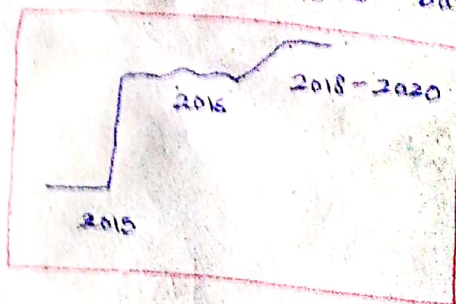
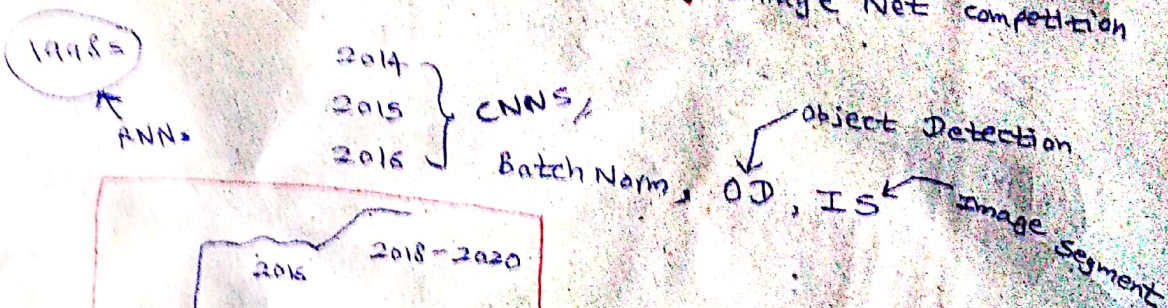


* Expected Test Error = Bias² + Variance + Noise

* **Optimizer** → Minimize the function of loss / error term

* we start with a random point on the function and move in the negative direction of the gradient of the function to reach the local / global minima

* **Revolutions** :- 2012 → Alex Net : Image Net competition



self-attention: Attending to the same sequence

Masking: To prevent undesired masking

$$\text{Attention} = \text{softmax} \left(\frac{q \cdot k^T}{\sqrt{d}} \right) \cdot v$$

Input Pipeline

- 1) **Tokenization** → Breaks a string into a sequence of tokens
- 2) **Lookup** → Assigns integer ID to each token using a vocab
- 3) **(Input) embedding** → Assigns a (learnable) vector to each integer ID in the vocab
- 4) **Positional encoding** → Encodes knowledge about the position in the sequence

Batch Gradient Descent → Fast

Stochastic Gradient Descent → slowest, Learning schedule

* **Mini-Batch Gradient Descent** → Fastest