

Q-1 classifiers

① K-Nearest Neighbor classifier is a non-parametric algorithm for classification purposes. The idea is

i) Loop over all training examples and find K nearest neighbours using some distance (Euclidean/Manhattan)

ii) Pick the class label for data point as the one with majority label in K neighbours.

* K is no of neighbours to select for predicting the label of test point.

* There is no training required in KNN

* Eg → If we are given weights and asked to classify as obese or not obese. Then for a test point, select (K=10) neighbours and select the major label in these 10 neighbours.

② SVM → Margin is the distance between support vectors and boundary hyperplane. Margin is meant to be maximized in SVM to pick the best classifier

$$L = \frac{\lambda}{2n} \|W\|^2 + \max(0, 1 - y_i w^T x_i)$$

$$\text{finding } \frac{\partial \max(0, a)}{\partial a} = \begin{cases} 0 & \text{for } a < 0 \\ 1 & \text{for } a > 0 \end{cases} = \Pi[a > 0]$$

$$\nabla_w L = \frac{\lambda W}{n} + \Pi[1 - y_i w^T x_i > 0] (-y_i x_i)$$

Next Page

$$\nabla_w L = \frac{\lambda w}{n} - y_i x_i \Pi[y_i w^T x_i < 1]$$

Sayam Kumar

S20180010158

Page 2

update rule

$$\cancel{w \leftarrow w + \eta \nabla_w L} \quad w \leftarrow w - \eta \nabla_w L$$

$$\cancel{w = w + \eta} \quad w = w - \eta \frac{\lambda w}{n} + y_i x_i \Pi[y_i w^T x_i < 1] \eta$$

$$= \left(1 - \eta \frac{\lambda}{n}\right) w + \eta \Pi[y_i w^T x_i < 1] y_i x_i$$

In other words

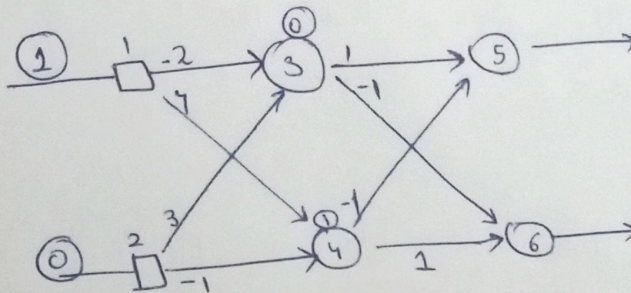
if $y_i w^T x_i < 1$

$$w = w + \eta \left(y_i x_i - \frac{\lambda}{n} w \right)$$

else

$$w = w - \eta \frac{\lambda}{n} w$$

Q-2 Neural Network



$$In_3 = -2 + 0 = 2 \Rightarrow out_3 = 0$$

$$In_4 = 4 + 0 = 4 \Rightarrow out_4 = 1$$

$$In_5 = 0 + -1 = -1 \Rightarrow out_5 = 0$$

$$In_6 = 0 + 1 = 1 \Rightarrow out_6 = 1$$

$$\boxed{\begin{matrix} y_5 = 0 \\ y_6 = 1 \end{matrix}}$$

$$y_5 = 0$$

$$y_6 = 1$$

Q-2 (b) Input Volume

$$30 \times 30 \times 256$$

Max
pool

$$2 \times 2 \quad \text{Stride } 2 \quad \text{Pad } 0$$

$$\text{Output size} = \frac{30-2}{2} + 1 = 15 \times 15$$

$$\text{Output size} = 15 \times 15, \text{ No of params} = 0$$

ii) Input $128 \times 128 \times 3$ Batch Norm Layer

No shape change

$$\text{Output shape } 128 \times 128 \times 3$$

$$\text{No of parameters} = \gamma, \beta \text{ (learnable)} = 2$$

(assuming full batch update)

iii) Input $16 \times 16 \times 16$

$$\text{filters } 5 \times 5$$

8 filters

$$\begin{array}{l} \text{2 stride} \\ \text{2 pad} \end{array}$$

$$W = \frac{16-5+2(2)}{2} + 1 = \frac{11+4}{2} + 1 = \frac{15}{2} + 1 = [8.5] = 8$$

$$\text{Output size} = 8 \times 8 \times 8$$

$$\text{Params} = (5 \times 5 \times 16) \times 8 \text{ weights} + 8 \text{ bias}$$

~~cannot apply this~~

$$\text{filter} = 3200 \text{ weights} + 8 \text{ bias}$$

Q-3 (a) Activation functions are used to introduce non linearities in a training and extracting features from data. It helps in activating and proper gradient flow during training

Next Page

Swish activation function is

$$f(x) = x \cdot \text{sigmoid}(\beta x)$$

where β is learnable parameter / hyperparameter

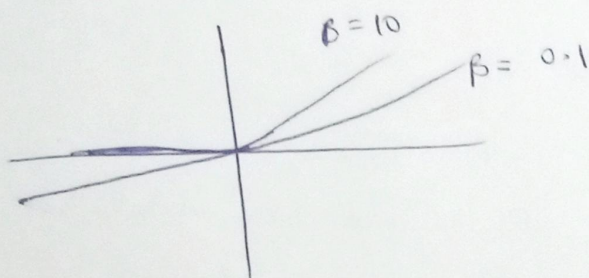
S20180010158

Sayam Kumar

Page 4

ReLU is a special case of swish (ReLU is with high values of β)

graph



(h) Adagrad optimizer \rightarrow Adagrad is for smoothing out the descent when reaching the optima (minima) value while True:

$$dx \leftarrow \text{gradient}(x)$$

$$\text{grad-sqr} \leftarrow \text{grad-sqr} + (dx) * (dx)$$

$$x \leftarrow x - \eta \times \frac{dx}{\sqrt{\text{grad-sqr} + 1e^{-7}}}$$

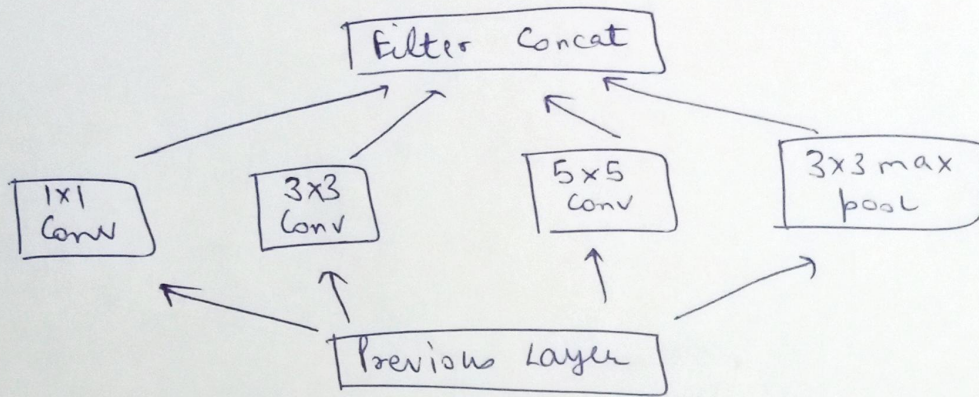
Drawback \rightarrow Over time the grad-sqr goes on increasing, thereby decreasing the learning rate before reaching minima. This is also called effective learning rate diminishing problem

RMS prop resolves this problem by a (decay-rate) parameter that accounts for weighted average of ∇f^2 and helps in mini batch training.

Q-4 Inception Module

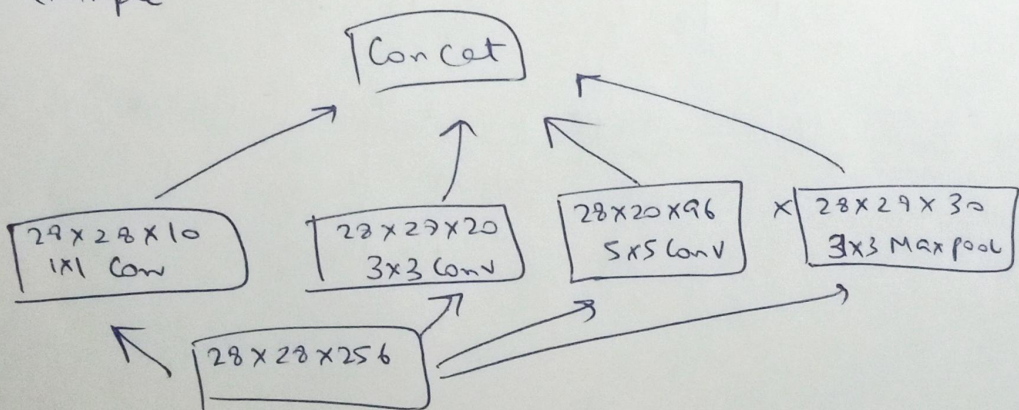
Sayam Kumar
S20180010153
Page 5

Inception is about using multiple filters with different sizes stacked together. This is a deeper network but a lot of computational efficiency.



It has a problem of computation with too many multiplication operations.

Taking example

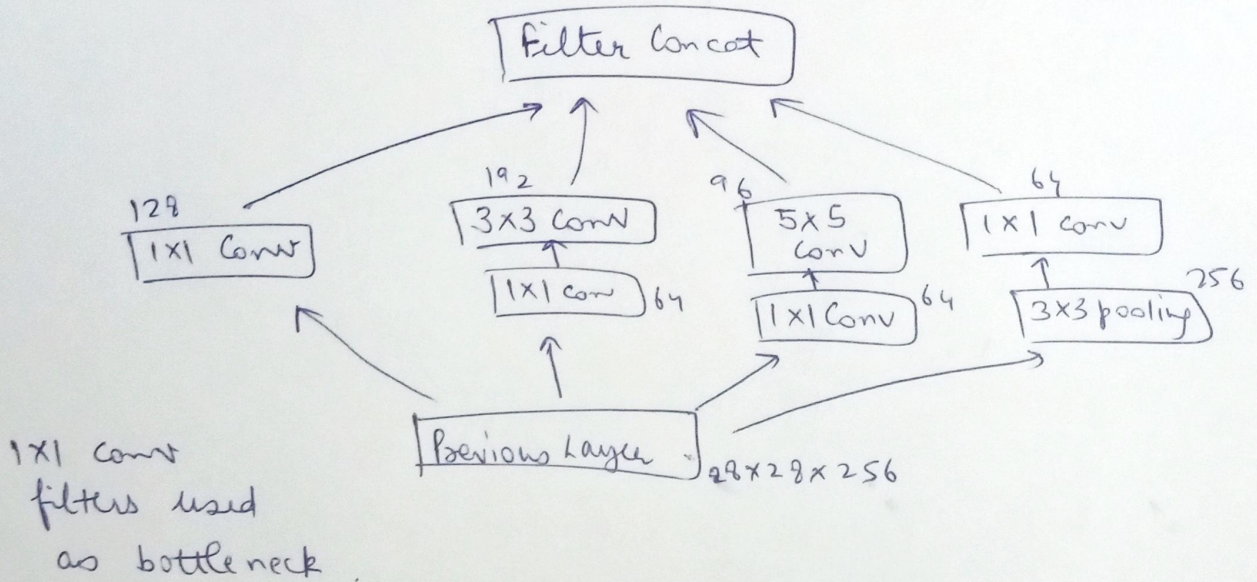


Total about 600M operations

Expensive

⑥ Bottleneck is used with 1×1 conv layer to reduce the computational complexity of Inception module

Sayam Kumar
S20180010158
Page 6



Operations are dramatically reduced from 850 M ops to 350 M operations.

1x1 conv filters helps in better weighted average for gradient flow and no of operations.