# Vanishing and Exploding Gradients
# Gradient Check

Deep Neural Networks

Session 11

Pramod Sharma
pramod.sharma@prasami.com

Neural Networks face problems of vanishing or exploding gradients

pra-sami

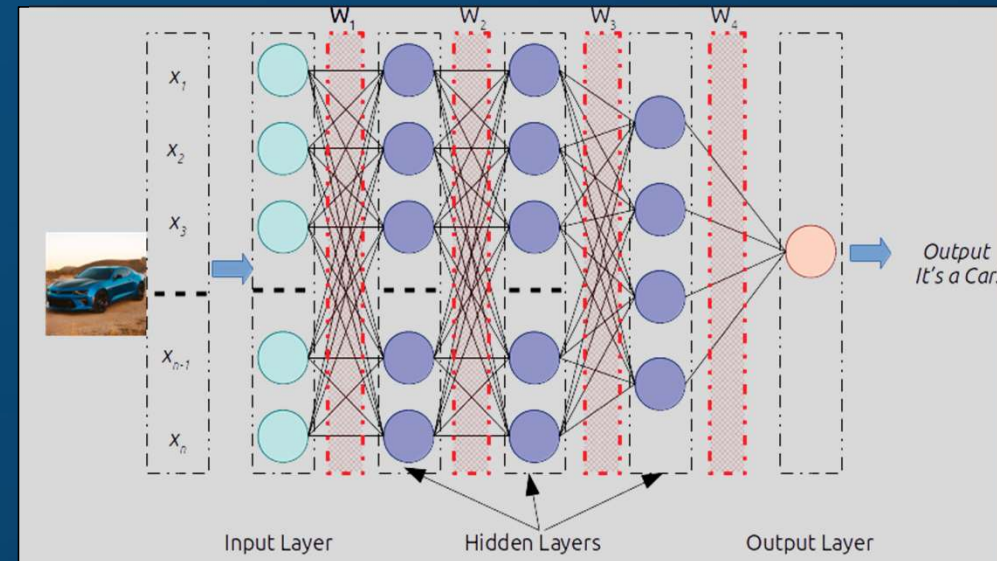Deeper the network more are the chances of the gradients becoming smaller and smaller or keep growing…

pra-sami

# Weights multiplied…

□ We know that :

- ❖ $z = X * W + b$
- ❖ $\hat{y} = a = \sigma(z)$
- ❖ $a_1 = \sigma(a_0 . W_1)$

□ That in multilayer network

- ❖ $\hat{y} = \sigma(\sigma(\sigma(\sigma(\sigma(a_0 . W_1).W_2).W_3).W_4)....$
- ❖ For explanation purpose assume σ (z) = z (say ReLU)
- ❖ $y = W_1 . W_2 . W_3 . W_4$
- ❖ so any change in y will result in $W_1 * W_2 * W_3 * W_4$ times y in layer 1
- ❖ Longer the chain, more $W_s$ will be multiplied.



Input Layer    Hidden Layers    Output Layer

Output
It's a Car!

pra-sami

# More Layers... More problems

❑ Assume we have 150 layers

❑ Also assume our weight is say 1.1

=> $1.1^{150}$ = 1.6 million

❑ On the other hand assume our weight is 0.9
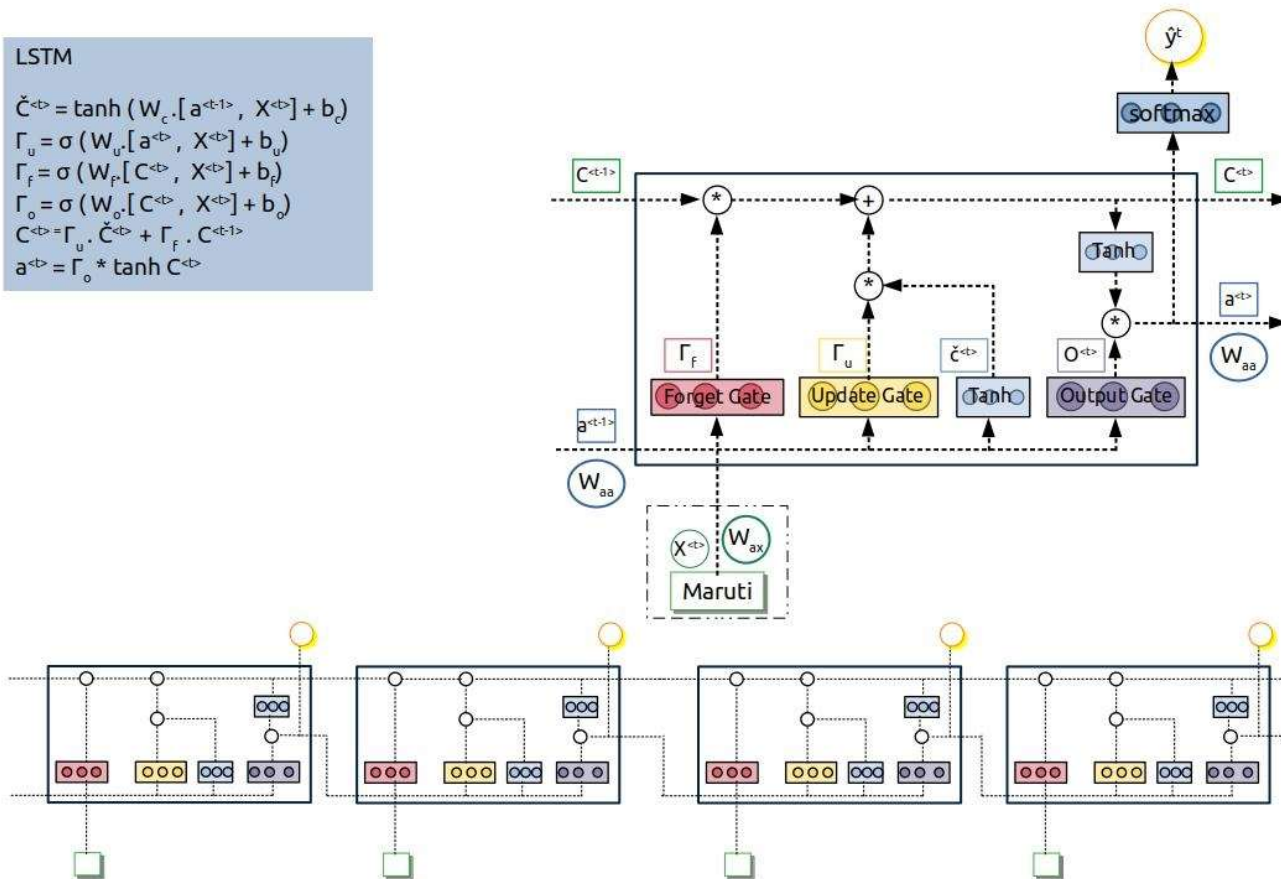
=> $0.9^{150}$ = 1.4 $e^{-7}$

pra-sami

# It's a Severe Problem…

- ❑ No silver bullet solution….

- ❑ There is multi-prong approach to it…

- ❑ First, Initialise your weights as close to 1 as possible ( not 1)

- ❑ It is found that for tanh activation function
  - ❖ Divide by $\sqrt{number\ of\ nodes\ in\ the\ previous\ layer}$

- ❑ Some cases : $\dfrac{2}{\sqrt{number\ of\ nodes\ in\ the\ previous\ layer}}$

- ❑ In ReLU , : $\dfrac{2}{\sqrt{number\ of\ nodes\ in\ the\ previous\ layer + number\ of\ nodes\ in\ current\ layer}}$

- ❑ Some literature suggest, even $\dfrac{K}{\sqrt{number\ of\ nodes\ in\ the\ previous\ layer}}$ ; κ is a another parameter to tune
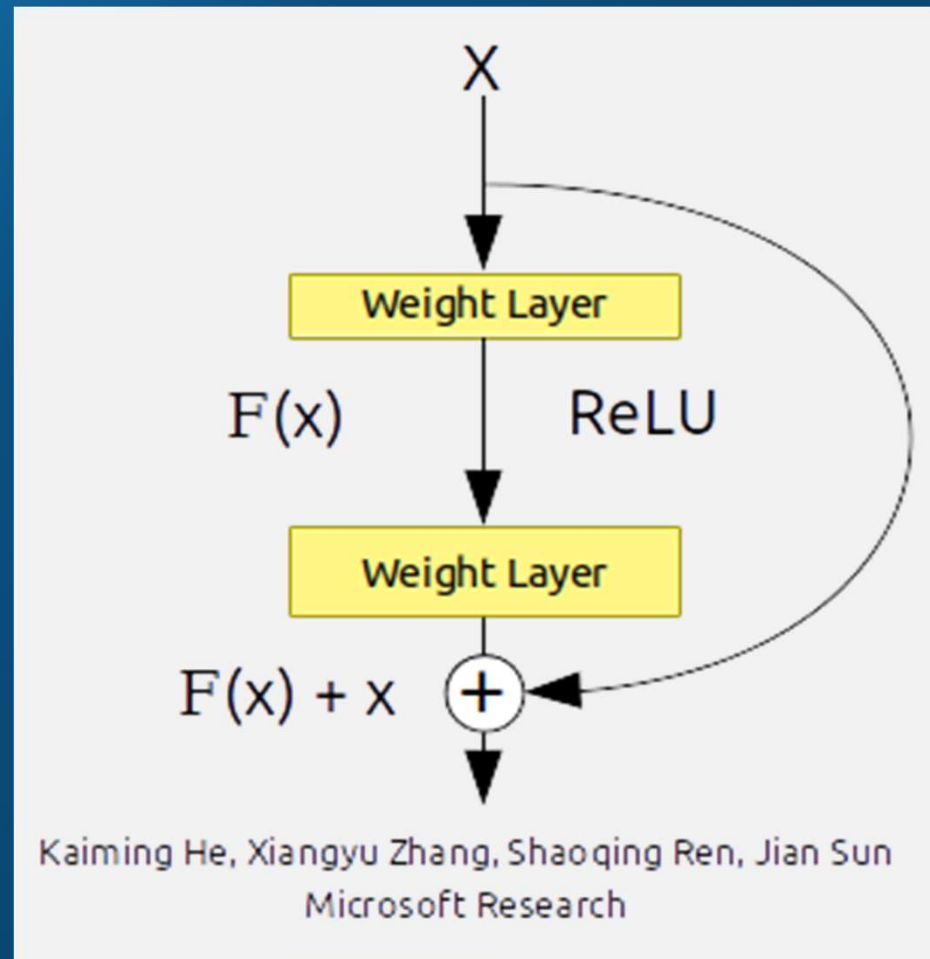
pra-sâmi

# Change Your Architecture - LSTM

# As in ResNet



Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun
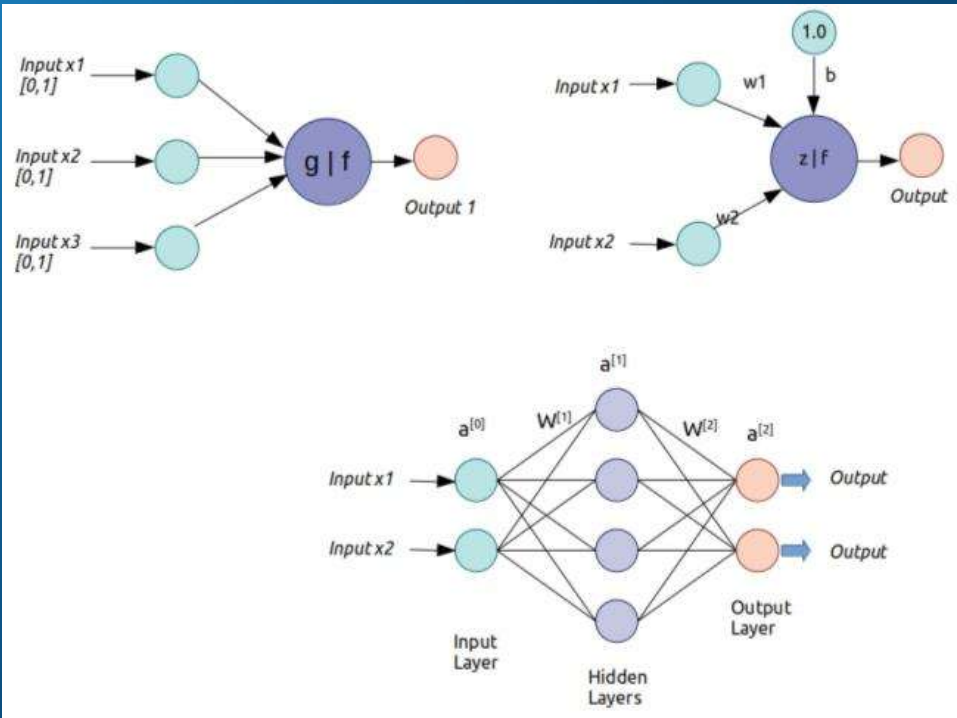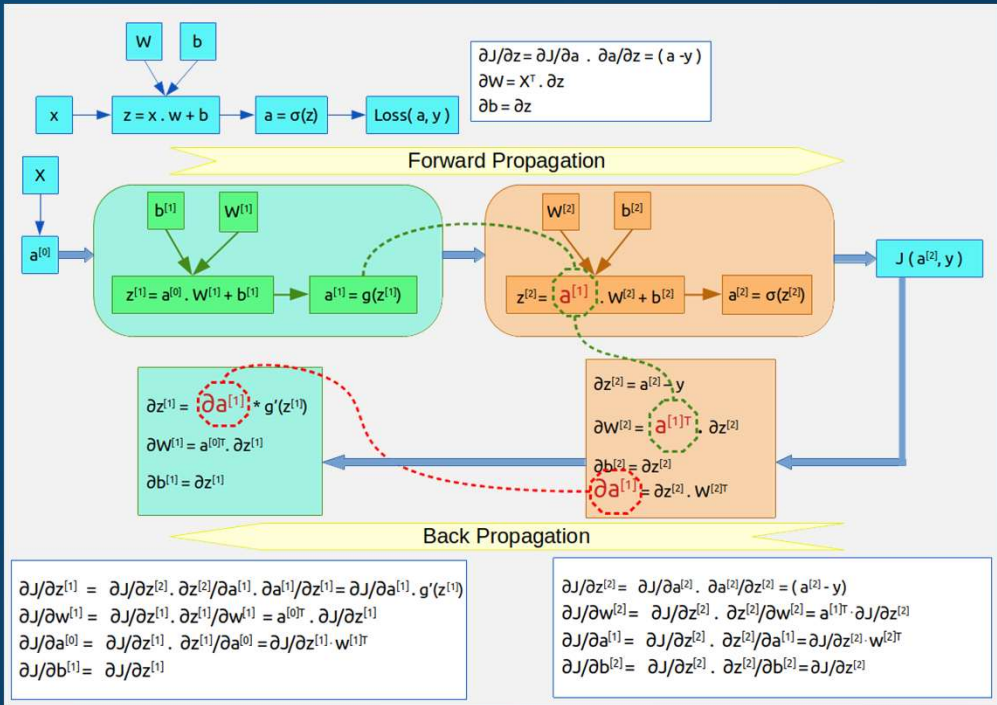Microsoft Research

# Gradient Check

pra-sami

# Gradient Checks

❏ Recap



❏ The math

pra-sami

# Loss Function and its Derivative

❑ The loss for our prediction $\hat{y}$ with respect to the true labels y is given by:

$$L(\hat{y}, y) = -y \cdot \log \hat{y} - (1-y) \cdot \log(1-\hat{y})$$

❑ For all samples:

$$J(\hat{y}, y) = -\frac{1}{m} \sum_{i \in m} y_i \log \hat{y}_i - (1-y_i) \cdot \log(1-\hat{y}_i)$$

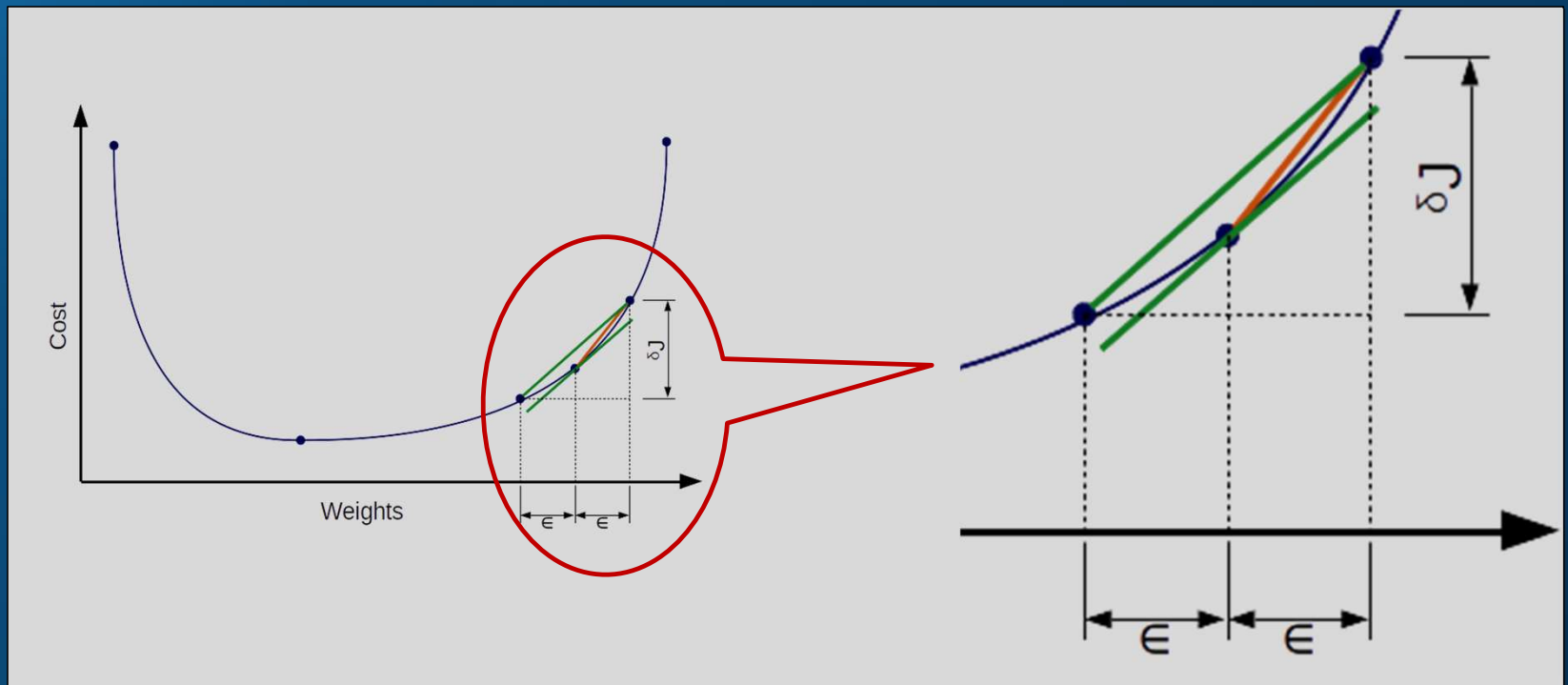Where $\hat{y} = \sigma(a \cdot W + b)$

❑ Therefore, we can say that:

$$J(\hat{y}, y) = J(W, b) = J(W_1, W_2, W_3, \ldots W_n, b_1, b_2, b_3, \ldots b_n)$$

pra-sami

# Calculation of derivative

❑ Use the centered formula

  ❖ The formula you may have seen for the finite difference approximation when evaluating the numerical gradient is not as good as centered formula

# Gradient Checking

- ❑ Also called "Grad Check"
- ❑ Do it to verify the model's math(Debug) only
  - ❖ Too heavy for training, switch off once the model is verified.
- ❑ For all values of Ws and Bs, we can calculate:

$$\delta\theta_{approx} = J(W_1, ..., W_i + \varepsilon, .... W_n, b_1, b_2, b_3, .... b_n) - J(W_1, ..., W_i - \varepsilon, .... W_n, b_1, b_2, b_3, .... b_n) / (2 * \varepsilon)$$

- ❑ To check if $\delta\theta_{approx}$ and $\delta\theta$ are close

$$\frac{\|\delta\theta_{approx} - \delta\theta\|_2}{\|\delta\theta_{approx}\|_2 + \|\delta\theta\|_2} \text{ is very small}$$

- ❑ For $\varepsilon$ = 1e-7
  - ❖ Relative error > 1e-2 usually means the gradient is probably wrong
  - ❖ 1e-2 > relative error > 1e-4 should make you feel uncomfortable
  - ❖ 1e-4 > relative error is usually okay for objectives with kinks. But if there are no kinks (e.g. use of tanh nonlinearities and softmax), then 1e-4 is too high.
  - ❖ 1e-7 and less you should be happy

pra-sâmi

# Grad Check Steps

❑ Recall: Our model has all weights and biases stored

- ❖ Model ={ "$W_1$": … , "$b_1$": …, "$W_2$": … , "$b_2$": …, … … "$W_n$": … , "$b_n$": …}
- ❖ We have implemented our forward prop and back prop

❑ Step 1 : Pick model and convert all weights and biases into a vector $\theta$

❑ Step 2: Similarly pick $\delta$W and $\delta$b and convert to a vector $\delta\theta$

❑ Step 3: for each of the value in the vector $\theta$

- ❖ Make copy of $\theta$ and $\delta\theta$
- ❖ Increase $\theta_i$ to $\theta_i + \varepsilon$
- ❖ Calculate J + ( Cost with increased $\theta$)
- ❖ Similarly calculate J- ( Cost with decreased $\theta$)
- ❖ Use J+ and J- to calculate $\delta\theta_{approx}$
- ❖ Calculate $\delta\theta$ as usual
- ❖ Find error

pra-sami

# Remember to Turn off Dropout/Augmentations

❑ When performing gradient check, remember to turn off any non-deterministic effects in the network, such as dropout, random data augmentations, etc.

❑ Otherwise these can clearly introduce huge errors when estimating the numerical gradient

❑ Downside ➜ you wouldn't be gradient checking them

  ❖ It might be that dropout isn't back-propagated correctly)

pra-sami

# Use Double Precision

❑ A common pitfall is using single precision floating point to compute gradient check. It is often that case that you might get high relative errors (as high as 1e-2) even with a correct gradient implementation

❑ There are cases when relative errors plummet from 1e-2 to 1e-8 by switching to double precision

pra-sami

# Stick Around Active Range of Floating Point

- What Every Computer Scientist Should Know About Floating-Point Arithmetic https://docs.oracle.com/cd/E19957-01/806-3568/ncg_goldberg.html

- In neural nets it is common to normalize the loss function over the batch

- If gradients per datapoint are very small, then *additionally* dividing them by the number of data points starts to give very small numbers

- Always print the raw numerical/analytic gradient,
  - ❖ Make sure that the numbers you are comparing are not extremely small (e.g. roughly 1e-10 and smaller in absolute value
  - ❖ If so, you may want to temporarily scale your loss function up by a constant to bring them to a "nicer" range where floats are more dense - ideally on the order of 1.0, where your float exponent is 0.

pra-sami

# Kinks in the Objective Function

❑ One source of inaccuracy to be aware of during gradient checking is the problem of *kinks*

❑ Kinks refer to non-differentiable parts of an objective function
  ❖ Such as ReLU max(0,x), or the SVM loss, Maxout neurons, etc.

❑ Consider gradient checking the ReLU function at x=−1e-6

❑ Since x<0, the analytic gradient at this point is exactly zero.

❑ However, the numerical gradient would suddenly compute a non-zero gradient because f( x + h ) might cross over the kink (e.g. if h>1e−6) and introduce a non-zero contribution.

❑ You might think that this is a pathological case, but in fact this case can be very common.

❑ Keeping track of the identities of all "winners" in a function of form max( x, y ); That is, was x or y higher during the forward pass. If the identity of at least one winner changes when evaluating f( x + h) and then f( x − h ), then a kink was crossed and the numerical gradient will not be exact

pra-sami

# Grad Check

❑ Phew… too lengthy calculations….

❑ Good News!

  ❖ Both **Tensorflow** and **Torch** have **autograd** implementation for us. So in real implementation we will be using those functions for gradient checking

❑ Caution!

  ❖ This check is resource hungry

  ❖ Once the verification is done, comment/switch off the code

❑ Deeper the network ➜ the higher the relative errors

  ❖ For the input data for a 10-layer network, a relative error of 1e-2 might be okay because the errors build up on the way

  ❖ Conversely, an error of 1e-2 for a single differentiable function likely indicates incorrect gradient

pra-sami

Next Session…
The process of learning the parameters…
Finding good hyper-parameters…