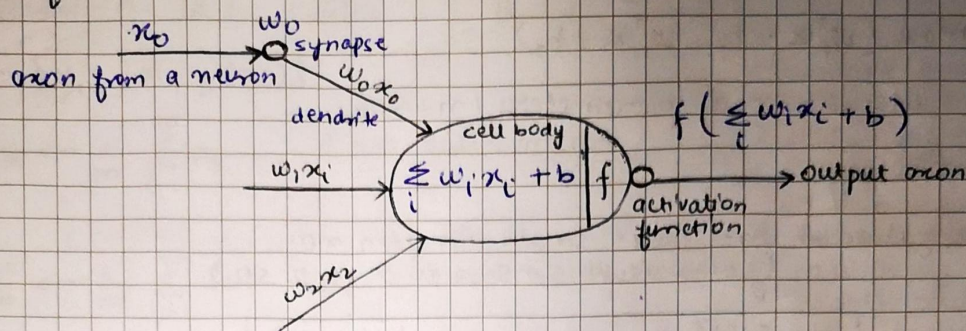


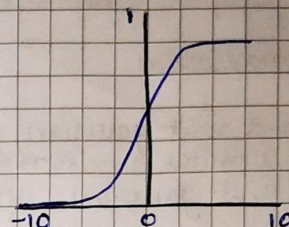
Lecture - 6 Training Neural Networks - I

→ Activation functions



① Sigmoid → $\sigma(x) = \frac{1}{(1 + e^{-x})}$

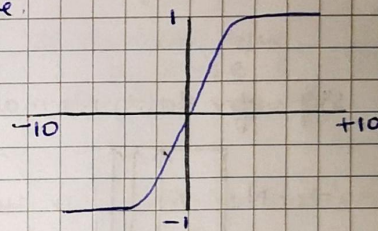
- Squashes no. to range $[0, 1]$
- Historically popular as they have nice interpretation as a saturating "firing rate" of a neuron



- Problems →
- Saturated neurons kill the gradients.
 - Sigmoid outputs are not zero-centered.
 - $\exp()$ is a bit computationally expensive.

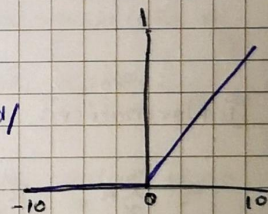
② Tanh(x) → Squashes no. to range $[-1, 1]$

- Zero-centered.
- Still kills the gradient when saturated.



③ ReLU (Rectified Linear Unit) → computes $f(x) = \max(0, x)$

- does not saturates in +ve region
- Computationally very efficient
- Converges much faster than sigmoid/tanh
- more biologically plausible than sigmoid



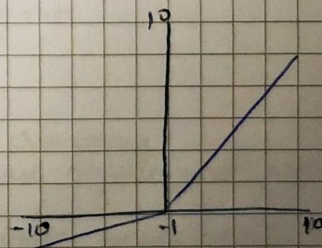
- Problems →
- output not zero-centered
 - gradient when $x < 0$

★ People like to initialize ReLU neurons with slightly positive bias (eg. 0.01)



- ④ Leaky ReLU →
- Does not saturate
 - computationally efficient
 - converges much faster than sigmoid/tanh in practice
 - will not die

$f(x) = \max(0.01x, x)$



- ⑤ max out "Neuron" → Does not have the basic form of dot product → Nonlinearity
 → Generalises ReLU and Leaky ReLU.
 → Linear Regime! Does not saturate! Does not die.

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

Problem → doubles the no. of parameters / neurons.

→ Data Preprocessing

→ Weight initialization → first idea: small random no.
 (Gaussian with 0 mean and $1e-2$ SD)

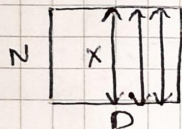
$$W = 0.01 * np.random.randn(D, H)$$

→ Batch Normalisation

"You want unit gaussian activations? Just make them so".
 Consider a batch of activations at some layer. To make each dimension unit gaussian, apply.

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}} \quad \rightarrow \text{This is a vanilla differentiable function.}$$

Steps for Batch Normalisation: ① compute the empirical mean and variance independently for each dimension.



② Normalize with the above formula.

③ usually inserted after fully connected or conv. layers and before non-linearity.

Problem: do we necessarily want a unit gaussian input to a tanh layer?

① Normalise:

$$\hat{x}^k = \frac{x^k - \mathbb{E}[x^k]}{\sqrt{\text{Var}[x^k]}}$$

and then allow the network to squash the range if it wants to.

$$y^k = \gamma^k \hat{x}^k + \beta^k$$

Squashing and scaling operation

$$\beta^k = \mathbb{E}[x^k]$$

$$\gamma^k = \sqrt{\text{Var}[x^k]}$$

It helps to recover the identity mapping.

Summary: Batch Normalization

Input: values of x over a mini-batch: $\beta = \{x_1, \dots, x_m\}$;
Parameter to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad // \text{ normalise}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta = \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift.}$$

- + Improves gradient flow through the network
- + Allows higher learning rate
- + Reduces the strong dependence on initialization
- + Acts as a form of Regularization in a funny way, and slightly reduces the need for a dropout.

→ Babysitting the Learning Process

- * start with the small regularization and find the learning rate that makes the loss go down.
 - * loss not going down \rightarrow learning rate too low.
 - * loss exploding \rightarrow learning rate too high.
- } somewhere b/w 10^{-3} to 10^{-5}

Hyperparameter optimization.

To choose the values of our hyperparameters, we do: "Cross-validation"

strategy. \rightarrow (i) coarse-fine cross-validation in stages

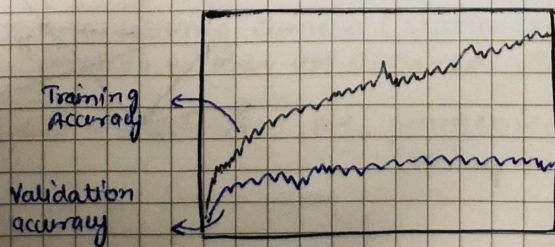
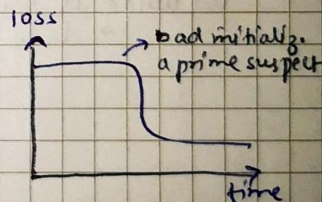
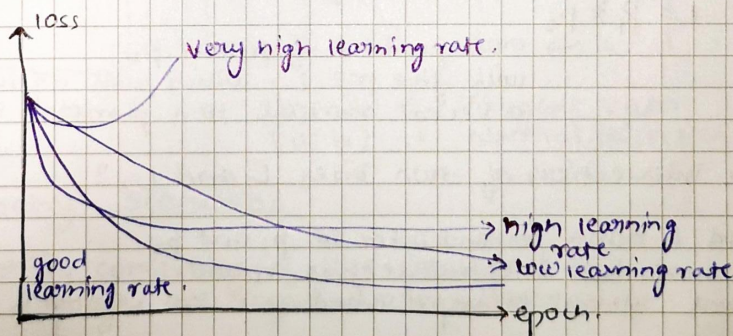
first stage \rightarrow only a few epochs to get a rough idea of what parameters work.

second stage \rightarrow longer running time, finer search.

Hyperparameters to play with \rightarrow Network architecture

- \rightarrow learning rate, its decay schedule, update type
- \rightarrow Regularization

loss curve:



} big gap = overfitting
= increase regularization strength.

No gap = increase model capacity