

① Deep learning, goodfellow (2016)

A feedforward network defines a mapping $y = f(\theta, x)$ & learns the value of the parameters θ that result in the best func p

Convolutional networks for object recog are a specialized kind of feedforward network

conceptual stepping stone to recurrent networks

FF called networks as represented by composing many different functions

The nets have a chain structure $f(x) = f^3(f^2(f'(x)))$

f' = first layer $f^{n/3}$ = last layer

the length of the chain is the depth of the model

$f(x)$ provides us with approx value of y $y \approx f^*(x)$

Dimensions of the network is the width

NN conceptually starts w/ linear models, i.e. lin reg, log reg, and adoptions are made to overcome non-linear limitations

To extend linear models to rep non-linear functions of x — we apply the lin model to a transform version of $x \rightarrow \phi(x)$

where ϕ = non-linear transformation

(2)

ϕ provides a new representation of x

How to choose the ϕ mapping?

- ① • use generic ϕ such as infinite-dimensional
 - that is implicit used by kernel machines or RBF Kernel
 - if $\phi(x)$ is high enough dimension, can always find a fit in the ff train data.
 - But generalization will be poor
 - Do not include enough info to solve advanced into
- ② • manually engineer ϕ
 - Dominant convention prior to DL
 - requires a lot of human effort for each task
 - specific to domains w/ little knowledge transfer between domains
- ③ • Deep learning strategy is to learn ϕ
 - we have a model: $y = f(x; \theta, w) = \phi(x; \theta)^T w$
 - θ paras , ϕ function , w weights/mappings
= Deep forward network
 - This gives up on convenience of problem
 - Parameterize $\phi(x; \theta)$ & use opt algo to find good θ that correspond to good output

③

6.2 Gradient Based Learning

training NN is not that much diff to any other ML model w/ gradient Descent

The largest diff between linear models & NNs is the non-linearity cause loss functions to become non-convex

This means NNs are trained using iterative gradient based optimizers

rather than linear equation solvers (lin reg)

or convex opt algos w/ global convergence guarantees used to train (log reg or SVM)

Convex optim converges from any initial parameters (in theory)

Stochastic gradient descent has no convergence guarantees & is sensitive to starting parameters (optim applied to non-convex)

for FF Nets, must set initial para to small random values

4

6.2.2 Output Units

Cost function chosen is tightly related to the output units of the NN

typically uses cross entropy of data dist & model dist

Output then determines form of cross ent function

Any hidden unit can be ea or used as an output unit

A FF NN provides hidden features given by
 $h = f(u; \theta)$

the role of the final output layer is to provide some additional transformation from the features to complete the network task

6.2.2.4 Other types of Output Units

NNs represent a function $f(u, \theta)$

the outputs are not direct predictions of y

instead $f(u, \theta) = w$ provides parameters for a distribution over y

loss func is then interpreted as: $-\log p(y; w(x))$

Therefore the model can be adapted to learn the variance or SD

(5)

6.3 Hidden Units

How to choose type of hidden unit to use in the hidden layers?

At time of book, was considered to be an active area of research w/o definitive theoretical principles

Rectified linear units are a popular choice

Picking the best in advance is not usually possible. Instead requires a design process of trial & error

Not all hidden units are differentiable at all input points

e.g. ReLU $g(z) = \max\{0, z\}$ is not diff @ $z=0$

This may seem like it invalidates g for use w/ a gradient-based learning algo

but Practice gd still works well. This is because NNs do not arrive @ the min of a loss func. — but g's reduce value significantly

Hidden units that are not diff usually only exist @ a small number of points

6

6.3.1 Rectified linear unit

active func where $g(z) = \max\{0, z\}$

easy to opt as similar to linear units. only diff = 0 across half the domain ($-z$)

Second the derivatives of the active unit (>0) are large & consistent

large Derivative =

- Small change in input = large Δ in output
- Allows for faster learning & bigger steps towards min the error
- Backprop gradients have a gradient that becomes very small in the layer — hence large gradient is a good qual to counter

Consistent Derivative:

- Stable = not wildly big changes per unit Δ of input
- reliable converge
- Avoid oscillations = hard to converge

Second Deriv of active $2=1$ & $2=0=0$

= gradient dir is far more useful vs act func that intro second-order effects

Rehs applied our top of affine transform

$$h = g(W^T x + b)$$

1

When initializing the parameters, it can be good to start w/ small values of b_f & 0.1

ensures most of the ~~hidden~~ units will remain active for most of the training

↳ Allows digerable units to flow through

There are some generalized versions of rect linear units which seek to address O issue

6.3.2 Logistic Sigmoid & Hyperbolic tangent

Prior to rectlin, most NNs used sig act func

$$g(z) = \sigma(z)$$

or hyperbolic tangent activ

$$g(z) = \tanh(z)$$

related as $\tanh(z) = 2\sigma(2z) - 1$

Issue w/ sig is that it saturates at high vs low (-) values

i.e. a big change in the input doesn't change much

only when value around 0 is the gradient significant

Saturate = learn slow or poorly

8

Wide-Spread Saturation makes gradient based learning difficult

for this reason, hidden sig units in FF Nets are now discouraged

Only approp when units are coupled w/ a loss function that undo sigmoid output
($\log \rightarrow \exp$)

(other than FFNN)
sig functions are common in other sets' FFNN
hence there inclusion may still happen

6.3.3 Other Hidden Units

In general a wide variety of hidden funcs work well

Many unpublished = just as good as popular

During research it is common to test many & find they test comparable

Newly discovered hidden types are considered to be so common to be considered as uninteresting

Using no activation can be considered as using an identity matrix

Not all activs can be linear otherwise the output of the model would have no non-lin

Acceptable for some to be linear

9

6.4. Architecture Design

Architecture = overall structure of the network

- * How many units it should have
 - * How they should be connected
- * Most NNs = organized into layers
- * these layers one in a chain structure
i.e. each a func of the prev func

Main considerations of network = Depth & width

Deep = less paras needed (width)
but harder to optimize

Ideal Network archi = Experiment &
guide via validation set

16

6.5 Backpropagation & Other Diff Algos

- * ff NN accepts an input x & produces \hat{y}
- * info flows forward through the network
- * initial info propagates up to the hidden units @ each layer
- * training forward prop continues until we create a scalar cost $J(\theta)$

Backprop Allows the info from the cost to the flow backward through the network in order to compute the network

Generally comp an analytical expression is easy, but evaluating is comp costly

Backprop does so using simple & inexpensive proc

Backprop = just the method for computing the gradient

→ Other algo, i.e., SGD, performs the learning using this gradient

We will work through how to compute the gradient for $\nabla_{\theta} f(x, y)$

x is set of vars whose derivs are desired

y = inputs to func but Derivs not required

the gradient we most require is that of the cost function w/ respect to the params $\nabla_{\theta} J(\theta)$

(11)

6.5.2 Chain Rule of Calculus

Backprop is the algo to compute chain rule, with a specific order of calculations

$n = \mathbb{R}$ $f, g: \text{functions } \mathbb{R} \rightarrow \mathbb{R}$

$$y = g(x), z = f(g(x)) = f(y)$$

then chain rule: $\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$

If x, y are vectors: $\mathbb{R}^n \times g: \mathbb{R}^n \rightarrow \mathbb{R}^m \quad f: \mathbb{R}^m \rightarrow \mathbb{R}$

$$\frac{\partial z}{\partial x_i} = \sum \frac{\partial z}{\partial y_i} \cdot \frac{\partial y_i}{\partial x_i}$$

$$\text{in vector notation: } \nabla_x z = \left(\frac{\partial g}{\partial x} \right)^T \nabla_y z$$

where $\frac{\partial g}{\partial x}$ is the $n \times m$ Jacobian matrix of g

gradient of value z obtained by multiplying Jacob matrix by the gradient $\nabla_y z$

Backprop algo = perform Jacob - gradient product for each operation in the graph

Usually Back prop is applied to n-dim tensor, not just vectors.

Conceptually, this is the same as vectors but rearranged

Backprop = Multi Jacob by gradients

(12)

6.5.3 Recursively applying the chain rule to obtain backprop

use chain rule to obtain gradient of a scalar (loss func) w/ respect to any node in the graph - that produced scalar

Computed gradients may be stored to reduce comp load

Start w/ version of Backprop that specifies the actual gradient directly

In the order that the transversed graph & chain rule ~~sho~~ direct

Consider a ~~fin~~ Comp graph desc how to comp a single scalar — $u^{(n)}$

the scalar = gradient we want to obtain w/ respect to the n_i input nodes $u^{(1)} \text{ to } u^{(n)}$

We wish to calc: $\frac{\partial u^{(n)}}{\partial u^{(i)}}$ for all $i \in \{1, 2, \dots, n\}$

$u^{(1)} \text{ to } u^{(n)}$ = Parameters of the model

Assume graph allow recursive calcs

Start @ $u^{(n_i + 1)}$ and go upto $u^{(n)}$

(B)

ff Calculates a graph G

Backprop creates a "new" graph B w/
one node per node of G

B order = Traverse of G

① each node of B we calc deriv $\frac{\partial u^{(n)}}{\partial u^{(l)}}$

assoc w/ ff node $u^{(l)}$

this is the chain rule w/ respect to scalar out

$$\frac{\partial u^{(n)}}{\partial u^{(l)}} = \sum_{i:j \in \text{Pa}(u^{(l)})} \frac{\partial u^{(n)}}{\partial u^{(l)}} \cdot \frac{\partial u^{(l)}}{\partial u^{(i)}} \quad (\text{Algob-2})$$

\uparrow edge

B = one edge from node $u^{(l)}$ to $u^{(i)}$ of G

Computation for Backprop scales to linear w/
number of edges in G

where the comp of each edge corresponds to
the partial derivative compute w/ respect
to one of its parameters. Parents

(14)

Algo 1 fast forward

Proc to perform mapping n_i inputs $u^{(1)}$ to $u^{(n)}$ to an output $u^{(n)}$

This begins a comp graph where each node computes numerical value $u^{(i)}$ by applying a function $f(i)$

to a set of args $A^{(i)}$ that comprises the values of the prev nodes $u^{(j)} : j < i$

the input of graph is vector x

and is set into the first n_i nodes $u^{(i)}$ to $u^{(n)}$

the output is read from the last node $u^{(n)}$

x vector = single data point

for $i = 1, \dots, n_i$ do

$$u^{(i)} \leftarrow x^i$$

Dimension
value holder

($u^{(i)}$ = i th input node)

end for

for $i = n_i + 1, \dots, n$ do

$$A^{(i)} \leftarrow \{u^{(j)} \mid j \in Pa(u^{(i)})\}$$

Set of values from parent nodes

$$u^{(i)} \leftarrow f^{(i)}(A^{(i)})$$

Parent nodes of $u^{(i)}$

end for

return $u^{(n)}$

function applied to A vector

Output

(15)

Algo 2 Backprop

- * Calculate derivatives of $u^{(n)}$
- * w/ respect to variables in the graph
- * simple example where all variables are scalars (not vec)
- * Compute derivs w/ respect to $u^{(1)}, \dots, u^{(n)}$
- * Computes Dev of all nodes in graph
- * Comp cost = proportional to num of edges
- * Assumes partial derivative w/ each edge requires const time
- * Each $\frac{\partial u^{(n)}}{\partial u^{(i)}}$ is a func of the parents $u^{(j)}$ of $u^{(i)}$

- Start with FF to obtain network - Algo 1
- init grad table - A data structure to store Devs that have been calced. $GT[u^{(i)}] = \text{value of } \frac{\partial u^{(n)}}{\partial u^{(i)}}$

$\text{grad-table}[u^{(n)}] \leftarrow 1$ using stored values

for $j = n-1$ down to 1
 compute $\frac{\partial u^{(n)}}{\partial u^{(j)}} = \sum_{i:j \in \text{Pa}(u^{(i)})} \frac{\partial u^{(n)}}{\partial u^{(i)}} \cdot \frac{\partial u^{(i)}}{\partial u^{(j)}}$

$\text{grad-table}[u^{(j)}] + \sum_{i:j \in \text{Pa}(u^{(i)})} \text{grad-table}[u^{(i)}] \frac{\partial u^{(i)}}{\partial u^{(j)}}$

end for

return $\{\text{grad-table}[u^{(i)}] | i = 1, \dots, n\}$

16

6.5.4 Backprop in full connect MLP

Algo 6.3 forward Prop to cost func

Req Network depth

Req $W^{(i)}$, $i \in \{1, \dots, l\}$ weight mats of model

Req $b^{(i)}$, $i \in \{1, \dots, l\}$

Req x input

Req y target output

$$h^{(0)} = \Theta(x)$$

for $k = 1, \dots, l$ do

$$a^{(k)} = b^{(k)} + W^{(k)} h^{(k-1)}$$

$$h^{(k)} = f(a^{(k)})$$

end for

$$\hat{y} = h^{(l)}$$

$$\frac{\partial J}{\partial \theta} = L(\hat{y}, y) + \lambda \Omega(\theta)$$

(17)

6.4 Algo Backprop for Deep NN

- Comp yields the gradients $\nabla_a^{(k)}$ on the activations $a^{(k)}$ for each layer k
- Start from output layer and go backwards to 1st hidden layer
- Gradients = how each layer should change to reduce error
- from these, obtain the ~~pa~~ gradient of the pars of each layer
- the grads on the w & b can be immediately used as part of a SGD update

After forward comp, compute gradient on the output layer:

$$g \leftarrow \nabla_{\hat{y}} L = \nabla_{\hat{y}} L(\hat{y}, y)$$

for $K=1, 1-1, \dots, 1$ do

Convert the grad on the layers output into a grad on the pre-linearity activation

$$g \leftarrow \nabla_{a^{(k)}} J = g \odot f'(a^{(k)})$$

Comp grads on weights & biases

$$\nabla_{b^{(k)}} J = g + \lambda \nabla_{b^{(k)}} \Omega(\theta)$$

$$\nabla_{W^{(k)}} J = g h^{(k-1)^T} + \lambda \nabla_{W^{(k)}} \Omega(\theta)$$

Propagate the grads to the next lower level layer actvns

$$g \leftarrow \nabla_{h^{(k-1)}} J = W^{(k)}^T g$$

end for