

Assignment 7

Tuesday, 3. January 2023 22:27

Prathrishi Mithmare - 7028692

Subrat Krishore Dutta - 7028082

Exercise 7.1

- a] From the paper we get to know that a feedforward network with given hidden layer, with given number of neurons can approximate continuous functions under appropriate activation functions. And also from the paper we get to know the reason why neural network performs well and can approximate functions.
- b] The main idea of the paper from part a is that since it uses appropriate activation function with given hidden layers it can approximate any continuous function. Approximation by superposition of a sigmoidal function like sigmoid or tanh activation function, thus we can say activation function plays an important role.
- c] Deeper networks can learn more complex relationships in the data and also they can generalize better which thus reduces overfitting. Deeper networks can also learn tree structures [ex: edges] better.

Exercise 7.2 L₁ and L₂ regularization

- (a) The bias term is responsible for only controlling one variable on the other hand weights control two variables i.e. it specifies how two variables interact. The bias term does not contribute much to the curvature of the surface that the model learns, thus regularizing it does not reduce the overfitting (high variance), or in other words does not induce much bias. On the contrary it can introduce significant underfitting as the position of the surface along the corresponding dimension can shift significantly.
- (b) As a result of using L₁ regularization the solution we obtain is sparse meaning the coefficients corresponding to a subset of the parameters can turn out to be zero. Using such an approach can be advantageous when we want a simple model which only takes into consideration the most important features in a dataset.

Mixing them together in some weighted combination can be beneficial as they will create a balance between the sparsity offered by L₁ regularization and the flexibility offered by L₂ regularization.

- (c) In general it is true as regularization is used to prevent overfitting and make a balance between minimizing the training error and minimizing the generalization error. When we regularize by adding penalty it normally has a tendency to increase the training error but it results in the model not overfitting and generalizing better resulting in better performance on the test data.

- (d) The penalty term added to the objective function is such that it penalizes the model if the weights are too large. The optimization algorithm will try to minimize the objective function by reducing the weight of the model, thus if the weight becomes too large the optimization algorithm may turn that to zero.

the weight is updated as

$$w_{\text{new}} = w_{\text{old}} - \alpha \frac{\partial L}{\partial w}$$

taking $\alpha = 0.5$ and $w_{\text{old}} = 1$

$$w_{\text{new}} = 1 - 0.5(1) = 0.5.$$

as we can see here the weight can be set to zero.

in case of L₂ with only one feature and in case of a L₂ regularized least-squares case the loss function can be approximated around the optimal value w^*

$$L(w) = \|w - w^*\|^2 + \frac{1}{2}(\alpha - w^*)H(w - w^*)$$

taking its derivative gives weight 0 at minima

$$\frac{\partial L(w)}{\partial w} = 0 + H(w - w^*) = 0$$

$$\Rightarrow \frac{\partial L(w)}{\partial w} H(w - w^*) = 0 \rightarrow 0$$

now the derivative of L₂ regularization gives αw adding this to eqn 0 and at minima.

$$\Rightarrow H(w - w^*) + \alpha w = 0$$

$$\Rightarrow Hw - Hw^* + \alpha w = 0$$

$$\Rightarrow (H + \alpha)w = Hw^*$$

$$\Rightarrow w = \frac{H}{(H + \alpha)}w^*$$

now in a least square problem $H > 0$ hence

number w^* is 0 so cannot be zero.

Exercise 7.3:

- a] $w_{t+1} = w_t - \eta \nabla_w J(w_t) \rightarrow \textcircled{1}$
For regularised loss the equation $\textcircled{1}$ becomes

$$w_{t+1} = w_t - \eta (\nabla_w J(w_t) + \lambda w_t)$$

$$= w_t - \eta \lambda w_t - \eta \nabla_w J(w_t)$$

$$= w_t(1 - \eta \lambda) - \eta \nabla_w J(w_t)$$

where η = learning parameter
 λ = regularization parameter.

$(1 - \eta \lambda)$ = weight decay.

- b] $\tilde{w} = (H + \alpha I)^{-1} H w^*$

$$\tilde{J}(w) = J(w^*) + \frac{1}{2}(w - w^*)^T H (w - w^*)$$

$$\nabla_w \tilde{J}(w) = H(w - w^*)$$

$$w^{(t+1)} = w^{(t+1)} - \varepsilon \nabla_w \tilde{J}(w^{(t+1)})$$

$$= w^{(t+1)} - \varepsilon H(w^{(t+1)} - w^*)$$

$$w^{(t+1)} - w^* = w^{(t+1)} - H(w^{(t+1)} - w^*) + \varepsilon H w^* - w^*$$

By taking common elements

$$w^{(t+1)} - w^* = (I - \varepsilon H)(w^{(t+1)} - w^*)$$

We know that $H = Q^T Q$

$$w^{(t+1)} - w^* = (I - \varepsilon Q^T Q)(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (Q^T - \varepsilon I)(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)(Q^T - \varepsilon I)(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)(I - \varepsilon Q^T)(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$

$$Q^T(w^{(t+1)} - w^*) = (I - \varepsilon Q^T)^2(w^{(t+1)} - w^*)$$