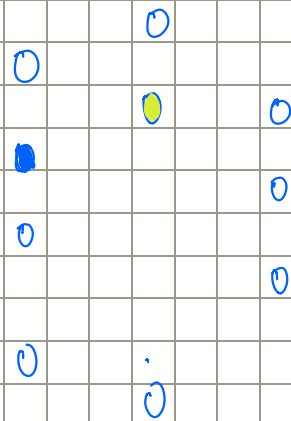


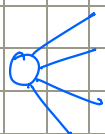
06/10/23

Deep Learning

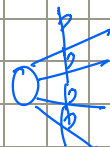
- Recap ✓
- Dropout ✓
- Weight Decay ✓
- ResNet ✓



- Dropout: Training
- For a given "thinned" network carry out SGD
 - Choose a new "thinned" network for each training batch.
 - Simply don't update weights corresponding to dropped nodes
- Test
- Multiply the weights at each node with the probability value used to not drop it.



Training: Node retained with prob. p .



Testing: Multiply weights with prob value p .

◦ Weight decay:

$$R(\theta) = \sum_{i=1}^n d(y^{(i)}, \hat{y}^{(i)}) \quad - (1)$$

$$R_{\text{norm}}(\theta) = R(\theta) + \frac{\lambda}{2} \|\theta\|_2^2 \quad - (2)$$

Without decay:

$$\theta^{(r)} = \theta^{(r-1)} - \eta \nabla_{\theta^{(r-1)}} R(\theta) \quad - (3) \text{ from } R(\theta) \text{ in } (1)$$

With decay:

$$\theta^{(r)} = \theta^{(r-1)} - \eta \nabla_{\theta^{(r-1)}} R_{\text{norm}}(\theta)$$

$$\theta^{(r)} = \theta^{(r-1)} - \eta \nabla_{\theta^{(r-1)}} \left(R(D) + \frac{\lambda}{2} \theta^T \theta \right)$$

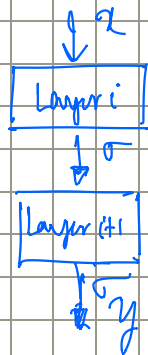
$$\theta^{(r)} = \underbrace{(1 - \eta \lambda)}_{\text{decay of weight}} \theta^{(r-1)} - \eta \nabla_{\theta^{(r-1)}} R(D)$$

decay of weight

Section 7.2 in
Goodfellow et al.

— x —

o Residual Network



Regular NN

$$y = \mathcal{F}(x, \theta)$$

\mathcal{F} is the model that captures the relation between x and y

Res Net

$\mathcal{F}(x, \theta) = y - x$ i.e. $\mathcal{F}(x, \theta)$ models the relation between x and the residual $y - x$.

$$\text{i.e. } \mathcal{F}(x, \theta) = \mathcal{F}(x; \theta) + x$$

