

卷之三

卷之三

A lot of important work on what happened in the 80's
and 90's. Not money itself but inflation

were found and thus completely disappeared in the 60's. Only in the last few years neural networks made a big comeback thanks to lots of data and fast GPU.

Supervised classification
Classification is the task of taking an input
and giving it a label

You have lots of examples where the
therapeutic actions have already been demonstrated

Given a new example, how do we know to which class it belongs?

Classification or prediction is the common building block of machine learning

Logistic Verifier (or linear)
condition

$$Y = b + X^T W \rightarrow \text{logistic}$$

input vector - biases etc.
" " + the best Y

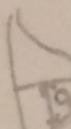
Training funds will be 10% less.

(R)
 S(yi) = $\frac{y_i}{\sum x_j y_j}$ \rightarrow softmax

If you multiply both sides by either $\frac{1}{x}$ or $\frac{1}{y}$, dividing by x or y respectively, we get

$$\frac{1 - \cos \theta}{2} = \frac{1 - \frac{128}{278}}{2} = \frac{149}{556}$$

This makes it
a bit easier for
the optimizer to do
its job

weight initialization

 draw weight randomizing
 from range & determine
 the order of increments
 of the outputs at the
 initial point of the optimization
 it's much better to begin
 with a low sigma (distribution
 around about things)

We need to find "W" and "b" such that
 we have a low deviation
 we have a low error and a high accuracy
 For the in correct case

$$f_i = \frac{1}{N} \sum_i D(s_i(w^T x_i + b), l_i)$$

 minimize D averaged over the entire training set
 we want the loss to be small
 To find the weights that cause the loss to be small
 we use gradient descent

Adding random error values to a very large
value can introduce errors in position
MEAN $X_1 = 0$ \leftarrow no limit on variables
possible

$$\text{VARIANCE } \sigma^2(X_i) = \sigma^2(X_j)$$

Measuring performance

Run classifier in your images with labels and see how many

it gets right.

That's my error measure.

Then you try your classifier on images it has never seen before

and your performance is not so good.

Imagine the classifier computes the new model with one of the other random class images and gives you back the idea

Every classifier you'll build will tend to the kind memory the training set, it will suddenly do that very well. You won't help it generalize to new data

Take a small subset of the training set. Don't use it in training. Then measure the error on that test data

Training a classifier is wonderful a process of trial and error.

Train a classifier, measure its performance. Then you try another one and measure again. And repeat by tweaking the model, exploring the parameter, measure until you have a nice classifier.

Then you'll deploy your system in a real production environment and you get more data, and you want your performance on that old data to do nearly as well.

Your classifier has memorized the test data through your decision about which classifier to use or which parameter to tune you some information to your classifier about the test set.

Take another chunk of your training set and never use it until you have made your final decision. The validation set to measure your actual error

The bigger your test set, the less noisy your accuracy measurement will be.

A change that affects 30 problems of your validation set, it's hardly significant and can be trusted.

Validation set size > 30 000

Model accuracy figures given relevant to the 10% data given you. enough reduction to see small improvements in longer tests.

If some important changes in data, if you have a small training set, cross-validation is a way to mitigate this problem, but it can be slow so getting more data is often better

The problem with scaling gradient descent is that you need $-\alpha \nabla f(w)$. If $f = \sum_i f_i$, takes n floating point operations, computing its gradient ($\rightarrow \nabla f(w) \rightarrow$) takes three times the compute. It depends on every element in the training set. This can be a lot of compute.

Since gradient descent is iterative, instead you have to compute that hundreds of times.

Instead of computing the loss we're going to compute an estimate of it (a noisy loss).

The estimate is computing the average loss for a small random fraction of the training data. Between 1 and 1000 training samples isn't random enough if the way you pick samples isn't random enough it no longer works.

So compute loss and derivative for that sample and update what the result is the right direction to move to do

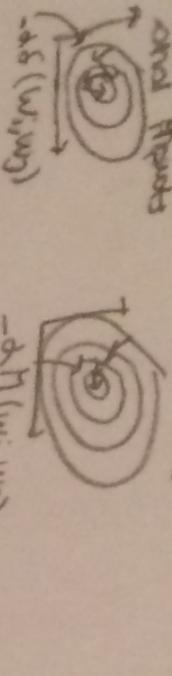
gradient descent

It is not the right direction and at times it might increase the real loss. But we're comparing by doing this many times. taking small steps each will be cheaper to compute. We take many more steps.

This is vastly more efficient than doing gradient descent. This is called stochastic gradient descent and it scales well with data and model size

Helping SAD you want 0 mean and lower the well to help SAD you want 0 mean and lower the variance for initial weights (picked at random)

but we could do want the model to be non-linear



8
WONDER
RUNNING
DANCE 26
1 1 9

For personal or
business model or web train

lower body

而其子也。故曰：「吾子」也。故曰：「吾子」也。

ADMITTED to the care of Dr. G. W. DODGE, superintendant

Deep Neural Networks

$$P \rightarrow X \rightarrow \Sigma \rightarrow \mathcal{G} \rightarrow Y \rightarrow S(X) \rightarrow \dots$$

$(N+1)^K$ PARAMETERS FOR THE MODEL, WHICH WOULD TAKE INTO ACCOUNT THE COMPLEXITY OF THE SYSTEM.

Also it's linear with matrix

Shows what GPUs are for

Linear operation: The results shows small changes in the input ion never result in big changes in the output. The derivative of a linear function is constant.

Diagram illustrating the relationship between a state transition diagram and a Mealy machine.

The top part shows a state transition diagram with three states: S_0 , S_1 , and S_2 . Transitions are labeled with inputs X and Y , leading to outputs Z and W respectively.

The bottom part shows a Mealy machine with two states. Each state has a self-loop transition. The top state has an input X leading to output Z and a transition to the bottom state. The bottom state has an input Y leading to output W and a transition back to the top state.

A large bracket connects the two diagrams, indicating they represent the same system.

```

graph TD
    input["input  
x"] --> sqrt["sqrt  
x"]
    sqrt --> minus1["minus 1  
- 1"]
    minus1 --> mult["mult  
* 8"]
    mult --> output["output  
y"]

```

chain rule!

Preston

人→□→□→□→X

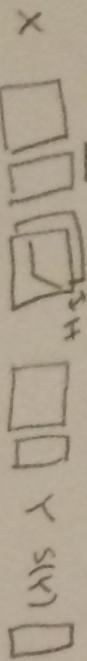
卷之三

To compute the derivatives

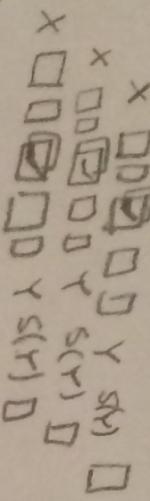
der & hot flow
out combined into
out from in.
hot open on to

Individual operations
in joint network
and founders will

Training a deep neural network



You can make this neural network more complex by increasing H . But doing this is not really efficient in general. One it gets really hard to train. Located you can add more layers and make your model deeper.



You can typically get much more performance by going deeper rather than wider.

Also, many natural phenomena have models with hierarchical structures which deep networks capture.

In order for images humans have lines up the objects longer and do many fine geometric steps. Complicated things further up the objects.

The model matches the abstraction you expect to see in your data and less on lower level features.

Regularization

When doing numerical optimization the network needs to find the right size for your data is very hard to optimize, so we always try networks that are way to big for our data, and then we try out best to prevent them from overfitting (leaving a zero problem)

Validation performance vs training time

With increasing training time

With increasing validation performance

With decreasing validation performance

With increasing validation time

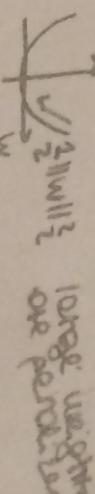
Another way is to apply regularization.

Regularizing means applying a fixed constraint on your network that implicitly reduces the number of free parameters while not making it more difficult to optimize.

Overfitting is when fit just on the training data, but doesn't generalize well to new data.

L2 regularization

$$f' = f + \beta \frac{1}{2} \|w\|_2^2$$



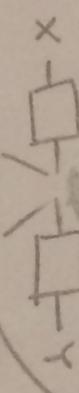
Because it's added to the structure of your network, didn't have to change.

$$\beta \frac{1}{2} \|w\|_2^2 \rightarrow \frac{1}{2} (w_1^2 + w_2^2 + w_3^2 + w_4^2)$$

whose derivative is w_j

Dropout

values that go from 0 (lower) to 1 (higher).



Random!

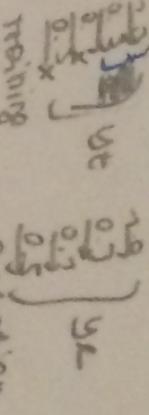
Set half of them to 0
so each sample
only sees half your network

The result will never rely on any given activation to be present. So it is forced to learn a distributed representation for everything to make sure that at least some of the information remain.

These things make it robust and prevents overfitting.

If also makes your network act as a consensus over an ensemble of networks

When evaluating the network trained with dropouts you longer want that random. To fix the consensus over these models by averaging the activations



$\hat{y}_t \sim E(y_t)$
to make sure the hidden state the remaining activation

be a factor of α^2 .

This was when backpropagation, you just remove these dropouts and leaving operations

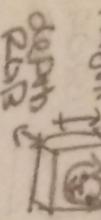
If your network doesn't have to work with the structure of your data from scratch, it is going to perform better

If you want your network to see an image has a cat in it, you can tell it that objects in images are largely the same whether they're in any position (translation invariance)

We get sharing when you know that the two inputs contain the same kind of information (thus you share their weights) only then the weights are the same. In moments like this, shared weights don't change on average over time or space.

(Convolutional Networks (or卷积网))
store parameters over space

height \downarrow width \downarrow



Given an image, it has several not unique outputs on it and didn't depend on it.

depth \downarrow width \downarrow

A convolution
The output is another image with different height, width, depth if your patch uses the size of 3×3 instead of 1×1 image it wouldn't be different but a lower or smaller size of kernel means fewer weights stored

A smaller patch means fewer weights stored

A layer will be a neural network which instead of stacks of matrix multiplying vectors will have stacks of convolutions.

It will form a pyramid of smaller but thinner layers going up

$256 \times 256 \times 3 \times 3 \times 16 \times 16 \times 16 \times 16$ units in total

The depth we're going to do the semantic complexity you've been squashed out and only

parameters that map to content of the image remain

- patch is good called feature map - patch is a patch in a feature map

if you're moving pixels you're not each time you move your filter. When 3×3 the output is roughly a third of the input.

Using padding - when you don't go out of the edge and pad

IP you go off the edge and pad with 0s such that the output image is the same size as the input image. It's zero padding

In all cases the depth is the same

An improvement is to bottleneck to reduce the spatial extent of your feature map through a convolutional pyramid method called pooling

Striding is very aggressive - it removes a lot of information

We can use a small stride and then take all the convolutions in a neighbourhood and combine them together

That is called pooling.

The most common is taking the max of all the neighbourhood

Max pooling (lower values) - Max parameters added

Pyramidal architecture is now pooling - for performing max pooling - not learning model

- learn more accurate (lower values)

Pyramidal architecture is now learning model

INTE \rightarrow CNN \rightarrow PReLU \rightarrow LayerNorm \rightarrow Comp

Image padding - providing a blurred view of the image

in convolution, taking a 3×3 image and 1 pixel

min network instead of a linear layer

Very inexpensive way to make your model deeper

Note your max parameters did not make

inception module

Deep model for text

often we different words to mean almost
the same thing

KITTY CAT when thinks are
similar will like
we have to learn to share parameters

that they are related between them
we would also like to we understand
between enough that we understand
from the meaning

that would require a lot of labeling.

We use unsupervised learning instead

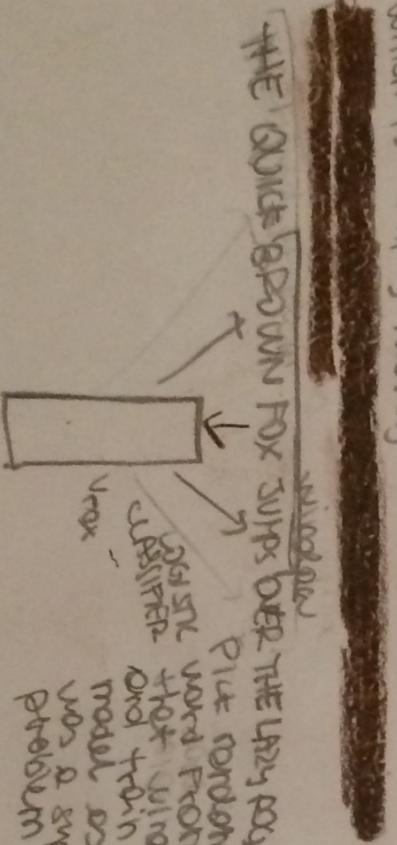
similar words tend to occur in similar contexts

THE CAT PURRS IT'S MEOWABLE
THIS CAT HUNTS WHILE TO USE KITTY

We use words to small vectors called embeddings
which are going to be close to each other
when words have similar meanings and apart
when they don't

Word2Vec

Given text, map each word to an embedding
initially a random one, and the user that
embeddings to the word predict context
which is simply nearby words



To see that the embeddings are clustering together we can
see we can do a scatter without looking at the words closest to
any given word

Proper way is to reduce the dimensionality of the embedding to two
dimensions and to plug the two dimensional representation
The naive way less local words too much information
t-SNE preserves the neighborhoods in input data

It's better to measure closeness with cosine instead of L2 as the
length of the vector is not relevant to the semantic values

JOAINE

VCA VTTEN
WATI MATTEN

it's good to normalize
embedding vectors, so
that norm (VCA+
WATI) / MATTEN

The softmax is really inefficient if long for every iteration
but you can compute the probabilities but you can compute the probabilities
but as if the original vector had been divided by them and then
normed out not there (scaled softmax)

Semantic Analogies

UPPER - UGLY + UGLY = UGLYNESS
FATHER - STOLE / STOLE + KIDNAP - KIDNAP

Recurrent Networks

Let's say the state s_t depends on the previous state s_{t-1} and the input x_t
and a function \rightarrow $s_t = f(s_{t-1}, x_t)$
or what time t depends on what time $t-1$

On has the same
classifier w

To compute the parameter gradients we backpropagate
through time to s_0 (beginning)

This gradient is the gradient to be propagated to
the form parameters (a lot of untrained weights)

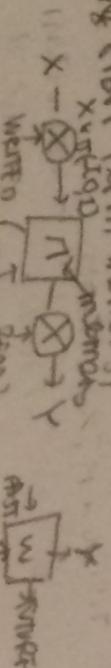
This is bad if set that never get updated weights
either the gradients slow or potentially go down to 0 really quickly (exploding
gradient) or exploding gradients

gradient clipping for exploding gradients

clip + Aw
max(1, Aw) / Aw

→

LSTM for vanishing gradients
(long short term memory)



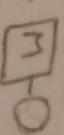
activation \rightarrow forget
 \rightarrow $\text{tanh}(Wx + b)$

\rightarrow $\text{tanh}(Wx + b)$ what need to lost

multiplies at each gate. $\beta = 0$ mean no input given

σ is continuous, differentiable, backprop possible

the softmax values gets controlled by a tiny logistic regression on the input



The softmax help the model to keep outputs within [0, 1]
the softmax longer when it needs to

regularization with L2 norm

drop out with L2 norm

input adapt

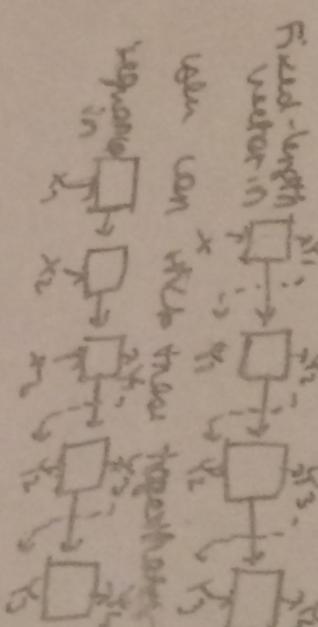
you can also run to predict next ties in your sequence
for that you take the sequence up to time t ,
generate the predictions and then you
compute from the predicted distribution, feed
that compute to the next step

just think so the next prediction is very steady.
but this will just compute more time, then you
can make requirement (at hypothesis) that
you could continue predicting from already this
by computing the total probability out of those
characters generated

that way you steered the sampling from many
one best choice and kept stuck with that

beam search. beam search only take the most
likely few candidates at every step and pruning next

RNN's not able length sequences to fixed length vectors
themselves to sequence generation stream then go the other way



outputs

whole length sequence out