

Module 1:

Inductive reasoning in DL

feature engineering → to lower down the need for
using all alter-
(color history etc.)

if we use softmax convert scores to probabilities,
the right loss function to use would be

Cross-entropy

tensor → multi-dimensional matrices

each instance is a vector size of m, batch is
of size $[B \times m]$

each instance is a matrix of size $W \times H$, our
batch is $[B \times W \times H]$

each instance is a multi-channel matrix of size
 $C \times W \times H$ our batch is $[B \times C \times W \times H]$

Deep Learning

Since deep learning takes in inputs as raw with no dimensionality reduction and learns feature representation for that data.

- Representation Learning
- Neural Networks

(Hierarchical) Compositionality:

cascade of non-linear transformations
multiple layers of representations

end-to-end learning
learning (goal-driven) representations
learning feature extraction

distributed representations
no single neuron "encodes" everything
groups of neurons well together

learn new features instead of being told what the best feature is.

compose into a complex function

Hierachical:

pic of a car → low-level feature

end-to-end learning

Shallow - learning.

pic of car

+ fixed
hand crafted
feature
extractor

+ learn
Simple
trainable
classifier

low-level feature



mid-level feature



high-level feature



trainable classifier



"car"

deep models:

pic of car



trained
feature
transform

trainable
feature
transform
classifier



layered
interval
representations

ML2:

linear classifier
input \rightarrow function of input \rightarrow loss function

input layer hidden layer output layer
 x w_1 w_2

$$f(x, w_1, w_2) = g(w_2 g(w_1 \cdot x))$$

\nwarrow sigmoid

Three layer NN can hypothetically represent any function.

two layer NN can hypothetically represent any continuous function.

Differentiability is important b/c of gradient descent.

backprop

the outputs of the model is called forward
calculate the gradients for each module
(backward pass)

backward pass → start at loss pass back through
the layers
the modules and end at input layer
(since it has no parameters)

$l_1 \rightarrow l_2 \rightarrow l_3$ forward pass

$$\frac{\partial L}{\partial h^{l-1}} \quad \boxed{\text{ }} \quad \frac{\partial L}{\partial h^l}$$



$$\frac{\partial L}{\partial w}$$

backward pass

$l_3 \leftarrow l_2 \leftarrow l_1$

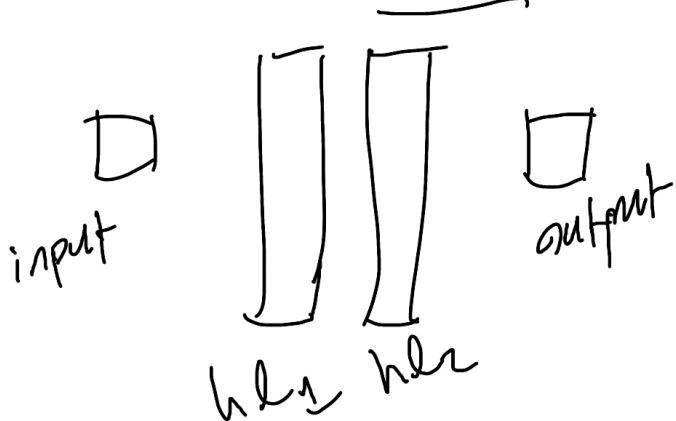
(you also need to
store intermediate
outputs of
all layers)

Since we might
need it for
gradient
calculator

R&LV has better gradient flow than sigmoid. Performed count-wise

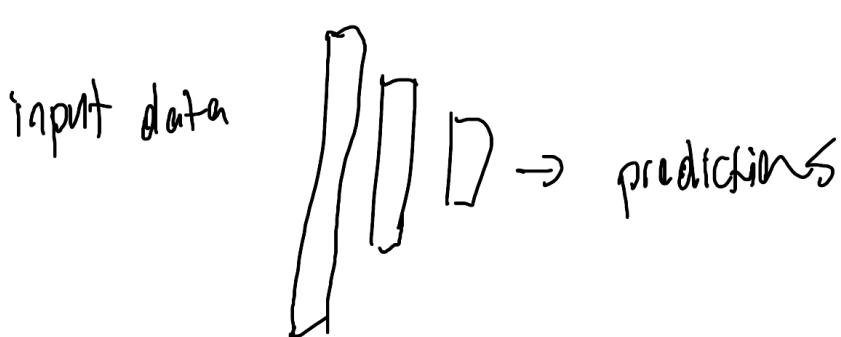
ML3:

a network with two or more hidden layers is considered deep



depth brings
parameter
efficiency
dimensionality
regularization

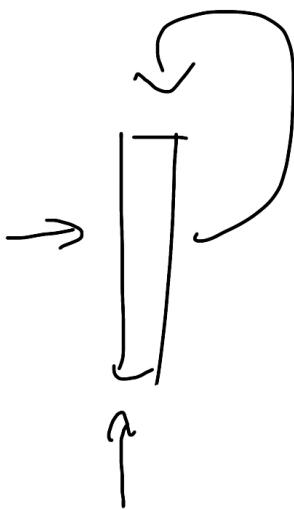
fully connected neural network





convolutional neural

networks



designing architecture:

understanding your data is the
first step.

vanishing gradients \Rightarrow
when gradients
become flat -

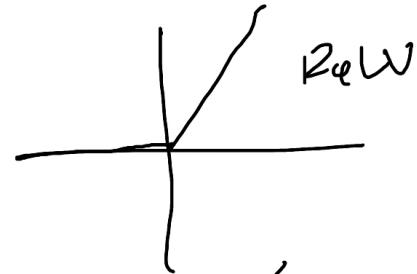
tanh \rightarrow central / balanced.

features at the end of its
derivatives are somewhat

computationally heavy.

ReLU
not central like tanh
but no saturation on positive end.

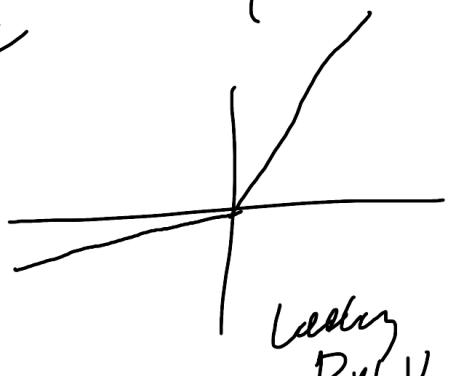
Leaky ReLU:



Learnable parameters

no saturation on negative side

no dead neurons



initializing values to
constant values leads to
an degenerate solution

- gradient will be the same
- all weights will be updated to same value.

common approach is small normally distributed random numbers

the deeper the network, the outputs of nodes become smaller, leading to larger standard deviations.

leads to small updates; layer lead to saturation.

Xavier Initialization
simple initialization rule:

Uniform $\left[-\frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}} \right]$

$n_j \rightarrow$ number of input nodes

$n_{j+1} \rightarrow$ number of output nodes

empirically well simpler version

$$N(0,1) \cdot \sqrt{\frac{1}{n_j}}$$

works for tanh, similar activations

ReLU:

(Xavier rule)

$$N(0,1) \cdot \sqrt{\frac{1}{n_j/2}}$$

deep learning involves a lot of complex, compositional, non-linear functions. so less landscape it needs
a lot of intuition & rule of thumb used.

a good way to effectively use the info is to take data in batches
but this causes issues of high variance
but it's unbiased. and has using descent.

dynamic learning rate

SGD can achieve similar results to RMSProp

Acknowledgment

Adm

L1 \rightarrow encourages sparsity,
 $\hookrightarrow |y - w_{x_i}|^2 + \lambda \cdot |w|$

$$L2 = \hookrightarrow |y - w_{x_i}|^2 + \lambda \cdot |w|^2$$

L1/L2 on weights encourage small values.

To prevent NN from overfitting
keep a D.T probability for each model
 \hookrightarrow allows you to check whether marking out a node
helps.

in practice \rightarrow it's a mask during test none
is dropped-out. But you scale output by p.

$$W_{\text{test}} : p \cdot W$$

if let's say the "important features" are selected
in dropout and the weights are updated in
backprop, this technique will allow you to
not rely on one particular feature

data augmentation \rightarrow performing a range of
transformations. It shouldn't change the meaning of
the data though.

for example, image flipping, random cropping;
color jitter

validation \rightarrow it's for understanding architecture
and what hyperparameter to use.

cross-validation \rightarrow if you don't have a lot
of data, a bit computationally expensive

check initial loss at small random weights
 $-\log(p)$ for cross entropy. where $p: 0.5$

or another example is to not do regularization
and make sure loss goes up when added.

you can also fit with a small dataset and see if
your model can overfit this small amount of data
before regularization.

tiny learning rate \rightarrow too small of a learning rate

overfitting \rightarrow validation loss starts to get worse
after a while

underfitting \rightarrow validation loss very close to
training loss or both are high

Validation loss has no regularization

hyper-parameters:

learning rate, weight decay

Momentum, others more stable

batch norm & dropout not needed together.

learning rate should be scaled proportionally
to batch size (higher batch size higher
learning rate)

TPR / FPR

tradeoff between the number of positive
predictions and correctness of predictions -

TPR : $\frac{tp}{tp + fn}$

FPR : $\frac{fp}{fp + tn}$

accuracy : $\frac{(tp + tn)}{(tp + tn + fn + fp)}$

Missing data imputation

for numerical data:

mean, mode, most frequent, zero, constant

for categorical data:

hot-deck imputation, k-NN, deep learned
embeddings