

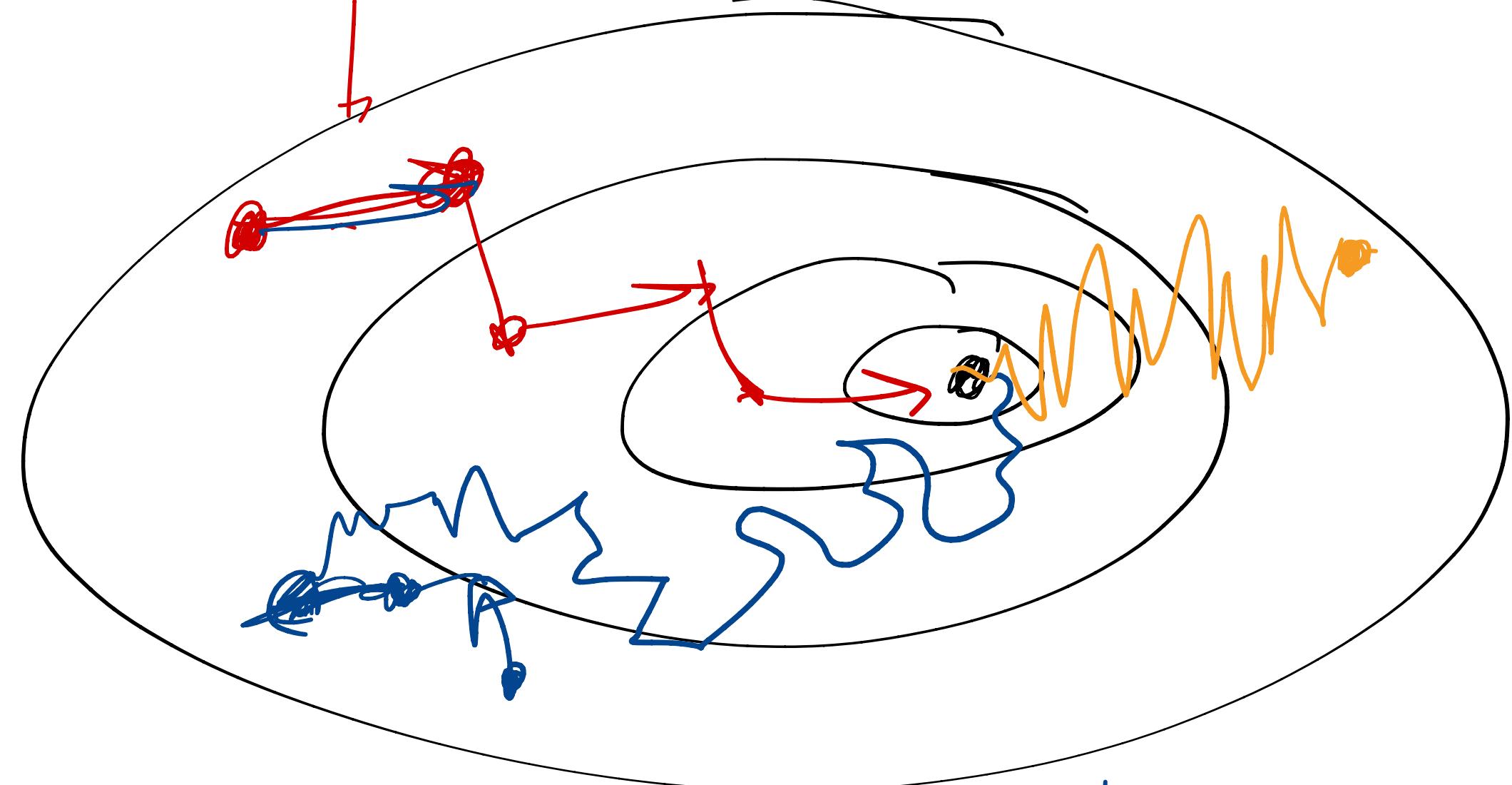
ICCS482 Deep Learning

Lecture 6: Optimization

Sunsern Cheamanunkul, Sep 24, 2020.

Gradient Descent

1-step of GD



X_B the entire data set

divide train data into
 $\{x^{(1)}\}, \{x^{(2)}\}, \dots, \{x^{(k)}\}$

doesn't make use of vectorization

- ① Stochastic GD
 - * update weights every example (noisy)
 - * should not pair with large α

- ② Batch GD
 - * update weights after seeing the entire train data.
 - * If data is not too big, this is ok.

- ③ Mini-batch GD
 - * small batches:
 - * update weight even batch

batch-size \rightarrow # of examples in mini-batch

of epoch \leftarrow # of passes on train data

$\hookrightarrow \frac{\# \text{ of examples}}{\text{batch-size}}$ weight updates

batch-size = 1 \leftarrow stochastic GD

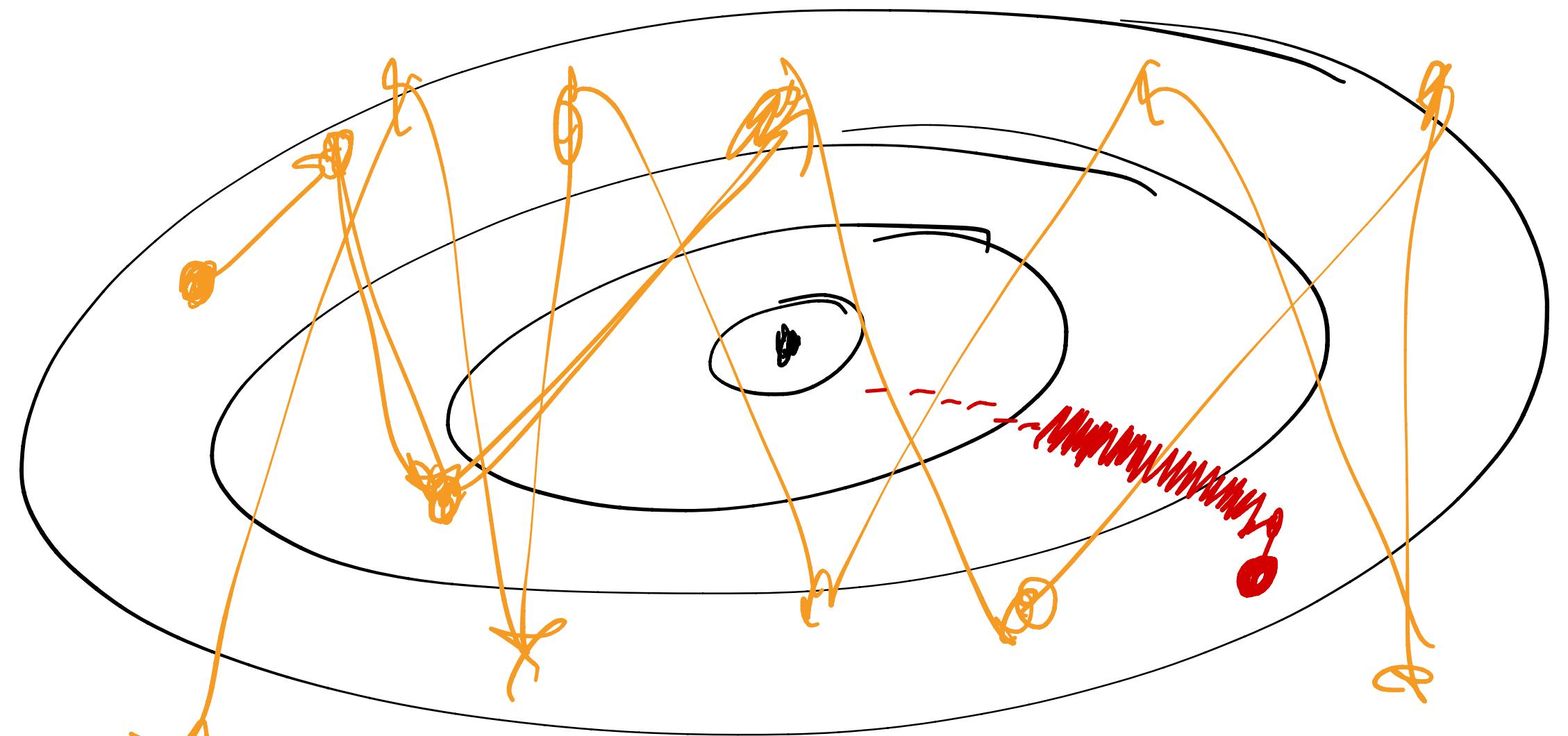
batch-size = M \leftarrow batch GD

2 < batch-size $\leq M-1$ \leftarrow minibatch GD

In practice,

$$\text{batch-size} = \left\{ \begin{array}{l} 32, 64, 128, 256, 512 \\ 1024 \end{array} \right\}$$

Problem with GD



$$\underline{w} = \underline{w} - \alpha \nabla w$$

learning rate

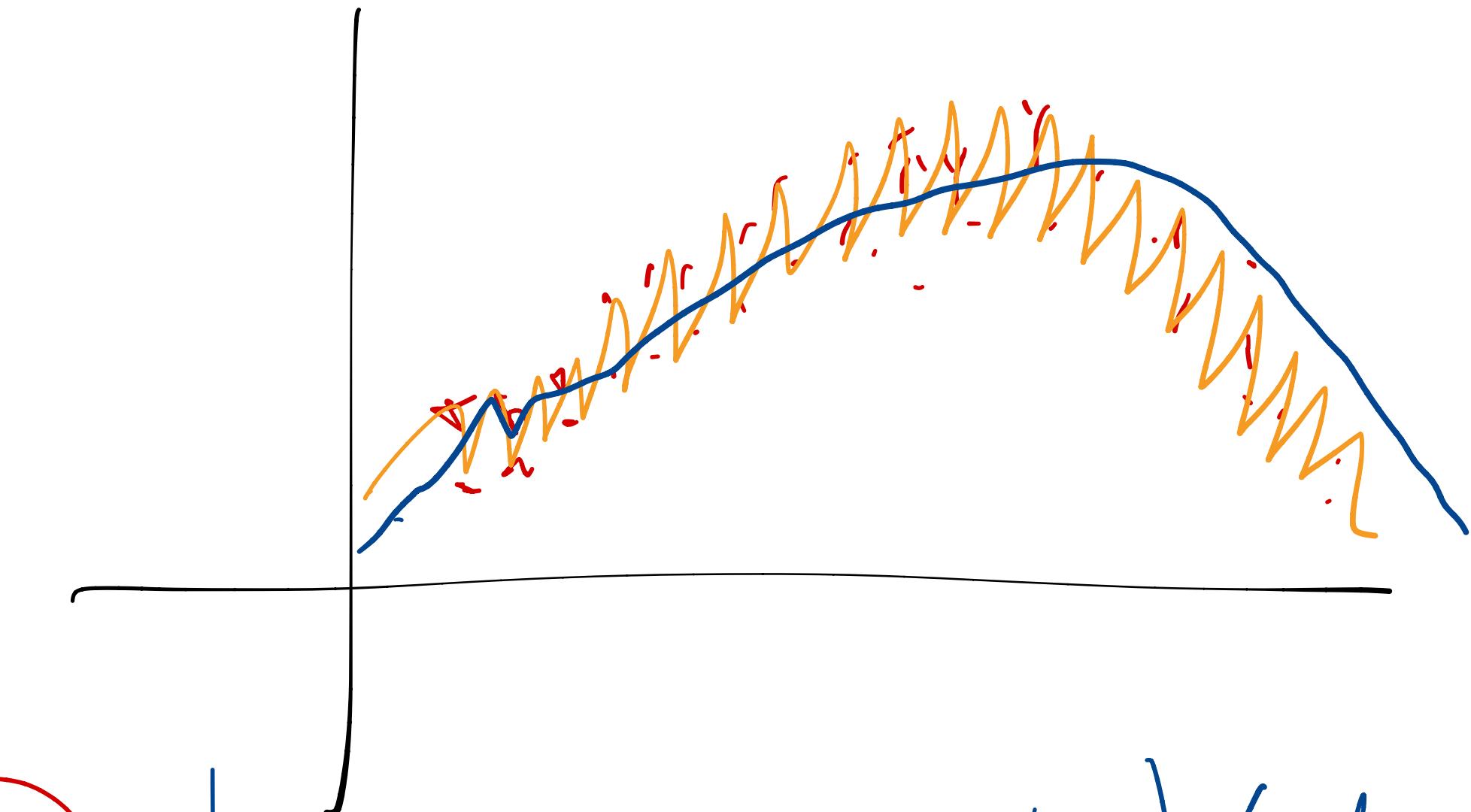
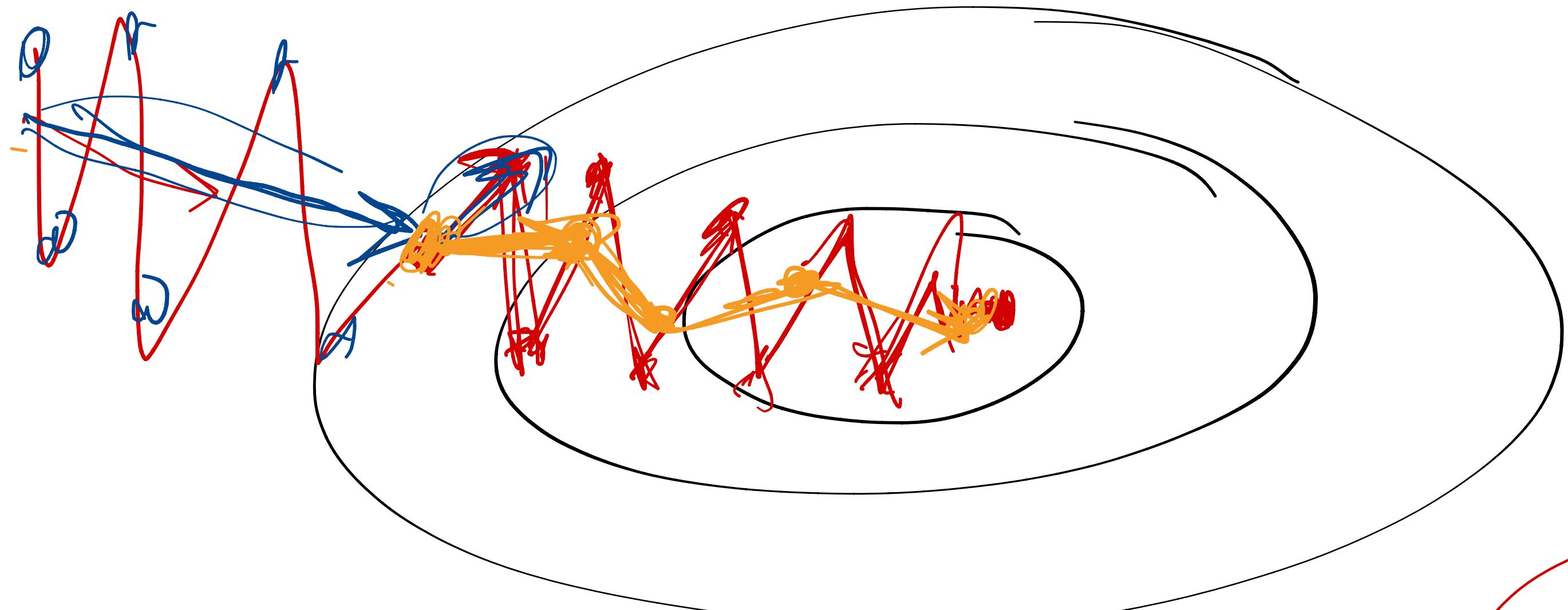
② too low

learning rate

problem of
not converge

① too high
learning rate

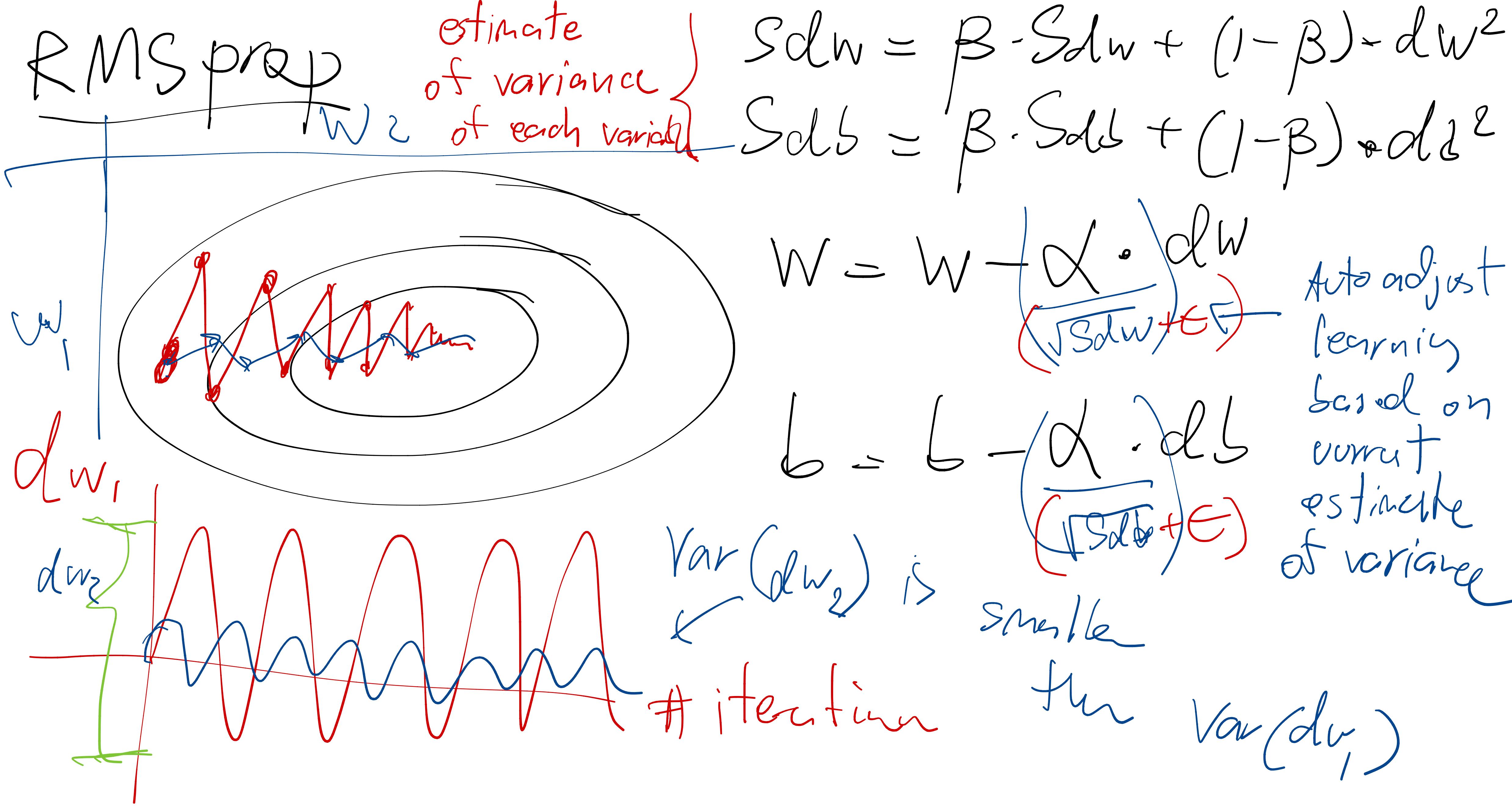
Gradient descent with Momentum



$$\begin{aligned} \overrightarrow{Vdw} &= \beta \overrightarrow{Vdw_t} + (1 - \beta) \overrightarrow{dw} \\ \overrightarrow{Vdb} &= \beta \overrightarrow{Vdb_t} + (1 - \beta) \overrightarrow{db} \end{aligned}$$

exponential weighted average of \overrightarrow{dw} over iterations.

$$\left. \begin{aligned} w &= w - \alpha \cdot \overrightarrow{Vdw} \\ b &= b - \alpha \cdot \overrightarrow{Vdb} \end{aligned} \right\}$$



Adam (Adaptive Moment Estimate)

$$Vdw = \beta_1 Vdw + (1 - \beta_1) dw$$

$$Vdb = \beta_1 Vdb + (1 - \beta_1) dl$$

$$Sdw = \beta_2 Sdw + (1 - \beta_2) dw^2$$

$$Sdb = \beta_2 Sdb + (1 - \beta_2) dl^2$$

$$w = w - \frac{\alpha \cdot Vdw}{Sdw + \epsilon}$$

1st moment

2nd moment

$$b = b - \frac{\alpha \cdot Vdb}{Sdb + \epsilon}$$

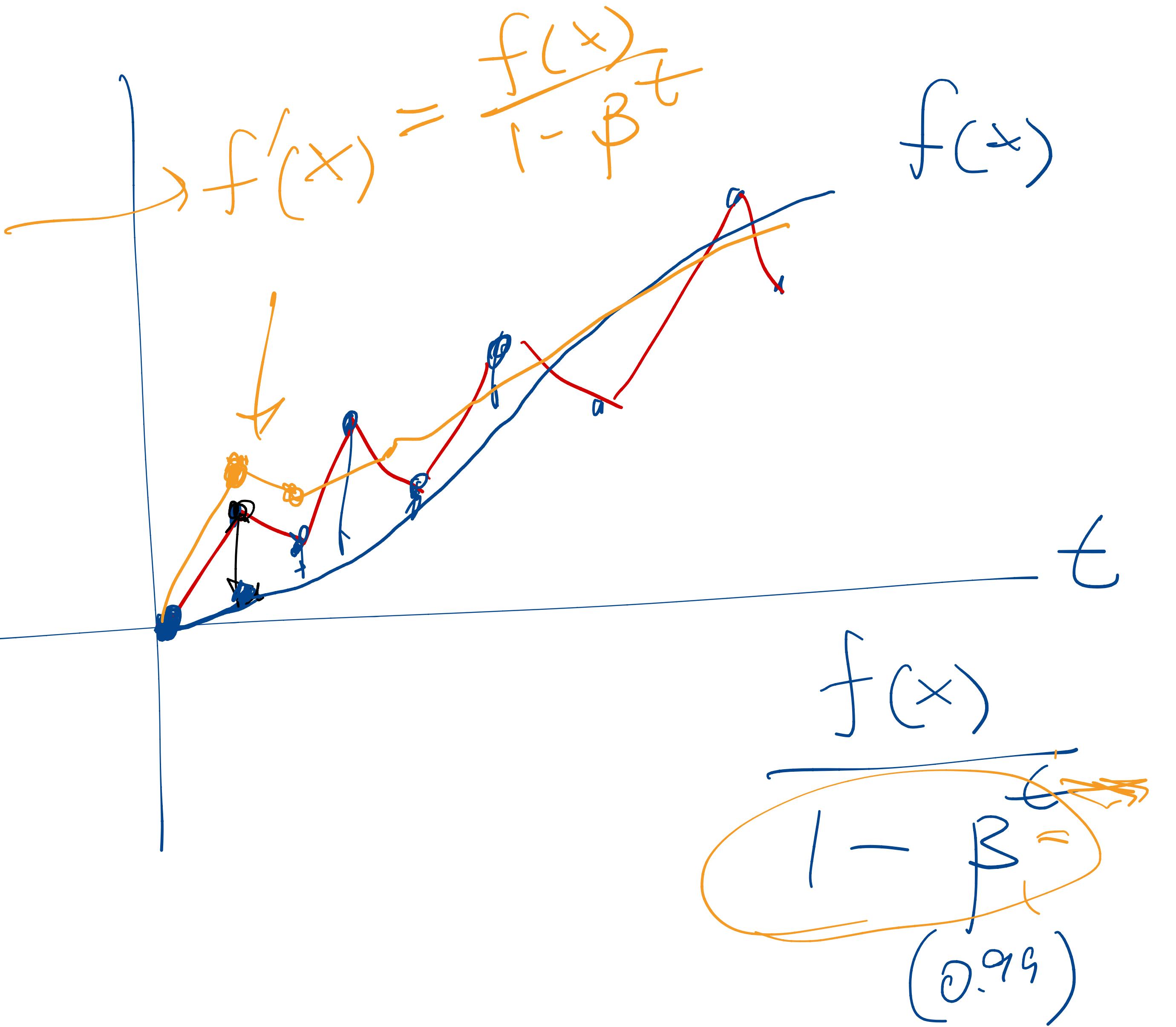
Bias correction

term



this corrects
odd
start
problem

after
applying/
bias
correction



Learning rate decay

$$\alpha \rightarrow \text{epoch} \cdot \alpha_0 \cdot \text{decay-rate}$$

initial learning rate

decay-rate

epoch

α_0

