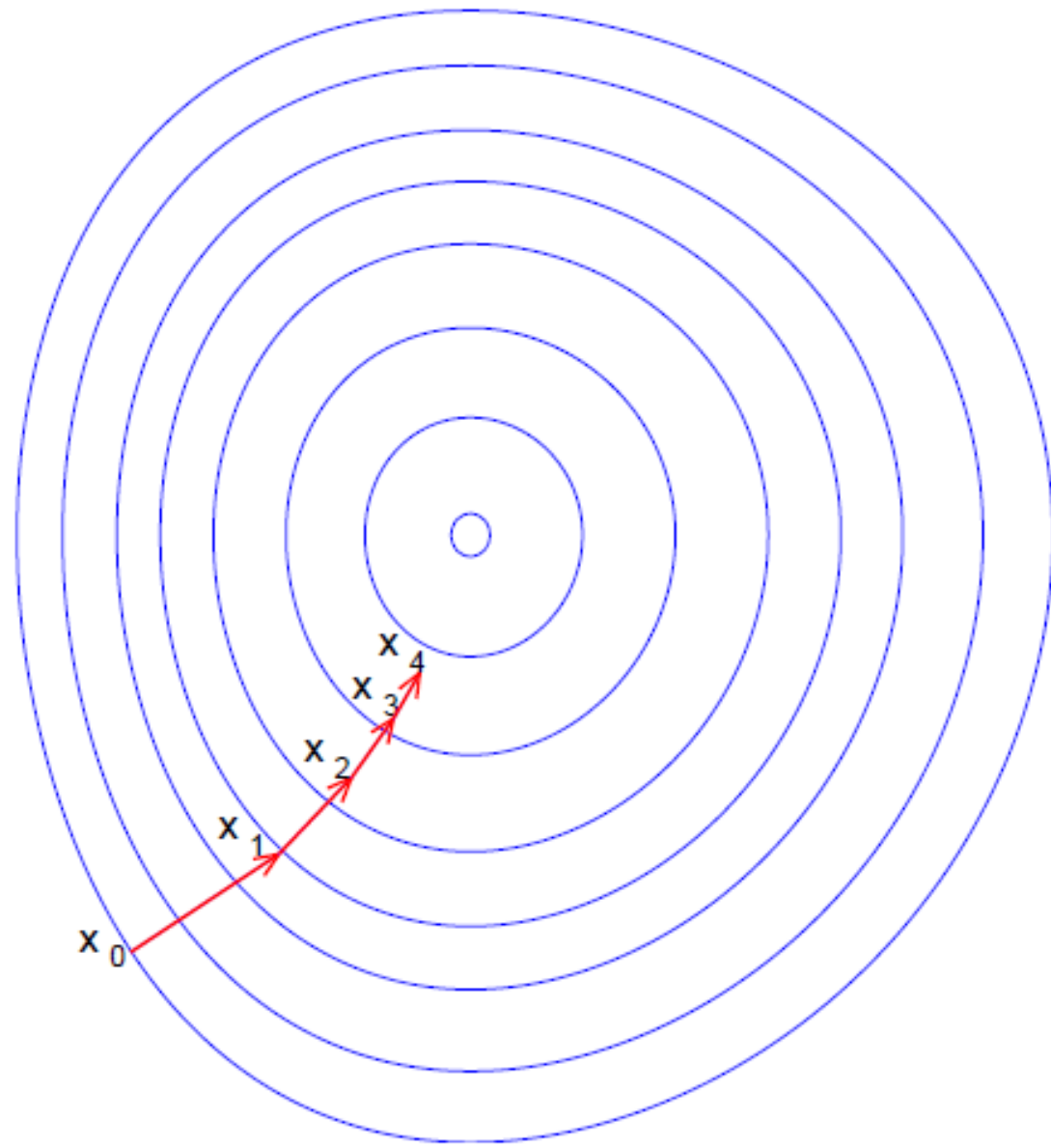Gradient descent
- Fixed and small rate
- Line search

$$f'(x) = 0 \qquad \text{Stationary point}$$

Grad estimation.  Backprop for NNs

# MLE optimization

Problem

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x},\mathbf{y} \sim \hat{p}_{\text{data}}} L(\boldsymbol{x}, y, \boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^{m} L(\boldsymbol{x}^{(i)}, y^{(i)}, \boldsymbol{\theta})$$

$$L(\boldsymbol{x}, y, \boldsymbol{\theta}) = -\log p(y \mid \boldsymbol{x}; \boldsymbol{\theta})$$

Gradient

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^{m} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{x}^{(i)}, y^{(i)}, \boldsymbol{\theta})$$
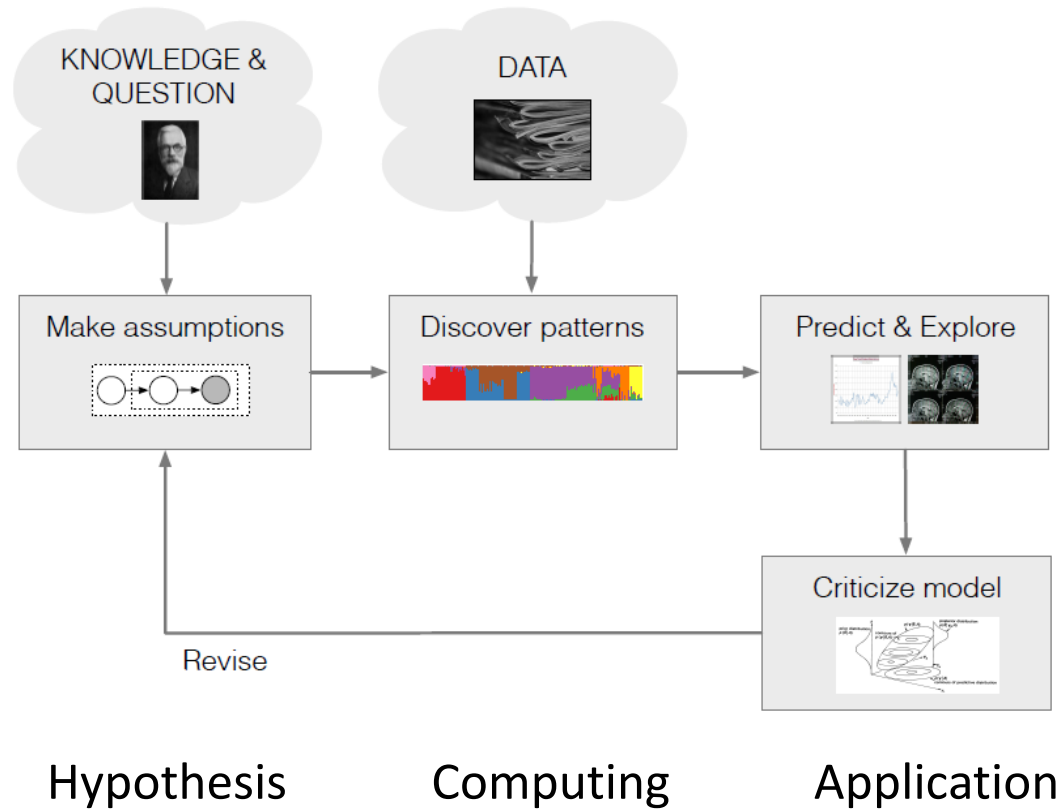
What if also regulariser?

What if  m is large?

Stochastic gradient descent….

# Optimization in ML

- Multimodality

- Large scale

- Gradients expensive

- Hessian superexpensive

- …

# Data flow in machine learning
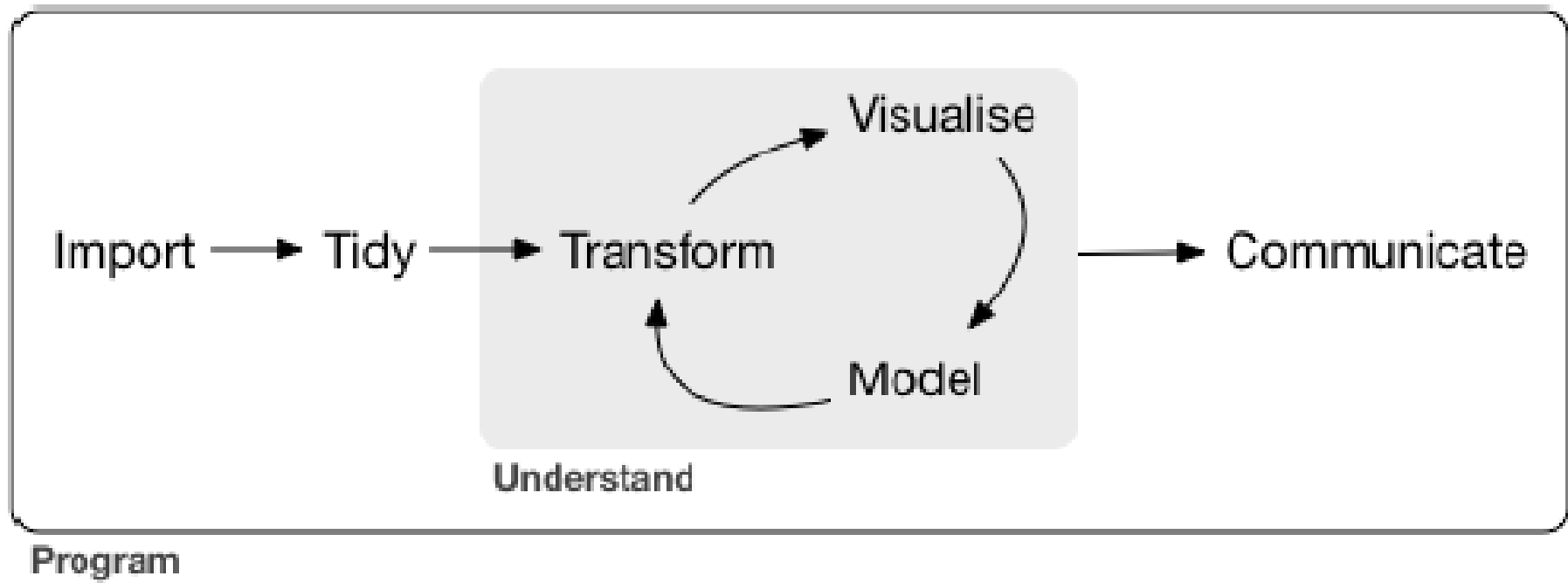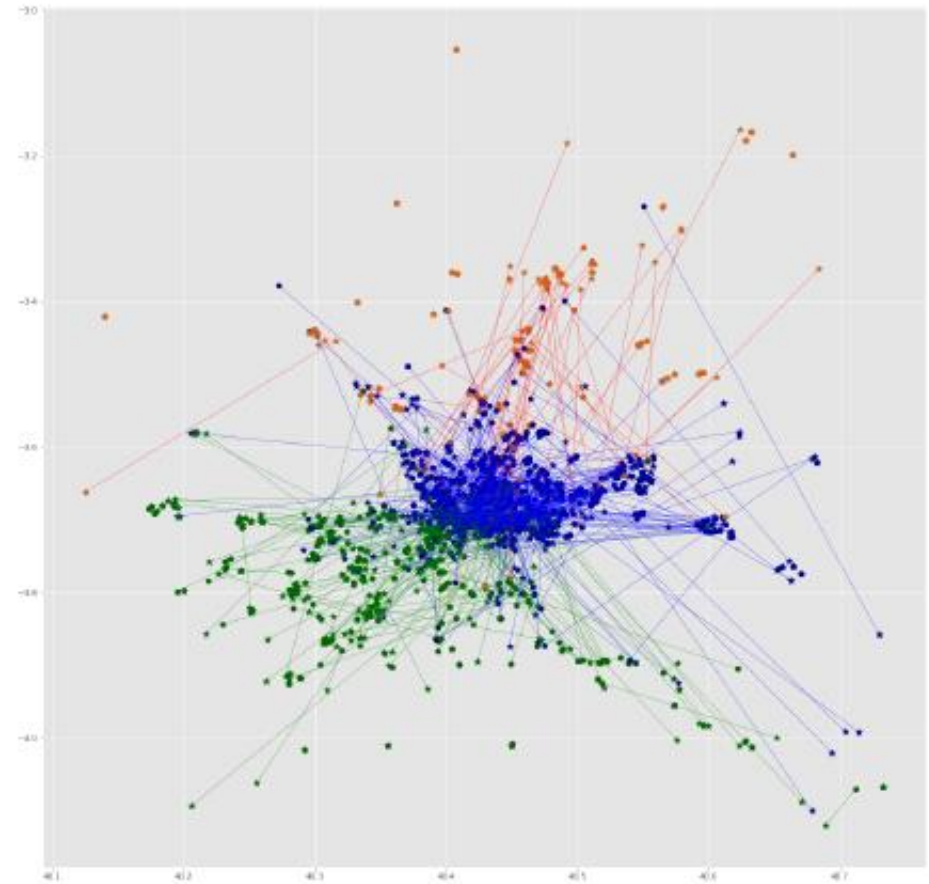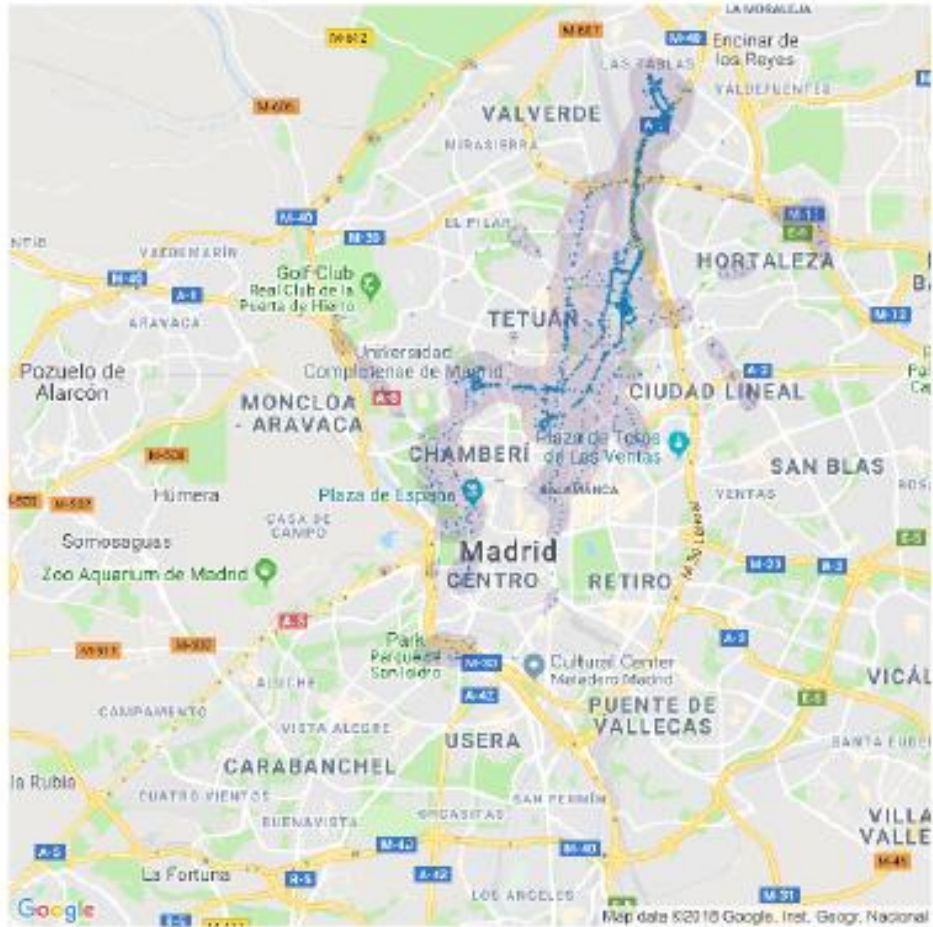
# Broad learning scheme



Inference

What does my
model say
about data?

Hypothesis       Computing       Application

General, Scalable

Just some ideas briefly. Many more during labs and at final part

# ML and BD

# First steps. Preprocessing

- Data from heterogeneous sources (social networks, sensors, samples,….) in different support (text, data bases, streams, images,…)

- First: identify problem to be solved, available variables that may provide information

- Combine available info in a coherent manner

- Final objective of preprocessing: organise data in tensorial/tabular  form

# Different types of data

- Not always trivial to transform data in numerical and/or categorical variables
- Extra preprocessing required
- Examples
  - Text (tweets, web pages,...) word2vec, bag-of-words, n-grams
  - Images: RGB values of pixels, grey intensities
  - Audio: Fourier transform, MFCC  (Mel Frequency Cepstral coeffs)
  - Video: sequences of frames
  - SMILE codes in chemoinformatics
  - Facebook likes

# Summing up

# Recap

- Supervised

- Unsupervised

- Reinforced

- ML vs Bayes
- Challenges due to BD