1. P(~JohnCalls, MaryCalls, Alarm,~Earthquake, Burglary)
   = P(~J, M, A, ~E, B)                                        //simplify variable names
   = P(~J|Parent(J))P(M|Parent(M))P(A|Parent(A))P(~E)P(B)
   = P(~J|A)P(M|A)P(A|~E^B)P(~E)P(B)
   = 0.1 * 0.7 * 0.94 * 0.998 * 0.001
   = 0.0000656684
2.



   Node ordering is Alarm, JohnCalls, Earthquake

| Edges | Explanation |
|---|---|
| Alarm→JohnCalls | P(JohnCalls\|Alarm) != P(JohnCalls)<br>In the original belief network, it is clear that the Alarm has an effect on whether or not JohnCalls. Therefore, this edge is required to satisfy conditional independence. |
| Alarm→Earthquake | Conditional Independence:<br>P(Earthquake\|Alarm, JohnCalls) = P(Earthquake\|Alarm)<br>Based on the original belief network, whether or not john calls is based on the alarm going off. So whether or not JohnCalls does not affect P(Earthquake\|Alarm). That is also why there is not an edge from JohnCalls to Earthquake. |

3. P(Burglary|JohnCalls) > P(Burglary|JohnCalls, Earthquake)
   => P(B|J) > P(B|J,E)                                        //simplify variable names
   => B is a cause, J is an effect, and E is another cause
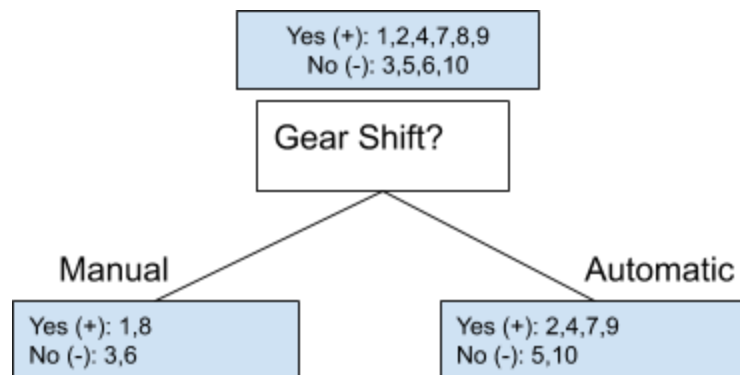   => This is Intercausal inferences (explaining away)
   Therefore, the inequality holds because the "cause is already found"
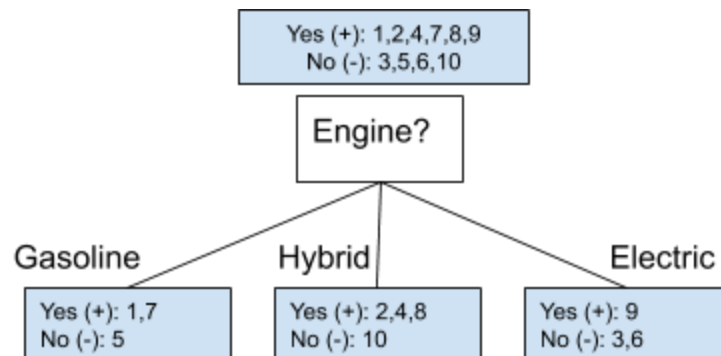4. Depth-one decision trees

a. Type

Yes (+): 1,2,4,7,8,9
No (-): 3,5,6,10

Type?

Sports

Yes (+): 1,7,9
No (-): 3,10

SUV

Yes (+): 2,4,8
No (-):

Mini Van

Yes (+):
No (-): 5, 6

b. Gear Shift

Yes (+): 1,2,4,7,8,9
No (-): 3,5,6,10

Gear Shift?

Manual

Yes (+): 1,8
No (-): 3,6

Automatic

Yes (+): 2,4,7,9
No (-): 5,10

c. Engine

Yes (+): 1,2,4,7,8,9
No (-): 3,5,6,10

Engine?

Gasoline

Yes (+): 1,7
No (-): 5

Hybrid

Yes (+): 2,4,8
No (-): 10

Electric

Yes (+): 9
No (-): 3,6

5. Information gain: Entropy(parent) - [average entropy(children)]

Entropy(parent) = sum for i in C [ $-P_i * \log_2(P_i)$]

= $-P_{yes} * \log_2(P_{yes}) + -P_{no} * \log_2(P_{no})$

= $-0.6 * \log_2(0.6) + -0.4 * \log_2(0.4) = 0.971$

   a. Gain(E, Type?) = 0.971 - sum($|E_v| * Entropy(E_v)/|E|$)

     = 0.971 -

     ($E_{sports} * Entropy(E_{sports}) + E_{suv} * Entropy(E_{suv}) + E_{van} * Entropy(E_{suv}))/|E|$

     = 0.971 - (5*0.971 + 3*0 + 2*0)/10

     = **0.486**

b. Gain(E, Gear Shift?) = 0.971 - sum(|Ev| * Entropy(Ev)/|E|)
= 0.971 - (4*1 + 6*0.918)/10
= **0.020**
c. Gain(E, Engine?) = 0.971 - sum(|Ev| * Entropy(Ev)/|E|)
= 0.971 - (3*0.918 + 4*0.811 + 3*0.918)/10
= **0.0958**

6. I would choose the "Type" attribute first because it has the maximal information gain out of all attributes as calculated above. In terms of a decision tree, we want to make as few tests before reaching a decision, so by maximizing information gain we decrease entropy the most with each test and ideally decrease the number of tests before reach a decision.

7. (1) Program output for "Info" was 0.485475, "Gear Shift" was 0.0199731, and "Engine" was 0.0954618. This confirms my calculations.
(2) Program output for attributes table slide06, page 9

| Attribute | Info Gain |
| --- | --- |
| Alt | 0.0000 |
| Bar | 0.0000 |
| Fri | 0.0207208 |
| Hun | 0.19571 |
| Pat | 0.540852 |
| Price | 0.19571 |
| Rain | 0.0000 |
| Res | 0.0207208 |
| Type | 0.0000 |
| Est | 0.207519 |

8. (1) Yes
(2)

Coordinates (x, y, z)
Red: Class = 1
Black: Class = 0
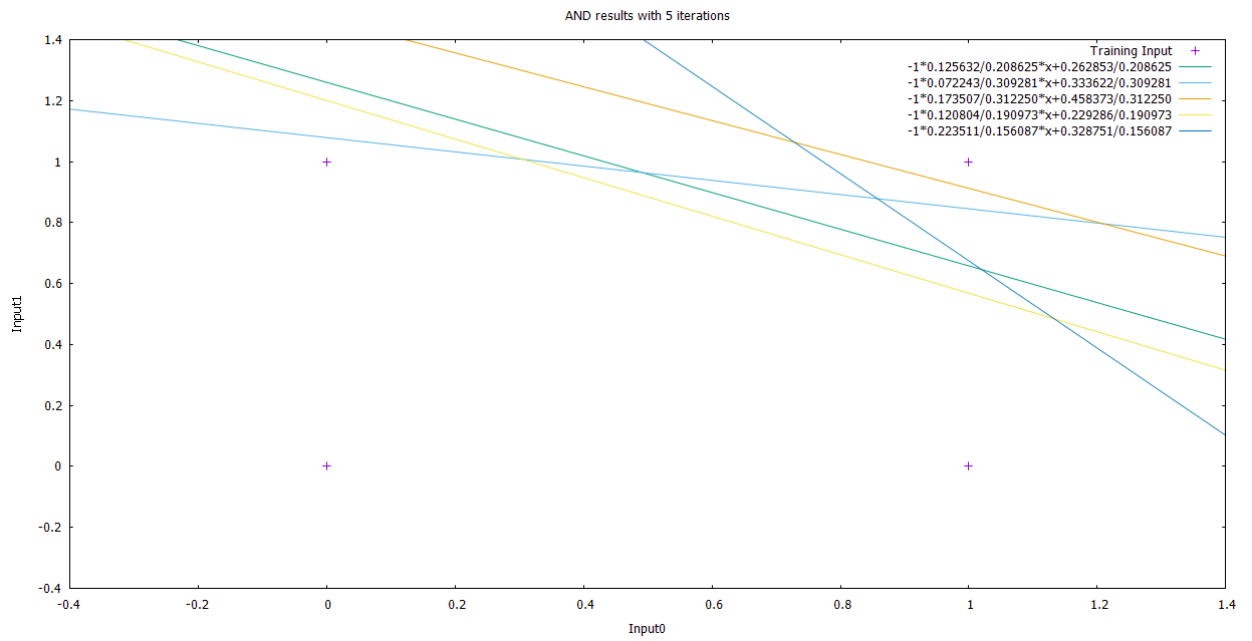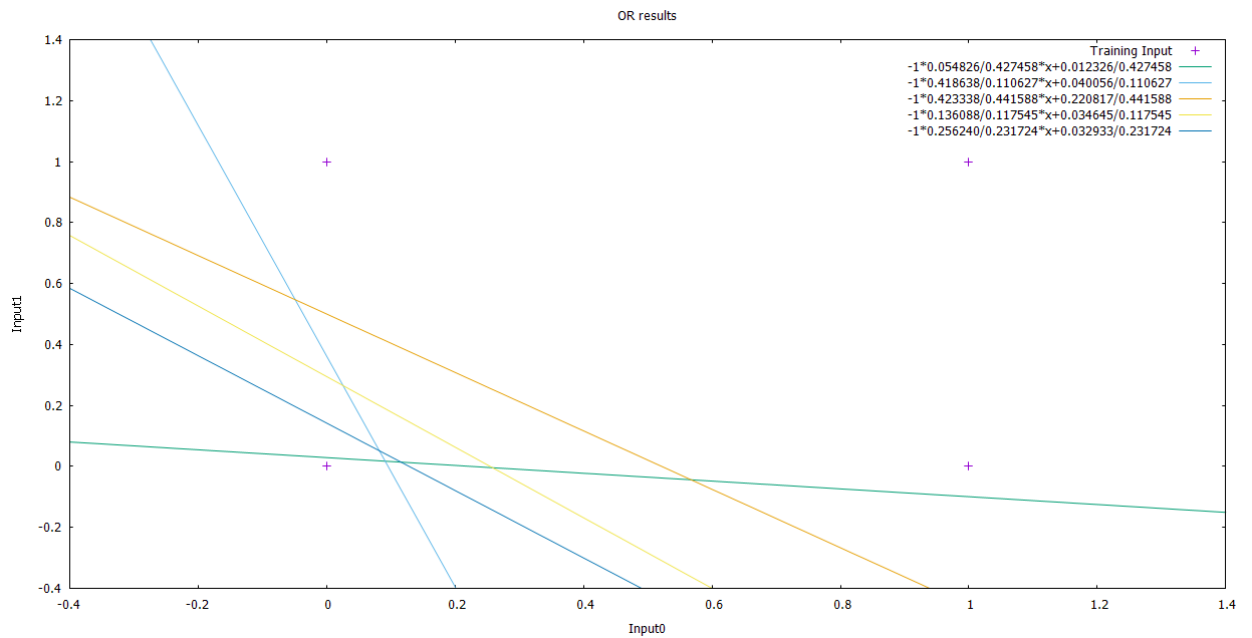
(3) A single perceptron unit can solve this classification problem because the points are linearly-separable. We can draw a flat plane to separate red and black dots. The plane would be parallel to the line created by the red dots, and at an angle. This drawing illustrates this example:



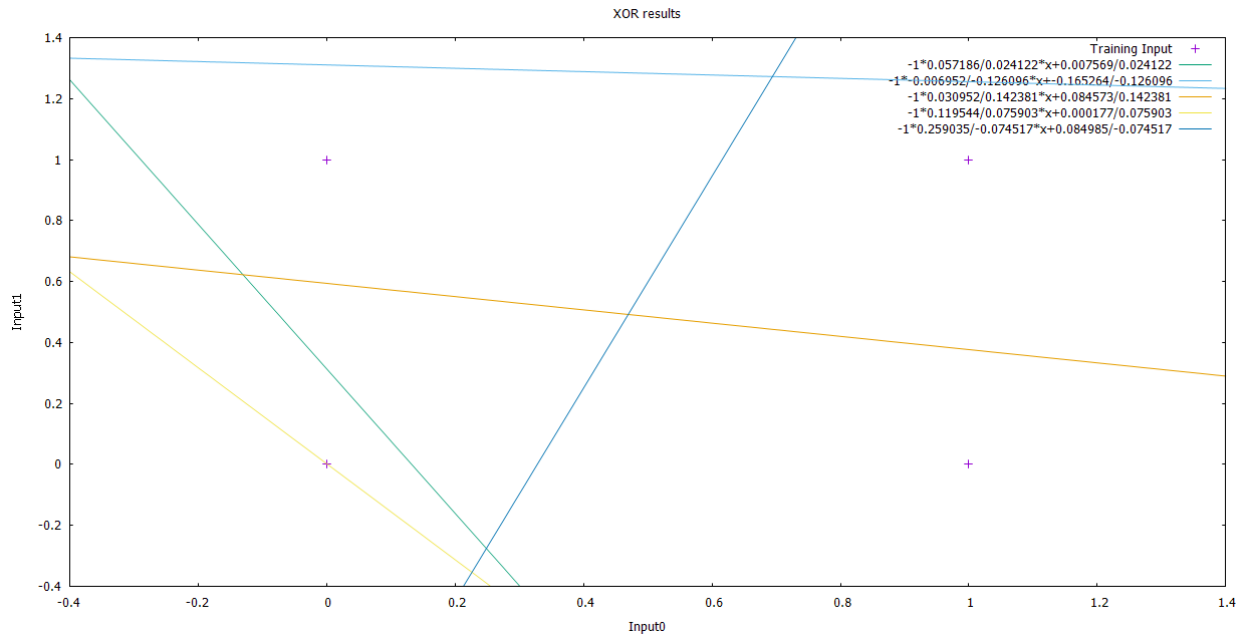Coordinates (x, y, z)
Red: Class = 1
Black: Class = 0

9. (2) You can also view these pictures under the project director "question9/results".
   - note that plotting was done with Gnuplot and on Windows machine. In order to run the C++ code you must have Gnuplot installed and using Windows OS.
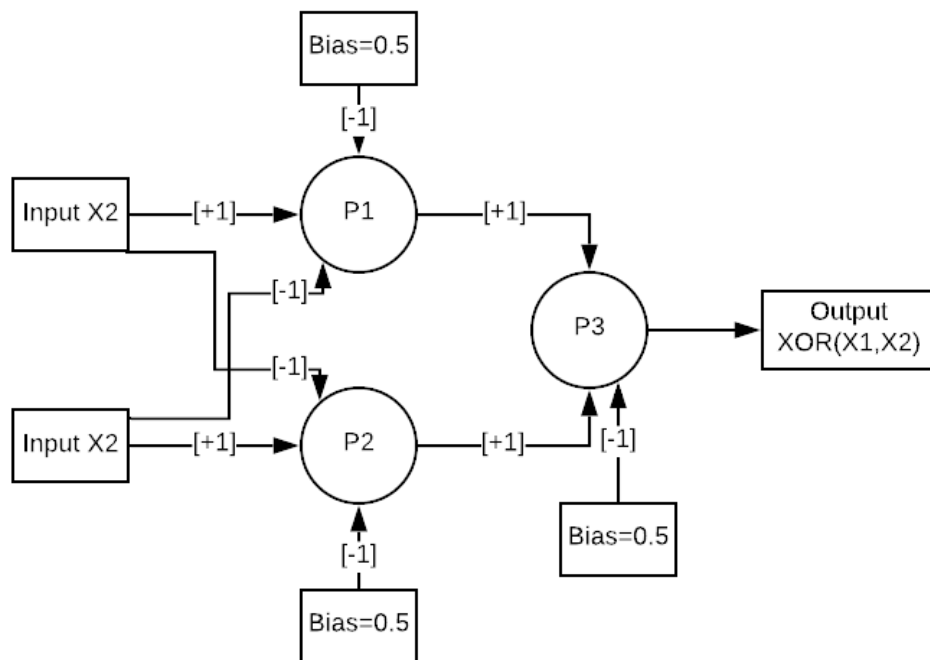   AND Results

AND results with 5 iterations

Training Input +
-1*0.125632/0.208625*x+0.262853/0.208625
-1*0.072243/0.309281*x+0.333622/0.309281
-1*0.173507/0.312250*x+0.458373/0.312250
-1*0.120804/0.190973*x+0.229286/0.190973
-1*0.223511/0.156087*x+0.328751/0.156087

OR results:



OR results

Training Input +
-1*0.054826/0.427458*x+0.012326/0.427458
-1*0.418638/0.110627*x+0.040056/0.110627
-1*0.423338/0.441588*x+0.220817/0.441588
-1*0.136088/0.117545*x+0.034645/0.117545
-1*0.256240/0.231724*x+0.032933/0.231724

XOR Results (100 iterations max)

XOR results

Legend:
- Training Input +
- -1*0.057186/0.024122*x+0.007569/0.024122
- -1*-0.006952/-0.126096*x+-0.165264/-0.126096
- -1*0.030952/0.142381*x+0.084573/0.142381
- -1*0.119544/0.075903*x+0.000177/0.075903
- -1*0.259035/-0.074517*x+0.084985/-0.074517
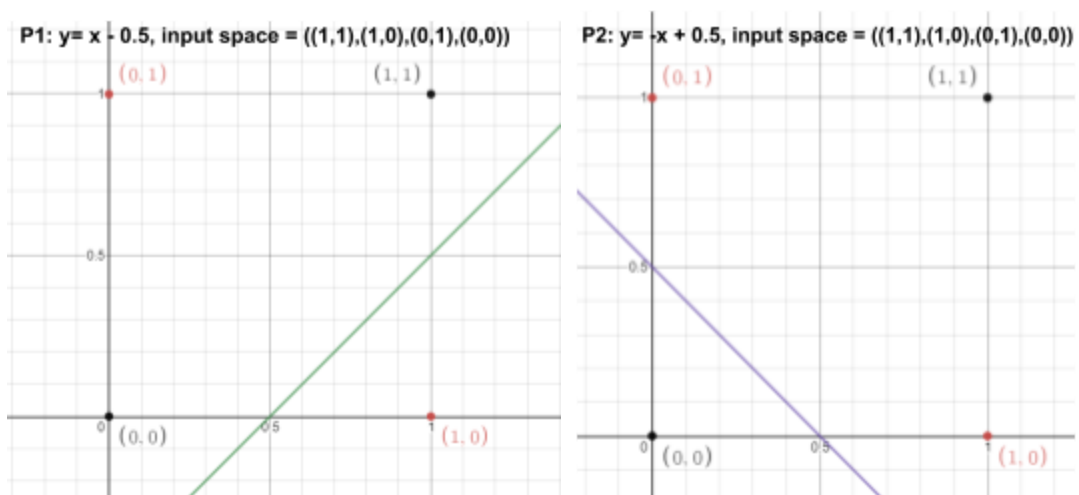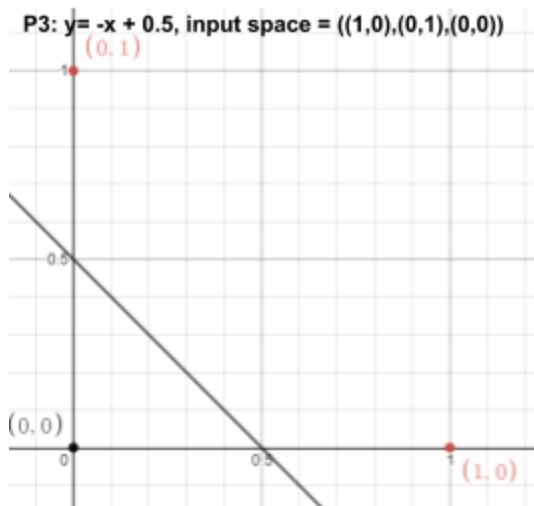
Axes: Input0 (horizontal), Input1 (vertical)

(3) For AND and OR, it is possible. The weights and bias are initialized randomly between -0.5 and 0.5. This means that the weights for AND can be randomly intialized as bias = -0.5 and weight1 = 0.3 and weight2 = 0.3, which will result in no errors. Similarly, if the weights and bias for OR are randomly initialized as bias = -0.3 and weight1 = 0.4 and weight2 = 0.4, which will also result in no errors.

This is heavily implementation-dependent. Note that if the weights and bias were only initialized randomly to a positive range like [0, 1], then it is not possible because all and biases weights will be greater than or equal to 0, which will always result in an output of 1!

## 10. (1) Neural Network Topology



## (2) Decision boundary

P3: y= -x + 0.5, input space = ((1,0),(0,1),(0,0))

11. Given error function E(w) = (w+1)(w-1)(w-3)(w-4)

    => E(w) = w^4 - 7 w^3 + 11 w^2 + 7 w - 12            // simplify

    => dE/dw = 4w^3 - 21w^2 + 22w + 7             // deriviative

    **=> Δw = −α * (4w^3 - 21w^2 + 22w + 7)**        // plug into Δw = −α*dE/dw

12. Plots:

    The pictures can also be found under "question12/results"

    w = 6.0 // this is a fixed value

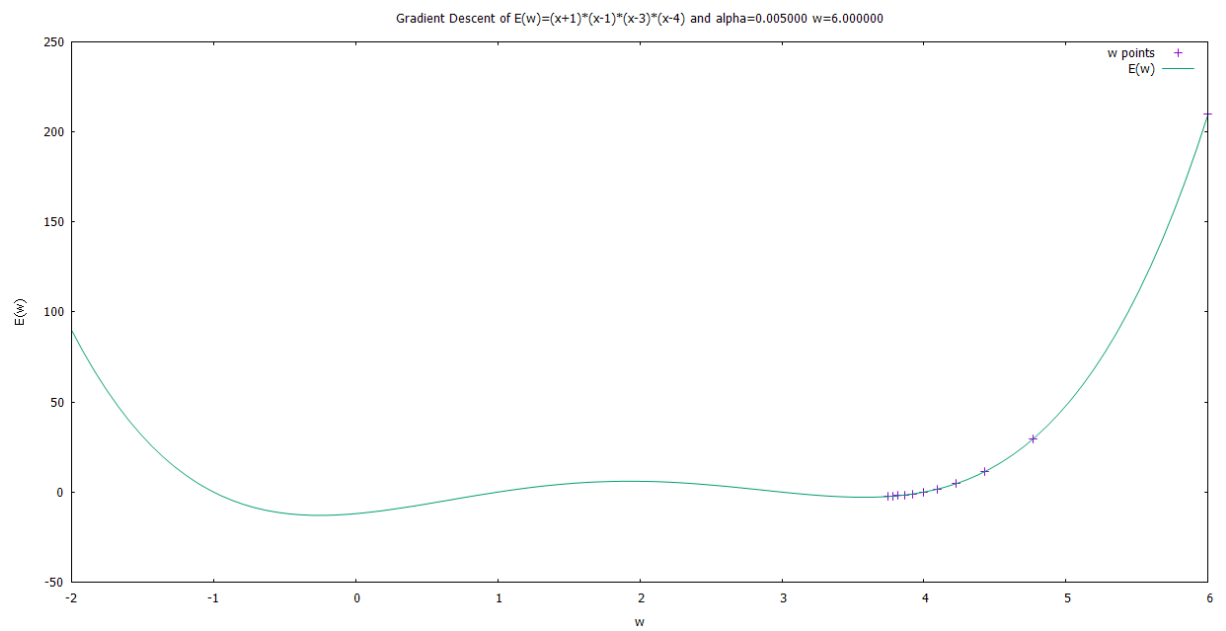    Experimented with different learning rates = [0.001, 0.002, 0.005, 0.01]

    Graph 1: alpha = 0.001
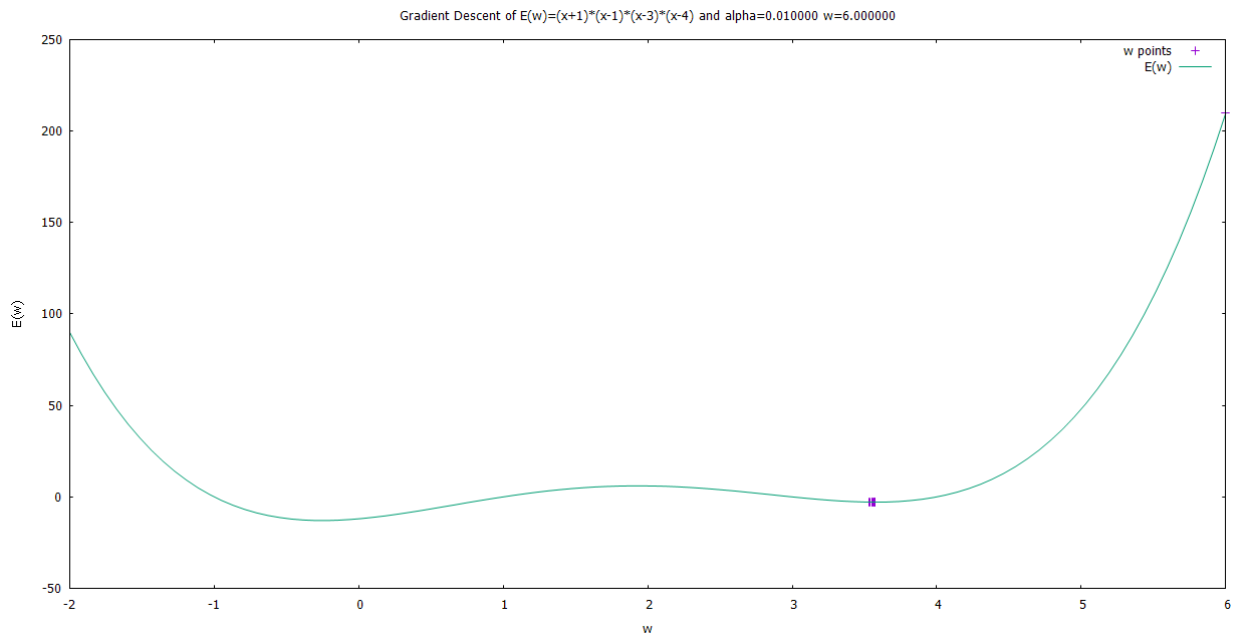


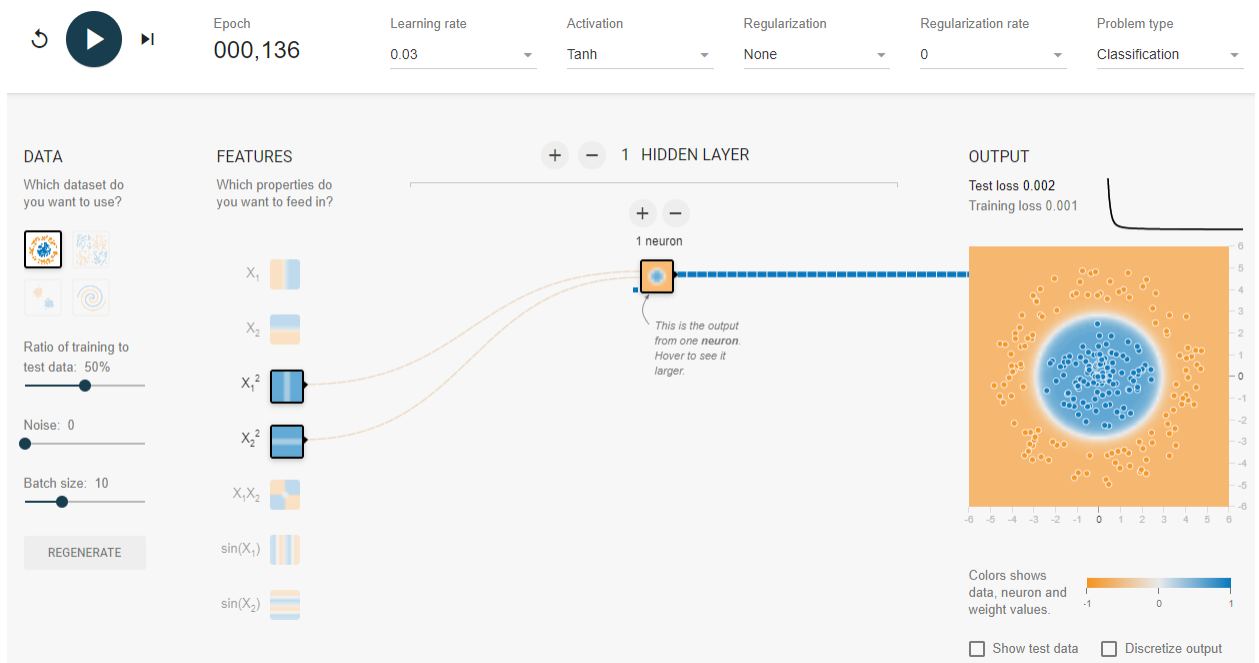Gradient Descent of E(w)=(x+1)*(x-1)*(x-3)*(x-4) and alpha=0.001000 w=6.000000

Graph 2: alpha = 0.002

Gradient Descent of E(w)=(x+1)*(x-1)*(x-3)*(x-4) and alpha=0.002000 w=6.000000



Graph 3: alpha = 0.005

Gradient Descent of E(w)=(x+1)*(x-1)*(x-3)*(x-4) and alpha=0.005000 w=6.000000

## Graph 4: alpha = 0.01



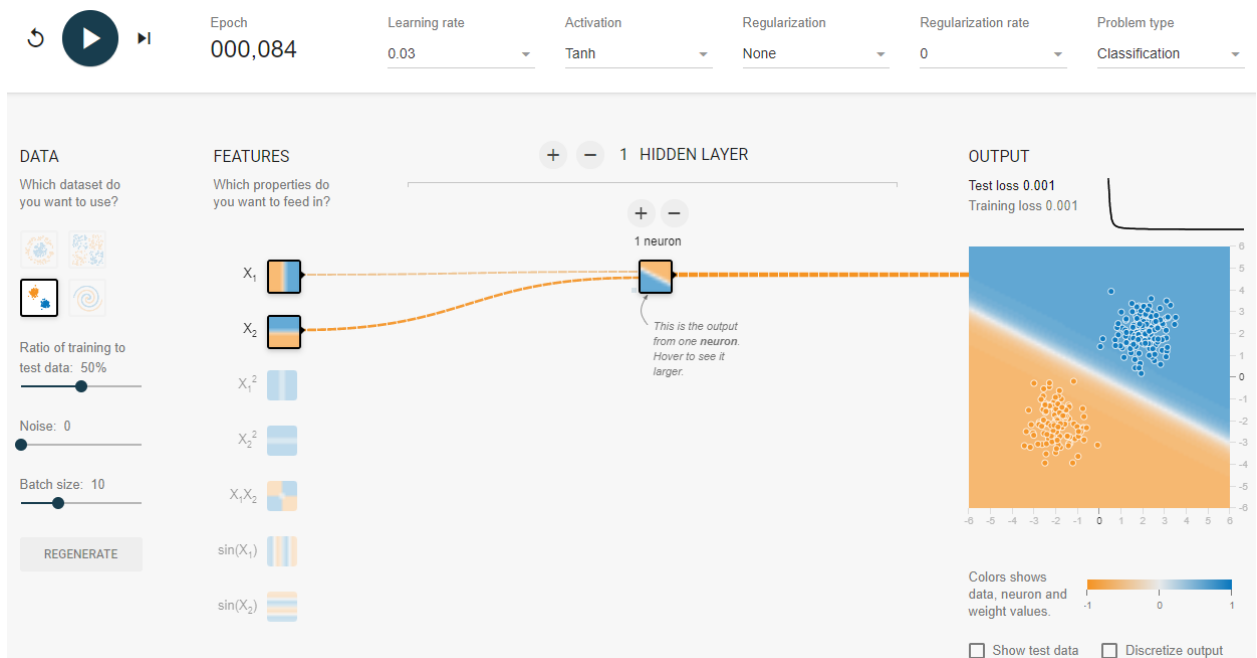Gradient Descent of E(w)=(x+1)*(x-1)*(x-3)*(x-4) and alpha=0.010000 w=6.000000

## 13. Dataset 1:



Dataset 2:

Dataset 3:



Dataset 4: Only x1, x2 used

Dataset 4: only x1,x2,x1^2,x2^2 used

14. (3):



1D SOM w1->w2->...->w7 Plotted

(4): input vector x = (4,1)
=> best matching is w6 (3,0)
    w6new = (3,0) + 1 * (1, 1) = (4, 1)
=> immidiate neighbors are w(6-1) and w(6+1) => w5, w7
    w5 = (5, 4) + ⅔ * (-1, -3) = (4.333, 2)
    w7 = (-1, 5) + ⅔ * (5, -4) = (2.333, 2.333)
=> second order neighbor is w(6-2) and w(6+2) => w4, w8
    w4 = (6, 2) + ⅓ * (-2, -1) = (5.333, 1.666)
    w8 doesn't exist
=> because rest h() value is 0, those weight vectors don't move!

The graph below shows the red line (before one iteration of training) and the blue line (after one iteration of training). The green dot is the input vector added.