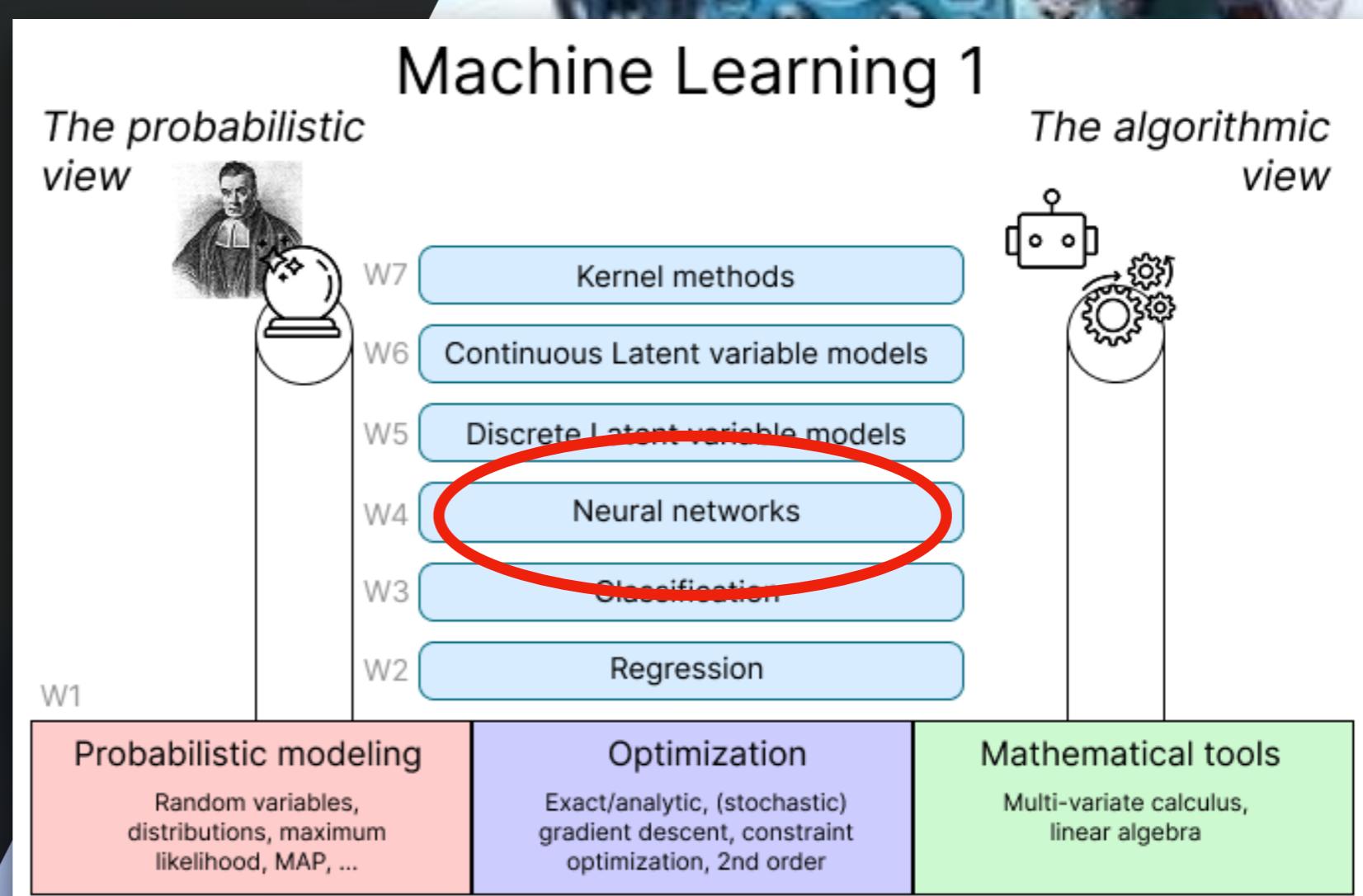


# Machine Learning 1

## Lecture 8 - Neural Networks

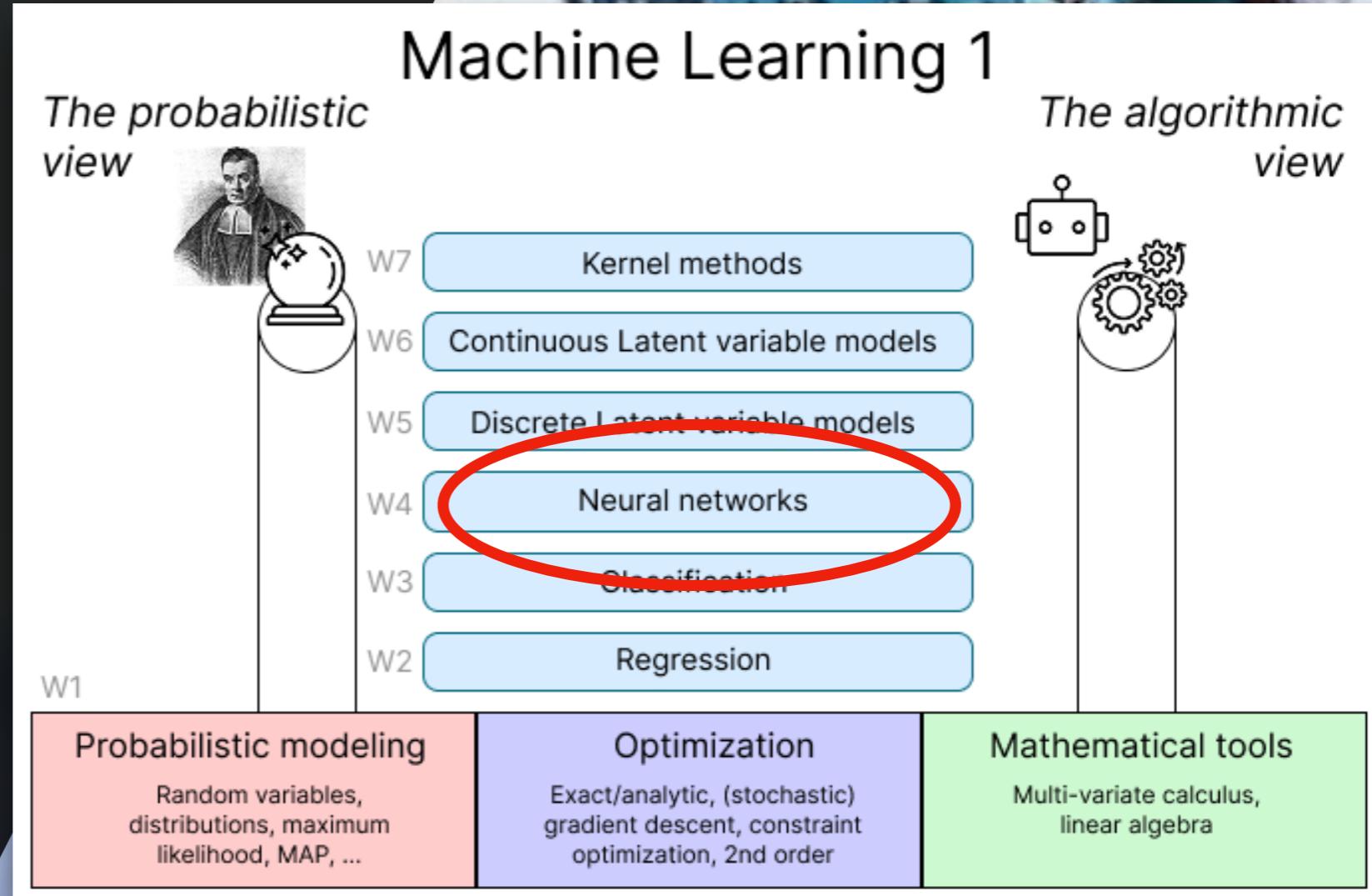
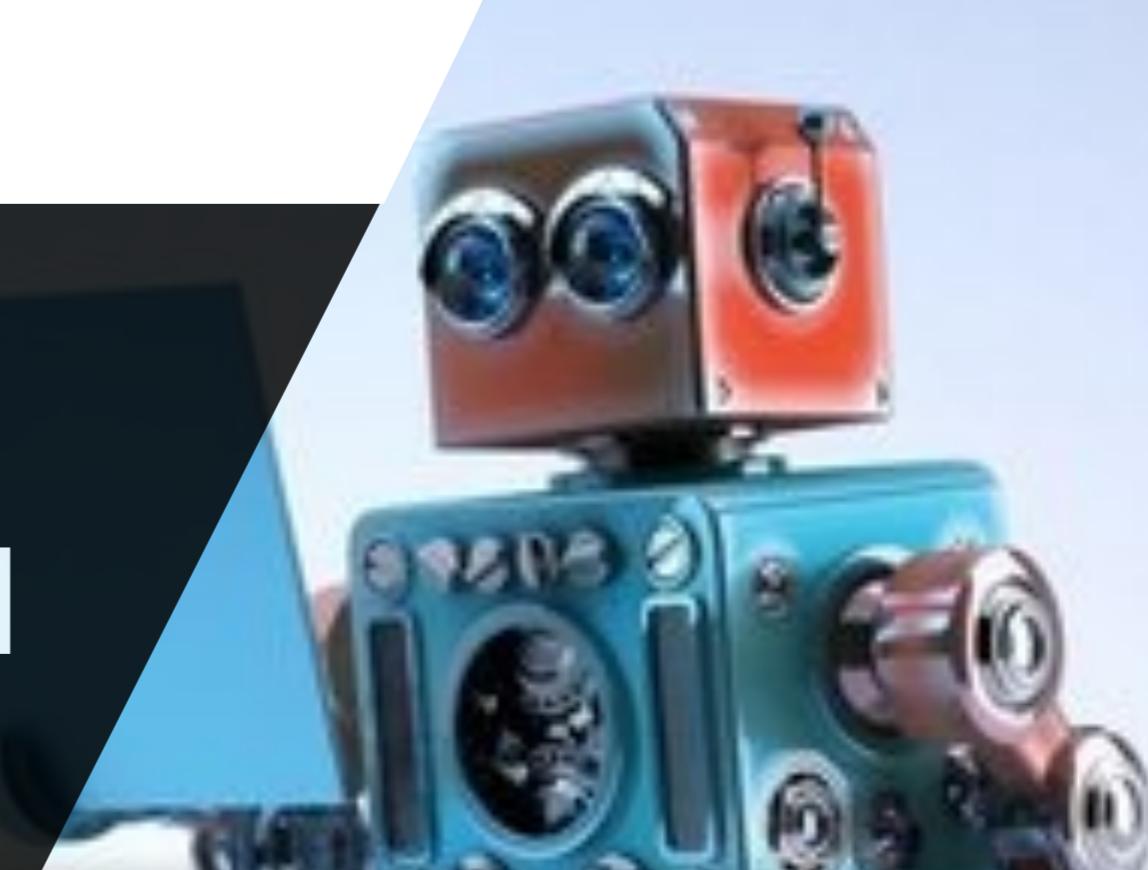
Erik Bekkers



# Machine Learning 1

Lecture 8.1 - Supervised Learning  
Neural Networks

Erik Bekkers  
(Bishop 5.1)











Why activation funct,  
are needed?

$$h\left(\underbrace{w^{(3)T} h(w^{(2)T} h(w^{(1)T} \underline{x}))}_{M'' \leftarrow M' \leftarrow M \leftarrow D}\right)$$

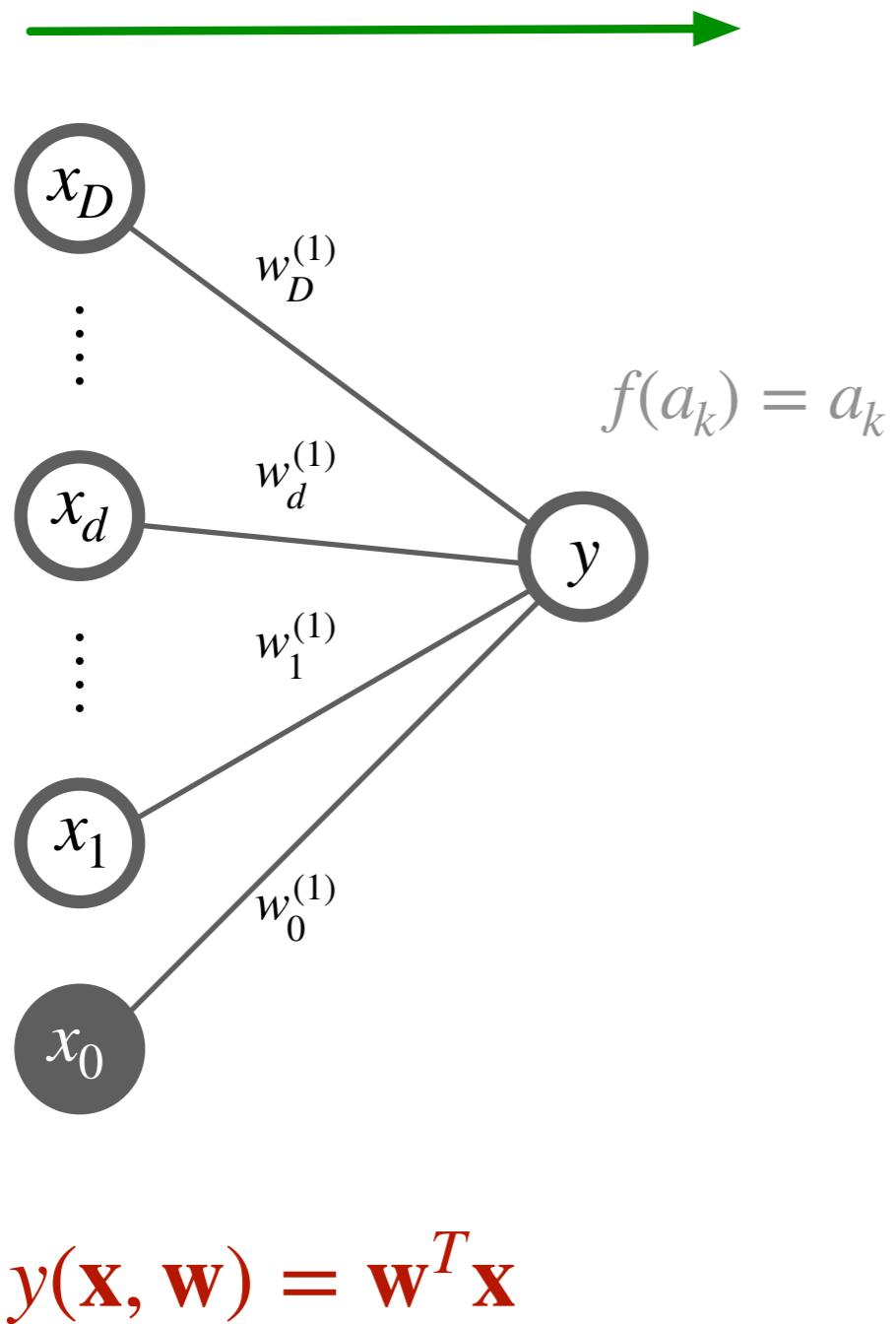
$$M'' \leftarrow M' \leftarrow M \leftarrow D$$

$$\underbrace{w^{(3,2,1)^T} \underline{x}}_{M'' \times D}$$

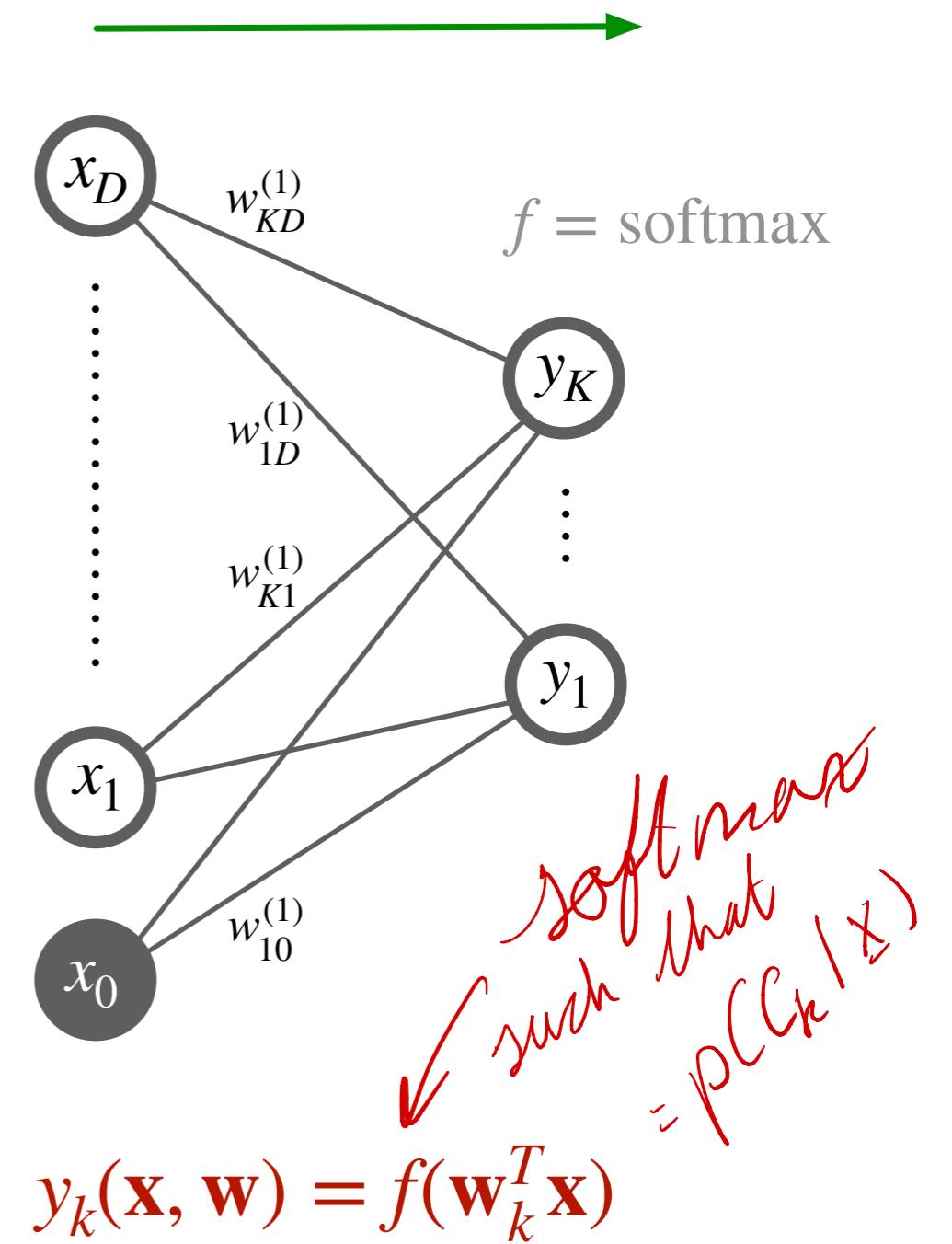
still  
linear

without act\_fn it is  
still effectively one matrix  
vector mult.

# Linear regression & classification as NN



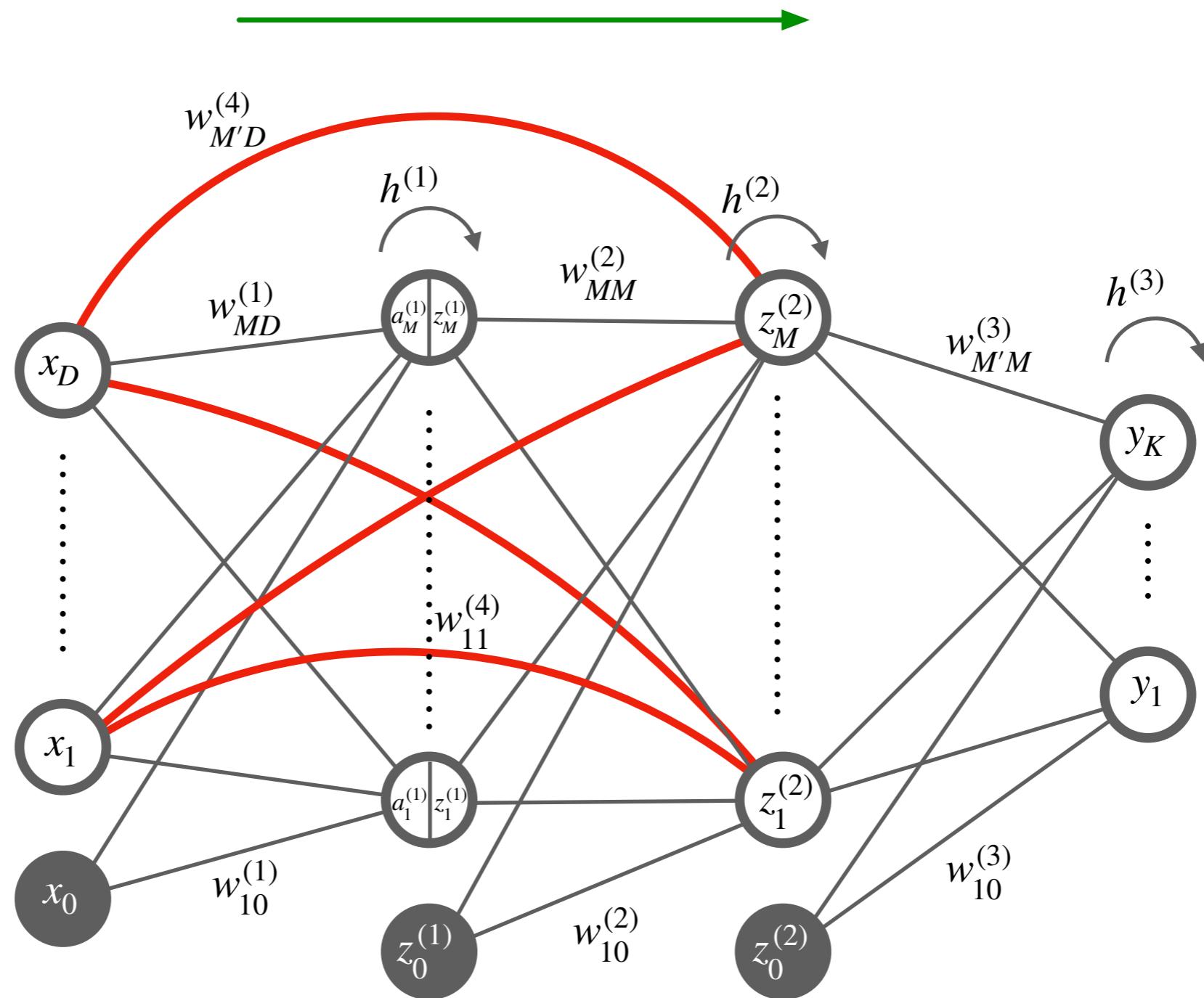
**Figure:** Linear regression as 1-layer NN



**Figure:** Linear Classification with  $K$  classes as 1-layer NN



# Feed-Forward Networks: Skip Connections



**Figure:** 3 layer feed-forward net with skip connections

# Feed-Forward Networks: Sparse Connections

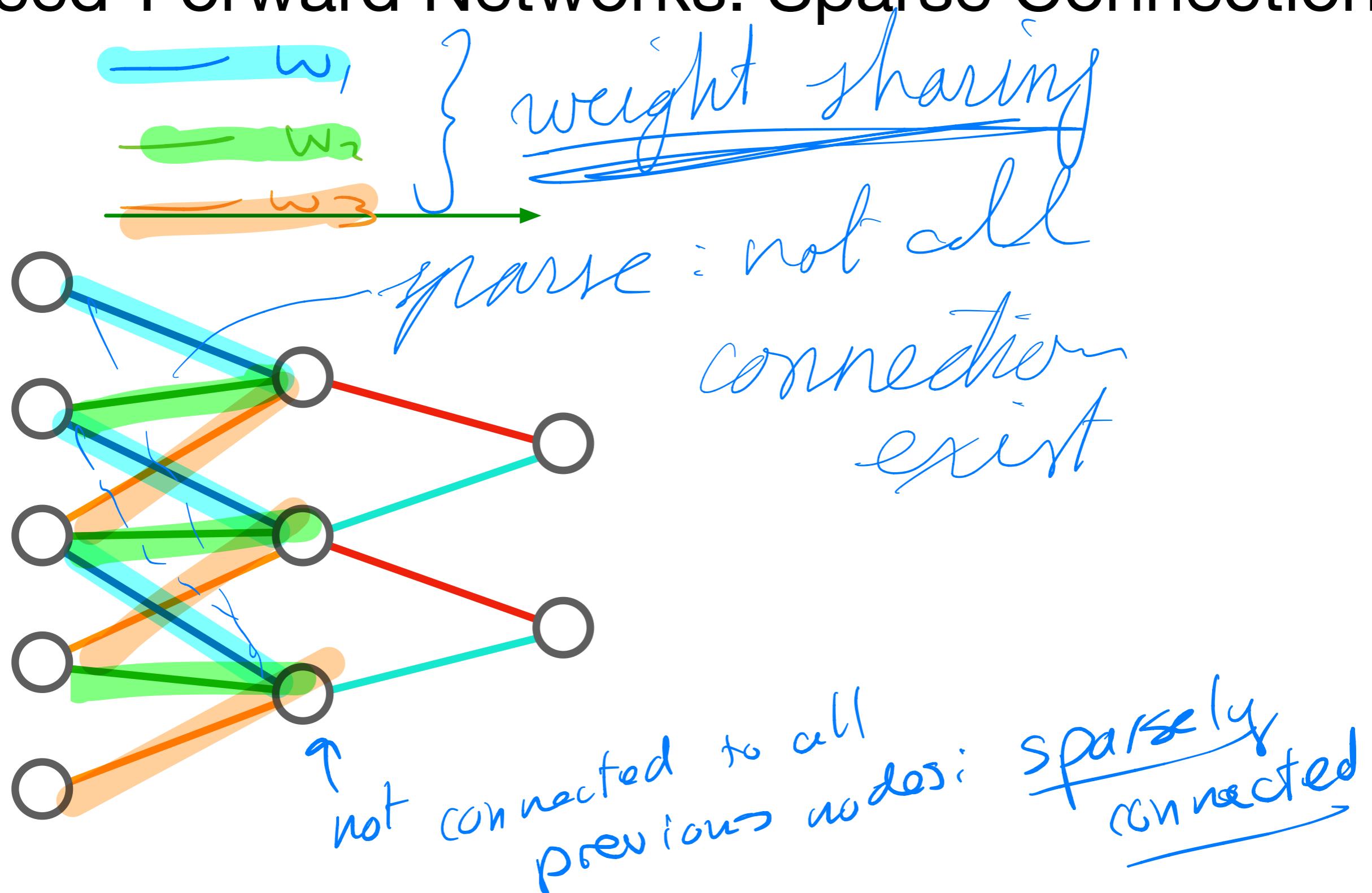


Figure: Feed-forward architecture with sparse connections. With "special" weight sharing --> Convolutional Neural Nets (Le Cun et al 1989)

layer  
is equivariant

# Feed-Forward Networks: Sparse Connections

*Here are  $g, x$  treated as signals*

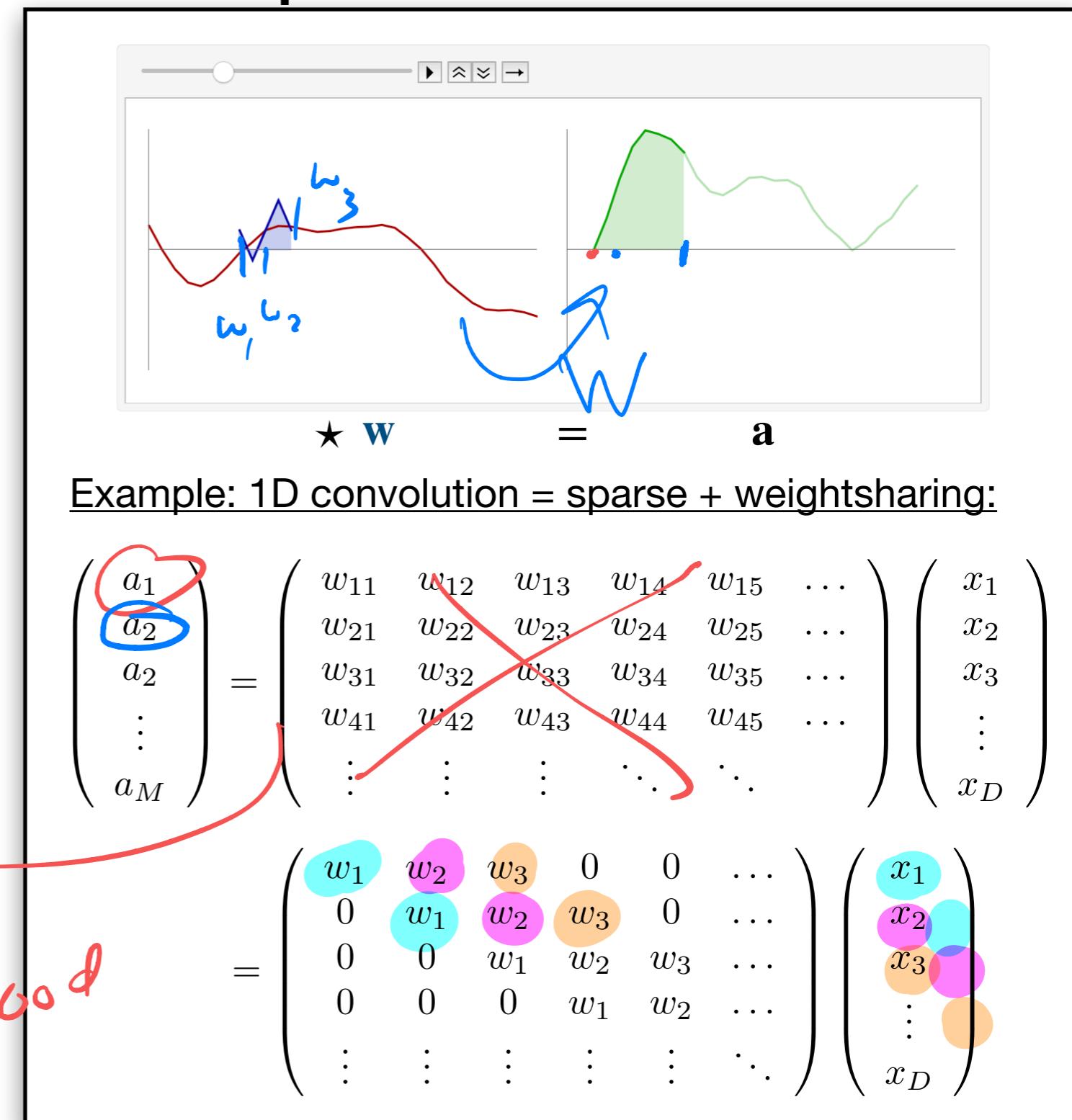
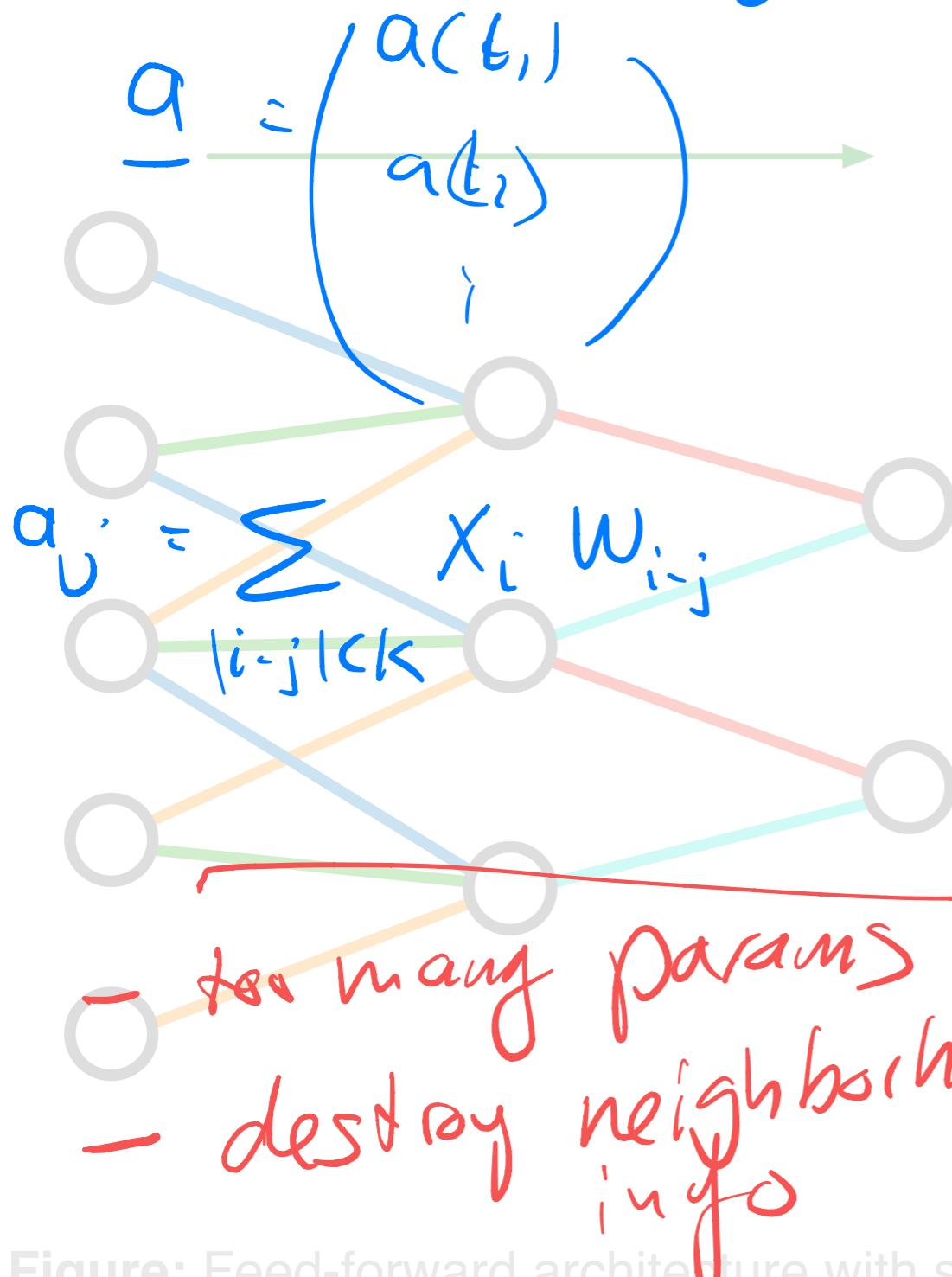


Figure: Feed-forward architecture with sparse connections. With special weight sharing --> Convolutional Neural Nets (Le Cun et al 1989)

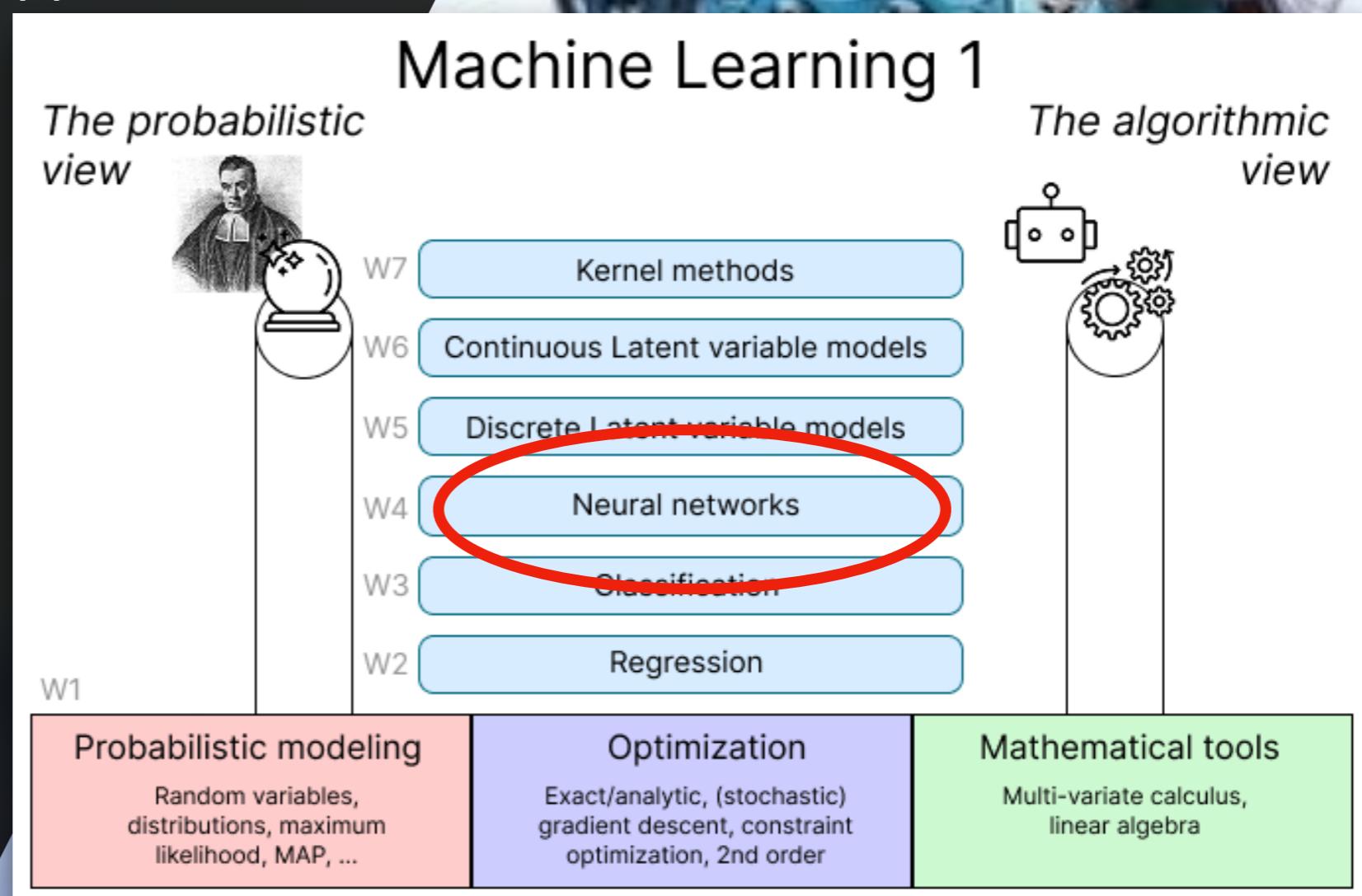


# Machine Learning 1

Lecture 8.2 - Supervised Learning  
Neural Networks - Universal Approximators



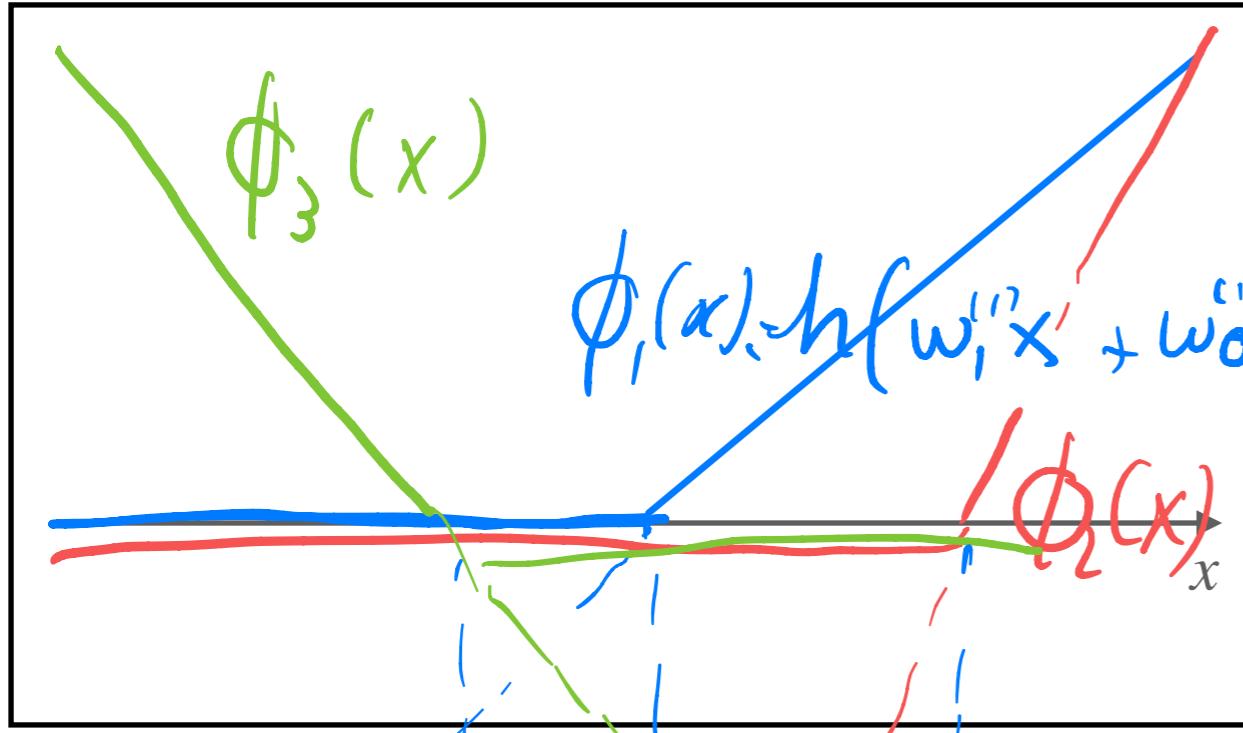
Erik Bekkers  
(Bishop 5.1)



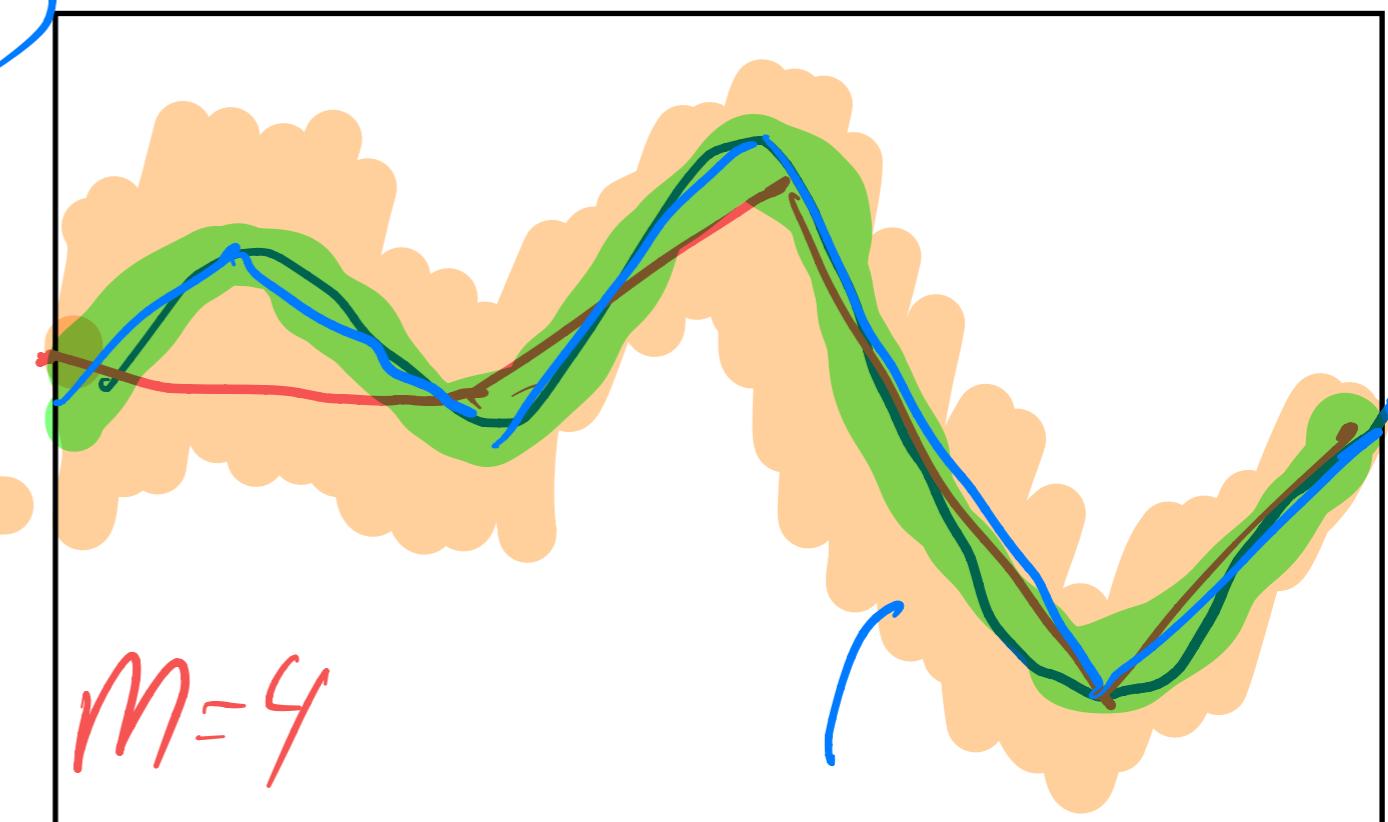


# Neural Networks with ReLU

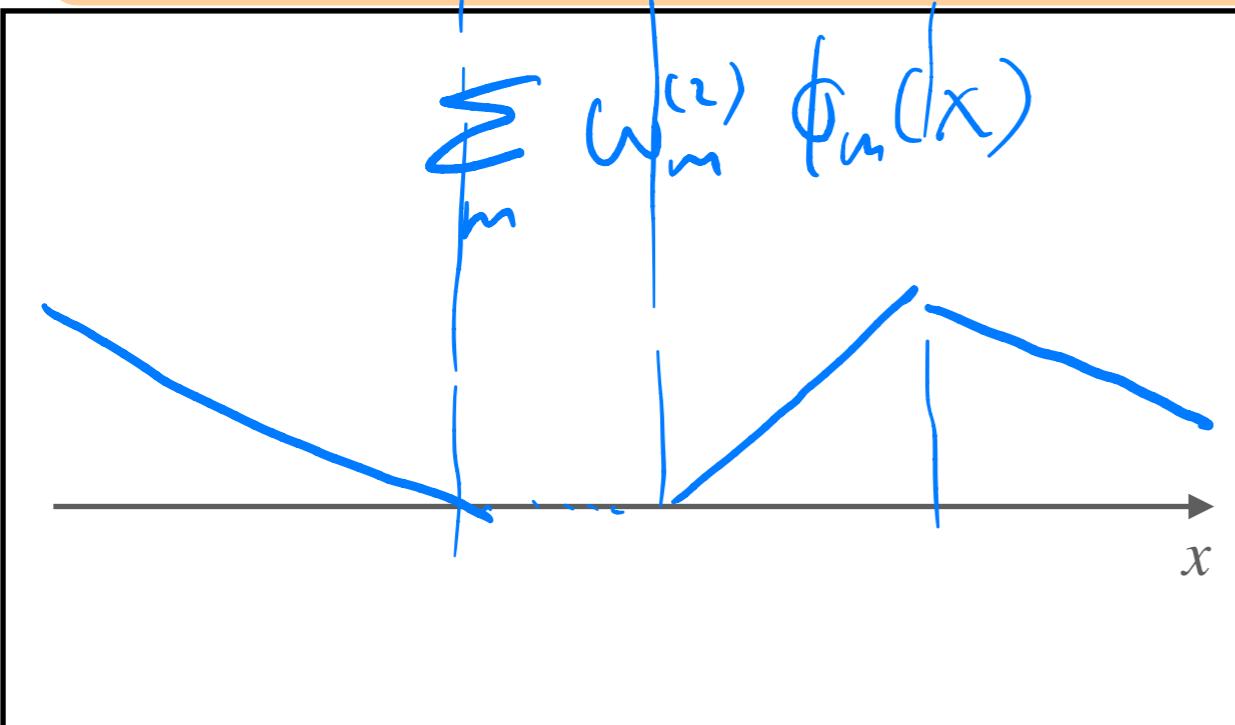
Learned basis functions with ReLU



Approximation with ReLU NNs/PWL functions



Leads to piece-wise linear (PWL) functions



$M=4$   
 $M=5$   
tolerance  
 $\epsilon$

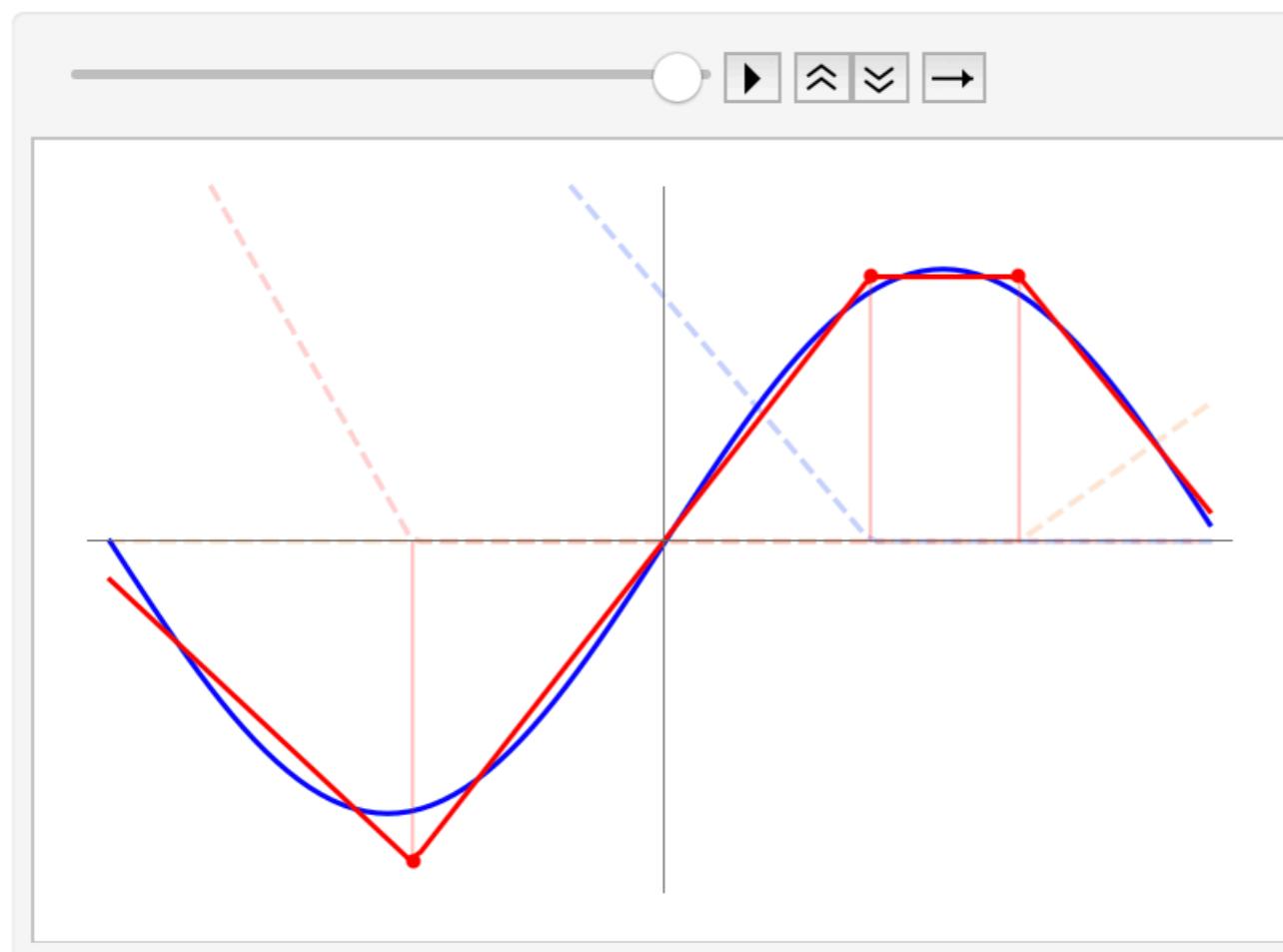




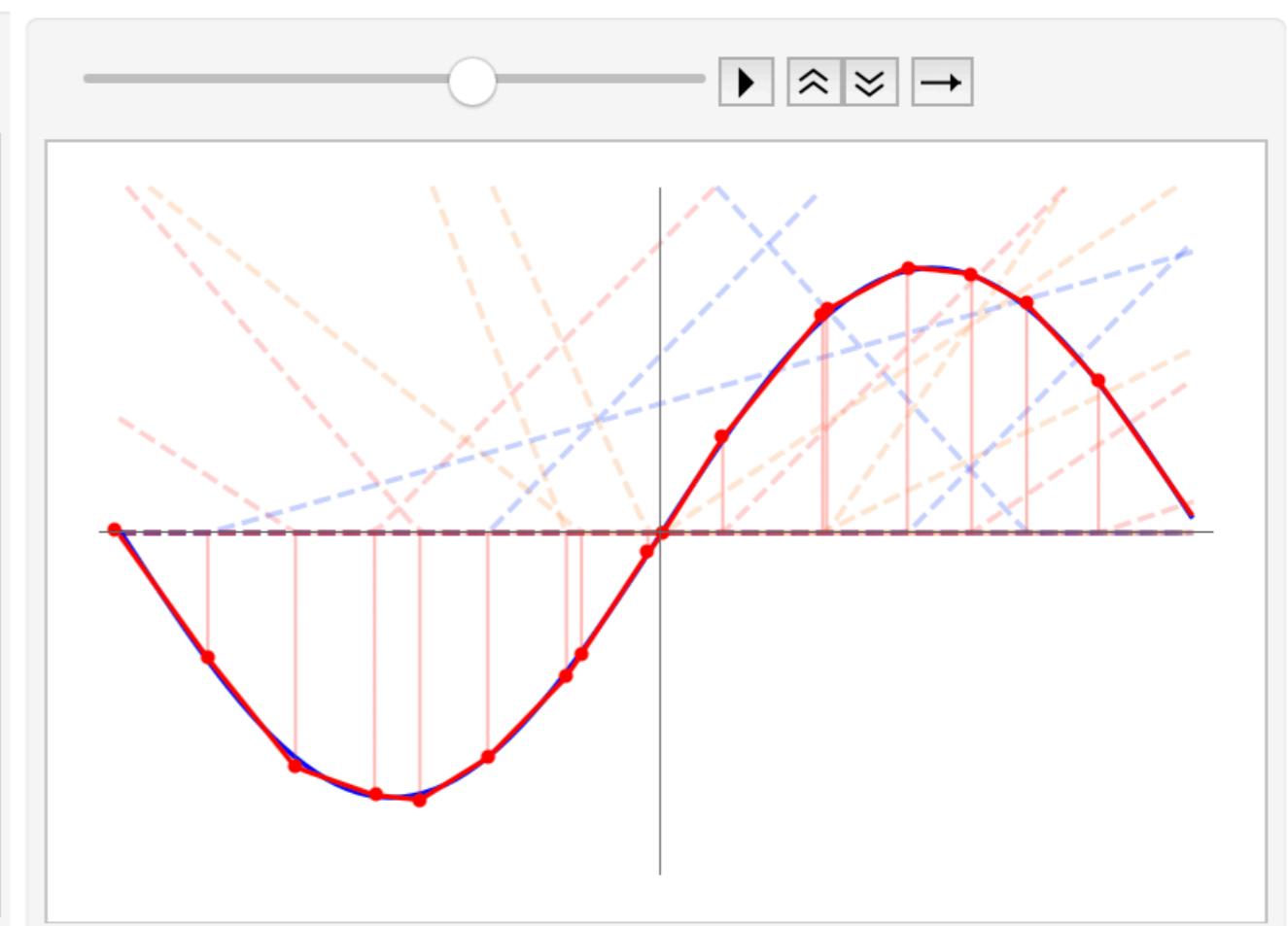


# Example: Function Approximations

Piece-wise linear approximation with 2-layer NN



**M=3 hidden units**

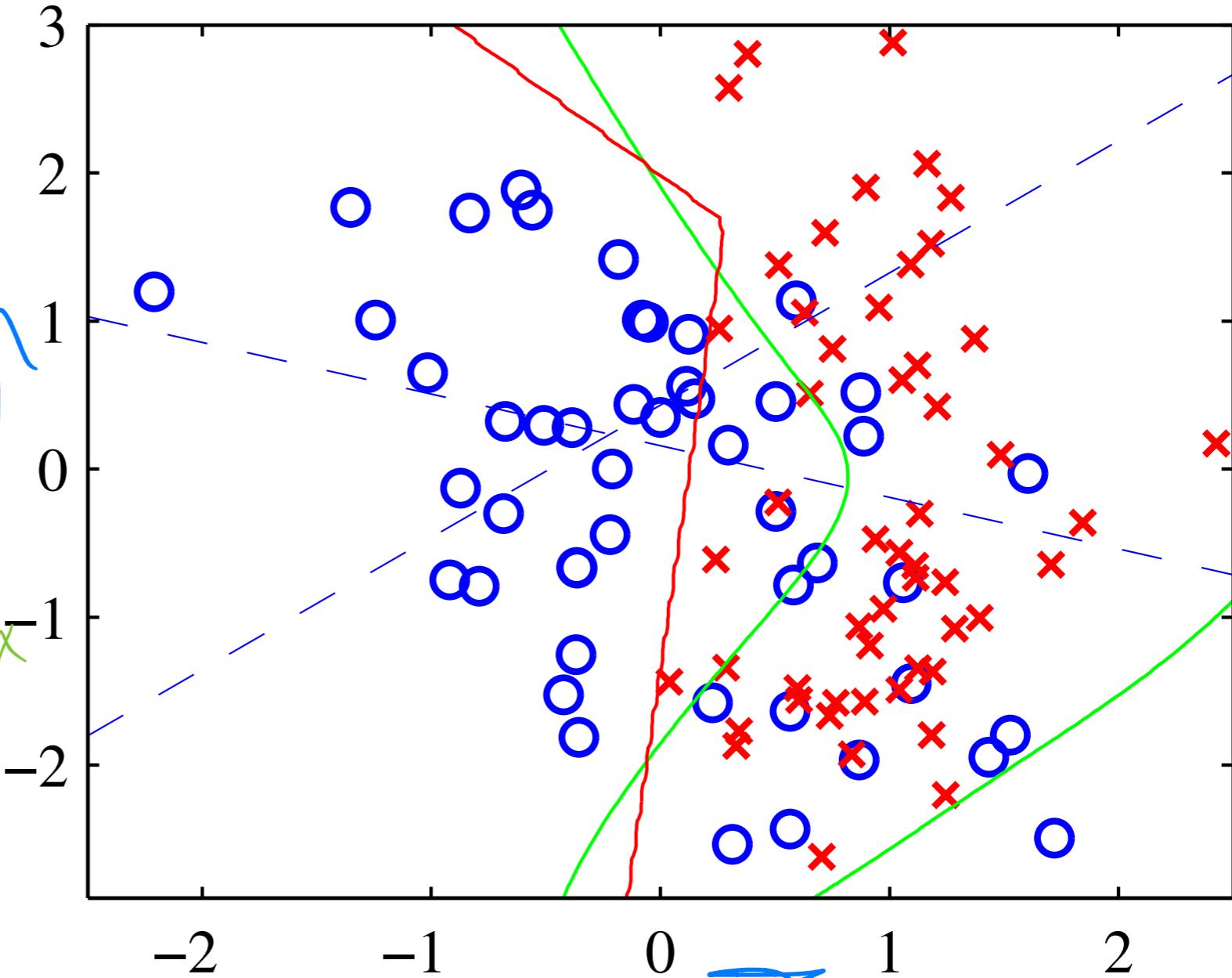


**M=20 hidden units**

# Example: Classification with Neural Nets

MLP:

- ▶ 2 layers
- ▶ # of inputs: 2
- ▶ 2 hidden units with tanh activation function
- ▶ # of outputs: 1 or 2  $x_2$
- ▶ Output activation function:  $\sigma$  softmax
- ▶ Red line: MLP decision boundary
- ▶ Green line: optimal decision boundary from synthetic data distribution

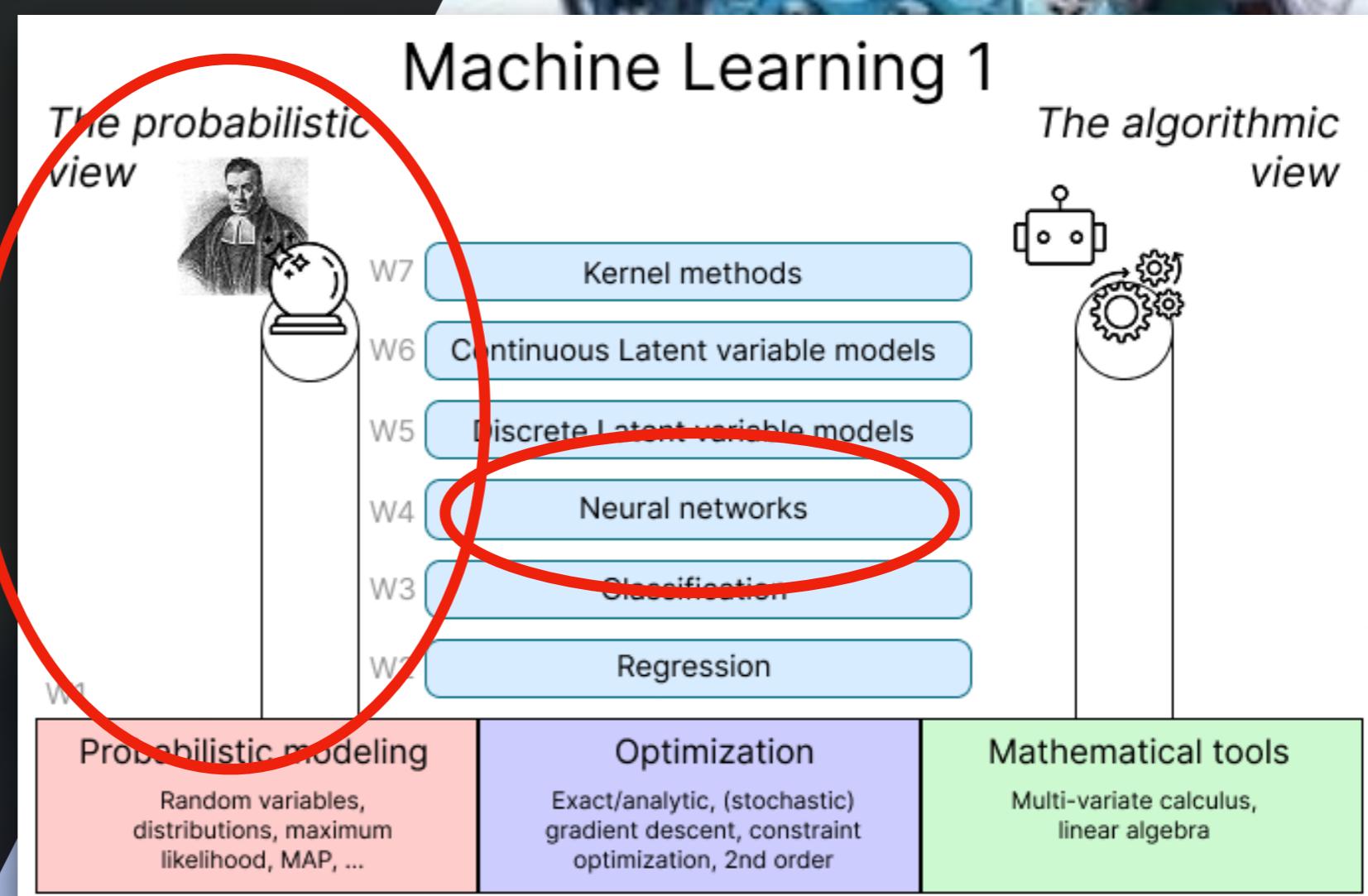
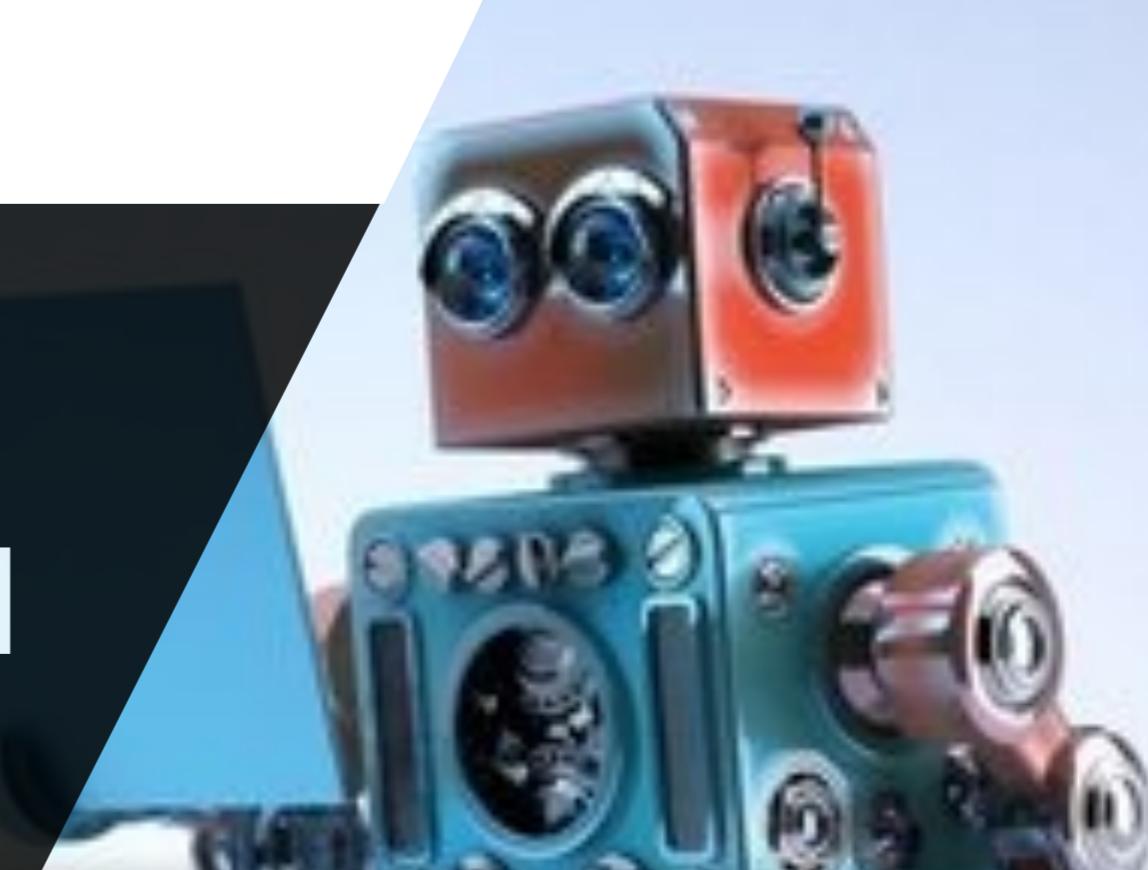


**Figure:** MLP for classification with 2 classes (Bishop5.4)

# Machine Learning 1

Lecture 8.3 - Supervised Learning  
Neural Networks - Losses

Erik Bekkers  
(Bishop 5.2.0)









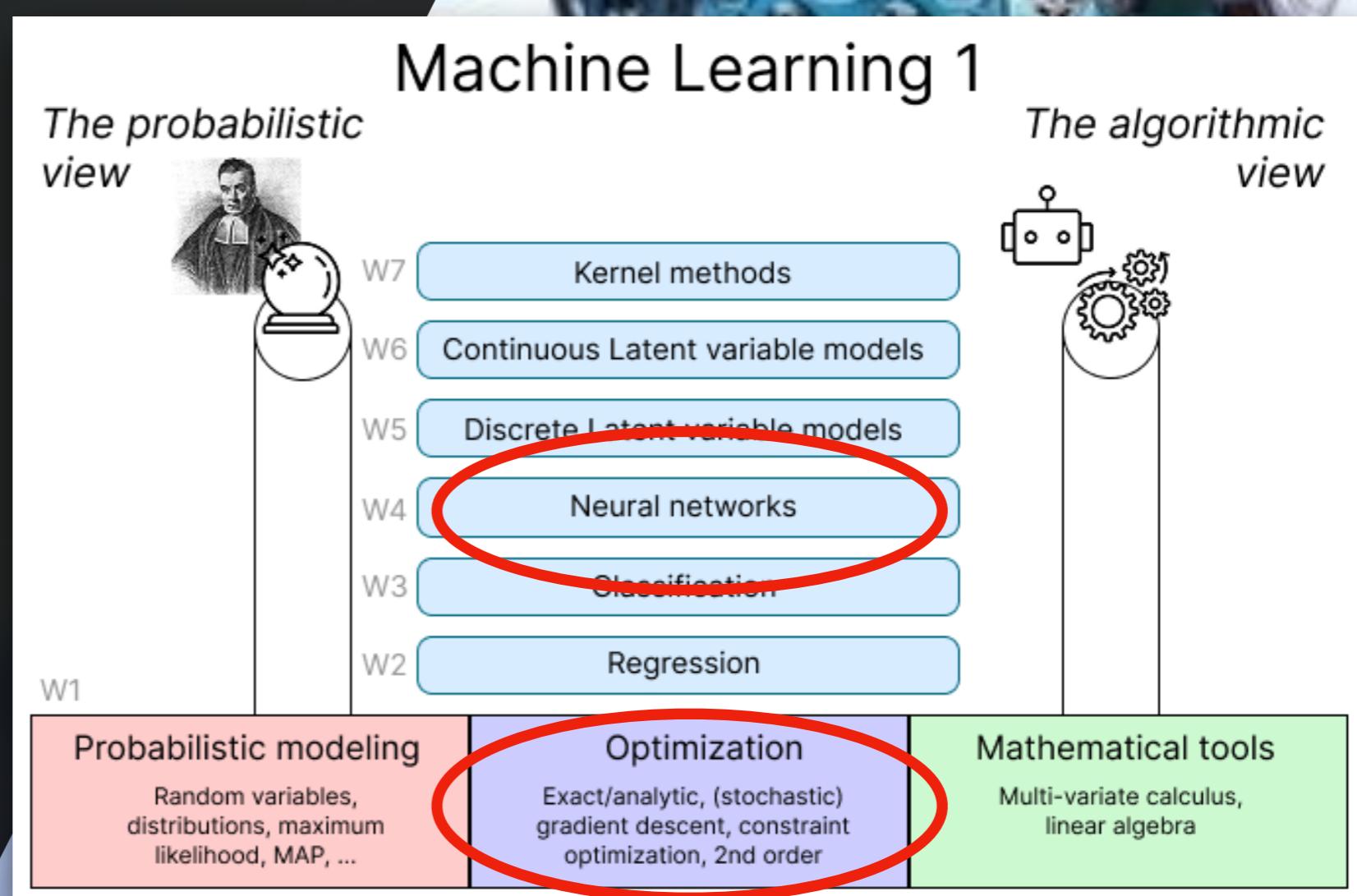




# Machine Learning 1

Lecture 8.4 - Supervised Learning  
Neural Networks - Training

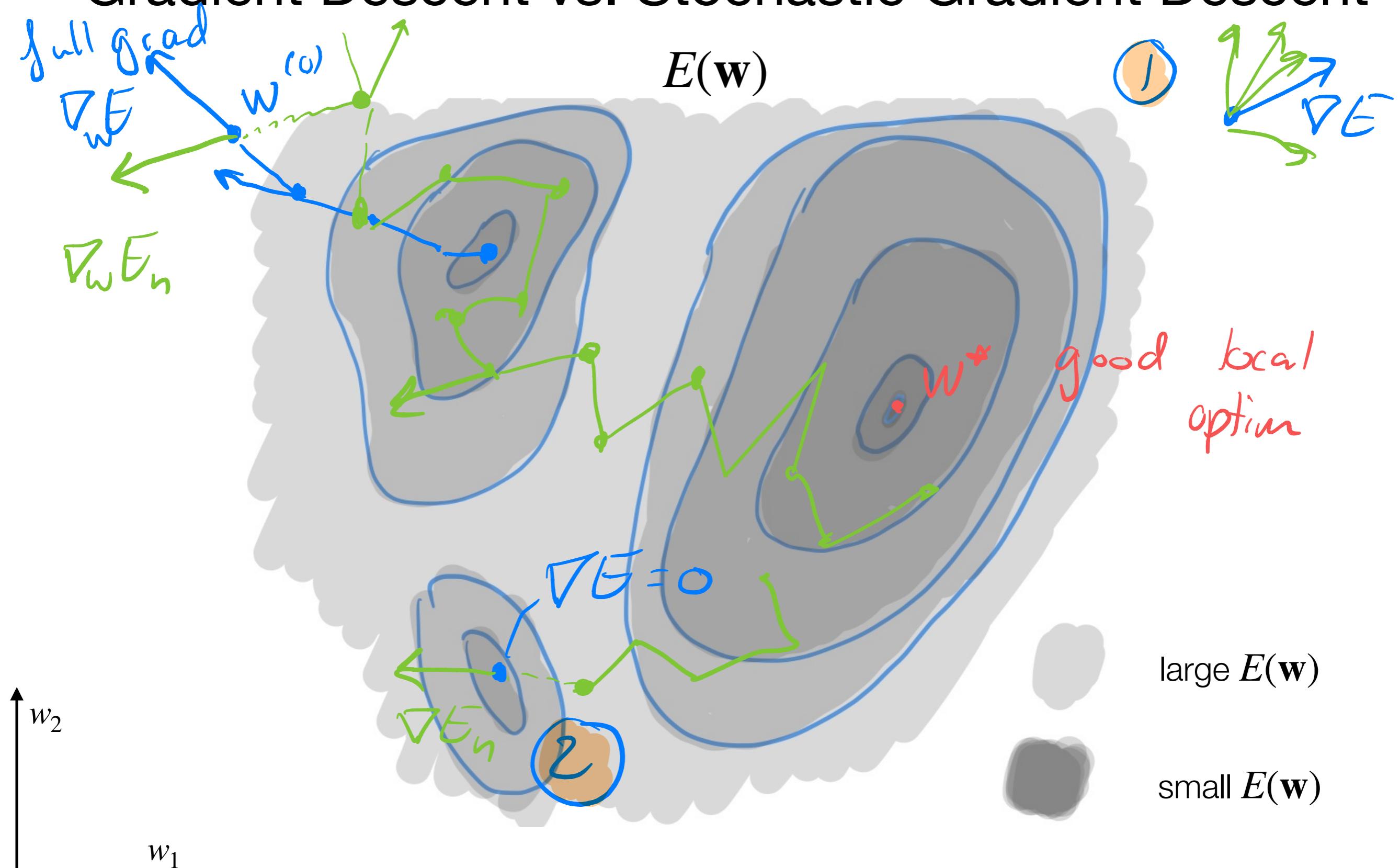
Erik Bekkers  
(Bishop 5.2)







# Gradient Descent vs. Stochastic Gradient Descent



# Gradient Descent vs. Stochastic Gradient Descent

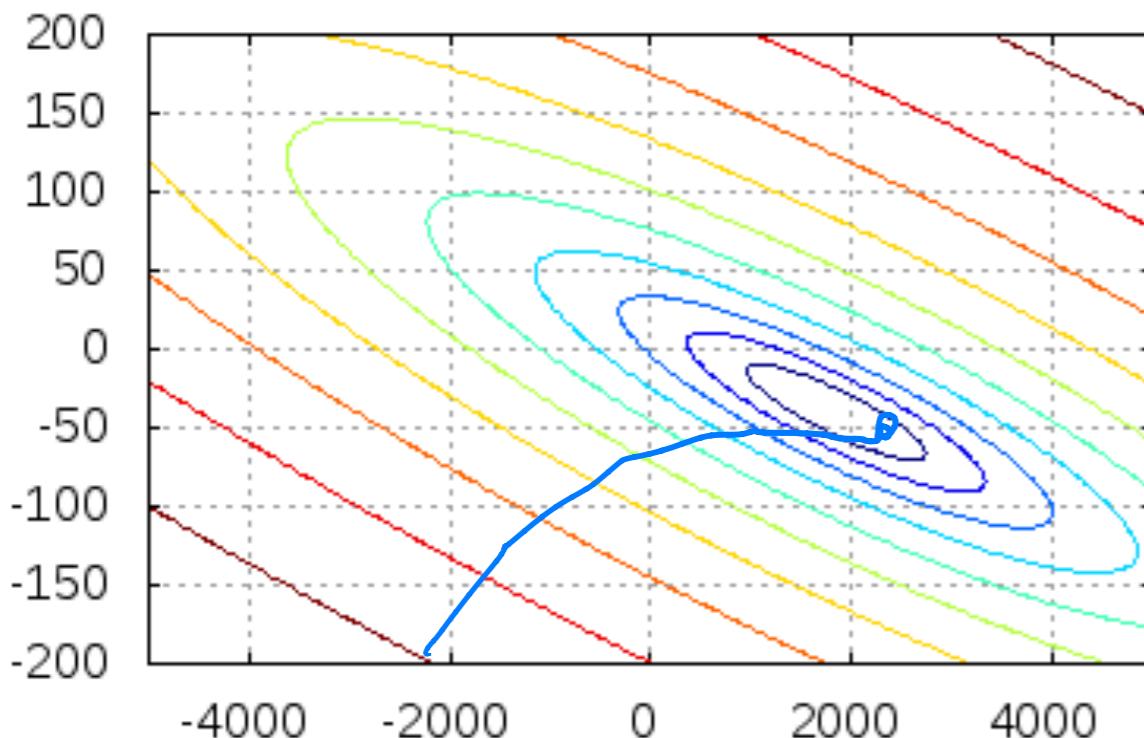
- ▶ About the learning rate:
  - ▶ If learning rate is too small: slow convergence
  - ▶ If learning rate is too large: oscillations around local minimum
  - ▶ Use **learning rate scheduling** with smaller learning rate over time
- ▶ Why SGD over full batch gradient descent
  - ▶ At the beginning of learning all gradients  $\nabla E_n(\mathbf{w})$  will roughly point in the same general direction. **Full batch gradient descent computes redundant number of gradients!** *Efficient + accurate enough*
  - ▶ SGD is more likely to escape a local minimum since  $\nabla E(\mathbf{w}) = \mathbf{0}$  does **not** necessarily imply  $\nabla E_n(\mathbf{w}) = \mathbf{0}$ !

1

2

# Gradient Descent vs. Stochastic Gradient Descent

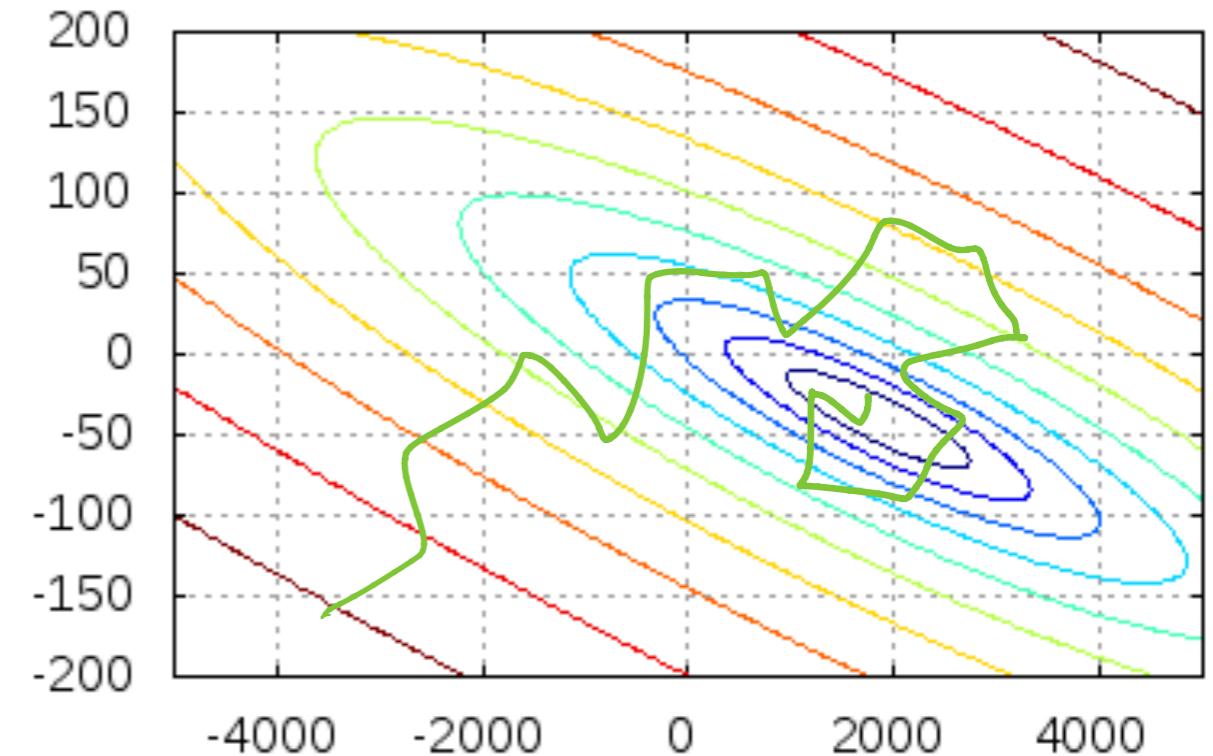
GD



$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E(\mathbf{w}^{(\tau)})$$

- + requires less iterations
- each iteration is slow to compute

SGD

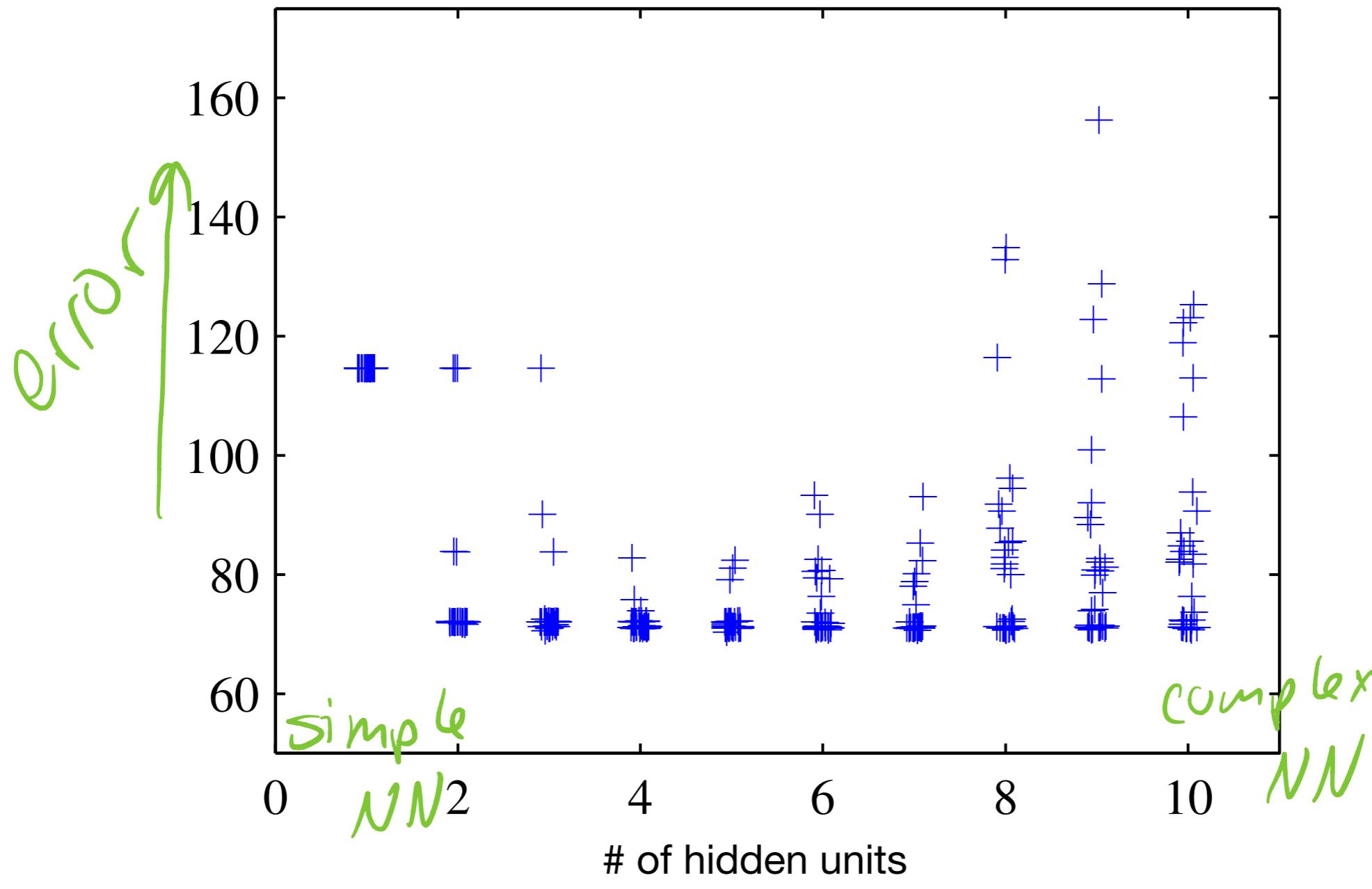


$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}^{(\tau)})$$

- requires more iterations
- + each iteration is fast to compute

# Example: Test Errors and Local Minima

Restart for different random initial  $\mathbf{w}^{(0)}$  to end up in different local minima



**Figure:** sum-of-squares test error vs. network size (# of hidden units) for 30 random starts each (Bishop 5.10)

