

② Ada Grad Optimization
(Adaptive Gradient)

$$r_t = r_{t-1} + (\nabla \theta_t)^2, \quad r_0 \geq 0$$

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\delta + \sqrt{r_t}} \nabla \theta_t$$

↑
the number
to avoid divide by 0.

⇒ As time progresses, r_t will get to a large value, hence, the moment does not happen.

③ ²⁰¹⁴ Ada M Optimization
(Adaptive Moments)

Adaptive + Momentum

① Momentum →
$$b_t = \rho_1 b_{t-1} + (1 - \rho_1) \nabla \theta_t \quad \text{--- ① } 0.9, 0.1$$

$$r_t = \rho_2 r_{t-1} + (1 - \rho_2) (\nabla \theta_t)^2 \quad \text{--- ② } 0.9, 0.1$$

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\delta + \sqrt{r_t}} \cdot \tilde{b}_t$$

divide by zero → (8) + √r_t ↑ ② Adaptive Gradient

$$\tilde{r}_t = \frac{r_t}{1 + \rho_2^t}, \quad \tilde{b}_t = \frac{b_t}{1 - \rho_1^t} \quad \text{--- ③ Bias Correction}$$

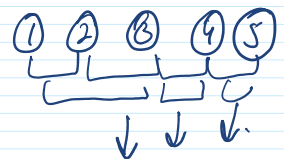
REGULARIZATION

→ Generalized Network.

→ Training Sample (a distribution)



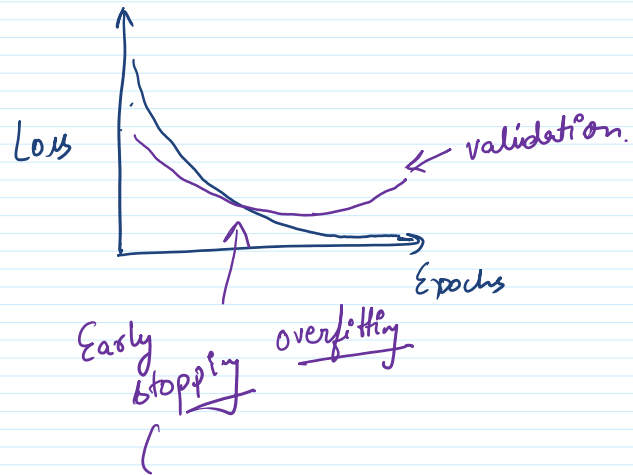
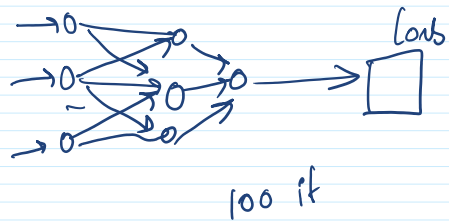
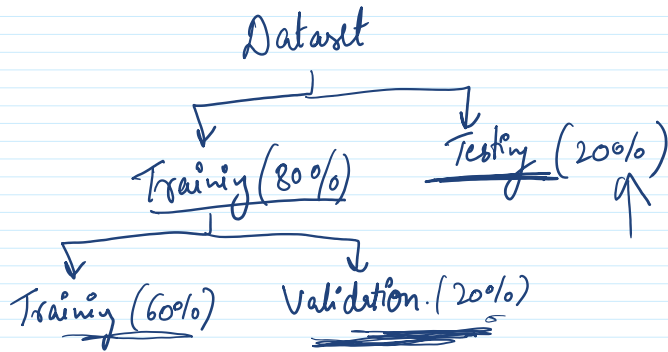
Noisy Training Samples



Testing

Domain-Shifting

Image Net
1.3 Billion.



① Early Stopping

② Data Augmentation

10 class classification.

→ ① Cat → 100 images

→ ② Dog → 100 images

→ ③ | → 100 images

!

1000 images dataset → 1000 X 5 (80%)

→ Geometric Transformations:

Rotation, Translation, Shear, Scale, ...

→ Photometric Transformations:

Noise, Blur, ...

(3) Adding regularization on weights.

$$\tilde{L}(w) = L(w) + \frac{\beta}{2} \|w\|^2.$$

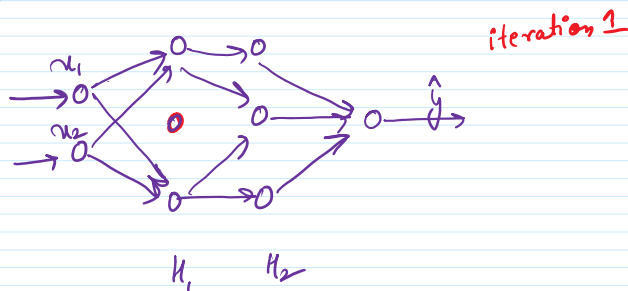
(4) Noise Injection :-
to the input

MSE Loss

$$\begin{aligned} y_{\text{noisy}} &= \sum w_i(x_i + \epsilon_i) \\ &= \sum w_i x_i + \sum w_i \epsilon_i \\ &= \underline{y} + \underbrace{\sum w_i \epsilon_i}_{\text{regularization}} \end{aligned}$$

Let: $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$
↑
Gaussian distribution.

(5) Drop out :-



'p' → probability of picking a neuron.
p = 0.5 50% of nodes of that hidden layer are dropped.

Mini-batch GD.

$$0 \rightarrow 0 \rightarrow 0$$

$$0 \rightarrow 0 \rightarrow 0 \rightarrow \dots$$

$$0 \rightarrow 0 \rightarrow \dots$$

$H_1 \quad H_2$

- ① If only some of the neurons take on the entire load of the task on hand, it can lead to overfitting.
- ② The rest of the neurons do not learn anything.
- ③ Typically, nodes are dropped only once for a minibatch

- of sample.
- ④ Nodes are dropped with a probability of 'p'.
 - ⑤ F/w & b/w pairs are done only through active neurons.

⑥ At the time of testing, weights are multiplied by probability.
weights will reduce by factor 'p'.

Algorithm.	No. of steps in 1 epoch	
Vanilla GD	1	
Stochastic GD	N	N → no. of training samples.
Minibatch.	N/B.	B → batch-size

	Outputs	
	Real Values (Regression)	Probability (Classification)
O/p activation	Linear	Sigmoid/ Softmax
Loss function	MSE	<u>CCE</u> .