NN → ASSIGNMENT - 07 , Tutor! - Redion
Priyanka Upadhyay (2581714), s8prupad@stud.uni-soarland.de
Gopal Bhattrai (7013597), gobho0001@stud-uni-sparland.de

| Problem → 7.1 | $L_1$ & $L_2$ Regularization |
| --- | --- |

**(1.a)**



L₁ Regularization Norm          L₂ Regularization

**(1.b)** In $L_1$ Regularization, the predictor' which does not have the significant impact on prediction, there coefficient becomes 0 (zero) which can result as less computational cost. whereas In $L_2$ Regularization, coefficient never becomes equal to 0 ( their value becomes close to 0).

So that's why we prefer $L_1$ Norm over $L_2$.

**(1.c)** Yes. $L_1$ Regularization can be used ↑for feature when the coefficient values is very small, and it does not hold much significance over prediction value.

$L_2$ Regularization should / can be used for features with very High value of coefficient and are highly significant for outcome.

**(1.d)** No, The main aim of training model is to perform well on unseen data and Hence Regularization is done to improve Generalization / Test error by introducing some bias and decreasing the variance.

**(1.e)** with 1 Feature:

$$L_2 = (y - wx)^2 + \alpha w^2 = y^2 + w^2 x^2 - 2wxy + \alpha w^2$$

Put $\dfrac{\partial L_2}{\partial w} = 0$

$$0 + 2wx^2 - 2xy + 2\alpha w = 0$$

$$\Rightarrow \boxed{w = \dfrac{xy}{x^2 + \alpha}}$$ w will become zero, when & only $\alpha \to \infty$.

**(1.f)** with 1 Feature

$$L_1 = (y - wx)^2 + \alpha w = y^2 + w^2 x^2 - 2wxy + \alpha w$$

$$\dfrac{\partial L_1}{\partial w} = 0 \to 0 + 2wx^2 - 2xy + \alpha = 0$$

$$\boxed{w = \dfrac{2xy - \alpha}{2x^2}}$$

w will become, when $\alpha = 2xy$

**Problem 7.2** multi-task learning and Data Augmentation

**(a)** Problems:
1. In some cases, Task performances decreases due to the task interference.
2. Lack of Analytical tool to predict task interference.
3. It becomes difficult to apply multi-task learning to Neural Networks, because it's feature space is given by layer-wise embeddings.

Solutions:
1. If the shared module's capacity is too large, then there is no interference between tasks.
2. Task variance is used to determine task interference.
3. Assign per-task weights for settings where different tasks share the same data but have different labels.

| Problem 7.3 | Early Stopping

(a)
$$E = E_0 + \frac{1}{2} (w - w^*)^T H (w - w^*)$$

After $\tau$ steps, the components of the weight vector parallel to the eigen vectors of $H$ is:

$$w_j^\tau = \left\{ 1 - (1 - \varepsilon \eta_j)^\tau \right\} w_j^*$$

$$\nabla E = H \left( w^{(\tau-1)} - w^* \right)$$

$\rightarrow$ $S_0$   $\tau \rightarrow \infty$    i.e   $\tau = 1, 2, 3, 4 \dots$

$\tau = 1 \rightarrow$   $\omega^{(0)} = 0$ (origin)

$$\omega_j^{(1)} = \omega_i^{(0)} - \varepsilon \, \eta_j \left( \omega_j^{(0)} - \omega_j^{*} \right)$$

$$= \varepsilon \, \eta_j \, \omega_j^{*} \longrightarrow \text{computing with given}$$
$$\text{value}$$

$$= 1 - \left( 1 - \varepsilon \, \eta_j \right) \, \omega_j^{*}$$

$\tau = N \rightarrow$ we   assume the   Result holds for

$$\tau = N - 1$$

$$w_j^N = w_j^{N-1} - \varepsilon \wedge_j \left( w_j^{(N-1)} - w_j^* \right)$$

$$= w_j^{(N-1)}(1 - \varepsilon \wedge_j) + \varepsilon \wedge_j w_j^*$$

$$= \left[ 1 - (1 - \varepsilon \wedge_j)^{N-1} \right] w_j^A (1 - \varepsilon \wedge_j)$$

$$+ \varepsilon \wedge_j w_j^*$$

$$= \left\{ (1 - \varepsilon \wedge_j) - (1 - \varepsilon \wedge_j)^N \right\} w_j^* + \varepsilon \wedge_j w_j^*$$

$$= \left\{ 1 - (1 - \varepsilon \wedge_j)^N \right\} w_j^*$$

which also Holds the Result

and provided $[1 - \varepsilon \wedge_j] < 1$

$(1 - \varepsilon \wedge_j)^\tau \to 0$ as $\tau \to \infty$ hence

$$\left\{ 1 - (1 - \varepsilon \wedge_j)^N \right\} \to 1 \text{ and we}$$

can prove that $\quad w^{(\tau)} \to w^*$

---

(b) If $\tau$ is finite but $\wedge_j \gg (\varepsilon \tau)^{-1}$
then $\tau$ must still be large, since
$\wedge_j \varepsilon \tau \gg 1$, even though $|1 - \varepsilon \wedge_j| < 1$

if $\tau$ is larger than from above we
can show that $\quad w_j^{(\tau)} \sim w_j^*$

if $\wedge_j \ll (\varepsilon \tau)^{-1}$, then it means $\varepsilon \wedge_j$
must be small, since $\varepsilon \wedge_j \tau \ll 1$
& $\tau$ is an integer greter than or
equal to $1$

$$w_j^{(\tau)} \ll |w_j^*|$$

(c) we know that, when regularization parameter $(\alpha)$ is much larger than one of the eigen values $(\Lambda_j)$ then corrosponding parameter value $w_i$ will close to 0.

In the same way when $\alpha$ is smaller than $\Lambda_j$, $w_i$, will be close to its maximum likelihood value.

Thus we can say that $\alpha$ is analogous to $\varepsilon\tau$

Problem 7.4    Universal Approximation Theorem

(a)   This paper contribute and validate the Idea that Multilayer Feedforward Architecture provides Neural Network the potential of being universal learning machines and not the choice of Activation Function.

(b)   The main idea of the paper is that activation function is more important for approximation rather than the depth of the Network. It states that Neural Networks with a Single hidden layer can accurately make predictions when provided with a continuous sigmoidal Non-Linear activation function.

(c)   A Deeper Network with multiple layers effectively learns various features at different levels of abstraction and hence provide better generalization.