



Weight finding in MLPs

- Although it has been known since the 1960's that Multi-Layered Perceptrons are not limited to linearly separable problems there remained a big problem which blocked their development and use
 - How do we find the weights needed to perform a particular function?
- The problem lies in determining an error at the hidden nodes
 - We have no desired value at the hidden nodes with which to compare their actual output and determine an error
 - We have a desired output which can deliver an error at the output nodes but how should this error be divided up amongst the hidden nodes?

MLP Learning Rule • In 1986 Rumelhart, Hinton and Williams proposed a Generalised Delta Rule • Also known as Error Back-Propagation or Gradient Descent Learning • This rule, as its name implies, is an extension of the good old Delta Rule $\Delta^p W_{ji} = \eta \delta^p_j O^p_i$ • The extension appears in the way we determine the δ values • For an output node we have - $\delta^p_j = (T^p_j - o^p_j) \cdot f'_j (net^p_j)$ • For a hidden node we have - $\delta^p_j = f'_j (net^p_j) \cdot \sum_k \delta^p_k w_{kj}$





Architectures How many hidden layers? How many nodes per hidden layer? There are no simple answers Kolmogorov's Mapping Neural Network Existence Theorem Due to Hecht-Nielsen A multi-layered perceptron with n inputs in 10,11 and m output nodes requires only 1 hidden layer of 2(n+1) nodes This is a theoretical result and, in practice, training times can be very long for such minimalist networks



Initial Weights • What size should they be? - No hard and fast rules - Since the common activation functions produce outputs whose magnitude doesn't exceed 1 a range of between -1 and +1 seems sensible - Some researchers believe values related to the fan-in of a node can improve performance and suggest magnitudes of around *1/sqrt(fan-in)* • Never use symmetric weight values - Symmetric patterns in the weights, once manifested, can be difficult to get rid of • So, use values between -1 and +1 and make sure there are no patterns in the weights

Problems with Gradient Descent

- The problems associated with gradient descent learning are the inverse of those present in classical hill-climbing search
- Local Minima
 - Getting stuck in a local minimum instead of reaching a global minimum
 - Detectable because weights don't change but the error remains unacceptable
- Plateaux
 - Moving around aimlessly because the error surface is flat
 - Detectable because although the weights keep changing the error doesn't
- Crevasses
 - Getting caught in a downwards spiral which doesn't lead to a global minimum
 - NOT detectable so dangerous but rare





Some More Problems

- Training with too high a learning rate can take longer or even fail
 - As a general rule the larger the learning rate, η, the faster the training. The weights are adjusted by larger amounts and so migrate towards a solution more rapidly
 - If the weight changes are too large though the training algorithm can keep "stepping over" the values needed for a solution rather than landing on them
- Networks with too many weights will not generalise well
 - The more weights there are in a network (the more degrees of freedom it has) the more arbitrary is the weight set discovered during training
 - One weight set chosen arbitrarily from many possible solutions that satisfy the requirements of the training set, is unlikely to satisfy data not used in training

Input Representations (I)

- The way in which the inputs to an ANN are represented can be crucial to the successful training and eventual performance of the system
- There is **no correct way** to select input representations since they are highly dependent on what the ANN is required to learn about the inputs
- A significant proportion of the design time for an ANN is spent on devising the input encoding scheme
- Consider the problem of representing some simple shapes such as *triangle*, *square*, *pentagon*, *hexagon* and *circle*
 - Possible schemes include
 - Bitmap images
 - Edge counts
 - · Shape-specific input nodes

Input Representations (II)

What are we seeking to do?

- Do we need to generalise about shapes?
 - If not then <u>shape-specific</u> input nodes should suffice because we won't need any more detailed information about the shapes
- Generalising about regular shapes
 - If we only need to be able to differentiate between and generalise about regular shapes then an <u>edge count</u> should suffice
- Generalising about irregular shapes
 - If we need to be able to differentiate between and generalise about irregular shapes then a <u>bitmap image</u> may be needed
 NB Angle sizes and edge lengths may suffice for
 - differentiating between different types of triangle or between squares, rectangles rhombuses, etc.
- Greater power => More refined data



Input Representations (IV)

- Another attribute Colour
 - Looks like a case for specific input nodes for each attribute value - one for each colour in this instance
 - All colours have a wavelength though so we might consider normalising the wavelengths and using a single input node to represent the wavelength
 - On the other hand, we know all colours can be generated from the three primaries, so we might use an encoding scheme with one input node for each colour but which allows a whole gamut of colours to be represented by treating the inputs like the colour guns in a television monitor
- Yet another attribute Number
 - Normalise values to avoid saturation
 - Quantise to use multiple discrete inputs
 - Don't employ clever encoding schemes



MLP Examples

• NETtalk

- Speech synthesis
- Sejnowski & Rosenberg (1987)

• ALVINN

- Steering a car along a road
- Pomerleau, et al. (1989)
- ZIP Codes
 - Recognising handwritten ZIP codes
 - Le Cun, et al. (1989)





• Claimed to be twice as fast an non-ANN rival systems



