Backpropagation

Designing, Visualizing and Understanding Deep Neural Networks

CS W182/282A

Instructor: Sergey Levine UC Berkeley



Neural networks

Drawing computation graphs



what **expression** does this compute? equivalently, what **program** does this correspond to?

 $||(x_1\theta_1 + x_2\theta_2) - y||^2$

this is a MSE loss with a linear regression model

neural networks are computation graphs

if we design **generic tools** for computation graphs, we can train **many kinds** of neural networks

Drawing computation graphs

a simpler way to draw the same thing:



I'll drop the decorator from now on...

what **expression** does this compute? equivalently, what **program** does this correspond to?

 $||(x_1\theta_1 + x_2\theta_2) - y||^2$

this is a MSE loss with a linear regression model

neural networks are computation graphs

if we design **generic tools** for computation graphs, we can train **many kinds** of neural networks

Logistic regression



let's draw the computation graph for **logistic regression** with the negative log-likelihood loss

 $p_{\theta}(y|x) = \frac{\exp(x^{T}\theta_{y})}{\sum_{y'} \exp(x^{T}\theta_{y'})}$ $-\log p_{\theta}(y|x) = -x^{T}\theta_{y} + \log \sum_{y'} \exp(x^{T}\theta_{y'})$

Logistic regression

a simpler way to draw the same thing:

$$p_{\theta}(y|x) = \frac{\exp(x^T \theta_y)}{\sum_{y'} \exp(x^T \theta_{y'})}$$

$$-\log p_{\theta}(y|x) = -x^{T}\theta_{y} + \log \sum_{y'} \exp(x^{T}\theta_{y'})$$

$$f_{\theta}(x) = \begin{bmatrix} x^{T} \theta_{y_{1}} \\ x^{T} \theta_{y_{2}} \\ \vdots \\ x^{T} \theta_{y_{m}} \end{bmatrix} \qquad f_{\theta}(x) = \theta x$$

matrix



$$p_{\theta}(y=i|x) = \operatorname{softmax}(f_{\theta}(x))[i] = \frac{\exp(f_{\theta,i}(x))}{\sum_{j=1}^{m} \exp(f_{\theta,j}(x))}$$



Drawing it even more concisely

Notice that we have **two types** of variables:

data (e.g., x, y), which serves as input or target output

parameters (e.g., θ)

the parameters usually affect one specific operation

(though there is often *parameter sharing*, e.g., conv nets – more on this later)



Neural network diagrams

(simplified) computation graph diagram

neural network diagram



Logistic regression with features





 $\operatorname{softmax}(x^T\theta)$



pop quiz: what is the dimensionality of θ ?

Learning the features

Problem: how do we represent the learned features?

Idea: what if each feature is a (binary) logistic regression output?

which layer $w_1^{(1)}$ $w_1^{(1)}$ which feature = rows of weight matrix



Let's draw this!

$$\phi(x) = \begin{pmatrix} \operatorname{softmax}(x^T w_1^{(1)}) \\ \operatorname{softmax}(x^T w_2^{(1)}) \\ \operatorname{softmax}(x^T w_3^{(1)}) \end{pmatrix} = \sigma(W^{(1)}x) \qquad p(y|x) = \operatorname{softmax}(\phi(x)^T \theta)$$





Simpler drawing



simpler way to draw the same thing:

even simpler:





Doing it multiple times





Activation functions

$$\phi_1(x) = \operatorname{softmax}(x^T w_1^{(1)}) = \frac{1}{1 + \exp(-x^T w_1^{(1)})}$$

we don't have to use a **sigmoid**!

a wide range of non-linear functions will work these are called **activation functions** we'll discuss specific choices later why **non-linear?**

$$a^{(2)} = \sigma(W^{(2)}\sigma(W^{(1)}x))$$

if
$$\sigma(z) = z$$
, then...
 $a^{(2)} = W^{(2)}W^{(1)}x = Mx$

multiple linear layers = one linear layer

enough layers = we can represent anything (so long as they're nonlinear)





Demo time!



Source: https://playground.tensorflow.org/

Aside: what's so neural about it?



Training neural networks

What do we need?

1. Define your model class



2. Define your loss function

negative log-likelihood, just like before

3. Pick your **optimizer**

stochastic gradient descen what do we need?

 $\nabla_{\theta} \mathcal{L}(\theta)$

4. Run it on a big GPU

$$dt = \begin{pmatrix} \frac{d\mathcal{L}(\theta)}{d\theta_1} \\ \frac{d\mathcal{L}(\theta)}{d\theta_2} \\ \vdots \\ \frac{d\mathcal{L}(\theta)}{d\theta_n} \end{pmatrix}$$



Chain rule for neural networks

A neural network is just a composition of functions

So we can use chain rule to compute gradients!



Does it work?

 $\frac{d\mathcal{L}}{dW^{(1)}} = \frac{dz^{(1)}}{dW^{(1)}} \frac{da^{(1)}}{dz^{(1)}} \frac{dz^{(2)}}{da^{(1)}} \frac{d\mathcal{L}}{dz^{(2)}}$

We can calculate each of these Jacobians!

Example:

Why might this be a **bad** idea?

 $z^{(2)} = W^{(2)}a^{(1)}$

 $\frac{dz^{(2)}}{da^{(1)}} = W^{(2)T}$

if each $z^{(i)}$ or $a^{(i)}$ has about n dims... each Jacobian is about $n \times n$ dimensions matrix multiplication is $O(n^3)$

do we care? AlexNet has layers with 4096 units...

Doing it more efficiently



this is **always** true because the loss is scalar-valued!

compute
$$\frac{da^{(1)}}{dz^{(1)}}\delta = \gamma$$

$$\frac{d\mathcal{L}}{dW^{(1)}} = \frac{dz^{(1)}}{dW^{(1)}}\gamma$$

this product is cheap: $O(n^2)$

The backpropagation algorithm



backward pass:

initialize $\delta = \frac{d\mathcal{L}}{dz^{(n)}}$ for each f with input x_f & params θ_f from end to start:



Let's walk through it...



backward pass:

initialize
$$\delta = \frac{d\mathcal{L}}{dz^{(n)}}$$

for each f with input x_f & params θ_f from end to start:

$$\begin{array}{c} \bullet & \frac{d\mathcal{L}}{d\theta_f} \leftarrow \frac{df}{d\theta_f} df \\ \bullet & \delta \leftarrow \frac{df}{dx_f} \delta \end{array}$$

Practical implementation

Neural network architecture details



Some things we should figure out:

How many layers?

How big are the layers?

What type of **activation function**?



Bias terms softmax sigmoid \mathcal{X} linear $z^{(1)}$ $a^{(1)}$ linear $z^{(2)}$ $\mathcal{L}(z^{(2)})$ layer layer 3x1 3x1 2x1 $W^{(1)}b^{(1)}$ $W^{(2)}b^{(2)}$

Linear layer:

 $z^{(i+1)} = W^{(i)}a^{(i)}$ problem: if $a^{(i)} = \vec{0}$, we always get 0...

Solution: add a "bias":

has nothing to do with bias/variance bias

$$z^{(i+1)} = W^{(i)}a^{(i)} + b^{(i)}$$

`additional parameters in each linear layer

What else do we need for backprop?



forward pass: calculate each $a^{(i)}$ and $z^{(i)}$ backward pass:

initialize $\delta = \frac{d\mathcal{L}}{dz^{(n)}}$

for each f with input x_f & params θ_f from end to start:

$$\frac{d\mathcal{L}}{d\theta_f} \leftarrow \frac{df}{d\theta_f} \delta$$
$$\delta \leftarrow \frac{df}{dx_f} \delta$$

for each function, we need to compute:

$$\frac{df}{d\theta_f}\delta \qquad \qquad \frac{df}{dx_f}\delta$$

linear layer softmax + cross-entropy sigmoid ReLU





for each function, we need to compute: $\frac{df}{d\theta_f}\delta = \frac{df}{dx_f}\delta$

linear layer: $z^{(i+1)} = W^{(i)}a^{(i)} + b^{(i)}$ z = Wa + b (just to simplify notation!)

$$\frac{dz}{db}\delta = \delta$$

$$z_i = \sum_k W_{ik}a_k + b_i \quad \frac{dz_i}{db_j} = \text{Ind}(i=j) \quad \frac{dz}{db} = \mathbf{I}$$



for each function, we need to compute: $\frac{df}{d\theta_f}\delta = \frac{df}{dx_f}\delta$

linear layer: $z^{(i+1)} = W^{(i)}a^{(i)} + b^{(i)}$ z = Wa + b (just to simplify notation!)

$$\frac{dz}{da}\delta = W^T\delta$$

$$z_i = \sum_k W_{ik}a_k + b_i \quad \frac{dz_i}{da_k} = W_{ik} \quad \frac{dz}{da} = W^T \qquad \left(\frac{dy}{dx}\right)_{ij} = \frac{dy_j}{dx_i}$$



for each function, we need to compute: $\frac{df}{d\theta_f}\delta = \frac{df}{dx_f}\delta$

linear layer: $z^{(i+1)} = W^{(i)}a^{(i)} + b^{(i)}$ z = Wa + b (just to simplify notation!)

$$\frac{dz}{da}\delta = W^T\delta \qquad \frac{dz}{dW}\delta = \delta a^T \qquad \frac{dz}{db}\delta = \delta$$

$$\frac{df}{dx_f}\delta \qquad \qquad \frac{df}{d\theta_f}\delta$$

Backpropagation recipes: sigmoid



Backpropagation recipes: ReLU



for each function, we need to compute:

$$\frac{df}{d\theta_f}\delta \quad \frac{df}{dx_f}\delta$$

$$f_i(z_i) = \max(0, z_i)$$
 $\frac{df_i}{dz_i} = \operatorname{Ind}(z_i \ge 0)$

$$\left(\frac{df}{dz}\delta\right)_i = \operatorname{Ind}(z_i \ge 0)\delta_i$$

Summary



forward pass: calculate each $a^{(i)}$ and $z^{(i)}$ backward pass:

initialize $\delta = \frac{d\mathcal{L}}{dz^{(n)}}$

for each f with input x_f & params θ_f from end to start:

$$\frac{d\mathcal{L}}{d\theta_f} \leftarrow \frac{df}{d\theta_f} \delta$$
$$\delta \leftarrow \frac{df}{dx_f} \delta$$

for each function, we need to compute:

$$\frac{df}{d\theta_f}\delta \qquad \qquad \frac{df}{dx_f}\delta$$

linear layer softmax + cross-entropy sigmoid ReLU