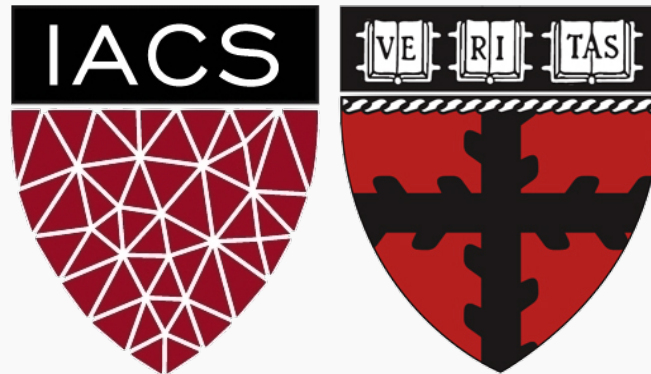# Gradient Descent

CS109B Data Science 2
Pavlos Protopapas, Mark Glickman
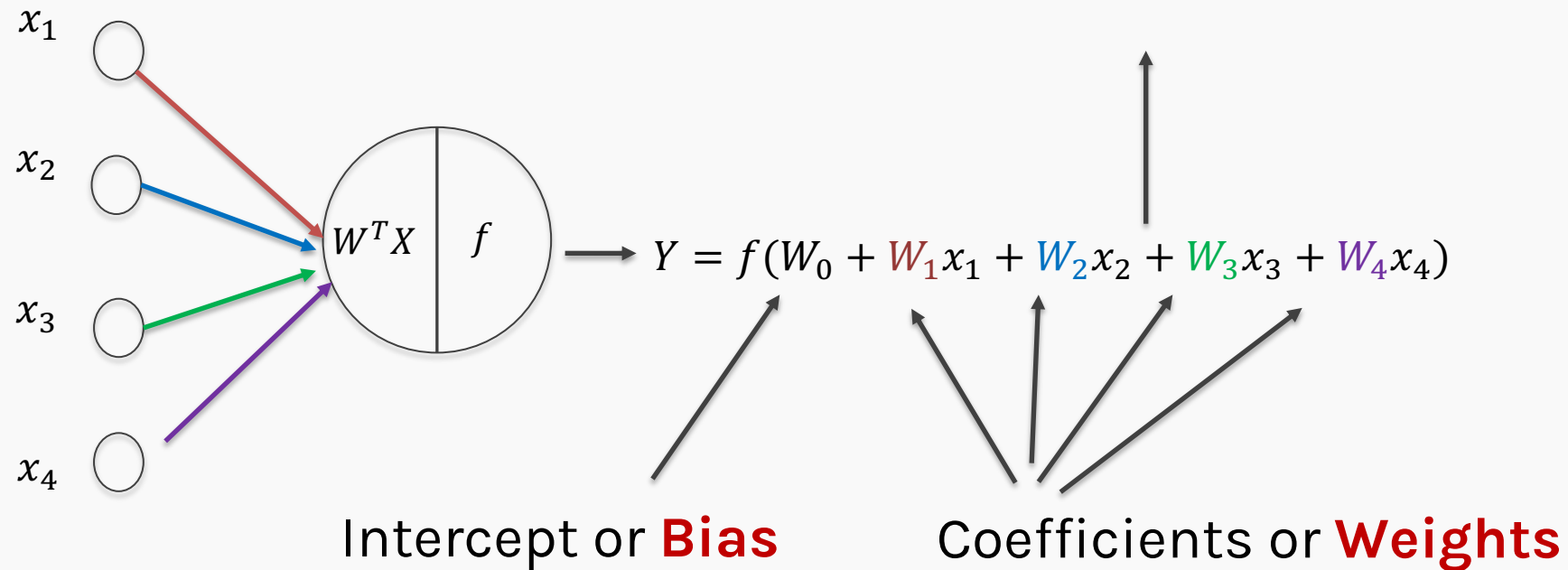
# Learning the coefficients

## Start with single neuron

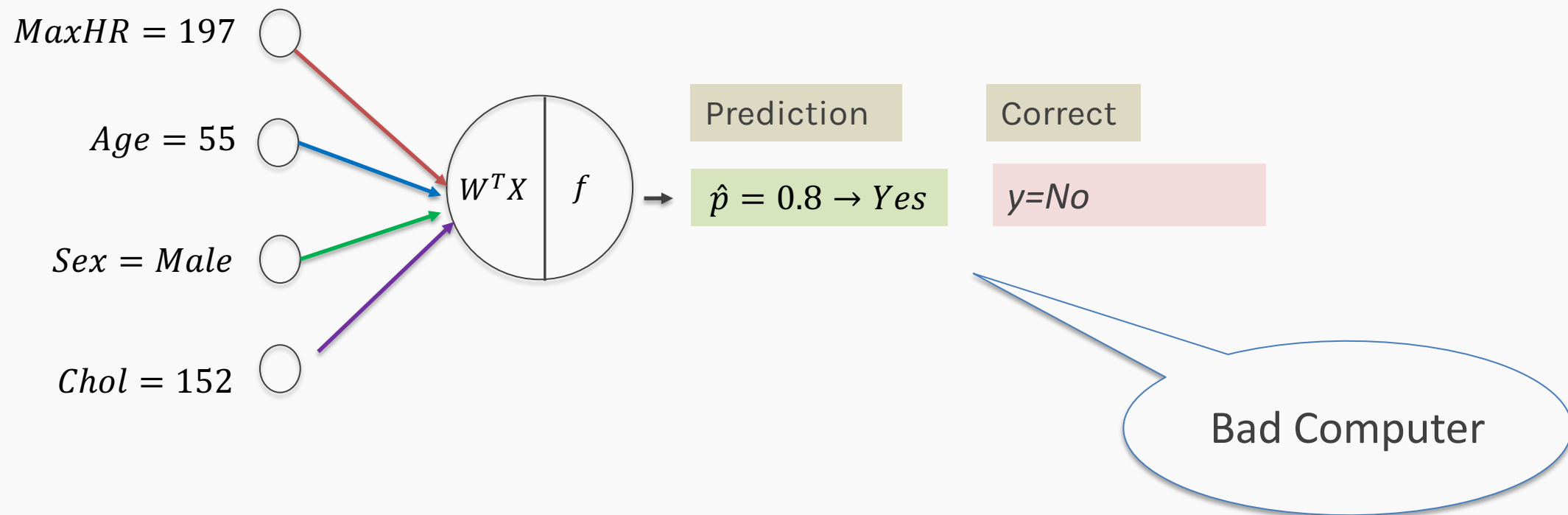Classification:
activation is sigmoid

$$f(X) = \frac{1}{1 + e^{-W^T X}}$$
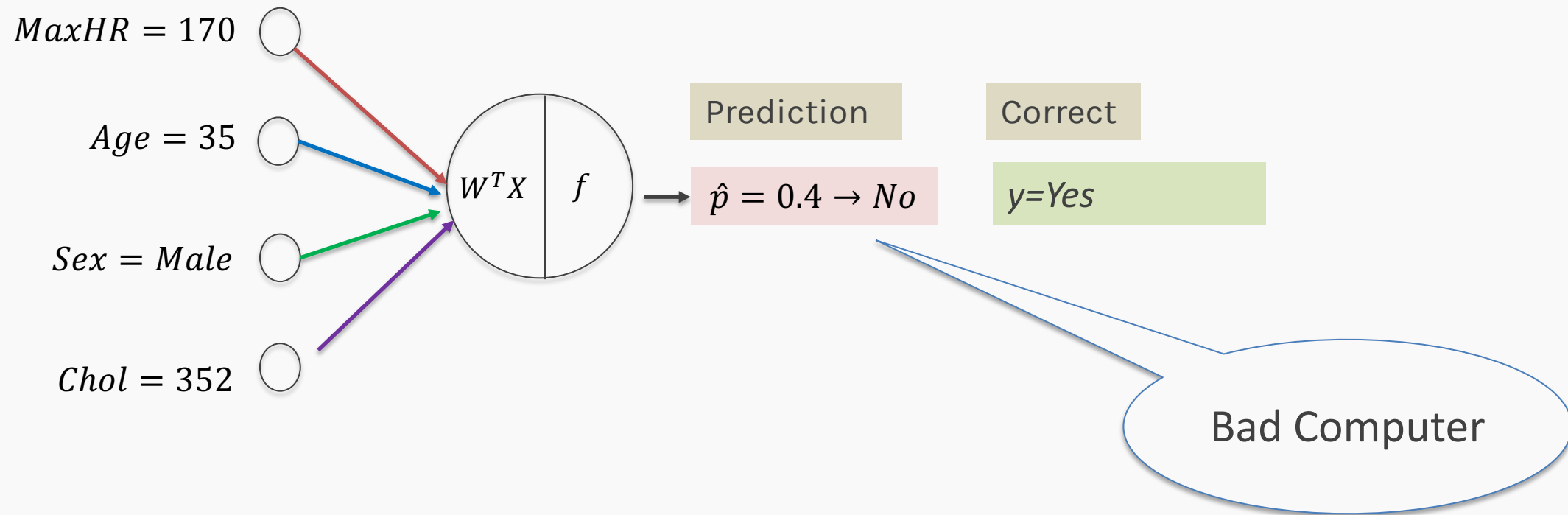
Regression:
activation is linear

$$f(X) = W^T X$$



$x_1$

$x_2$

$x_3$

$x_4$

$W^T X$  |  $f$

$Y = f(W_0 + W_1 x_1 + W_2 x_2 + W_3 x_3 + W_4 x_4)$

Intercept or **Bias**

Coefficients or **Weights**

# But what is the idea?

Start with all randomly selected weights. Most likely it will perform horribly. For example, in our heart data, the model will be giving us the wrong answer.

$MaxHR = 197$

$Age = 55$

$W^T X \mid f$

$Sex = Male$

$Chol = 152$

Prediction

$\hat{p} = 0.8 \rightarrow Yes$

Correct

$y=No$

Bad Computer

# But what is the idea?

Start with all randomly selected weights. Most likely it will perform horribly. For example, in our heart data, the model will be giving us the wrong answer.



$MaxHR = 170$

$Age = 35$

$Sex = Male$

$Chol = 352$

$W^TX$ | $f$

Prediction

$\hat{p} = 0.4 \rightarrow No$

Correct

$y=Yes$

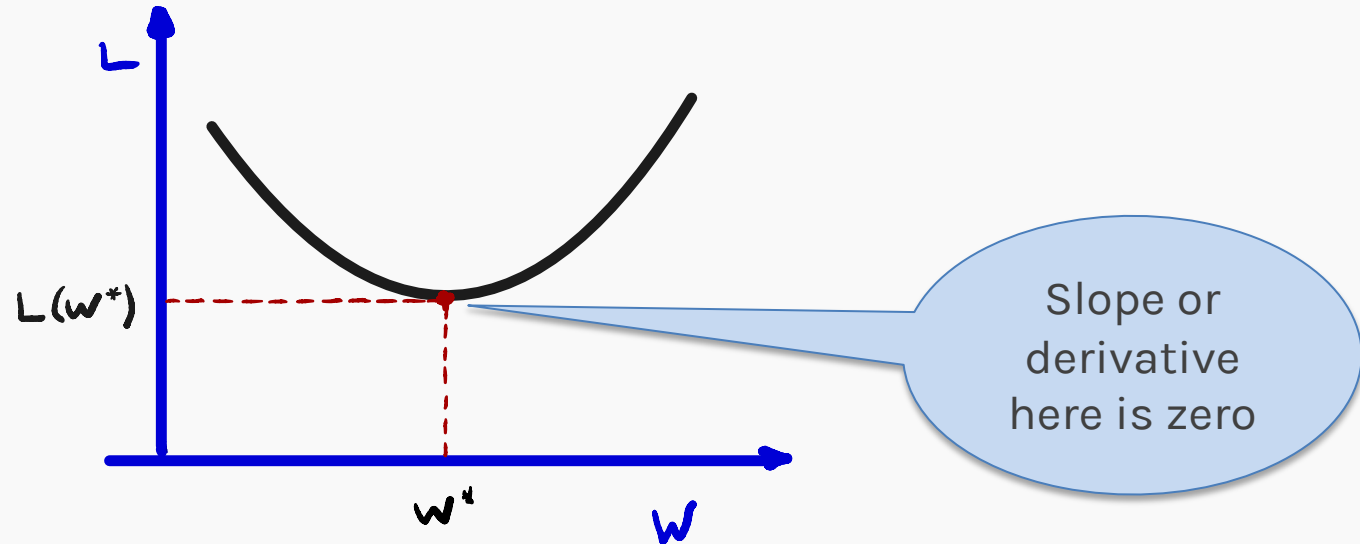Bad Computer

# But what is the idea?

- The **Loss Function** takes all these results and averages them and tells us how bad or good the computer or those weights are.

- Telling the computer how **bad** or **good** it is, does not help.

- You want to tell it how to change those weights, so it gets better.

Loss function: $\mathcal{L}(w_0, w_1, w_2, w_3, w_4)$

For now, let's only consider a single weight, $\mathcal{L}(w_1)$

# Minimizing the Loss function

Ideally, we want to know the value of $W$ that gives the minimal $\mathcal{L}(W)$
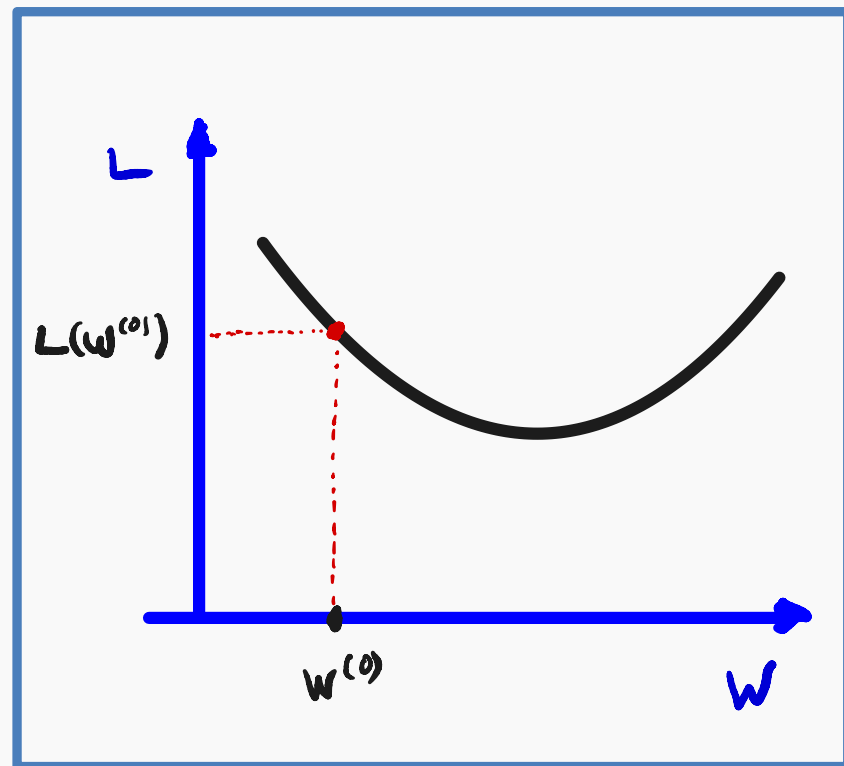


To find the optimal point of a function $\mathcal{L}(W)$, we take the derivative wrt to the weight:
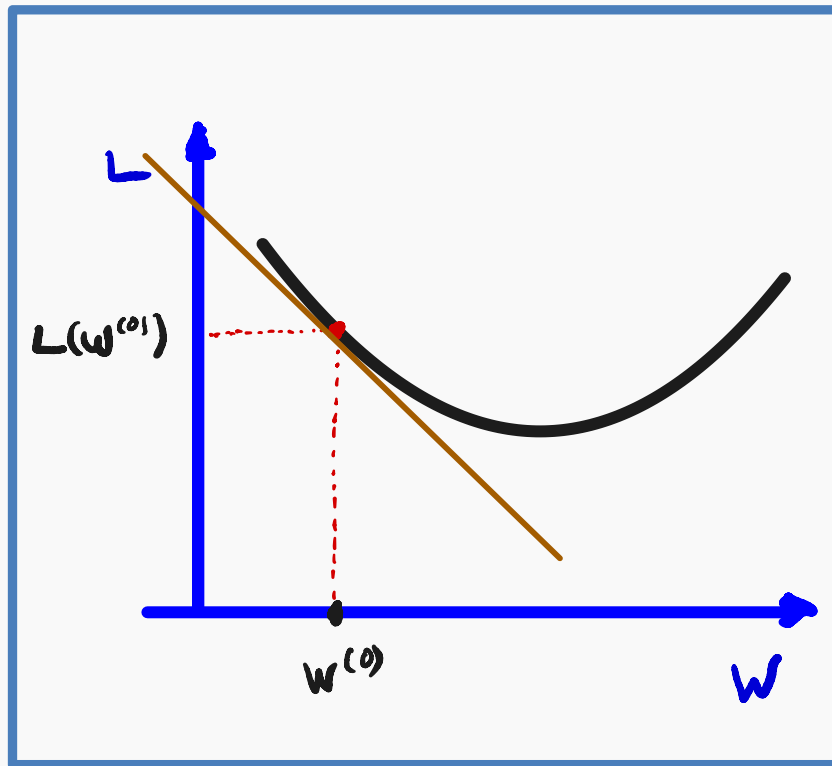
$$\frac{d\mathcal{L}(W)}{dW} = 0$$

And find the $W$ that satisfies that equation. **Sometimes** there is no explicit solution for that.

# Estimate of the regression coefficients: gradient descent
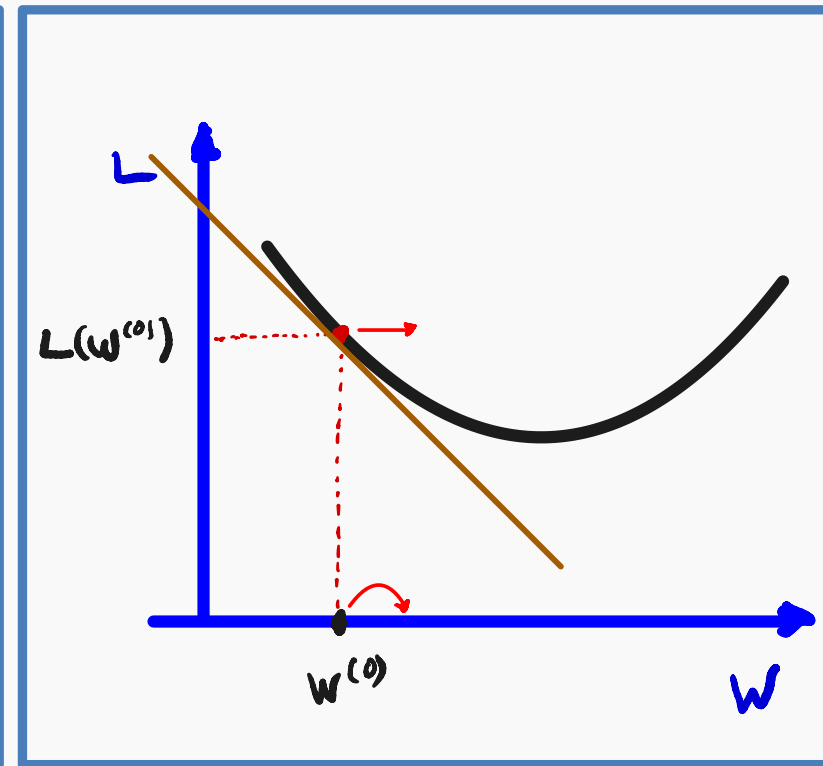
A more flexible method would be



Start from a random point

Compute the slope/derivative at this point

Step to the opposite direction of the derivative

A more flexible method would be



Compute the slope/derivative at $W^{(1)}$ and step again in the opposite direction of the derivative.

Continue,

Stop when no more improvement or after a certain number of iterations.

# Estimate of the regression coefficients:  gradient descent

**Question**: How do we generalize this to more than one weight?

**Take the gradient:**

$$\nabla_W L(W) = \left[ \frac{\partial L}{\partial W_1}, \ \frac{\partial L}{\partial W_2}, \ \dots, \frac{\partial L}{\partial W_p} \right]$$

**Question:**  What do you think is a good approach for telling the model how to change (what is the step size) to become better?

# Gradient Descent (cont.)

If the step is proportional to the slope, then you avoid overshooting the minimum. How?

# Gradient Descent (cont.)

If the step is proportional to the slope, then you avoid overshooting the minimum. How?

$$\frac{d\mathcal{L}(W)}{dW}$$

$$\frac{d\mathcal{L}(W)}{dW}$$

$$\frac{d\mathcal{L}(W)}{dW}$$

# Gradient Descent

We know that we want to go in the opposite direction of the derivative, and we know we want to be making a step proportional to the derivative.

Making a step means:

$$w^{new} = w^{old} + step$$

Step size is proportional to derivative

Opposite direction of the derivative and proportional to the derivative means:

$$w^{new} = w^{old} - \eta \frac{d\mathcal{L}}{dw}$$

Learning Rate

Change to more conventional notation:

$$w^{(i+1)} = w^{(i)} - \eta \frac{d\mathcal{L}}{dw}$$

# Gradient Descent

- Algorithm for optimization of first order to finding a minimum of a function.

- It is an iterative method.

- $L$ is decreasing much faster in the direction of the negative derivative.

- The learning rate is controlled by the magnitude of $\eta$.

$$w^{(i+1)} = w^{(i)} - \eta \frac{d\mathcal{L}}{dw}$$

# Gradient Descent Considerations

- We still need to calculate the derivatives.

- We need to set the learning rate.

- Local vs global minima.

- The full likelihood function includes summing up all individual 'errors'. Sometimes this includes hundreds of thousands of examples.

# Gradient Descent Considerations

- **We still need to calculate the derivatives.**

- We need to set the learning rate.

- Local vs global minima.

- The full likelihood function includes summing up all individual '*errors*'. Sometimes this includes hundreds of thousands of examples.

# Calculate the Derivatives

Can we do it? Can we calculate the derivative of  any loss function?

**Wolfram Alpha** can do it for us!

We need a formalism to deal with these derivatives.

# Chain Rule

Chain rule for computing gradients:

$$y = g(x) \qquad z = f(y) = f(g(x))$$

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

$$\boldsymbol{y} = g(\boldsymbol{x}) \qquad z = f(\boldsymbol{y}) = f(g(\boldsymbol{x}))$$

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

For longer chains:

$$\frac{\partial z}{\partial x_i} = \sum_{j_1} \cdots \sum_{j_m} \frac{\partial z}{\partial y_{j_1}} \cdots \frac{\partial y_{j_m}}{\partial x_i}$$

# Logistic Regression derivatives

For logistic regression, the –ve log of the likelihood is:

$$\mathcal{L} = \sum_i \mathcal{L}_i = -\sum_i \log L_i = -\sum_i [y_i \log p_i + (1 - y_i) \log(1 - p_i)]$$

$$\mathcal{L}_i = -y_i \log \frac{1}{1 + e^{-W^T X}} - (1 - y_i) \log(1 - \frac{1}{1 + e^{-W^T X}})$$

To simplify the analysis let us split it into two parts,

$$\mathcal{L}_i = \mathcal{L}_i^A + \mathcal{L}_i^B$$

So the derivative with respect to $W$ is:

$$\frac{\partial \mathcal{L}}{\partial W} = \sum_i \frac{\partial \mathcal{L}_i}{\partial W} = \sum_i (\frac{\partial \mathcal{L}_i^A}{\partial W} + \frac{\partial \mathcal{L}_i^B}{\partial W})$$

$$\mathcal{L}_i^A = -y_i \log \frac{1}{1 + e^{\boxed{-W^T X}}}$$

| Variables | Partial derivatives | Partial derivatives |
|---|---|---|
| $\xi_1 = -W^T X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ |
| $\xi_2 = e^{\xi_1} = e^{-W^T X}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{-W^T X}$ |
| $\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ |
| $\xi_4 = \dfrac{1}{\xi_3} = \dfrac{1}{1 + e^{-W^T X}} = p$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\xi_3^2}$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\left(1 + e^{-W^T X}\right)^2}$ |
| $\xi_5 = \log \xi_4 = \log p = \log \dfrac{1}{1 + e^{-W^T X}}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = \dfrac{1}{\xi_4}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = 1 + e^{-W^T X}$ |
| $\mathcal{L}_i^A = -y\xi_5$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ |
| $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = \dfrac{\partial \mathcal{L}_i}{\partial \xi_5} \dfrac{\partial \xi_5}{\partial \xi_4} \dfrac{\partial \xi_4}{\partial \xi_3} \dfrac{\partial \xi_3}{\partial \xi_2} \dfrac{\partial \xi_2}{\partial \xi_1} \dfrac{\partial \xi_1}{\partial W}$ | | $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = -yX e^{-W^T X} \dfrac{1}{\left(1 + e^{-W^T X}\right)}$ |

$$\mathcal{L}_i^A = -y_i \log \frac{1}{1 + \boxed{e^{-W^T X}}}$$

| Variables | Partial derivatives | Partial derivatives |
|---|---|---|
| $\xi_1 = -W^T X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ |
| $\xi_2 = e^{\xi_1} = e^{-W^T X}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{-W^T X}$ |
| $\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ |
| $\xi_4 = \dfrac{1}{\xi_3} = \dfrac{1}{1 + e^{-W^T X}} = p$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\xi_3^2}$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\left(1 + e^{-W^T X}\right)^2}$ |
| $\xi_5 = \log \xi_4 = \log p = \log \dfrac{1}{1 + e^{-W^T X}}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = \dfrac{1}{\xi_4}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = 1 + e^{-W^T X}$ |
| $\mathcal{L}_i^A = -y \xi_5$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ |
| $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = \dfrac{\partial \mathcal{L}_i}{\partial \xi_5} \dfrac{\partial \xi_5}{\partial \xi_4} \dfrac{\partial \xi_4}{\partial \xi_3} \dfrac{\partial \xi_3}{\partial \xi_2} \dfrac{\partial \xi_2}{\partial \xi_1} \dfrac{\partial \xi_1}{\partial W}$ | | $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = -yX e^{-W^T X} \dfrac{1}{\left(1 + e^{-W^T X}\right)}$ |

$$\mathcal{L}_i^A = -y_i \log \boxed{\frac{1}{1 + e^{-W^T X}}}$$

| Variables | Partial derivatives | Partial derivatives |
|---|---|---|
| $\xi_1 = -W^T X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ |
| $\xi_2 = e^{\xi_1} = e^{-W^T X}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{-W^T X}$ |
| $\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ |
| $\xi_4 = \dfrac{1}{\xi_3} = \dfrac{1}{1 + e^{-W^T X}} = p$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\xi_3^2}$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\left(1 + e^{-W^T X}\right)^2}$ |
| $\xi_5 = \log \xi_4 = \log p = \log \dfrac{1}{1 + e^{-W^T X}}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = \dfrac{1}{\xi_4}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = 1 + e^{-W^T X}$ |
| $\mathcal{L}_i^A = -y \xi_5$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ |
| $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = \dfrac{\partial \mathcal{L}_i}{\partial \xi_5} \dfrac{\partial \xi_5}{\partial \xi_4} \dfrac{\partial \xi_4}{\partial \xi_3} \dfrac{\partial \xi_3}{\partial \xi_2} \dfrac{\partial \xi_2}{\partial \xi_1} \dfrac{\partial \xi_1}{\partial W}$ | | $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = -y X e^{-W^T X} \dfrac{1}{\left(1 + e^{-W^T X}\right)}$ |

$$\mathcal{L}_i^A = -y_i \log \boxed{\frac{1}{1 + e^{-W^T X}}}$$

| Variables | Partial derivatives | Partial derivatives |
|---|---|---|
| $\xi_1 = -W^T X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ |
| $\xi_2 = e^{\xi_1} = e^{-W^T X}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{-W^T X}$ |
| $\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ |
| $\xi_4 = \dfrac{1}{\xi_3} = \dfrac{1}{1 + e^{-W^T X}} = p$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\xi_3^2}$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\left(1 + e^{-W^T X}\right)^2}$ |
| $\xi_5 = \log \xi_4 = \log p = \log \dfrac{1}{1 + e^{-W^T X}}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = \dfrac{1}{\xi_4}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = 1 + e^{-W^T X}$ |
| $\mathcal{L}_i^A = -y \xi_5$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ |
| $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = \dfrac{\partial \mathcal{L}_i}{\partial \xi_5} \dfrac{\partial \xi_5}{\partial \xi_4} \dfrac{\partial \xi_4}{\partial \xi_3} \dfrac{\partial \xi_3}{\partial \xi_2} \dfrac{\partial \xi_2}{\partial \xi_1} \dfrac{\partial \xi_1}{\partial W}$ | | $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = -yX e^{-W^T X} \dfrac{1}{\left(1 + e^{-W^T X}\right)}$ |

$$\mathcal{L}_i^A = -y_i \boxed{\log \frac{1}{1 + e^{-W^T X}}}$$

| Variables | Partial derivatives | Partial derivatives |
|---|---|---|
| $\xi_1 = -W^T X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ |
| $\xi_2 = e^{\xi_1} = e^{-W^T X}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{-W^T X}$ |
| $\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ |
| $\xi_4 = \dfrac{1}{\xi_3} = \dfrac{1}{1 + e^{-W^T X}} = p$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\xi_3^2}$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\left(1 + e^{-W^T X}\right)^2}$ |
| $\xi_5 = \log \xi_4 = \log p = \log \dfrac{1}{1 + e^{-W^T X}}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = \dfrac{1}{\xi_4}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = 1 + e^{-W^T X}$ |
| $\mathcal{L}_i^A = -y\xi_5$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ |
| $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = \dfrac{\partial \mathcal{L}_i}{\partial \xi_5} \dfrac{\partial \xi_5}{\partial \xi_4} \dfrac{\partial \xi_4}{\partial \xi_3} \dfrac{\partial \xi_3}{\partial \xi_2} \dfrac{\partial \xi_2}{\partial \xi_1} \dfrac{\partial \xi_1}{\partial W}$ | | $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = -yX e^{-W^T X} \dfrac{1}{\left(1 + e^{-W^T X}\right)}$ |

$$\mathcal{L}_i^A = \boxed{-y_i \log \frac{1}{1 + e^{-W^T X}}}$$

| Variables | Partial derivatives | Partial derivatives |
|---|---|---|
| $\xi_1 = -W^T X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ |
| $\xi_2 = e^{\xi_1} = e^{-W^T X}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{-W^T X}$ |
| $\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ |
| $\xi_4 = \dfrac{1}{\xi_3} = \dfrac{1}{1 + e^{-W^T X}} = p$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\xi_3^2}$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\left(1 + e^{-W^T X}\right)^2}$ |
| $\xi_5 = \log \xi_4 = \log p = \log \dfrac{1}{1 + e^{-W^T X}}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = \dfrac{1}{\xi_4}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = 1 + e^{-W^T X}$ |
| $\mathcal{L}_i^A = -y\xi_5$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ |
| $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = \dfrac{\partial \mathcal{L}_i}{\partial \xi_5} \dfrac{\partial \xi_5}{\partial \xi_4} \dfrac{\partial \xi_4}{\partial \xi_3} \dfrac{\partial \xi_3}{\partial \xi_2} \dfrac{\partial \xi_2}{\partial \xi_1} \dfrac{\partial \xi_1}{\partial W}$ | | $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = -yX e^{-W^T X} \dfrac{1}{\left(1 + e^{-W^T X}\right)}$ |

$$\mathcal{L}_i^A = -y_i \log \frac{1}{1 + e^{-W^T X}}$$

| Variables | Partial derivatives | Partial derivatives |
|---|---|---|
| $\xi_1 = -W^T X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ | $\dfrac{\partial \xi_1}{\partial W} = -X$ |
| $\xi_2 = e^{\xi_1} = e^{-W^T X}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$ | $\dfrac{\partial \xi_2}{\partial \xi_1} = e^{-W^T X}$ |
| $\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ | $\dfrac{\partial \xi_3}{\partial \xi_2} = 1$ |
| $\xi_4 = \dfrac{1}{\xi_3} = \dfrac{1}{1 + e^{-W^T X}} = p$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\xi_3^2}$ | $\dfrac{\partial \xi_4}{\partial \xi_3} = -\dfrac{1}{\left(1 + e^{-W^T X}\right)^2}$ |
| $\xi_5 = \log \xi_4 = \log p = \log \dfrac{1}{1 + e^{-W^T X}}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = \dfrac{1}{\xi_4}$ | $\dfrac{\partial \xi_5}{\partial \xi_4} = 1 + e^{-W^T X}$ |
| $\mathcal{L}_i^A = -y \xi_5$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ | $\dfrac{\partial \mathcal{L}}{\partial \xi_5} = -y$ |
| $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = \dfrac{\partial \mathcal{L}_i}{\partial \xi_5} \dfrac{\partial \xi_5}{\partial \xi_4} \dfrac{\partial \xi_4}{\partial \xi_3} \dfrac{\partial \xi_3}{\partial \xi_2} \dfrac{\partial \xi_2}{\partial \xi_1} \dfrac{\partial \xi_1}{\partial W}$ | | $\dfrac{\partial \mathcal{L}_i^A}{\partial W} = -y X e^{-W^T X} \dfrac{1}{\left(1 + e^{-W^T X}\right)}$ |

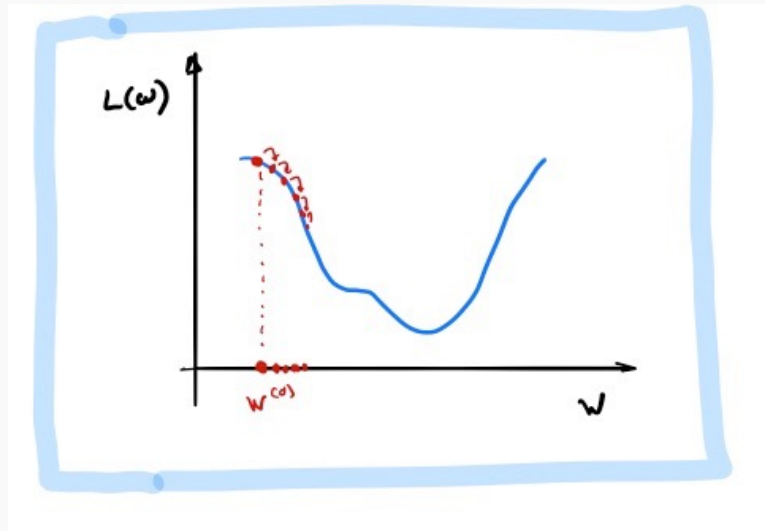$$\mathcal{L}_i^B = -(1 - y_i) \log[1 - \frac{1}{1 + e^{-W^T X}}]$$

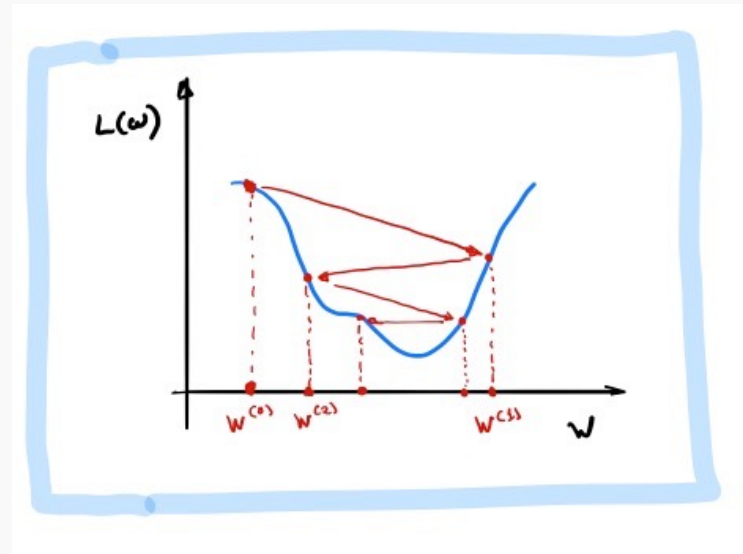| Variables | derivatives | Partial derivatives wrt to X,W |
|---|---|---|
| $\xi_1 = -W^T X$ | $\frac{\partial \xi_1}{\partial W} = -X$ | $\frac{\partial \xi_1}{\partial W} = -X$ |
| $\xi_2 = e^{\xi_1} = e^{-W^T X}$ | $\frac{\partial \xi_2}{\partial \xi_1} = e^{\xi_1}$ | $\frac{\partial \xi_2}{\partial \xi_1} = e^{-W^T X}$ |
| $\xi_3 = 1 + \xi_2 = 1 + e^{-W^T X}$ | $\frac{\partial \xi_3}{\partial \xi_2} = 1$ | $\frac{\partial \xi_3}{\partial 2} = 1$ |
| $\xi_4 = \frac{1}{\xi_3} = \frac{1}{1 + e^{-W^T X}} = p$ | $\frac{\partial \xi_4}{\partial \xi_3} = -\frac{1}{\xi_3^2}$ | $\frac{\partial \xi_4}{\partial \xi_3} = -\frac{1}{\left(1 + e^{-W^T X}\right)^2}$ |
| $\xi_5 = 1 - \xi_4 = 1 - \frac{1}{1 + e^{-W^T X}}$ | $\frac{\partial \xi_5}{\partial \xi_4} = -1$ | $\frac{\partial \xi_5}{\partial \xi_4} = -1$ |
| $\xi_6 = \log \xi_5 = \log(1 - p) = \log \frac{1}{1 + e^{-W^T X}}$ | $\frac{\partial \xi_6}{\partial \xi_5} = \frac{1}{\xi_5}$ | $\frac{\partial \xi_6}{\partial \xi_5} = \frac{1 + e^{-W^T X}}{e^{-W^T X}}$ |
| $\mathcal{L}_i^B = (1 - y)\xi_6$ | $\frac{\partial \mathcal{L}}{\partial \xi_6} = 1 - y$ | $\frac{\partial \mathcal{L}}{\partial \xi_6} = 1 - y$ |
| $\frac{\partial \mathcal{L}_i^B}{\partial W} = \frac{\partial \mathcal{L}_i^B}{\partial \xi_6} \frac{\partial \xi_6}{\partial \xi_5} \frac{\partial \xi_5}{\partial \xi_4} \frac{\partial \xi_4}{\partial \xi_3} \frac{\partial \xi_3}{\partial \xi_2} \frac{\partial \xi_2}{\partial \xi_1} \frac{\partial \xi_1}{\partial W}$ | | $\frac{\partial \mathcal{L}_i^B}{\partial W} = (1 - y)X \frac{1}{\left(1 + e^{-W^T X}\right)}$ |

# Considerations

- We still need to calculate the derivatives.

- **We need to set the learning rate.**

- Local vs global minima.

- The full likelihood function includes summing up all individual '*errors'.* Sometimes this includes hundreds of thousands of examples.
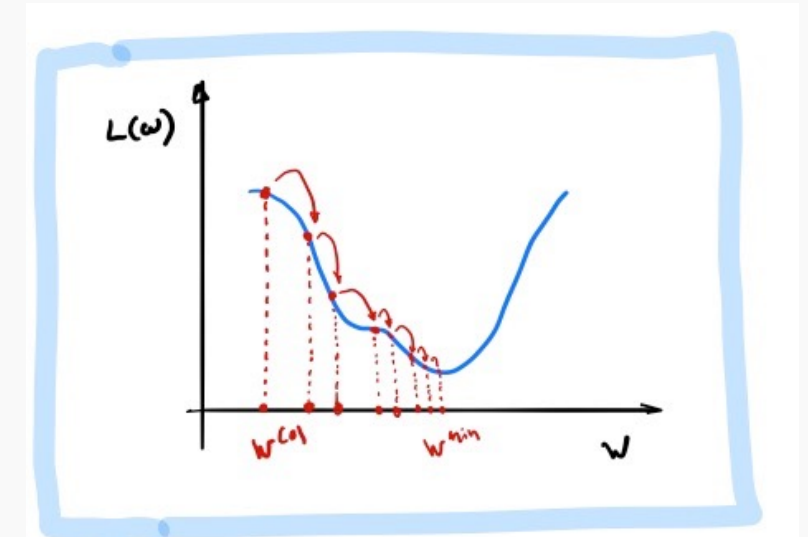
# Learning Rate

Our choice of the learning rate has a significant impact on the performance of gradient descent.



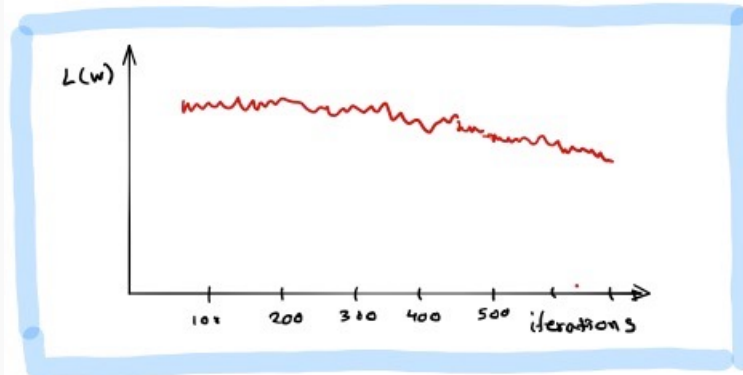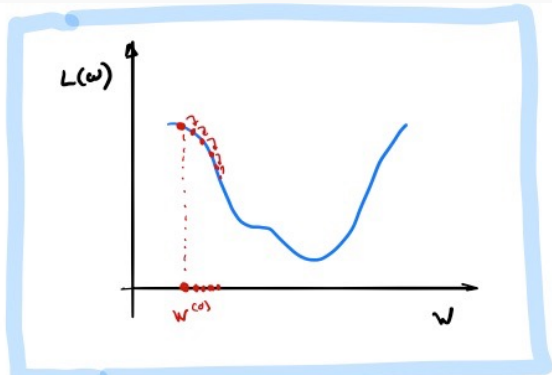When $\eta$ is too small, the algorithm makes very little progress.



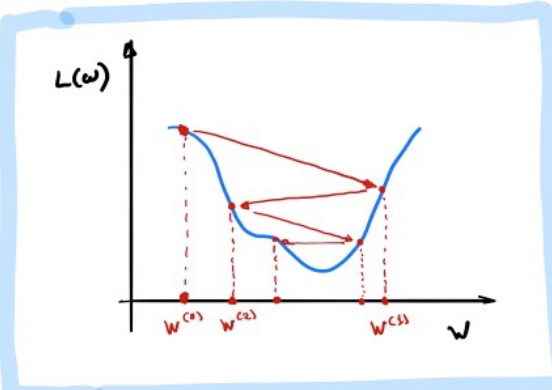When $\eta$ is too large, the algorithm may overshoot the minimum and has crazy oscillations.



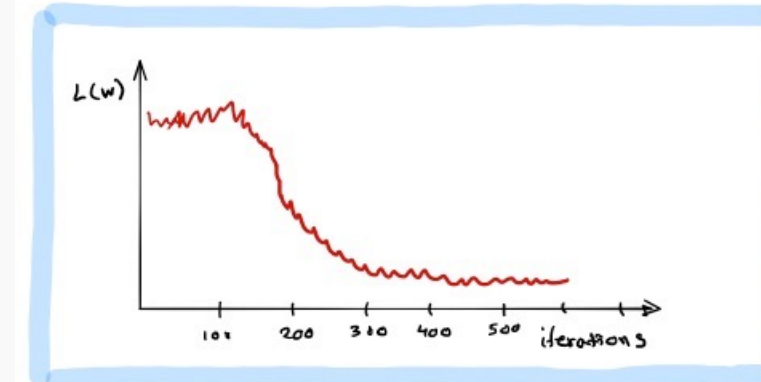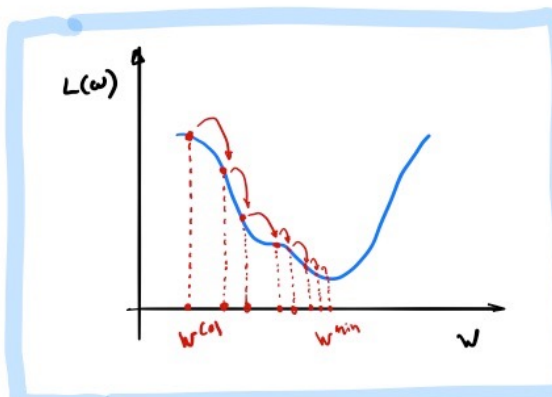When $\eta$ is appropriate, the algorithm will find the minimum. The algorithm **converges**!

How can we tell when gradient descent is converging? We visualize the loss function at each step of gradient descent. This is called the **trace plot**.



While the loss is decreasing throughout training, it does not look like descent hit the bottom.



Loss is mostly oscillating between values rather than converging.



The loss has decreased significantly during training. Towards the end, the loss stabilizes and it can't decrease further.

# Learning Rate

There are many alternative methods which address how to set or adjust the learning rate, using the derivative or second derivatives and or the momentum.

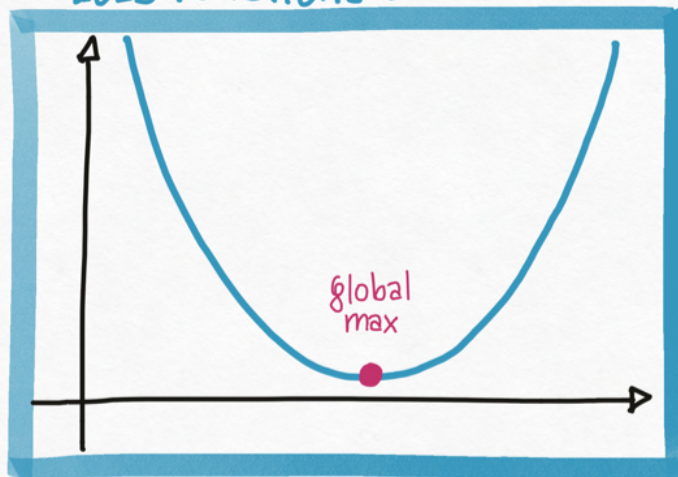**More on this later.**

# Considerations

- We still need to calculate the derivatives.

- We need to set the learning rate.

- **Local vs global minima.**

- The full likelihood function includes summing up all individual '*errors'.* Sometimes this includes hundreds of thousands of examples.

# Local vs Global Minima

If we choose $\eta$ correctly, then gradient descent will converge to a stationary point. But will this point be a **global minimum**?

If the function is convex then the stationary point will be a global minimum.
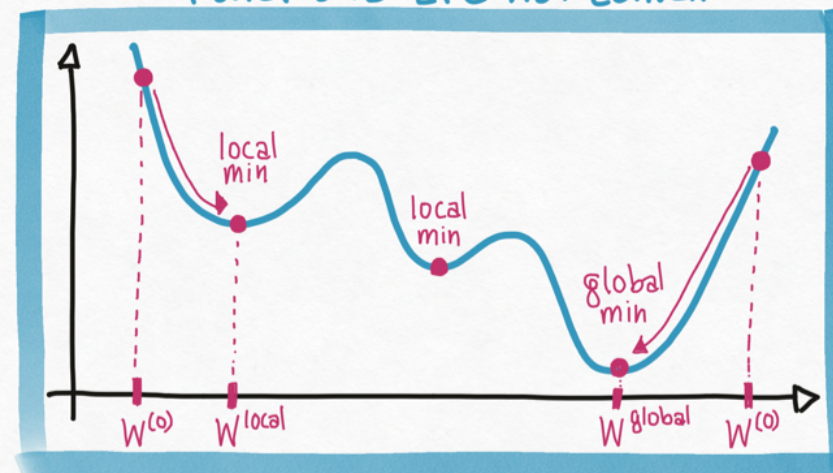


Linear & Polynomial Regression Loss Functions are Convex

global max

Hessian (2nd Derivative) positive semi-definite everywhere.
Every stationary point of the gradient is a global min.

Neural Network Regression Loss Functions are not Convex

local min

local min

global min

$W^{(0)}$    $W^{local}$    $W^{global}$    $W^{(0)}$

Neural networks with different weights can correspond to the same function.

Most stationary points are local minima but not global optima.

# Local vs Global Minima

No guarantee that we get the global minimum.

**Question:** What would be a good strategy?

- Random restarts
- Add noise to the loss function