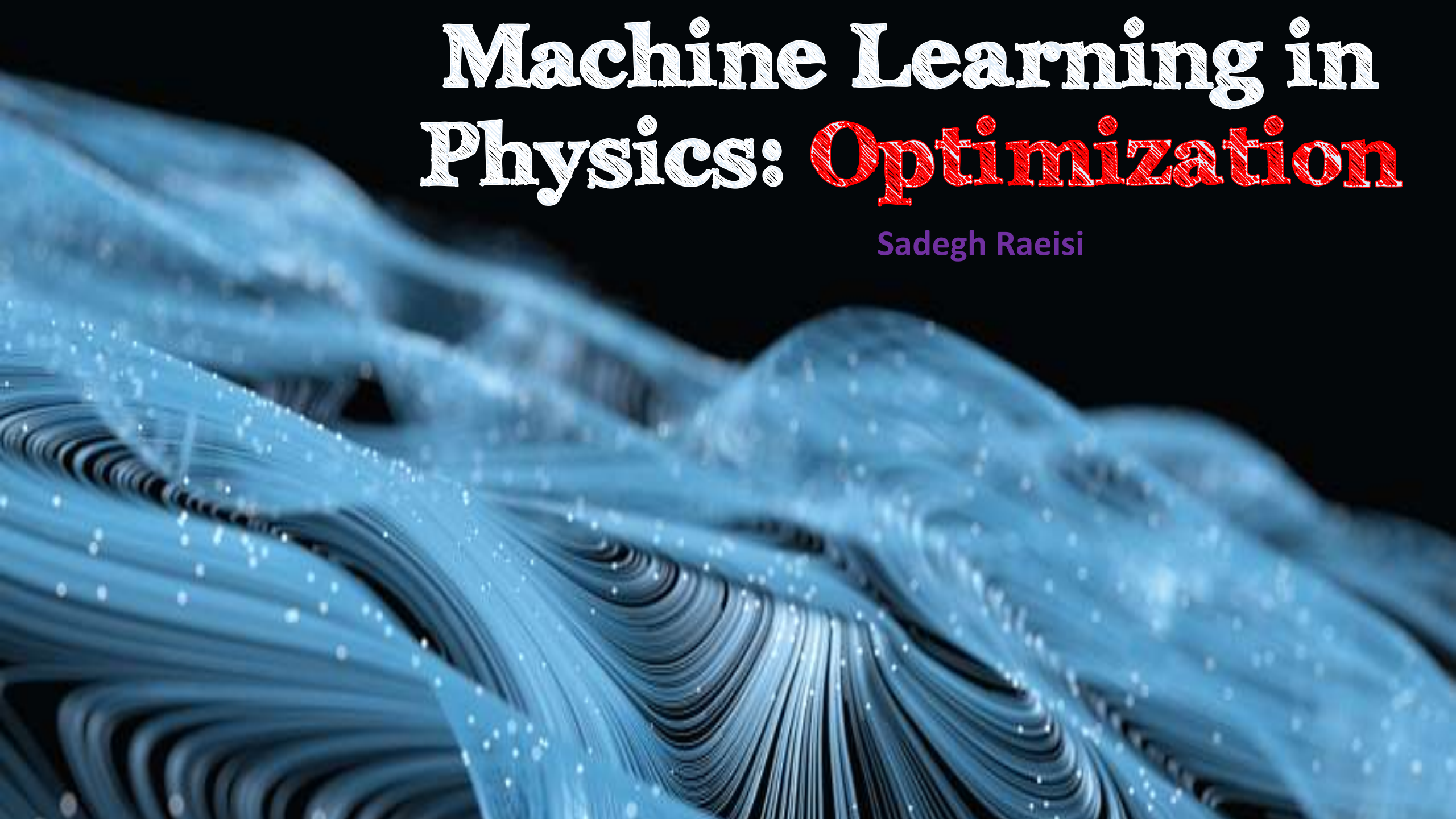
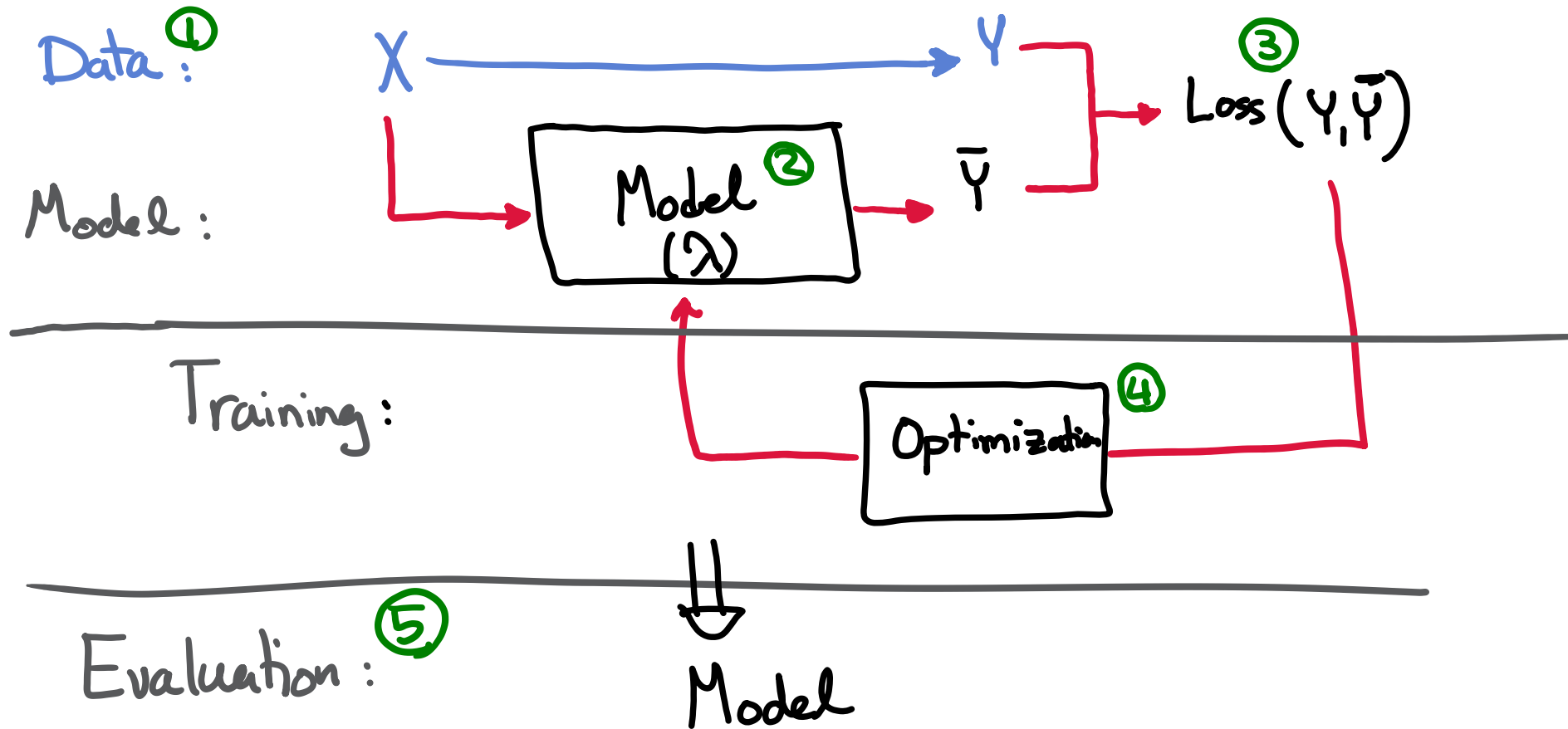


Machine Learning in Physics: **Optimization**

Sadegh Raeisi



Supervised: Ingredients



Outline

Concept

Gradient descent

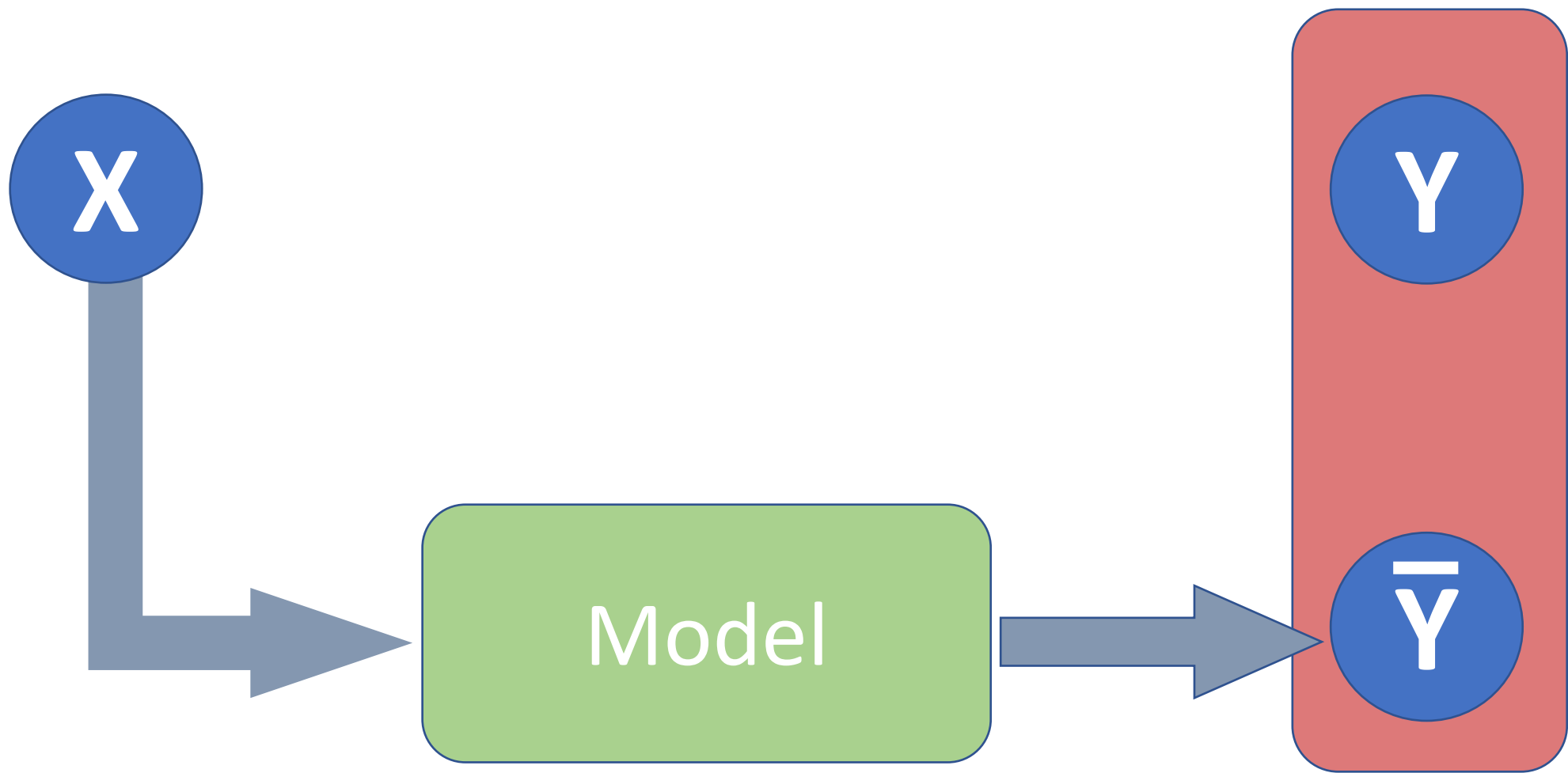
Challenges

Stochastic gradient descent

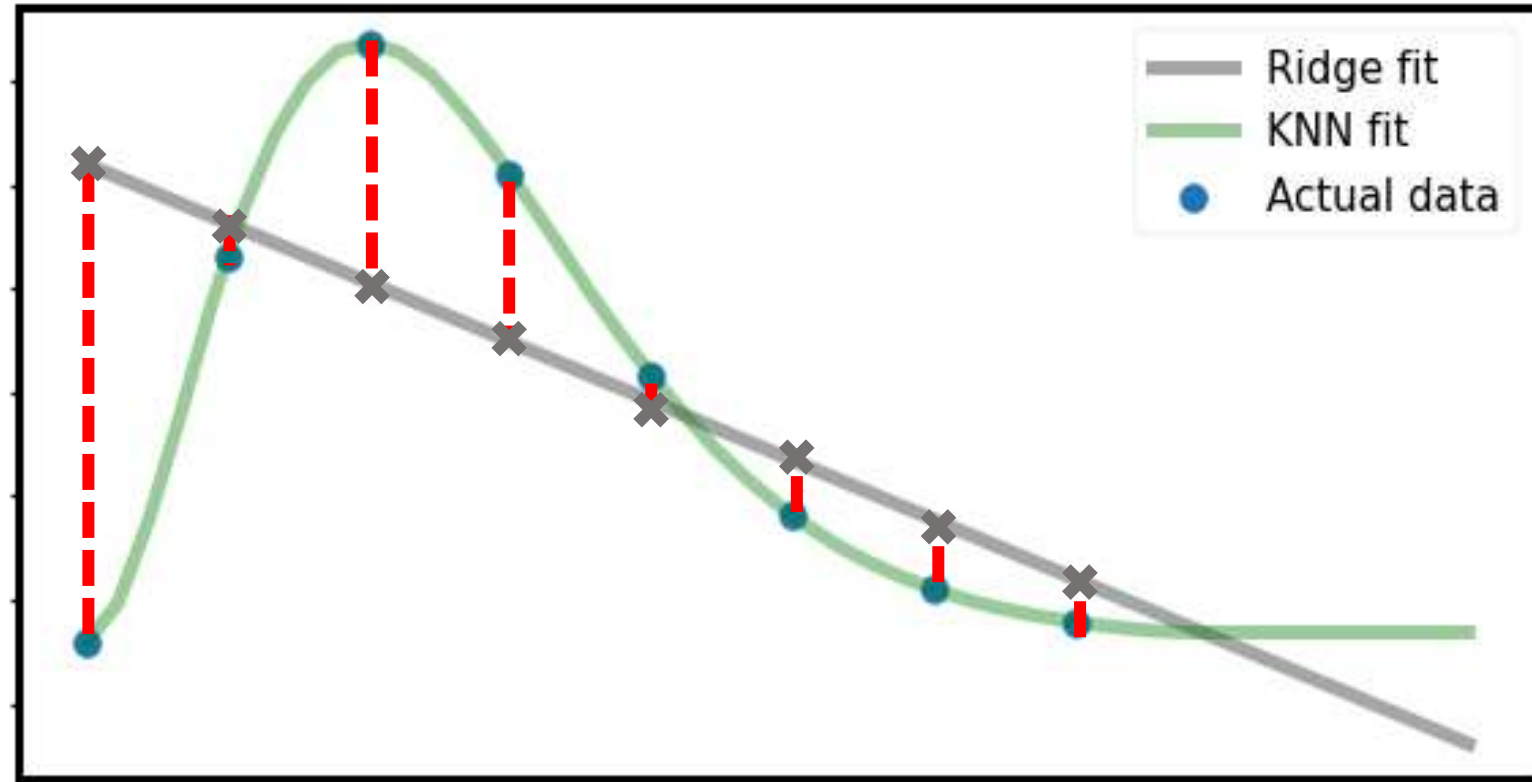
Weight decay

Adaptive learning rate


Concept



How close are the predictions?

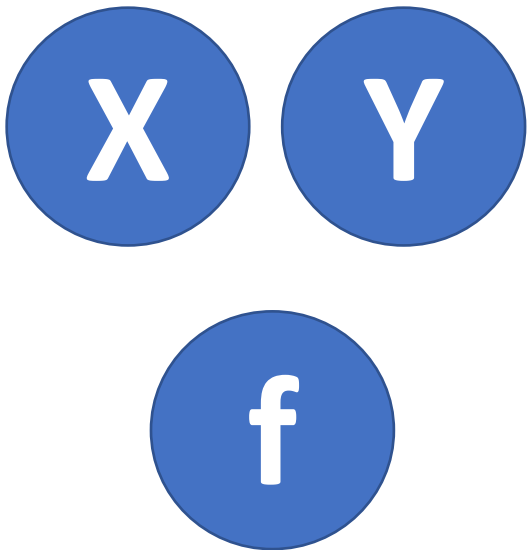


How can we quantify the difference?

$$\text{Dist}(\textcircled{Y}, \textcircled{\bar{Y}})$$

$$f_w(X)$$

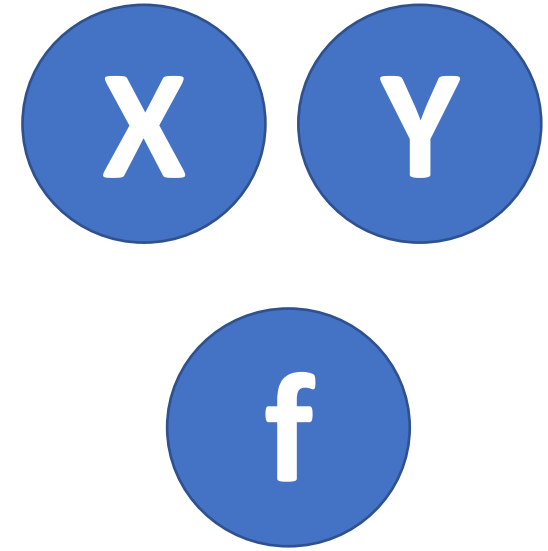
How can we quantify the difference?

$$\mathcal{L}_w (\textcircled{Y}, \textcircled{\bar{Y}})$$



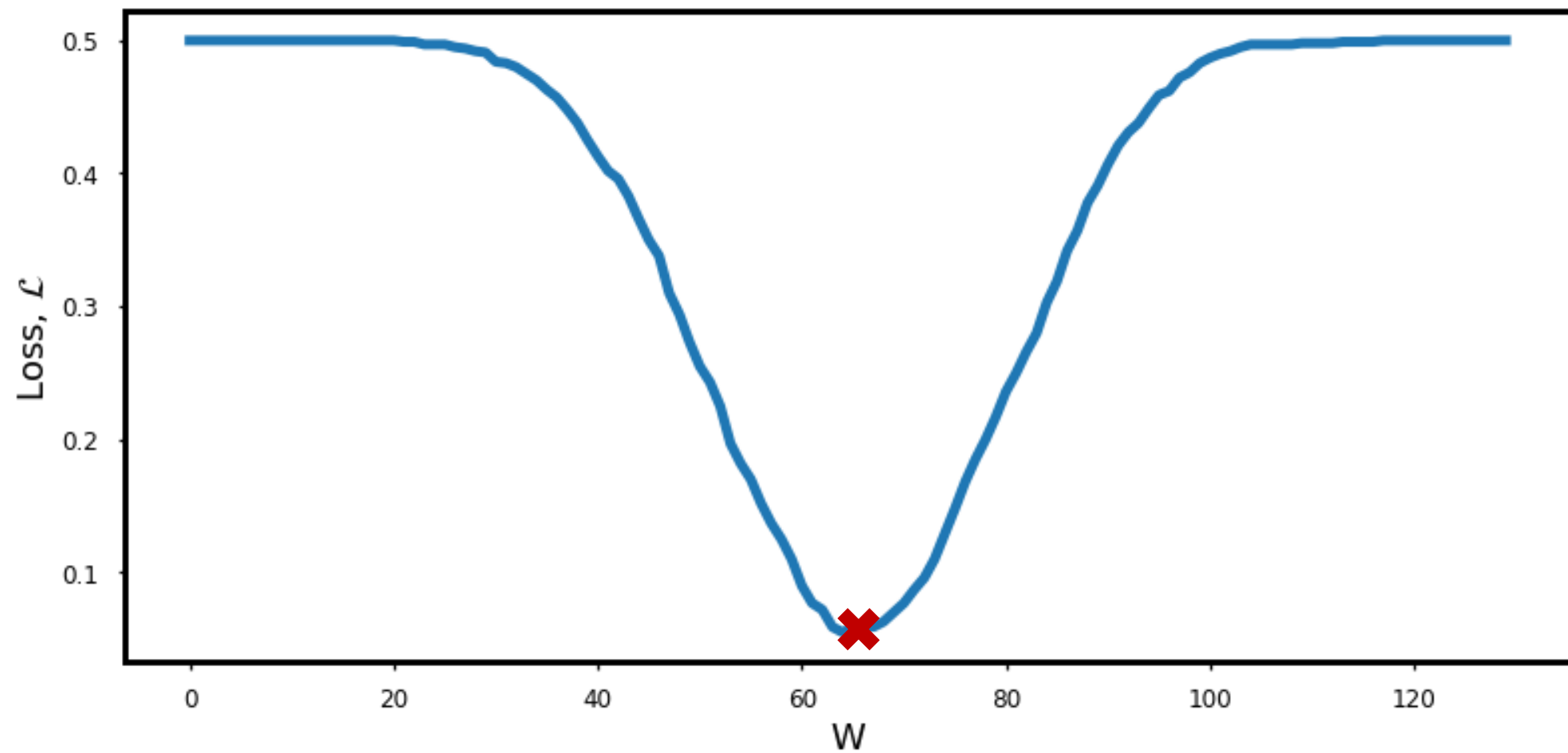
$$\mathcal{L}(w)$$

Problem statement



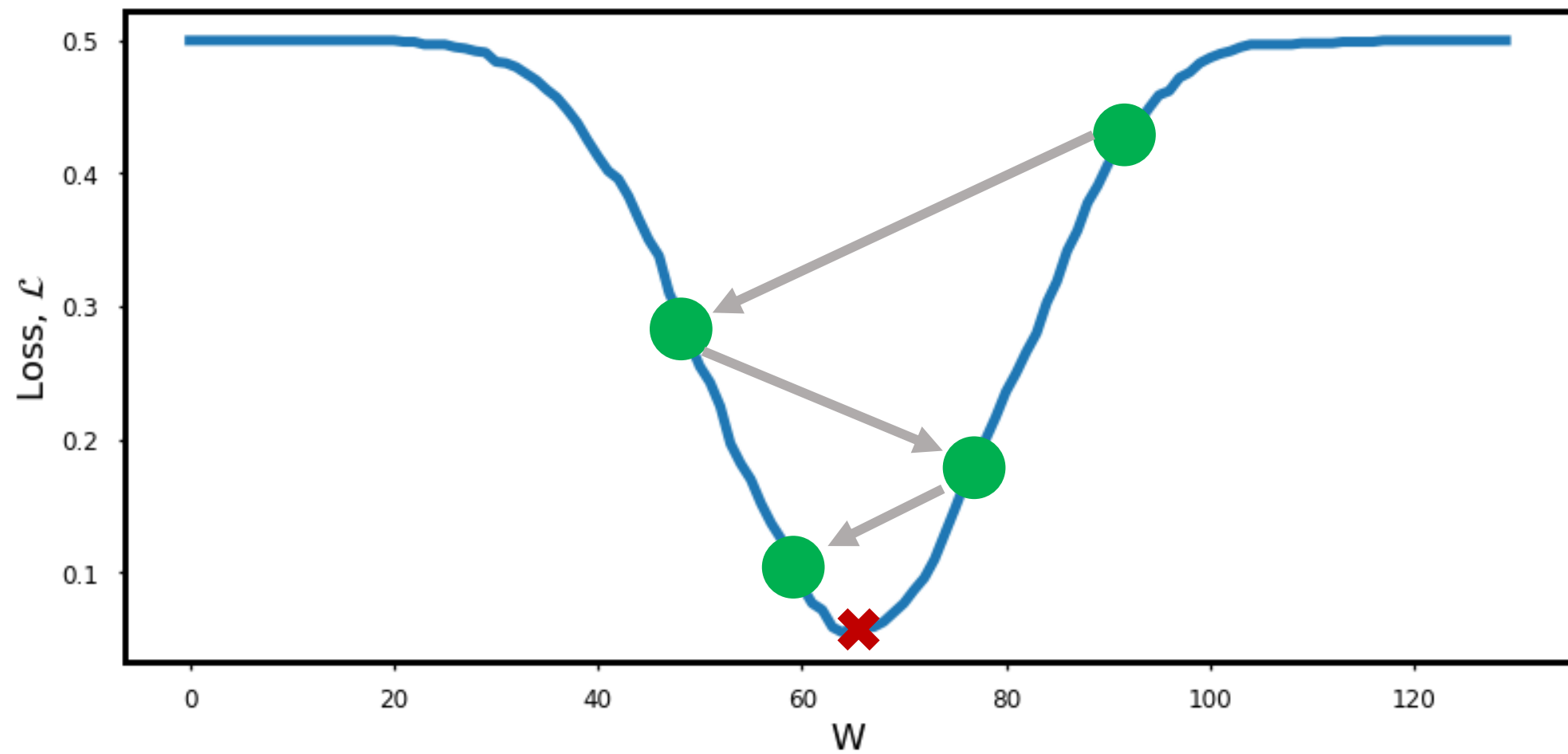
Minimize $\mathcal{L}(w)$
 w

Loss landscape



Gradient descent

Idea



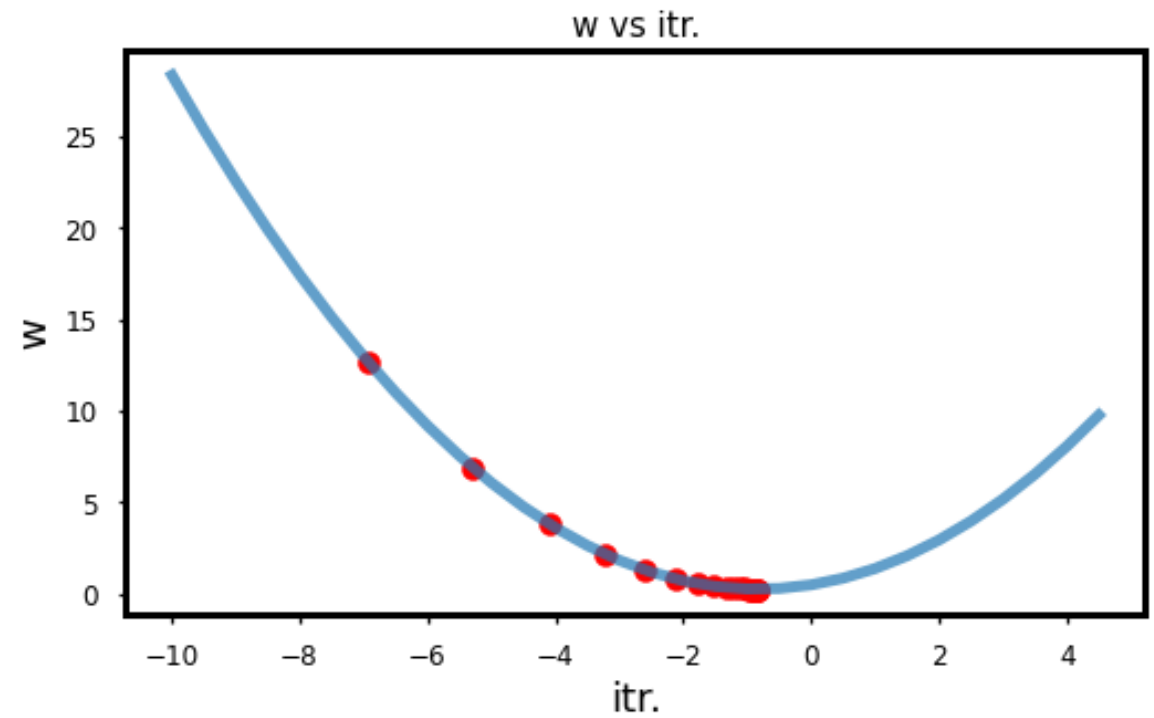
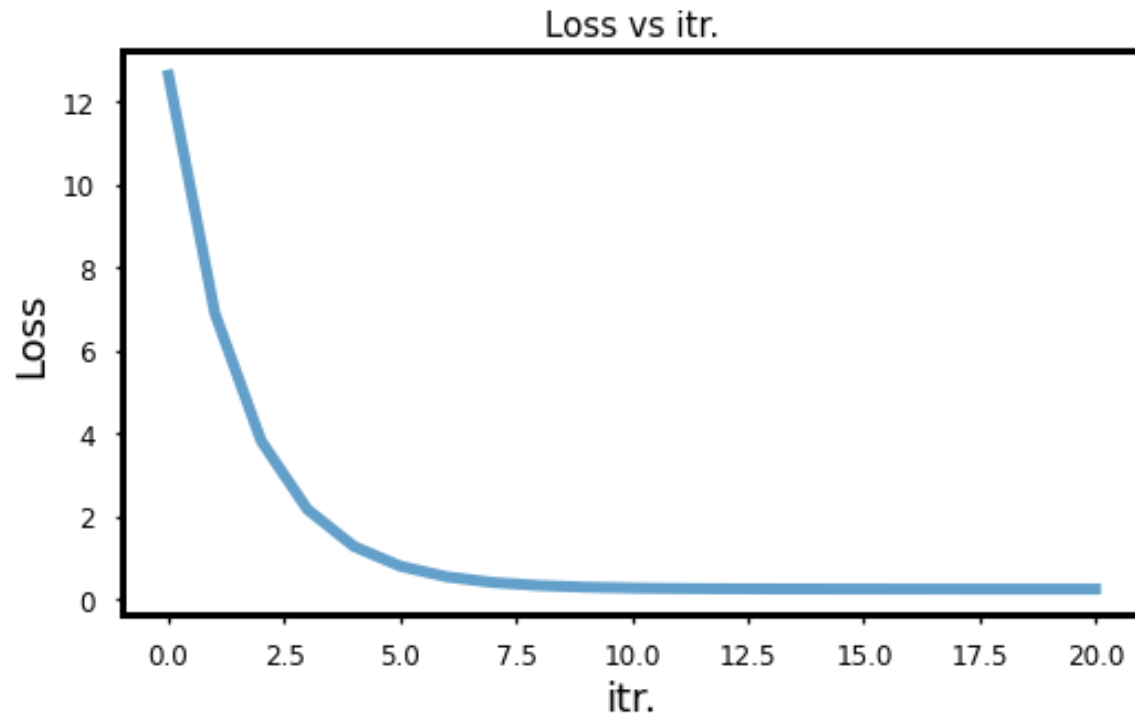
Algorithm

- Start with a random w .
- Calculate the derivative $\frac{\partial \mathcal{L}}{\partial w}$.
- Update $w \rightarrow w - \eta \frac{\partial \mathcal{L}}{\partial w}$.
- Repeat

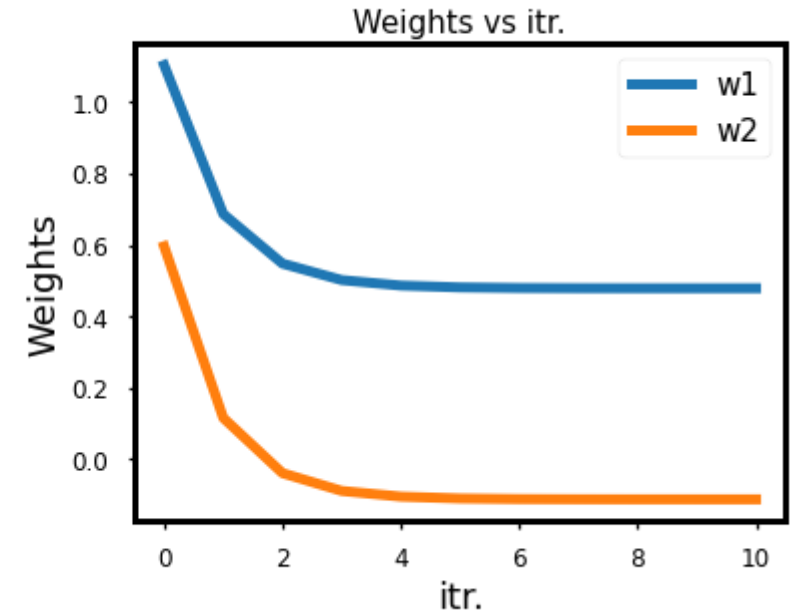
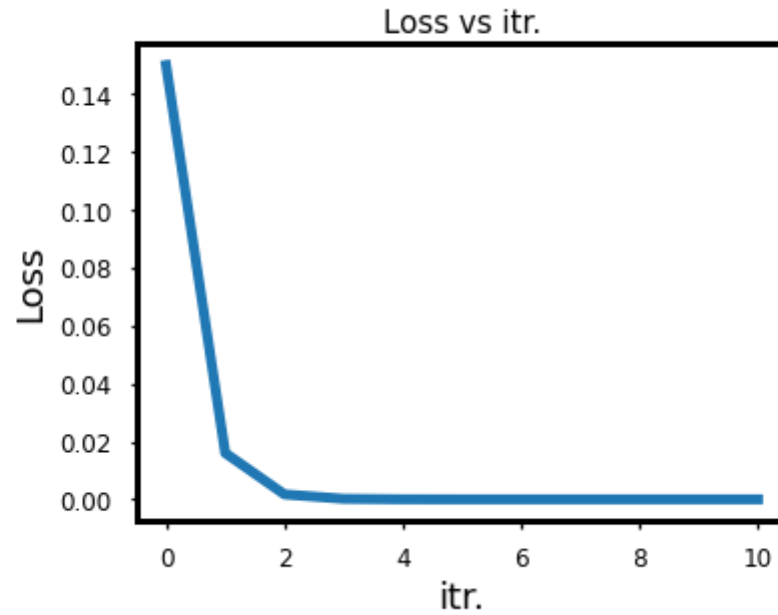
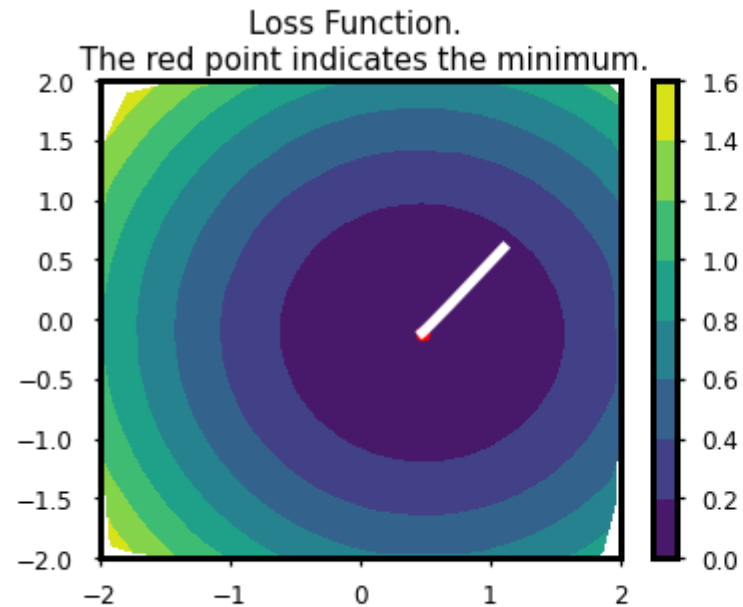
- Start with a random w .
- Calculate the derivative $\frac{\partial \mathcal{L}}{\partial w}$.
- Update $w \rightarrow w - \eta \frac{\partial \mathcal{L}}{\partial w}$.
- Repeat

η : Learning rate

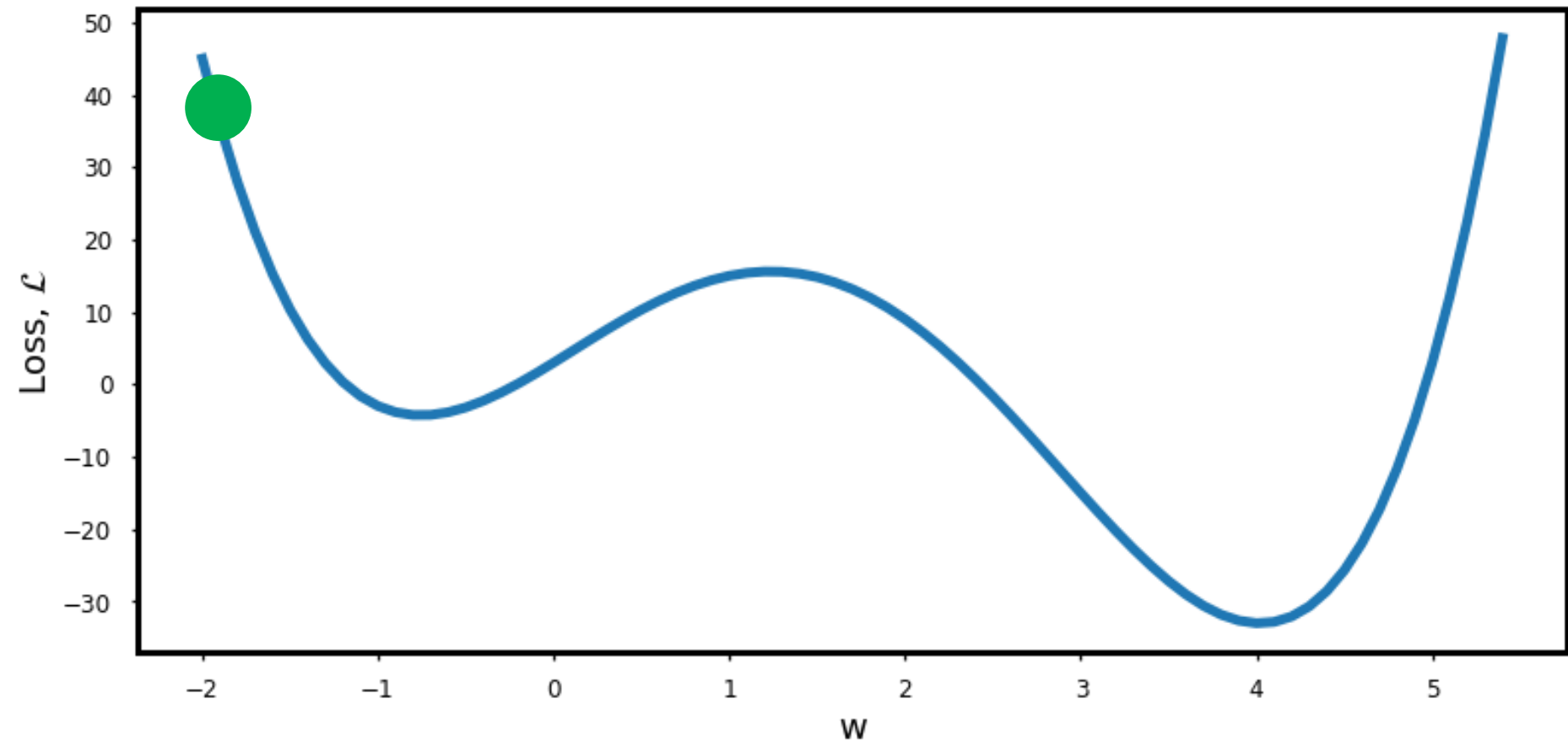
How does it work?



How does it work?

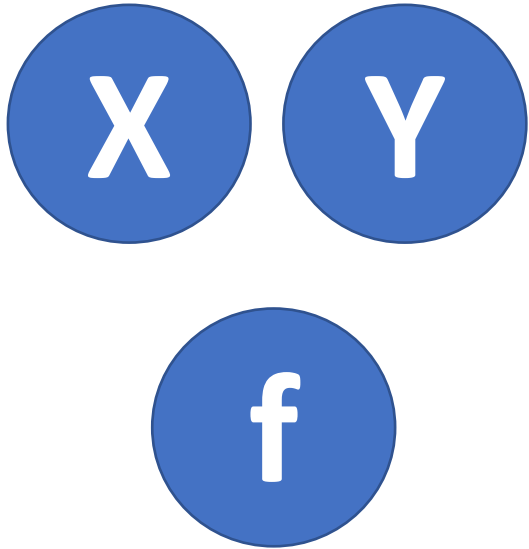


What if there is a local minimum?



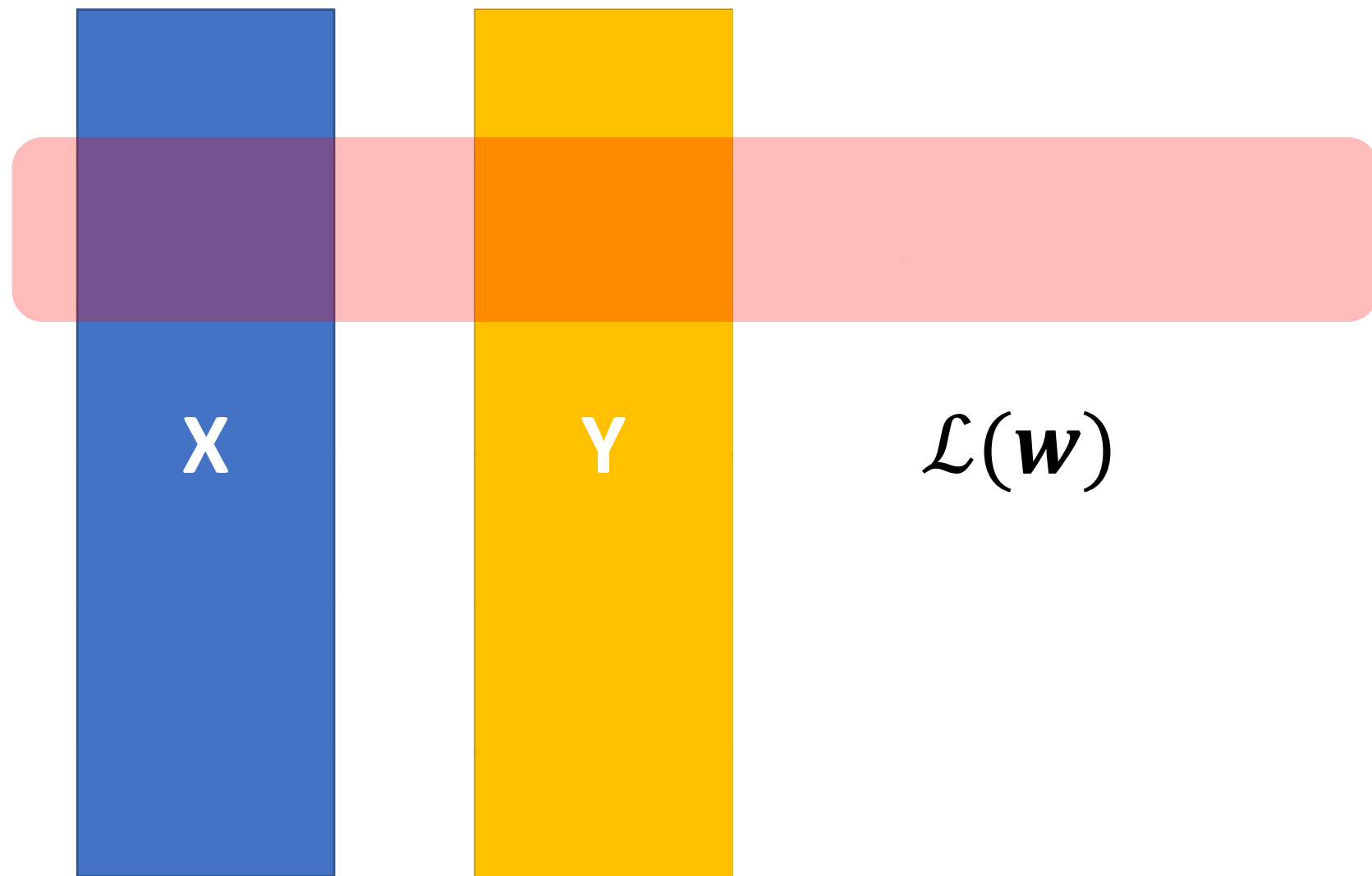
Batch Optimization

Reminder



$$\mathcal{L}(\mathbf{w})$$

Batch Optimization

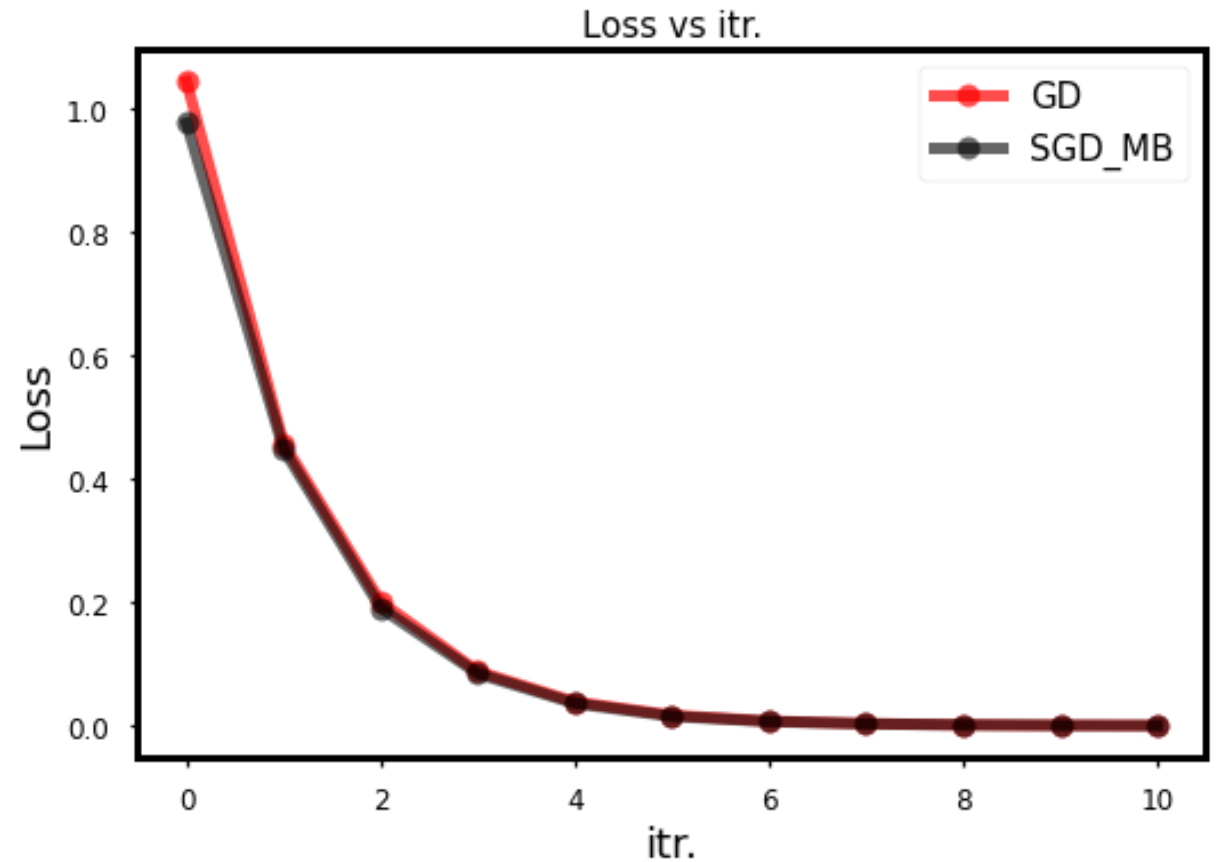
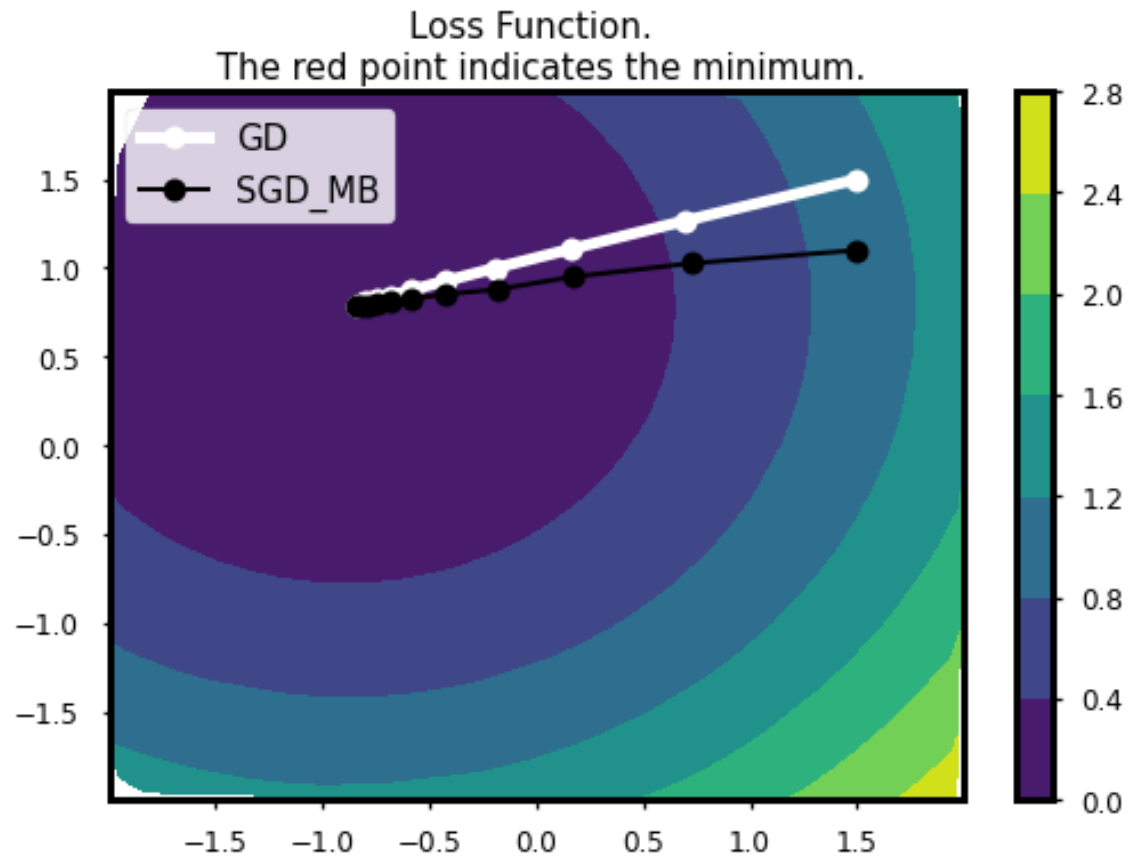


Hyper-parameter

Batch-size

Learning rate

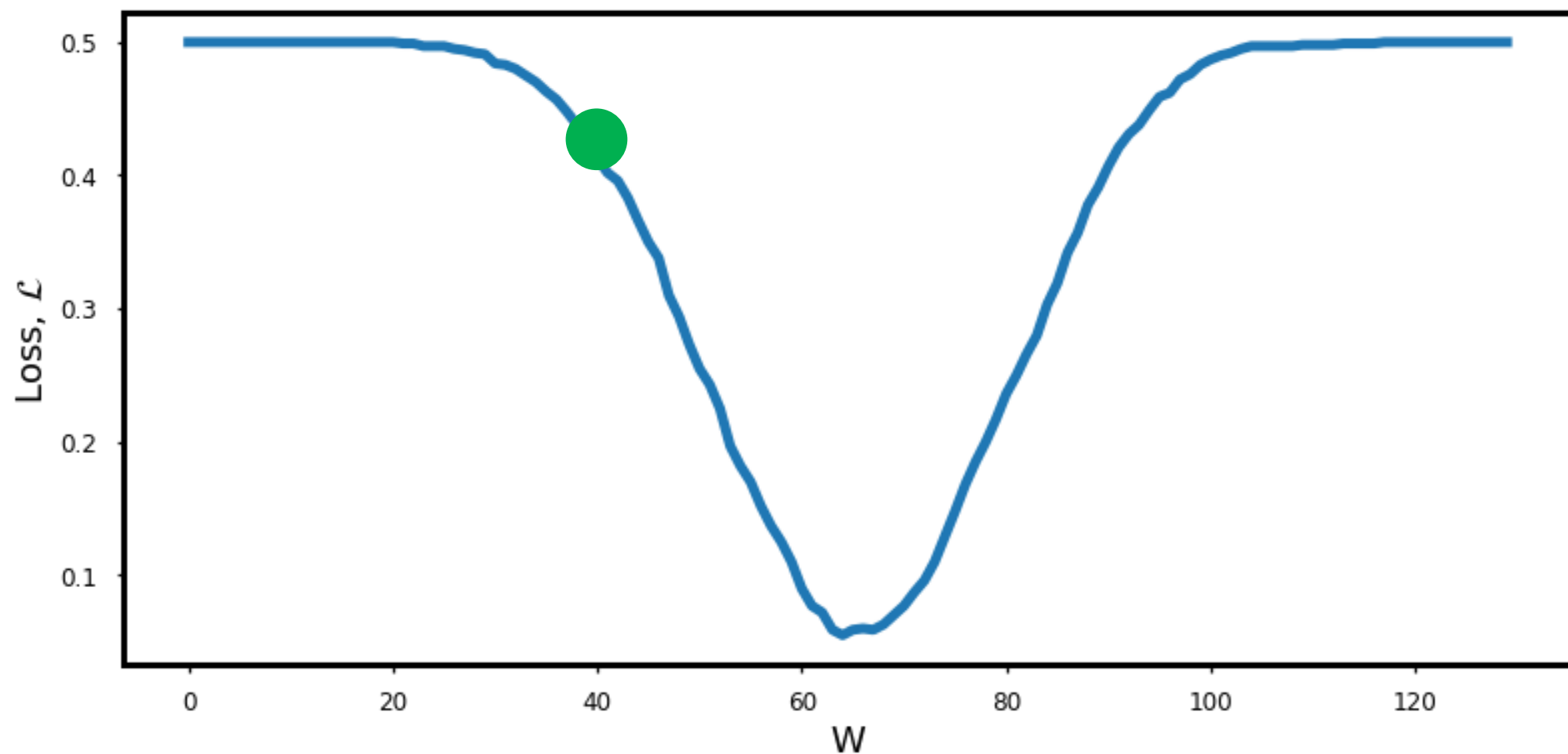
How does it look like?



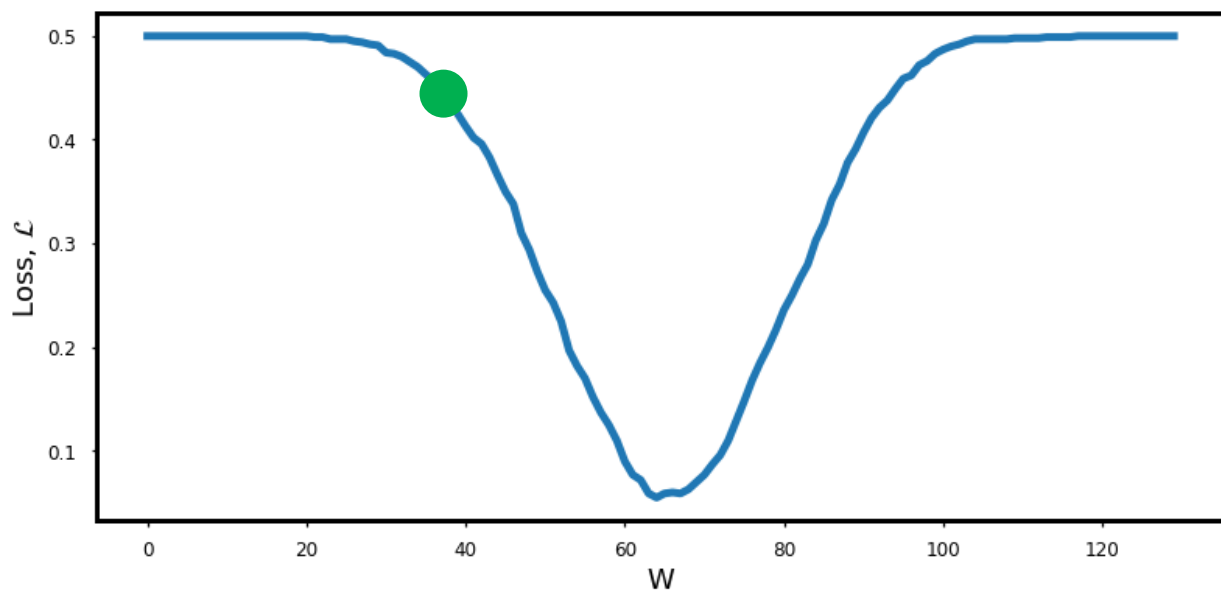
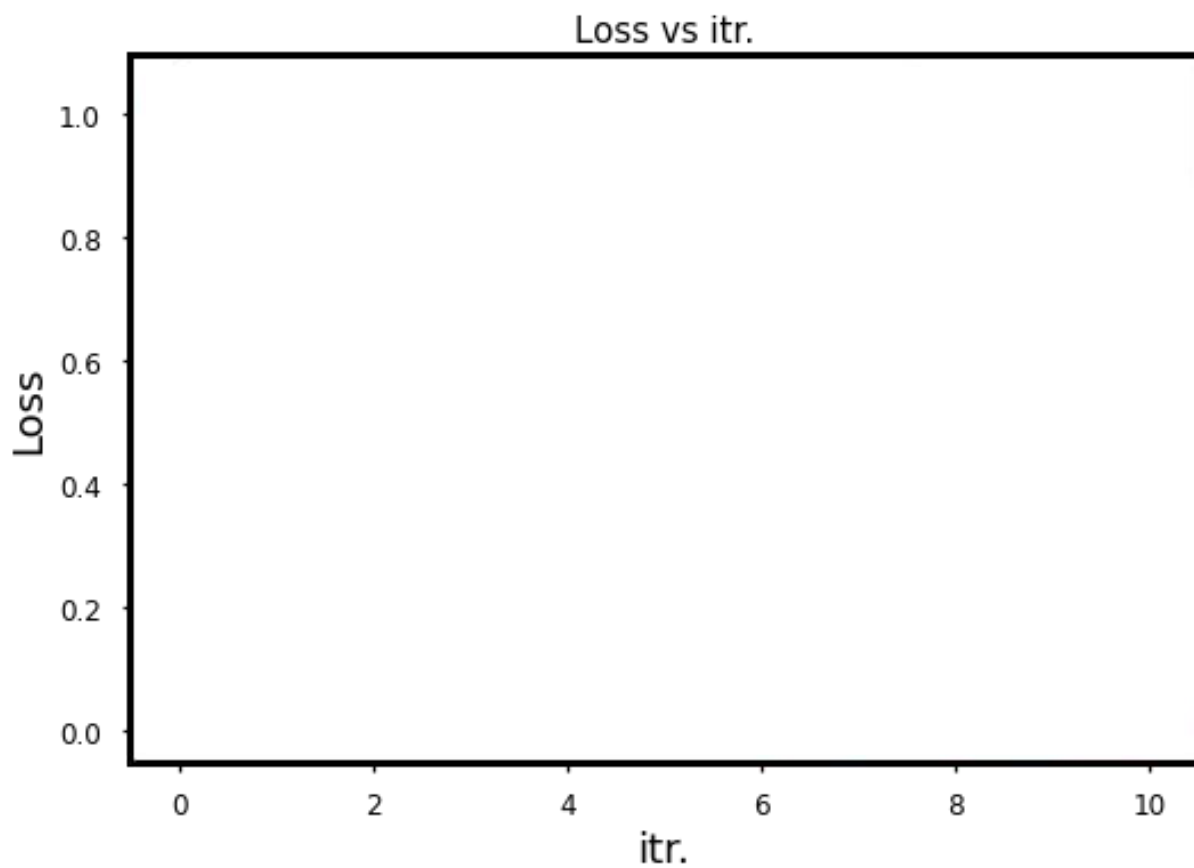
How does batch optimization help?

- Faster
- Can escape the local minimum

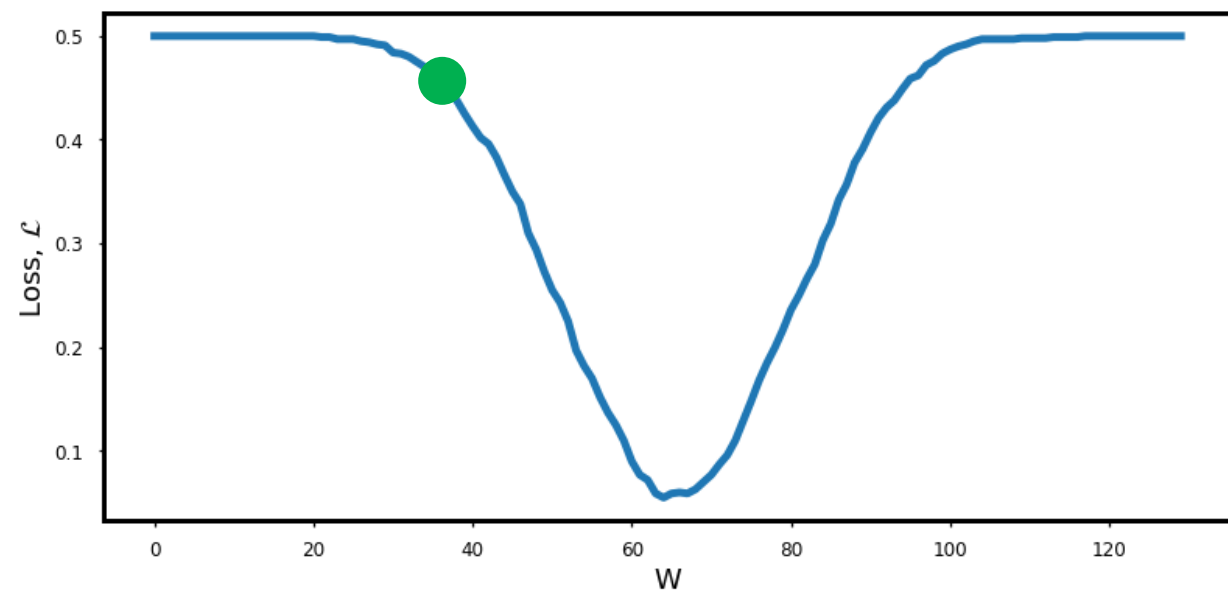
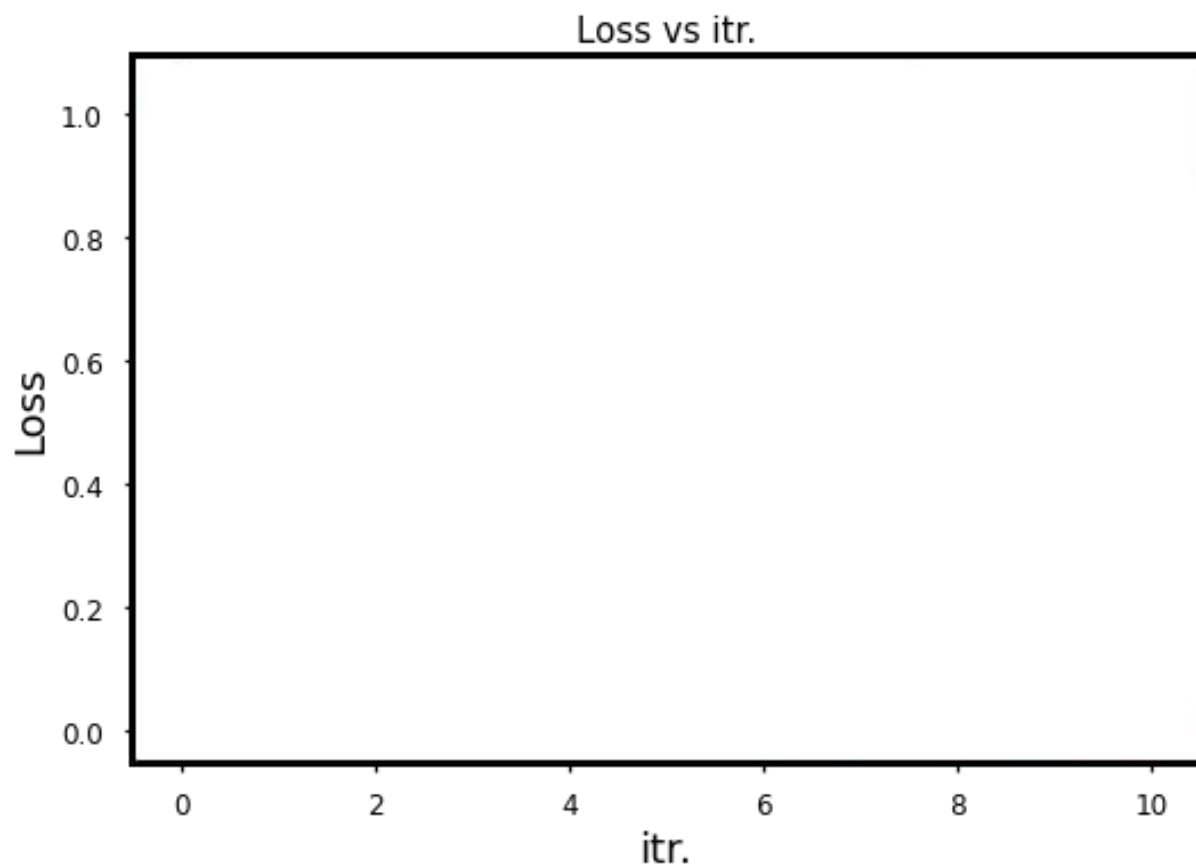
What happens if l_r is small?



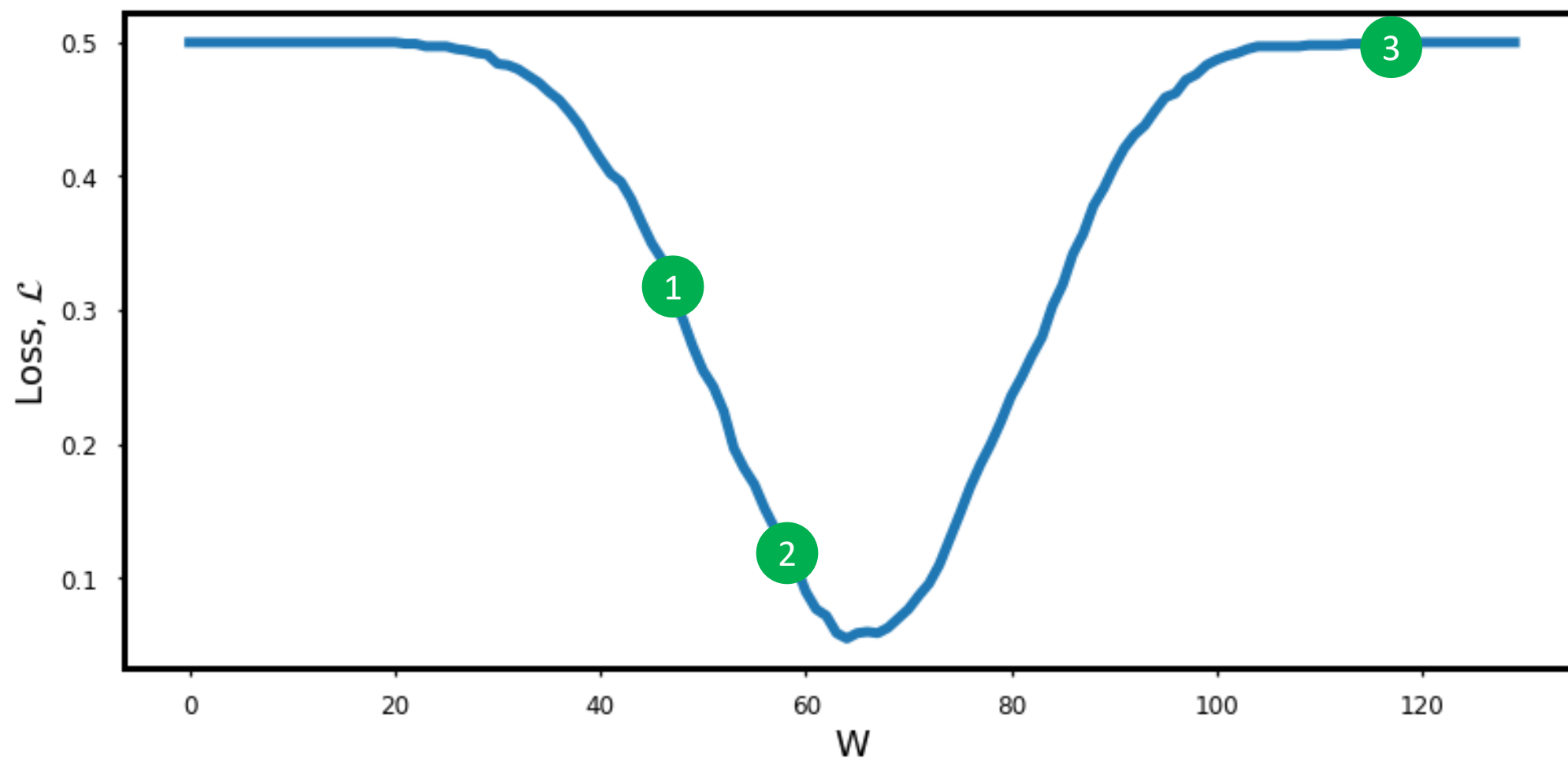
What happens if l_r is big?



What happens if l_r is small?



What is a good choice for the lr ?



Variable
learning rate

Decaying learning rate

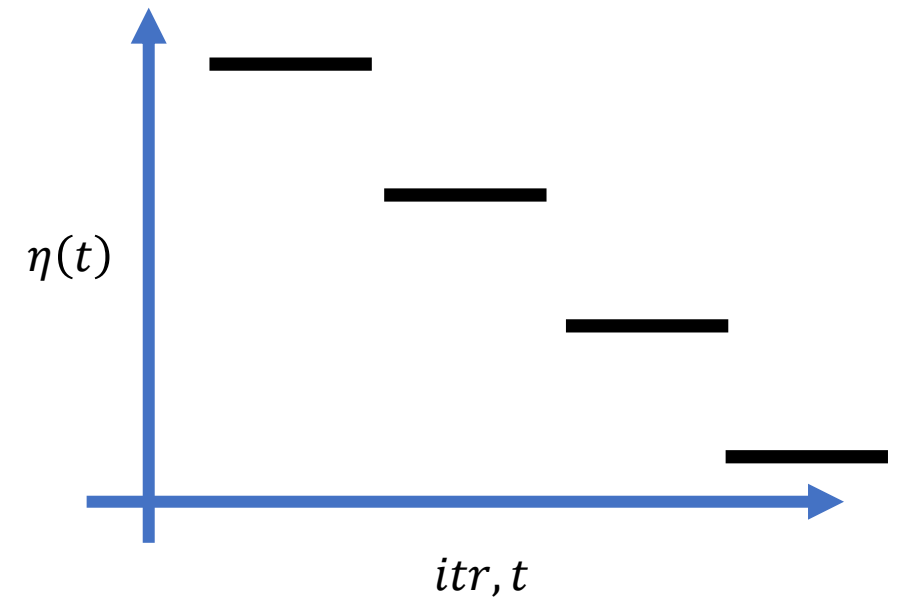
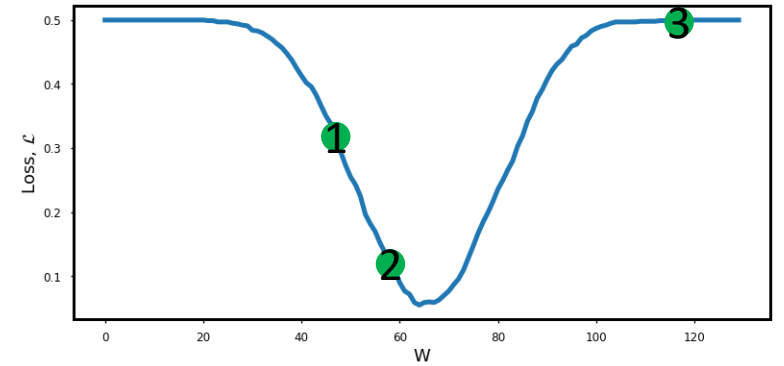
- $\eta(t) = \frac{1}{\sqrt{t}} \eta_0$

- $\eta(t) = e^{-kt} \eta_0$

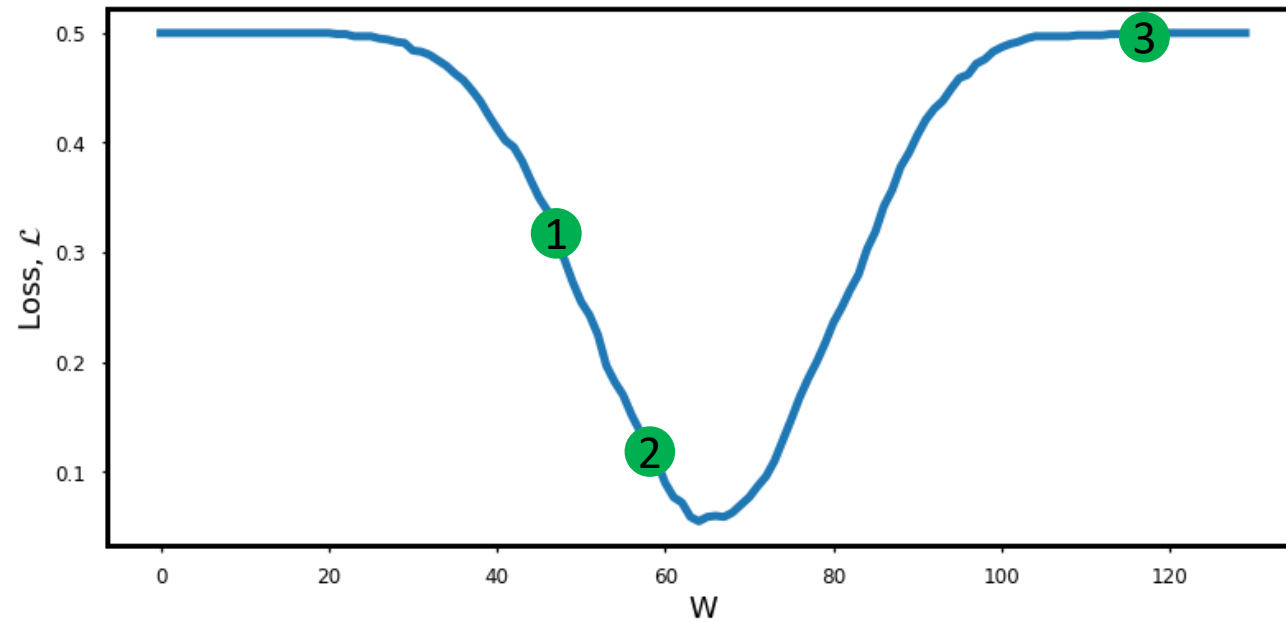
t is the number of itr.

η_0 & k are a hyper-parameter.

- Staircase

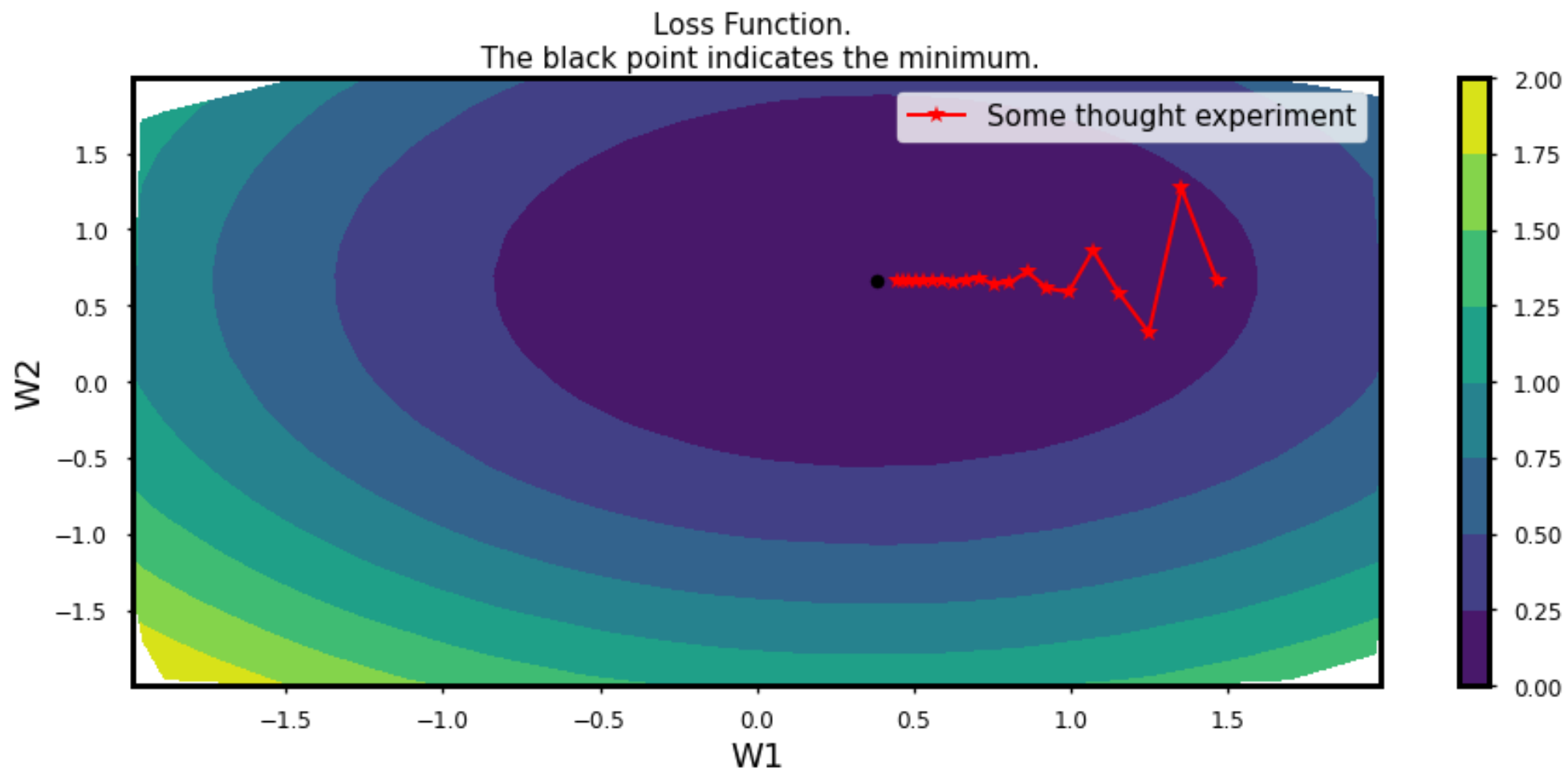


Can we do better?

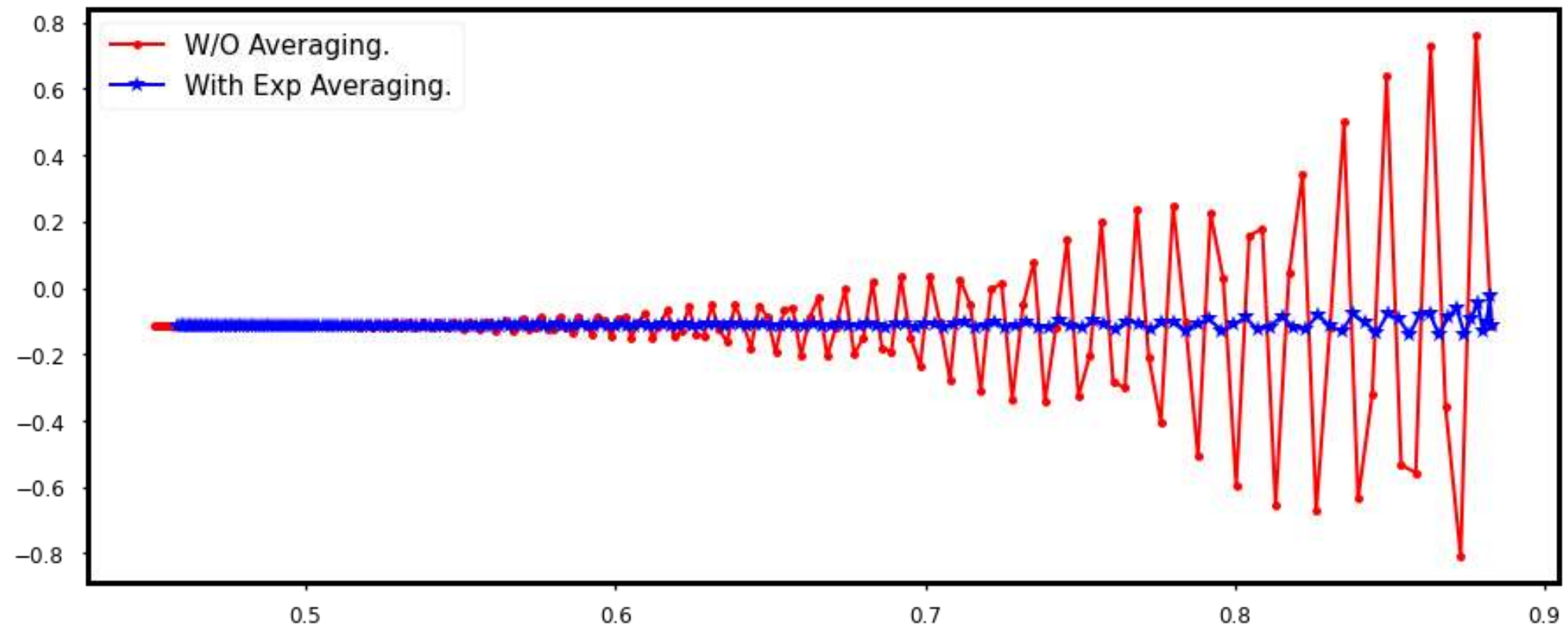


Adaptive learning rate

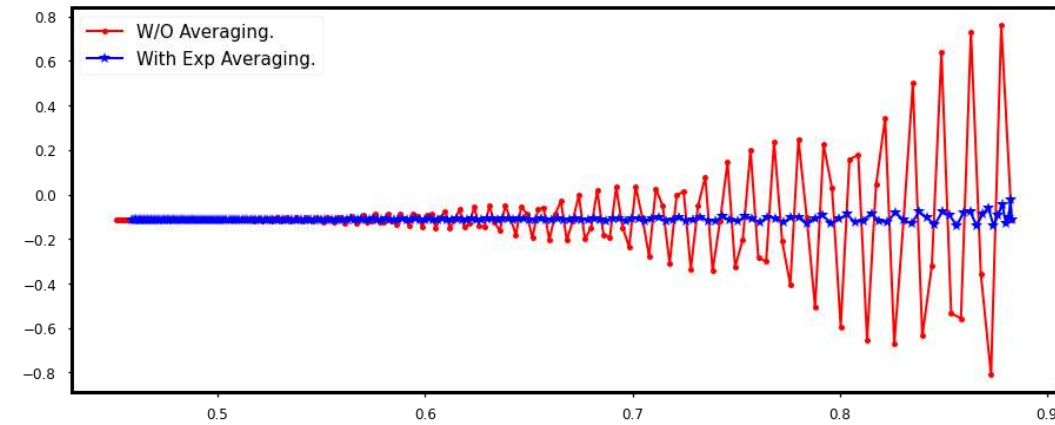
Imagine ...



Momentum



Momentum



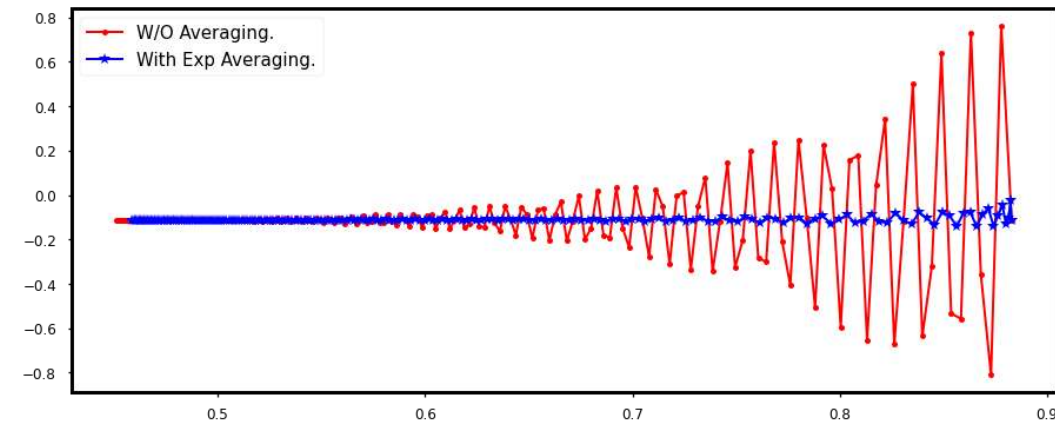
Update rule:

$$w \rightarrow w - \eta \frac{\partial \mathcal{L}}{\partial w}$$



$$w \rightarrow w - \eta \left(\frac{\partial \mathcal{L}}{\partial w} \right)_{avg}$$

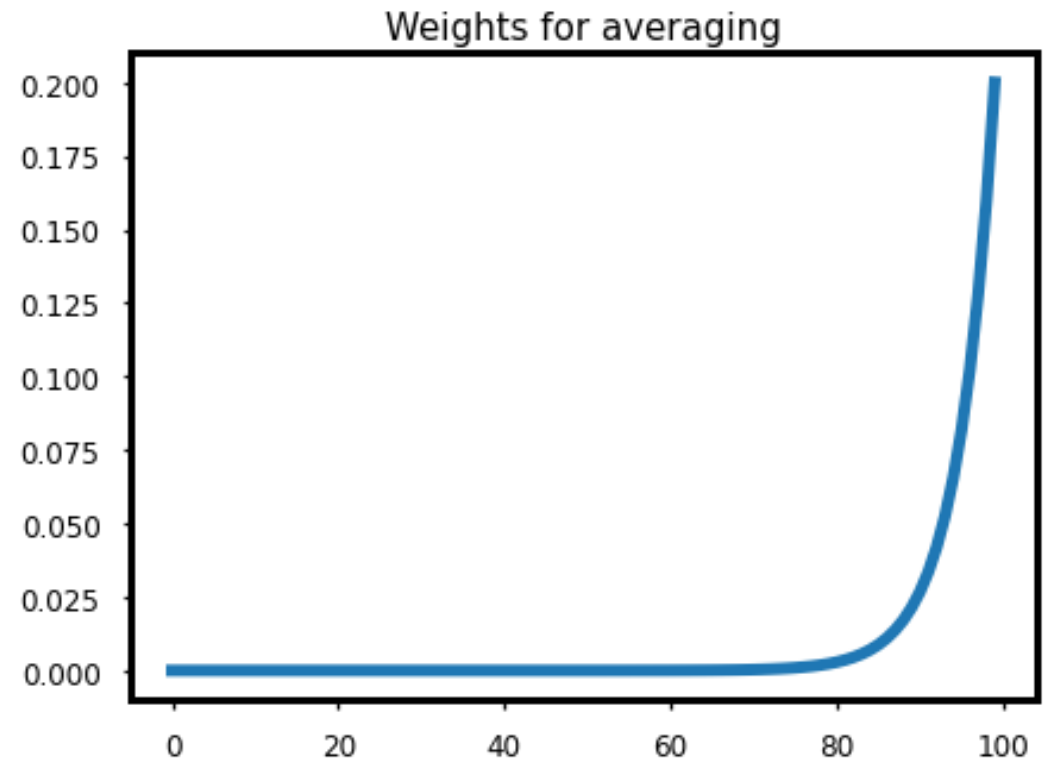
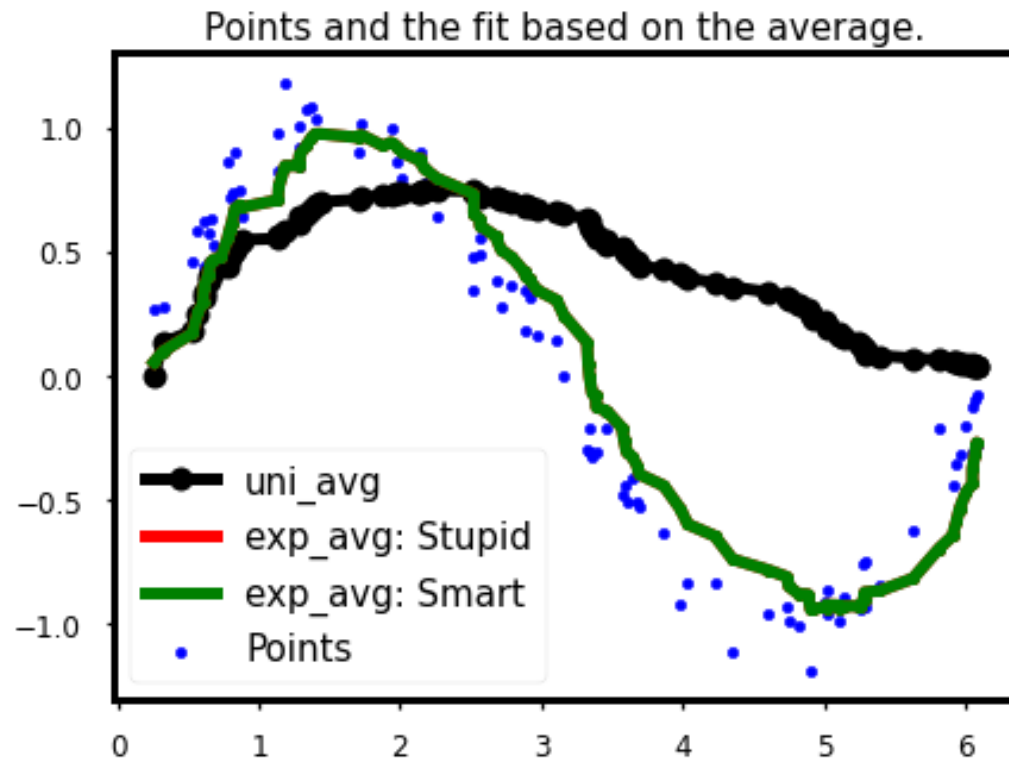
What kind of average?



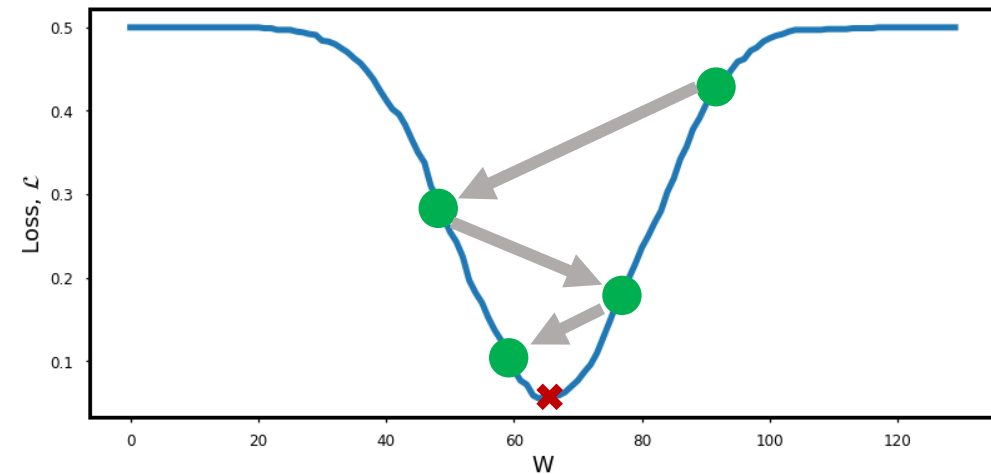
$$\left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}}\right)_{avg}^{t+1} = \beta \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}}\right)_{avg}^t + (1 - \beta) \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}}\right)^{t+1}$$

Exponentially weighted average

$$\left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}}\right)_{avg}^{t+1} = \beta \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}}\right)_{avg}^t + (1 - \beta) \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}}\right)^{t+1}$$



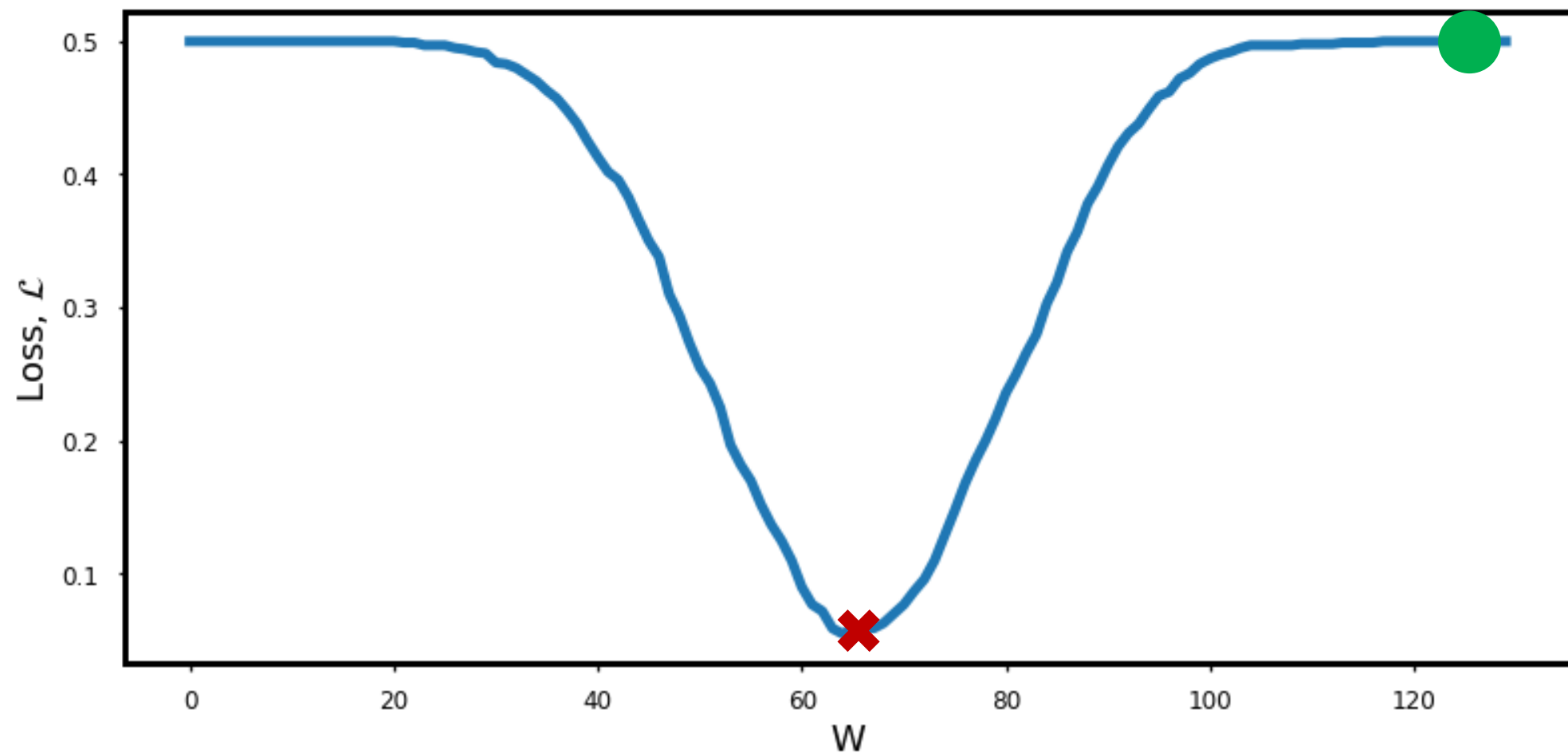
Momentum



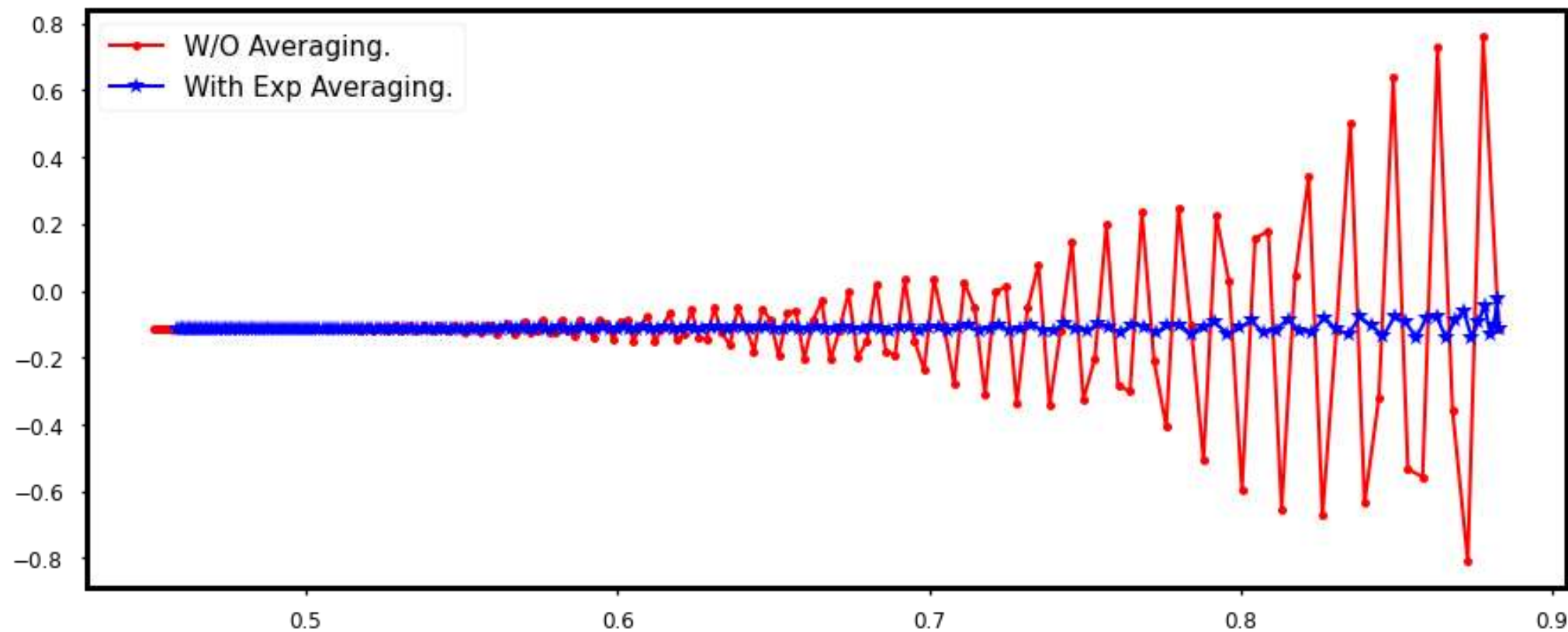
$$w \rightarrow w - \eta \left(\frac{\partial \mathcal{L}}{\partial w} \right)_{avg}$$

$$\left(\frac{\partial \mathcal{L}}{\partial w} \right)_{avg}^{t+1} = \beta \left(\frac{\partial \mathcal{L}}{\partial w} \right)_{avg}^t + (1 - \beta) \left(\frac{\partial \mathcal{L}}{\partial w} \right)^{t+1}$$

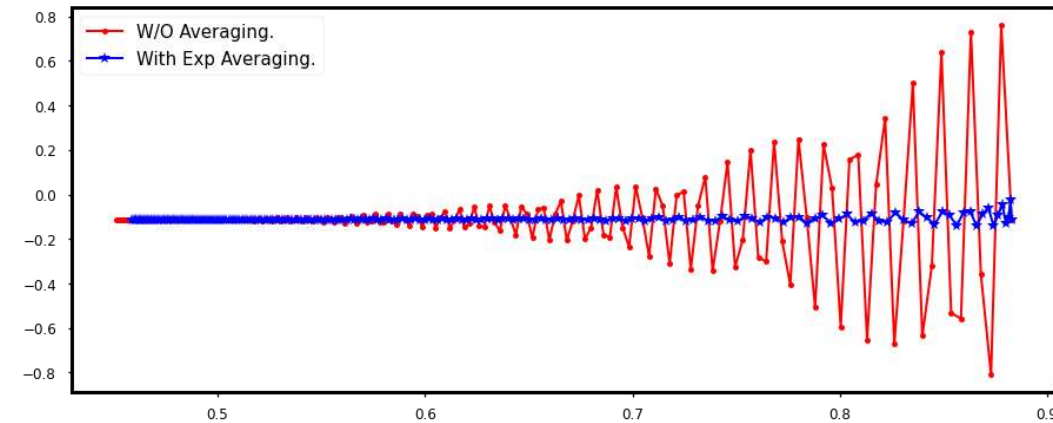
How could make it increase the l_r ?



Variance



Variance



Update rule:

$$w \rightarrow w - \eta \frac{\partial \mathcal{L}}{\partial w}$$



$$w \rightarrow w - \frac{\eta}{\sqrt{ms}} \left(\frac{\partial \mathcal{L}}{\partial w} \right)$$

ms: **Mean of squared of the gradients.**

Average of variance

$$(\text{ms})_{avg}^{t+1} = \beta (\text{ms})_{avg}^t + (1 - \beta) \left(\left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right)^2 \right)^{t+1}$$

RMS prop

$$\mathbf{w} \rightarrow \mathbf{w} - \frac{\eta}{\sqrt{\text{ms} + \epsilon}} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right)$$

$$(\text{ms})_{avg}^{t+1} = \beta' (\text{ms})_{avg}^t + (1 - \beta') \left(\left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right)^2 \right)^{t+1}$$

Now let's put both of them together ...

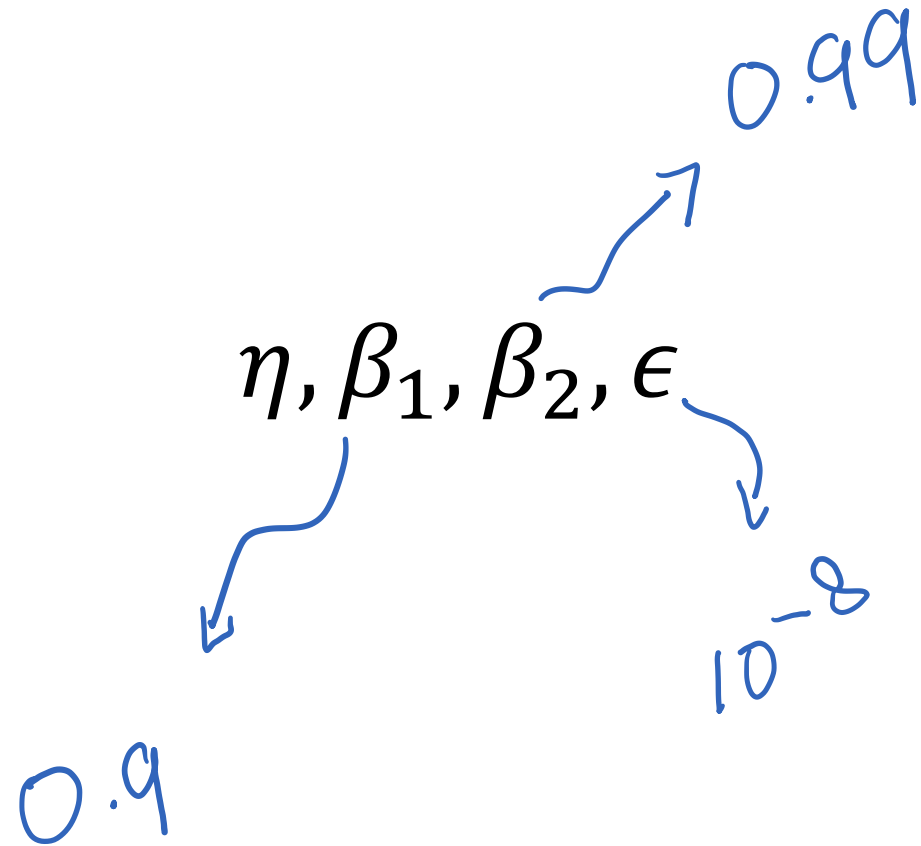
Adam

$$\mathbf{w} \rightarrow \mathbf{w} - \frac{\eta}{\sqrt{m\mathbf{s} + \epsilon}} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right)_{avg}$$

$$\left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right)_{avg} = \beta_1 \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right)_{avg}^t + (1 - \beta_1) \left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right)^{t+1}$$

$$m\mathbf{s} = \beta_2 (m\mathbf{s})_{avg}^t + (1 - \beta_2) \left(\left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}} \right)^2 \right)^{t+1}$$

Hyper parameters



So far ...

