



The
University
Of
Sheffield.

Non-linear Modelling by Adaptive Pre-processing

Rob Harrison
Automatic Control & Systems Engineering

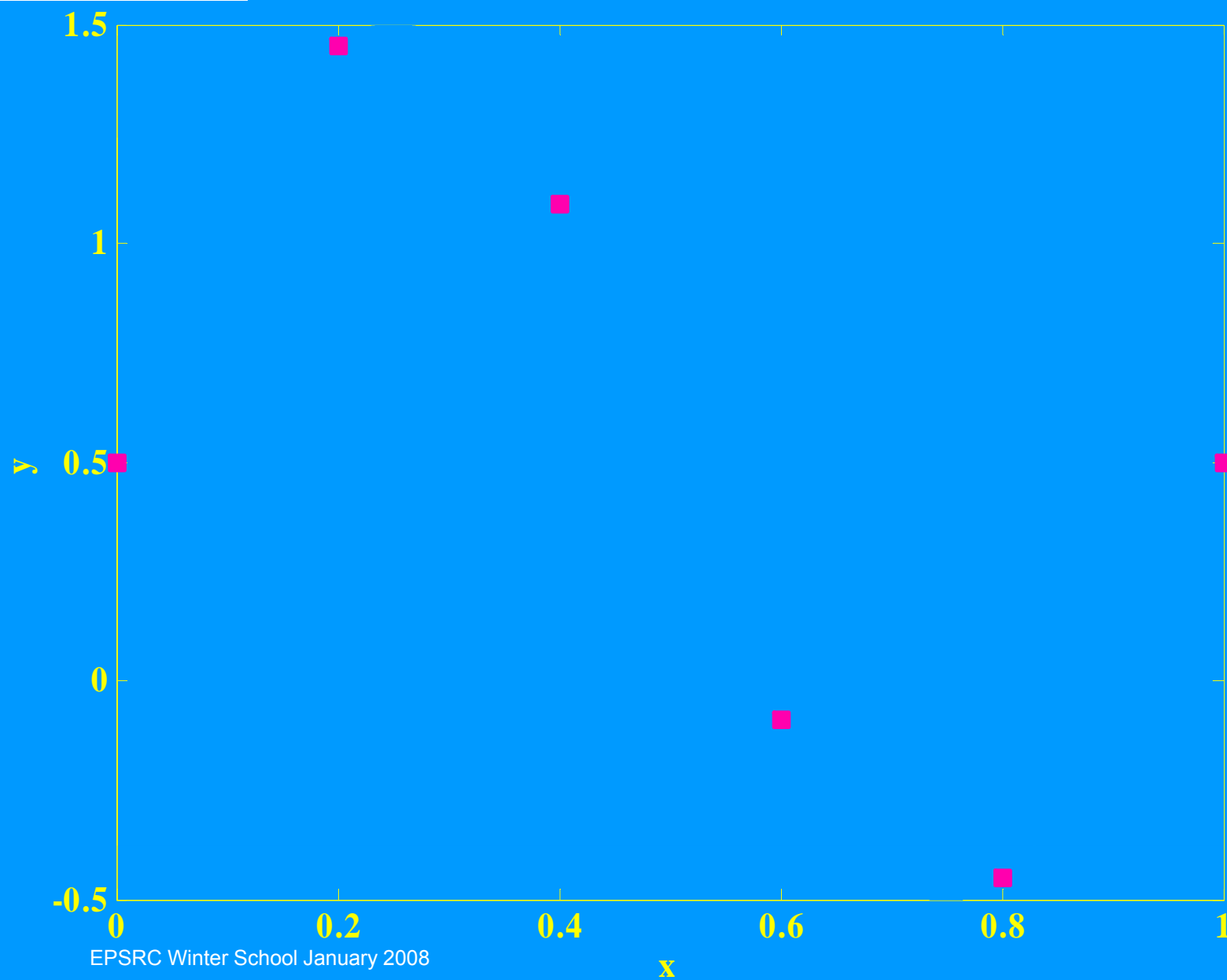


The Data Modelling Problem

- $y = f(x)$ $z = y + e$
- multivariate & non-linear
 - measurement errors
- $\{x_i, z_i\} \ i = 1:N$ $z_i = f(x_i) + e_i$
- infer behaviour *everywhere* from a few examples
 - little or no prior information on $f(x)$
- \hat{y} etc. indicates *estimate*

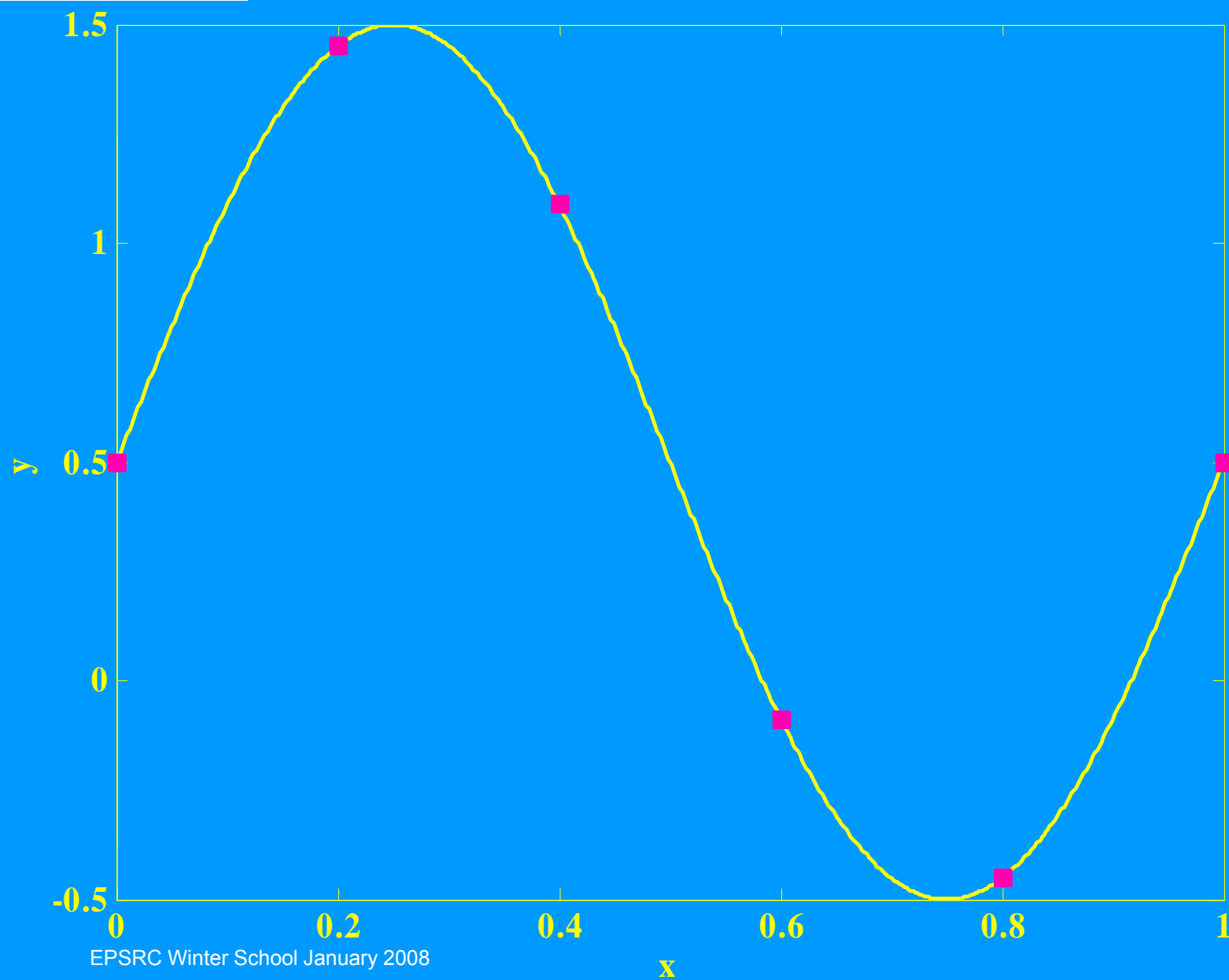


The
University
Of
Sheffield.



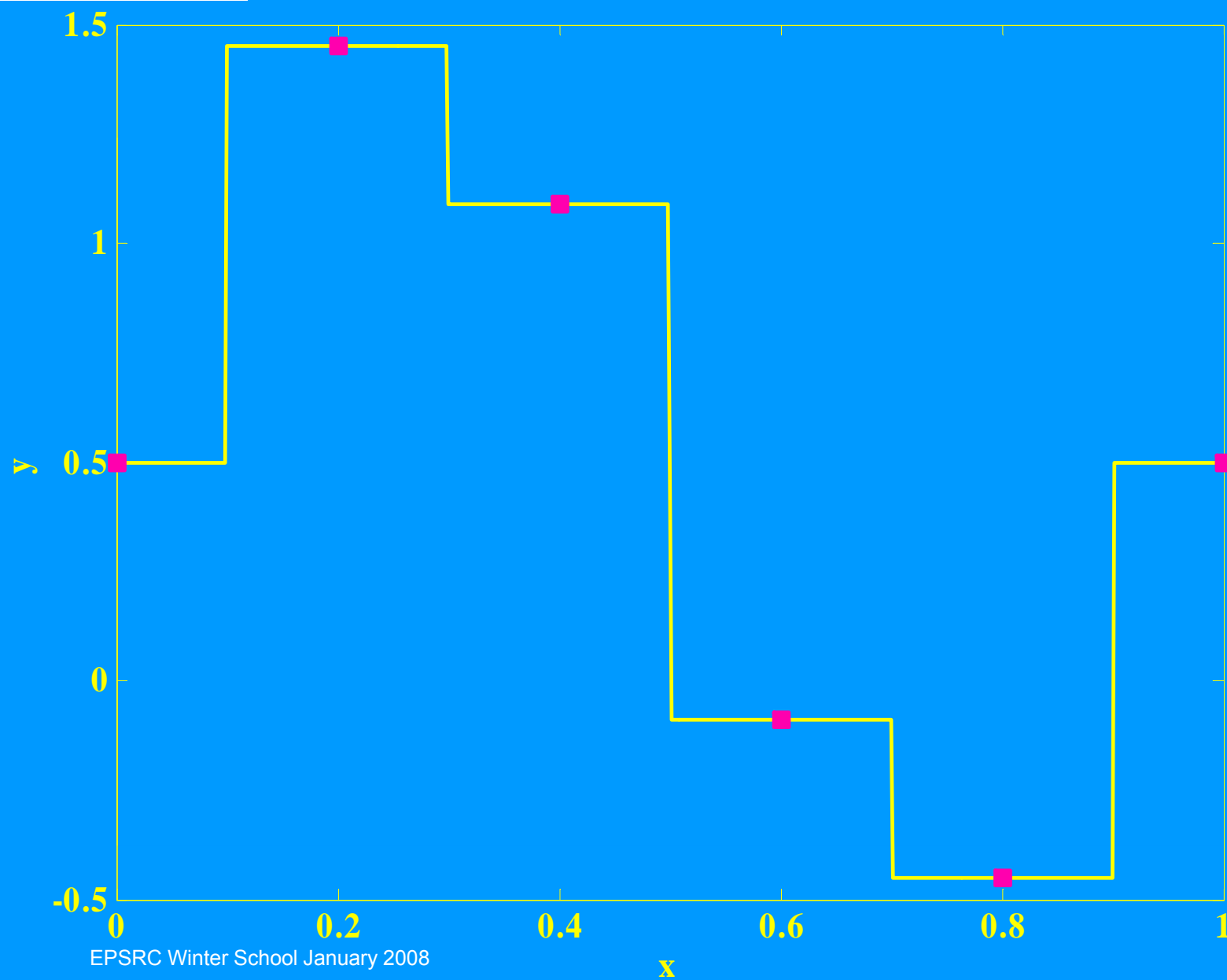


The
University
Of
Sheffield.



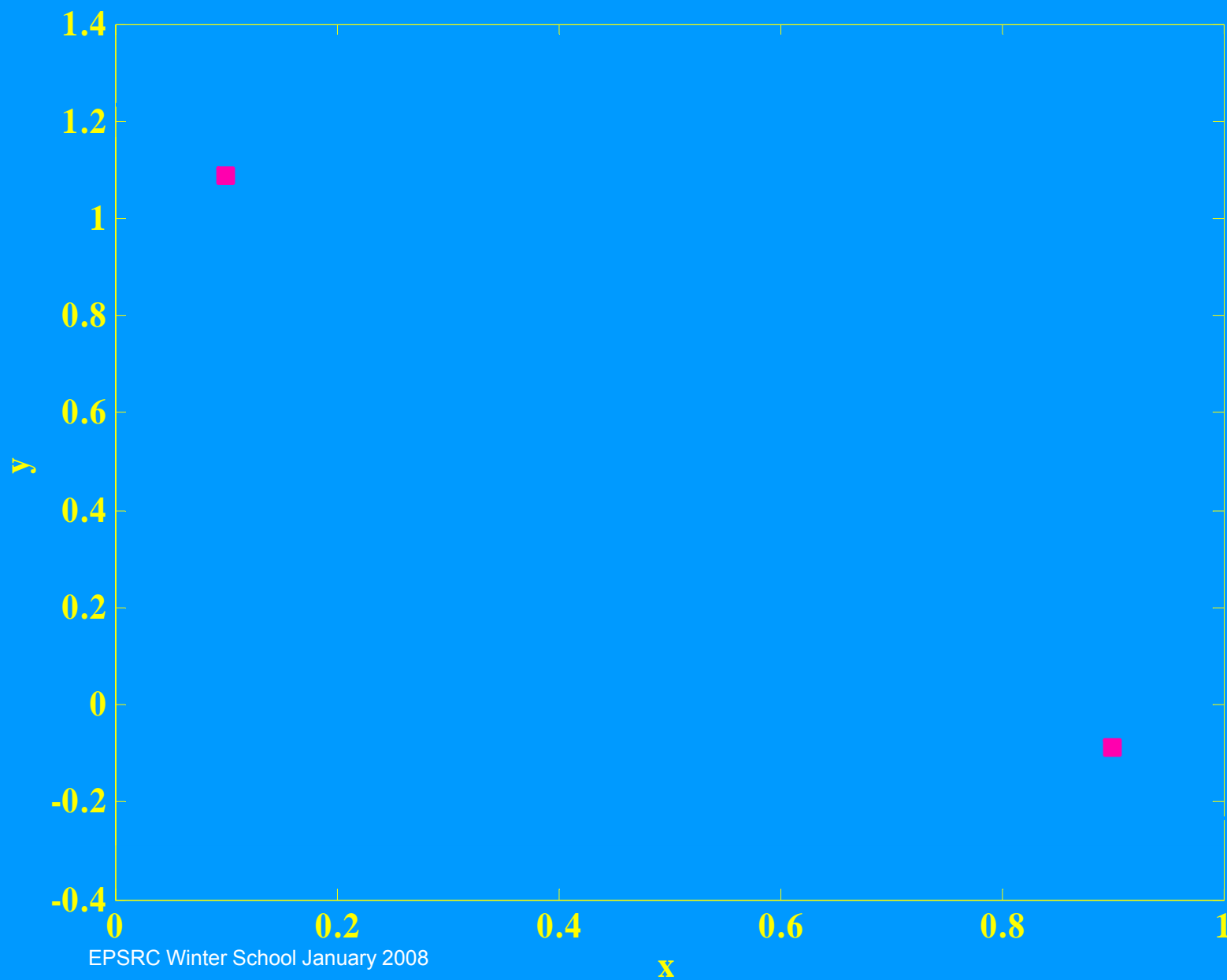


The
University
Of
Sheffield.



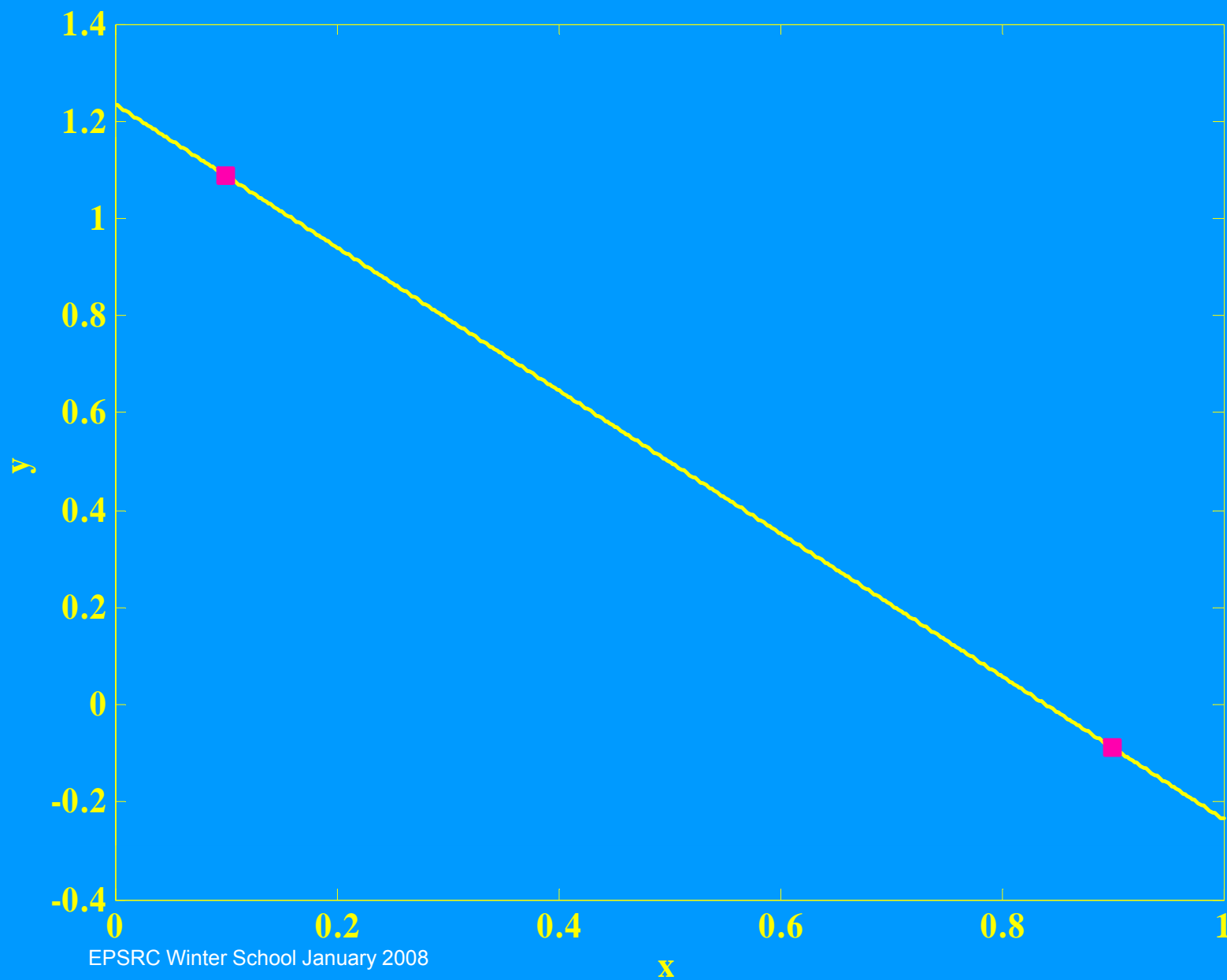


The
University
Of
Sheffield.



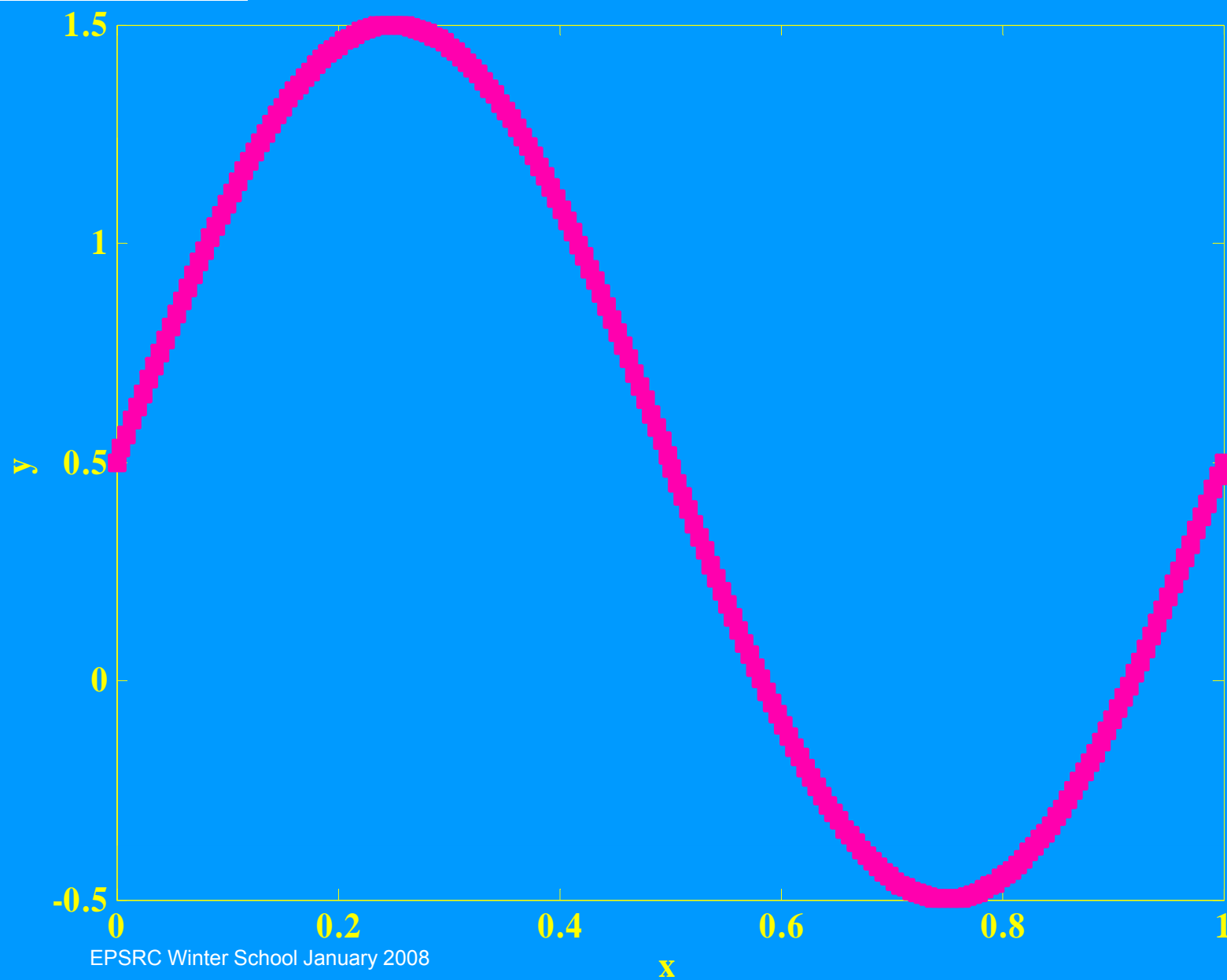


The
University
Of
Sheffield.





The
University
Of
Sheffield.





Dimensionality

- lose ability to see the shape of $f(x)$
 - try it in 13-D
- number of samples exponential in d
 - if N OK in 1-D, N^d needed in d -D
- how do we know if “well-spaced”?
 - how can we sample where the action is?
 - observational vs experimental data!
- ALWAYS undersampled!

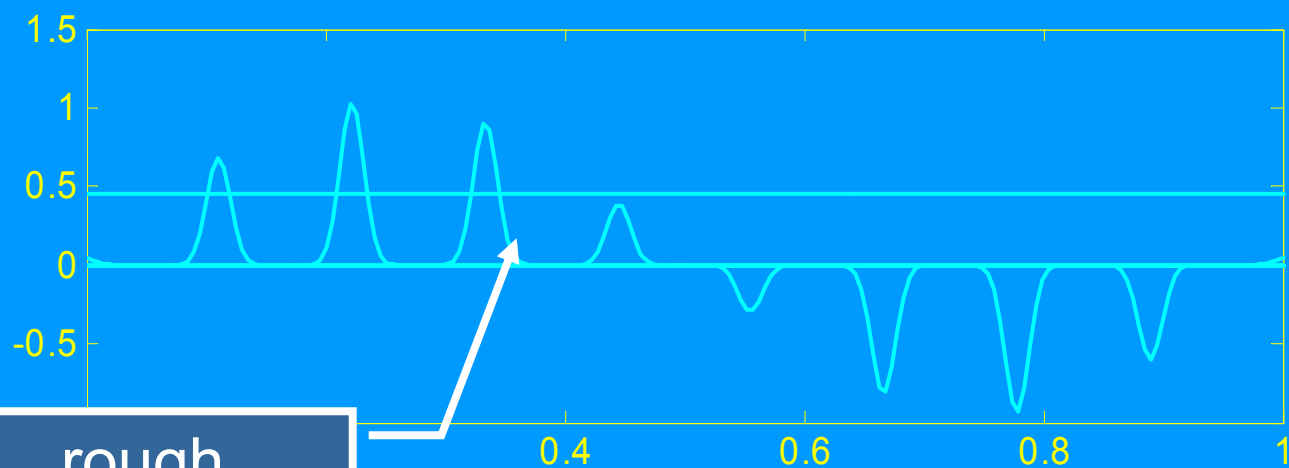
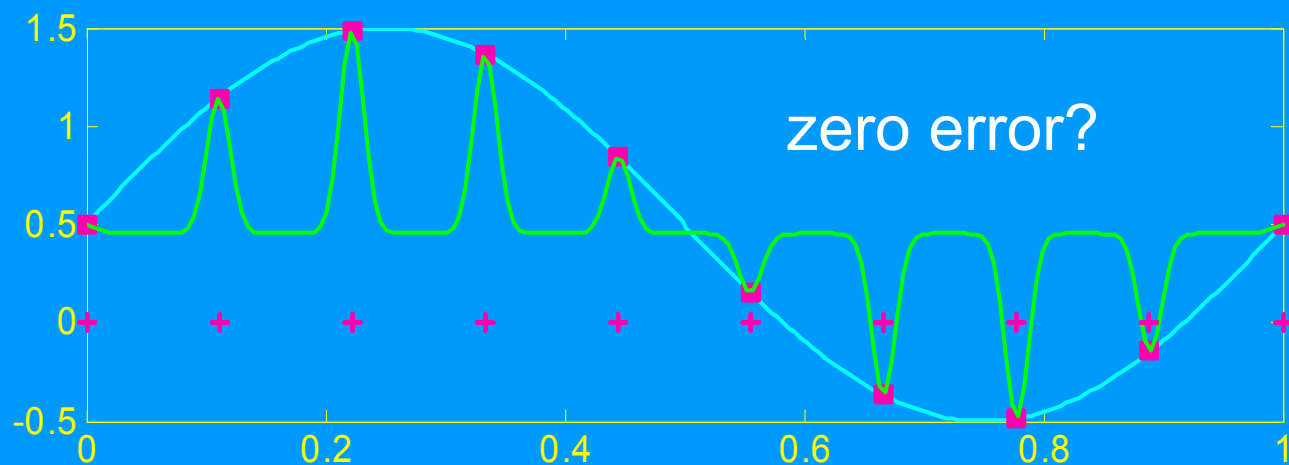


What goes on in the Gaps?

- Universal Approximation
- Advantage
 - can bend to (almost) any shape
- Disadvantage
 - can bend to (almost) any shape
- Training data is all we have to go on



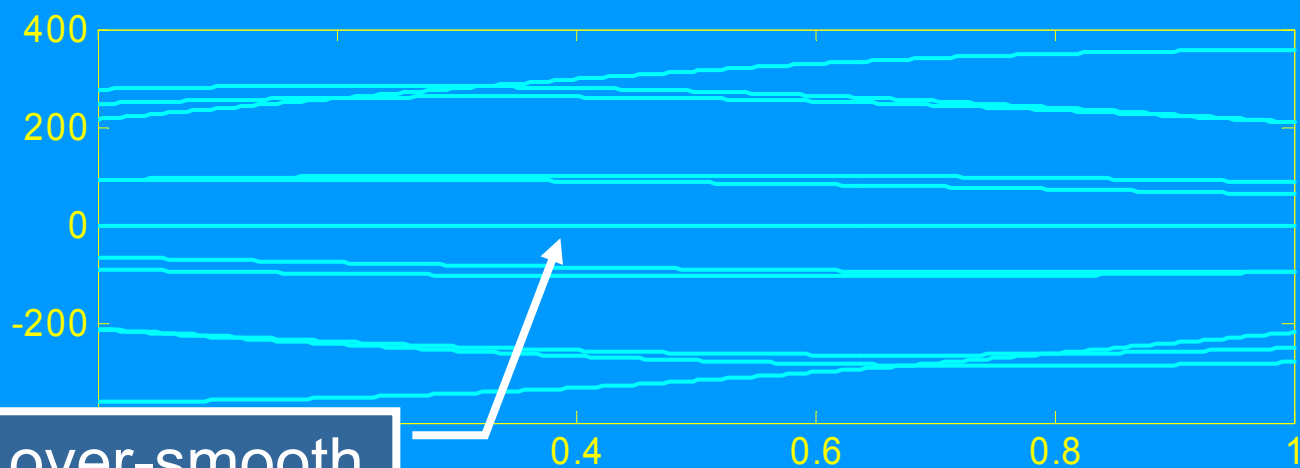
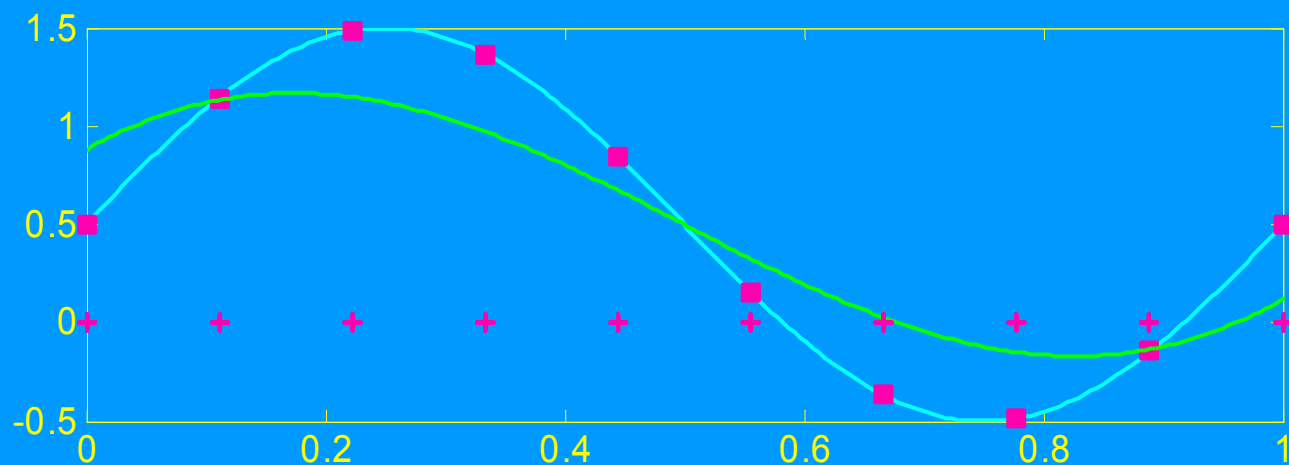
Overfitting (sample data)



rough
components



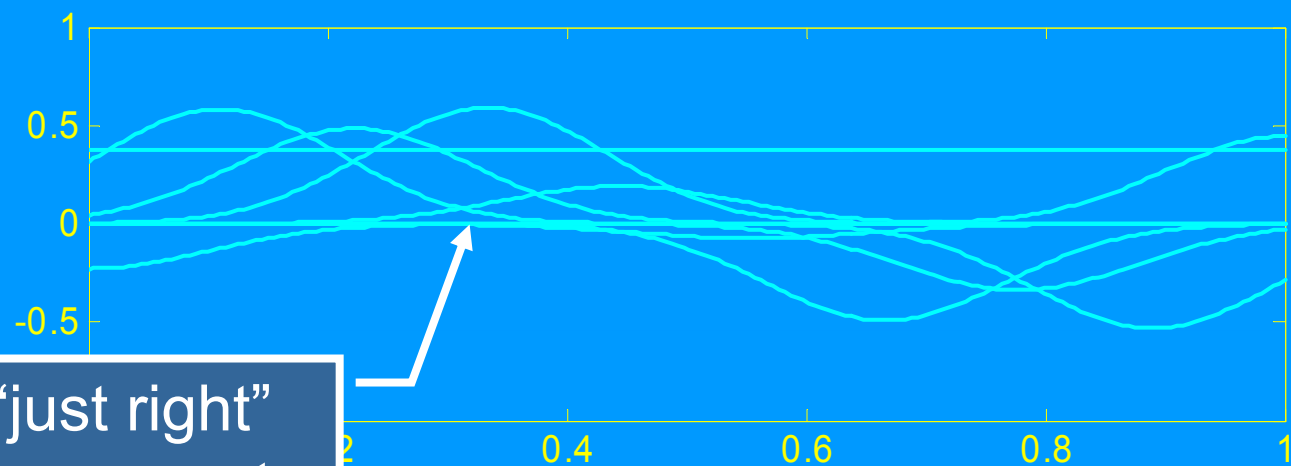
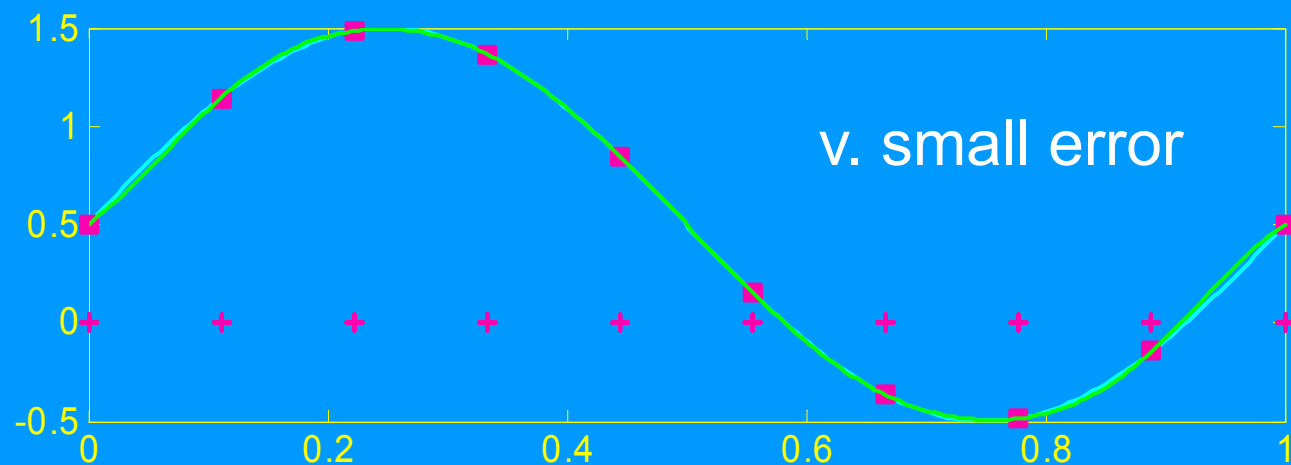
Underfitting (sample data)



over-smooth
components



Goldilocks



"just right"
components

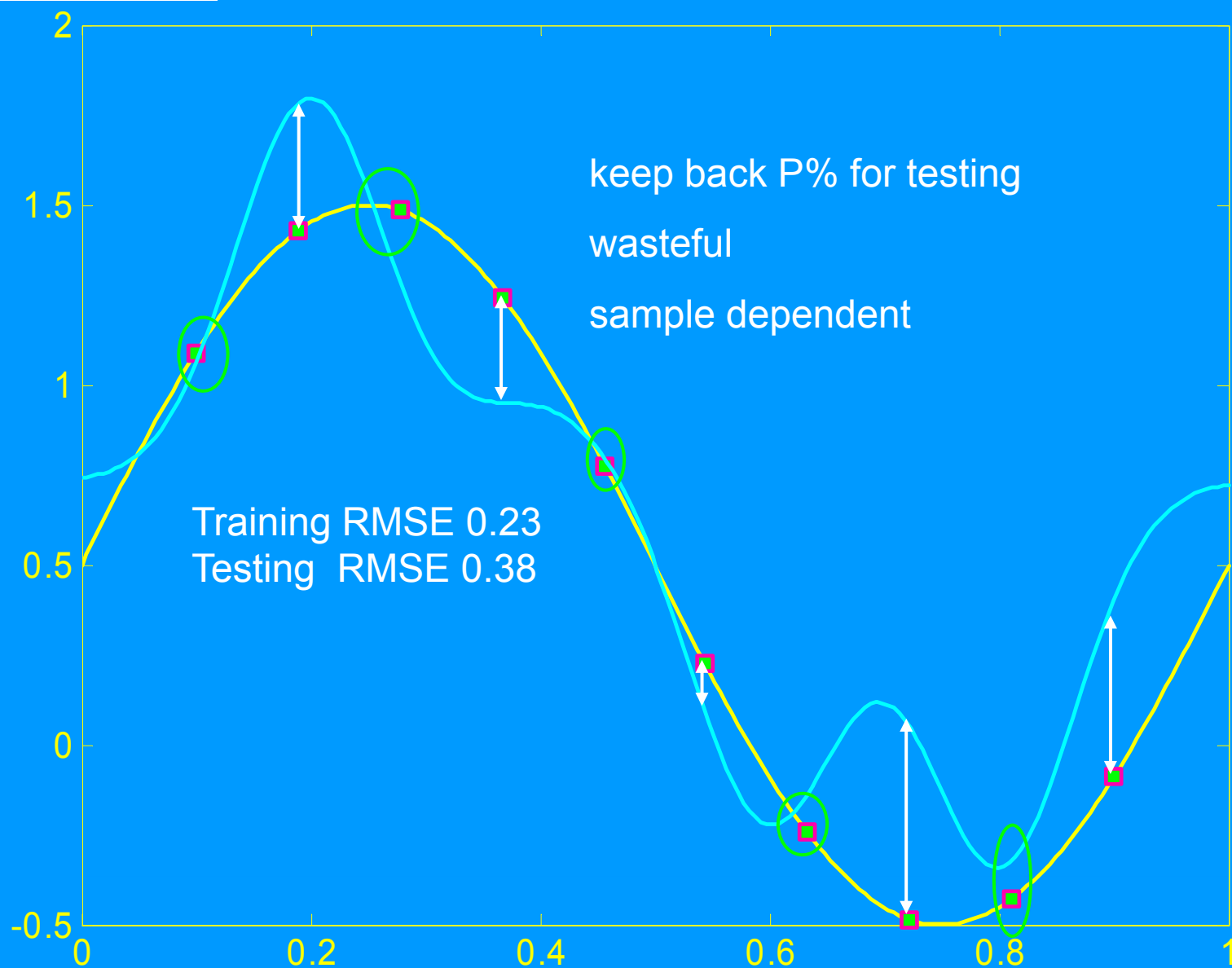


Restricting “Flexibility”

- use data to tell the estimator how to behave
- regularization/penalization
 - penalize “roughness”
 - e.g. $SSE + \rho Q$
 - $Q = \sum w_{ij}^2 \rightarrow \hat{w} = (\Phi^T \Phi + \rho I)^{-1} \Phi^T z$
- use potentially complex structure
 - data constrains where it can
 - Q constrains elsewhere



Hold-out Method





Cross Validation

- leave-one-out CV
 - train on all but one
 - test that one
 - repeat N times
 - compute performance
- m-fold CV
 - divide sample into m non-overlapping sets
 - proceed as above
- all data used for training and testing
 - more work but realistic performance estimates
- used to choose “hyper-parameters”
 - e.g. ρ , number, width ...



Training Data

X_1	Z_1
X_2	Z_2
X_3	Z_3
X_4	Z_4
X_5	Z_5

train

test

Y_4

Y_1

Y_2

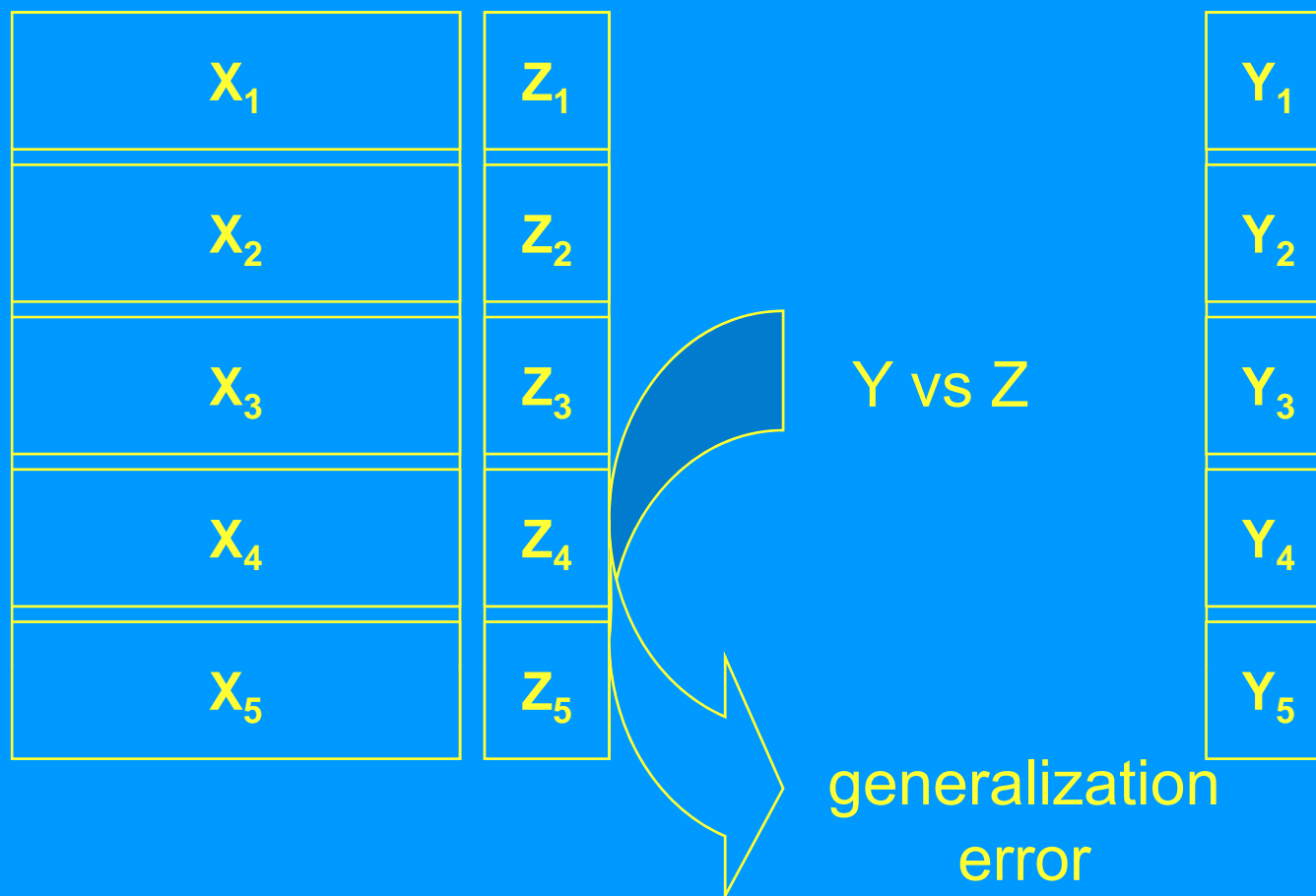
Y_3



The
University
Of
Sheffield.

X_1	Z_1
X_2	Z_2
X_3	Z_3
X_4	Z_4
X_5	Z_5

Y_1
Y_2
Y_3
Y_4





Adaptive Basis Functions

- “linear” models
 - fixed pre-processing
 - parameters \rightarrow cost “benign”
 - easy to optimize

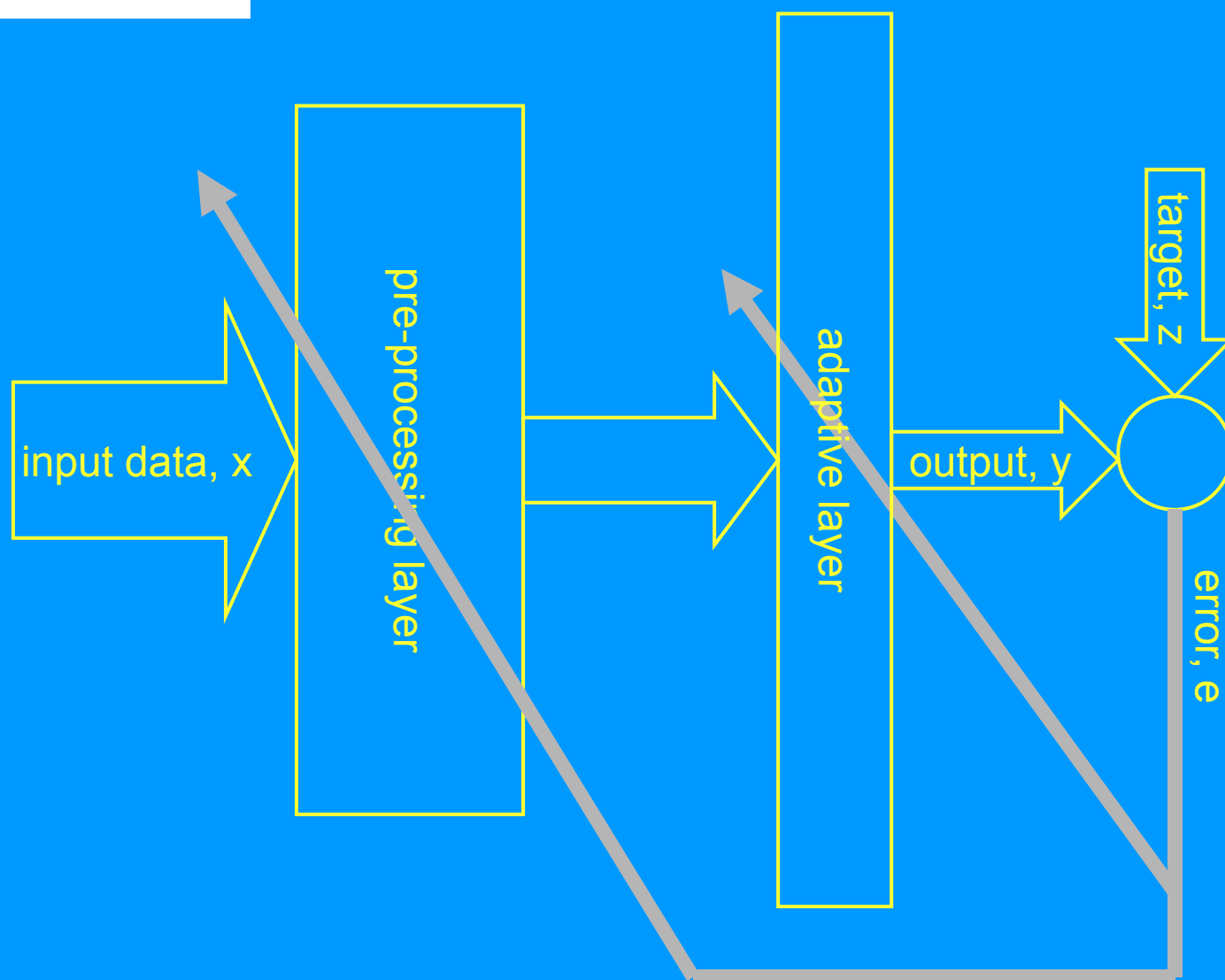
but

- combinatorial
- arbitrary choices

what is best pre-processor to choose?



The
University
Of
Sheffield.





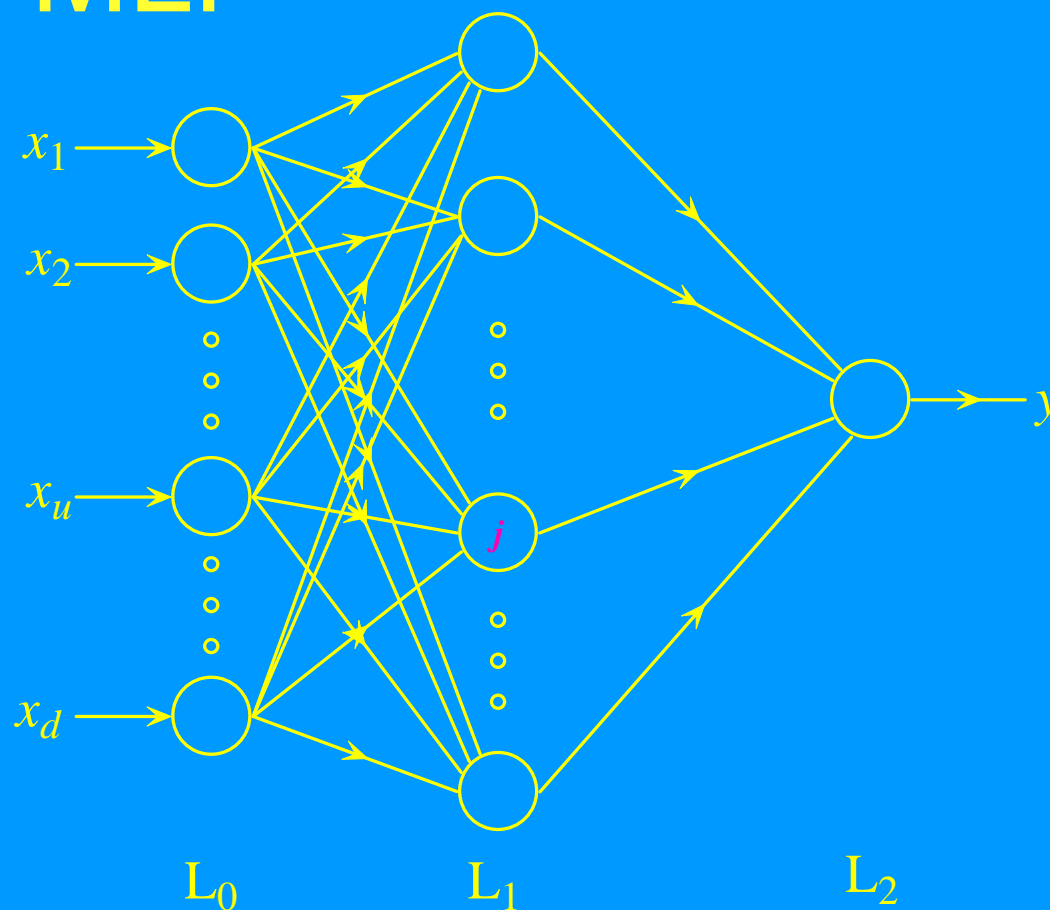
The Multi-Layer Perceptron

- formulated from loose biological principles
- popularized mid 1980s
 - Rumelhart, Hinton & Williams 1986
 - Werbos 1974, Ho 1964
- “learn” pre-processing stage from data
- layered, feed-forward structure
 - sigmoidal pre-processing
 - task-specific output

non-linear model



Two-Layer MLP



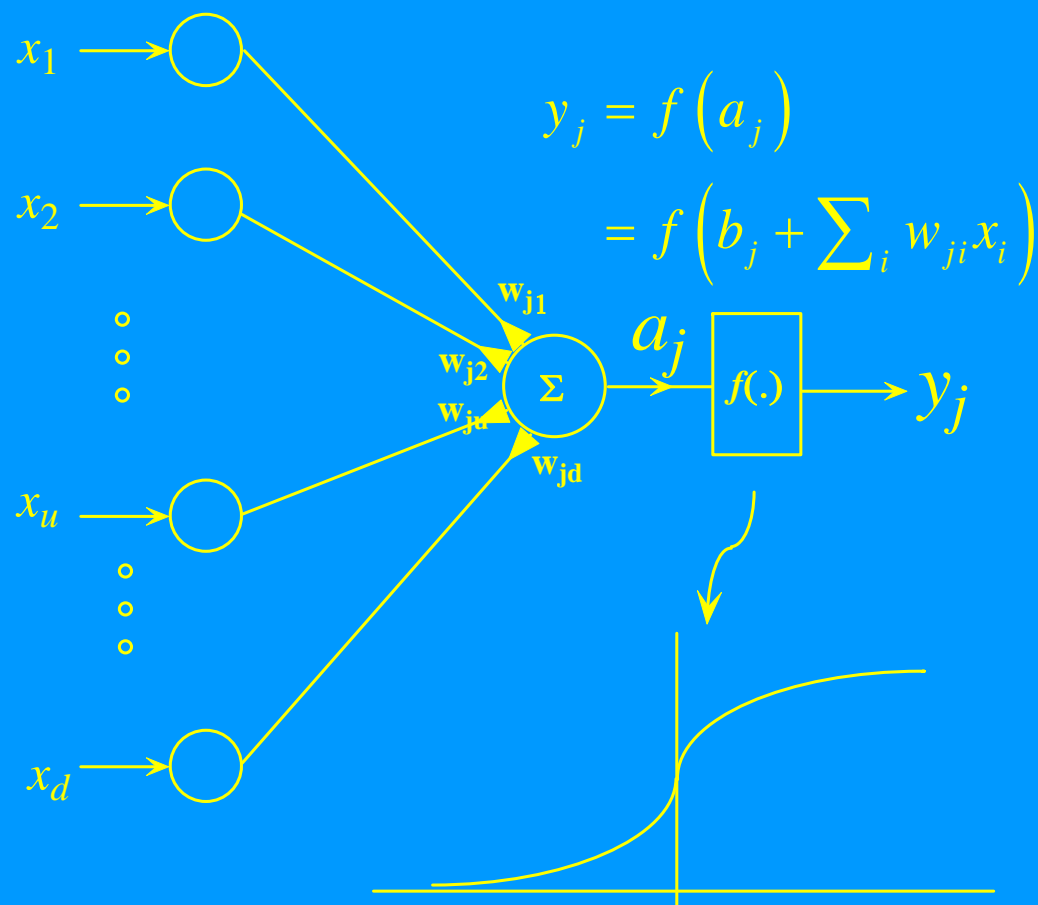
$$y = \theta(\underline{w}^T \underline{v} + b)$$

$$v_j = \sigma(\underline{w}_j^T \underline{x} + b_j)$$

$$y = \theta\left(\sum_i w_i \sigma\left(\sum_j w_j x_j + b_j\right) + b\right)$$



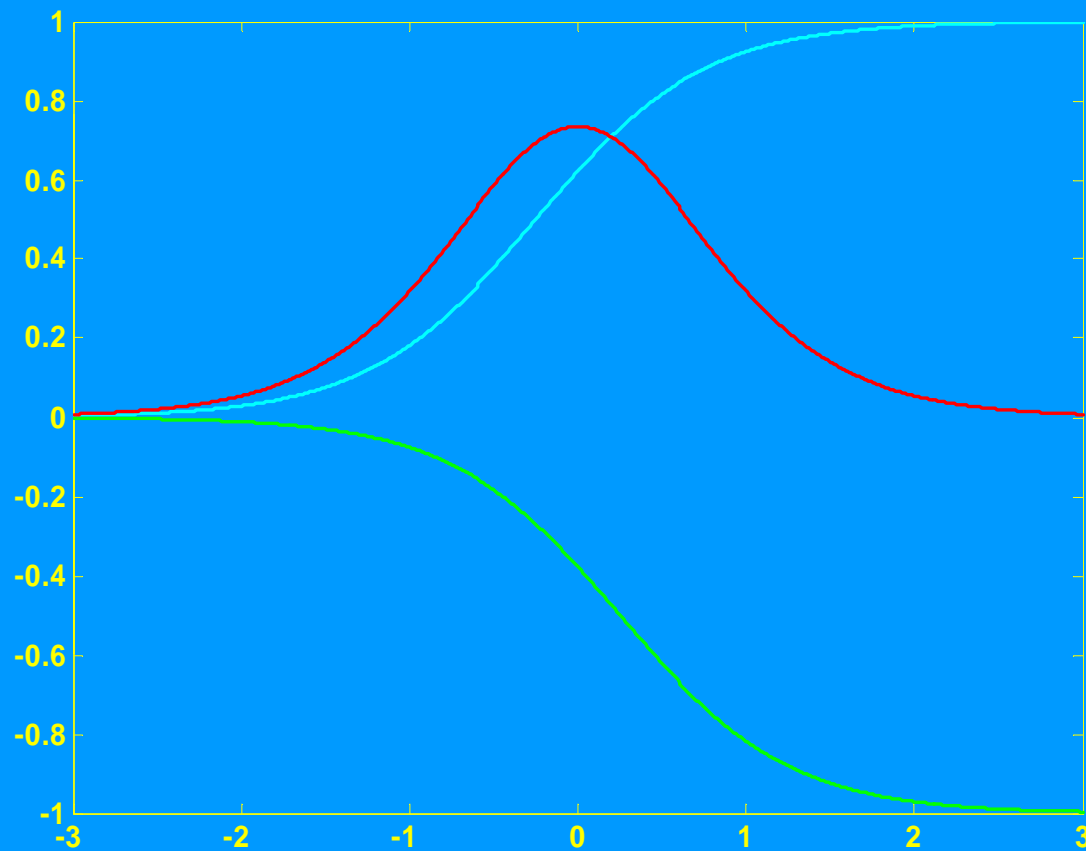
A Sigmoidal Unit





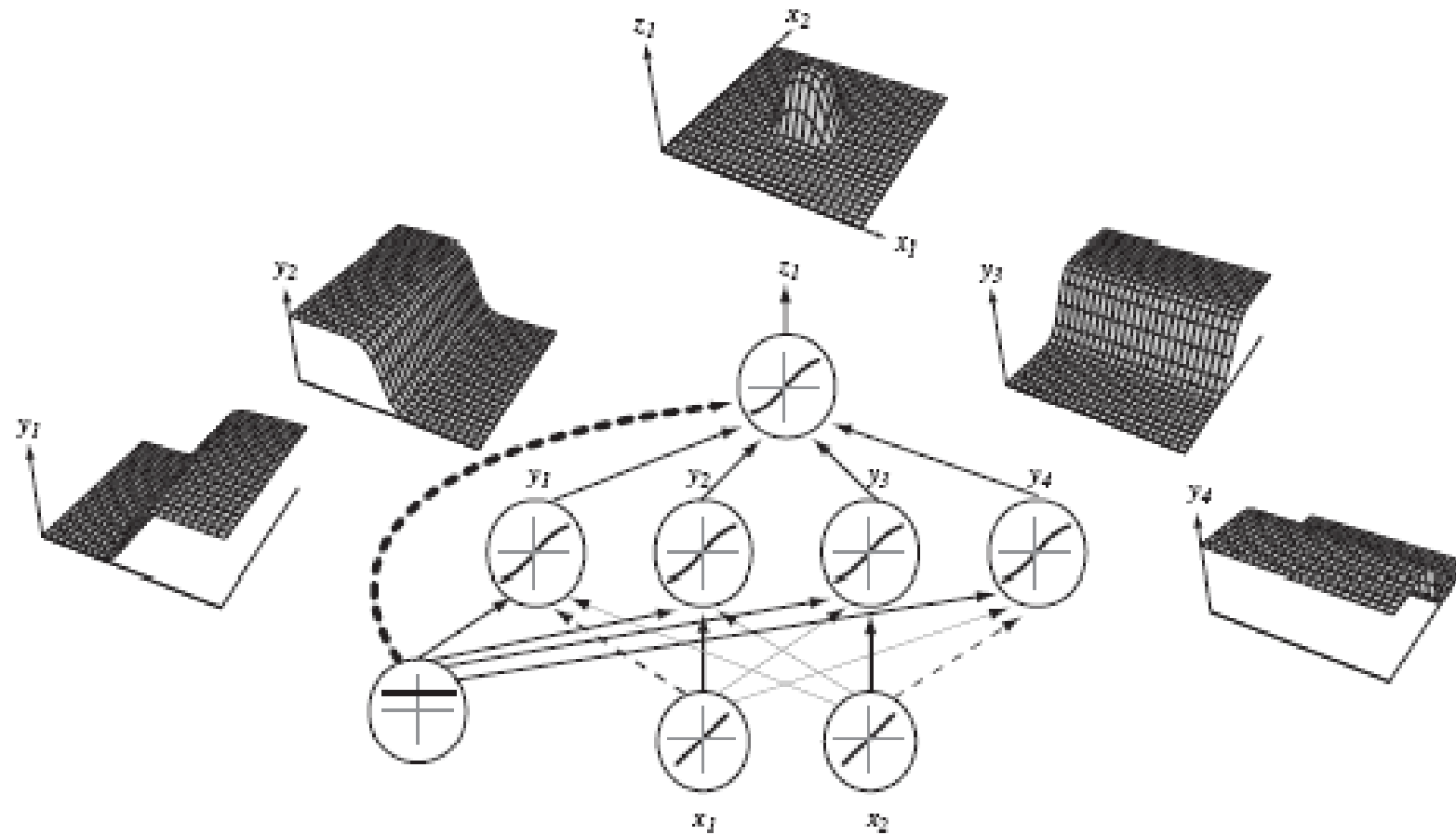
The
University
Of
Sheffield.

Combinations of Sigmoids





The
University
Of
Sheffield.






Universal Approximation

- linear combination of “enough” sigmoids
 - Cybenko, 1990
 - single hidden layer adequate
 - more may be better
 - choose hidden parameters (w , b) optimally
- problem solved?**



Interpretation

- Minimising SSE equivalent to finding *conditional mean* of target data
 - infinite sample / global minimum

$$J_{\infty} = \frac{1}{2} \sum_{i=1}^n \int_{\underline{x} \in \mathbb{R}^d} \left(E[z_i | \underline{x}] - \hat{f}_i(\underline{x}; \mathcal{W}) \right)^2 p(\underline{x}) d\underline{x}$$


MLP

$$+ \frac{1}{2} \sum_{i=1}^n \int_{\underline{x} \in \mathbb{R}^d} \left(z_i - E[z_i | \underline{x}] \right)^2 p(\underline{x}) d\underline{x}$$

does not depend
on \mathcal{W}



$$\hat{f}_i(\underline{x}; \mathcal{W}^*) = E[z_i | \underline{x}]$$

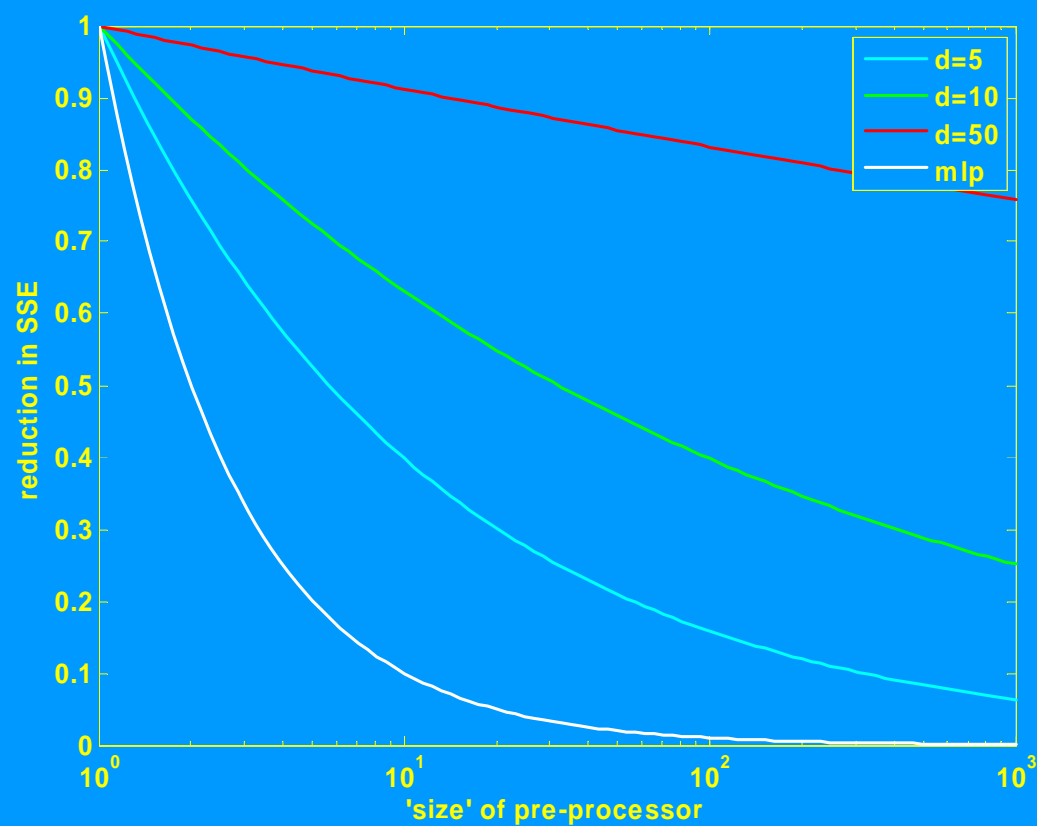


Pros

- compactness
 - potential to obtain same veracity with much smaller model
 - c.f. sparsity/complexity control in linear models
- “simple” training algorithm



Compactness of Model



$$MLP \ O(1/H)$$

$$SER \ O\left(1/H^{2/d}\right)$$



Backpropagation Algorithm

Gradient Descent

$$w_{jr}(t+1) = w_{jr}(t) + \eta(t) \delta_j(t) y_r(t) \quad j \in L_m, r \in L_{m-1} \quad \text{Update Rule}$$

Generalised Delta Rule

$$\delta_j(t) = \theta'(net_j(t)) e_j(t) \quad m = M \quad \text{output layer}$$

$$\delta_j(t) = \sigma'(net_j(t)) \sum_{i \in L_m} w_{ij}(t) \delta_i(t) \quad j \in L_{m-1}, m \in \text{hidden layers}$$

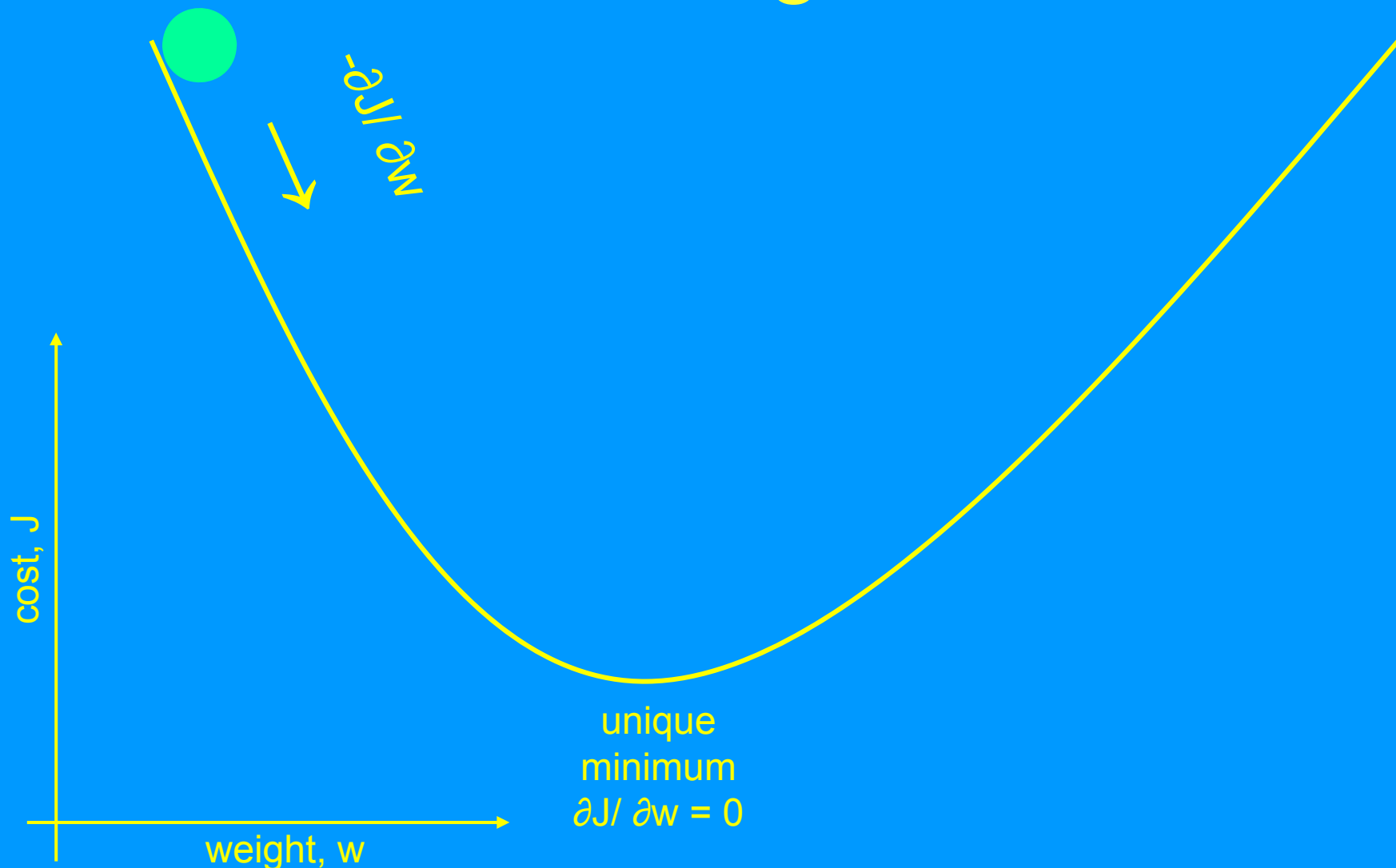


& Cons

- parameter \rightarrow cost “malign”
 - optimization difficult
 - many solutions possible
 - effect of hidden weights in output non-linear

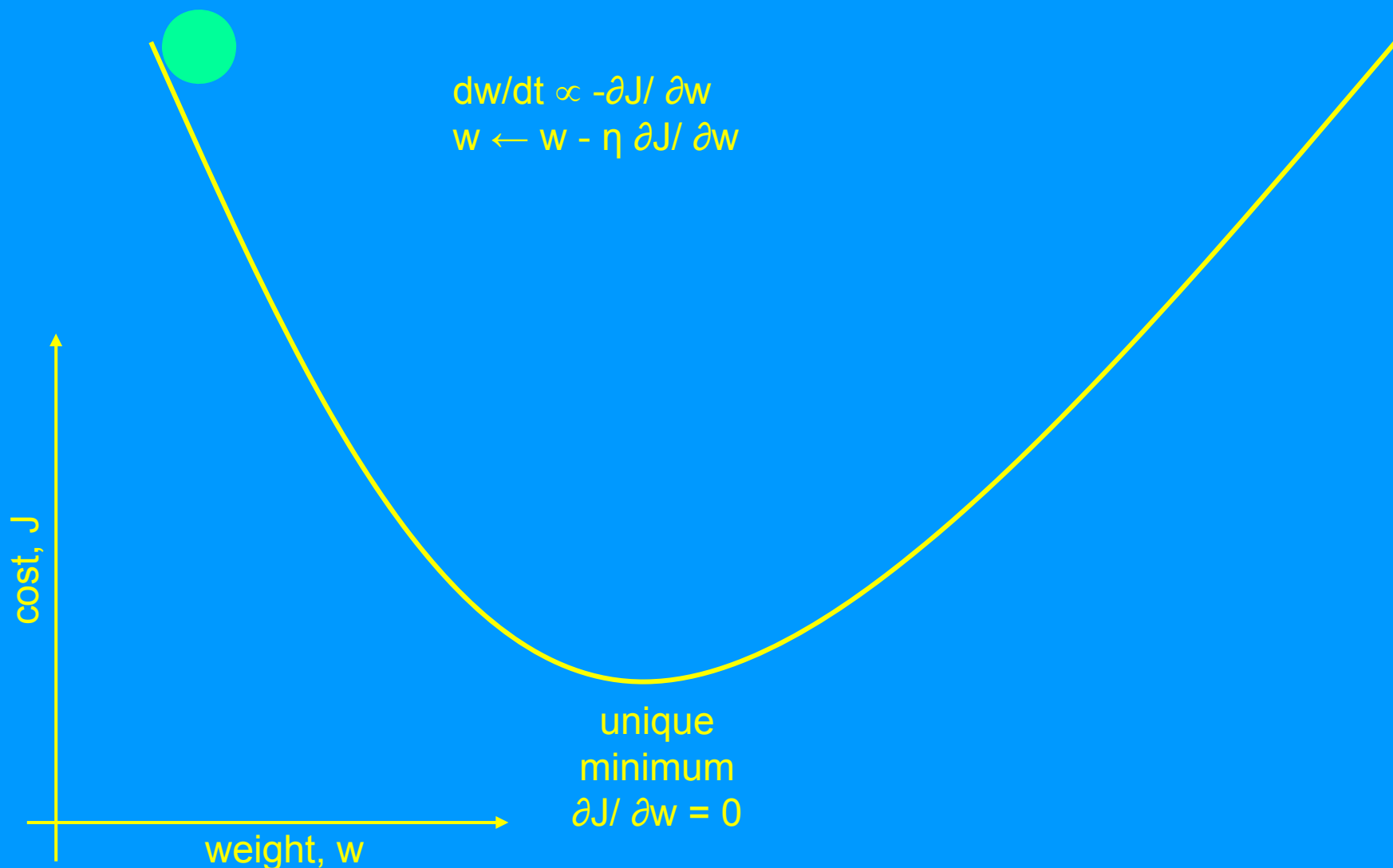


Rolling Ball



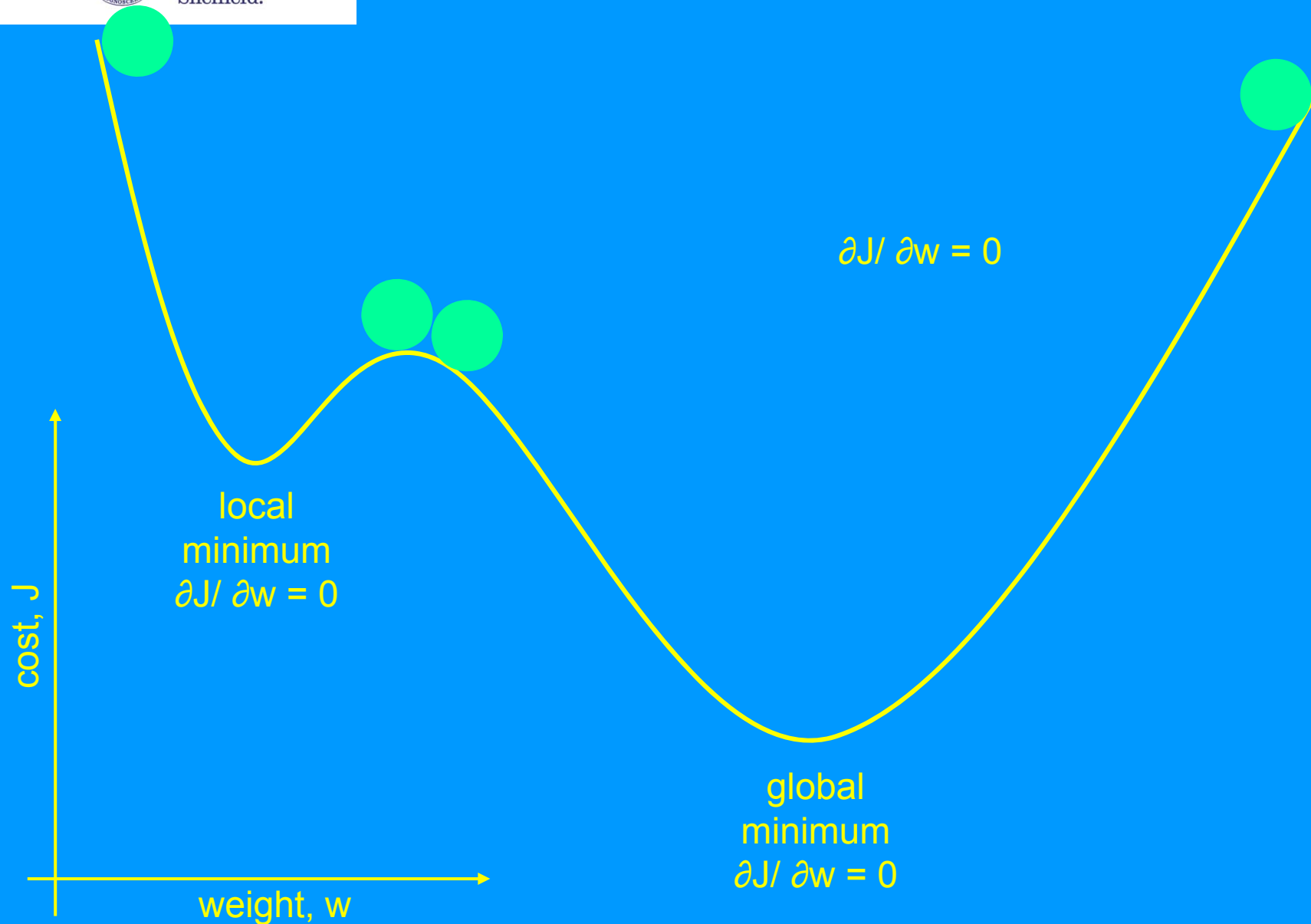


Gradient Descent





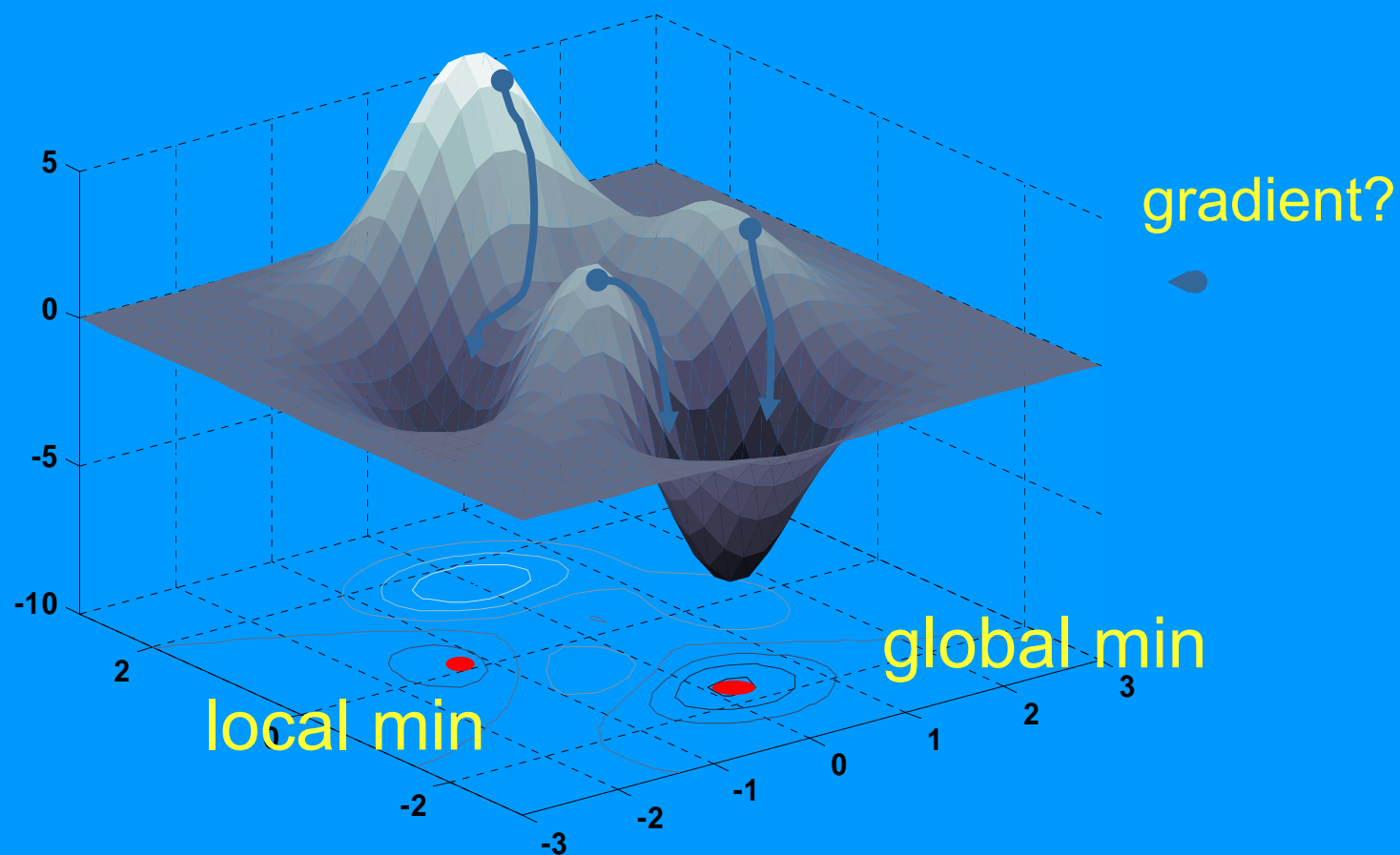
The
University
Of
Sheffield.





The
University
Of
Sheffield.

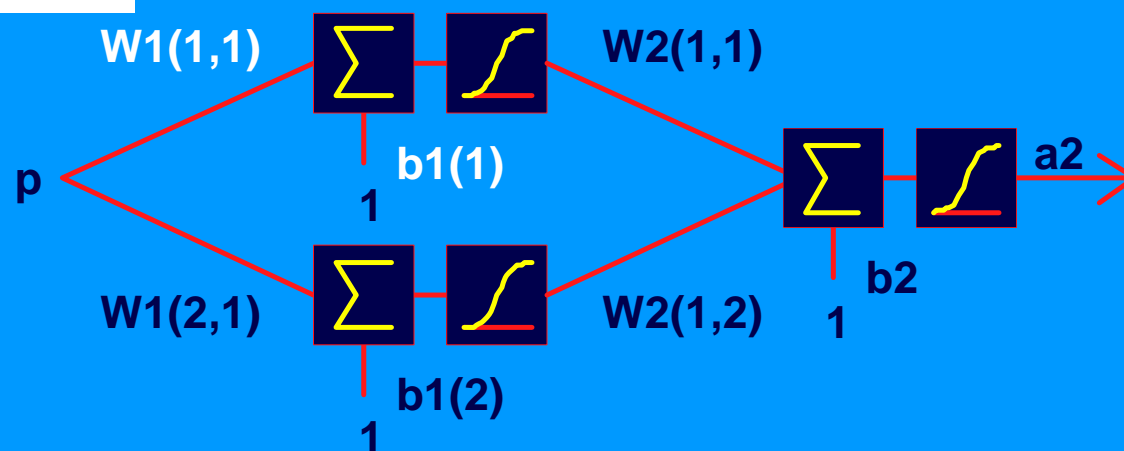
Multi-Modal Cost Surface



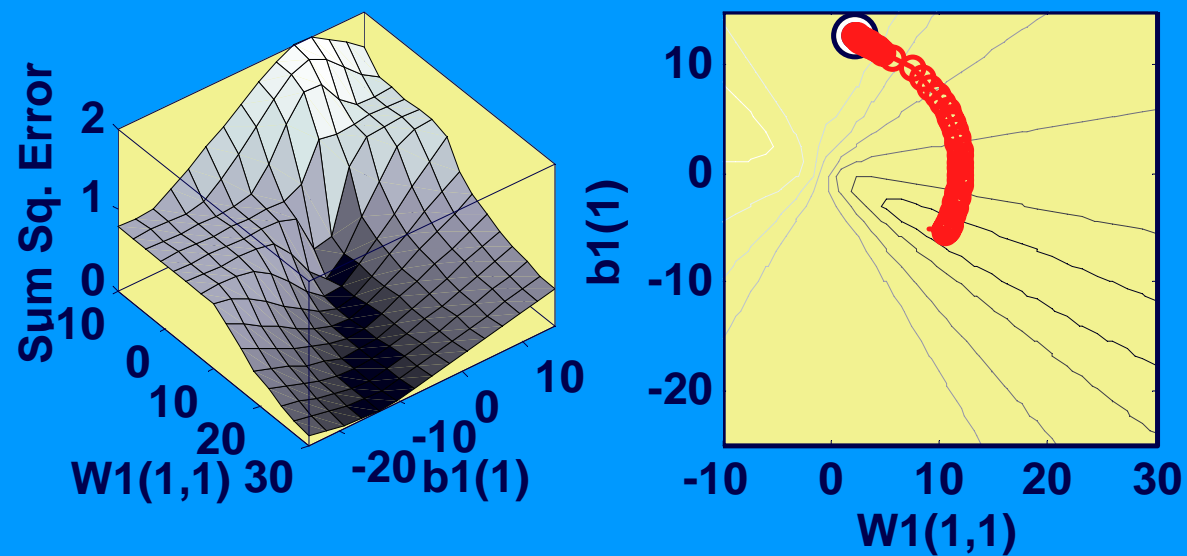


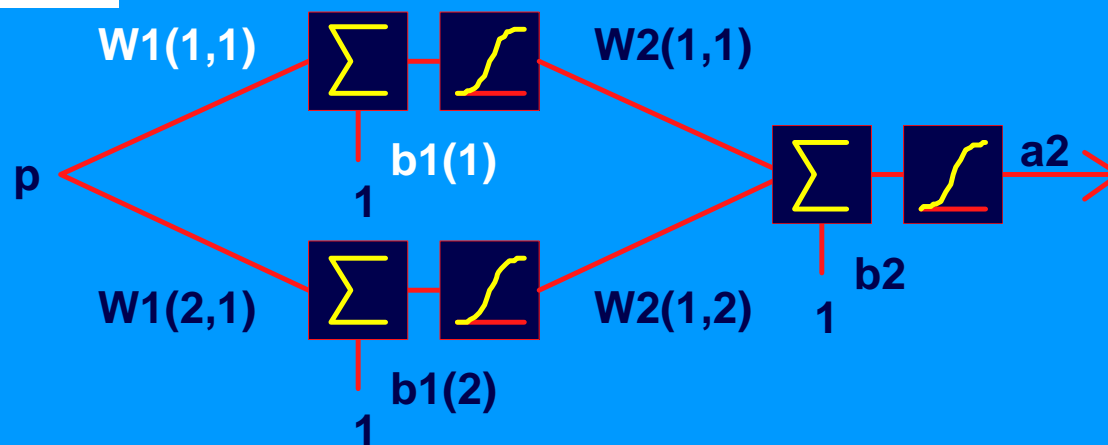
Heading Downhill

- assume
 - minimization (e.g. SSE)
 - analytically intractable
- step parameters downhill
- $W_{\text{new}} = W_{\text{old}} + \text{step in right direction}$
- backpropagation (of error)
 - slow but efficient
- conjugate gradients, Levenburg/Marquardt
 - for preference

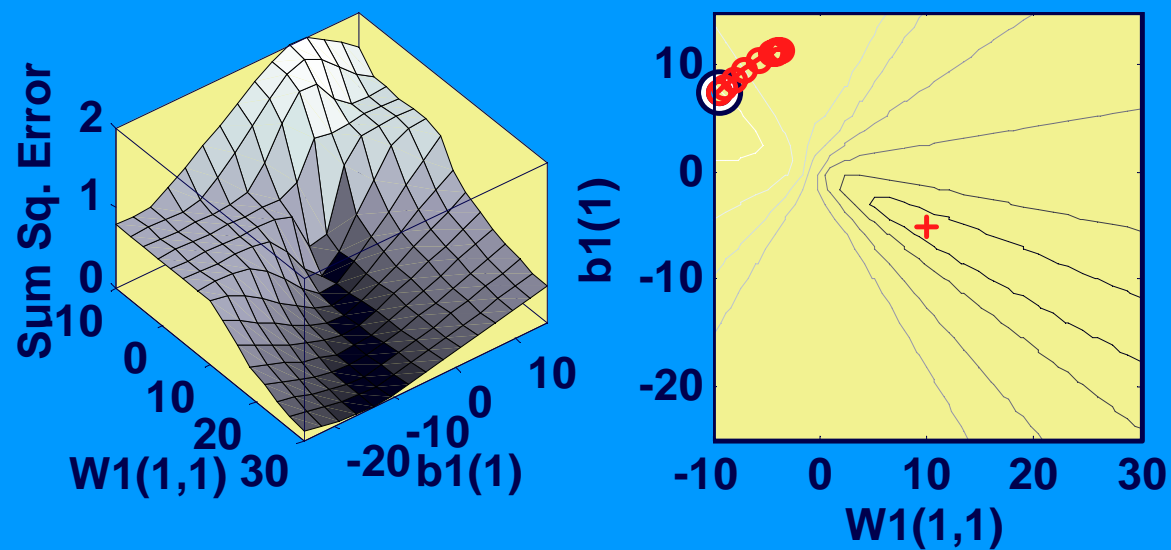


☐ $W1(1,1), W2(1,1)$ ☒ $W1(1,1), b1(1)$ ☐ $b1(1), b1(2)$





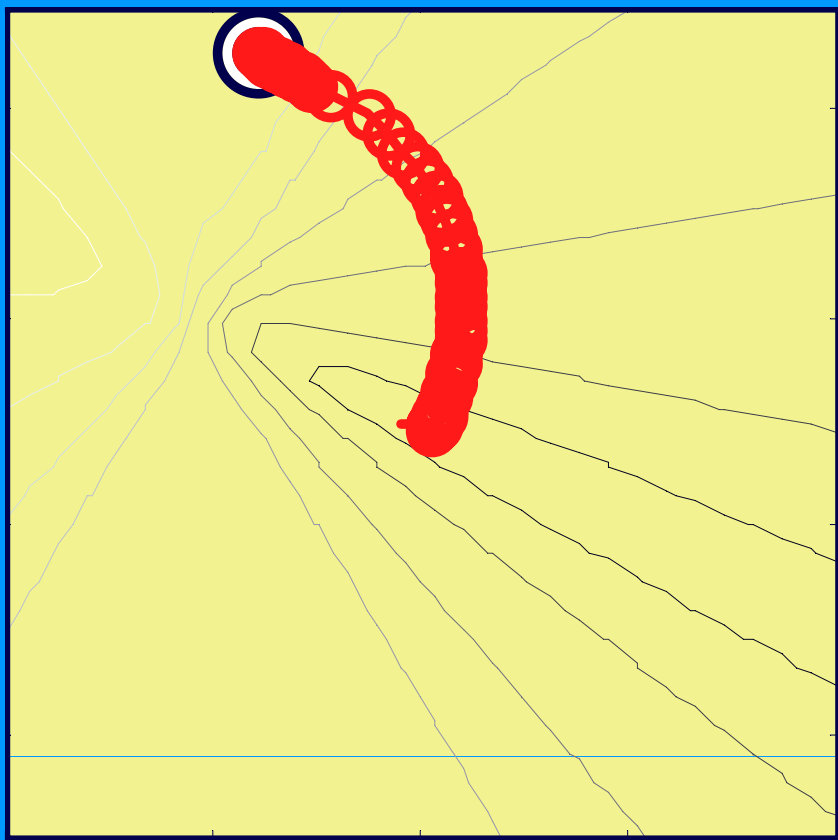
☐ $W1(1,1), W2(1,1)$ ☒ $W1(1,1), b1(1)$ ☐ $b1(1), b1(2)$



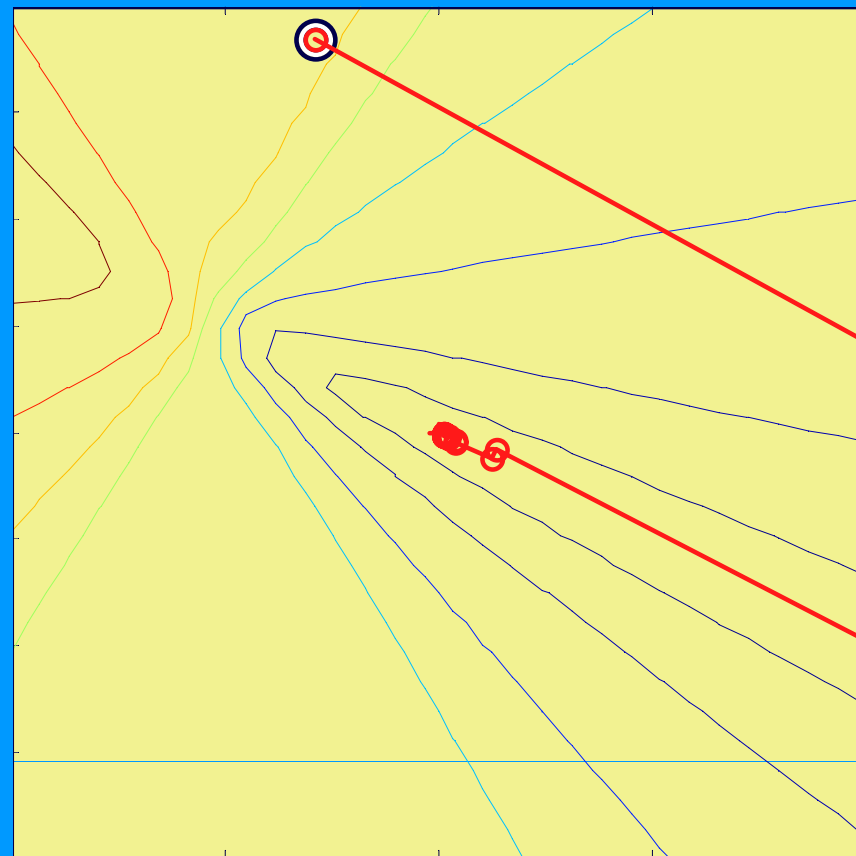


The
University
Of
Sheffield.

backprop



conjugate gradients





Implications

- correct structure can get “wrong” answer
 - dependency on initial conditions
 - might be good enough
- train / test (cross-validation) required
 - is poor behaviour due to
network structure?
ICs?

additional dimension in development



RBF NN Warning!

- RBF NNs claimed to have unique solution
- BUT**
- who picks the pre-processing layer?
 - direct optimisation of centres and widths
 - some other method
- L.I.P. models have “non-linear parameters” to select \rightarrow multi-modal cost \equiv MLP



Are multiple minima a problem?

- pros seem to outweigh cons
- good solutions often arrived at quickly
- all previous issues apply
 - sample density & distribution
 - lack of prior knowledge



How to Use

- to “generalize” a GLM
 - linear regression – curve-fitting
linear output + SSE
 - logistic regression – classification
logistic output + cross-entropy (deviance)
extend to multinomial, ordinal
e.g. softmax output + cross entropy
 - Poisson regression – count data



What is Learned?

- the right thing
 - in a maximum likelihood sense

theoretical

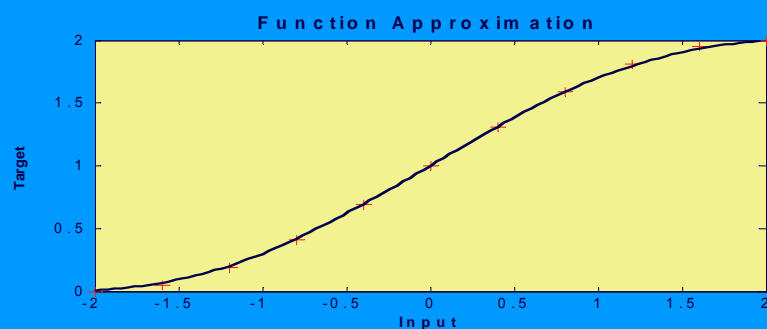
- conditional mean of target data, $E(z|x)$
 - implies probability of class membership for classification $P(C_i|x)$

estimated

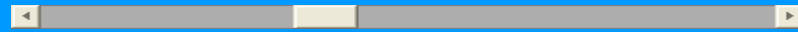
- if good estimate then $y \rightarrow E(z|x)$



Simple 1-D Function



Number of Hidden Neurons S1:



1

4

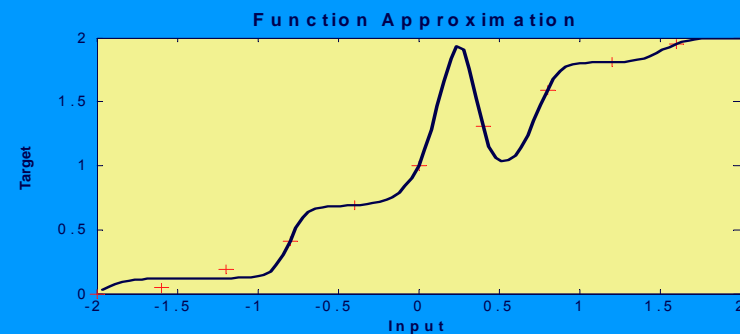
9

Difficulty Index:



1

9



Number of Hidden Neurons S1:



1

9

9

Difficulty Index:

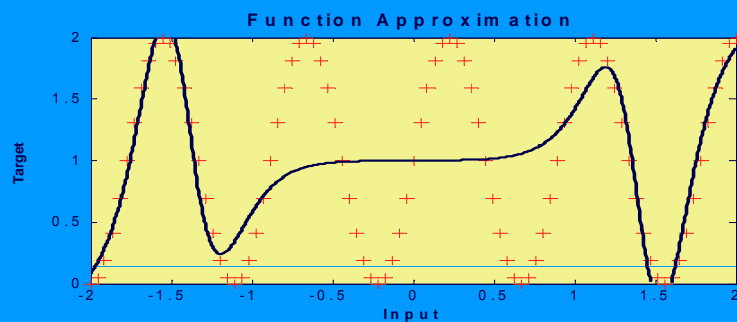


1

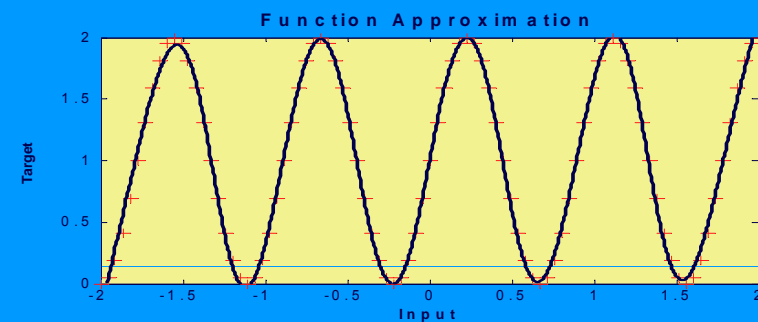
9



More Complex 1-D Example



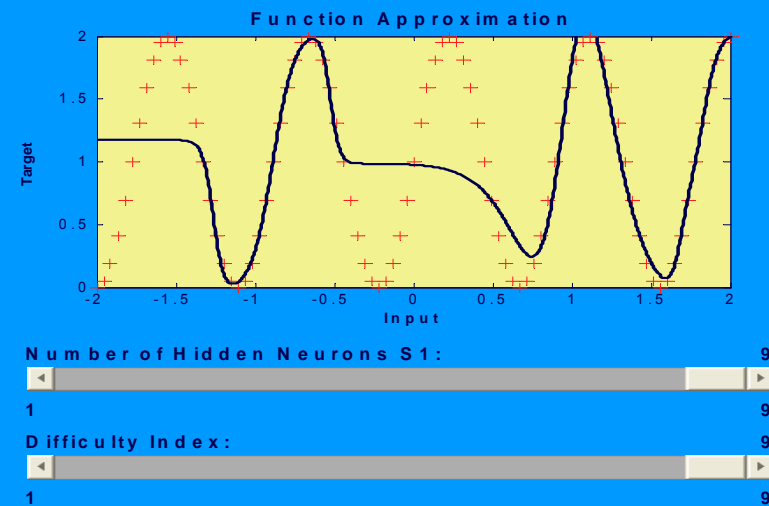
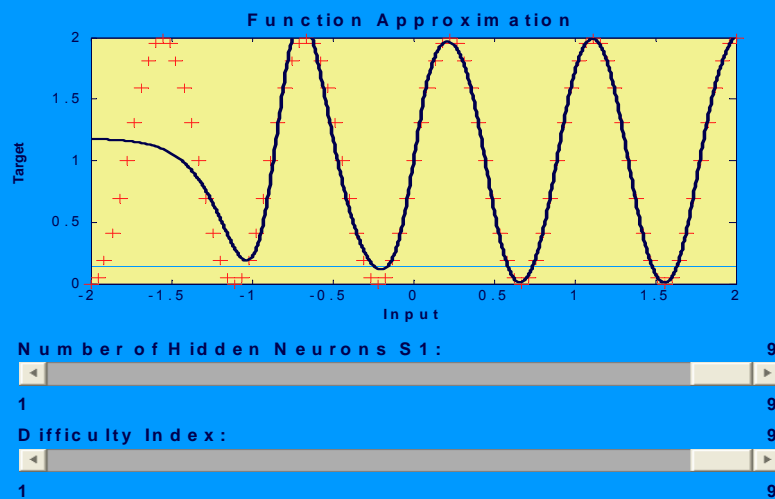
Number of Hidden Neurons S1: 4
1 9
Difficulty Index: 9
1 9



Number of Hidden Neurons S1: 7
1 9
Difficulty Index: 9
1 9



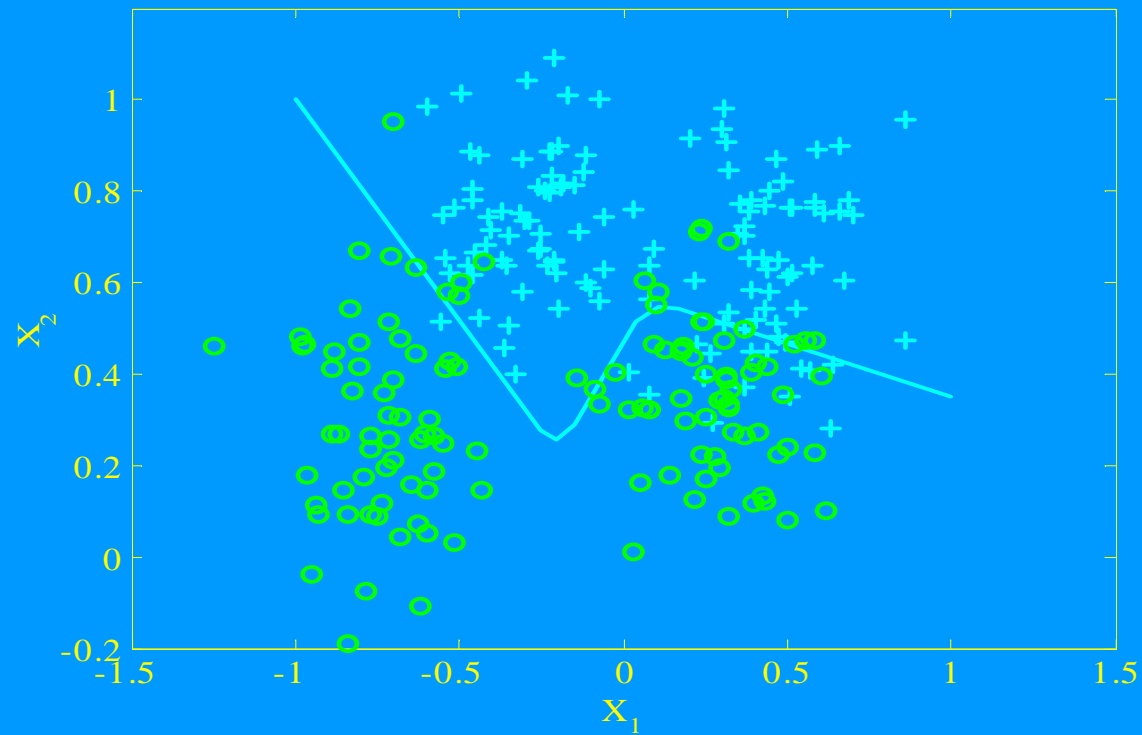
Local Solutions





The
University
Of
Sheffield.

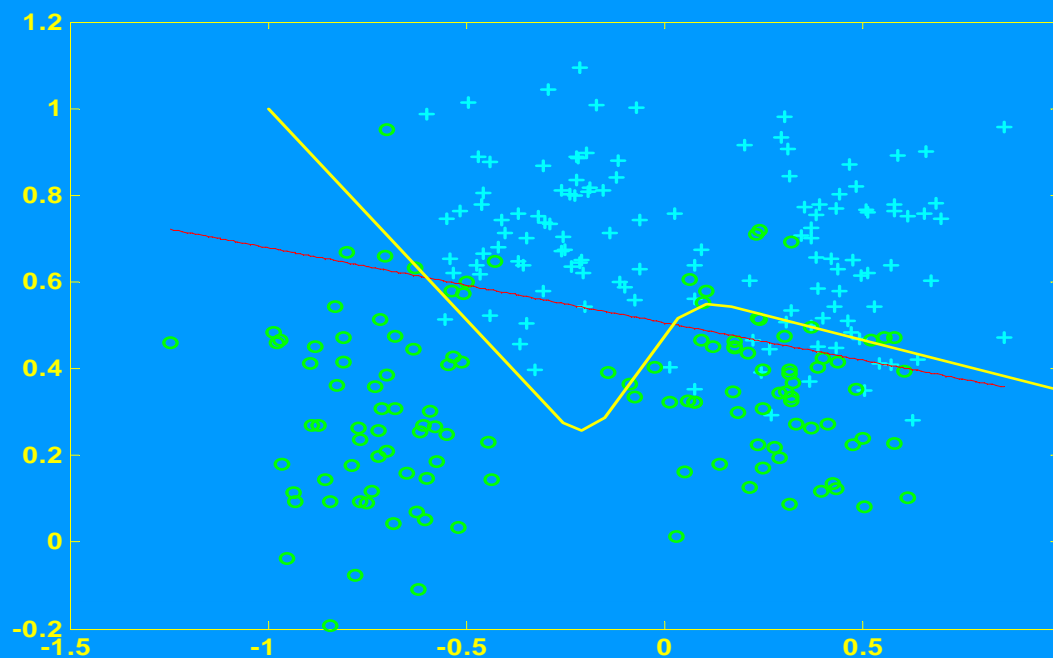
A Dichotomy



data courtesy B Ripley



Linear Decision Boundary

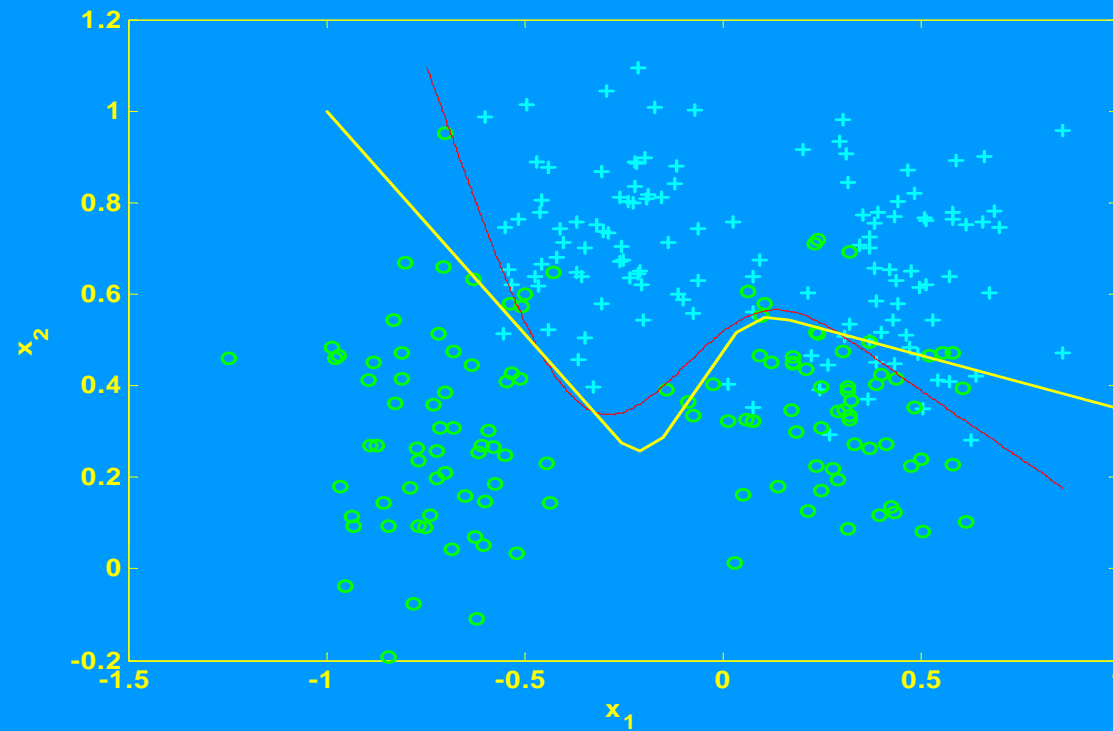


induced by $P(C_1|x) = P(C_2|x)$



The
University
Of
Sheffield.

Non-Linear Decision Boundary





The
University
Of
Sheffield.

Over-fitted Decision Boundary

