

Machine Learning 1 - Homework 4

Pascal Mattia Esser

October 11th, 2017

1 Lagrange Multipliers: Warm-up

1.1

1.2

1.3

$$\max : x_1 + 2x_2 - 2x_3 \quad (1)$$

$$\text{const} : x_1^2 + x_2^2 + x_3^2 = 1 \quad (2)$$

$$L(x_1, x_2, x_3, \lambda) = x_1 + 2x_2 - 2x_3 + \lambda(x_1^2 + x_2^2 + x_3^2 - 1) \quad (3)$$

Constraints:

$$\nabla_{x_1, x_2, x_3, \lambda} L(x_1, x_2, x_3, \lambda) = (2\lambda x_1 + 1, 2\lambda x_2 + 2, 2\lambda x_3 - 2, x_1^2 + x_2^2 + x_3^2 - 1) = 0 \quad (4)$$

we get: $\lambda = -1/(2x_1)$, $x_2 = 2x_1$, $x_3 = -2x_1$, $9x_1^2 = 1$ gives:

$$x_1 = 1/3 \quad (5)$$

$$x_2 = 2/3 \quad (6)$$

$$x_3 = -2/3 \quad (7)$$

which gives us $f(x_1, x_2, x_3) = 3$ or

$$\lambda = 1.5 \quad (8)$$

$$x_1 = -1/3 \quad (9)$$

$$x_2 = -2/3 \quad (10)$$

$$x_3 = 2/3 \quad (11)$$

which gives us $f(x_1, x_2, x_3) = -3$. therefore the first alternative maximizes the function.

1.4

$$\max : 1 - x_1^2 + x_2^2 \quad (12)$$

$$\text{const} : -x_1 - x_2 - 1 \geq 0 \quad (13)$$

$$L(x_1, x_2, \lambda) = 1 - x_1^2 + x_2^2 + \lambda(-x_1 - x_2 - 1) \quad (14)$$

$$\nabla_{x_1, x_2, \lambda} L(x_1, x_2, \lambda) = (-\lambda - 2x_1, -\lambda - 2x_2, x_1 + x_2 - 1) = 0 \quad (15)$$

we get $x_1 = x_2$ if $\lambda \neq 0$. we would get that $\lambda = -1$, which could contradict the constrain $\lambda \geq 0$, which gives us:

$$\lambda = 0 \quad (16)$$

$$x_1 = 0 \quad (17)$$

$$x_2 = 0 \quad (18)$$

as the maximized valued. Alternatively:

1.5

$$max : 6x^{2/3}y^{1/2} \quad (19)$$

$$const : 4x + 3y \leq 7000 \quad (20)$$

$$L(x, y, \lambda) = 6x^{2/3}y^{1/2} + \lambda(4x + 3y - 7000) \quad (21)$$

$$\nabla_{x, y, \lambda} L(x, y, \lambda) = \left(\frac{4\sqrt{y}}{\sqrt[3]{x}} + 4\lambda, \frac{3x^{2/3}}{\sqrt{y}} + 3\lambda, 4x + 3y - 7000 \right) = 0 \quad (22)$$

we get $\lambda = \frac{z^{1/2}}{x^{1/3}}$. this gives us $x = y$

$$\lambda = \frac{1}{\sqrt{10}} \quad (23)$$

$$x = 1000 \quad (24)$$

$$y = 1000 \quad (25)$$

as the maximized valued.

2 Kernel Outlier Detection

2.1

define the primal program for the circle as: From

$$\forall i : \|\mathbf{x}_i - \mathbf{a}\|^2 \leq R^2 + \xi_i \leq 0, \xi_i \geq 0 \quad (26)$$

rewrite as:

$$\forall i : \|\mathbf{x}_i - \mathbf{a}\|^2 - R^2 + \xi_i \leq 0, -\xi_i \leq 0 \quad (27)$$

primal Lagrangian:

$$L(\mathbf{a}, R, \boldsymbol{\xi}, \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\mu}) = R^2 + C \sum_{i=1}^N \xi_i + \sum_{i=1}^N \alpha_i (\|\mathbf{x}_i - \mathbf{a}\|^2 - R^2 - \xi_i) - \sum_{i=1}^N \mu_i \xi_i \quad (28)$$

2.2

let in general $N_{\setminus z}$ stand for terms that are not depending on z . Use this to shorten the notation for derivatives.

2.2.1

$$\nabla_{R^2} L = \nabla_{R^2} \left(R^2 + \sum_{i=1}^N \alpha_i (-R^2) + N_{\setminus R^2} \right) \quad (29)$$

$$= 1 - \sum_{i=1}^N \alpha_i = 0 \quad (30)$$

$$\Leftrightarrow \sum_{i=1}^N \alpha_i = 1 \quad (31)$$

2.2.2

$$\nabla_{\mathbf{a}} L = \nabla_{\mathbf{a}} \left(\sum_{i=1}^N \alpha_i (\|\mathbf{x}_i - \mathbf{a}\|^2) + N_{\setminus \mathbf{a}} \right) \quad (32)$$

$$= \nabla_{\mathbf{a}} \left(\sum_{i=1}^N \alpha_i ((\mathbf{x}_i - \mathbf{a})^T (\mathbf{x}_i - \mathbf{a})) + N_{\setminus \mathbf{a}} \right) \quad (33)$$

$$= \nabla_{\mathbf{a}} \left(\sum_{i=1}^N \alpha_i (\mathbf{x}_i^T \mathbf{x}_i - 2\mathbf{x}_i^T \mathbf{a} + \mathbf{a}^T \mathbf{a}) + N_{\setminus \mathbf{a}} \right) \quad (34)$$

$$= \sum_{i=1}^N \alpha_i (-2\mathbf{x}_i + 2\mathbf{a}) = 0 \quad (35)$$

$$\Leftrightarrow \sum_{i=1}^N \alpha_i \mathbf{x}_i = \mathbf{a} \sum_{i=1}^N \alpha_i \quad (36)$$

$$\text{Using Equation 31} \quad (37)$$

$$\sum_{i=1}^N \alpha_i \mathbf{x}_i = \mathbf{a} * 1 = \mathbf{a} \quad (38)$$

2.2.3

$$\forall i : \nabla_{\boldsymbol{\xi}} L = \nabla_{\boldsymbol{\xi}} \left(C \sum_{i=1}^N \xi_i + \sum_{i=1}^N \alpha_i (-\xi_i) - \sum_{i=1}^N \mu_i \xi_i \right) \quad (39)$$

$$C - \alpha_i - \mu_i = 0 \quad (40)$$

$$\Leftrightarrow \forall i : \mu_i = \alpha_i + C \quad (41)$$

2.2.4

KKT for slack term

$$\forall i : \alpha_i (||\mathbf{x}_i - \mathbf{a}||^2 - R^2 - \xi_i) = 0 \quad (42)$$

$$\mu_i \xi_i = 0 \quad (43)$$

$$\forall i : ||\mathbf{x}_i - \mathbf{a}||^2 - R^2 - \xi_i \leq 0 \quad (44)$$

$$\alpha_i \geq 0 \quad (45)$$

$$\mu_i \geq 0 \quad (46)$$

$$\xi_i \geq 0 \quad (47)$$

2.3

Derive which data-cases \mathbf{x}_i will have $\alpha_i > 0$ and which ones will have $\mu_i > 0$. NOTE: I solved this question by starting from the given equations [Equation 41](#) to [Equation 45](#) and derive from them the relation of the point to the radius:

For $\alpha_i > 0$: From [Equation 42](#) and [Equation 45](#) follows, that $\alpha_i > 0 \Rightarrow ||\mathbf{x}_i - \mathbf{a}||^2 - R^2 - \xi_i = 0 \Leftrightarrow ||\mathbf{x}_i - \mathbf{a}||^2 = R^2 + \xi_i$. from $\alpha_i > 0$ and [Equation 41](#) follows $\mu_i \neq 0 \therefore \mu_i = \mathbf{a}_i + C$ with $\xi_i = 0 \therefore \mu_i \xi_i = 0$ ([Equation 43](#)).

What follows is: $||\mathbf{x}_i - \mathbf{a}||^2 = R^2$, assuming $C \neq a_i$. meaning that \mathbf{x}_i is a support vector and therefore defines the radius.

If we now consider $C = a_i$, ξ can also be $\neq 0$ because $\mu_i \xi_i = 0$ is no longer required to have $\xi_i = 0$ to be fulfilled. This gives us $||\mathbf{x}_i - \mathbf{a}||^2 \geq R^2$, meaning, the point can be outside of the radius.

For $\mu_i > 0$: $\xi_i = 0 \therefore \mu_i \xi_i = 0$ ([Equation 43](#)). From [Equation 44](#) follows: $||\mathbf{x}_i - \mathbf{a}||^2 \leq R^2$ meaning that \mathbf{x}_i is in or on the radius.

2.4

Derive the dual Lagrangian and specify the dual optimization problem. Kernelize the problem, i.e. write the dual program only in terms of kernel entries and Lagrange multipliers.

$$D_{\mathbf{x}}(\alpha, \mu) = \min_{\mathbf{a}, R, \xi} \left(R^2 + C \sum_{i=1}^N \xi_i + \sum_{i=1}^N \alpha_i (\|\mathbf{x}_i - \mathbf{a}\|^2 - R^2 - \xi_i) - \sum_{i=1}^N \mu_i \xi_i \right) \quad (48)$$

$$= \min_{\mathbf{a}, R, \xi} \left(R^2 + \sum_{i=1}^N \alpha_i (\|\mathbf{x}_i - \mathbf{a}\|^2) - \sum_{i=1}^N \alpha_i R^2 \right) \quad (49)$$

$$= \min_{\mathbf{a}, R, \xi} \left(\sum_{i=1}^N \alpha_i (\mathbf{x}_i^T \mathbf{x}_i - \mathbf{x}_i^T \mathbf{a} - \mathbf{a}^T \mathbf{x}_i + \mathbf{a}^T \mathbf{a}) \right) \quad (50)$$

$$= \min_{\mathbf{a}, R, \xi} \left(\sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{x}_i - \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{a} - \mathbf{a}^T \sum_{i=1}^N \alpha_i (\mathbf{x}_i - \mathbf{a}) \right) \quad (51)$$

$$= \min_{\mathbf{a}, R, \xi} \left(\sum_{i=1}^N \alpha_i k(\mathbf{x}_i^T \mathbf{x}_i) - \sum_{i=1}^N \alpha_i \mathbf{x}_i^T \sum_{j=1}^N \alpha_j \mathbf{x}_j \right) \quad (52)$$

$$= \min_{\mathbf{a}, R, \xi} \left(\sum_{i=1}^N \alpha_i k(\mathbf{x}_i, \mathbf{x}_i) - \sum_{ij}^N (\alpha_i \alpha_j) (\mathbf{x}_i^T \mathbf{x}_j) \right) \quad (53)$$

$$\text{rewriting with a kernel} \quad (54)$$

$$= \sum_{i=1}^N \alpha_i k(\mathbf{x}_i, \mathbf{x}_i) - \sum_{ij}^N (\alpha_i \alpha_j) k(\mathbf{x}_i, \mathbf{x}_j) \quad (55)$$

2.5

Optimize Values for $\{\alpha_i\}$, $0 < \alpha < c$. For this, we have to maximize over α :

$$\min_{\alpha} D_{\mathbf{x}}(\alpha, \mu) = \min_{\alpha} \left(\sum_{i=1}^N \alpha_i k(\mathbf{x}_i, \mathbf{x}_i) - \sum_{ij}^N (\alpha_i \alpha_j) k(\mathbf{x}_i, \mathbf{x}_j) \right) \quad (56)$$

in terms on α_i compute μ_i , use $\alpha_i \geq 0$ (Equation 45), $\sum \alpha_i = 1$ (Equation 31) and $\mu_i = C - \alpha_i$ (Equation 41).

Solving for \mathbf{a}^* (Equation 38)

$$\sum_{i=1}^N \alpha_i \mathbf{x}_i = \mathbf{a}^* \quad (57)$$

Solve for R : for $\alpha_i > 0 \rightarrow \|\mathbf{x}_i - \mathbf{a}^*\|^2 = R^2$ for this $x_i = 0$ holds. write the radius as:

$$R = \|\mathbf{x}_r - \mathbf{a}^*\| \quad (58)$$

solve for ξ_i

$$\|\mathbf{x}_r - \mathbf{a}^*\|^2 < R \Rightarrow \xi_i = 0 \quad (59)$$

$$\|\mathbf{x}_r - \mathbf{a}^*\|^2 > R \Rightarrow \xi_i = \|\mathbf{x}_r - \mathbf{a}^*\|^2 - R \quad (60)$$

2.6

let \mathbf{x}_o being an out layer. we need a way to write R as described in the last section for \mathbf{x}_r , being a support vector.

$$R = \|\mathbf{x}_r - \mathbf{a}\| \quad (61)$$

Now write an outlayer, \mathbf{x}_o , as:

$$\|\mathbf{x}_o - \mathbf{a}\| > R \quad (62)$$

$$\|\mathbf{x}_o - \mathbf{a}\| > \|\mathbf{x}_r - \mathbf{a}\| \quad (63)$$

$$(\mathbf{x}_o^T \mathbf{x}_o - 2\mathbf{x}_o^T \mathbf{a} + \mathbf{a}^T \mathbf{a}) > (\mathbf{x}_r^T \mathbf{x}_r - 2\mathbf{x}_r^T \mathbf{a} + \mathbf{a}^T \mathbf{a}) \quad (64)$$

rewrite with kernels and Lagrangian parameters:

$$k(\mathbf{x}_o, \mathbf{x}_o) - 2 \sum_i^N \alpha_i k(\mathbf{x}_o^T \mathbf{x}_i) + \sum_{ij}^N (\alpha_i \alpha_j) k(\mathbf{x}_i^T \mathbf{x}_j) \quad (65)$$

$$> k(\mathbf{x}_r, \mathbf{x}_r) - 2 \sum_i^N \alpha_i k(\mathbf{x}_r^T \mathbf{x}_i) + \sum_{ij}^N (\alpha_i \alpha_j) k(\mathbf{x}_i^T \mathbf{x}_j) \quad (66)$$

2.7

now let $C = 0$. this means, no penalty for out layers. we can

$$\min_{\alpha, \mu, R} L(\mathbf{a}, R, \boldsymbol{\xi}, \mathbf{x}, \boldsymbol{\alpha}, \mathbf{R}, \boldsymbol{\xi}) \quad (67)$$

if we set $R^2 = 0$

$$\forall i : \|\mathbf{x}_i - \mathbf{a}\|^2 - R^2 - \xi_i \leq 0, -\xi_i \leq 0 \quad (68)$$

$$\forall i : \|\mathbf{x}_i - \mathbf{a}\|^2 \leq \xi_i \leq 0, -\xi_i \leq 0 \quad (69)$$

Therefore all points are out of the circle. the circle can become a point.

$C = \infty$ means, that all points are in the circle, because

$$\lim_{C \rightarrow \infty} \mu_i \Rightarrow \xi_i \rightarrow 0 \Rightarrow \forall i : \|\mathbf{x}_i - \mathbf{a}\|^2 \leq R^2 \quad (70)$$

2.8

Gaussian Kernel

$$k(x_i, x_j) = \exp\left(-\frac{1}{2\sigma^2}\right) \|\mathbf{x}_i - \mathbf{x}_j\|^2 \quad (71)$$

smaller bandwidth, so for $\sigma^2 \rightarrow 0$ two points the higher the kernel value. Data with less divination is weighted higher. As kernels express a function for similarity, this makes sense, that data with less divination, that are near by have higher 'weights'. Very small σ can lead to over fitting, because the Gaussian kernel is because of its construction likely to create areas around single points, which if they are outlays are over fitted.

2.9

$$\min_{\alpha, \mu, R} R^2 + C \sum_i^N \xi_i \quad (72)$$

$$\forall i : y_i(\|\mathbf{x}_i - \mathbf{a}\|^2 - R^2 - \xi_i) \leq 0, \xi_i \geq 0 \quad (73)$$

3 Neural Network

3.1

Do the forward pass to calculate the activations of all nodes and then determine the total error of the neural network. Let in the following $\sigma(\cdot)$ be the sigmoid function.

$$z1_{out} = \sigma(x1 * w1 + x2 * w2) \quad (74)$$

$$= 0.5664790559676278 \quad (75)$$

$$z2_{out} = \sigma(x1 * w3 + x2 * w4) \quad (76)$$

$$= 0.513746534902355 \quad (77)$$

$$y_{out} = \sigma(w5 * z1_{out} + w6 * z2_{out}) \quad (78)$$

$$= 0.6295403049912803 \quad (79)$$

$$error = E = -t * \ln(y_{out}) + (1 - t) \ln(1 - (y_{out})) \quad (80)$$

$$= 0.4627654005556673 \quad (81)$$

3.2

Apply back-propagation to obtain the partial derivatives of the error function with respect to the weights between the hidden layer and the output layer (w5 and w6). Then, use these derivatives to update the weights between the hidden layer and the output layer of the network using stochastic gradient descent (SGD) with a learning rate $\eta = 0.05$

$$w5' = w5 - lr * \frac{\partial E}{\partial w5} \quad (82)$$

$$= w5 - lr * \frac{\partial E}{\partial y_{out}} \frac{\partial y_{out}}{\partial y} \frac{\partial y}{\partial w5} \quad (83)$$

$$= w5 - lr * (-t/(y_{out}) * y_{out} * (1 - y_{out}) * z1_{out}) \quad (84)$$

$$= 0.8095161192315857 \quad (85)$$

$$w6' = w6 - lr * (-t/(y_{out}) * y_{out} * (1 - y_{out}) * z2_{out}) \quad (86)$$

$$= 0.16049288291512973 \quad (87)$$

3.3

Further apply back-propagation to obtain the partial derivatives of the error function with respect to the weights between the input layer and the hidden layer (w1, w2, w3 and w4). Again, use SGD to find the updates for these weights. Finally, calculate the total error for the neural network with the updated weights.

$$w1' = w1 - lr * \frac{\partial E}{\partial w1} \quad (88)$$

$$= w1 - lr * \frac{\partial E}{\partial z1_{out}} \frac{\partial z1_{out}}{\partial z1} \frac{\partial z1}{\partial w1} \quad (89)$$

$$\frac{\partial E}{\partial z1_{out}} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial z1_{out}} \quad (90)$$

$$\frac{\partial E}{\partial y} = \frac{\partial E}{\partial y_{out}} \frac{\partial y_{out}}{\partial y} \quad (91)$$

$$= -t/(y_{out}) * y_{out} * (1 - y_{out}) \quad (92)$$

$$\frac{\partial y}{\partial z1_{out}} = \frac{\partial(w5 * z1_{out} + w6 * z2_{out})}{\partial z1_{out}} \quad (93)$$

$$= w5 \quad (94)$$

$$\frac{\partial z1_{out}}{\partial z1} = z1_{out}(1 - z1_{out}) \quad (95)$$

$$\frac{\partial z1}{\partial w1} = \frac{\partial(w1 * x1 + w2 * x2)}{\partial w1} \quad (96)$$

$$= x1 \quad (97)$$

$$w1' = w1 - lr * -t/(y_{out}) * y_{out} * (1 - y_{out}) * w5 * z1_{out}(1 - z1_{out}) * x1 \quad (98)$$

$$= 0.4003639107605591 \quad (99)$$

$$\Rightarrow \text{doing this for all other weights gives us:} \quad (100)$$

$$w2' = w2 - lr * -t/(y_{out}) * y_{out} * (1 - y_{out}) * w5 * z2_{out}(1 - z2_{out}) * x1 \quad (101)$$

$$= 0.6503701796760513 \quad (102)$$

$$w3' = w3 - lr * -t/(y_{out}) * y_{out} * (1 - y_{out}) * w6 * z1_{out}(1 - z1_{out}) * x2 \quad (103)$$

$$= 0.2002388164366169 \quad (104)$$

$$w4' = w4 - lr * -t/(y_{out}) * y_{out} * (1 - y_{out}) * w6 * z2_{out}(1 - z2_{out}) * x2 \quad (105)$$

$$= 0.10024293041240866 \quad (106)$$

Apply feed forward again with updated values:

$$z1'_{out} = \sigma(x1 * w1' + x2 * w2') \quad (107)$$

$$= 0.566519810580869 \quad (108)$$

$$z2'_{out} = \sigma(x1 * w3' + x2 * w4') \quad (109)$$

$$= 0.5137737411036912 \quad (110)$$

$$y'_{out} = \sigma(w5' * z1'_{out} + w6' * z2'_{out}) \quad (111)$$

$$= 0.632059884290427 \quad (112)$$

$$error' = -t * \ln \sigma(y'_{out}) + (1 - t) \ln(1 - \sigma(y'_{out})) \quad (113)$$

$$= 0.4587711357000273 \quad (114)$$

This gives an improve of error of: 0.003994264855639973