Spring 2023!

Introduction to Machine Learning



Neural Networks I

Review: Non-linear classifiers



Review: Non-linear classifiers









$$z = \theta^{\top} \phi(x) + \theta_0$$

= $\theta_1 \phi_1(x) + \theta_2 \phi_2(x)$
+ θ_0
= $1 \cdot \phi_1(x) + 1 \cdot \phi_2(x)$
+ (-0.5)





















$$z = \theta^{\top} \phi(x) + \theta_0$$

= $\theta_1 \phi_1(x) + \theta_2 \phi_2(x)$
+ $\theta_3 \phi_3(x) + \theta_0$
= $1 \cdot \phi_1(x) + 1 \cdot \phi_2(x)$
+ $1 \cdot \phi_3(x) + (-0.5)$





These weights decide the how the input is transformed

These weights decide how the new features are combined



These weights decide the how the input is transformed

These weights decide how the new features are combined



- A rich hypothesis class of parameterized functions with multiple layers of "hidden" units that can learn to represent complex features of the input.
- Trained with gradient descent ("backpropagation").
- Originally inspired by the brain, but connection is tentative.

$$\mathbf{a} = \mathbf{f}(z) = \mathbf{f}\left(\left(\sum_{j=1}^{m} x_j w_j\right) + w_0\right) = \mathbf{f}(w^{\mathsf{T}} \mathbf{x} + w_0)$$



$$\mathbf{a} = \mathbf{f}(z) = \mathbf{f}\left(\left(\sum_{j=1}^{m} \mathbf{x}_{j} \mathbf{w}_{j}\right) + \mathbf{w}_{0}\right) = \mathbf{f}(\mathbf{w}^{\mathsf{T}} \mathbf{x} + \mathbf{w}_{0})$$



 $f(\mathbf{x}) = \sigma(\mathbf{x})$ $J(w, w_0)$ $= \sum_{i} L\left(NN(\mathbf{x}^{(i)}; w, w_0), \mathbf{y}^{(i)}\right)$ $L(g^{(i)}, y^{(i)}) =$ $y^{(i)} \log g^{(i)} + (1 - y^{(i)}) \log(1 - g^{(i)})$

$a = f(z) = f\left(\left(\sum_{j=1}^{m} x_{j}w_{j}\right) + w_{0}\right) = f(w^{T}x + w_{0})$ This is just logistic regression!!

 χ_1 $\chi_{\rm m}$ Neural Network with 0 hidden layers

 $= \sum L\left(NN(x^{(i)}; w, w_0), y^{(i)}\right)$

 $f(x) = \sigma(x) :$



$$A = f(Z) = f(W^T X + W_0)$$

- *W* is an $m \times n$ matrix,
- W_0 is an $n \times 1$ column vector,
- *X*, the input, is an m × 1 column vector,
- $Z = W^T X + W_0$, the *pre-activation*, is an $n \times 1$ column vector,
- A, the *activation*, is an n × 1 column vector,



Neural Network Activation Function



Neural Network Activation Function



Neural Network Activation Function















Two Layer Neural Network for Classification



Two Layer Neural Network for Classification

J



$$g = \sigma(U^{\top}A + U_0)$$
$$= \sigma(U^{\top}f(W^{\top}x + W_0) + U_0)$$

$$(W, W_0, U, U_0) =$$

= $\frac{1}{n} \sum_{i=1}^n y^{(i)} \log g^{(i)} + (1 - y^{(i)}) \log(1 - g^{(i)})$

This objective is a differentiable function of model parameters and hence we can perform gradient-based optimization as in Logistic regression





$$Z^{l} = W^{l} A^{l-1} + W^{l}_{0}$$
$$A^{l} = f^{l}(Z^{l})$$



Last activation function depends on problem at hand



Last activation function and loss function depends on problem at hand









Gradient Descent

- All model parameters are differentiable functions of the loss ⇒ can train with (stochastic) GD
- Backpropagation: an approach to efficiently obtain gradients in multilayer neural networks.
- An instance of **dynamic programming**. Avoid redundant computation by breaking down the gradient expression into shared subproblems.

One layer neural network with linear activation: Input dimension = 4, output dimension = 3

$$\begin{array}{ll} A^{(1)} = \underbrace{W^{(1)}}_{4\times3} \\ x + W^{(1)}_{0} \\ x + W^{(1)}$$

One layer neural network with linear activation: Input dimension = 4, output dimension = 3

$$\begin{array}{ll} A^{(1)} = \underbrace{W^{(1)}}_{4\times3} \\ x + W^{(1)}_{0} \\ x + W^{(1)}$$



Matrix chain rule

One layer neural network with linear activation: Input dimension = 4, output dimension = 3

$$\begin{array}{ll} A^{(1)} = \underbrace{W^{(1)}}_{4\times3} + W^{(1)}_{0} \\ 3\times1 & 4\times3 & 4\times1 & 3\times1 \end{array} \qquad \begin{array}{ll} \underline{loss} = L(A^{(1)}, y) = \sum_{k=1}^{3} (A^{(1)}_{k} - y_{k})^{2} \\ 1\times1 & 1\times1 \end{array}$$



Matrix chain rule Scalar chain rule

One layer neural network with linear activation: Input dimension = 4, output dimension = 3

$$\begin{array}{ll} A^{(1)} = \underbrace{W^{(1)}}_{4\times3} + W^{(1)}_{0} \\ 3\times1 & 4\times3 & 4\times1 & 3\times1 \end{array} \qquad \qquad \underbrace{loss}_{1\times1} = L(A^{(1)}, y) = \sum_{k=1}^{3} (A^{(1)}_{k} - y_{k})^{2} \\ \end{array}$$

$$\frac{\partial \text{loss}}{\partial W_{0,j}^{(1)}} = \frac{\partial A^{(1)}}{\partial W_{0,j}^{(1)}} \frac{\partial \text{loss}}{\partial A^{(1)}}$$

$$\frac{\partial W_{0,j}^{(1)}}{1 \times 3} \frac{\partial A^{(1)}}{3 \times 1}$$

Matrix chain rule

$$\frac{\partial \text{loss}}{\partial W_0^{(1)}} = \frac{\partial A^{(1)}}{\partial W_0^{(1)}} \frac{\partial \text{loss}}{\partial A^{(1)}}$$

3×1 3×3 3×1

Can calculate for all W0 via matrix multiplications!

One layer neural network with linear activation: Input dimension = 4, output dimension = 3

$$\begin{array}{ll}
A^{(1)} &= \underbrace{W^{(1)}}_{4\times3} \stackrel{\top}{}_{4\times1} \stackrel{+}{}_{3\times1} \stackrel{W^{(1)}_{0}}{}_{1\times1} &= L(A^{(1)}, y) = \sum_{k=1}^{3} (A^{(1)}_{k} - y_{k})^{2} \\
\times 1 & \frac{\partial \log s}{\partial W^{(1)}_{j,k}} = \sum_{p=1}^{3} \frac{\partial A^{(1)}_{p}}{\partial W^{(1)}_{j,k}} \frac{\partial \log s}{\partial A^{(1)}_{p}} = \frac{\partial A^{(1)}_{k}}{\partial W^{(1)}_{j,k}} \frac{\partial \log s}{\partial A^{(1)}_{k}} = x_{j} \frac{\partial \log s}{\partial A^{(1)}_{k}}
\end{array}$$

One layer neural network with linear activation: Input dimension = 4, output dimension = 3

$$\begin{array}{ll}
A^{(1)} &= \underbrace{W^{(1)}}_{4\times 3} \stackrel{\times}{}_{4\times 1} \stackrel{\times}{}_{3\times 1} \stackrel{W^{(1)}_{0}}{}_{1\times 1} &= L(A^{(1)}, y) = \sum_{k=1}^{3} (A^{(1)}_{k} - y_{k})^{2} \\
\times^{1} & \frac{\partial \log s}{\partial W^{(1)}_{j,k}} = \sum_{p=1}^{3} \frac{\partial A^{(1)}_{p}}{\partial W^{(1)}_{j,k}} \frac{\partial \log s}{\partial A^{(1)}_{p}} = \frac{\partial A^{(1)}_{k}}{\partial W^{(1)}_{j,k}} \frac{\partial \log s}{\partial A^{(1)}_{k}} = x_{j} \frac{\partial \log s}{\partial A^{(1)}_{k}}
\end{array}$$

$$\frac{\partial \text{loss}}{\partial W^{(1)}} = x \left(\frac{\partial \text{loss}}{\partial A^{(1)}}\right)^{\top}$$

$$4 \times 3 \qquad 4 \times 1 \qquad 3 \times 1$$

Again, can calculate everything at once via matrix multiplications (outer products)

One layer neural network with linear activation: Input dimension = 4, output dimension = 3

$$\begin{array}{l}
A^{(1)} = \underbrace{W^{(1)}}_{4\times3} \stackrel{\top}{}_{4\times1} + \underbrace{W^{(1)}}_{3\times1} & \underbrace{loss}_{1\times1} = L(A^{(1)}, y) = \sum_{k=1}^{3} (A^{(1)}_{k} - y_{k})^{2} \\
\times 1 \quad \frac{\partial loss}{\partial W^{(1)}_{j,k}} = \sum_{p=1}^{3} \frac{\partial A^{(1)}_{p}}{\partial W^{(1)}_{j,k}} \frac{\partial loss}{\partial A^{(1)}_{p}} = \frac{\partial A^{(1)}_{k}}{\partial W^{(1)}_{j,k}} \frac{\partial loss}{\partial A^{(1)}_{k}} = x_{j} \frac{\partial loss}{\partial A^{(1)}_{k}} \\
\frac{\partial loss}{\partial A^{(1)}_{k}} = \left(\frac{\partial loss}{\partial A^{(1)}_{k}} \right)^{\top} \quad \text{Eucrements of the second se$$

Exercise: convince yourselves that this is correct!!

Backpropagation Intuition

Forward propagation to obtain the output (model's guess)



Backpropagation Intuition

Forward propagation to obtain the output (model's guess)



Backpropagation to obtain gradients with respect to the loss

Backpropagation Intuition

Forward propagation to obtain the output (model's guess)



Backpropagation to obtain gradients with respect to the loss

$$\begin{array}{c} \underbrace{X = A^{0}}_{\substack{W_{0}^{1}\\W_{0}^{1}\\\hline\\W_{0}^{1$$



$$\begin{array}{c} \underbrace{X = A^{0}}_{\substack{W_{1}^{0}\\ W_{0}^{1}\\ \hline{W_{0}^{1}\\ \hline{W_{0}^{1}$$

$$\begin{array}{c} \underbrace{X = A^{0}}_{\substack{W_{1}^{1}\\W_{0}^{1}\\\vdots\\W_{0}^{1$$

$$\begin{array}{c} \overset{\forall y}{\underset{\substack{\lambda = A^{0} \\ W_{0}^{1} \\ \vdots \\ \partial U_{0}^{2}}}} \underbrace{Z^{1}}_{\substack{\lambda = 0 \\ \partial U_{0}^{2}}} \underbrace{f^{1}}_{\substack{\lambda = 0 \\ \partial U_{0}^{2}}} \underbrace{Z^{2}}_{\substack{W_{0}^{0} \\ \partial U_{0}^{2}}} \underbrace{f^{2}}_{\substack{\lambda = 0 \\ \partial U_{0}^{2}}} \underbrace{A^{2}}_{\substack{\partial U_{0}^{2} \\ \partial U_{0}^{2}}} \underbrace{A^{L-1}}_{\substack{W_{0}^{1} \\ \partial U_{0}^{2}}} \underbrace{Z^{L}}_{\substack{\partial U_{0}^{2} \\ \partial U_{0}^{2}}} \underbrace{f^{L}}_{\substack{\lambda = 0 \\ \partial U_{0}^{2}}} \underbrace{A^{L}}_{\substack{\partial U_{0}^{2} \\ \partial U_{0}^{2}}} \underbrace{A^{L-1}}_{\substack{\partial U_{0}^{2} \\ \partial U_{0}^{2}}} \underbrace{Z^{L}}_{\substack{\partial U_{0}^{2} \\ \partial U_{0}^{2}}} \underbrace{f^{L}}_{\substack{\partial U_{0}^{2} \\ \partial U_{0}^{2} \\ \partial U_{0}^{2}}} \underbrace{A^{L-1}}_{\substack{\partial U_{0}^{2} \\ \partial U_{0}^{2} \\ \partial U_{0}^{2} \\ \partial U_{0}^{2} \\ i \end{bmatrix}} \underbrace{A^{\ell}(\ell) = W^{\ell}(\ell)^{\top} A^{\ell}(\ell-1) + W_{0}^{\ell}}_{0} \\ \hline \underbrace{\partial U_{0}^{2} \\ \partial W_{j,k}^{\ell}} = \frac{\partial Z_{k}^{(\ell)}}{\partial Z_{k}^{\ell}} \frac{\partial U_{0}^{2} \\ \partial Z_{k}^{\ell}} = A_{j}^{\ell}(\ell-1) \frac{\partial U_{0}^{2} \\ \partial Z_{k}^{\ell}} \underbrace{\partial U_{0}^{2} \\ \partial W^{\ell}(\ell)} = A^{\ell}(\ell-1) \left(\frac{\partial U_{0}^{2} \\ \partial Z^{\ell}(\ell)}\right)^{\top} \\ \underbrace{\partial U_{0}^{2} \\ i \end{bmatrix} \underbrace{A^{\ell}(\ell) \\ i \end{bmatrix} \underbrace{A^{\ell}(\ell) \\ \partial Z^{\ell}(\ell) \\ \partial Z^{\ell}(\ell) \\ \partial Z^{\ell}(\ell) \\ \partial A^{\ell}(L) \\ \partial A^{\ell}(L$$

$$\begin{split} & \underbrace{X = A^{0}}_{\substack{W_{0}^{l} \\ \downarrow 0}} \underbrace{Z^{1}}_{\substack{\partial \log s \\ \partial ds}} \underbrace{f^{1}}_{\substack{\partial \log s \\ \partial ds}} \underbrace{Z^{2}}_{\substack{\partial \log s \\ \partial ds}} \underbrace{f^{2}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{Z^{2}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{f^{2}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{Z^{L}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{f^{L}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{f^{L}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{Z^{L}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{f^{L}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{Z^{L}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{f^{L}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{Z^{L}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{f^{L}}_{\substack{\partial \log s \\ \partial dss}} \underbrace{Z^{L}}_{\substack{\partial dss}} \underbrace$$







Can calculate $\frac{\partial loss}{\partial Z^{(l)}}$ for all layers in one "backward" pass due to shared computations

$$\frac{\partial \text{loss}}{\partial Z^{(l)}} \text{ needed to calculate } \frac{\partial \text{loss}}{\partial W^{(\ell)}} \text{ and } \frac{\partial \text{loss}}{\partial W_0^{(l)}}$$
$$\frac{\partial \text{loss}}{\partial W^{(\ell)}} = A^{(\ell-1)} \left(\frac{\partial \text{loss}}{\partial Z^{(\ell)}}\right)^{\top}$$

That's it!