

Machine Learning: Regression and Neural Nets

1 Regularized Linear Regression

Recall the problem of linear regression and the derivation of its various learning algorithms. We are given a dataset $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$, where $x^{(i)} \in \mathbb{R}^p$ (i.e., each data point has p features) and $y^{(i)} \in \mathbb{R}$. The hypothesis class we consider in linear regression is, for $w \in \mathbb{R}^{p+1}$:

$$h_w(x) = w_0 + w_1x_1 + \dots + w_px_p = w_0 + \sum_{j=1}^p w_jx_j$$

Consider the following alternative error function to *minimize* (while using squared-error loss):

$$J(w) = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - h_w(x^{(i)}))^2 + \lambda \sum_{j=1}^p w_j^2$$

1. Compared to standard linear regression, there is an extra *regularization* term: $\lambda \sum_{j=1}^p w_j^2$. From an optimization perspective, what is the role of this term, i.e., which weights (parameters) w does it prefer, and which does it penalize, assuming $\lambda > 0$? *Regularization helps to avoid the overfitting.*

2. λ is known as a *regularization constant*; it is a *hyperparameter* that is chosen by the algorithm designer and fixed during learning. What do you expect to happen to the optimal w when $\lambda = 0$? $\lambda \rightarrow +\infty$? $\lambda \rightarrow -\infty$? *LL likes the small value, since it penalize the large one.*

3. Find the partial derivatives $\frac{\partial J}{\partial w_0}$ and $\frac{\partial J}{\partial w_j}, j \in \{1, \dots, p\}$. (These two cases are different; for the second one, only derive once for an arbitrary index j .) *$\lambda=0$, no regularization applied. $\lambda \rightarrow -\infty$ can get optimal w | all $w \rightarrow 0$. $\lambda \rightarrow \infty$, penalty grow.*

4. Write out the update rule for gradient descent applied to the error function $J(w)$ above. Compare with the gradient descent update rule for standard linear regression; what is the difference? How does this difference affect gradient descent, assuming $\lambda > 0$?

Consider what happens both when w_j is positive and when it is negative.

5. In order for the problem to be well-defined, an appropriate regularization constant λ must be chosen for the given problem (dataset). How should the designer choose λ ?

$$3. \frac{\partial J}{\partial w_0} = \frac{1}{n} \cdot \sum_{i=1}^n \frac{\partial}{\partial w_0} (y_i - h_{w_0}(x_i))^2 + \lambda \frac{\partial}{\partial w_0} (w_0^2 + w_1^2 + w_2^2)$$

$$= -\frac{1}{n} \sum_{i=1}^n 2(y_i - h_{w_0}(x_i))$$

$$\frac{\partial J}{\partial w_i} = -\frac{1}{n} \sum_{i=1}^n 2(y_i - h_{w_0}(x_i)) \cdot x_i + 2\lambda w_i$$

$$4. w_0 = w_0 + 2 \left[\frac{1}{n} \sum_{i=1}^n (y_i - h_{w_0}(x_i)) \right]$$

$$w_j = w_j + 2 \left[\frac{1}{n} \sum_{i=1}^n (y_i - h_{w_0}(x_i)) \cdot x_{ij} - \lambda \cdot 2w_j \right]$$

5. The regularization we saw can't be too scalable, since it doesn't avoid the overfit, it has too much bias. The designer can use the cross fold method, to average some parameters and plot the results to find a great range for it.

2 Perceptron

You are given the following training set:

Sample	x_1	x_2	x_3	y
1	1	4	1	1
2	1	2	3	1
3	0	0	1	1
4	-1	4	0	1
5	1	0	-2	0
6	-1	-1	1	0
7	0	-4	0	0
8	1	0	-3	0

Run two iterations of the perceptron learning algorithm on this data set. Start with initial weights of 0, and use a learning rate of $\alpha = 0.5$. Don't forget to add a bias neuron $x_0 = 1$.

$$i=1$$

$$f(x) = x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3$$

$$\rightarrow 1 = 1 + 0 + 0 + 0$$

$$w_i = w_i + 0.5(1 - 1)$$

$$1 = w_i + \dots$$

$$0 = (1 + 0 + 0 + 0) = 1$$

$$\Delta w_1 = 0.5 \cdot (-1) \cdot 1 = -0.5$$

$$w_1 = 0 - 0.5 = -0.5$$

$$\Delta w_3 = 0.5 \cdot (-1) \cdot -2 = 1$$

$$w_3 = 0 + 1 = 1$$

$$6) 1 + (-0.5 + 0.5) + 0 + 1 \cdot 1 = 2.5 > 1 = 1$$

$$\Delta w_1 = 0 - 0.5 + 0.5 = 0$$

$$\Delta w_2 = 0.5 = 0.5$$

$$\Delta w_3 = 1 + 0.5 \cdot (0 - 1) \cdot 1 = 0.5$$

$$7) w_0 = -1 + 0.5 = -0.5$$

$$\Delta w_1 = 0.5 \cdot 0 = 0$$

$$w_1 = 0$$

$$\Delta w_2 = 0 \quad w_2 = 0.5$$

$$\Delta w_3 = 0 \quad w_3 = 0.5$$

$$8) \uparrow \text{ same}$$

$$I=2$$

$$f(x) = -1x_1 + 0x_2 + 0.5x_3 + 0.5x_4$$

$$1) f(x) = 1 = y = 1$$

$$2)$$

$$5) \Delta w_0 = 0.5 \quad \Delta w_1 = 0$$

$$w_0 = -0.5 \quad \Delta w_2 = 0.5$$

$$w_3 = 0.5 + 0.5 = 1$$

$$3) w_1 = 0$$

$$6) w_0 = -0.5 \quad \Delta w_0 = -0.5$$

$$w_0 = -0.5 + (-0.5) = -1$$

$$w_1 = 0 + 0.5 = 0.5$$

$$w_2 = 1$$

$$\Delta w_3 = 0.5 \quad w_3 = 0.5$$

$$7)$$

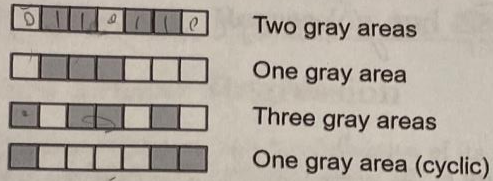
$$f(x): \text{practice} = \text{label}$$

$$8) \text{ so no change}$$

$$f(x) = -1x_1 + 0.5x_2 + x_3 + 0.5x_4$$

3 Neural Networks

You are given the following patterns:



Each pattern consists of N cells, that can be white or gray. Treat the patterns as cyclic (i.e., the first and the last cell are considered to be adjacent).

Design a neural network that identifies if there are more than G gray areas in the pattern. Show the architecture of your network and its weights for $N = 7$ and $G = 2$.

MORE than $N \rightarrow 1$

