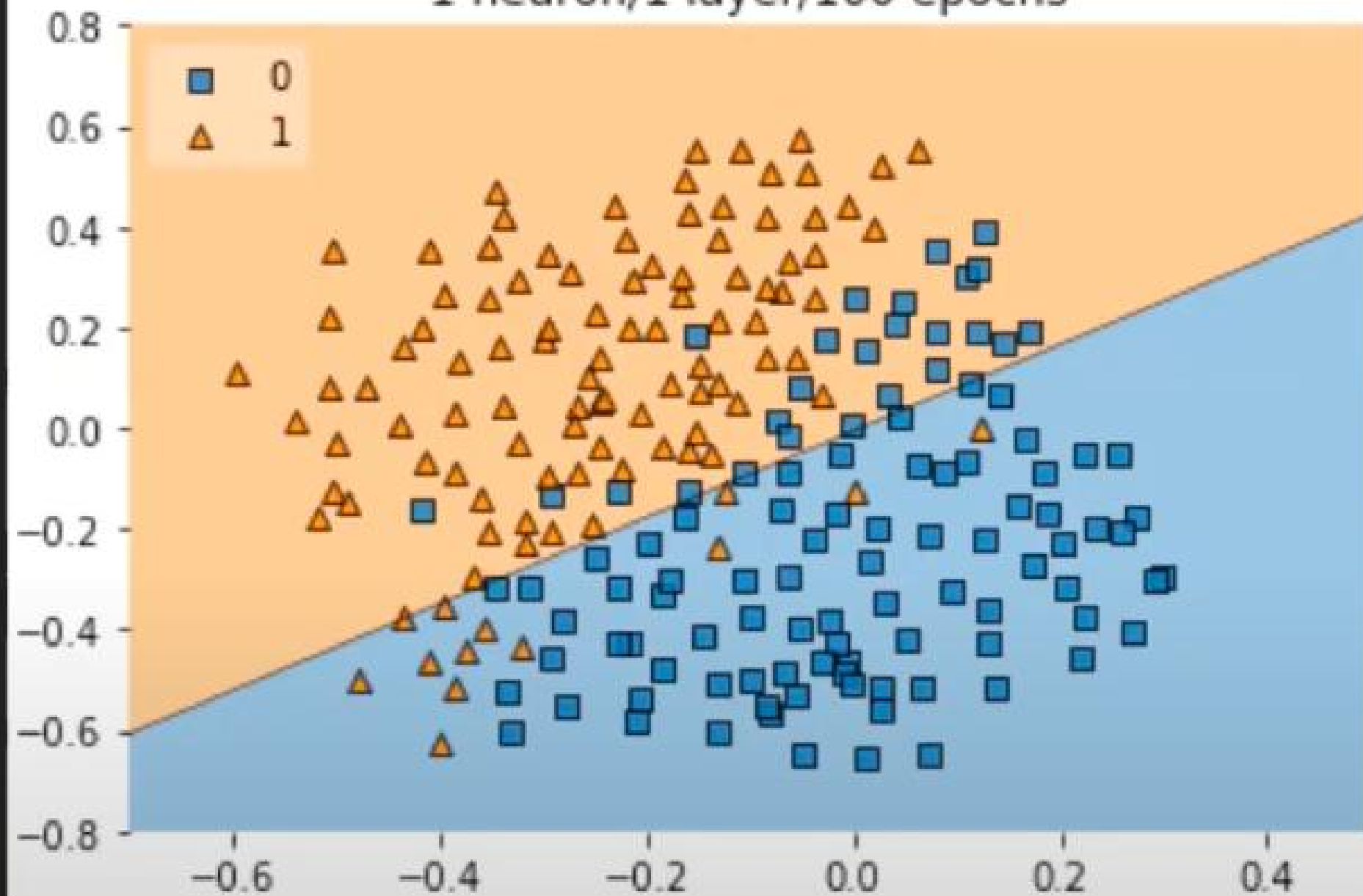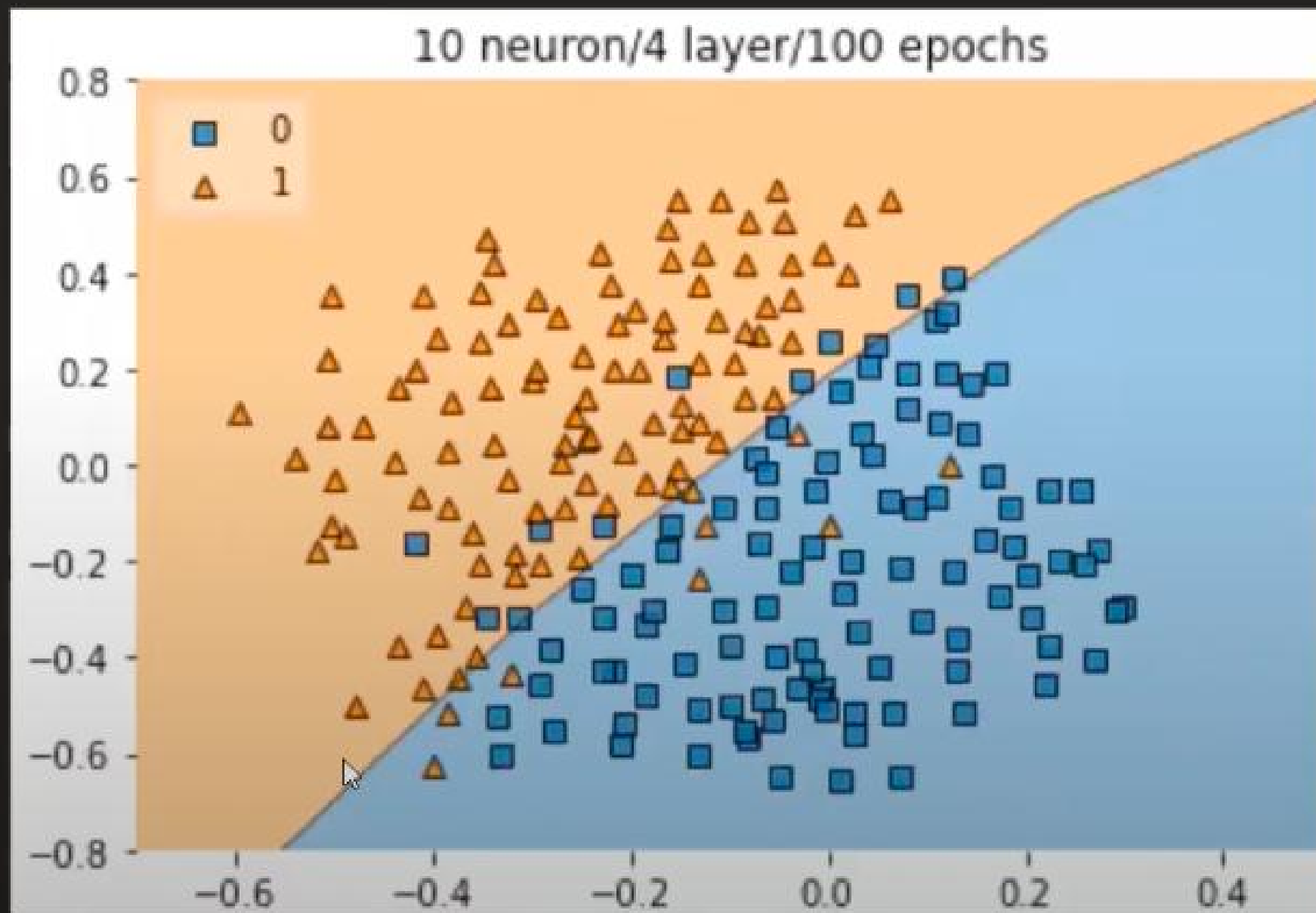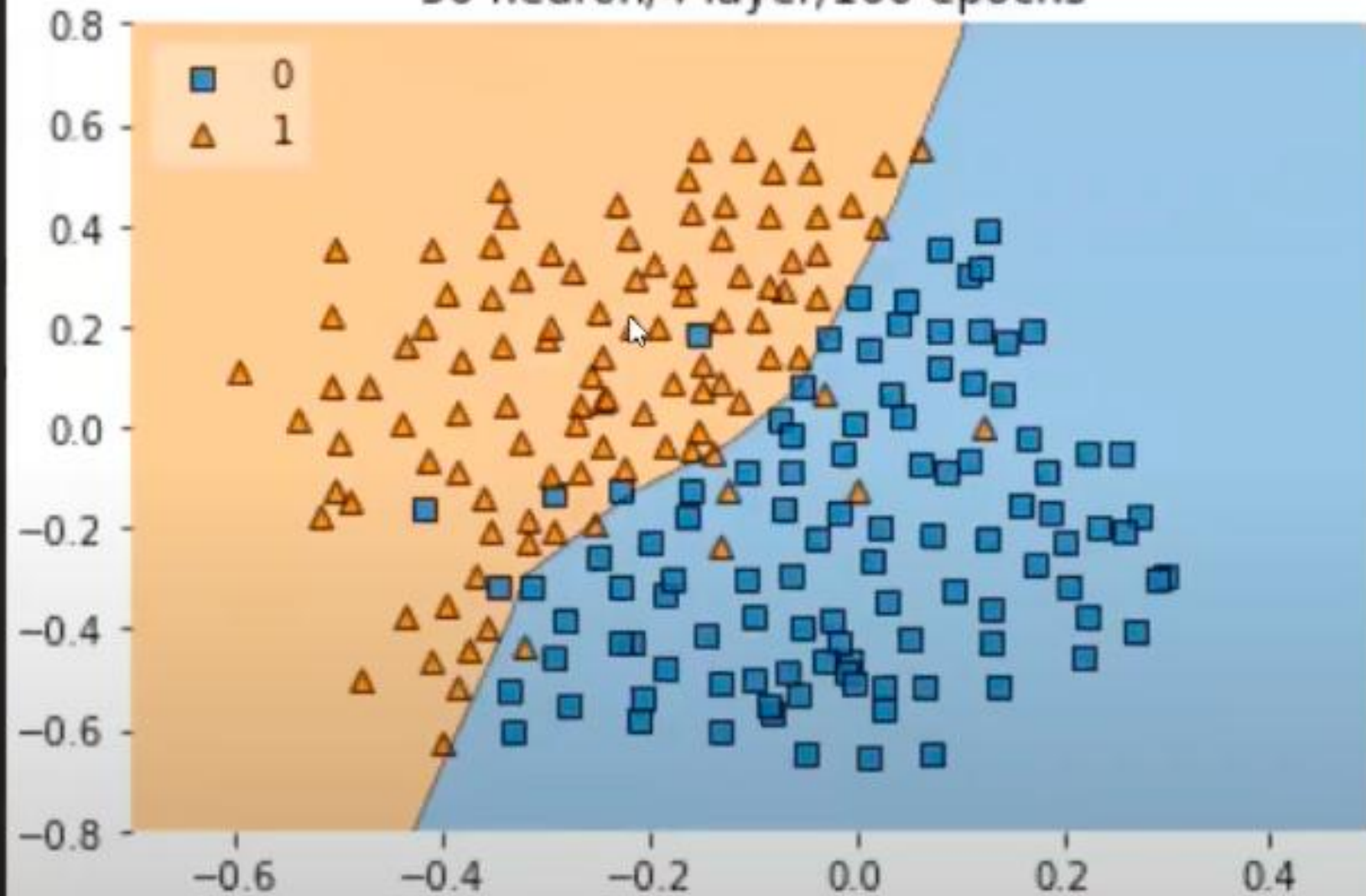1 neuron/1 layer/100 epochs
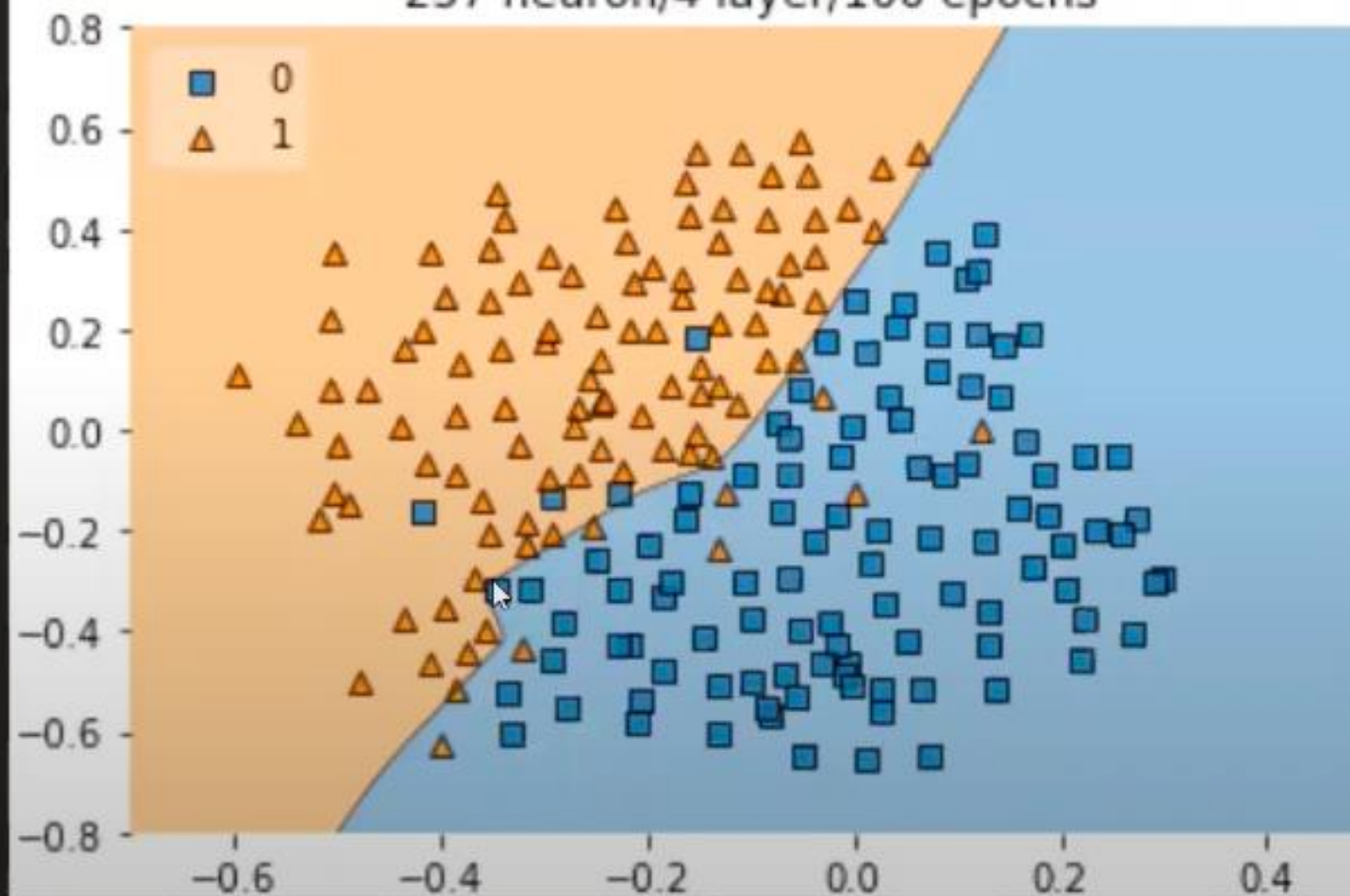
10 neuron/4 layer/100 epochs
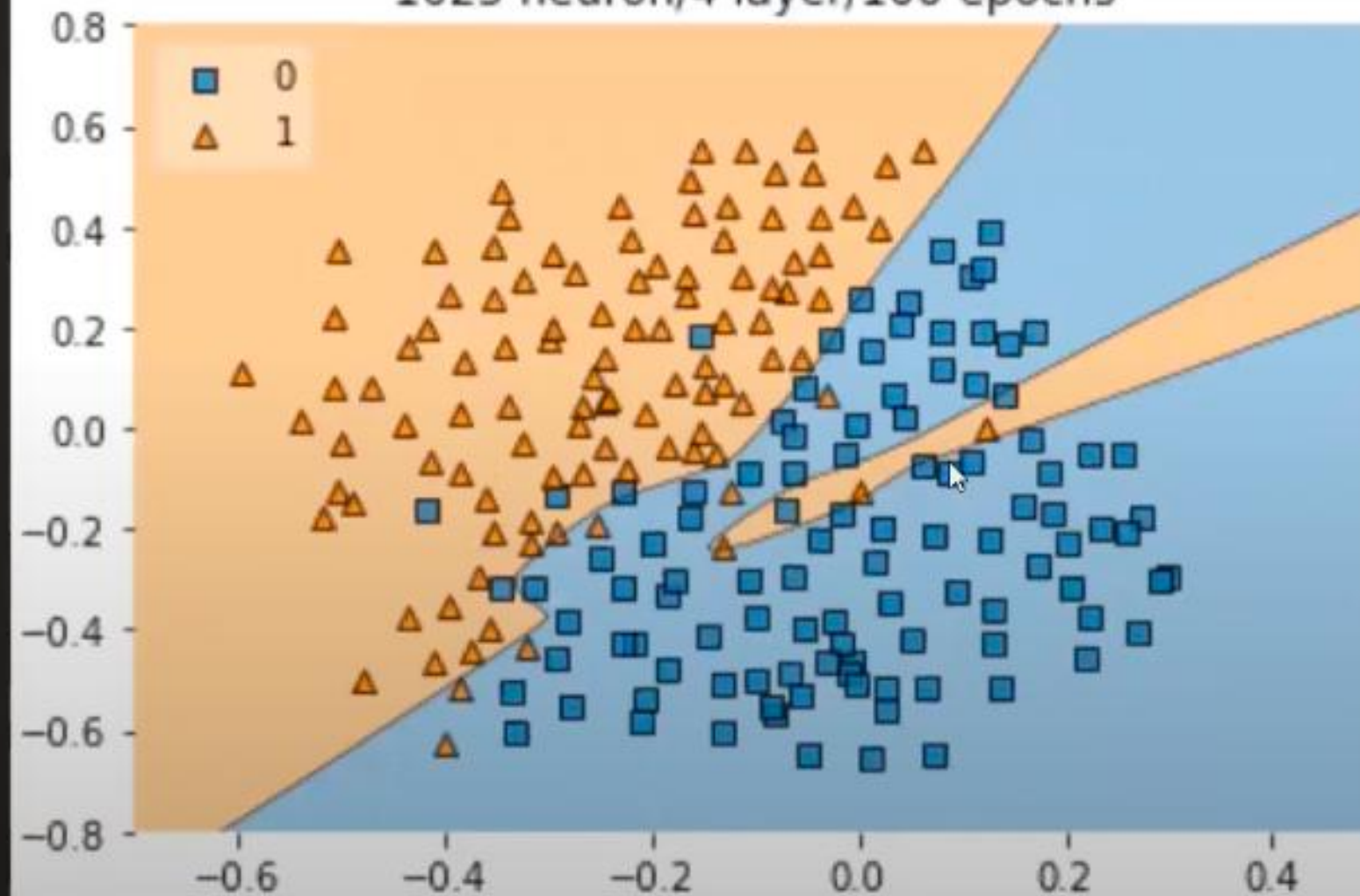
50 neuron/4 layer/100 epochs

257 neuron/4 layer/100 epochs

1025 neuron/4 layer/100 epochs

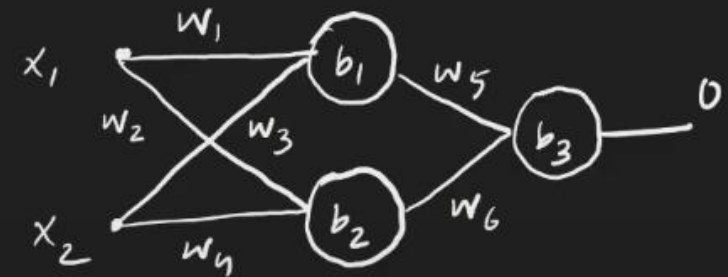# What are Activation Functions?

In **artificial neural networks**, each neuron forms a weighted sum of its inputs and passes the resulting scalar value through a function referred to as an activation function or transfer function. If a neuron has n inputs then the output or activation of a neuron is

$$a = g(w_1 x_1 + w_2 x_2 + w_3 x_3 + \ldots w_n x_n + b)$$
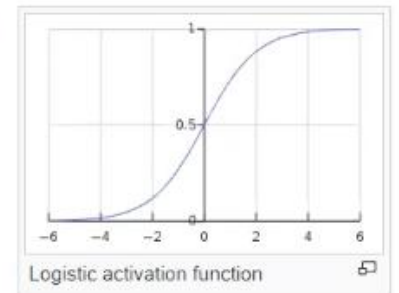
This function g is referred to as the activation function.

# Activation function

*For the formalism used to approximate the influence of an extracellular electrical field on neurons, see activating function. For a linear system's transfer function, see transfer function.*

In artificial neural networks, each neuron forms a weighted sum of its inputs and passes the resulting scalar value through a function referred to as an activation function or transfer function. If a neuron has n inputs $x_1, x_2, \ldots x_n$ then the output or activation of a neuron is $a = g(w_1 x_1 + w_2 x_2 + w_3 x_3 + \ldots w_n x_n + b)$. This function g is referred to as the activation function. If the function g is taken as the linear function $g(z) = z$ then the neuron performs linear regression or classification. In general g is taken to be a nonlinear function to do nonlinear regression and solve classification problems that are not linearly separable. When g is taken to be a sigmoidal or 's' shaped function varying from 0 to 1 or -1 to 1, the output value of the neuron can be interpreted as a YES/NO answer or binary decision. However saturating activation function can cause the vanishing gradient problem in deep networks. Replacing saturating sigmoidal activation functions with activation functions like ReLU that have larger derivative values allowed deeper networks to be trained for the first time. Non-monotonic and oscillating activation functions that significantly outperform ReLU have since been found.[1] In particular oscillating activation functions improve gradient flow, speedup training and allow single neurons to learn the XOR function like certain human cerebral neurons.[2][3]
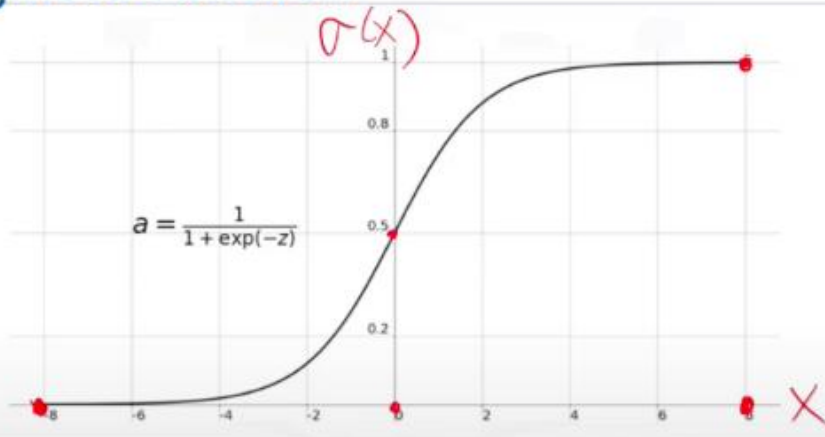


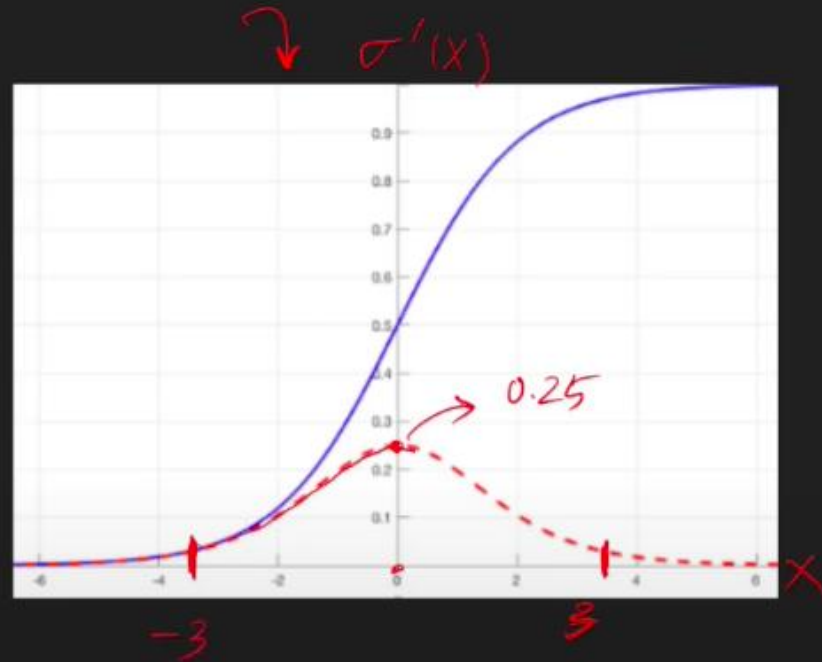Logistic activation function

**Contents** [hide]

# Sigmoid Activation Function

31 May 2022     14:50



Sigmoid Function

$\sigma(x)$

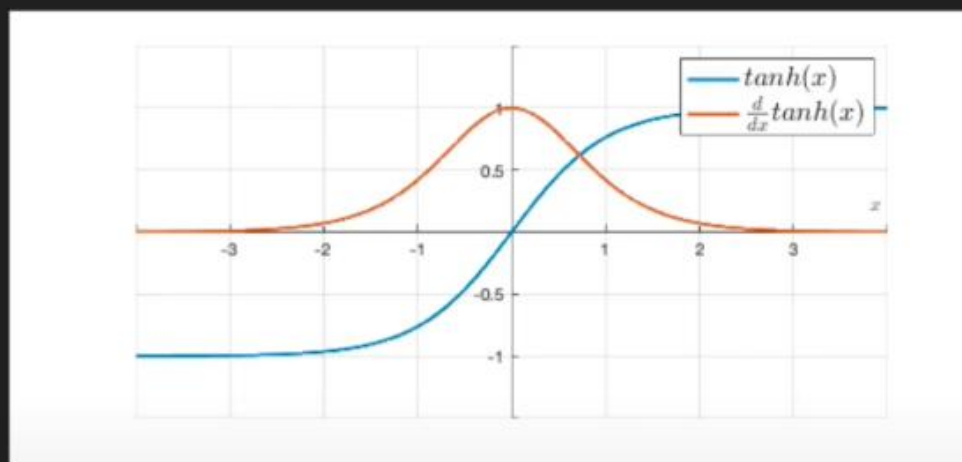$$a = \frac{1}{1 + \exp(-z)}$$

$\sigma'(x)$

0.25

−3     3
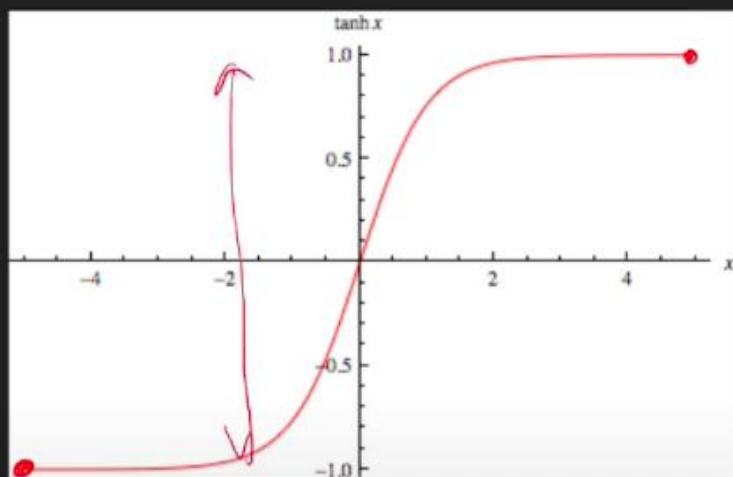
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

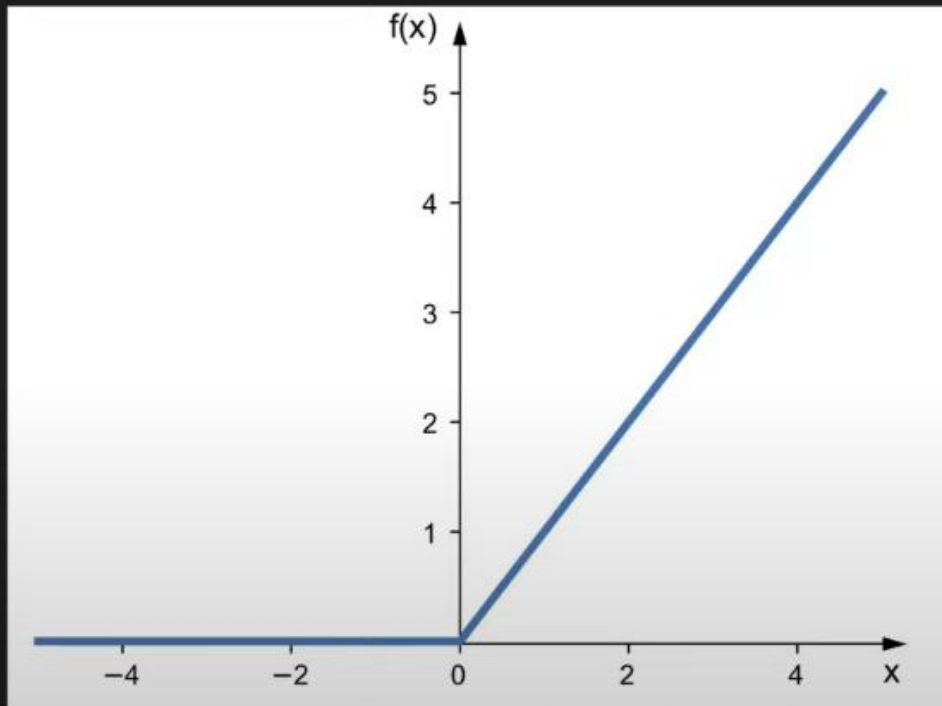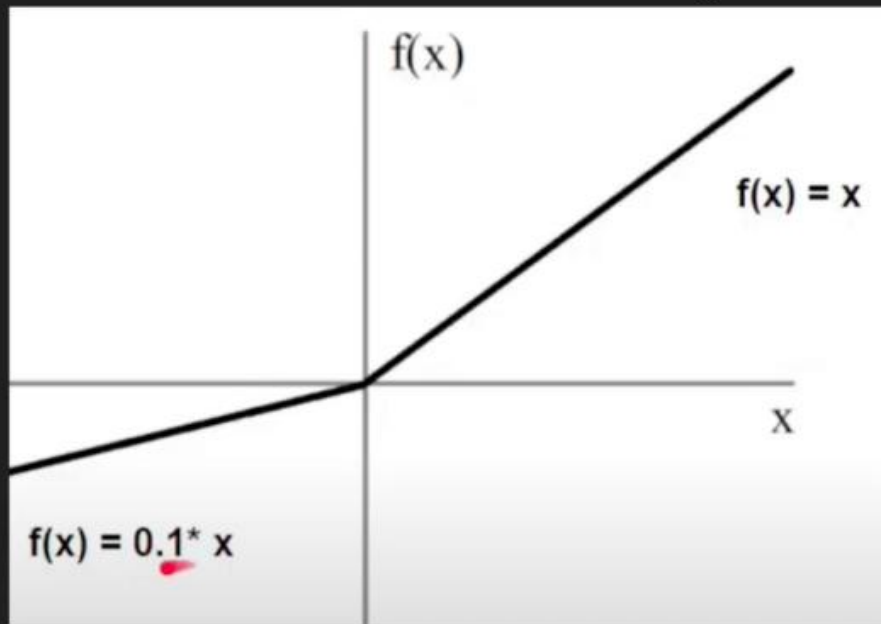# Tanh Activation Function

31 May 2022    14:50

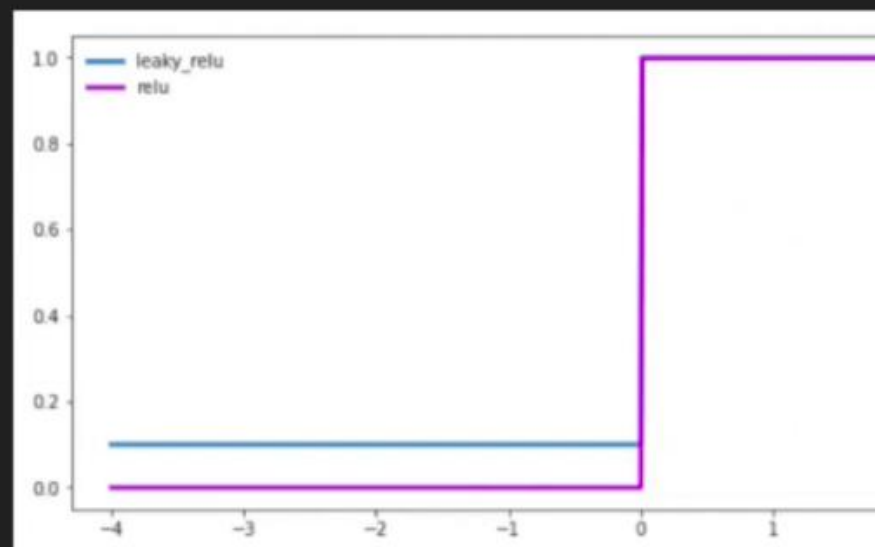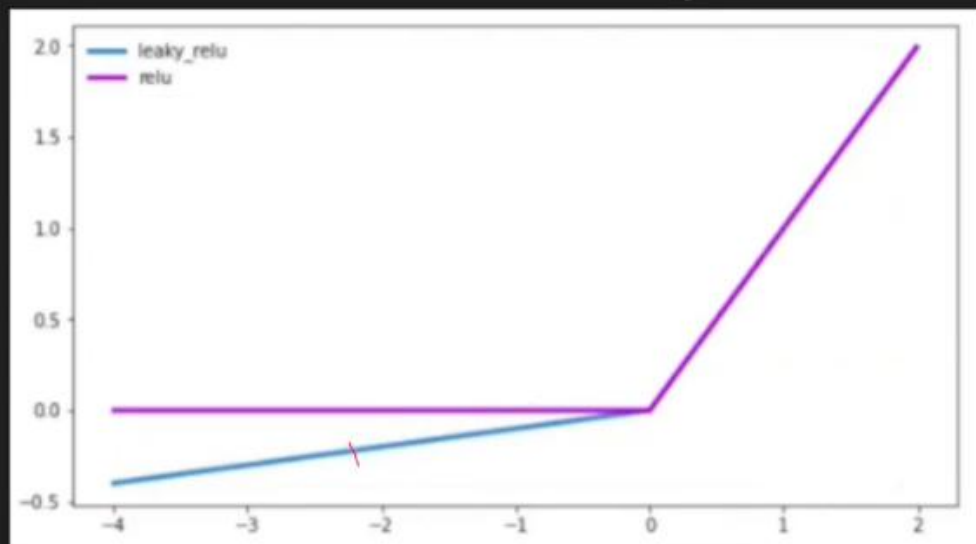# Relu Activation Function

# Leaky Relu

08 June 2022    22:53



$$f(z) = \max(0.01z, z)$$

$$z \geqslant 0 \rightarrow z$$

$$z < 0 \rightarrow \frac{1}{100}z \quad (\text{fraction of } z)$$
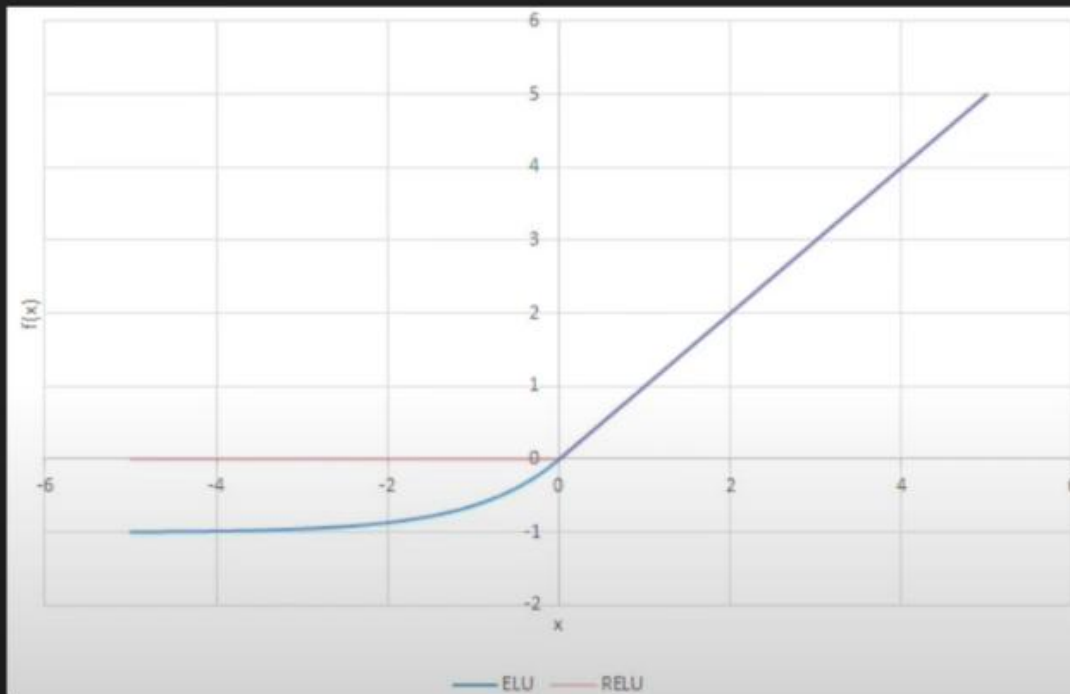
# Parametric Relu

f(x)

f(x) = x

x

f(x) = α*x

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ ax & \text{otherwise} \end{cases}$$

# Elu - Exponential Linear Unit

09 June 2022      00:29



$$ELU(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } \overline{x} < 0 \end{cases}$$

$$ELU'(x) = \begin{cases} 1 & \text{if } x > 0 \\ ELU(x) + \alpha & \text{if } x \leq 0 \end{cases}$$

ELU ——— RELU

# Selu - Scaled Exponential Linear Unit

09 June 2022    00:29

$$\text{SELU}(x) = \lambda \begin{cases} x & \text{if } x > 0 \\ \alpha e^x - \alpha & \text{if } x \leq 0 \end{cases}$$

$$a \approx 1.6732632423543772848170429916717$$

$$\lambda \approx 1.0507009873554804934193349852946$$

$$\text{SELU}'(x) = \lambda \begin{cases} 1 & \text{if } x > 0 \\ \alpha e^x & \text{if } x \leq 0 \end{cases}$$

# Why Weight Initialization is Important?

22 June 2022    12:48



1. Initialize the parameters
2. Choose an *optimization algorithm*
3. Repeat these steps:

    1. Forward propagate an input
    2. Compute the cost function
    3. Compute the gradients of the cost with respect to parameters using backpropagation
    4. Update each parameter using the gradients, according to the optimization algorithm

$$\boxed{\frac{\partial L}{\partial w'_{11}}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{11}} \frac{\partial a_{11}}{\partial z_{11}} X_1 \qquad \boxed{\frac{\partial L}{\partial w'_{12}}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{12}} \frac{\partial a_{12}}{\partial z_{12}} X_1$$

$$\frac{\partial L}{\partial w'_{21}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{11}} \frac{\partial a_{11}}{\partial z_{11}} X_2 \qquad \frac{\partial L}{\partial w'_{22}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{12}} \frac{\partial a_{12}}{\partial z_{12}} X_2$$

Random Initialization (Large values)

Batch-Normalization (BN) is an algorithmic method which makes the training of Deep Neural Networks (DNN) faster and more stable.

It consists of normalizing activation vectors from hidden layers using the mean and variance of the current batch. This normalization step is applied right before (or right after) the nonlinear function.

# Covarite Shift

# Keras Implementation

```python
model = Sequential()

model.add(Dense(3,activation='relu',input_dim=2))
model.add(BatchNormalization())
model.add(Dense(2,activation='relu'))
model.add(BatchNormalization())
model.add(Dense(1,activation='sigmoid'))
```

# Role of Optimizer

| cgpa | iq | placed |
|------|-----|--------|
| 8 | 80 | 1 |
| 9 | 90 | 1 |
| 7 | 70 | 0 |

$x_1$

$x_2$

4    2    2    1

$$\boxed{9 \text{ parameters}}$$

Optimization

$$\dfrac{L \ min}{y \quad \hat{y}} \quad \boxed{(w, b)}$$

$\hat{y}$

## Mathematical Intuition

$$V_t = \beta V_{t-1} + (1-\beta)\theta_t$$

$$V_0 = 0$$

$$V_1 = (1-\beta)\theta_1$$

$$V_2 = \beta V_1 + (1-\beta)\theta_2$$

$$= \beta(1-\beta)\theta_1 + (1-\beta)\theta_2$$

$$V_3 = \beta V_2 + (1-\beta)\theta_3$$

$$= \beta^2(1-\beta)\theta_1 + \beta(1-\beta)\theta_2 + (1-\beta)\theta_3$$

$$\begin{cases} V_4 = \beta V_3 + (1-\beta)\theta_4 \\ \\ V_4 = \beta^3(1-\beta)\theta_1 + \beta^2(1-\beta)\theta_2 + \beta(1-\beta)\theta_3 + (1-\beta)\theta_4 \\ \\ = (1-\beta)\left[\beta^3\theta_1 + \beta^2\theta_2 + \beta\theta_3 + \theta_4\right] \\ \qquad\qquad \uparrow \qquad\quad \uparrow \qquad\quad \uparrow \qquad \uparrow \end{cases}$$

$$\beta^3 < \beta^2 < \beta \qquad\qquad 0 < \beta < 1$$

# Understanding Graphs

20 July 2022    12:03



2D



3D



Contour

# meshc() combines mesh/contour plots

gnuplot - how to obtain contour lines with the same level color of the 3d plot - Stack Overflow

Visit

**Related images**

See more

An experimental investigation on mechanical properties of Fe2O3 microparticles reinforced

function $f(x, y)$

- ⦿ $ax^2 + by^2 + cxy + dx$
- ○ $\cos(axy + bx + cy)$
- ○ $\log(axy + bx + cy + 1)$
- ○ $\dfrac{ax^2 + by^2 + cx^2 y}{x^2 + y^2}$

constants

a ——◯———— ⊡ −1
b ————◯— ⊡ 1
c ———◯—— ⊡ 0
d —◯——— ⊡ −1

3D and contour plots of $f(x, y) = -x^2 - x + y^2$

608 × 518

cocula.gob.mx

**Wolfram Grapher Discount, 60% OFF | cocula.gob.mx**

Visit

# Convex Vs Non-Convex Optimization

function $f(x, y)$

- $a x^2 + b y^2 + c x y + d x$
- $\cos(a x y + b x + c y)$
- $\log(a x y + b x + c y + 1)$
- $\dfrac{a x^2 + b y^2 + c x^2 y}{x^2 + y^2}$

constants

a ——— -1
b ——— 1
c ——— 0
d ——— -1

3D and contour plots of $f(x, y) = -x^2 - x + y^2$

608 × 518

cocula.gob.mx

**Wolfram Grapher Discount, 60% OFF | cocula.gob.mx**

Visit

Small Radius
High Curvature

Large Radius
Low Curvature

Radius of
Curvature
Comparison

SGD                    SGD momentum

Momentum Optimizer(decay = 0.9)
epoch number: = 49

Momentum Optimizer(decay = 0.8)
epoch number: = 36

Batch Gradient Descent
epoch number: = 40

Momentum Optimizer(decay = 0.9)
epoch number: = 49

Momentum(NAG) Optimizer(decay = 0.9)
epoch number: = 24

# Keras Code

```
tf.keras.optimizers.SGD(
    learning_rate=0.01, momentum=0.0, nesterov=False, name="SGD", **kwargs
)
```

SGD

# Keras Code

```
tf.keras.optimizers.SGD(
    learning_rate=0.01, momentum=0.0, nesterov=False, name="SGD", **kwargs
)
```

SGD {

Momentum

momentum = 0.9

nesterov = False
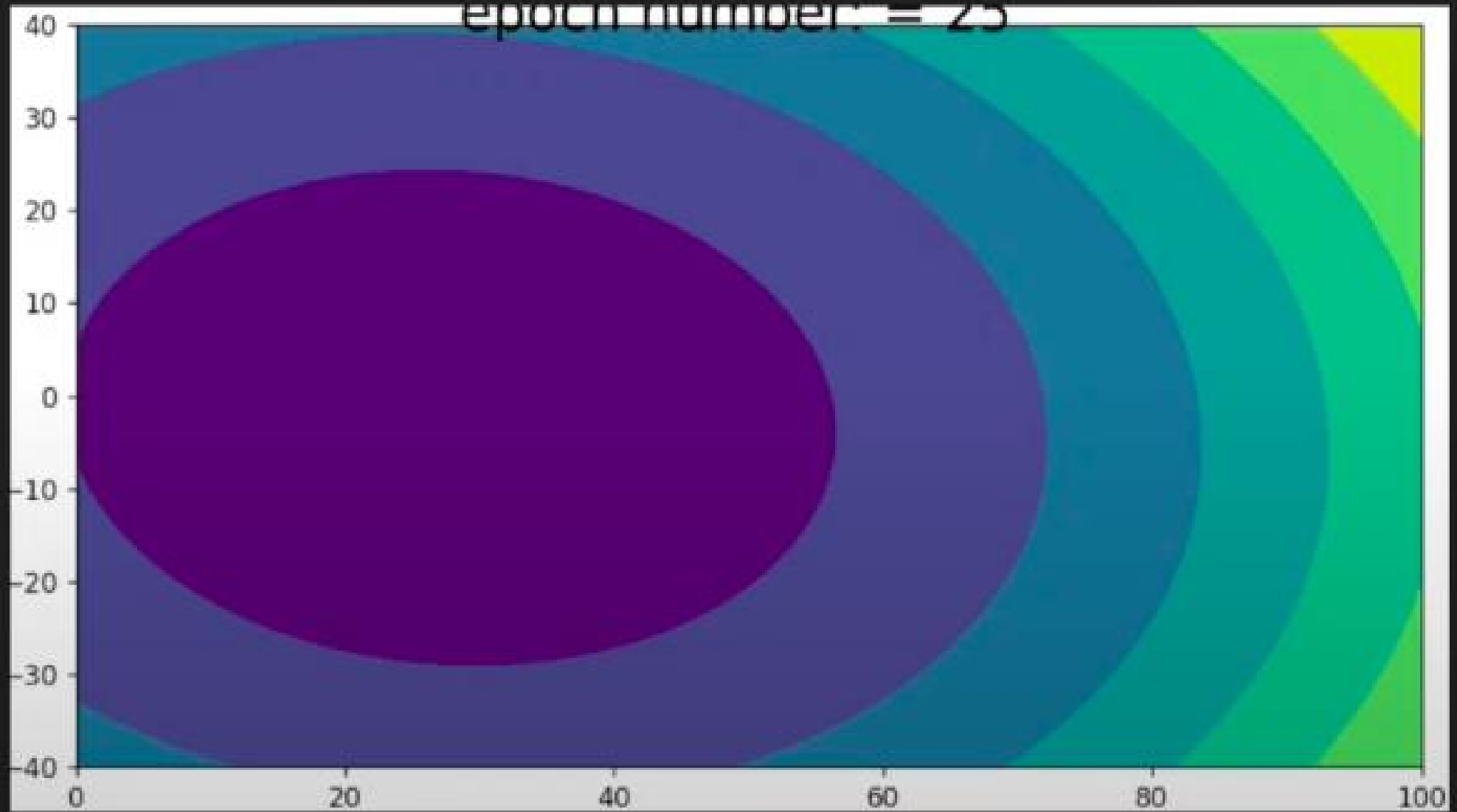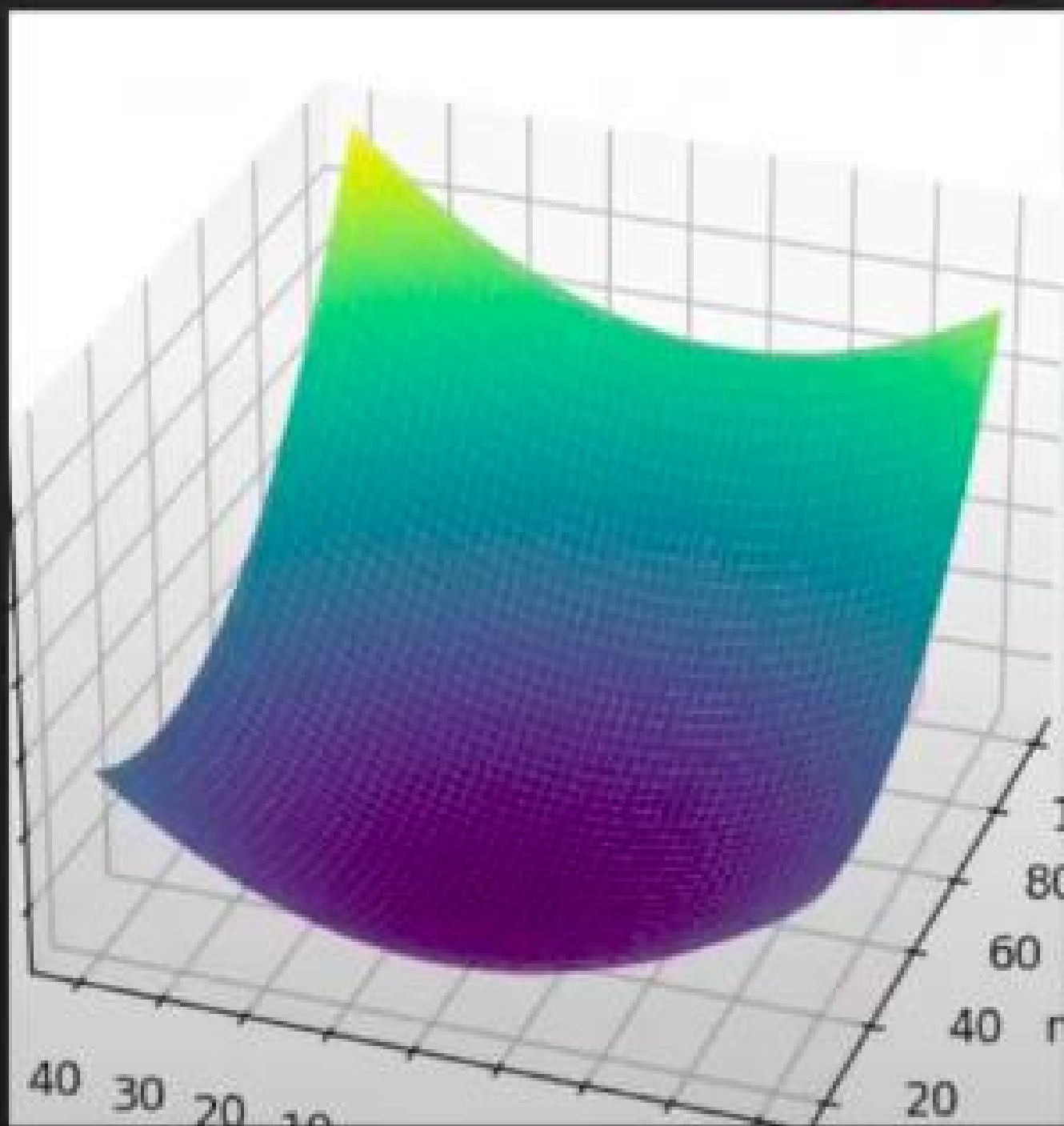
# Keras Code

```
tf.keras.optimizers.SGD(
    learning_rate=0.01, momentum=0.0, nesterov=False, name="SGD", **kwargs
)
```

SGD

Momentum
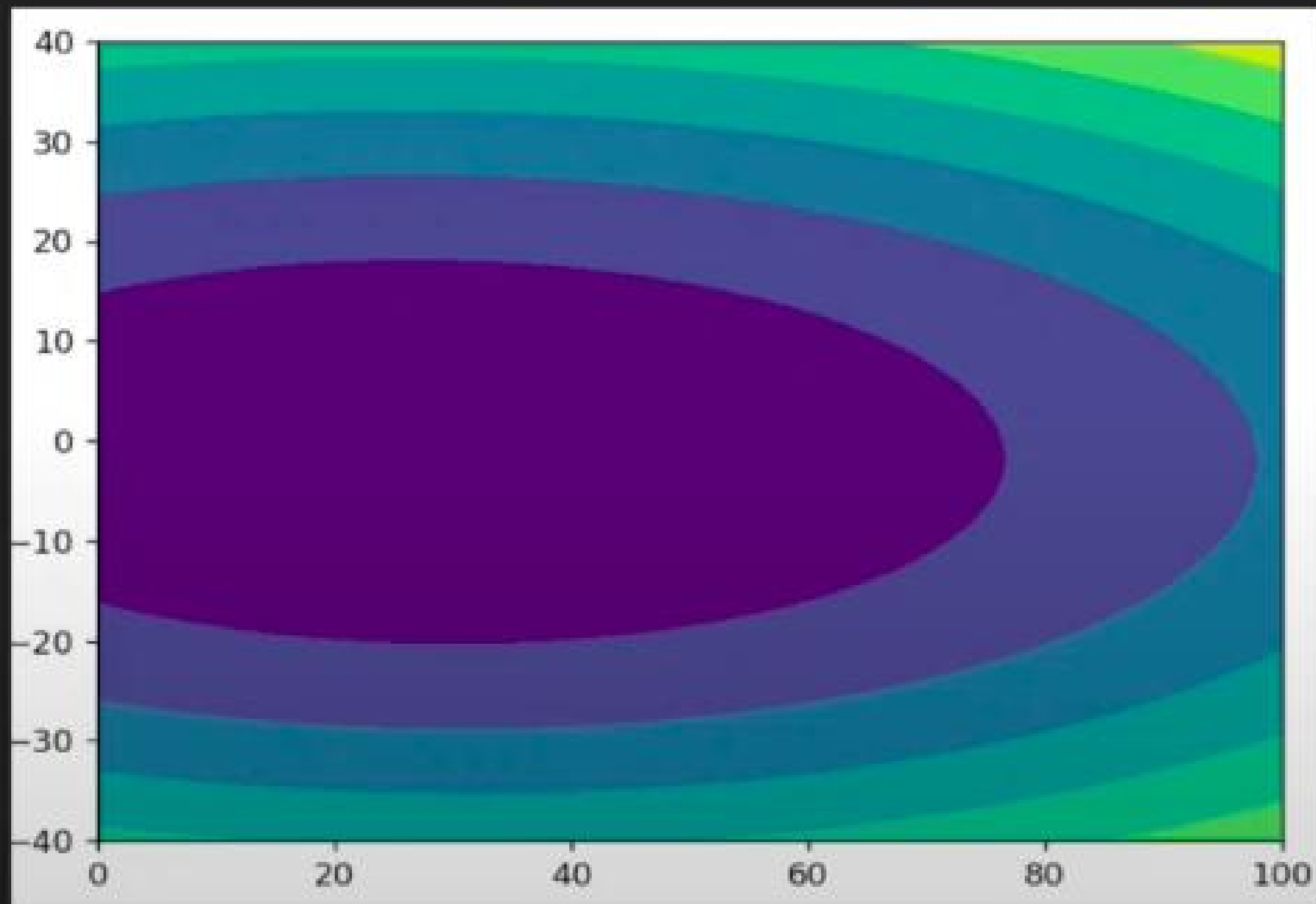
momentum = 0.9
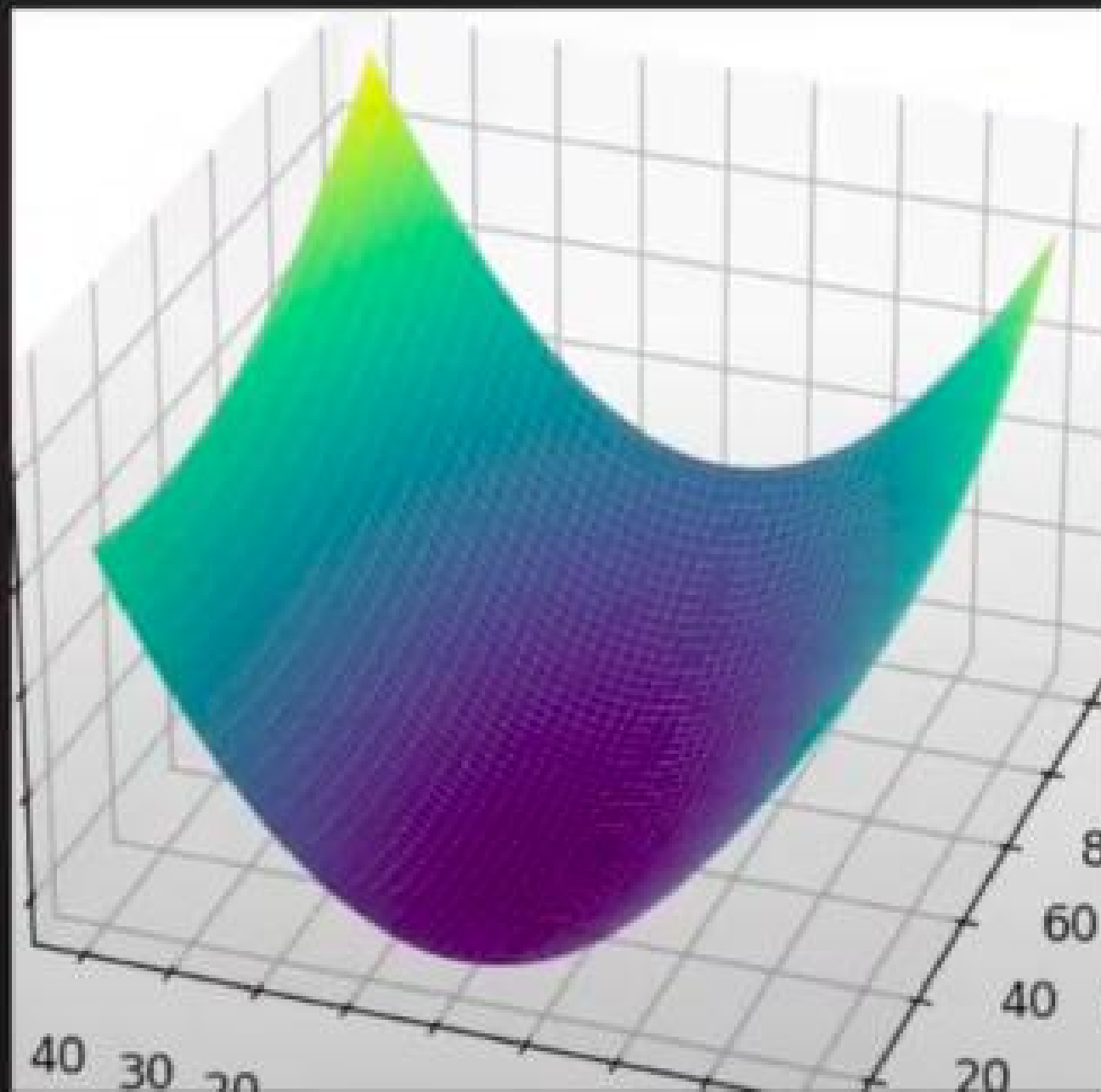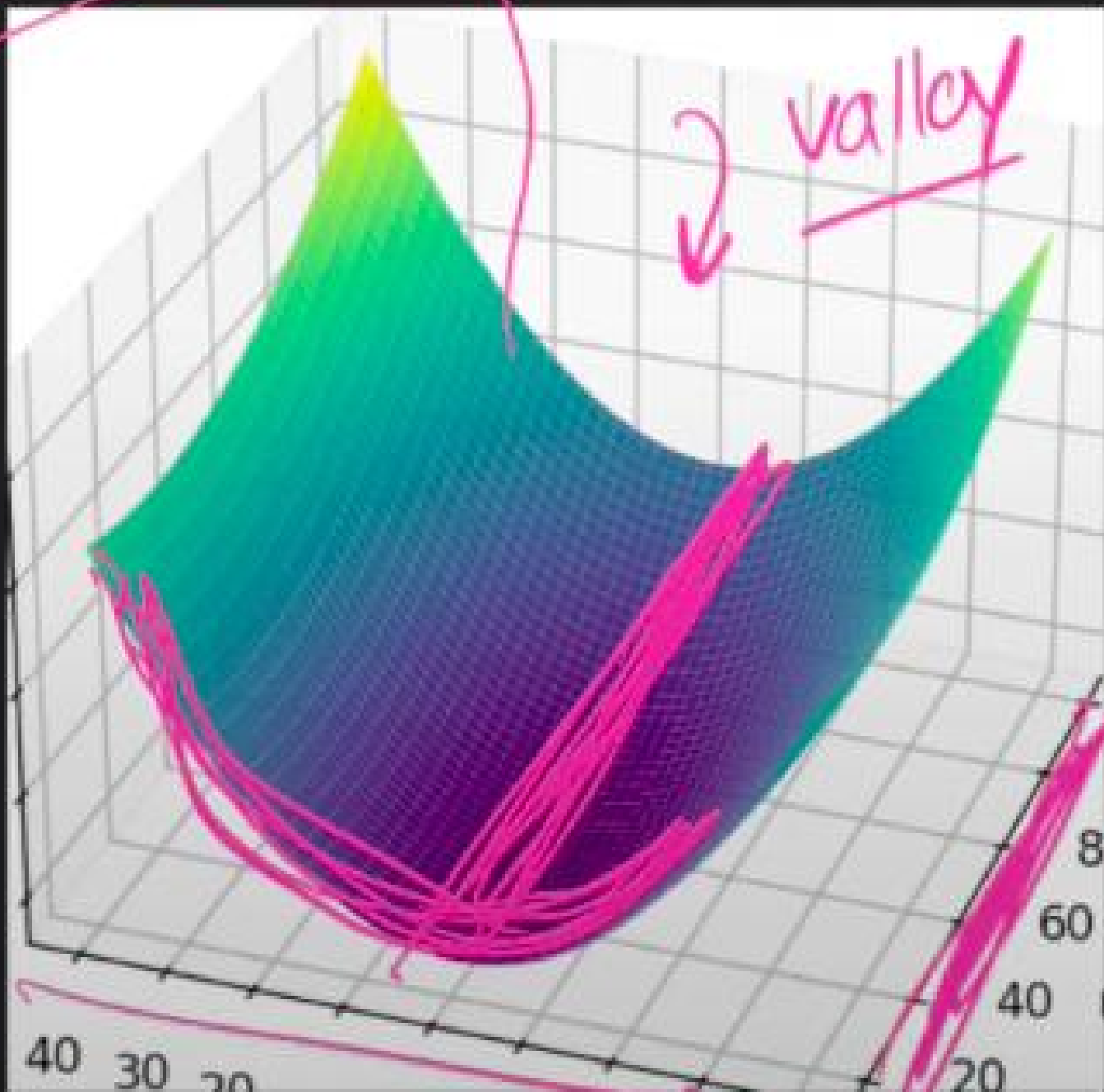
nesterov = False

NAG

momentum = 0.9

nesterov = True

epoch number = 23

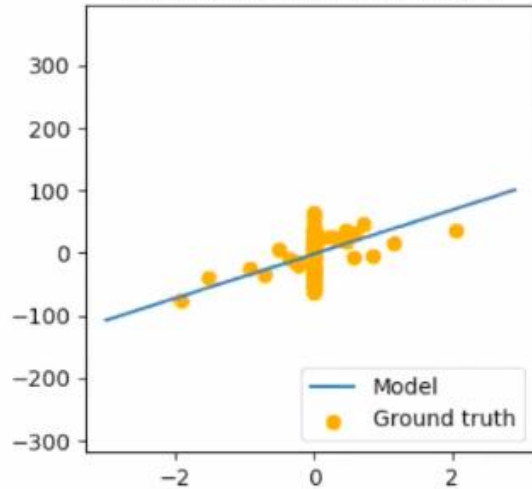Batch Gradient Descent
epoch number: = 75

# Momentum Optimizer(decay = 0.5)
## epoch number: = 42

# How optimizers behave(Why?)

Batch GD

Momentum

$(y-\hat{y})^2$

$x$  $w$ $\longrightarrow$ Linear

$\longrightarrow \hat{y}$

$b$

$\boxed{\dfrac{\partial L}{\partial w}} = \dfrac{\partial L}{\partial \hat{y}} \dfrac{\partial \hat{y}}{\partial w}$

$= -2(y-\hat{y})\boxed{X}$ add

$\dfrac{\partial L}{\partial b} = \boxed{-2(y-\hat{y}) * ①}$ ← bignum

$BGD \longrightarrow$ sparspe

for i in epochs:

$Ⓦ = w - \eta \boxed{\dfrac{\partial L}{\partial w}}$

$b = b - \eta \boxed{\dfrac{\partial L}{\partial b}}$ → big update in.

$\boxed{100 \text{ rows}}$  0  $\underline{X\_sporse}$

↓ →

small number

small in every epoch
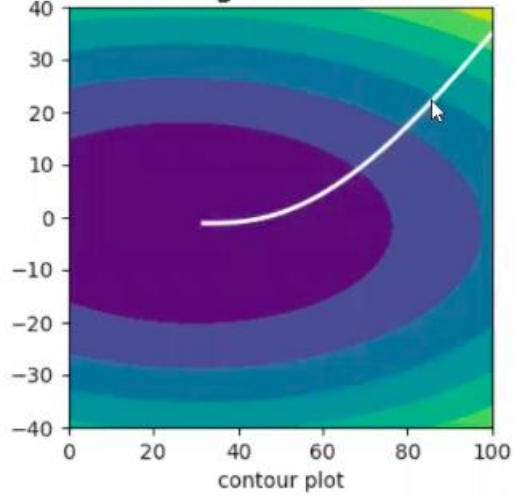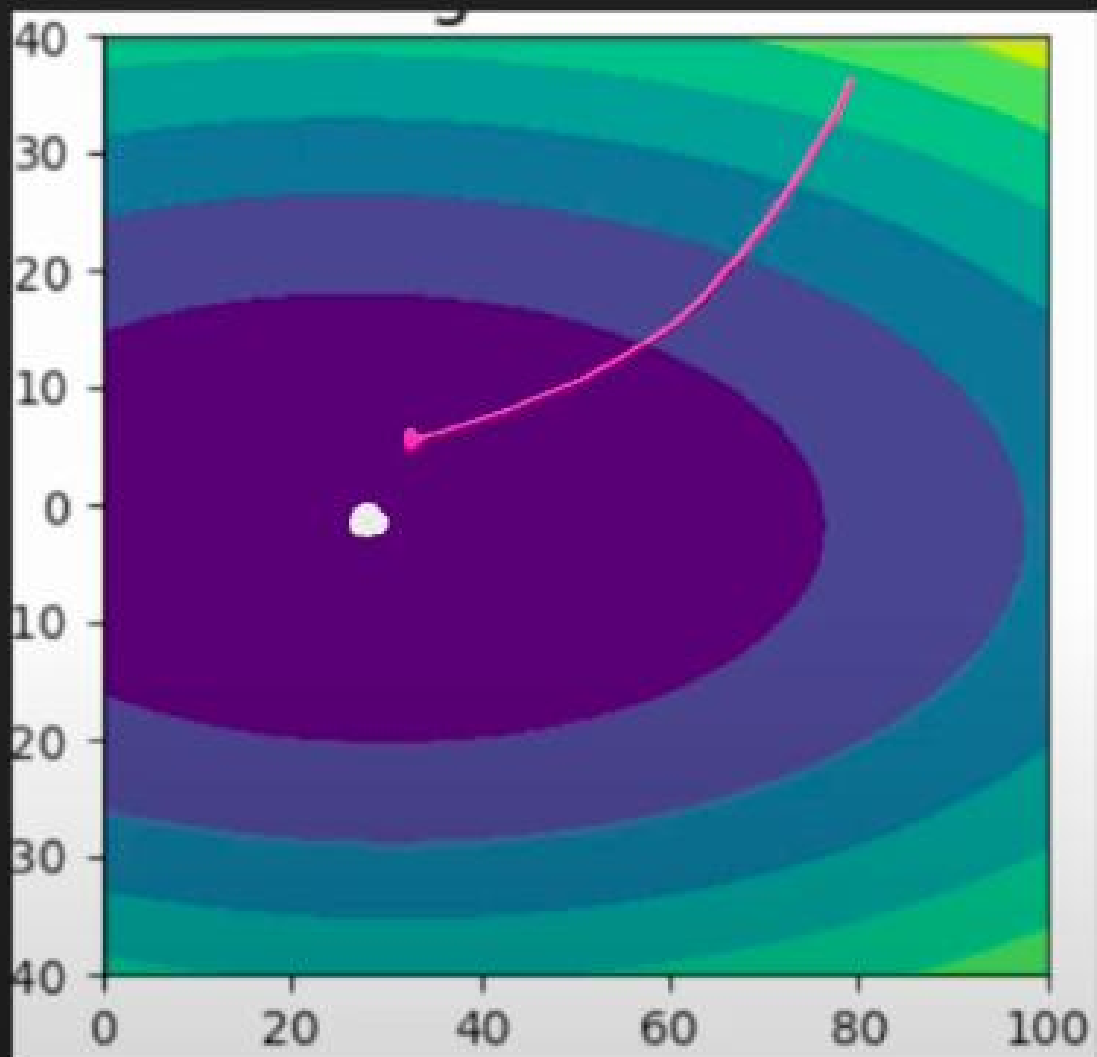
Adagrad Optimizer
epoch number: = 61

Ground truth & Model

Weights & Loss
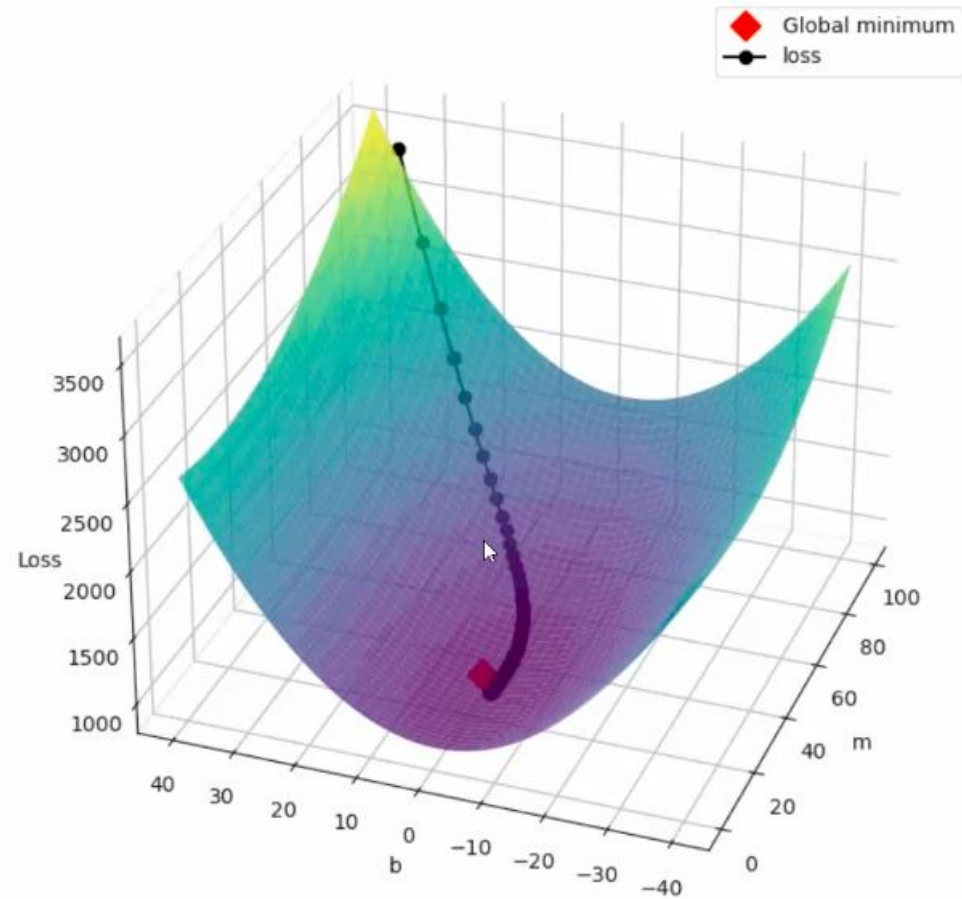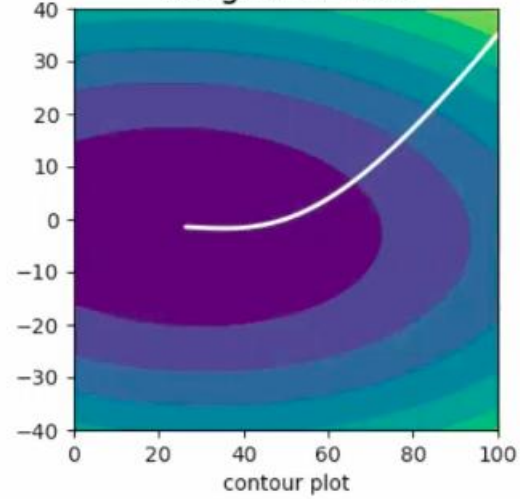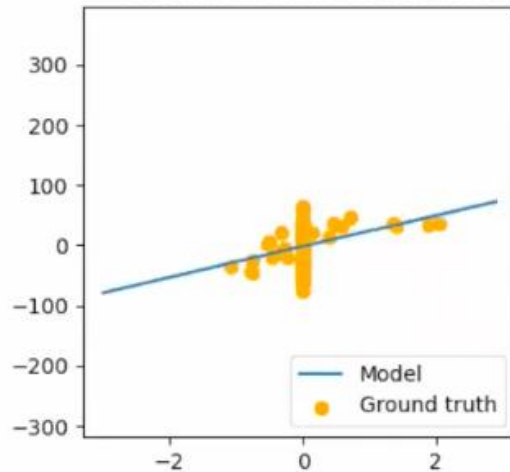contour plot

Global minimum
loss

Loss

# RMSProp
## epoch number: = 98

### Ground truth & Model

### Weights & Loss
contour plot