

Optimization

Noa Lubin & Lior Sidi



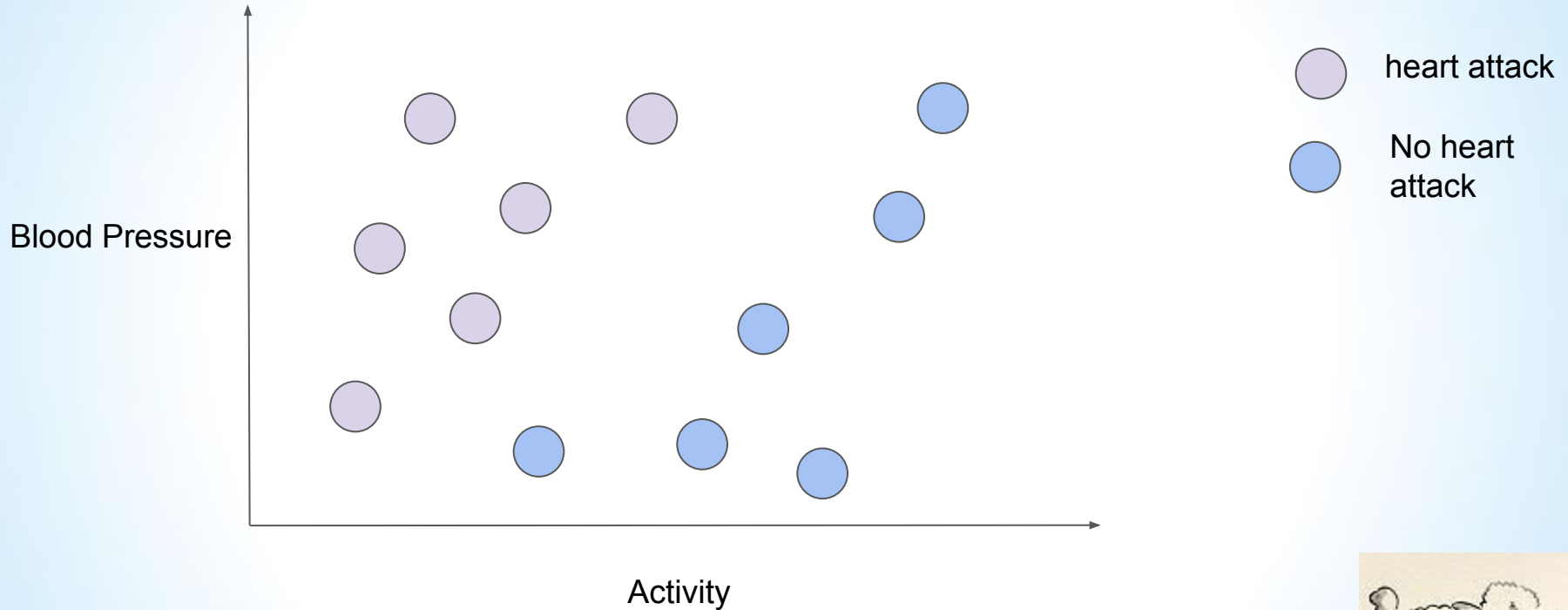
Agenda

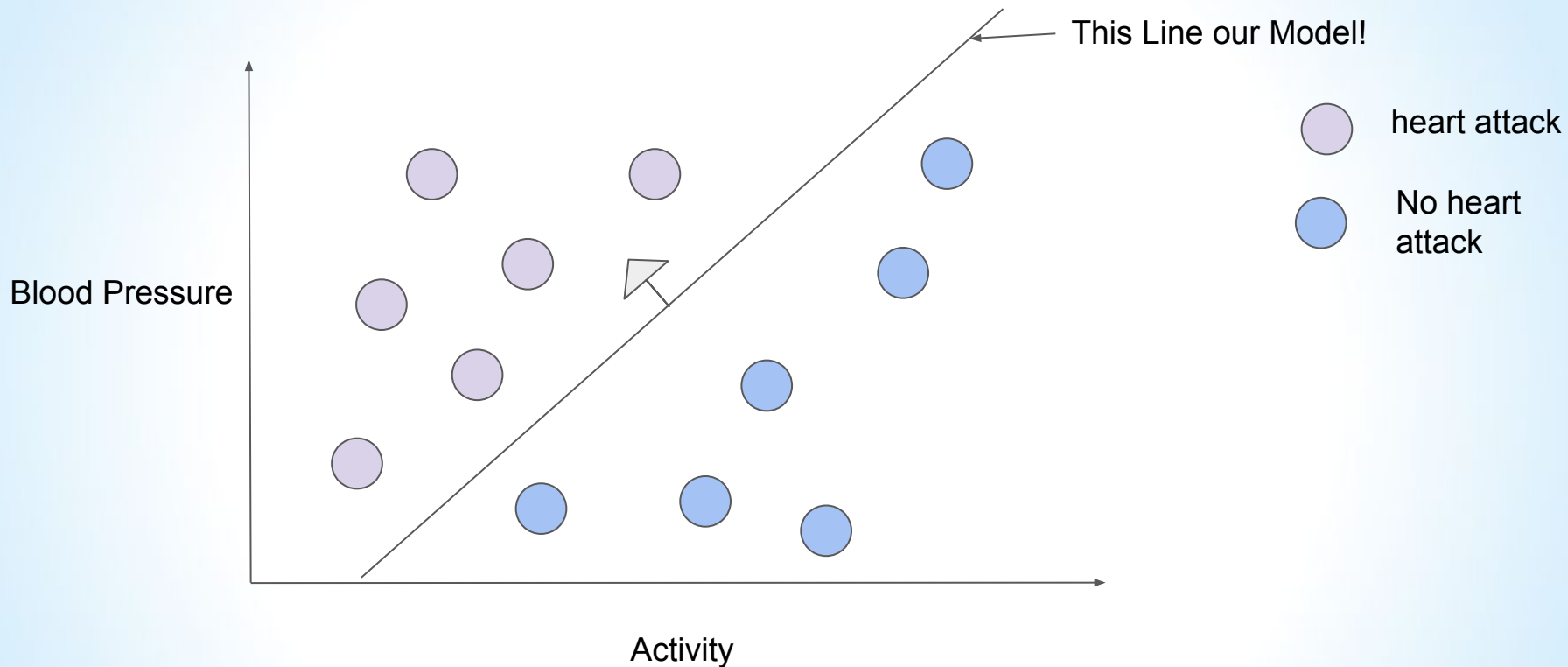
- Motivation
- Optimization Theory
- Gradient Descent Toy Example

Motivation

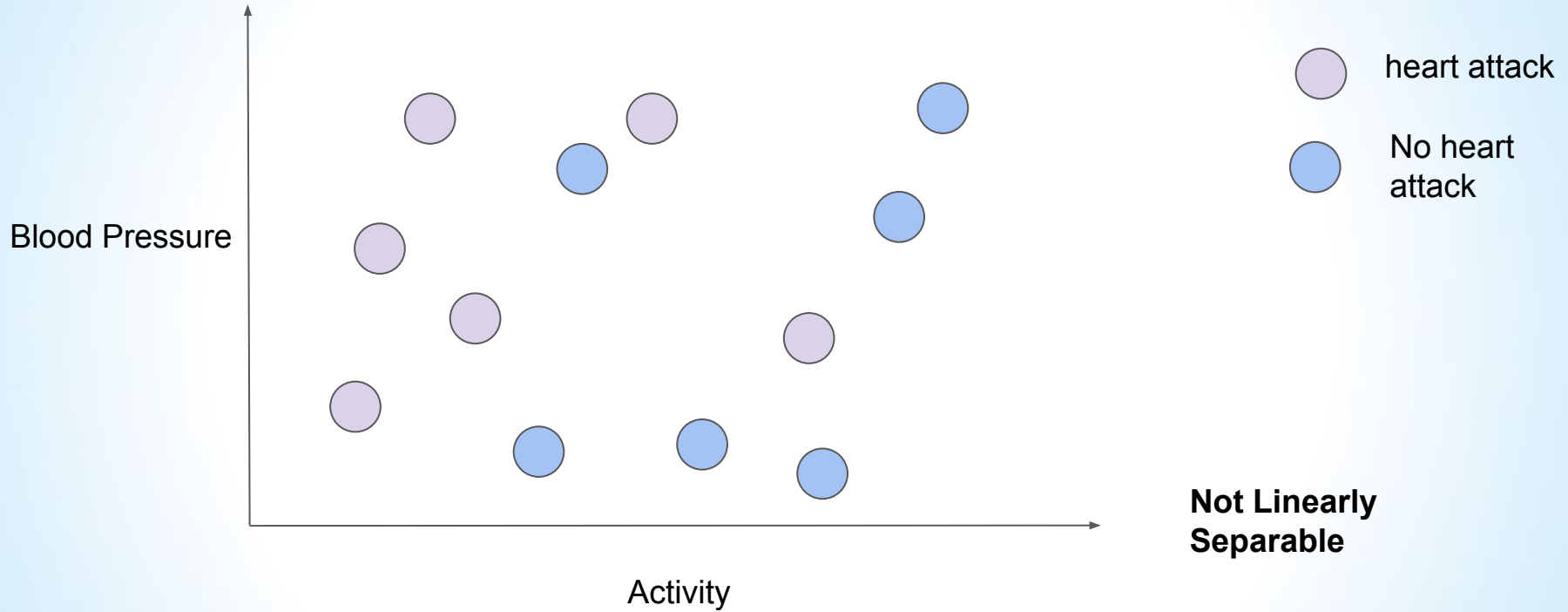


Draw Me a Model





In Real Life



Mathematical Formulation

Features:

x_1 - blood pressure

x_2 - activity

$f(x_1, x_2) =$ Heart Attack if $w_0 + w_1 * x_1 + w_2 * x_2 \geq 0$

No Heart Attack otherwise

Mathematical Formulation - Matrix Form

Features:

x_1 - blood pressure

x_2 - activity



$$x = \langle 1, x_1, x_2 \rangle$$

$$w = \langle w_0, w_1, w_2 \rangle$$

$$wx = 1 \cdot w_0 + x_1 \cdot w_1 + x_2 \cdot w_2$$

$f(x_1, x_2) =$ Heart Attack if $wx \geq 0$

No Heart Attack otherwise

This is a linear
classifier!

How can we find w ?

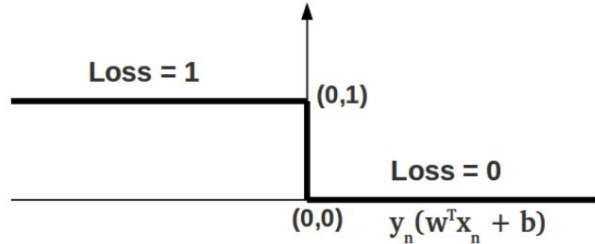
Loss Function

A loss function or cost function indicates “how wrong” we are with our current hypothesis

For example:

0 - 1 loss: 1 If the classification is incorrect

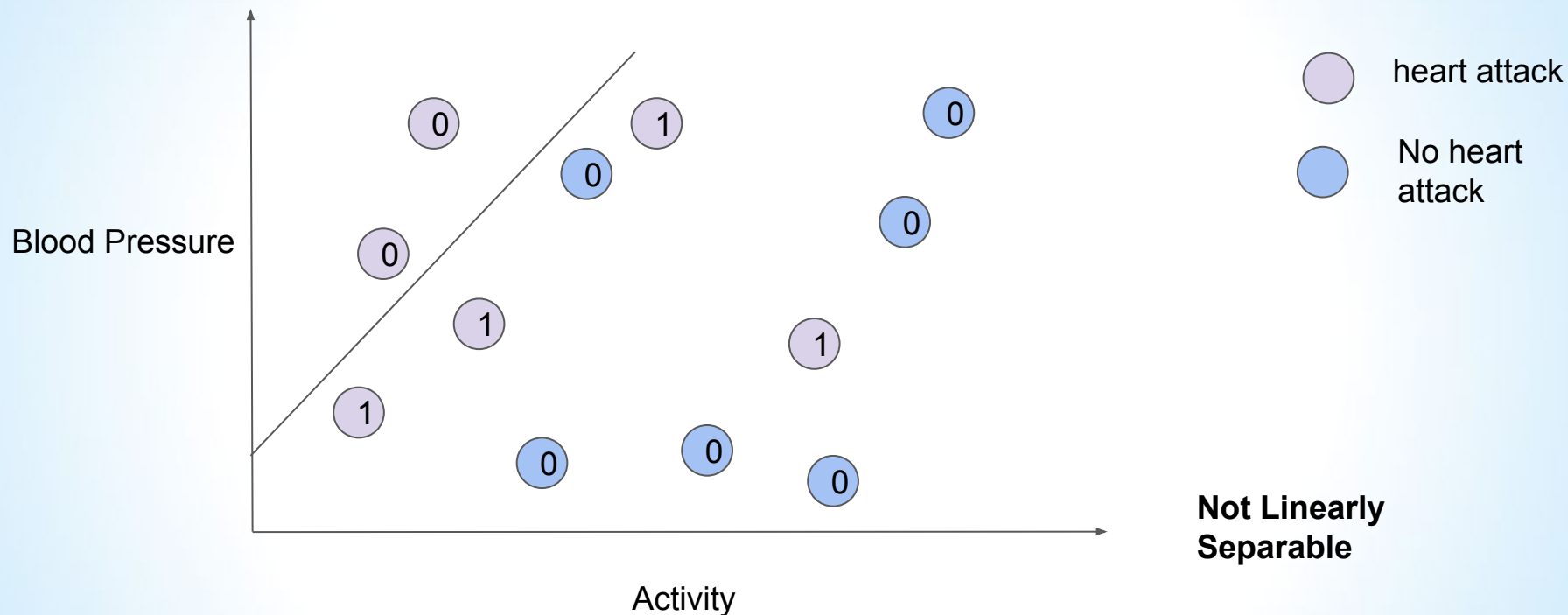
0 If the classification is correct



$$l(\hat{y}, y) = \begin{cases} 0 & y = \hat{y} \\ 1 & y \neq \hat{y} \end{cases} = I\{y \neq \hat{y}\}$$

1-0 loss

Let's guess a model

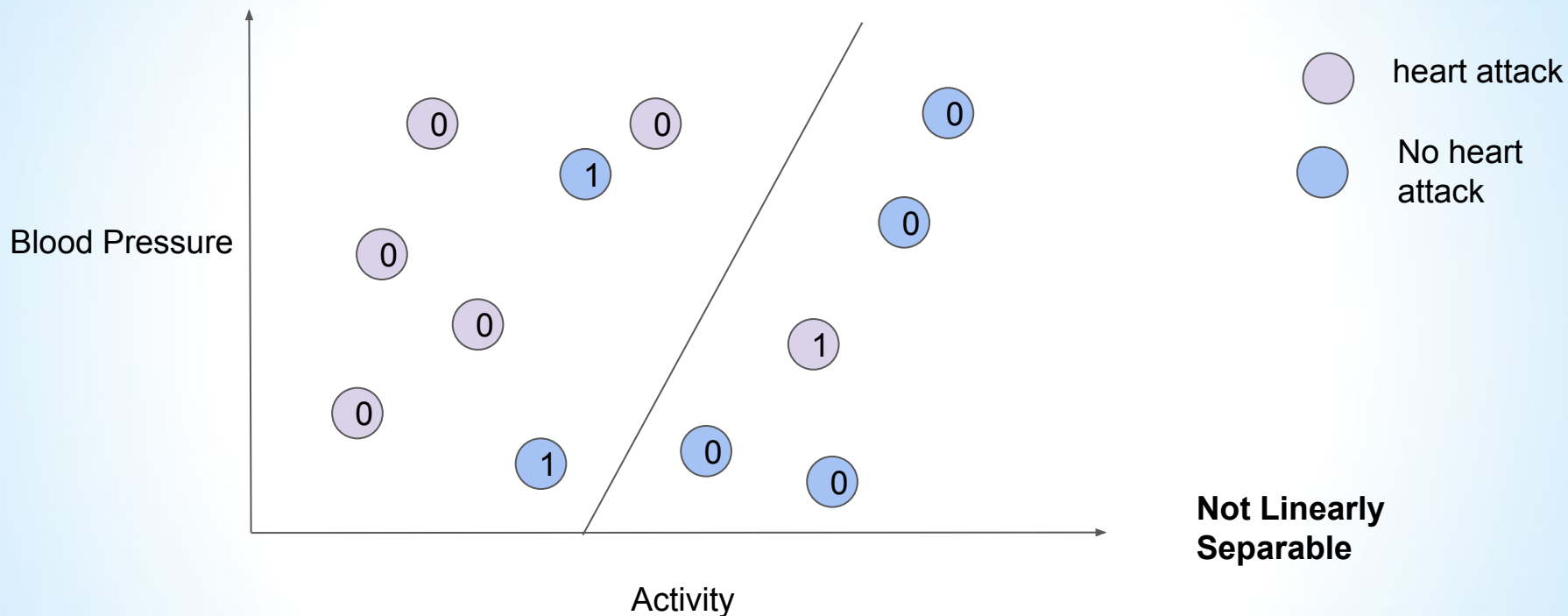


Loss = 4

Not Linearly
Separable

1-0 loss

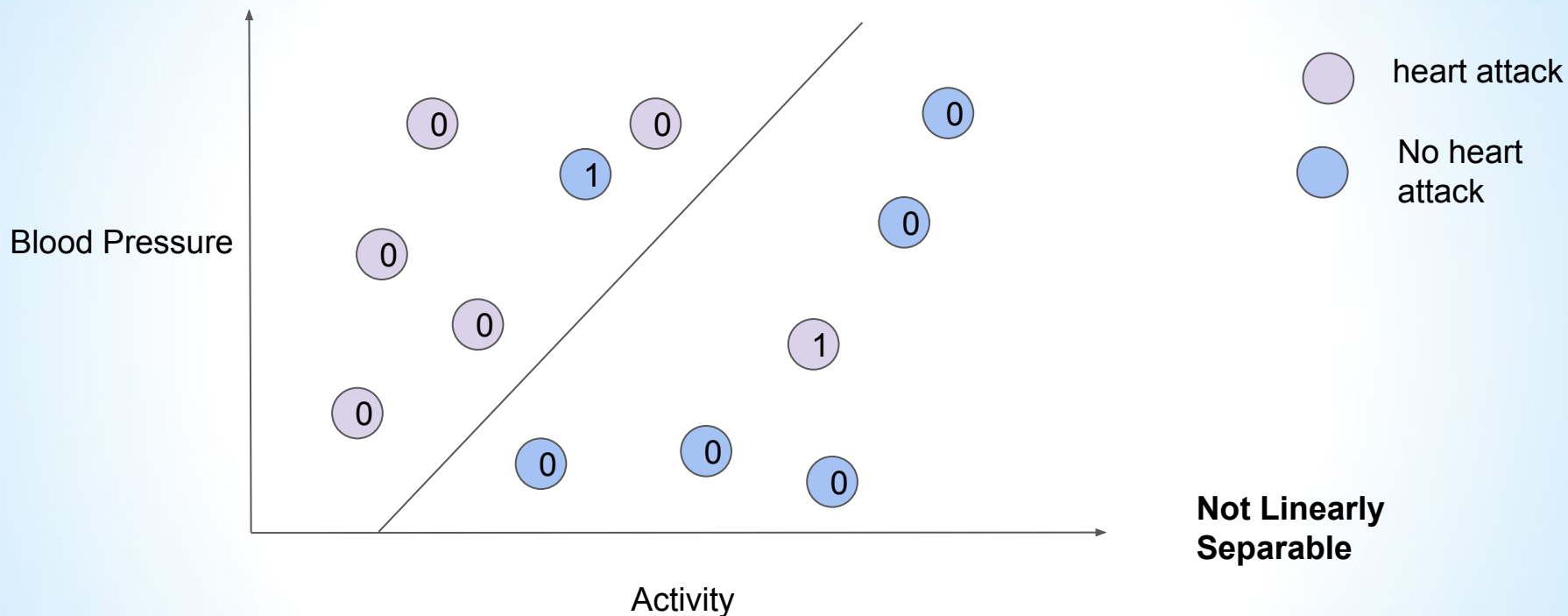
Let's guess a model



Loss = 3

1-0 loss

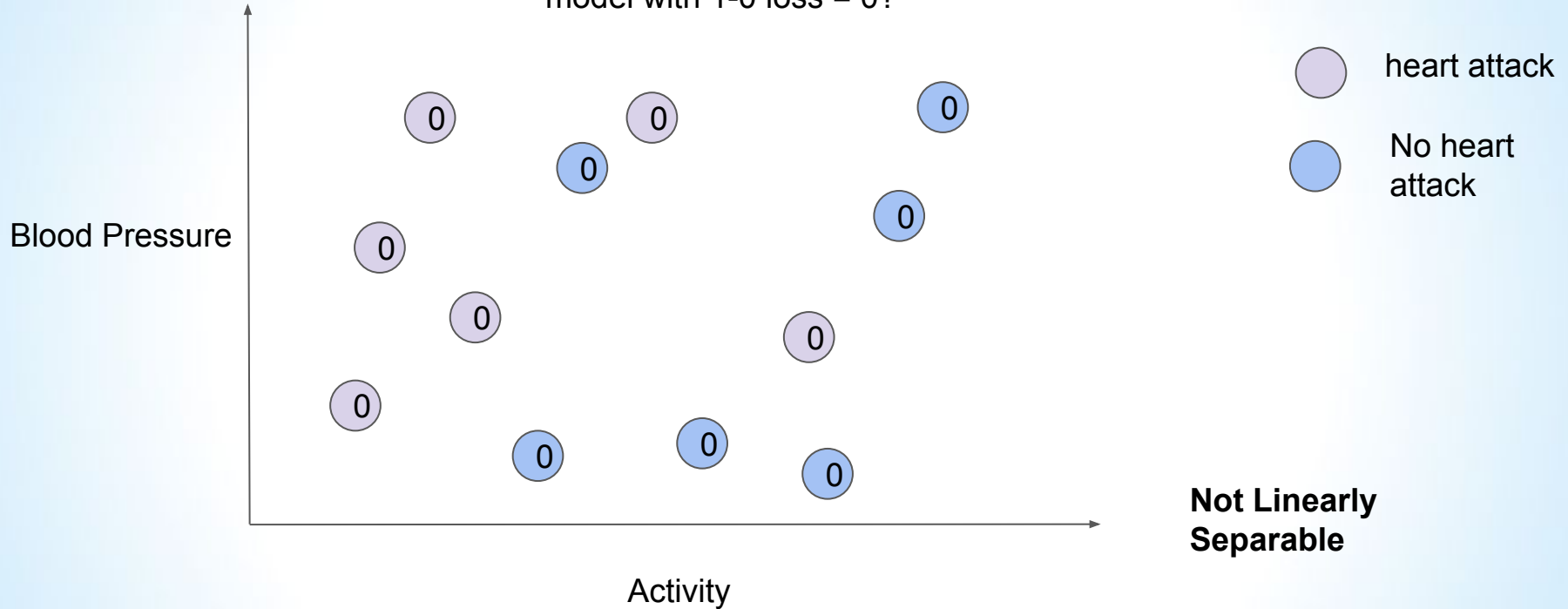
Let's guess a model



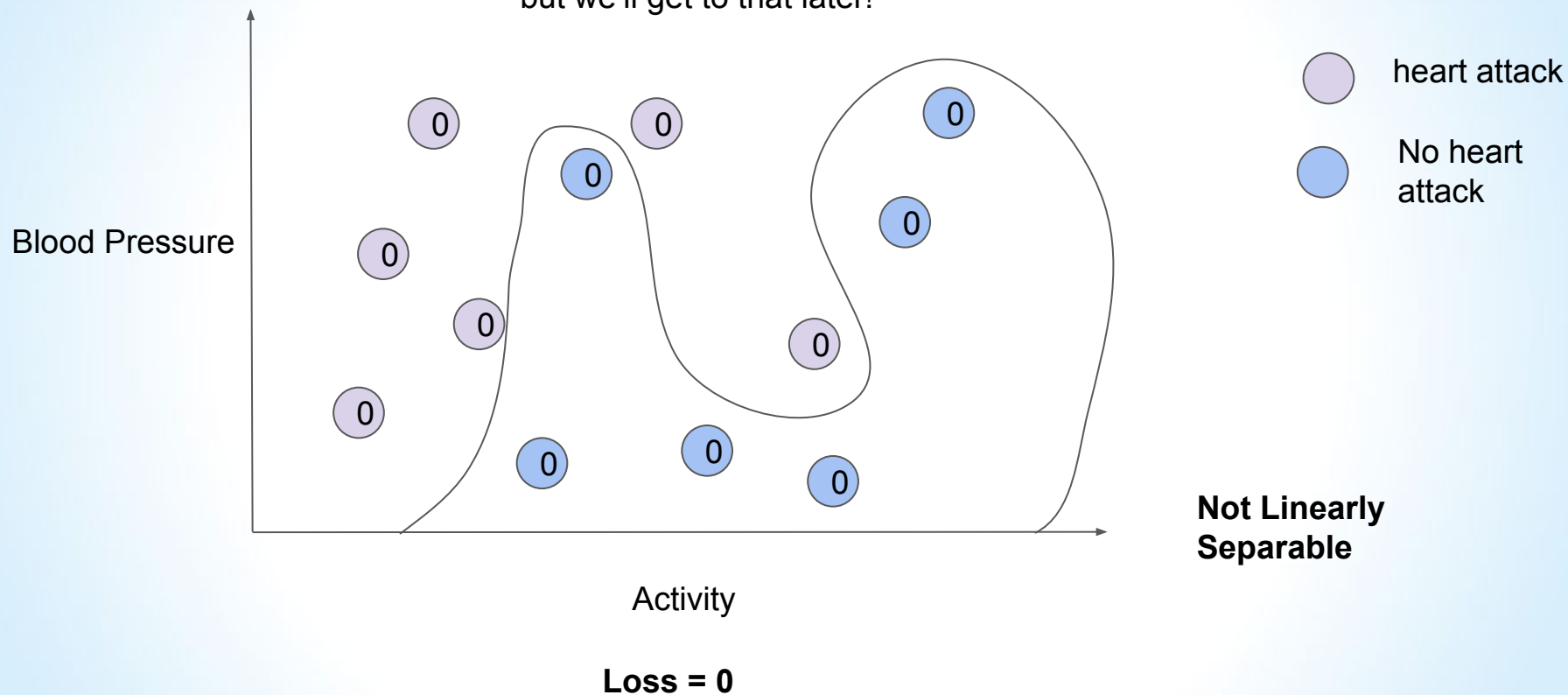
Loss = 2

1-0 loss

Can you guess a non - linear
model with 1-0 loss = 0?



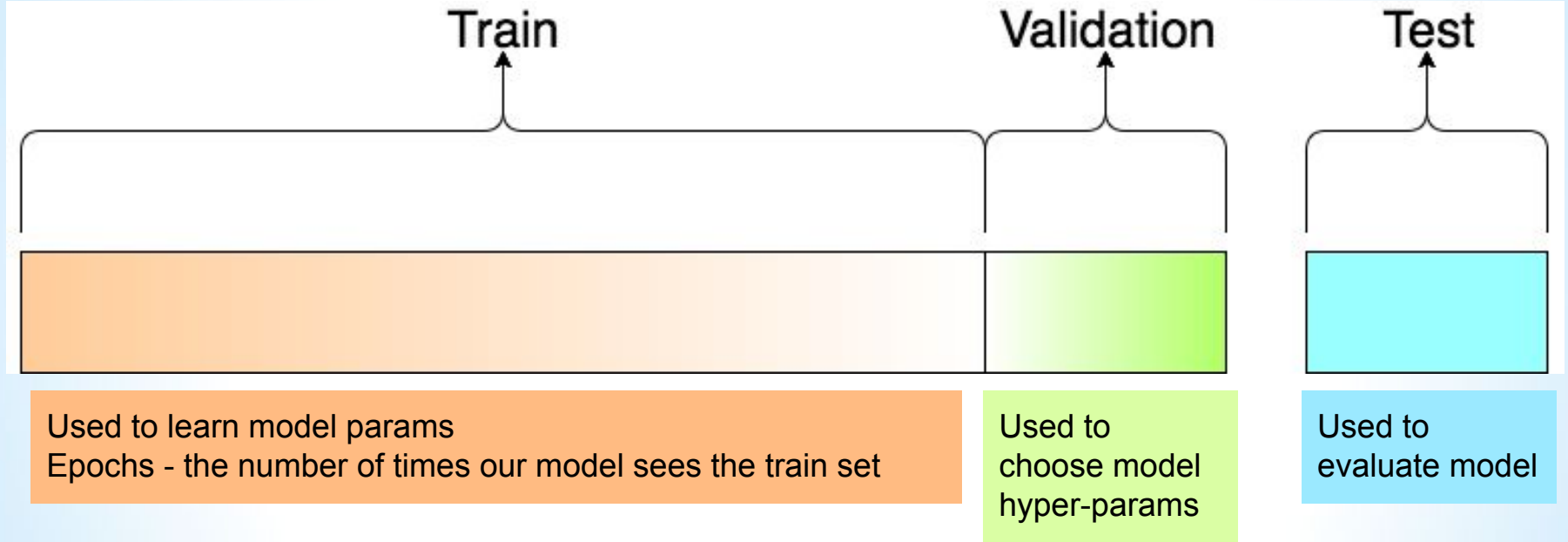
This what we call overfitting
but we'll get to that later!



Optimization Theory



Reminder:



ERM principal

- Minimize your loss on the observed data (training set)

$$\theta^* = \operatorname{argmin}_{\theta} \frac{1}{m} \sum_{j=1}^m l(y_j, f_{\theta}(x_j))$$

Optimization

The optimization method is independent of the loss function.

Loss - Where do we want to get?

Optimization - How do we want to get there?

We'll now discuss optimization methods.

Types Of Optimization

- **Gradient Descent**
- **Stochastic Gradient Descent**
- **Mini-batch Gradient Descent**
- Momentum
- AdaGrad
- RMSprop
- Adam
- ...

Gradient Descent


The idea of gradient descent is:

Take iterative steps to update parameters in the direction of the gradient

$$w^t \leftarrow w^{t-1} - \eta \cdot \frac{\partial \ell}{\partial w^{t-1}}$$

Diagram illustrating the gradient descent update formula:

- Previous weights** points to w^{t-1} .
- gradient** points to $\frac{\partial \ell}{\partial w^{t-1}}$.
- New weights** points to w^t .
- Step size / learning rate** points to η .



Gradient

A collection of partial derivatives

$$\nabla f(x_1, x_2, x_3) = \left(\frac{\partial f(x_1, x_2, x_3)}{\partial x_1} \frac{\partial f(x_1, x_2, x_3)}{\partial x_2} \frac{\partial f(x_1, x_2, x_3)}{\partial x_3} \right)$$

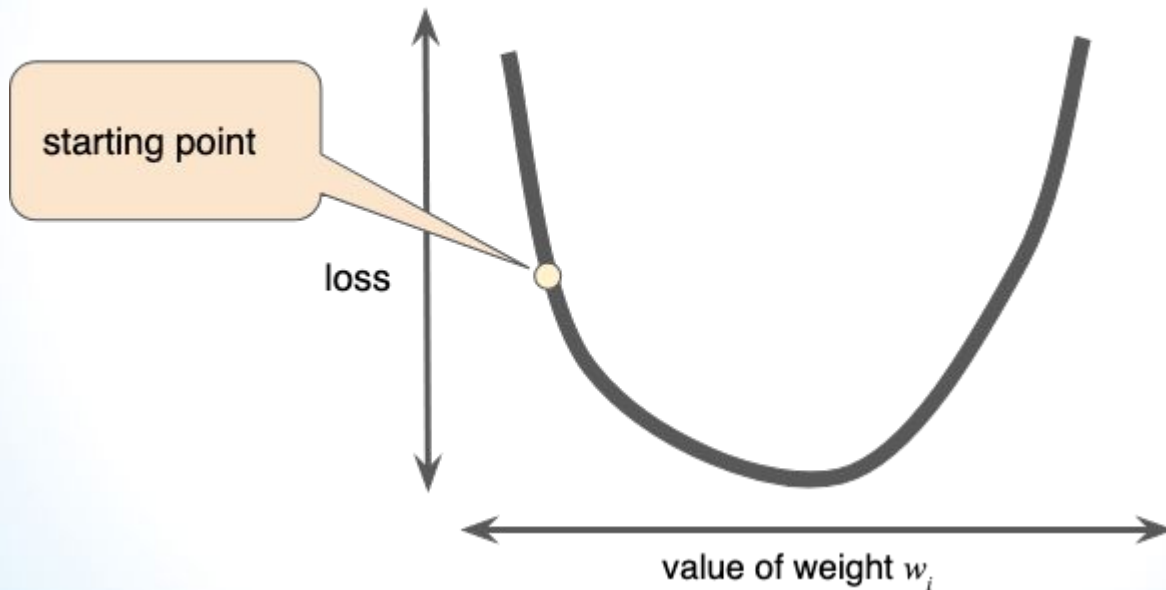
Gradient for multiple inputs

$$\frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{W}} = \begin{pmatrix} \frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{W}_{[1,1]}} & \frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{W}_{[1,2]}} & \dots & \frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{W}_{[1,n]}} \\ \frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{W}_{[2,1]}} & \frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{W}_{[2,2]}} & \dots & \frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{W}_{[2,n]}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{W}_{[m,1]}} & \frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{W}_{[m,2]}} & \dots & \frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{W}_{[m,n]}} \end{pmatrix}$$

$$\frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{b}} = \left(\frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{b}_{[1]}} \quad \frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{b}_{[2]}} \quad \dots \quad \frac{\partial L(\hat{\mathbf{y}}, \mathbf{y})}{\partial \mathbf{b}_{[n]}} \right)$$

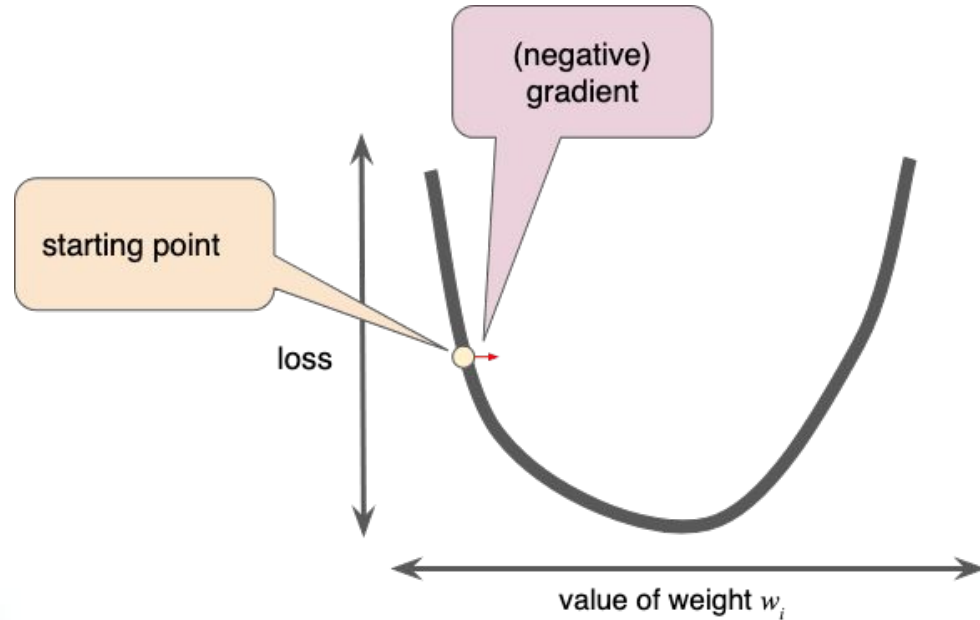
Gradient Descent

- Random starting point



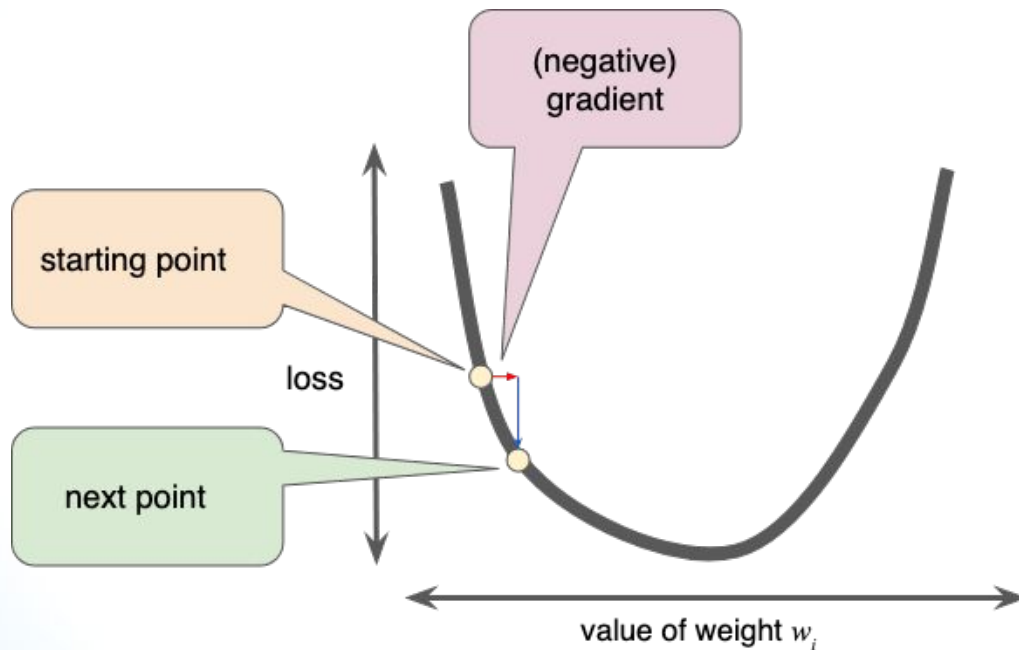
Gradient Descent

- A gradient is a vector so it has both magnitude and direction

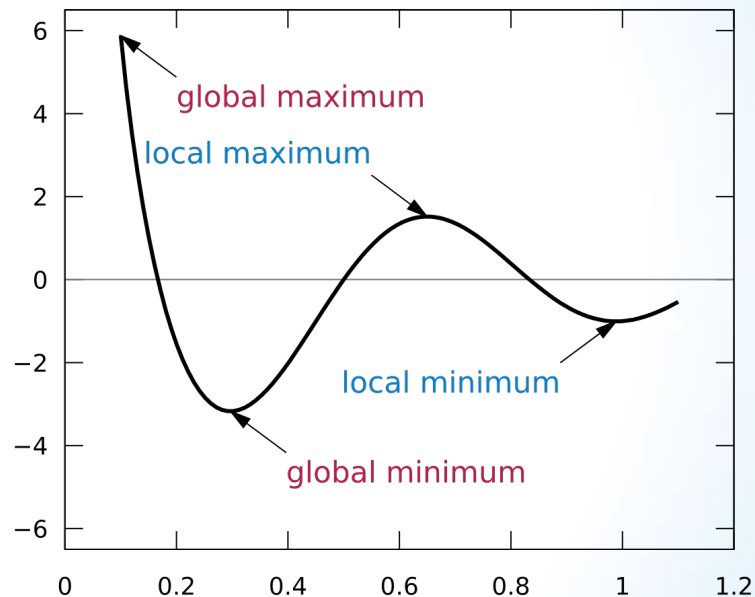


Gradient Descent

Until we meet stop criteria



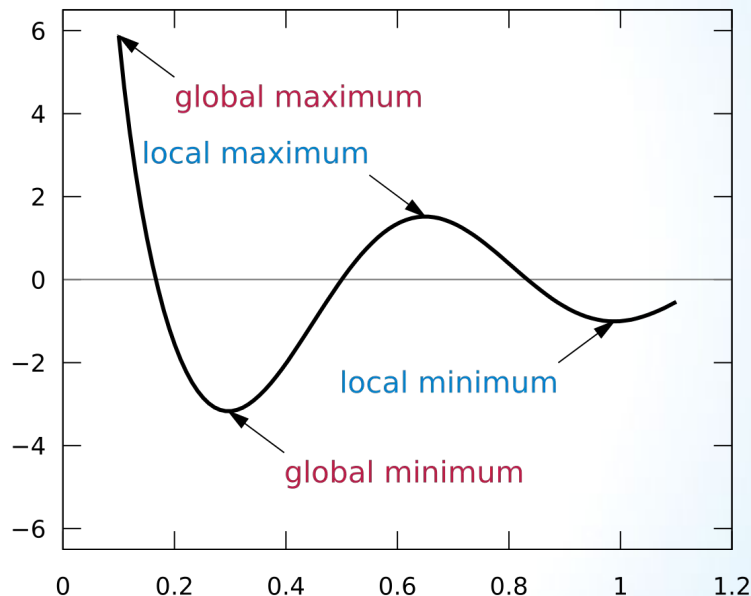
Does Gradient Descent Promise Global Minima?



Does Gradient Descent Promise Global Minima?

No!

Only when $f(x)$ is convex.



Gradient Descent

1. When do we stop?

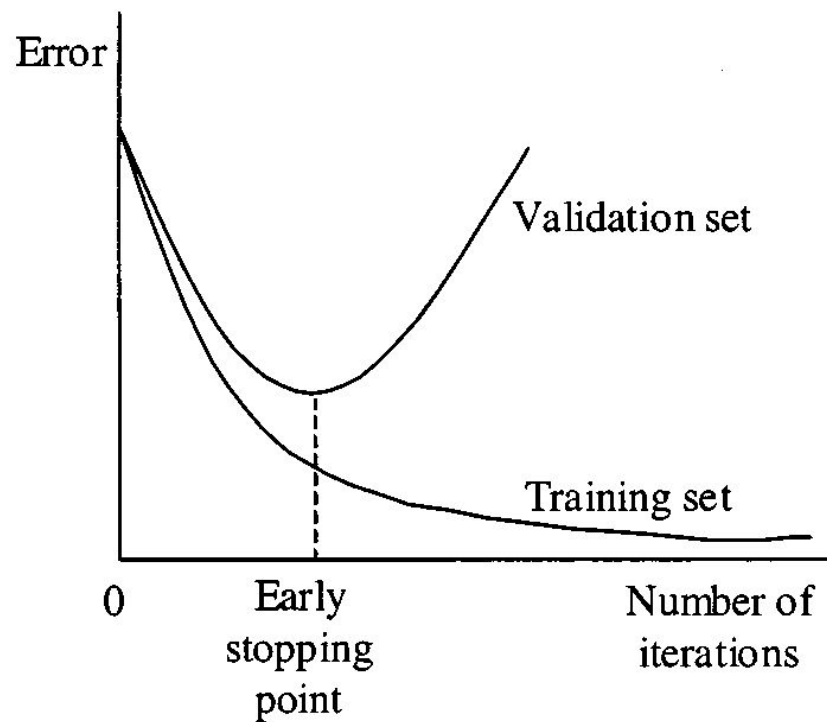


a. Fixed

b. Early Stopping

- i. When the gradient gets close to zero
- ii. When the loss stops changing much
- iii. When the parameters stop changing much
- iv. When performance on held-out validation set plateaus

Early Stopping



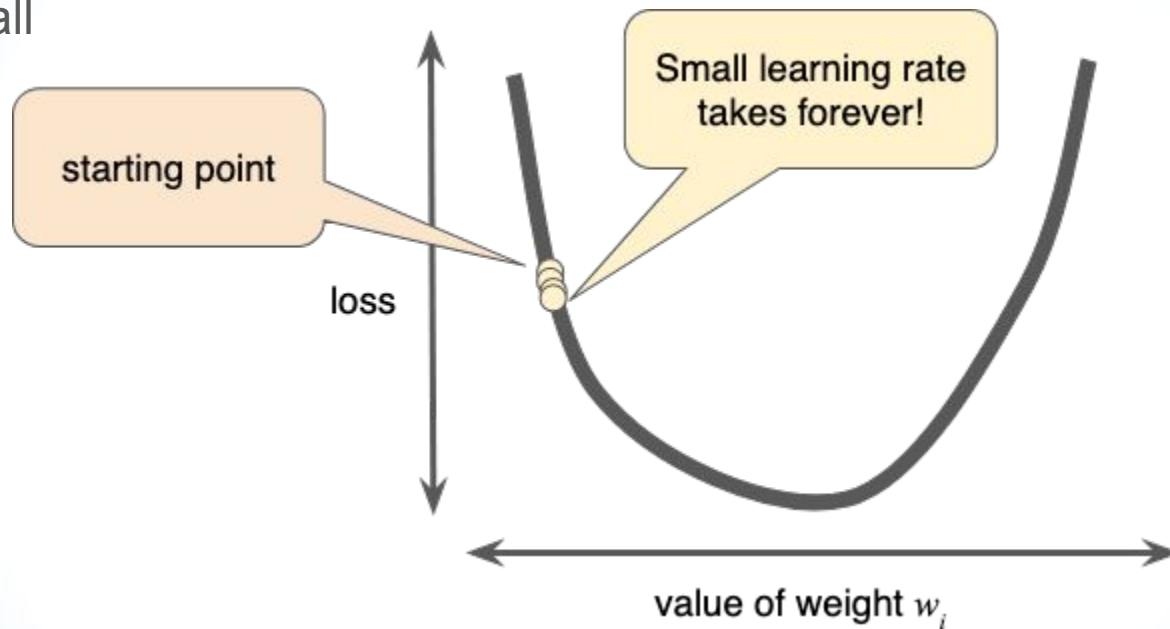
Gradient Descent

2. How do we choose step size aka learning rate?



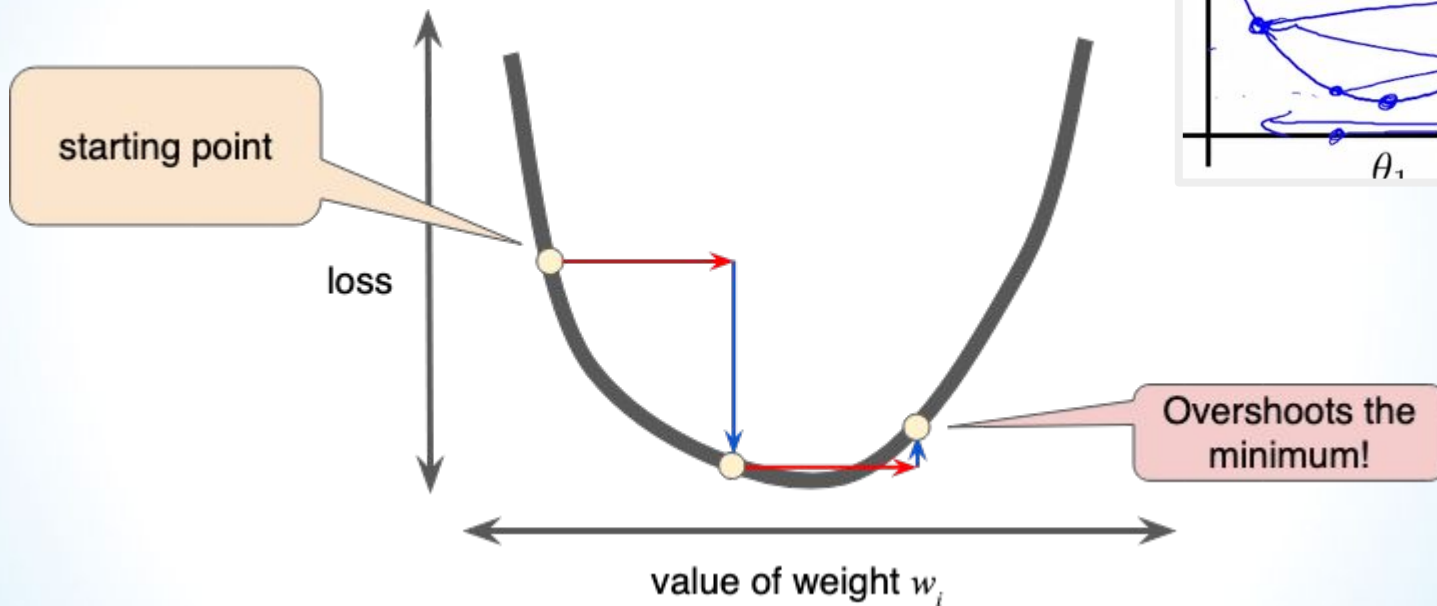
Learning Rate

LR is too small



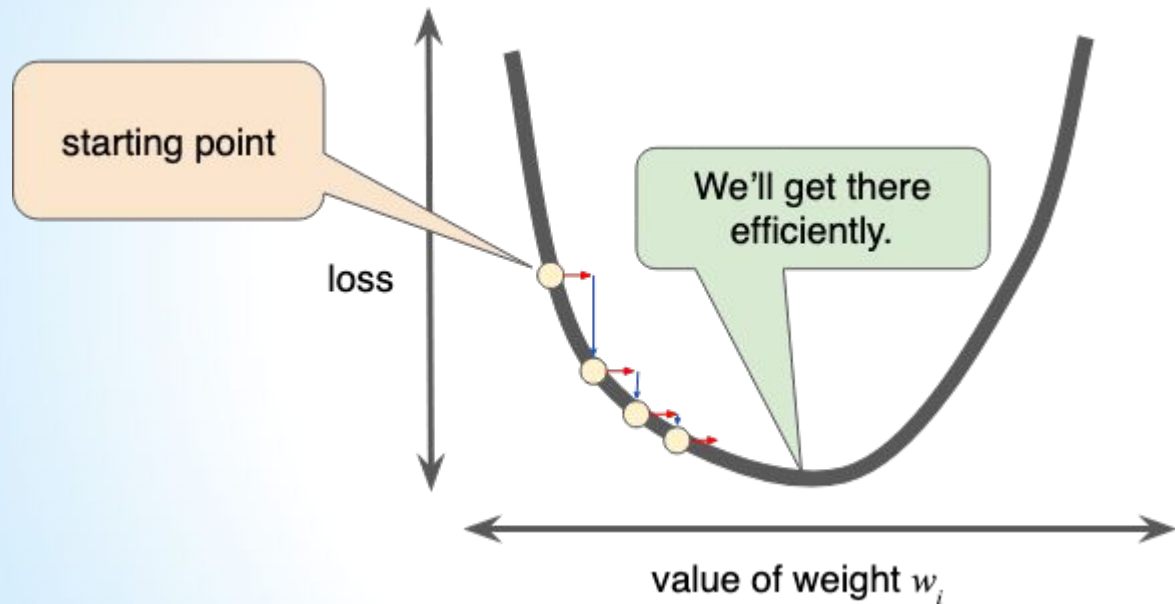
Learning Rate

LR is too large



Learning Rate

Learning rate is optimal



How can we choose the best LR?

1. Start with Large LR and reduce in each step
2. Use validation set to empirically choose the LR

Gradient Descent Training

Algorithm 1 Gradient Descent Training

Input:

- Function $f(\mathbf{x}; \Theta)$ parameterized with parameters Θ .
- Training set of inputs $\mathbf{x}_1, \dots, \mathbf{x}_n$ and desired outputs $\mathbf{y}_1, \dots, \mathbf{y}_n$.
- Loss function L .

-
- 1: **while** stopping criteria not met **do**
 - 2: Compute the loss $\mathcal{L}(\Theta) = \sum_i L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$ **<-- slow! goes over all data.**
 - 3: $\hat{\mathbf{g}} \leftarrow$ gradients of $\mathcal{L}(\Theta)$ w.r.t Θ
 - 4: $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$
 - 5: **return** Θ
-

Stochastic Gradient Descent Training



Algorithm 2 Online Stochastic Gradient Descent Training

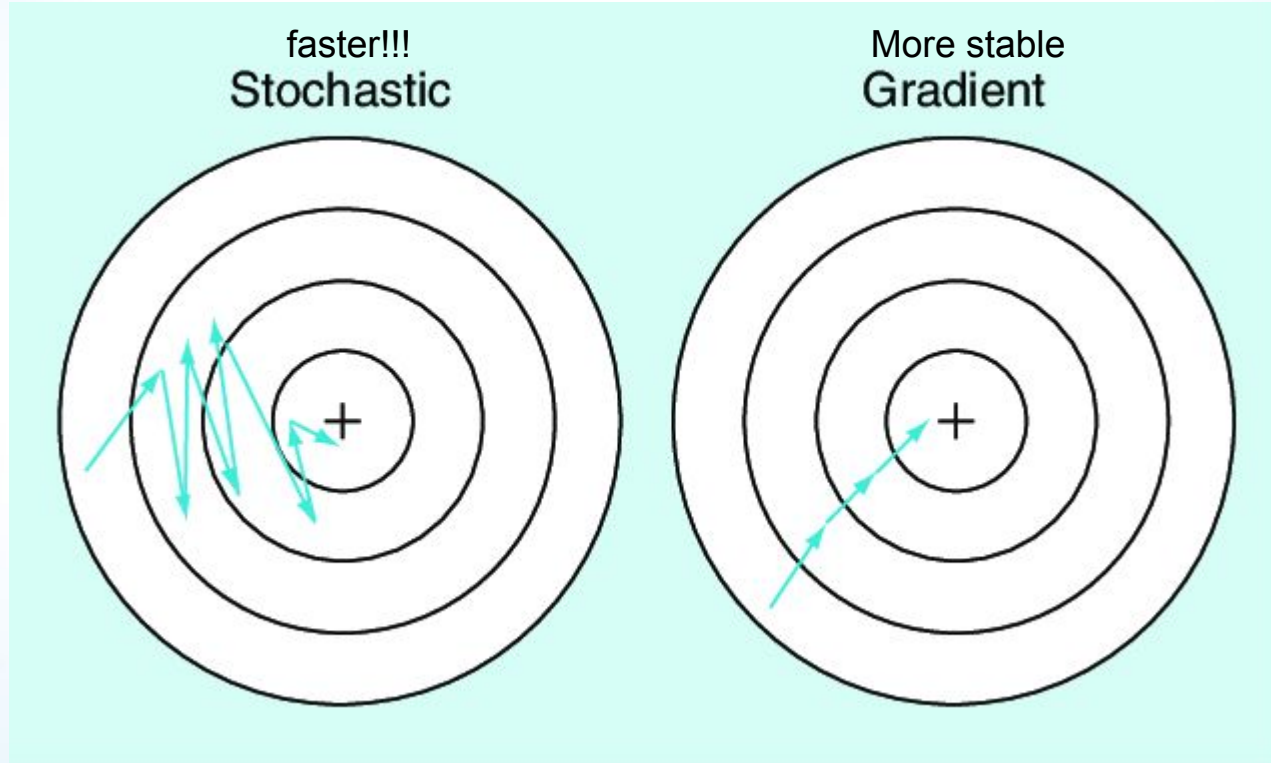
Input:

- Function $f(\mathbf{x}; \Theta)$ parameterized with parameters Θ .
- Training set of inputs $\mathbf{x}_1, \dots, \mathbf{x}_n$ and desired outputs $\mathbf{y}_1, \dots, \mathbf{y}_n$.
- Loss function L .

```

1: while stopping criteria not met do
2:   Sample a training example  $\mathbf{x}_i, \mathbf{y}_i$ 
3:   Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$ 
4:    $\hat{\mathbf{g}} \leftarrow$  gradients of  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$  w.r.t  $\Theta$ 
5:    $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$ 
6: return  $\Theta$ 
    
```

Gradient Descent vs Stochastic Gradient Descent



How can we combine both?

Mini-batch Gradient Descent Training

Algorithm 3 Minibatch Stochastic Gradient Descent Training

Input:

- Function $f(\mathbf{x}; \Theta)$ parameterized with parameters Θ .
 - Training set of inputs $\mathbf{x}_1, \dots, \mathbf{x}_n$ and desired outputs $\mathbf{y}_1, \dots, \mathbf{y}_n$.
 - Loss function L .
-

```

1: while stopping criteria not met do
2:   Sample a minibatch of  $m$  examples  $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_m, \mathbf{y}_m)\}$ 
3:    $\hat{\mathbf{g}} \leftarrow \mathbf{0}$ 
4:   for  $i = 1$  to  $m$  do
5:     Compute the loss  $L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i)$ 
6:      $\hat{\mathbf{g}} \leftarrow \hat{\mathbf{g}} + \text{gradients of } \frac{1}{m}L(f(\mathbf{x}_i; \Theta), \mathbf{y}_i) \text{ w.r.t } \Theta$ 
7:    $\Theta \leftarrow \Theta - \eta_t \hat{\mathbf{g}}$ 
8: return  $\Theta$ 

```

Gradient Descent Toy Example



Problem - Two Iterations of GD

3 inputs: x_1, x_2, x_3 (and $x_0=1$)

1 output: y

Initial weights: $w_0 = -0.7, w_1 = 0.2, w_2 = 0.7, w_3 = 0.9$

The true output: 0

The learning rate: 0.2

Let's say the gradient is based on the perceptron rule: $(y_{\text{pred}} - y_{\text{true}}) x$

Show 2 iterations of gradient descent

Solution First Iteration

1. Original output:

$$Wx = w_0 \cdot x_0 + w_1 \cdot x_1 + w_2 \cdot x_2 + w_3 \cdot x_3$$

$$= -0.7 \cdot 1 + 0.2 \cdot 0 + 0.7 \cdot 0 + 0.9 \cdot 1 = 0.2 > 0 \rightarrow Y_{\text{observed}} = 1$$

2. Let's calc the gradient:

$$a. \quad g_0 = (o - t) \cdot x_0 = (1 - 0) \cdot 1 = 1$$

$$g_1 = (o - t) \cdot x_1 = (1 - 0) \cdot 0 = 0$$

$$g_2 = (o - t) \cdot x_2 = (1 - 0) \cdot 0 = 0$$

$$g_3 = (o - t) \cdot x_3 = (1 - 0) \cdot 1 = 1$$

3. Update weights:

$$a. \quad w_0 = -0.7 - (0.2 \cdot 1) = -0.9$$

$$w_1 = 0.2 - (0.2 \cdot 0) = 0.2$$

$$w_2 = 0.7 - (0.2 \cdot 0) = 0.7$$

$$w_3 = 0.9 - (0.2 \cdot 1) = 0.7$$

$$X_0 = 1, x_1 = 0, x_2 = 0, x_3 = 1$$

$$Y_{\text{true}} = 0$$

Iteration 0:

$$W_0 = -0.7, w_1 = 0.2, w_3 = 0.7, w_4 = 0.9$$

$$w^t \leftarrow w^{t-1} - \eta \cdot \frac{\partial \ell}{\partial w^{t-1}}$$

0.2

Solution Second Iteration

1. Original output:

....

$X_0 = 1, x_1 = 0, x_2 = 0, x_3 = 1$
 $Y_{\text{true}} = 0$

Iteration 1:

$W_0 = -0.9, w_1 = 0.2, w_3 = 0.7, w_4 = 0.7$

4. New output:

$wx = 1 * (-0.9) + 0 * 0.2 + 0 * 0.7 + 1 * 0.7 = -0.2 < 0 \rightarrow Y_{\text{observed}} = 0$

$$w^t \leftarrow w^{t-1} - \eta \cdot \frac{\partial \ell}{\partial w^{t-1}}$$

0.2

Solution Second Iteration → No Change

1. Original output:

$$wX = 1 * (-0.9) + 0 * 0.2 + 0 * 0.7 + 1 *$$

$$0.7 = -0.2 < 0 \rightarrow \mathbf{Y \text{ observed} = 0}$$

2. Let's calc the gradient:

$$a. \quad g_0 = (o - t) * x_0 = (0 - 0) * 1 = 0$$

$$g_1 = (o - t) * x_1 = (0 - 0) * 0 = 0$$

$$g_2 = (o - t) * x_2 = (0 - 0) * 0 = 0$$

$$g_3 = (o - t) * x_3 = (0 - 0) * 1 = 0$$

3. Update weights:

$$a. \quad w_0 = \mathbf{-0.9} - (0.2 * 0) = \mathbf{-0.9}$$

$$w_1 = \mathbf{0.2} - (0.2 * 0) = \mathbf{0.2}$$

$$w_2 = \mathbf{0.7} - (0.2 * 0) = \mathbf{0.7}$$

$$w_3 = \mathbf{0.7} - (0.2 * 0) = \mathbf{0.7}$$

$$X_0 = 1, x_1 = 0, x_2 = 0, x_3 = 1$$

$$Y_{\text{true}} = 1$$

Iteration 2:

$$W_0 = -0.9, w_1 = 0.2, w_3 = 0.7, w_4 = 0.7$$

Convergence ! Gradient and weights didn't change.

Next Week

- Regularization
- Linear Regression
- Logistic Regression
- Code
- Summary