

5.4.2 Batch normalization

Input normalization : done once for the dataset

Hidden layers : data distribution changes due to the update of $\underline{\theta}$

- over layers
 - over time during training
- } internal covariate shift

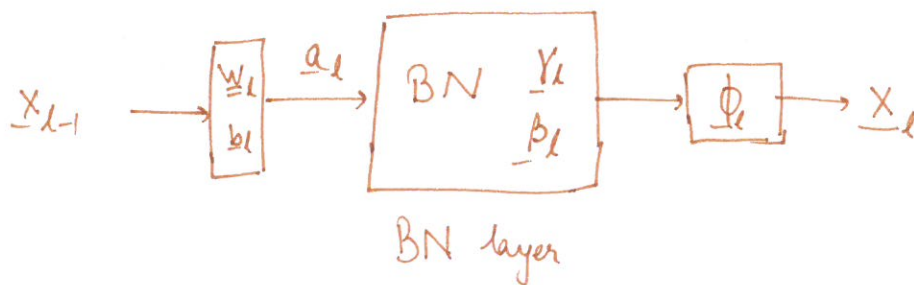
→ slow down training if nothing is done

Soln:

Batch normalization (BN):

like Input normalization, but

- for hidden layers (we have a choice on which layer to apply)
- for each mini batch → stabilize training



Slide S20

S21

Q
Ilias

If we first do a zero mean unit variance scaling and then again boost the features then we still obtain the dynamic range with a different μ and σ . So what is the use of using an affine transform of the activation o/p?

Q

Mathematically, we can write this whole term as one term in place of \underline{w} and \underline{b} , so how does it make a difference?

5.5

Parameter initialization

SGD does a local search, so we need a good starting point. The solⁿ depends on the initial values

$$\underline{\theta^0} = \underline{w_0}$$

Different parameter ~~it~~ initializations :

Slide 5-24

2. Random initialization:

Only done for $\underline{w_1}$. Bias are set to 0. eg

• normal distribution $N()$

$$[\underline{w_1}]_{ij} \sim \text{iid } \underset{\uparrow}{\delta} \cdot N(0, 1) = N(0, \delta^2)$$

to control the data dynamic range.

This leads to asymmetry.

• uniform distribution $U()$

$$[\underline{w_1}]_{ij} \sim \text{iid } \delta U(-1, 1) = U(-\delta, \delta)$$

$\delta = \text{const } \forall i, j, l \rightarrow \text{not optimum}$

3.

He initialization

③

as 2), but $\sigma_l \sim \frac{1}{\sqrt{M_{l-1}}} \rightarrow$ constant
activation flow after
initialization

layer l after initialization:

$$\underline{a}_l = \underline{W}_l \cdot \underline{x}_{l-1} + \underline{b}_l \quad \underline{b}_l = \underline{0} \quad \underline{W}_l \cdot \underline{x}_{l-1}$$

or

$$\begin{array}{ccc} \underline{a} & = & \underline{W} \cdot \underline{x} \\ M_{l \times 1} & & M_{l \times M_{l-1}} \quad M_{M_{l-1} \times 1} \end{array}$$

$$a_i = \sum_{j=1}^{M_{l-1}} w_{ij} x_j$$

Assumption:

- X_j iid - zero mean, variance σ_x^2
- w_{ij} iid " " " σ_w^2
- X_j and w_{ij} independent.

$$\Rightarrow \bullet E(a_i) = \sum_j \underbrace{E(w_{ij})}_0 \underbrace{E(x_j)}_0 = 0$$

$$\bullet \text{Var}(a_i) = E(a_i^2) - \underbrace{(E(a_i))^2}_0$$

$$= E \left[\left(\sum w_{ij} x_j \right)^2 \right]$$

$$= E \left[\sum_j \sum_k w_{ij} w_{ik} x_j x_k \right]$$

$$= E \sum_j \sum_k E[w_{ij} w_{ik}] E[x_j x_k] = \sum_{j=1}^{M_{l-1}} \sigma_w^2 \sigma_x^2 = M_{l-1} \sigma_w^2 \sigma_x^2$$

Therefore constant activation flow in fwd pass:

$$\text{Var}(a_i) = \text{const} \quad \forall i, l$$

$$\Rightarrow \delta_{w,l} \sim \frac{1}{\sqrt{M_{l-1}}} \quad M_{l-1}: \text{fan in}$$

4. Glorot initialization

constant gradient flow in backward pass

$$\left\| \frac{\partial L}{\partial a_l} \right\| = \text{const} \quad \forall l$$

$$\Rightarrow \delta_{w,l} \sim \frac{1}{\sqrt{M_l}} \quad M_l: \text{fan-out}$$

$$\text{Compromise: } \delta_{w,l} \sim \frac{1}{\sqrt{M_{l-1} + M_l}}$$

5. 6

Improved model

Model - architecture of NN

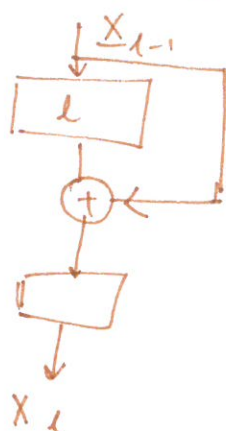
A. A better activations function $\phi(x)$:

- ReLU instead of sigmoid.
- leaky ReLU instead of ReLU as we have small gradient for -ve values.

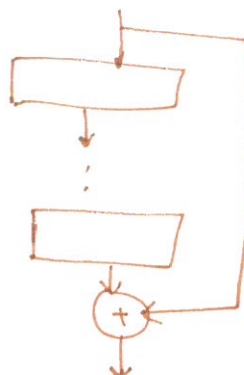
B. Skip-connections, shortcuts, residual n/w

(5)

over one layer



over multiple layer



- forward pass:

Combine low level features of shallow layer with high level features of deep layer

see ch 7, ch 10

- backward pass

$$\underline{x}_l = \phi_l (\underline{w}_l \cdot \underline{x}_{l-1} + \underline{b}_l) + \underline{x}_{l-1}$$

$$\frac{\partial \underline{x}_l}{\partial \underline{x}_{l-1}} = \underbrace{\frac{\partial \phi_l}{\partial \underline{x}_{l-1}}}_{\text{This can have}} + \underline{I}$$

a vanishing gradient

term, but \underline{I} guarantees that propagation will always happen. Hence we need shortcuts.

(?) Many other architecture improvements see ch 10.

6

Overfitting and regularization

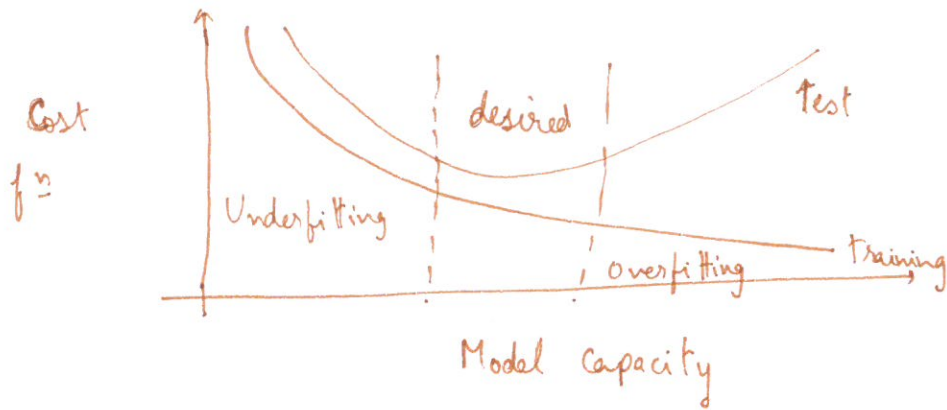
⑥

6.1

Model capacity, underfitting and overfitting

Slide 6-1

Slide 6-2



Slide 6-3