

Neural Computation

General Loss Function = (predicted output - original output)².

Underfitting = Training ↓ Testing ↓ Performance.

Oversetting = Training ↑ Testing ↓ Performance

Optimum = Training ↑ Testing ↑ Performance.

Data (100%)	
Training (80%)	Testing (20%)
Train (80%)	Valid (20%)
	Testing (20%)

Linear Model

$$C(\omega) = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \omega^\top x^{(i)})^2 \quad C'(\omega) = \frac{1}{n} \sum_{i=1}^n (x^{(i)})^\top \omega - \frac{1}{n} \sum_{i=1}^n y^{(i)} x^{(i)}$$

$$= \frac{1}{2n} (\bar{\omega}^\top x^\top \omega - 2\omega^\top x^\top y + y^\top y)$$

ω^* is the minimum point.

First order optimal condition : $\nabla C(\omega^*) = 0$.

$$\text{One Dimension case : } \epsilon'(\omega) \quad \omega^* = \frac{\sum_{i=1}^n y^{(i)} x^{(i)}}{\sum_{i=1}^n (x^{(i)})^2}$$

$$(x\omega)^\top = \omega^\top x^\top \quad \text{Closed form: } \omega^* = \arg \min_{\omega} C(\omega) \Rightarrow \omega^* = (x^\top x)^{-1} x^\top y$$

Polynomial regression

$$f(x) = \omega^\top \phi(x) = \omega_0 + \omega_1 x + \omega_2 x^2 + \omega_3 x^3.$$

$$\omega^* = (x^\top x)^{-1} x^\top y$$

$$\hat{y} = x\omega^*$$

Regularized Regression.

$$C(\omega) = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \omega^\top x^{(i)})^2 + \frac{\lambda}{2} \|\omega\|_2^2 \quad \omega^* = \left(\frac{1}{n} (x^\top x) + \lambda I \right)^{-1} \left(\frac{1}{n} x^\top y \right)$$

Gradient Descent

$$\omega^{t+1} = \omega^t - \eta \nabla C(\omega^t)$$

GD For Regression

$$C(\omega) = \frac{1}{2n} (\omega^T x^T x \omega - 2\omega^T x^T y + y^T y)$$

$$\nabla C(\omega) = \frac{1}{n} (x^T x \omega - x^T y)$$

$$\omega^{(t+1)} = \omega^{(t)} - \frac{\eta}{n} (x^T x \omega - x^T y)$$

Stochastic Gradient Descent

$$C(\omega) = \frac{1}{n} \sum_{i=1}^n c_i(\omega)$$

Learning rate $\eta_t = c/\sqrt{t}$, c is a parameter needed to tune.

Minibatch SGD

$$\omega^{(t+1)} = \omega^{(t)} - \frac{\eta_t}{b} \sum_{i \in B_t} \nabla c_i(\omega^{(t)})$$

Linear Model for Classification

$$\hat{y} = \text{sgn}(\omega^T x) = \begin{cases} 1 & \omega^T x > 0 \\ -1 & \text{otherwise} \end{cases}$$

$$\text{Margin} = y \omega^T x$$

$$C(\omega) = \frac{1}{2n} \sum_{i=1}^n (\max \{0, 1 - y^{(i)} \omega^T x^{(i)}\})^2$$

$$\nabla C(\omega) = \begin{cases} 0 & \text{if } y^{(i)} \omega^T x^{(i)} \geq 1 \\ (\omega^T x^{(i)} - y^{(i)}) x^{(i)} & \text{otherwise} \end{cases}$$

SGD for Linear Model

$$c_i(\omega) = \begin{cases} 0 & \text{if } y^{(i)} \omega^\top x^{(i)} \geq 1 \\ \frac{1}{2} (1 - y^{(i)} \omega^\top x^{(i)})^2 & \text{otherwise} \end{cases}$$

$$\omega^{t+1} = \omega^t - \eta \nabla C(\omega^t) = \begin{cases} \omega^{(t)} & \text{if } y^{(i)} (\omega^{(t)})^\top x^{(i)} \geq 1 \\ \omega^{(t)} - \eta y^{(i)} (y^{(i)} \omega^{(t)} - 1) x^{(i)} & \text{otherwise} \end{cases}$$

Perception

if $y\omega^\top x \leq 0$ then $\omega = \omega + yx$.

$$y = \sigma(\sum w_i x_i + b)$$

$$C(\omega) = \frac{1}{2n} \sum_{i=1}^n (\sigma(\omega^\top x^{(i)} + b) - y^{(i)})^2$$

Activation Functions

$$\text{Sigmoid } \sigma(x) = \frac{1}{1 + e^{-x}} \quad (-\infty, \infty) \rightarrow (0, 1)$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (-\infty, \infty) \rightarrow (-1, 1)$$

$$\text{ReLU}(x) = \max(0, x) \quad (-\infty, \infty) \rightarrow (0, \infty)$$

(Not Differentiable)

$$x \rightarrow y \rightarrow z$$

$$\frac{dz}{dx} = \frac{dy}{dx} \times \frac{dz}{dy} \rightarrow \text{local gradient}$$

$$\text{Backpropagation gradient: } (\delta_j^l = \frac{\partial C}{\partial z_j^l})$$

Backpropagation

$$\frac{\partial C}{\partial w^l} = \delta^l (a^{l-1})^\top \quad \frac{\partial C}{\partial b^l} = \delta^l$$

Backpropagated Gradient for Output Layer

$$\delta_j^L = \sigma'(z_j^L) (a_j^L - g_j)$$

For Hidden Layers

$$\delta_j^L = \sigma'(z_j^L) \sum_k \delta_k^{L+1} \omega_{kj}^{L+1}$$

GD with momentum

$$v^{(t+1)} = \alpha \cdot v^{(t)} - \eta \nabla C(\omega^{(t)}) \quad \alpha \in [0, 1]$$

$$\omega^{(t+1)} = \omega^{(t)} + v^{(t+1)}$$

Nesterov Momentum (Nesterov Accelerated Gradient)

$$\omega^{(\text{ahead})} = \omega^{(t)} + \alpha v^{(t)}$$

$$v^{(t+1)} = \alpha v^{(t)} - \eta \nabla C(\omega^{(\text{ahead})})$$

$$\omega^{(t+1)} = \omega^{(t)} + v^{(t+1)}$$

SGD with Momentum

$$g(t) = \nabla C_i(\omega_t)$$

$$v^{(t+1)} = \alpha v^{(t)} - \eta g(t)$$

$$\omega^{(t+1)} = \omega^{(t)} + v^{(t+1)}$$

SGD with Nesterov Momentum

$$\omega^{(\text{ahead})} = \omega^{(t)} + \alpha v^{(t)}$$

$$g^{(\text{ahead})} = \nabla C_i(\omega^{(\text{ahead})})$$

$$v^{(t+1)} = \alpha v^{(t)} - \eta g^{(\text{ahead})}$$

$$\omega^{(t+1)} = \omega^{(t)} + v^{(t+1)}$$

AdaGrad

$$g^{(t)} = \nabla C_i(\omega^{(t)})$$

accumulated gradient norm square $\delta^{(t+1)} = \delta^{(t)} + g^{(t)} \odot g^{(t)}$

$$\omega^{(t+1)} = \omega^{(t)} - \frac{\eta}{\delta + \sqrt{\delta^{(t+1)}}} \odot g^{(t)}$$

RMSProp

$$\gamma^{(t+1)} = \beta \cdot \gamma^{(t)} + (1 - \beta) g^{(t)} \odot g^{(t)}$$

Adam

$$s^{(t+1)} = \beta_1 \cdot s^{(t)} + (1 - \beta_1) g^{(t)}$$

$$\gamma^{(t+1)} = \beta_2 \gamma^{(t)} + (1 - \beta_2) g^{(t)} \odot g^{(t)}$$

$$s^{(t+1)} = \frac{s^{(t+1)}}{1 - \beta_1^{t+1}}$$

$$\gamma^{(t+1)} = \frac{\gamma^{(t+1)}}{1 - \beta_2^{t+1}}$$

$$\omega^{(t+1)} = \omega^{(t)} - \frac{\eta}{\delta + \sqrt{\gamma^{(t+1)}}} \odot s^{(t+1)}$$

CNN Formula

$$\text{Output} = \frac{\text{Input} - \text{Kernel} + 2 \times \text{Padding}}{\text{Stride}} + 1$$

$$\text{Depth} = \text{Number of filters}$$

$$\begin{aligned} \text{Number of Parameters} &= \text{Kernel} \times \text{Kernel} \times \text{Depth} \times \text{no. of filters} \\ &\quad + \text{no. of filters} \end{aligned}$$

Upsampling

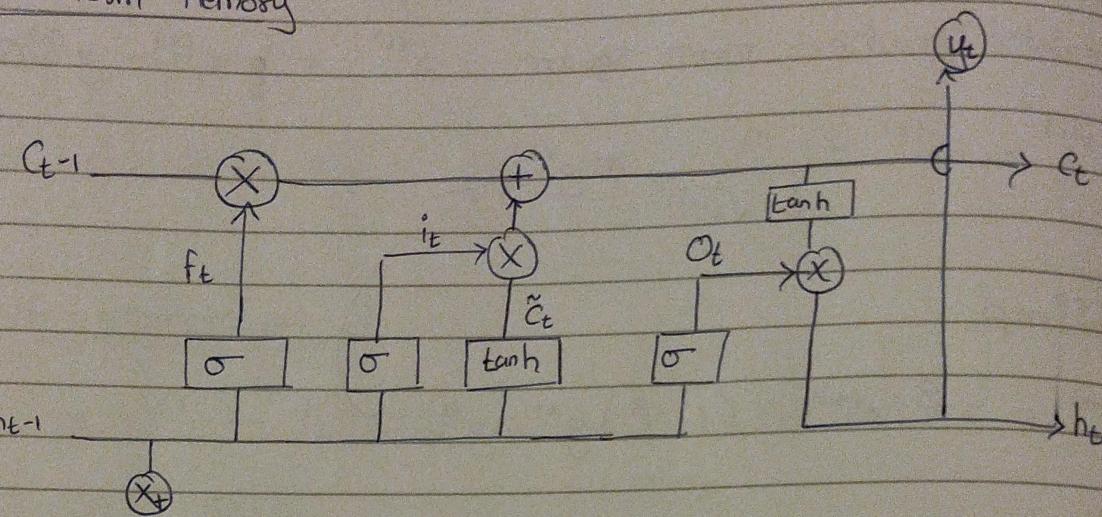
$$\text{Output} = (\text{Input} - 1) \times \text{Stride} - 2 \times \text{Padding} + (\text{Kernel} - 1) + 1 + \text{output padding}$$

$$\text{Depth} = \text{No. of filters}$$

$$\text{Min-Max} = \frac{\text{value} - \text{min}}{\text{max} - \text{min}}$$

$$z\text{-score} = \frac{\text{value} - \text{mean}}{\text{std}}$$

Long Short Term Memory



① forget

$$f_t = \sigma(\omega_{fb}^T h_{t-1} + \omega_{fx}^T x_t)$$

② Store

$$\tilde{i}_t = \sigma(\omega_{ih}^T h_{t-1} + \omega_{ix}^T x_t)$$

$$\tilde{c}_t = \tanh(\omega_{ch}^T h_{t-1} + \omega_{cx}^T x_t)$$

③ Update

$$c_t = f_t \times c_{t-1} + \tilde{i}_t \times \tilde{c}_t$$

④ Output

$$o_t = \sigma(\omega_{oh}^T h_{t-1} + \omega_{ox}^T x_t)$$

$$h_t = o_t \times \tanh(c_t)$$

$$y_t = w_{hy}^T h_t$$

CNN Problem

$$\text{Input} = 11 \times 11$$

$$\text{Output} = 5 \times 5 \times 10$$

$$\text{Kernel} = 8 \times 8$$

$$N = 10$$

$$\text{Stride} = 2$$

$$\text{Padding} = 1$$

$$\text{Output} = 4 \times 4 \times 5$$

$$\text{Kernel} = 1 \times 1 \times 5$$

$$4 = \frac{5 - K + 2}{2} + 1$$

$$11 = 1 \quad \omega = 1 \quad O = 10 \quad M = 5$$

Auto Encoders

$$\text{Loss function} = \frac{1}{d} \sum_{j=1}^d (x^{(j)} - g_\phi^{(j)}(f_\phi(x)))^2$$

Probabilistic Encoder

$$P_\phi(z|x) = N(\mu_\phi(x), \sigma_\phi^2(x)^2)$$

VAE

$$L_{\text{VAE}} = L_{\text{rec}} + \lambda L_{\text{reg}}$$

↑
Reconstruction loss Regularizer.

$$L_{\text{rec}} = \frac{1}{d} \sum_{j=1}^d (x^{(j)} - g_\phi(\tilde{z}^{(j)})^2)$$

Reparameterization Trick

$$\tilde{z} = N(\mu_\phi(x), \sigma_\phi^2(x)^2) = \mu_\phi(x) + \sigma_\phi(x) \odot \epsilon . \quad \epsilon \sim N(0, I)$$

Practise

$$C(\omega) = \frac{1}{2} (\omega^\top x - y)^2$$

$$\nabla C(\omega) = (\omega^\top x - y) x \quad \text{margin} = y \omega^\top x$$

For Linear classification, $\nabla C(\omega) = (1 - \text{margin}) g x$

Paper

- Q1. a) No. of neurons in Input Layer = 2.
 Hidden = 2
 Output = 1.

This is a fully connected NN.

$$\text{No. of parameters} = [(2 \times 2) + 2] + [(2 \times 1) + 1] = 6 + 3 = 9$$