UNIVERSITY OF CAMBRIDGE

# Advanced Data Science

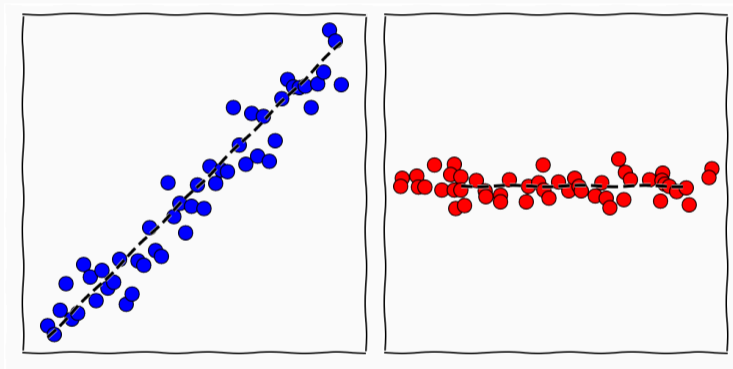Lecture 9 : Statistical Learning Outlook

Carl Henrik Ek - che29@cam.ac.uk

16th of November, 2022

http://carlhenrik.com

# Introduction

**Principal Component Analysis** diagonalises a $D \times D$ matrix

**Principal Component Analysis** diagonalises a $D \times D$ matrix

- finds a geometrical representation where the covariance is diagonal

**Principal Component Analysis** diagonalises a $D \times D$ matrix

- finds a geometrical representation where the covariance is diagonal

**Multi-Dimensional-Scaling** diagonalises a $N \times N$ matrix

**Principal Component Analysis** diagonalises a $D \times D$ matrix

- finds a geometrical representation where the covariance is diagonal

**Multi-Dimensional-Scaling** diagonalises a $N \times N$ matrix

- finds a geometrical representation that "matches" a distance matrix

**Principal Component Analysis** diagonalises a $D \times D$ matrix

- finds a geometrical representation where the covariance is diagonal

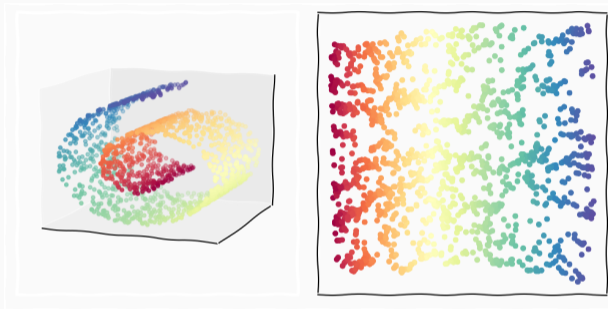**Multi-Dimensional-Scaling** diagonalises a $N \times N$ matrix

- finds a geometrical representation that "matches" a distance matrix
- equivivalent to PCA with euclidian distance

**Principal Component Analysis** diagonalises a $D \times D$ matrix

- finds a geometrical representation where the covariance is diagonal
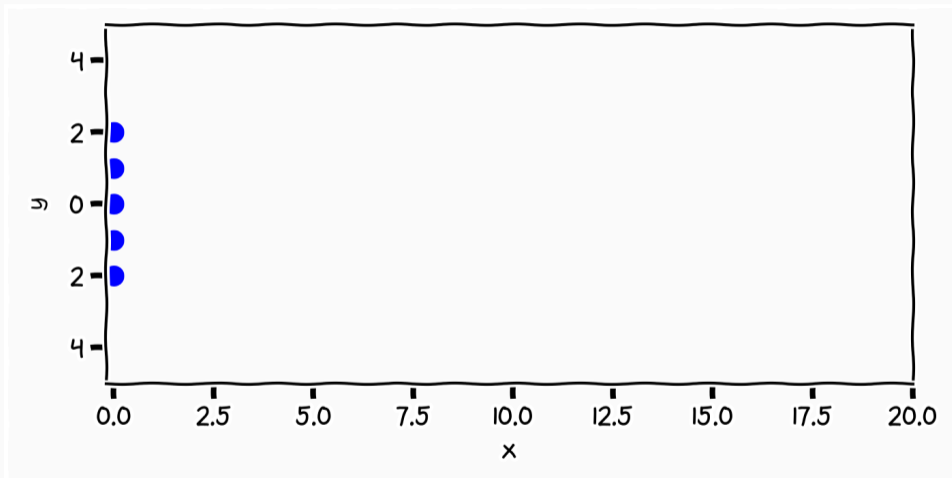
**Multi-Dimensional-Scaling** diagonalises a $N \times N$ matrix

- finds a geometrical representation that "matches" a distance matrix
- equivivalent to PCA with euclidian distance
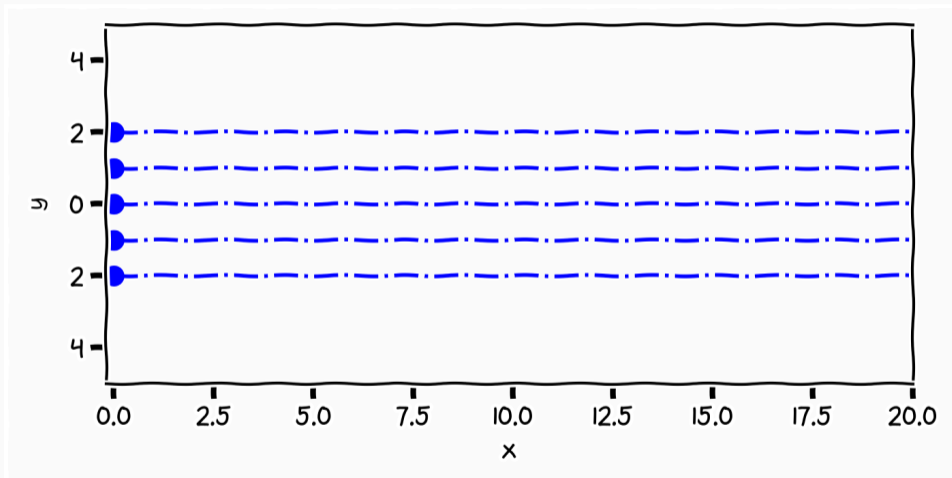- can be non-linearised with a non-linear distance measure

$$\mathbf{y}_i = f(\mathbf{x}_i)$$
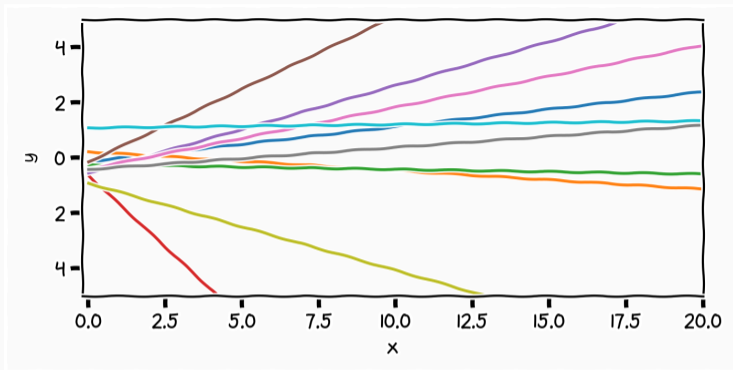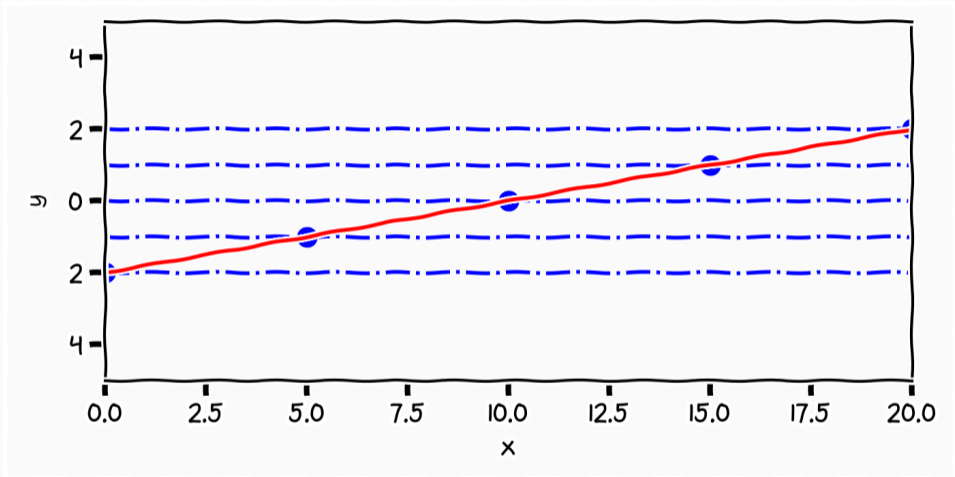
- This problem is very ill-posed
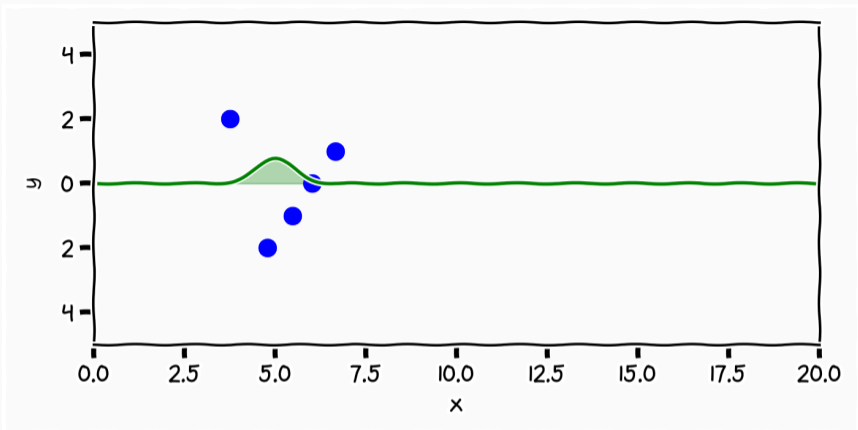- We have to encode a preference towards the solution that we want

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmax}} \prod_{i=1}^{N} p(y_i \mid \boldsymbol{\beta}, \mathbf{x}_i)$$

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmax}} \prod_{i=1}^{N} p(y_i \mid \boldsymbol{\beta}, \mathbf{x}_i) + \lambda \left( \sum_{j=1}^{d} \beta_j^p \right)^{\frac{1}{p}}$$
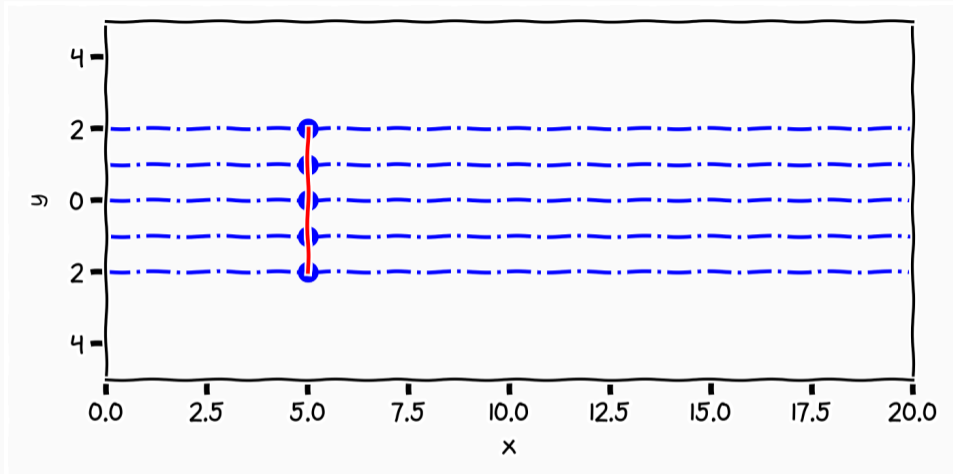
$$y_i = \mathbf{w}^{\mathrm{T}}\mathbf{x}$$

$$p(\mathbf{w}) \sim \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \alpha\mathbf{I})$$
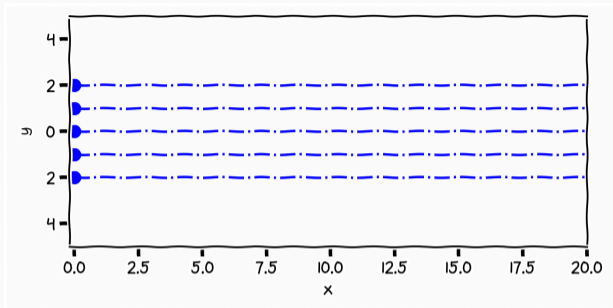
$$p(\mathbf{X}) \sim \mathcal{N}(\mathbf{X} \mid \mathbf{0}, \alpha_2 \mathbf{I})$$

$$p(\mathbf{X}) \sim \mathcal{N}(\mathbf{0}, \alpha_2 \mathbf{I})$$

$$\{\hat{\mathbf{X}}, \hat{\mathbf{w}}\} = \underset{\hat{\mathbf{X}}, \hat{\mathbf{w}}}{\operatorname{argmax}} \left( \underbrace{\mathcal{L}(\mathbf{Y}, \mathbf{X}, \mathbf{w})}_{\log p(\mathbf{Y}|\mathbf{X}, \mathbf{w})} + \gamma_1 \log p(\mathbf{w}) + \gamma_2 \log p(\mathbf{X}) \right)$$

$$p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, \alpha \mathbf{I})$$

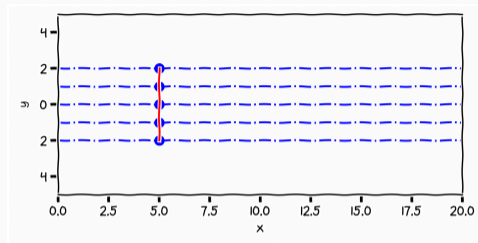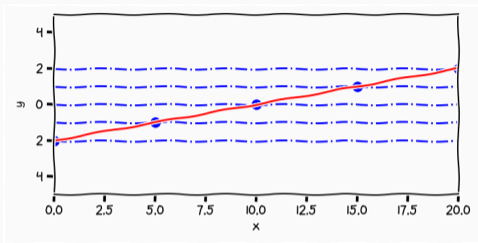$$p(\mathbf{X}) \sim \mathcal{N}(\mathbf{0}, \alpha_2 \mathbf{I})$$



13

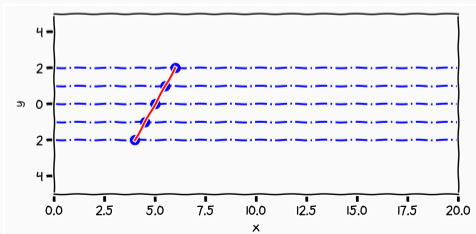**Statistical learning** machine learning is inherently ill-posed

**Statistical learning** machine learning is inherently ill-posed

**Interpretation** our "predictions" can only ever be interpreted in light of the knowledge we put in

**Statistical learning** machine learning is inherently ill-posed

**Interpretation** our "predictions" can only ever be interpreted in light of the knowledge we put in

**Knowledge** how can we incorporate knowledge in a principled manner

**access** what data did I acquire

**assess** how did I prepare/treat the data

**address** which model to choose, how did I set the parameters of the model

# Learning

16

$$p(t) \sim \mathcal{N}(15, 1.5)$$

Deterministic Variable

```
Code
int x = 3;
float y = 3.14;
```

Stochastic Variable

$$x \sim p(x)$$

$$y \sim \mathcal{N}(0, 1)$$

$$\tilde{f}(x) = \int f(x,t)p(t)\mathrm{d}t$$

Sum Rule

$$p(x) = \sum_{\forall y \in \mathcal{Y}} p(x, y)$$

Product Rule

$$p(x, y) = p(x \mid y)p(y)$$

$$p(\mathcal{D}) = \int p(\mathcal{D} \mid \theta) p(\theta) \mathrm{d}\theta$$

$$p(\mathcal{D}) = \int p(\mathcal{D} \mid \theta)p(\theta)\mathrm{d}\theta$$

$$p(\mathcal{D}) = \int p(\mathcal{D} \mid \theta)p(\theta)\mathrm{d}\theta$$

$$p(\mathcal{D}) = \int p(\mathcal{D} \mid \theta) \underbrace{p(\theta)\mathrm{d}\theta}_{\mathrm{d}\hat{\theta}}$$

3?

$$p(x, y) = p(y|x)p(x)$$

$$p(x, y) = p(y|x)p(x)$$
$$p(x, y) = p(x|y)p(y)$$

$$p(x,y) = p(y|x)p(x)$$
$$p(x,y) = p(x|y)p(y)$$
$$p(x|y)p(y) = p(y|x)p(x)$$

$$p(x, y) = p(y|x)p(x)$$
$$p(x, y) = p(x|y)p(y)$$
$$p(x|y)p(y) = p(y|x)p(x)$$
$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

$$p(x, y) = p(y|x)p(x)$$
$$p(x, y) = p(x|y)p(y)$$
$$p(x|y)p(y) = p(y|x)p(x)$$
$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$
$$= \frac{p(y|x)p(x)}{\sum_x p(y|x)p(x)}$$

$$p(\theta \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \theta)p(\theta)}{p(\mathcal{D})}$$

**Likelihood** for a specific parameter setting how does the observation manifest itself

**Prior** what do I believe/know about the parameters

**Evidence** what is the probability of a specific set of data

**Posterior** what is the probability for different parameter settings given a set of data

**Ad-hoc Regularisation** maximum likelihood or regularised error

$$\{\hat{\mathbf{X}}, \hat{\mathbf{w}}\} = \underset{\hat{\mathbf{X}}, \hat{\mathbf{w}}}{\operatorname{argmax}} \left( \underbrace{\mathcal{L}(\mathbf{Y}, \mathbf{X}, \mathbf{w})}_{\log p(\mathbf{Y}|\mathbf{X}, \mathbf{w})} + \gamma_1 \log p(\mathbf{w}) + \gamma_2 \log p(\mathbf{X}) \right)$$

**Principled Regularisation** posterior distribution

$$p(\theta \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \theta)p(\theta)}{\int p(\mathcal{D} \mid \theta)p(\theta)\mathrm{d}\theta}$$

*Integration is a key step in inference, where it is encountered when averaging over the many states of the world consistent with observed data. Indeed, a provocative Bayesian view is that integration is the single challenge separating us from systems that fully automate statistics. More speculatively still, such systems may even exhibit artificial intelligence (ai)*
*– Universal Artificial Intelligence - M. Hutter*

*"Nature laughs at the difficulties of integrations"*
*– Simon Laplace*

DIFFERENTIATION

START

TRY APPLYING
CHAIN POWER
RULE RULE
QUOTIENT PRODUCT
RULE RULE
ETC

DONE? — NO
YES
DONE!

INTEGRATION

START

TRY APPLYING
INTEGRATION SUBSTITUTION
BY PARTS

DONE?

HAHA, NOPE!

CAUCHY'S FORMULA ??? ???!? ???

PARTIAL FRACTIONS

? RIEMANN INTEGRATION

??? INSTALL MATHEMATICA

? ???

STOKES' THEOREM ??? RISCH ALGORITHM

????? ??? WHAT THE HECK IS A BESSEL FUNCTION??

PHONE CALLS TO MATHEMATICIANS OH NO BURN THE EVIDENCE

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) = \sum_{i}^{N} p(\mathbf{y}|\mathbf{x}_i)p(\mathbf{x}_i)$$

- $\mathbf{x}_i$ is a specific binary images

2290593203500326442498254071102 8779924646158308390547680551234
5054431338510774037915738775865 8057318635099533562444284837656
6408900340661545734126916095393 4651531316272895970961099648619
5486636741656944283948869330648 4701733713508133208092688099524
0707971539803921050200955733579 4366205566676730638553849508752
9677470990968153918788613785751 3890052212385415364000233552517
9230941551480812783648467474496 1578781252261713953420063416790
7552057630497077601674681891226 1453204962575441115371836944715
6895505073882545721273943517481 6507334054019330445298798029650
8746618030728963410359112463410 9184832439049686890853942279882

9655406361370980789697504759416 7461331023628146001054998291892
8850448033966038407878196527044 7157474368533868315778800203562
1474121034155871572968019805251 8982409725023084881200238736500
2027283572275248844963488736471 3943526031912848227248826190464
8476965948928382396693052519124 1687725175533908692952453783598
2837023543516588536916371046489 4220310701508827933380526429979
2599815801920922903898158871712 8926097153382729134531621865313
9786085815417055159827515344471 3326325034781836776513703100360
9793889758575377908303501066776 6548311999605347475370343426743
8253400053810997864187276609708 2093090380663944422789696913654
8900202322285082544979530967870 6304437009833849217731493021674

## Number of terms iii

```
25506248717508338594766791895095   5680602732346712939153259990811
48939130328420650376019730541966   1524092173016464047938013691439
66718432036059811187775136277555   7250792266837423597968228683403
40891384751547673727271229322226   8878852083218796660305975797728
87782987686468159942599573254086   8749600987758158350339985951647
51217086975807460294738428018336   8592485796034133919973077413533
68694919563685166113776742372086   1780419191068702807890339161440
99126661387307752660057804524226   5302437317858452782485229505751
37610939444647228055539117717166   4315059230286413698788578331540
17822394957907816501100598872746   5959467831004471989549305375741
90738099064718222518825147478496   0657161167548497523333968812279
```

43

49114751199656354594624473392897828672753085721621023943443062
01449072780844668538929442057198697060107876495003418069047901
81420256733072612769503473201816461274039931292984401423199725
43409301707634660377253374196629143599599348813527131013125346
35085302320378163021153281388668643014293963947674718567131663
50435955804654725436951706056632361702749907044372801683830358
69913652994643262056428393431504053504888101754720253838078891
92539392721103826349328251385543816977282386956487514065578882
34747518138465426828255208381310069117625217360239526199430454
34643503384285930316545135079767510717638042435127189839307791
20937657434512013867455548820224148073627378623609980111113076

0640189547044207203761774747082 0243516866198003957569584101060
8046613562965001201466456771415 5778664863093617634553900426210
9110167208910075825348801584001 7224071067971558665492397885347
6607256313817084019127947685341 8537351879721277733449450507730
3189505040470344922506903873556 9656865708529073446623478695245
6543122517479114466613670208736 0842313671545657762822696089905
6802168279902278674508669673834 7816102210900054189076993778672
7705964820658607375143364171301 1744511704016132334906338900377
1777472580944833242545989973822 5646744609738390155521757096422
2619375692340966923479020630115 9076383049447801135255878205328
2752643299087648267991015324907 4963538068771014944040060242262

3804497742682401904233153226013 9373317250133351983527123955504
2292211010517136771541981666250 0131430427440349387764312765762
4870317305687566284108475166000 1324414350620739304183073837766
8972502903711649967733818943578 9237255328232566165426546313829
11359993958629376

- Possible black and white 3 Megapixel images

$$2^{3145728}$$

- Possible black and white 3 Megapixel images

$$2^{3145728}$$

- Number of atoms in the universe

$$10^{80} \approx (2^{\frac{10}{3}})^{80} \approx 2^{267}$$

## Numbers

- Possible black and white 3 Megapixel images

$$2^{3145728}$$

- Number of atoms in the universe

$$10^{80} \approx (2^{\frac{10}{3}})^{80} \approx 2^{267}$$

- Age of the universe in seconds

$$4.35 \cdot 10^{17} \approx 2^{59}$$

- Computational intractability: there are too many states to sum over (image segmentation)
- Analytic: no closed form exists for the distribution (unsupervised learning)

- Computational intractability: there are too many states to sum over (image segmentation)
- Analytic: no closed form exists for the distribution (unsupervised learning)
- The double annoyance: machine learning is not just ill-posed, the computations needed for making it well posed is intractable

- There exists no universal learner

- There exists no universal learner
- For every learner there exist a task on which it fails

- There exists no universal learner
- For every learner there exist a task on which it fails
- Every algorithm that learns something useful does so by assumptions

- There exists no universal learner
- For every learner there exist a task on which it fails
- Every algorithm that learns something useful does so by assumptions
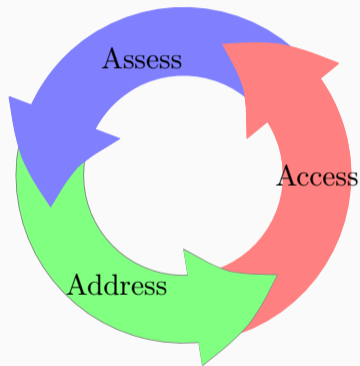- *There is no free lunch algorithm*

Imagen

HOW TO MAKE AN **OMELETTE**

3 EGGS WHISK HEAT OLIVE OIL

POUR EGGS INTO PAN

USE A SPATULA TO PULL SOLIDIFIED EDGES TOWARDS THE MIDDLE

KEEP PULLING UNTIL THERE'S JUST A BIT OF FLUID LEFT

FOLD IN THREE AND SERVE
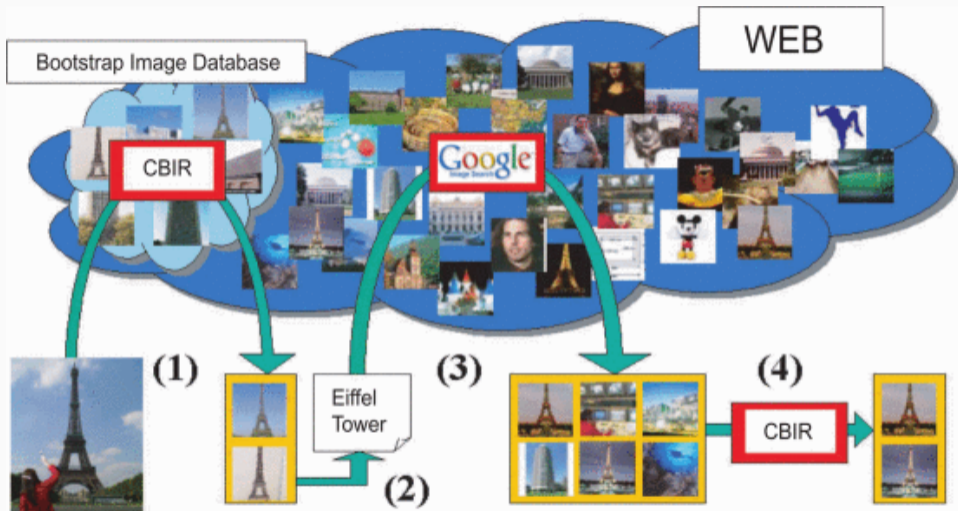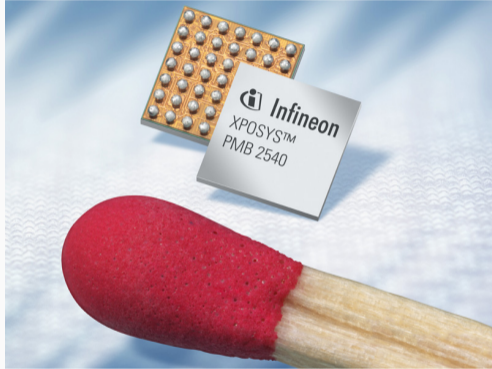
51

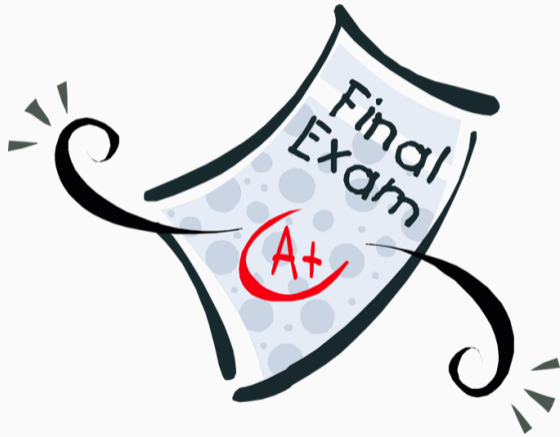*"You need to put Machine Learning in the context of data (and humans)"*

2010

55

# Summary

- Machine learning problems are inherently ill-posed

- Machine learning problems are inherently ill-posed
- We need to introduce knowledge/assumptions

- Machine learning problems are inherently ill-posed
- We need to introduce knowledge/assumptions
- The results can only be interpreted in light of the knowledge/assumptions

- Machine learning problems are inherently ill-posed
- We need to introduce knowledge/assumptions
- The results can only be interpreted in light of the knowledge/assumptions
- Use methods that you can explain as you need to communicate with domain experts

*Alongside your implementation you will provide a* short repository overview *describing how you have implemented the different parts of the project and where you have placed those parts in your code repository. You will submit* your code *alongside a version of* this notebook *that will allow your examiner to understand and reconstruct the thinking behind your analysis.*

*Remember the notebook you create should tell a story, any code that is not critical to that story can safely be placed into the associated analysis library and imported for use (structured as given in the Fynesse template)*

**Lack of narrative** why are you doing what you are doing?

**Lack of narrative** why are you doing what you are doing?

**Spaghetti code** encapsulate code, clean up code

**Lack of narrative** why are you doing what you are doing?

**Spaghetti code** encapsulate code, clean up code

**The perfect prediction** what does this even mean?

**Lack of narrative** why are you doing what you are doing?

**Spaghetti code** encapsulate code, clean up code

**The perfect prediction** what does this even mean?

**ML Ninjas** we will not give additional marks for "advanced methods"

eof