Intro to Deep Neural Networks: Multi-layer Perceptrons

CS771: Introduction to Machine Learning Pivush Rai

Limitation of Linear Models

Linear models: Output produced by taking a linear combination of input features



• A basic unit of the form $y = f(w^T x)$ is known as the "Perceptron" (not to be confused with the Perceptron "algorithm", which learns a linear classification model) Although can kernelize to

This can't however learn nonlinear functions or nonlinear decision boundaries

make them nonlinear

CS771: Intro to ML

Neural Networks: Multi-layer Perceptron (MLP)

- An MLP is a network containing several Perceptron units across many layers
- An MLP consists of an input layer, an output layer, and one or more hidden layers



Illustration: Neural Net with Single Hidden Layer \hat{y}_n • Compute K pre-activations for each input \boldsymbol{x}_n Last layer activation function *o* A linear model with learnable weight vec \boldsymbol{w}_k $\boldsymbol{z}_{nk} = \boldsymbol{w}_k^{\mathsf{T}} \boldsymbol{x}_n = \sum_{d=1}^{D} \boldsymbol{w}_{dk} \boldsymbol{x}_{nd} \quad (k = 1, 2, ..., K)$ can be an identify function too (e.g., for regression, $\hat{y}_n = s_n$) or sigmod/softmax/sign etc Apply nonlinear activation on each pre-act $h_{n^{\uparrow}}$ h_{n2} Called a hidden unit $h_{nk} = g(z_{nk})$ (k = 1, 2, ..., K)Hidden layer • Apply a linear model with h_n acting as features activation function gmust be nonlinear A linear model with Z_{n2} Z_{n1} learnable weight vec \boldsymbol{v} Score of $\boldsymbol{s}_n = \boldsymbol{v}^{\mathsf{T}} \boldsymbol{h}_n = \sum_{k=1}^{K} \boldsymbol{v}_k \boldsymbol{h}_{nk}$ the input w_{11} Finally, output is produced as w_{12} W_{32} w_{21} w_{22} w_{31} $\hat{y}_n = \mathbf{o}(s_n)$ Score converted to the actual prediction • Loss: $\mathcal{L}(\boldsymbol{W}, \boldsymbol{v}) = \sum_{n=1}^{N} \ell(y_n, \hat{y}_n) \mathcal{X}_n$ x_n ?

Neural Nets: A Compact Illustration

Will denote a linear combination of inputs followed by a nonlinear operation on the result

Note: Hidden layer pre-act z_{nk} and post-act h_{nk} will be shown together for brevity



• Denoting $W = [w_1, w_2, ..., w_K]$, $w_k \in \mathbb{R}^D$, $h_n = g(W^T x_n) \in \mathbb{R}^K$ (K = 2, D = 3 above). Note: g applied elementwise on pre-activation vector $z_n = W^T x_n$ cs771: Intro to ML

Activation Functions: Some Common Choices





Superposition of two linear models = Nonlinear model



CS771: Intro to ML

Examples of some basic NN/MLP architectures



Single Hidden Layer and Single Output

One hidden layer with K nodes and a single output (e.g., scalar-valued regression or binary classification)



CS771: Intro to ML

10

Single Hidden Layer and Multiple Outputs

• One hidden layer with K nodes and a vector of C outputs (e.g., vector-valued regression or multi-class classification or multi-label classification)



CS771: Intro to ML

Multiple Hidden Layers (One/Multiple Outputs)

Most general case: Multiple hidden layers with (with same or different number of hidden nodes in each) and a scalar or vector-valued output





12

The Bias Term

• Each layer's pre-activations $\boldsymbol{z}_n^{(\ell)}$ have a an add bias term $\boldsymbol{b}^{(\ell)}$ (has the same size as $\boldsymbol{z}_n^{(\ell)}$ and $\boldsymbol{h}_n^{(\ell)}$) as well

$$\boldsymbol{z}_n^{(\ell)} = \boldsymbol{W}^{(\ell)^{\mathsf{T}}} \boldsymbol{h}_n^{(\ell-1)} + \boldsymbol{b}^{(\ell)}$$
$$\boldsymbol{h}_n^{(\ell)} = g(\boldsymbol{z}_n^{(\ell)})$$

- Bias term increases the expressiveness of the network and ensures that we have nonzero activations/preactivations even if this layer's input is a vector of all zeros
- Note that the bias term is the same for all inputs (does not depend on n)
- \blacksquare The bias term $\boldsymbol{b}^{(\ell)}$ is also learnable





Neural Nets are Feature Learners

Hidden layers can be seen as <u>learning</u> a feature rep. $\phi(x_n)$ for each input x_n



14

Kernel Methods vs Neural Nets

Also note that neural nets are faster than kernel methods at test time since kernel methods need to store the training examples at test time whereas neural nets do not



Recall the prediction rule for a kernel method (e.g., kernel SVM)

$$y_n = \mathbf{w}^{\mathsf{T}} \phi(\mathbf{x}_n)$$
 OR $y_n = \sum_{i=1}^{m} \alpha_n k(\mathbf{x}_i, \mathbf{x}_n)$

- It's like a one hidden layer NN with
 - Pre-defined N features $\{k(x_i, x_n)\}_{i=1}^N$ acting as feature vector h_n
 - $\{\alpha_i\}_{n=1}^N$ are learnable output layer weights
- It's also like a one hidden layer NN with
 - Pre-defined M features $\phi(x_n)$ (M being size of feature mapping ϕ) acting as feature vector h_n
 - $w \in \mathbb{R}^M$ are learnable output layer weights
- Both kernel methods and deep neural networks extract new features from the inputs
 - For kernel methods, the features are pre-defined via the kernel function
 - For deep NN, the features are learned by the network
- Note: Kernels can also be learned from data ("kernel learning") in which case kernel methods and deep neural nets become even more similar in spirit ^(C) CS771: Intro to ML



Why Neural Networks Work Better: Another View^{1/}

- Linear models tend to only learn the "average" pattern
 - E.g., Weight vector of a linear classification model represent average pattern of a class
- Deep models can learn multiple patterns (each hidden node can learn one pattern)
 - Thus deep models can learn to capture more subtle variations that a simpler linear model



Neural Nets: Some Aspects

- Much of the magic lies in the hidden layers
- Hidden layers learn and detect good features
- Need to consider a few aspects
 - Number of hidden layers, number of units in each hidden layer
 - Why bother about many hidden layers and not use a single very wide hidden layer (recall Hornik's universal function approximator theorem)?
 - Complex networks (several, very wide hidden layers) or simpler networks (few, moderately wide hidden layers)?
 - Aren't deep neural network prone to overfitting (since they contain a huge number of parameters)?

Choosing the right NN architecture is important and a research area in itself. Neural Architecture Search (NAS) is an automated technique to do this





Representational Power of Neural Nets

• Consider a single hidden layer neural net with K hidden nodes





- Recall that each hidden unit "adds" a simple function to the overall function
- Increasing K (number of hidden units) will result in a more complex function
- Very large K seems to overfit (see above fig). Should we instead prefer small K?
- No! It is better to use large K and regularize well. Reason/justification:
 - Simple NN with small K will have a few local optima, some of which may be bad
 - Complex NN with large K will have many local optimal, all equally good (theoretical results on this)
- We can also use multiple hidden layers (each sufficiently large) and regularize well CS771: Intro to ML

Wide or Deep?

While very wide single hidden layer can approx. any function, often we prefer many, less wide, hidden layers



 Higher layers help learn more directly useful/interpretable features (also useful for compressing data using a small number of features)