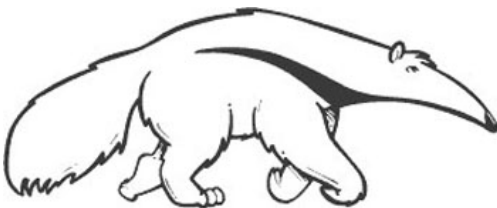


+

Machine Learning and Data Mining

Multi-layer Perceptrons & Neural Networks: Basics

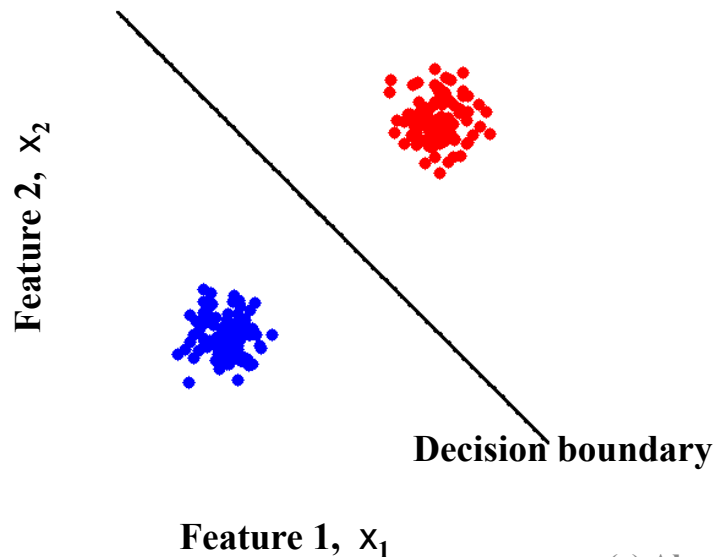
Prof. Alexander Ihler



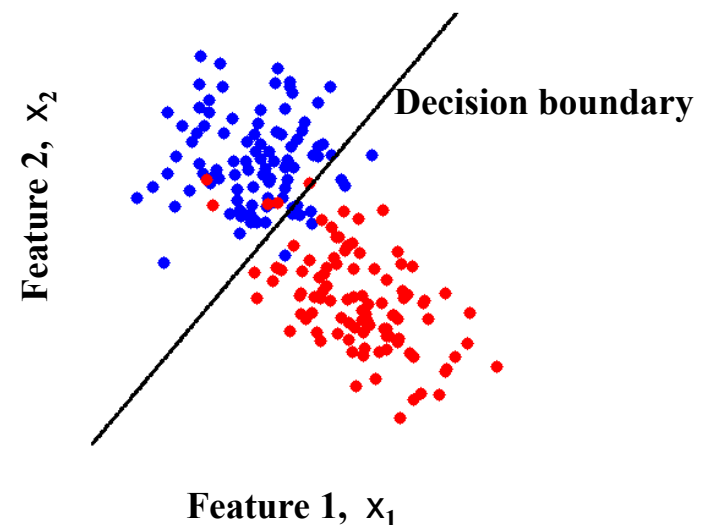
Linear Classifiers (Perceptrons)

- Linear Classifiers
 - a linear classifier is a mapping which partitions feature space using a linear function (a straight line, or a hyperplane)
 - separates the two classes using a straight line in feature space
 - in 2 dimensions the decision boundary is a straight line

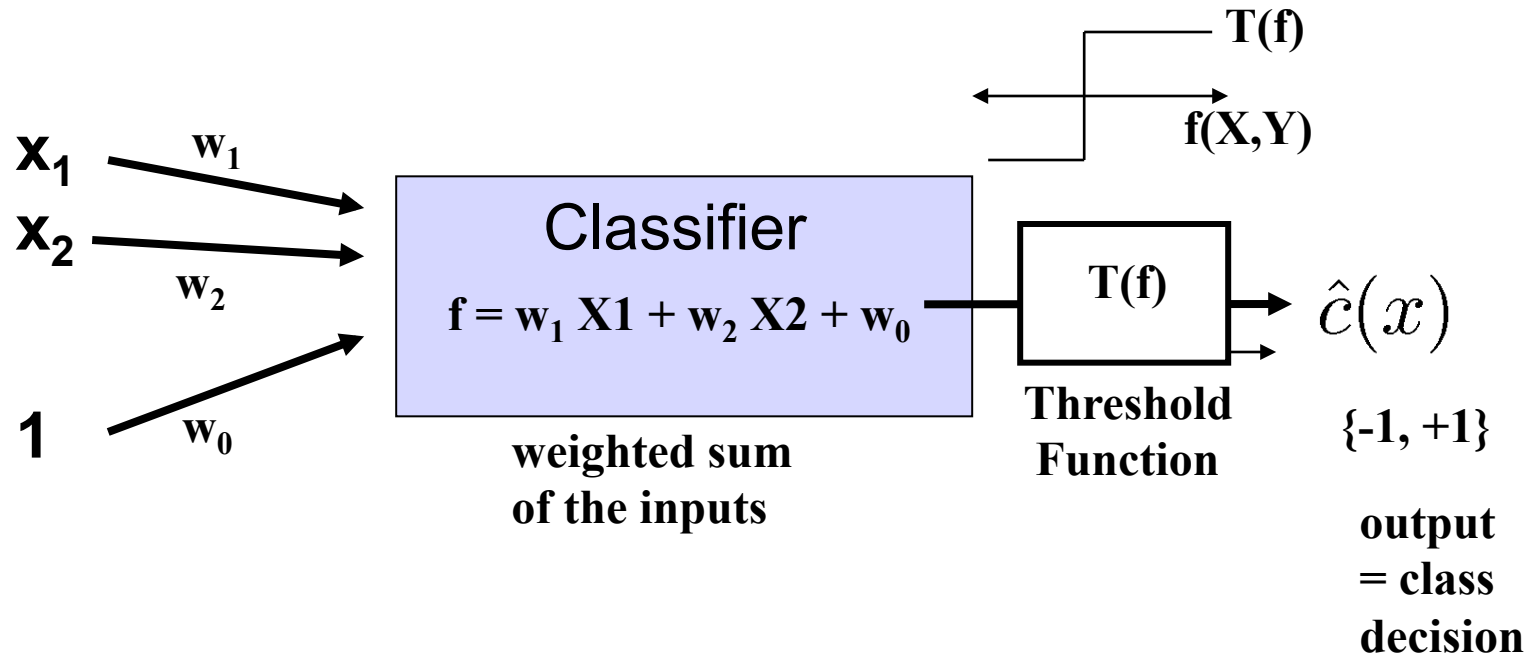
Linearly separable data



Linearly non-separable data



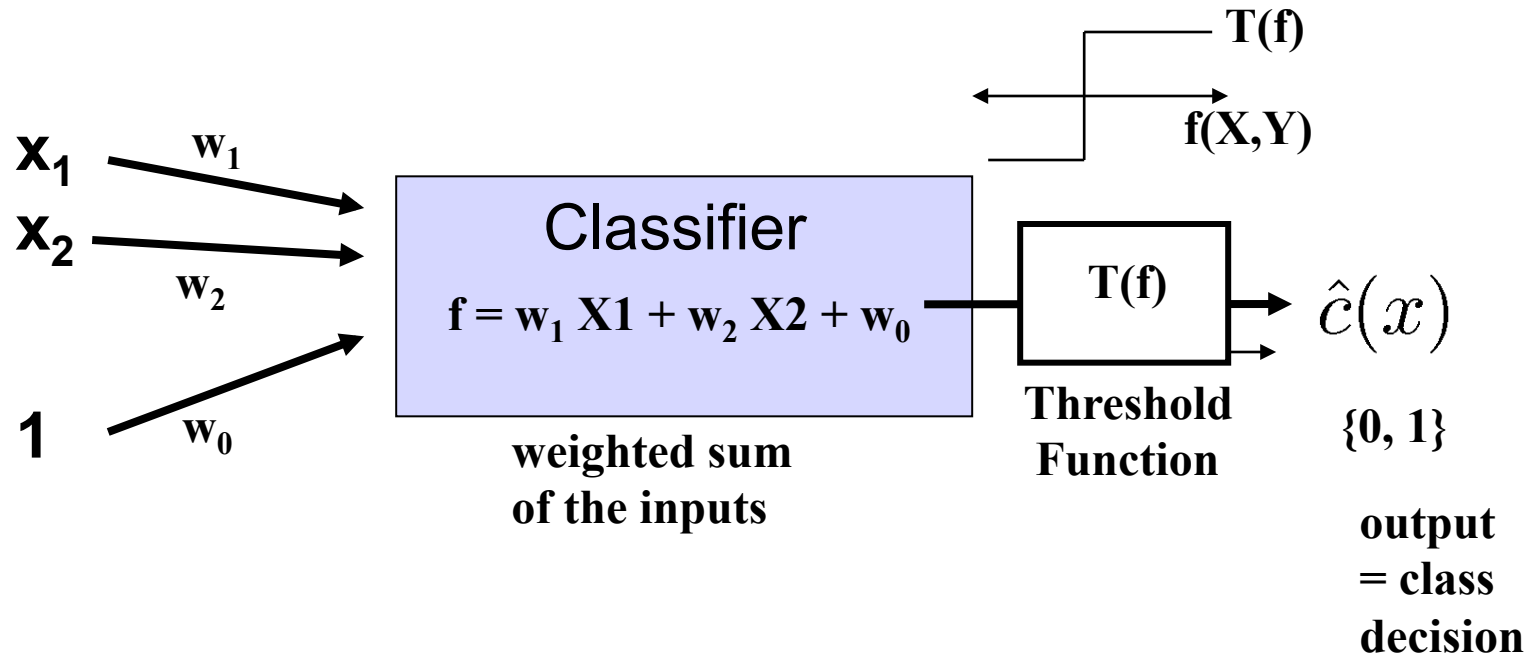
Perceptron Classifier (2 features)



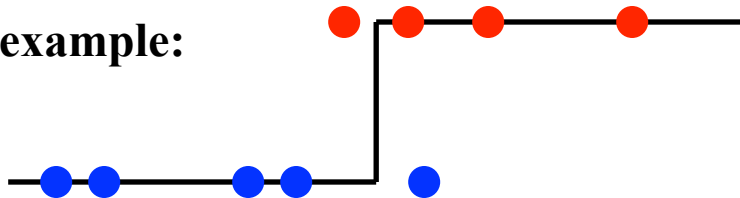
Decision Boundary at $f(x) = 0$

Solve: $X_2 = -w_1/w_2 X_1 - w_0/w_2$ (Line)

Perceptron (Linear classifier)



1D example:

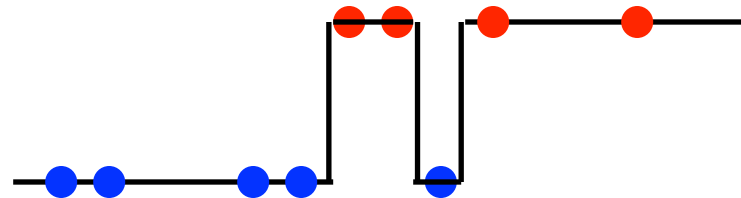
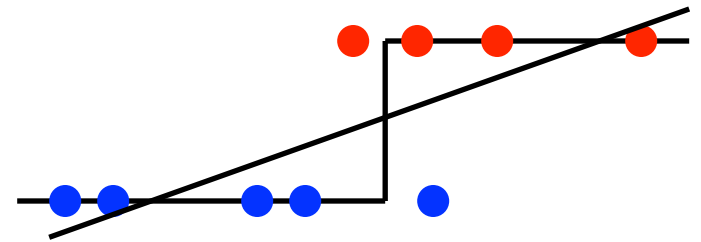


$$\begin{aligned} T(f) &= 0 \text{ if } f < 0 \\ T(f) &= 1 \text{ if } f > 0 \end{aligned}$$

Decision boundary = "x such that $T(w_1 x + w_0)$ transitions"

Features and perceptrons

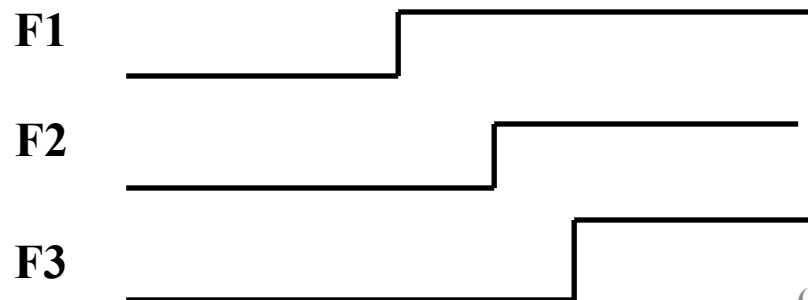
- Recall the role of features
 - We can create extra features that allow more complex decision boundaries
 - Linear classifiers
 - Features $[1, x]$
 - Decision rule: $T(ax+b) = ax + b > /< 0$
 - Boundary $ax+b=0 \Rightarrow$ point
 - Features $[1, x, x^2]$
 - Decision rule $T(ax^2+bx+c)$
 - Boundary $ax^2+bx+c = 0 = ?$
 - What features can produce this decision rule?



(c) Alexander Ihler

Features and perceptrons

- Recall the role of features
 - We can create extra features that allow more complex decision boundaries
 - For example, polynomial features
$$\Phi(x) = [1 \ x \ x^2 \ x^3 \dots]$$
- What other kinds of features could we choose?
 - Step functions?



Linear function of features
 $a F1 + b F2 + c F3 + d$

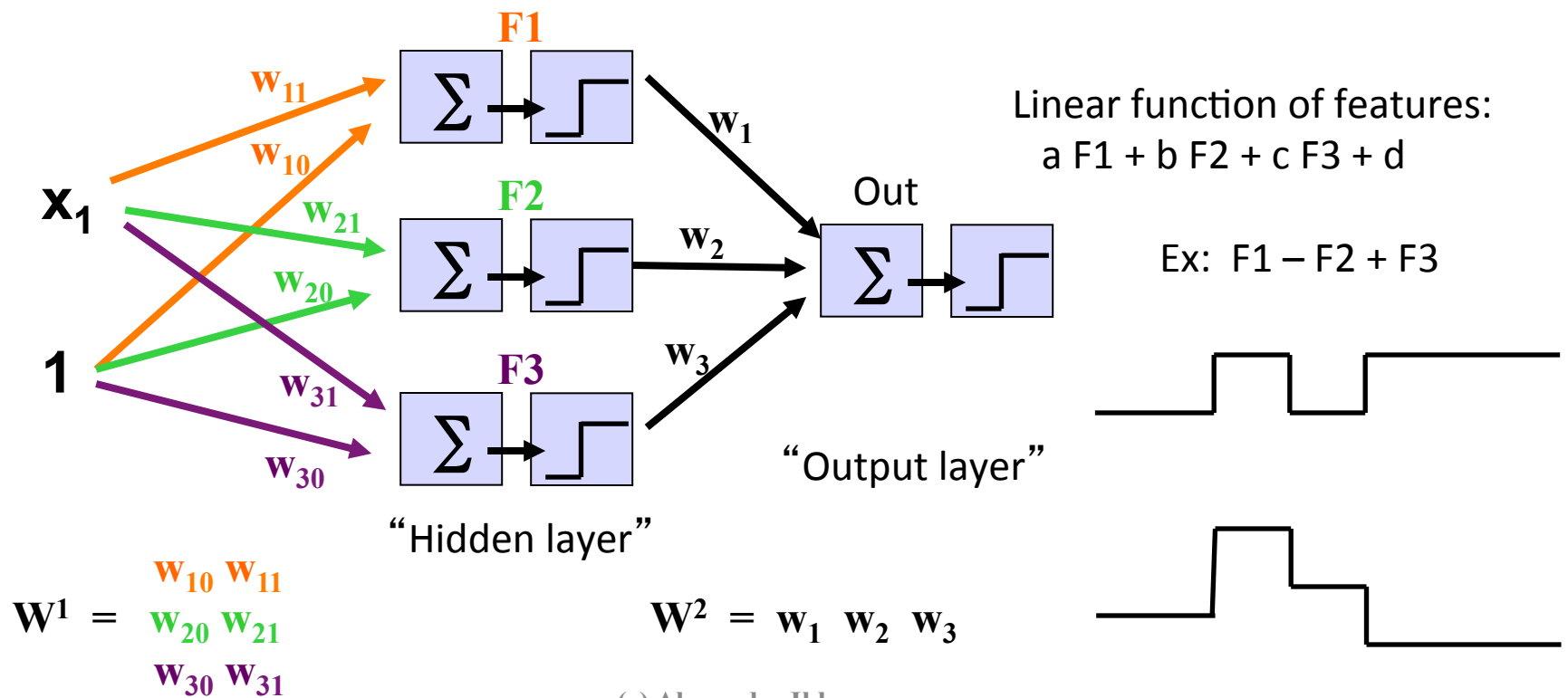
Ex: $F1 - F2 + F3$



(c) Alexander Ihler

Multi-layer perceptron model

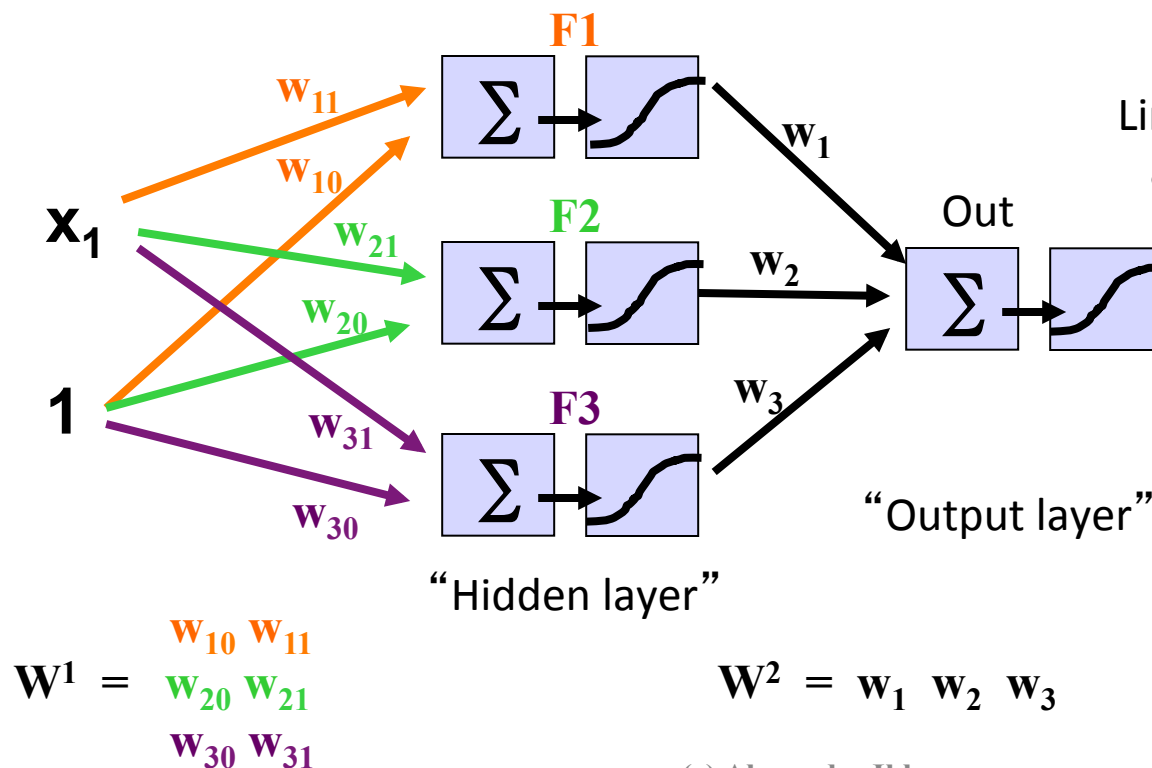
- Step functions are just perceptrons!
 - “Features” are outputs of a perceptron
 - Combination of features output of another



Multi-layer perceptron model

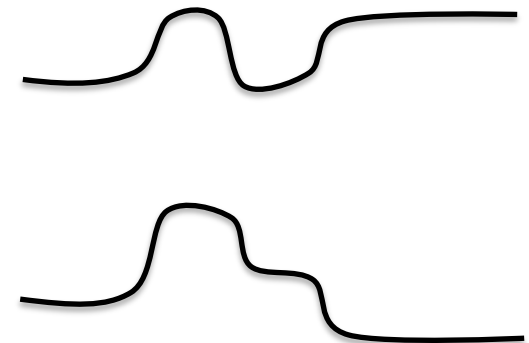
- Step functions are just perceptrons!
 - “Features” are outputs of a perceptron
 - Combination of features output of another

Regression version:
Remove activation
function from output



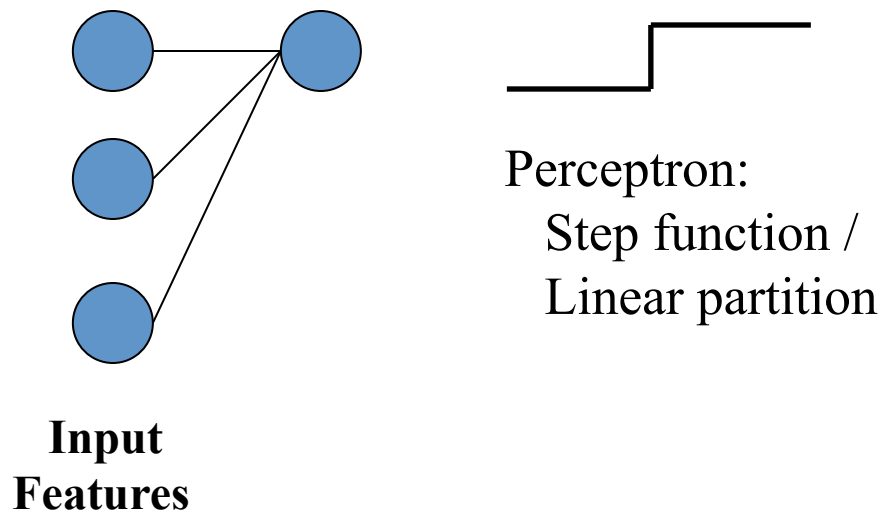
Linear function of features:
 $a F1 + b F2 + c F3 + d$

Ex: $F1 - F2 + F3$



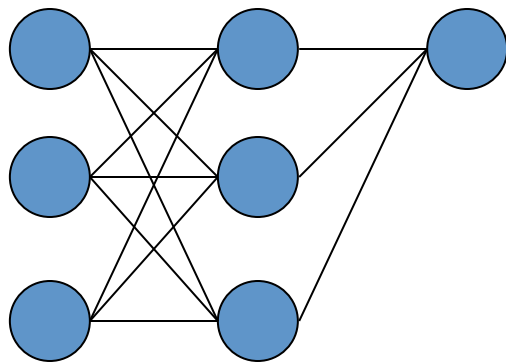
Features of MLPs

- Simple building blocks
 - Each element is just a perceptron f' n
- Can build upwards



Features of MLPs

- Simple building blocks
 - Each element is just a perceptron f' n
- Can build upwards



Input
Features



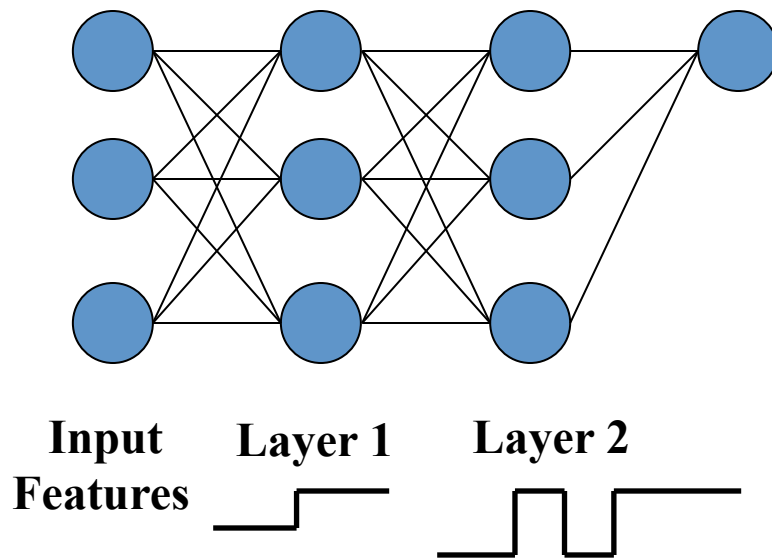
2-layer:

“Features” are now partitions

All linear combinations of those partitions

Features of MLPs

- Simple building blocks
 - Each element is just a perceptron f' n
- Can build upwards

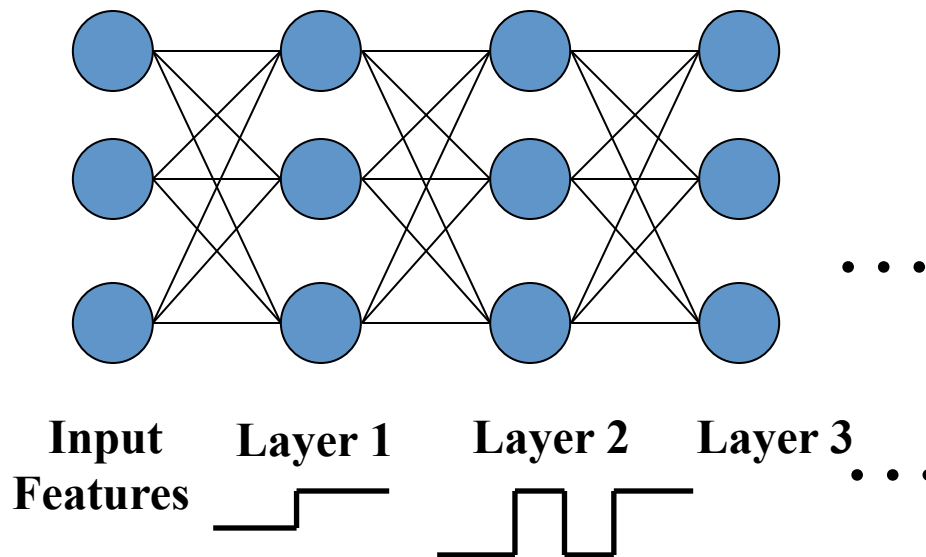


3-layer:

“Features” are now complex functions
Output any linear combination of those

Features of MLPs

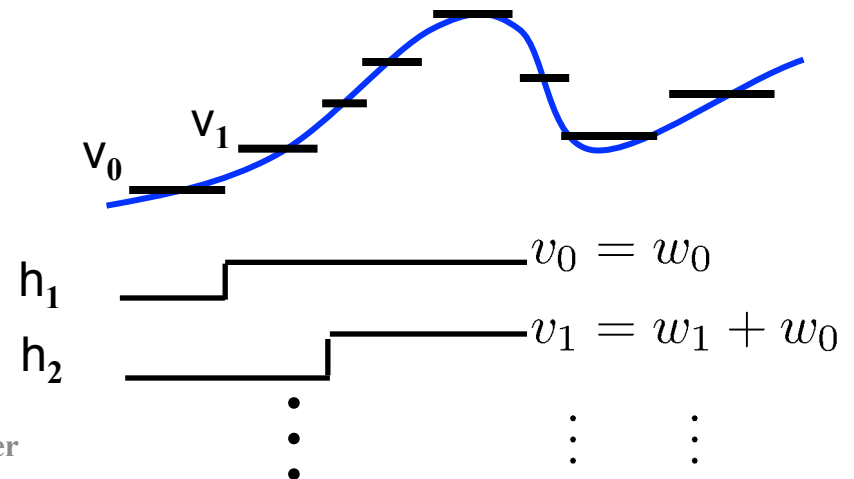
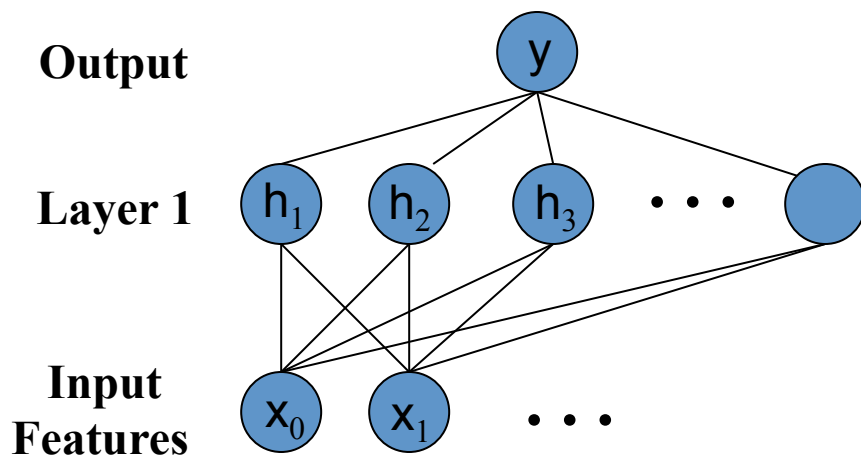
- Simple building blocks
 - Each element is just a perceptron f' n
- Can build upwards



Current research:
“Deep” architectures
(many layers)

Features of MLPs

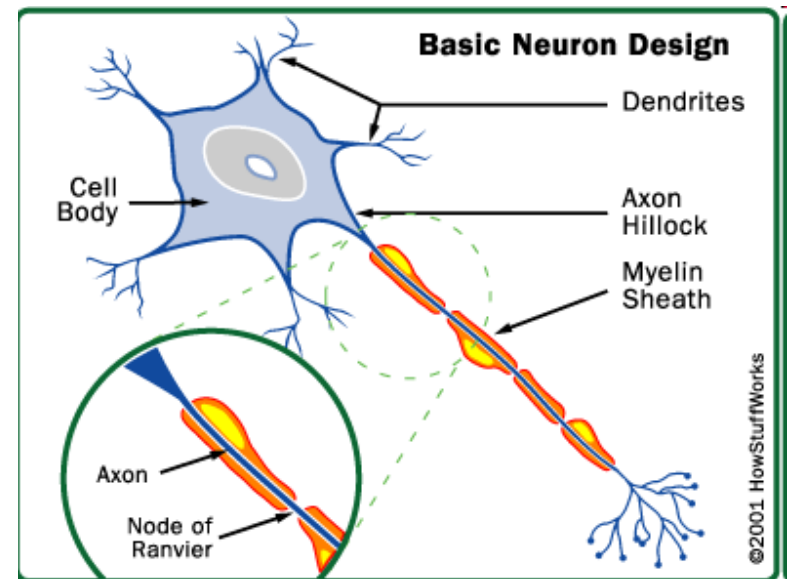
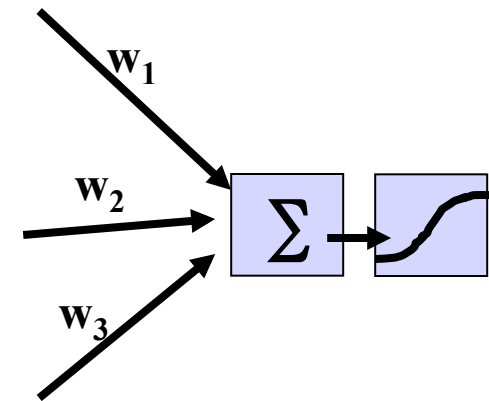
- Simple building blocks
 - Each element is just a perceptron function
- Can build upwards
- Flexible function approximation
 - Approximate arbitrary functions with enough hidden nodes



(c) Alexander Ihler

Neural networks

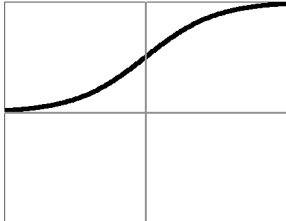
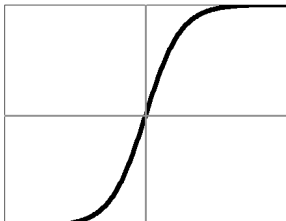
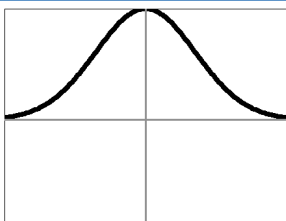
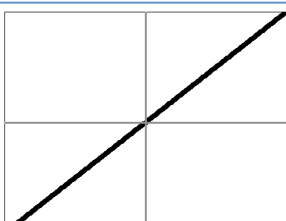
- Another term for MLPs
- Biological motivation
- Neurons
 - “Simple” cells
 - Dendrites sense charge
 - Cell weighs inputs
 - “Fires” axon



(c) Alexander Ihler

“How stuff works: the brain”

Activation functions

Logistic	$\sigma(z) = \frac{1}{1 + \exp(-z)}$		$\frac{\partial \sigma}{\partial z}(z) = \sigma(z)(1 - \sigma(z))$
Hyperbolic Tangent	$\sigma(z) = \frac{1 - \exp(-2z)}{1 + \exp(-2z)}$		$\frac{\partial \sigma}{\partial z}(z) = 1 - (\sigma(z))^2$
Gaussian	$\sigma(z) = \exp(-z^2/2)$		$\frac{\partial \sigma}{\partial z}(z) = -z\sigma(z)$
Linear	$\sigma(z) = z$		$\frac{\partial \sigma}{\partial z}(z) = 1$

And many others...

Feed-forward networks

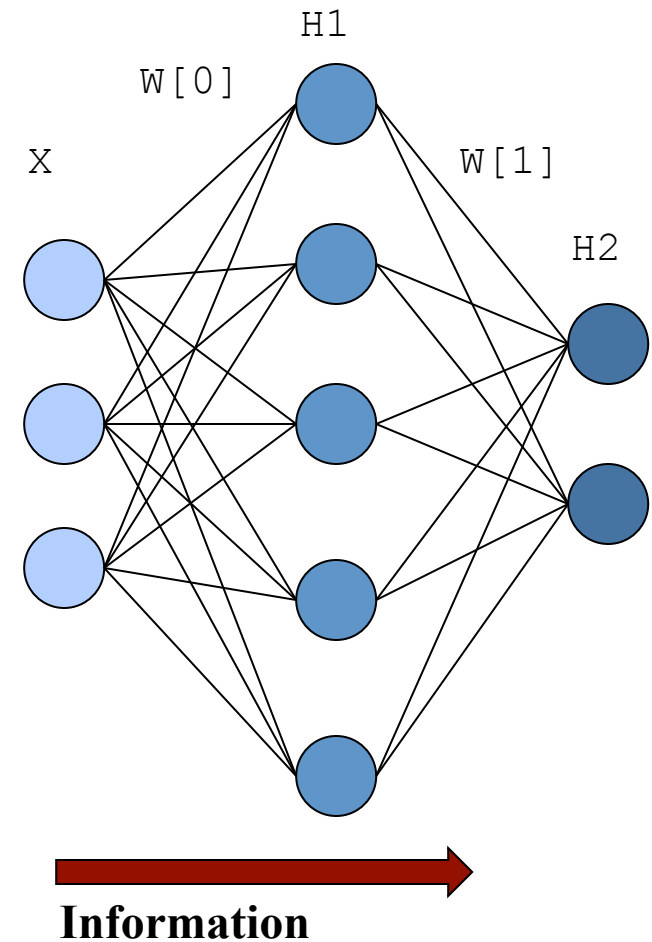
- Information flows left-to-right
 - Input observed features
 - Compute hidden nodes (parallel)
 - Compute next layer...

```
X1 = _add1(X);      # add constant feature
T  = X1.dot(W[0].T); # linear response
H  = Sig( T );      # activation f'n

H1 = _add1(H);      # add constant feature
S  = H1.dot(W[1].T); # linear response
H2 = Sig( S );      # activation f'n

% ...
```

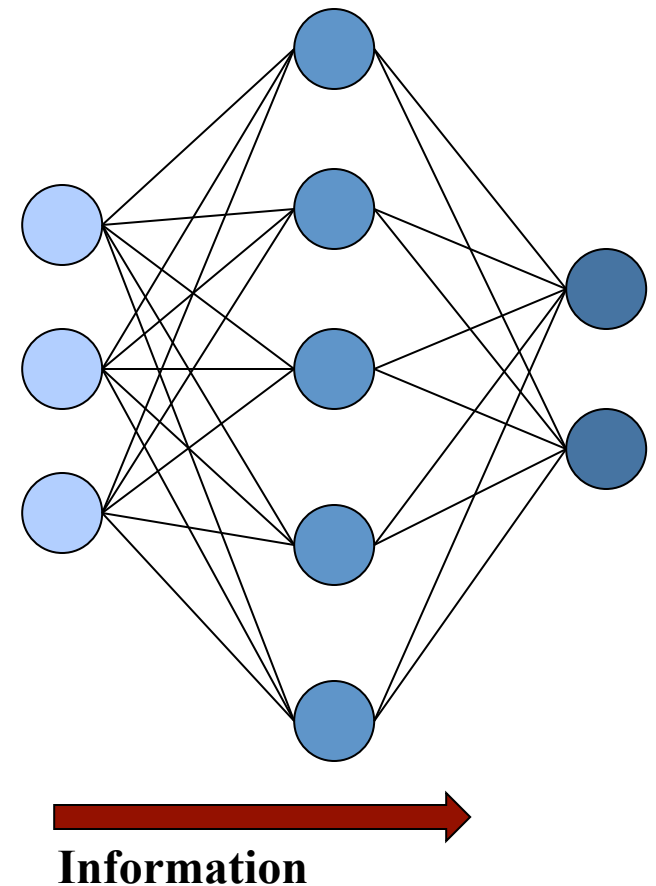
- Alternative: recurrent NNs...



Feed-forward networks

A note on multiple outputs:

- Regression:
 - Predict multi-dimensional y
 - “Shared” representation
= fewer parameters
- Classification
 - Predict binary vector
 - Multi-class classification
 $y = 2 = [0 \ 0 \ 1 \ 0 \ \dots]$
 - Multiple, joint binary predictions
(image tagging, etc.)
 - Often trained as regression (MSE),
with saturating activation

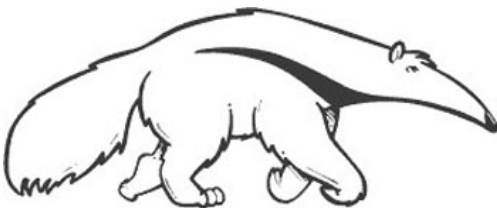


+

Machine Learning and Data Mining

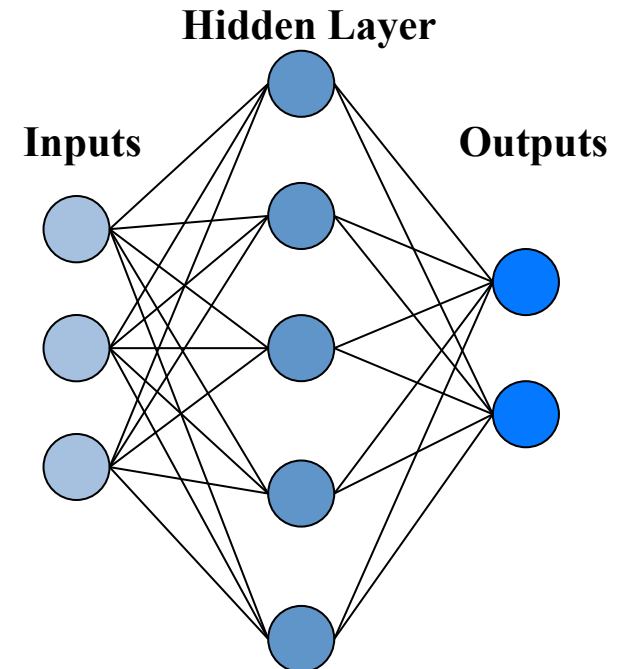
Multi-layer Perceptrons & Neural Networks: Backpropagation

Prof. Alexander Ihler



Training MLPs

- Observe features “x” with target “y”
 - Push “x” through NN = output is “ \hat{y} ”
 - Error: $(y - \hat{y})^2$ (Can use different loss functions if desired...)
 - How should we update the weights to improve?
-
- Single layer
 - Logistic sigmoid function
 - Smooth, differentiable
 - Optimize using:
 - Batch gradient descent
 - Stochastic gradient descent



Backpropagation

- Just gradient descent...
- Apply the chain rule to the MLP

$$\frac{\partial J}{\partial w_{kj}^2} = -2 \sum_{k'} (y_{k'} - \hat{y}_{k'}) (\partial \hat{y}_{k'})$$

$$= -2(y_k - \hat{y}_k) \sigma'(s_k) h_j \quad (\text{Identical to logistic mse regression with inputs "h}_j\text{"})$$

Forward pass

Loss function

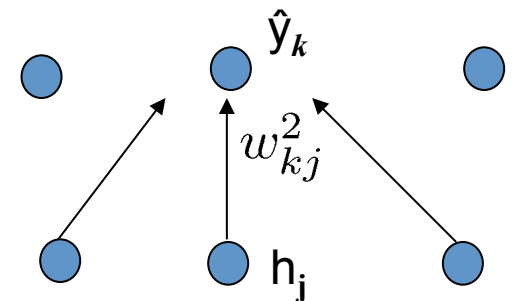
$$J_i(W) = \sum_k (y_k^{(i)} - \hat{y}_k^{(i)})^2$$

Output layer

$$\hat{y}_k = \sigma(s_k) = \sigma(\sum_j w_{kj}^2 h_j)$$

Hidden layer

$$h_j = \sigma(t_j) = \sigma(\sum_i w_{ji}^1 x_i)$$



Backpropagation

- Just gradient descent...
- Apply the chain rule to the MLP

$$\frac{\partial J}{\partial w_{kj}^2} = -2 \sum_{k'} (y_{k'} - \hat{y}_{k'}) (\partial \hat{y}_{k'})$$

$$= \underbrace{-2(y_k - \hat{y}_k) \sigma'(s_k)}_{\beta_k^2} h_j \quad (\text{Identical to logistic mse regression with inputs "h}_j\text{"})$$

$$\frac{\partial J}{\partial w_{ji}^1} = \sum_k -2(y_k - \hat{y}_k) (\partial \hat{y}_k)$$

$$= \sum_k -2(y_k - \hat{y}_k) \sigma'(s_k) w_{kj}^2 \partial h_j$$

$$= \sum_k \underbrace{-2(y_k - \hat{y}_k) \sigma'(s_k)}_{\beta_k^2} w_{kj}^2 \sigma'(t_j) x_i$$

(c) Alexander Ihler

Forward pass

Loss function

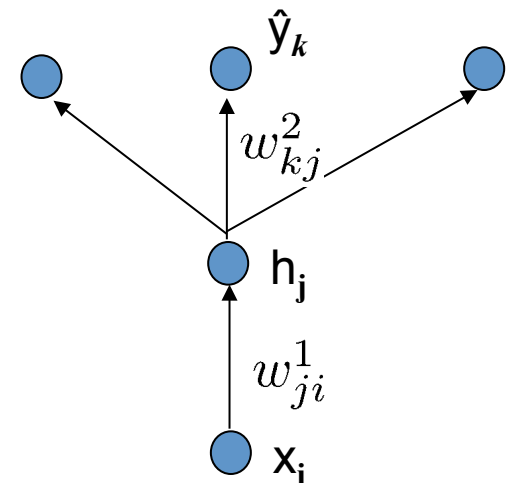
$$J_i(W) = \sum_k (y_k^{(i)} - \hat{y}_k^{(i)})^2$$

Output layer

$$\hat{y}_k = \sigma(s_k) = \sigma(\sum_j w_{kj}^2 h_j)$$

Hidden layer

$$h_j = \sigma(t_j) = \sigma(\sum_i w_{ji}^1 x_i)$$



Backpropagation

- Just gradient descent...
- Apply the chain rule to the MLP

$$\frac{\partial J}{\partial w_{kj}^2} = -2(y_k - \hat{y}_k) \sigma'(s_k) h_j$$

$$\frac{\partial J}{\partial w_{ji}^1} = \sum_k \beta_k^2 -2(y_k - \hat{y}_k) \sigma'(s_k) w_{kj}^2 \sigma'(t_j) x_i$$

Forward pass

Loss function

$$J_i(W) = \sum_k (y_k^{(i)} - \hat{y}_k^{(i)})^2$$

Output layer

$$\hat{y}_k = \sigma(s_k) = \sigma(\sum_j w_{kj}^2 h_j)$$

Hidden layer

$$h_j = \sigma(t_j) = \sigma(\sum_i w_{ji}^1 x_i)$$

```
% X : (1xN1)
H = Sig(X1.dot(W[0]))
% W1 : (N2 x N1+1)
% H : (1xN2)
Yh = Sig(H1.dot(W[1]))
% W2 : (N3 x N2+1)
% Yh : (1xN3)
```

```
B2 = (Y-Yhat) * dSig(S) # (1xN3)

G2 = B2.T.dot( H ) # (N3x1) * (1xN2) = (N3xN2)

B1 = B2.dot(W[1])*dSig(T) # (1xN3) . (N3*N2) * (1xN2)

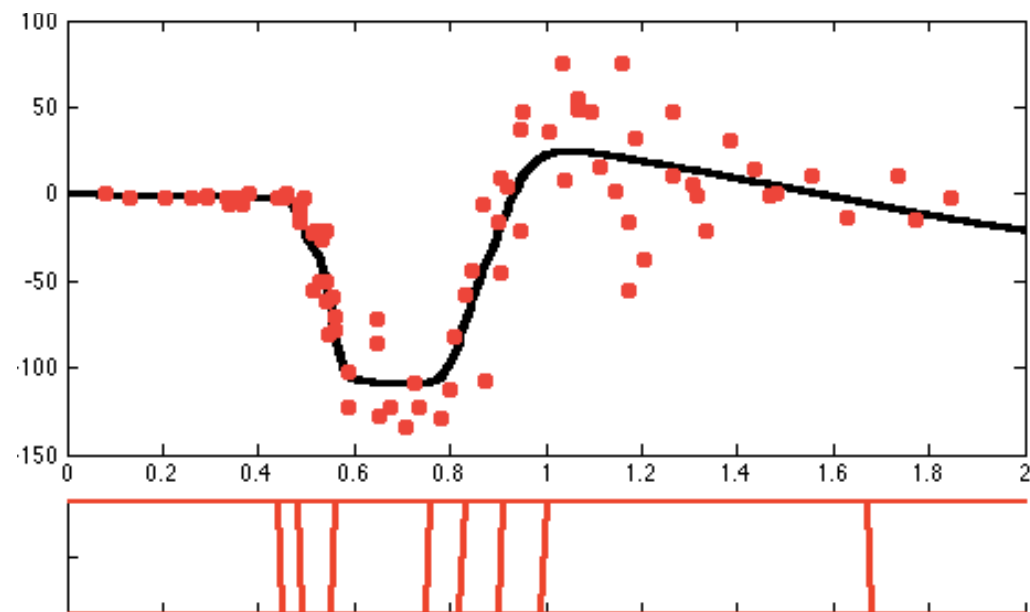
G1 = B1.T.dot( X ) # (N2 x N1+1)
```

Example: Regression, MCycle data

- Train NN model, 2 layer
 - 1 input features => 1 input units
 - 10 hidden units
 - 1 target => 1 output units
 - Logistic sigmoid activation for hidden layer, linear for output layer

Data:
+
learned prediction f' n:

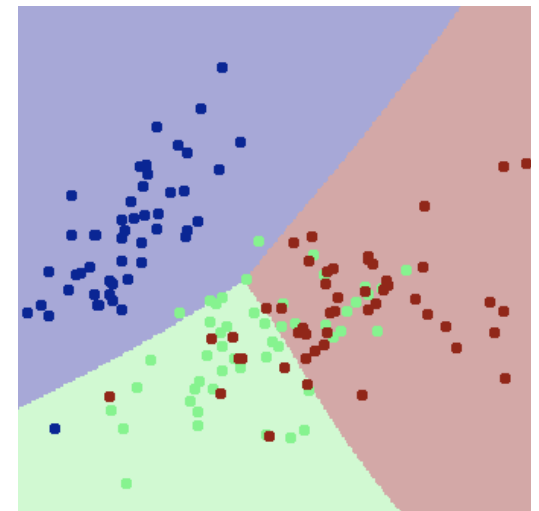
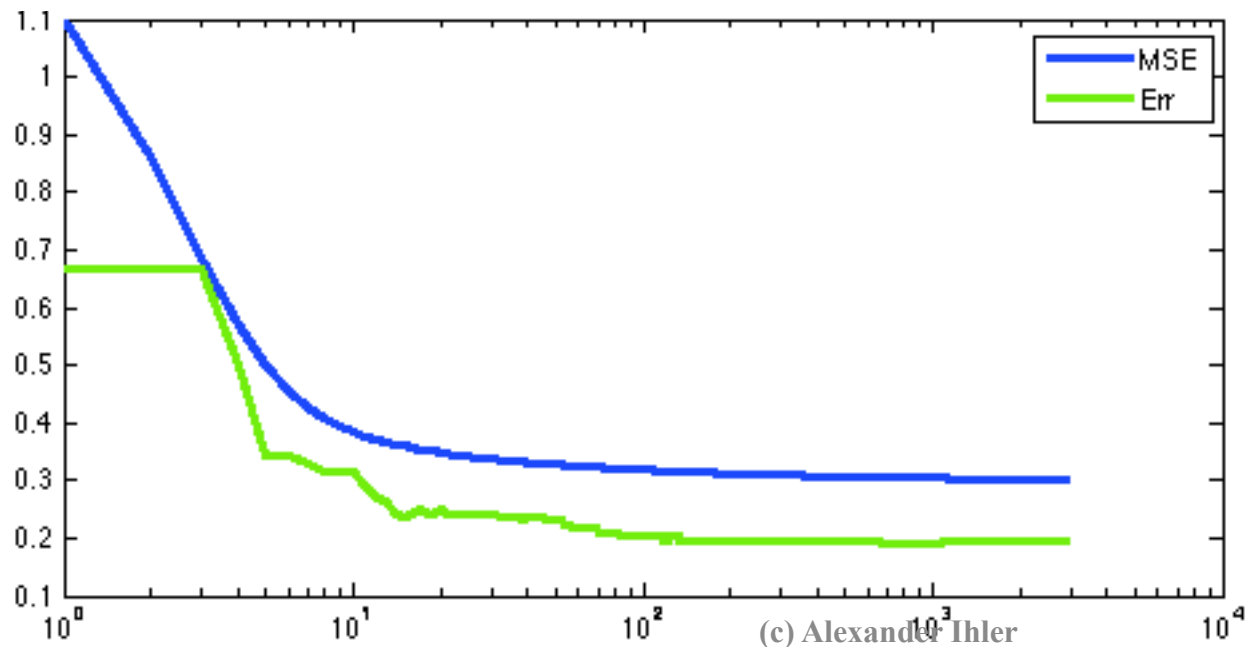
Responses of hidden nodes
(= features of linear regression):
select out useful regions of "x"



(c) Alexander Ihler

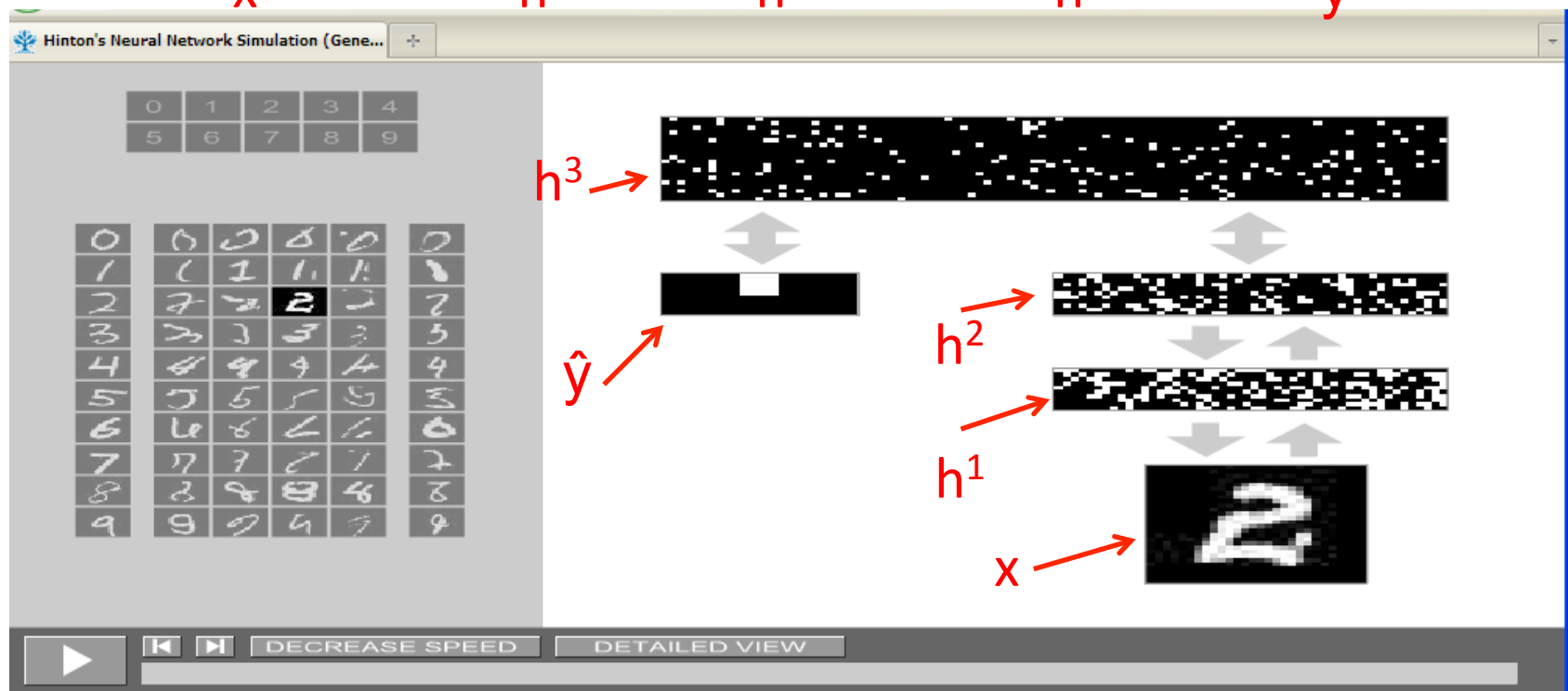
Example: Classification, Iris data

- Train NN model, 2 layer
 - 2 input features => 2 input units
 - 10 hidden units
 - 3 classes => 3 output units ($y = [0 \ 0 \ 1]$, etc.)
 - Logistic sigmoid activation functions
 - Optimize MSE of predictions using stochastic gradient



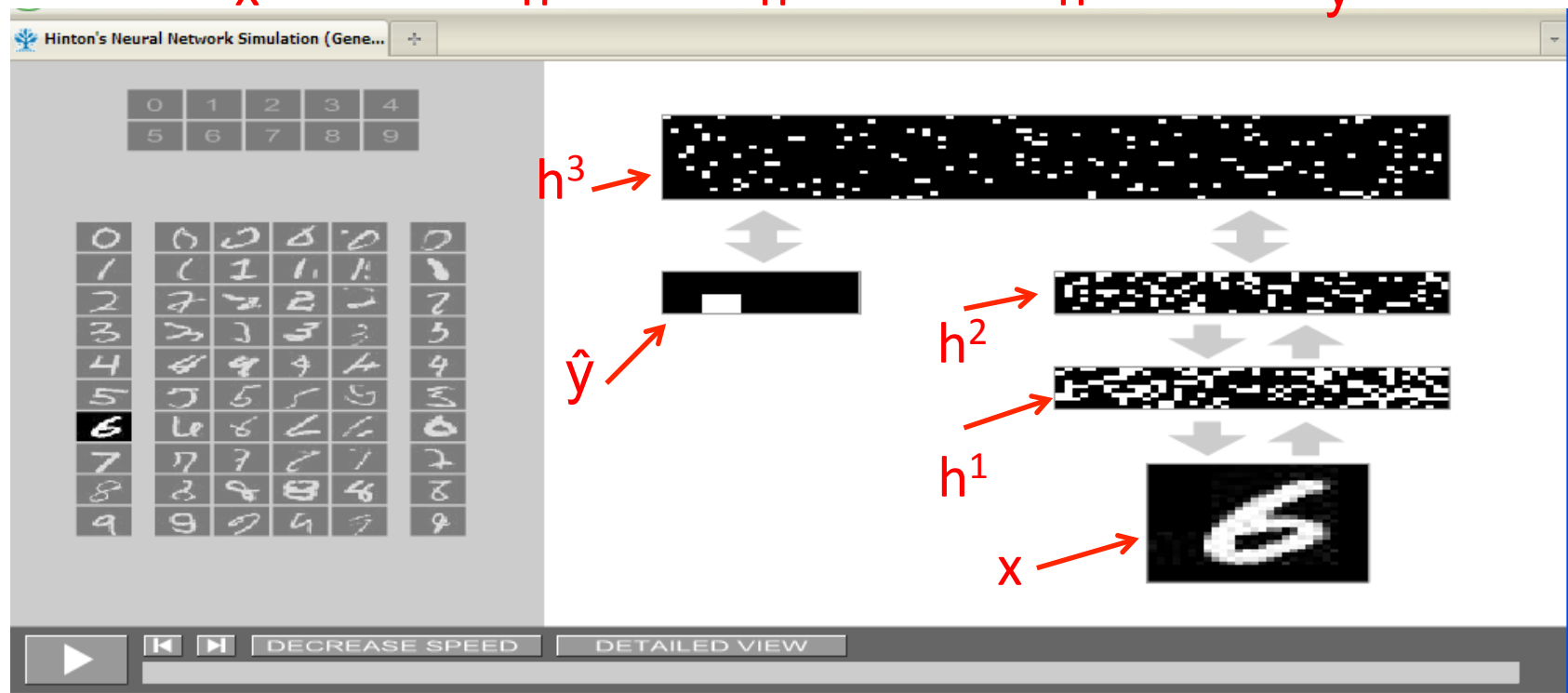
MLPs in practice

- Example: Deep belief nets (Hinton et al. 2007)
 - Handwriting recognition
 - Online demo
 - 784 pixels \Leftrightarrow 500 mid \Leftrightarrow 500 high \Leftrightarrow 2000 top \Leftrightarrow 10 labels



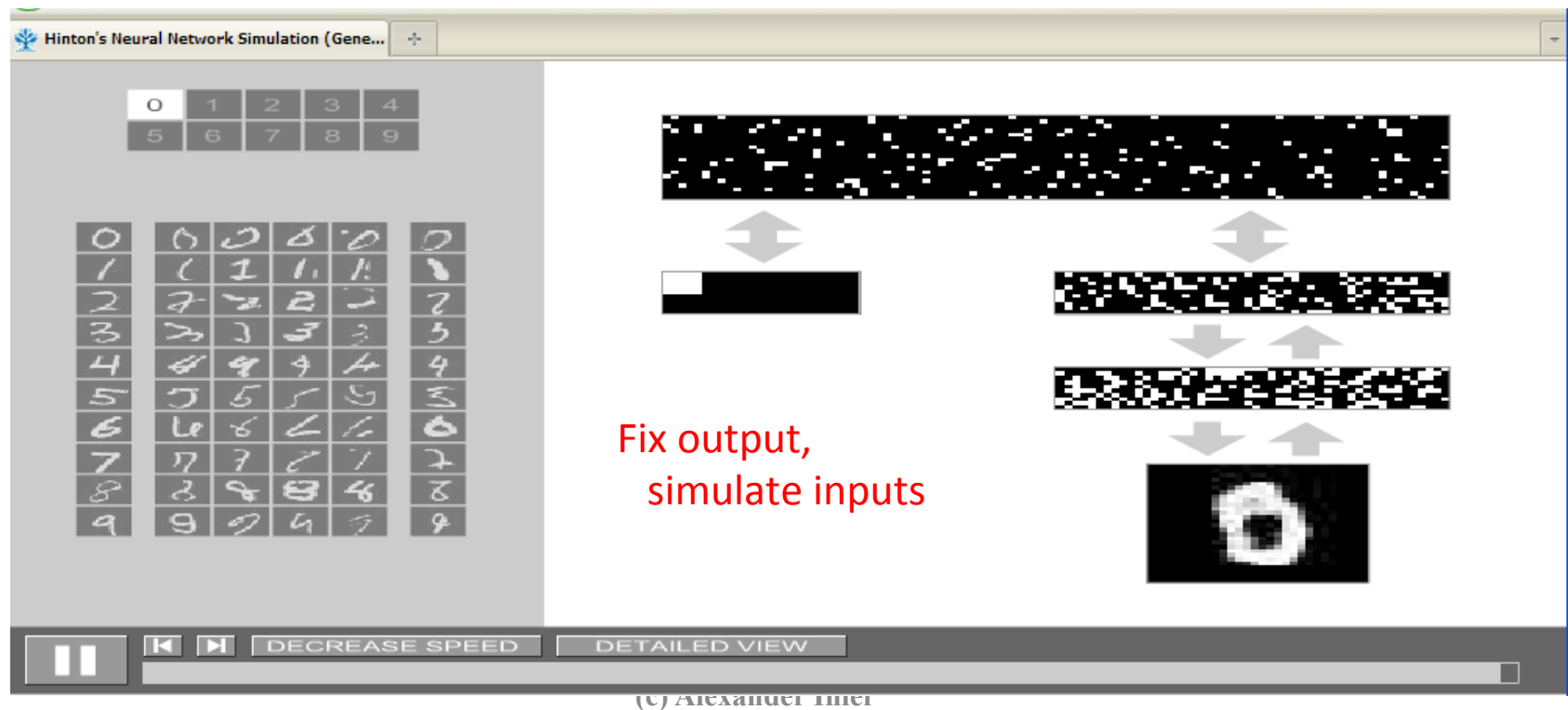
MLPs in practice

- Example: Deep belief nets (Hinton et al. 2007)
 - Handwriting recognition
 - Online demo
 - 784 pixels \Leftrightarrow 500 mid \Leftrightarrow 500 high \Leftrightarrow 2000 top \Leftrightarrow 10 labels



MLPs in practice

- Example: Deep belief nets (Hinton et al. 2007)
 - Handwriting recognition
 - Online demo
 - 784 pixels \Leftrightarrow 500 mid \Leftrightarrow 500 high \Leftrightarrow 2000 top \Leftrightarrow 10 labels



Neural networks & DBNs

- Want to try them out?
- Matlab “Deep Learning Toolbox”
<https://github.com/rasmusbergpalm/DeepLearnToolbox>



Matlab/Octave toolbox for deep learning. Includes Deep Belief Nets, Stacked Autoencoders, Convolutional Neural Nets, Convolutional Autoencoders and vanilla Neural Nets. Each method has examples to get you started.

- PyLearn2
<https://github.com/lisa-lab/pylearn2>
- TensorFlow

Summary

- Neural networks, multi-layer perceptrons
- Cascade of simple perceptrons
 - Each just a linear classifier
 - Hidden units used to create new features
- Together, general function approximators
 - Enough hidden units (features) = any function
 - Can create nonlinear classifiers
 - Also used for function approximation, regression, ...
- Training via backprop
 - Gradient descent; logistic; apply chain rule