

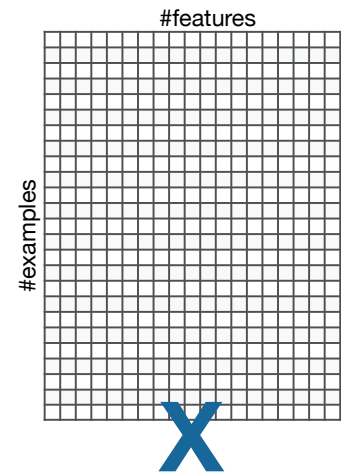
COS234: INTRODUCTION TO MACHINE LEARNING

Prof. Yoram Singer



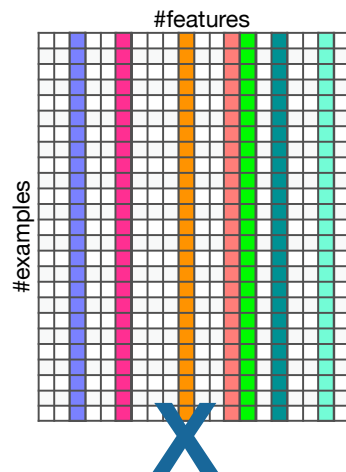
Topic: Gradient-Based Learning - Part I

© 2020 YORAM SINGER



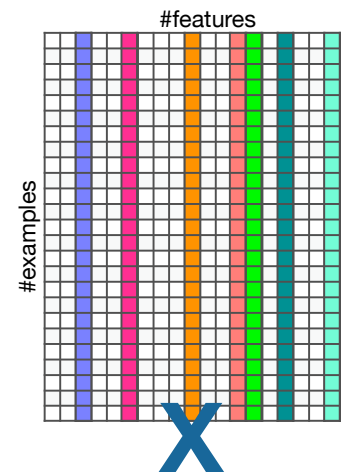
© 2020 YORAM SINGER

2



© 2020 YORAM SINGER

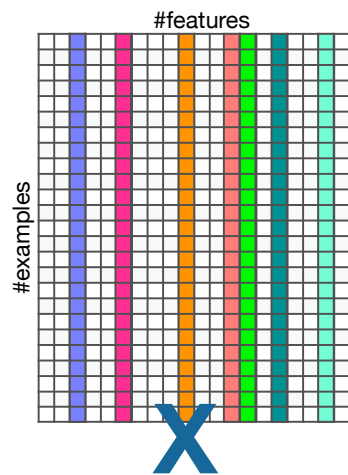
2



© 2020 YORAM SINGER

2

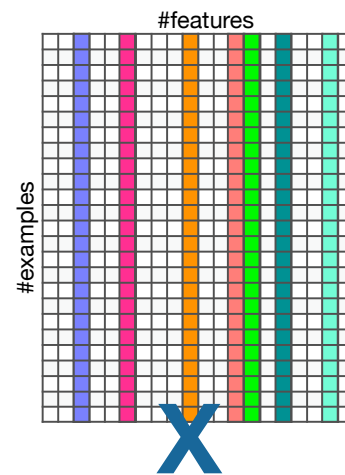
- Sequential: a column at a time



- Sequential: a column at a time
- Oblivious to “similar” examples

© 2020 YORAM SINGER

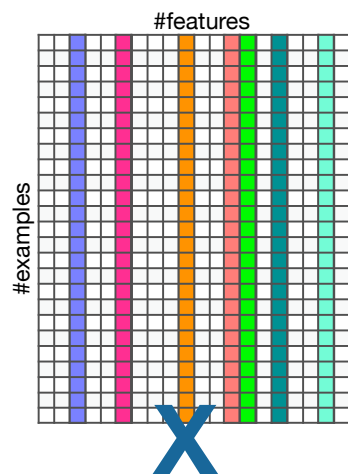
2



- Sequential: a column at a time
- Oblivious to “similar” examples
- Difficult to parallelize

© 2020 YORAM SINGER

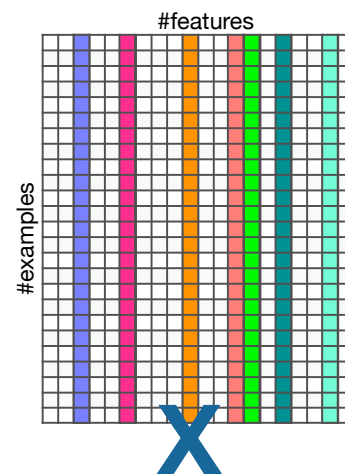
2



- Sequential: a column at a time
- Oblivious to “similar” examples
- Difficult to parallelize
- Requires dedicated update per loss

© 2020 YORAM SINGER

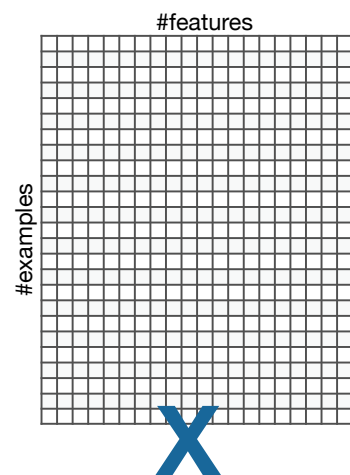
2



- Sequential: a column at a time
- Oblivious to “similar” examples
- Difficult to parallelize
- Requires dedicated update per loss
- Fails to work in non-linear settings

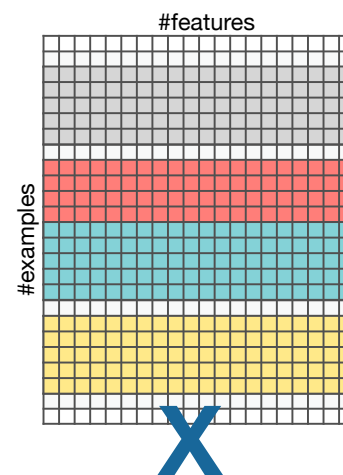
© 2020 YORAM SINGER

2



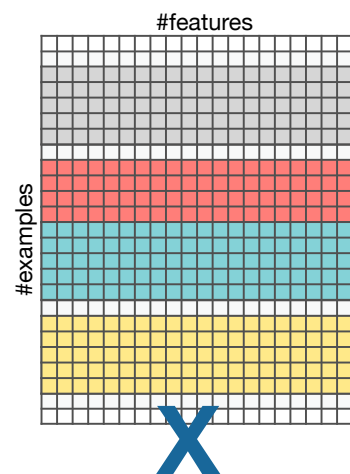
© 2020 YORAM SINGER

3



© 2020 YORAM SINGER

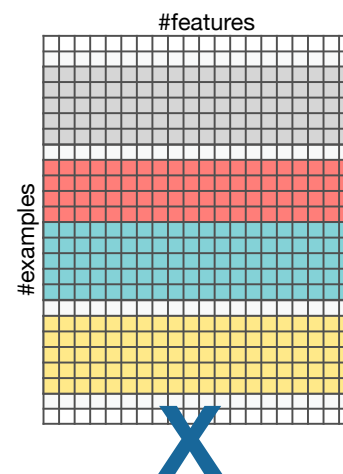
3



- Pick a small subset of rows of X to use

© 2020 YORAM SINGER

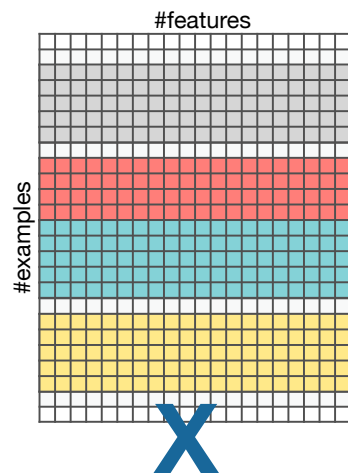
3



- Pick a small subset of rows of X to use
- Subset not necessarily consecutive rows

© 2020 YORAM SINGER

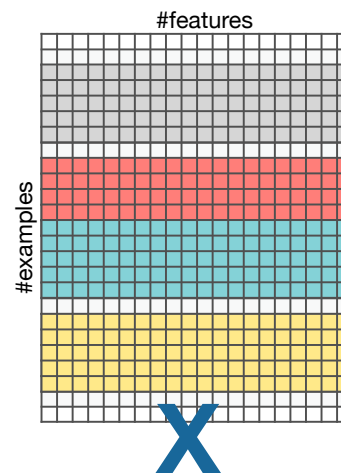
3



- Pick a small subset of rows of X to use
- Subset not necessarily consecutive rows
- Subset used to update all weights

© 2020 YORAM SINGER

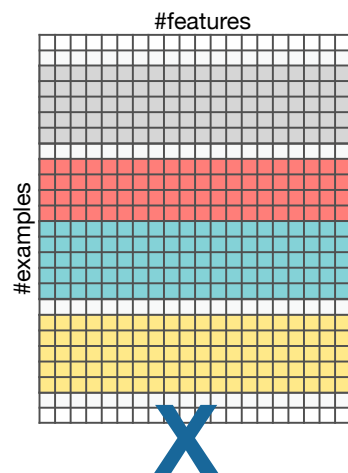
3



- Pick a small subset of rows of X to use
- Subset not necessarily consecutive rows
- Subset used to update all weights
- Update “local” to subset selected

© 2020 YORAM SINGER

3



- Pick a small subset of rows of X to use
- Subset not necessarily consecutive rows
- Subset used to update all weights
- Update “local” to subset selected
- Albeit local, effect is “global”

© 2020 YORAM SINGER

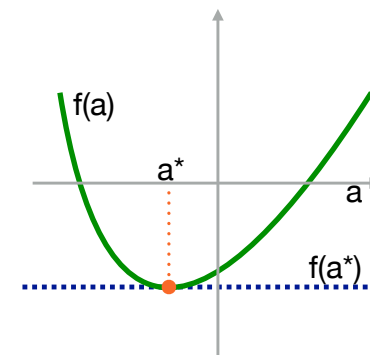
3

- a^* minimum point of f :

$$f(a^*) \leq f(b) \text{ for all } b \neq a^*$$
- Derivate of $f(a)$ at a^* is zero

$$\left. \frac{df}{da} \right|_{a=a^*} \equiv f'(a^*) = 0$$
- Often no closed form solution for:

$$\frac{df}{da} = 0$$



© 2020 YORAM SINGER

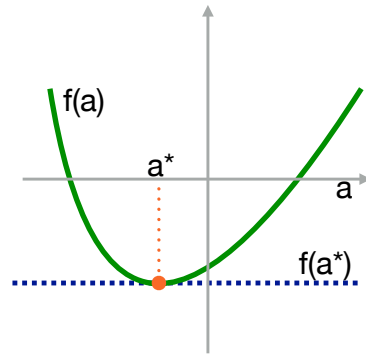
4

we can still make use of f' ...

- a^* minimum point of f :
 $f(a^*) \leq f(b)$ for all $b \neq a^*$
- Derivate of $f(a)$ at a^* is zero

$$\left. \frac{df}{da} \right|_{a=a^*} \equiv f'(a^*) = 0$$
- Often no closed form solution for:

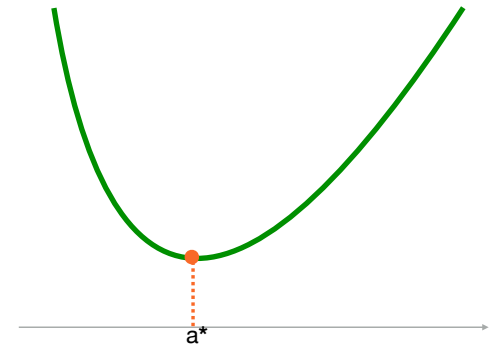
$$\frac{df}{da} = 0$$



© 2020 YORAM SINGER

4

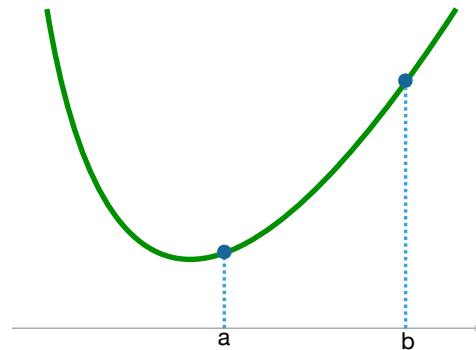
(Strict) Convexity



© 2020 YORAM SINGER

5

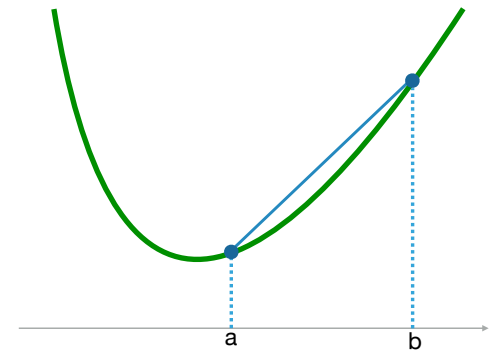
(Strict) Convexity



© 2020 YORAM SINGER

5

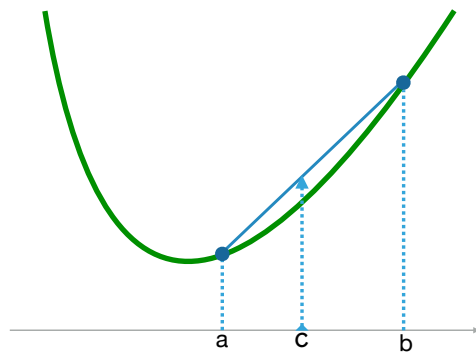
(Strict) Convexity



© 2020 YORAM SINGER

5

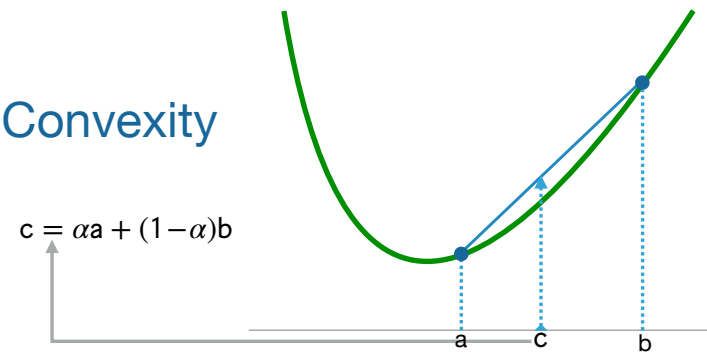
(Strict) Convexity



© 2020 YORAM SINGER

5

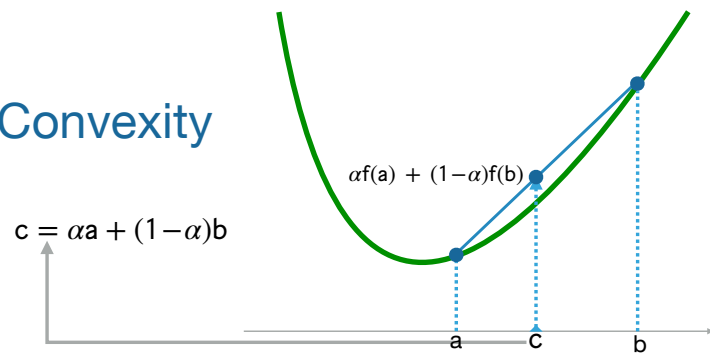
(Strict) Convexity



© 2020 YORAM SINGER

5

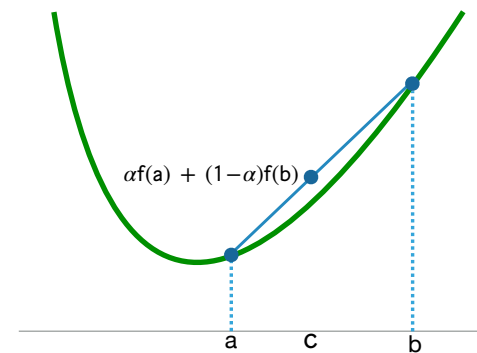
(Strict) Convexity



© 2020 YORAM SINGER

5

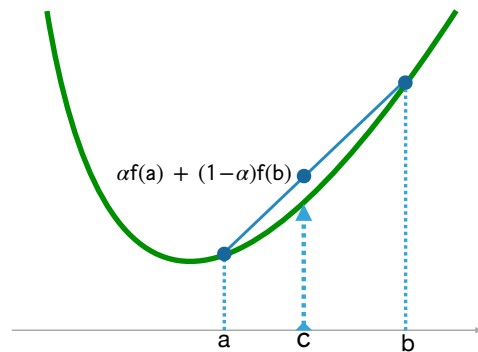
(Strict) Convexity



© 2020 YORAM SINGER

5

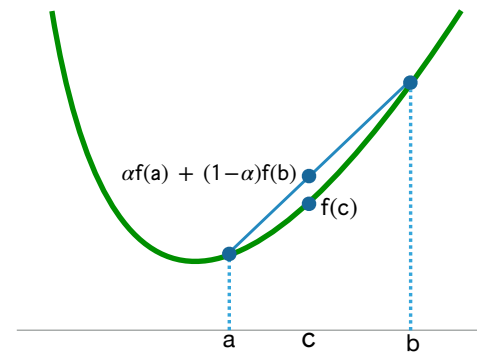
(Strict) Convexity



© 2020 YORAM SINGER

5

(Strict) Convexity

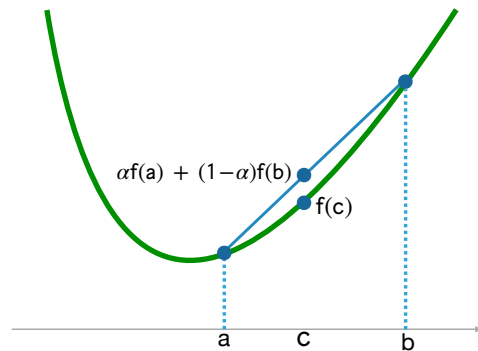


© 2020 YORAM SINGER

5

$$\alpha \in (0, 1) : f(\alpha a + (1 - \alpha)b) < \alpha f(a) + (1 - \alpha)f(b)$$

(Strict) Convexity



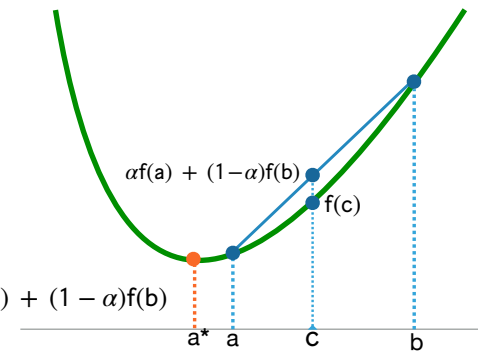
© 2020 YORAM SINGER

5

(Strict) Convexity

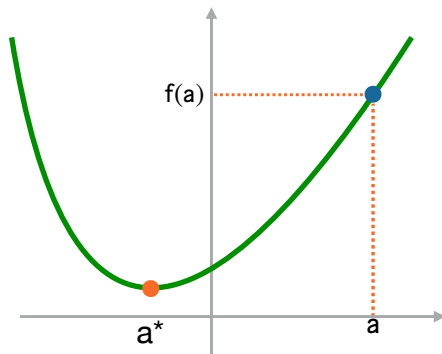
- a^* is a unique minimum
- derivative is < 0 left to a^*
- derivative is > 0 right to a^*
- derivative of $f(a)$ at a^* is 0
- second derivative $f''(a) > 0$

$$\alpha \in (0, 1) : f(\alpha a + (1 - \alpha)b) < \alpha f(a) + (1 - \alpha)f(b)$$



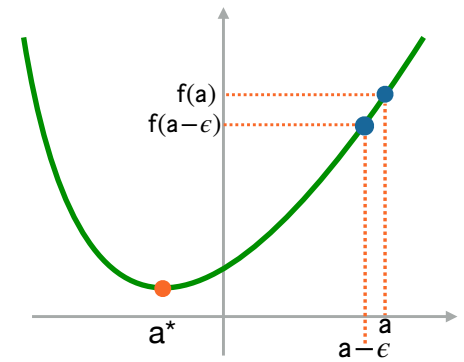
© 2020 YORAM SINGER

6



© 2020 YORAM SINGER

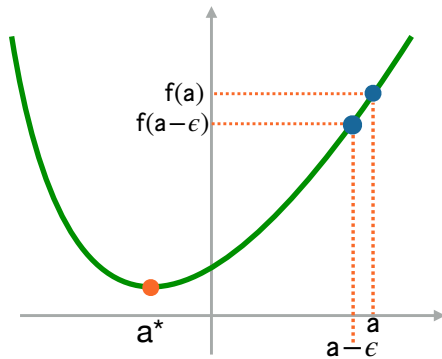
7



© 2020 YORAM SINGER

7

- if a to the right ($a > a^*$) of a^* then
 $f(a) > f(a - \epsilon)$



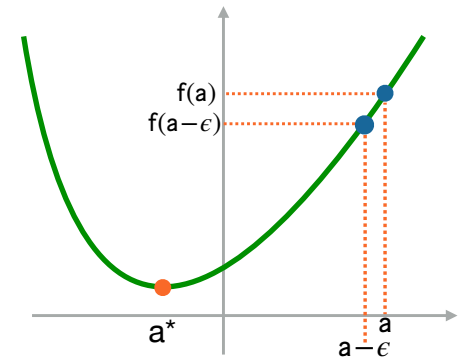
© 2020 YORAM SINGER

7

- if a to the right ($a > a^*$) of a^* then
 $f(a) > f(a - \epsilon)$

- Therefore

$$\frac{f(a) - f(a - \epsilon)}{\epsilon} > 0$$



© 2020 YORAM SINGER

7

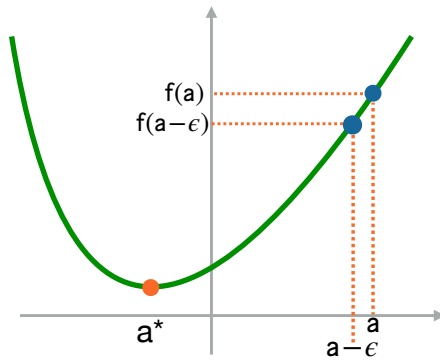
- if a to the right ($a > a^*$) of a^* then

$$f(a) > f(a - \epsilon)$$

- Therefore

$$\frac{f(a) - f(a - \epsilon)}{\epsilon} > 0$$

- Taking $\epsilon \rightarrow 0$ we get $f'(a) > 0$



© 2020 YORAM SINGER

7

- if a to the right ($a > a^*$) of a^* then

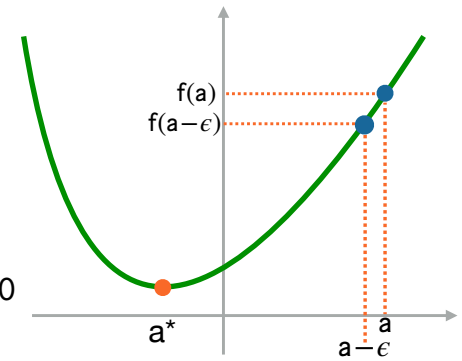
$$f(a) > f(a - \epsilon)$$

- Therefore

$$\frac{f(a) - f(a - \epsilon)}{\epsilon} > 0$$

- Taking $\epsilon \rightarrow 0$ we get $f'(a) > 0$

- Similarly for $a < a^*$ we get $f'(a) < 0$



© 2020 YORAM SINGER

7

- if a to the right ($a > a^*$) of a^* then

$$f(a) > f(a - \epsilon)$$

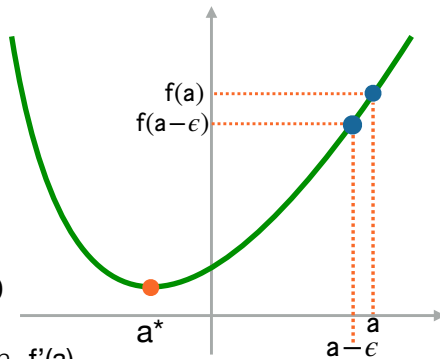
- Therefore

$$\frac{f(a) - f(a - \epsilon)}{\epsilon} > 0$$

- Taking $\epsilon \rightarrow 0$ we get $f'(a) > 0$

- Similarly for $a < a^*$ we get $f'(a) < 0$

- Get closer to a^* by going in direction $-f'(a)$



© 2020 YORAM SINGER

7

- if a to the right ($a > a^*$) of a^* then

$$f(a) > f(a - \epsilon)$$

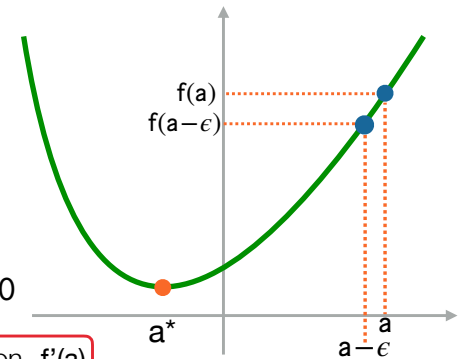
- Therefore

$$\frac{f(a) - f(a - \epsilon)}{\epsilon} > 0$$

- Taking $\epsilon \rightarrow 0$ we get $f'(a) > 0$

- Similarly for $a < a^*$ we get $f'(a) < 0$

- Get closer to a^* by going in direction $-f'(a)$



© 2020 YORAM SINGER

7

- if a to the right ($a > a^*$) of a^* then

$$f(a) > f(a - \epsilon)$$

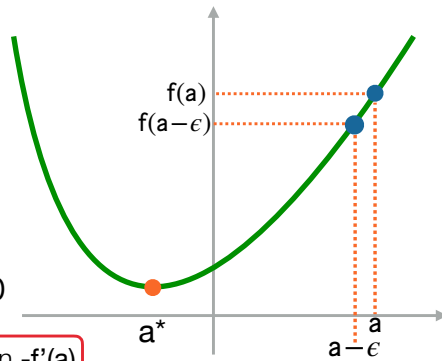
- Therefore

$$\frac{f(a) - f(a - \epsilon)}{\epsilon} > 0$$

- Taking $\epsilon \rightarrow 0$ we get $f'(a) > 0$

- Similarly for $a < a^*$ we get $f'(a) < 0$

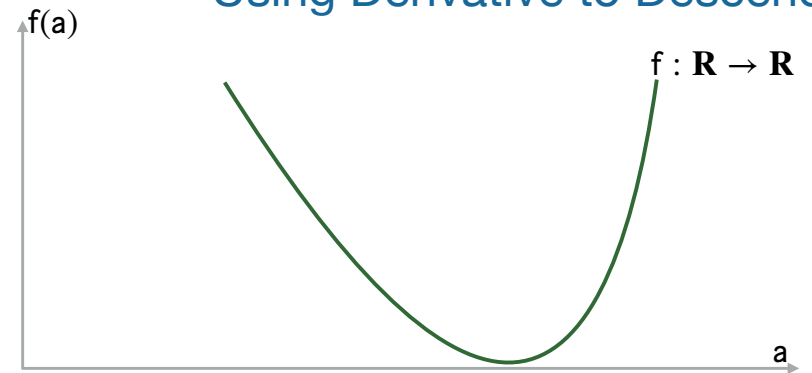
- Get closer to a^* by going in direction $-f'(a)$



© 2020 YORAM SINGER

7

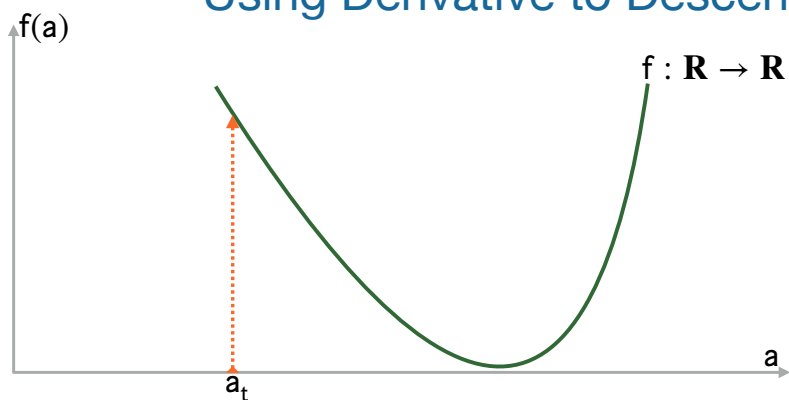
Using Derivative to Descend



© 2020 YORAM SINGER

8

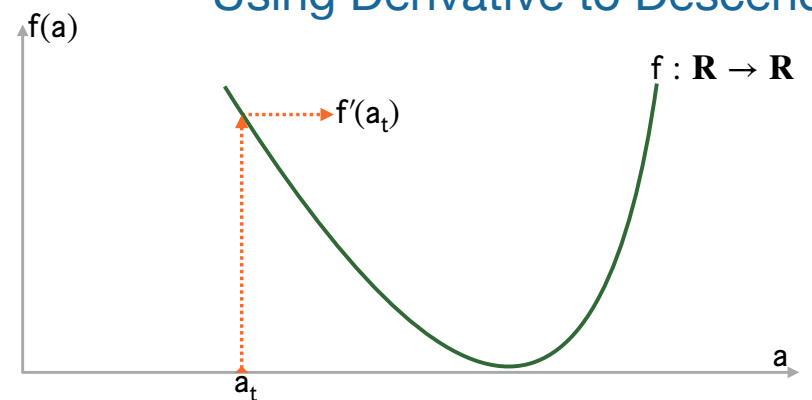
Using Derivative to Descend



© 2020 YORAM SINGER

8

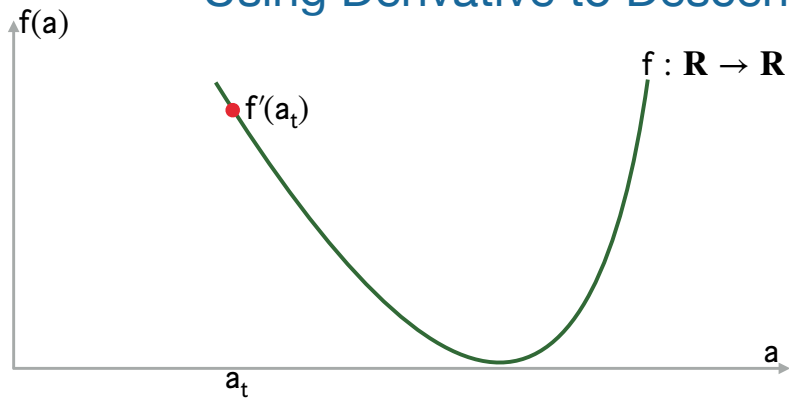
Using Derivative to Descend



© 2020 YORAM SINGER

8

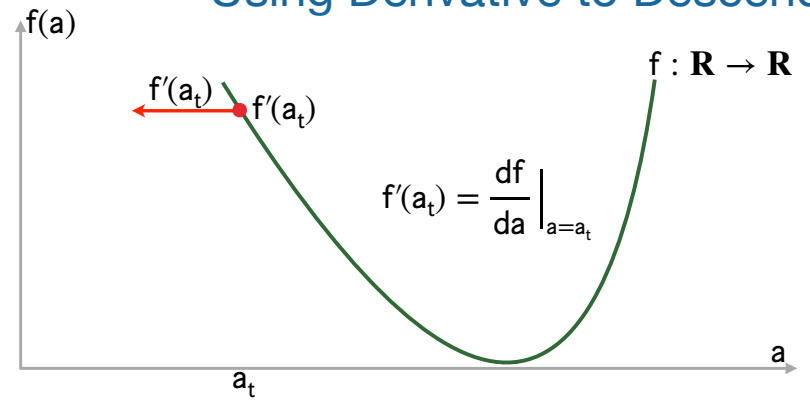
Using Derivative to Descend



© 2020 YORAM SINGER

8

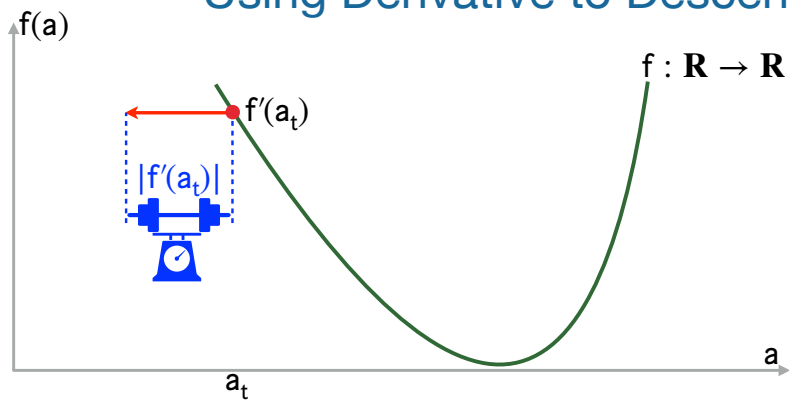
Using Derivative to Descend



© 2020 YORAM SINGER

8

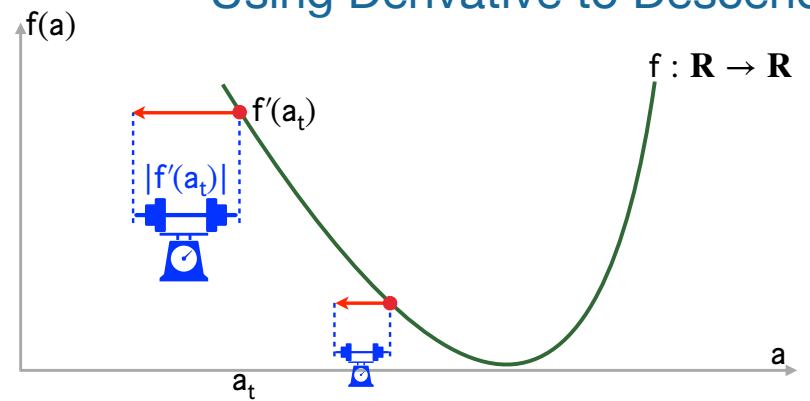
Using Derivative to Descend



© 2020 YORAM SINGER

8

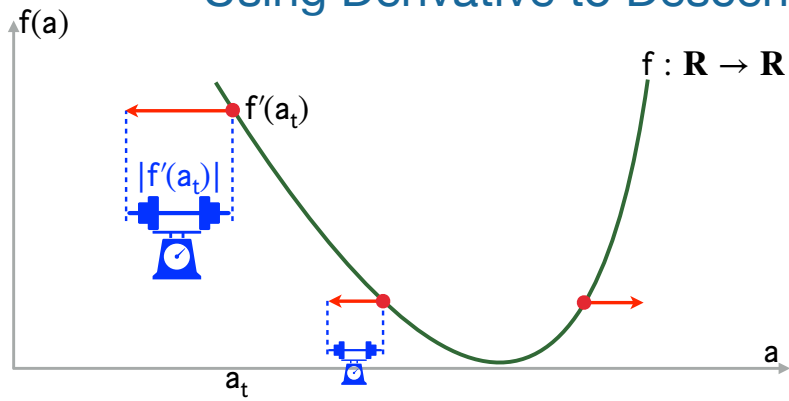
Using Derivative to Descend



© 2020 YORAM SINGER

8

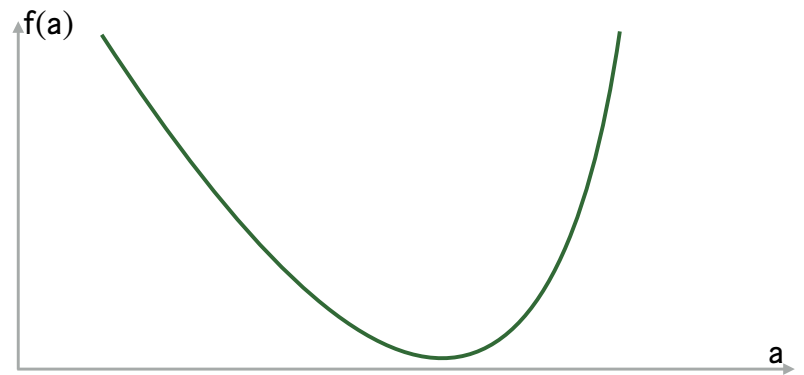
Using Derivative to Descend



© 2020 YORAM SINGER

8

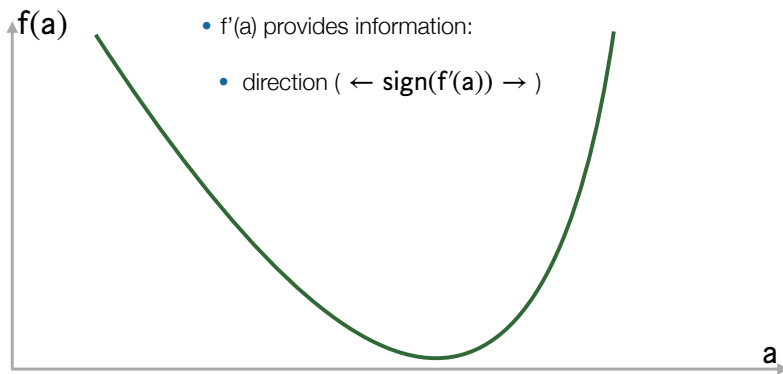
Iterative Derivative Procedure



© 2020 YORAM SINGER

9

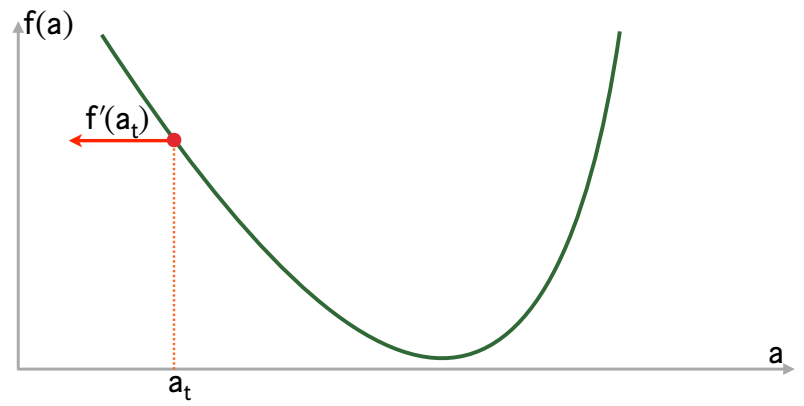
Iterative Derivative Procedure



© 2020 YORAM SINGER

9

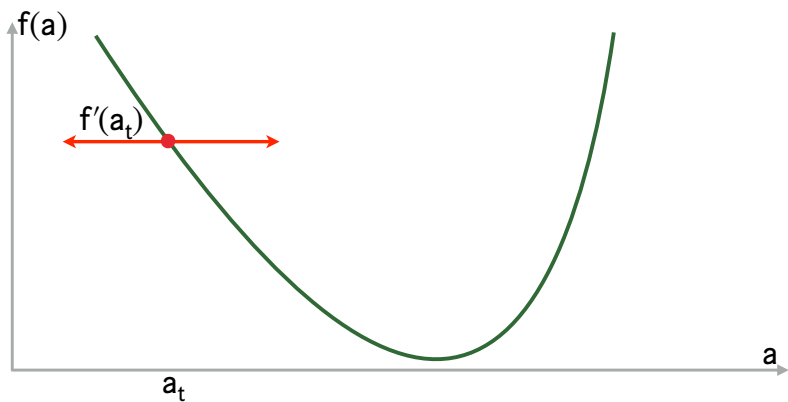
Iterative Derivative Procedure



© 2020 YORAM SINGER

9

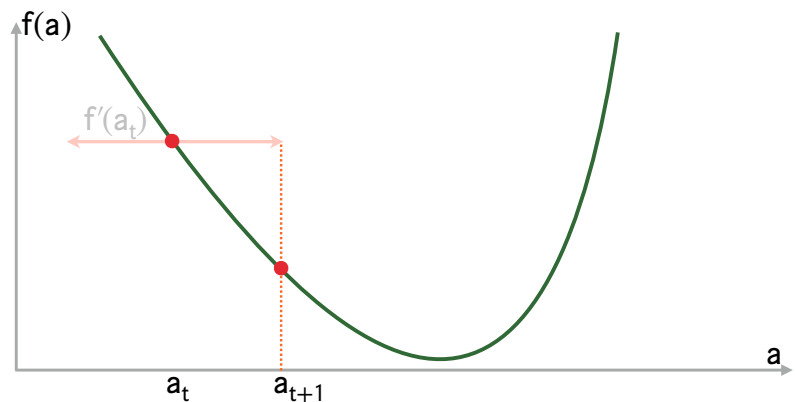
Iterative Derivative Procedure



© 2020 YORAM SINGER

9

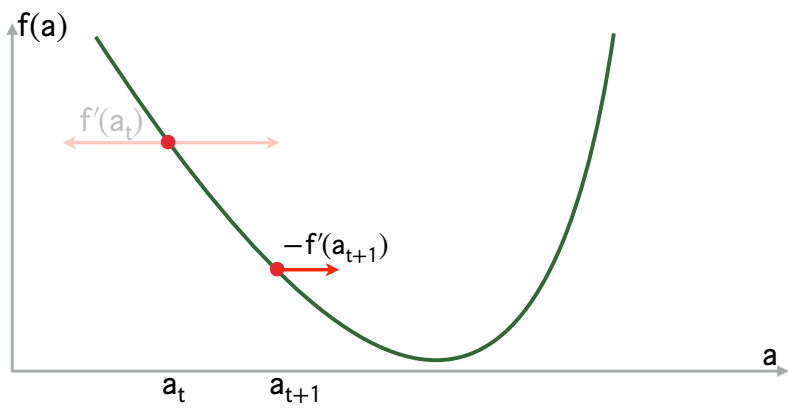
Iterative Derivative Procedure



© 2020 YORAM SINGER

9

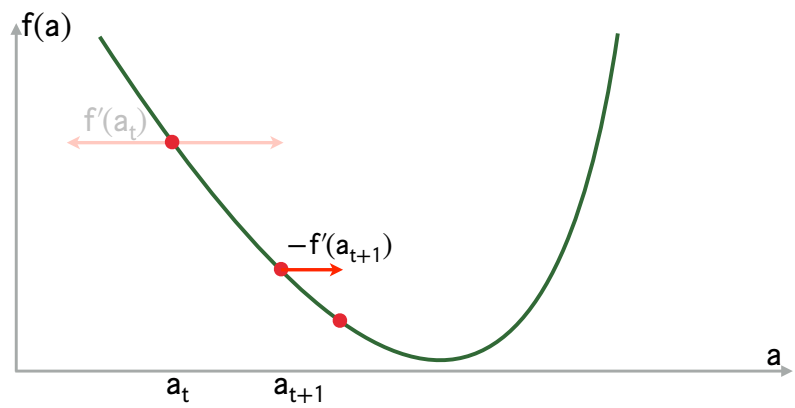
Iterative Derivative Procedure



© 2020 YORAM SINGER

9

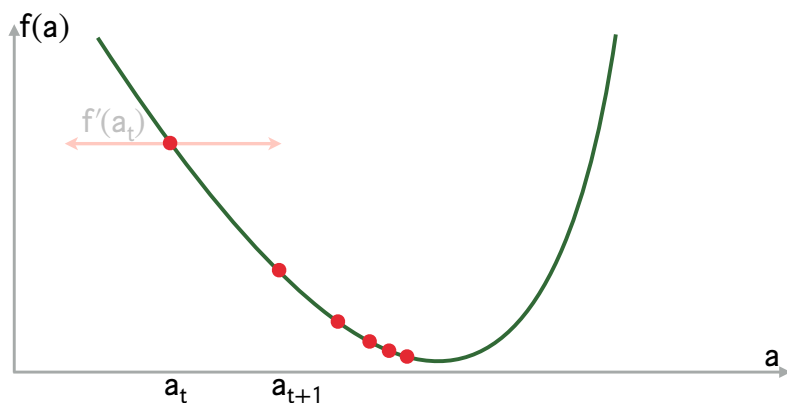
Iterative Derivative Procedure



© 2020 YORAM SINGER

9

Iterative Derivative Procedure



© 2020 YORAM SINGER

9

Iterative Derivative Procedure

- Input: function f
- Goal: find \hat{a} such that $|f(\hat{a}) - f(a)| \leq \epsilon$
- Choose initial value a_1
- Loop:
 - $a_{t+1} \leftarrow a_t - \eta f'(a_t)$
- Until ...

© 2020 YORAM SINGER

10

Iterative Derivative Procedure

- Input: function f
- Goal: find \hat{a} such that $|f(\hat{a}) - f(a)| \leq \epsilon$
- Choose initial value a_1
- Loop:
 - $a_{t+1} \leftarrow a_t - \eta f'(a_t)$
- Until ...

★ learning rate
★ step size

© 2020 YORAM SINGER

10

Learning Rate

- Crucial in many learning problems
- Fixed learning-rate can be used in certain circumstances
- Self-tuning procedure of learning-rate exist, notably AdaGrad
- In many applications:
 - Linear decrease $\eta_t = \frac{\eta_0}{b + st}$ where $\eta_0 \in [0.1, 1]$,
 - Sub-linear decrease $\eta_t \sim \frac{\eta_0}{\sqrt{t}}$

© 2020 YORAM SINGER

11

Example

- Find minimum of

$$f(x) = \log(1 + e^{x-a-\delta}) + \log(1 + e^{b-a-\delta})$$

- Derivative is

$$f'(x) = \frac{1}{1 + e^{a+\delta-x}} - \frac{1}{1 + e^{x+\delta-b}}$$

- $a=1$ $b=2$ $\delta=4.5$
- Run with different initializations and learning rates ($\eta \in \{0.01, 4, 40\}$)

© 2020 YORAM SINGER

12

Implementation of $f(x)$

```
import numpy as np

def smooth_l1(a, b, delta):
    def smooth_l1_close(x):
        lloss = np.log(1 + np.exp(x - a - delta))
        rloss = np.log(1 + np.exp(b - x - delta))
        return 0.5 * (rloss + lloss)
    return smooth_l1_close
```

© 2020 YORAM SINGER

13

Implementation of $f'(x)$

```
def smooth_l1_deriv(a, b, delta):
    def smooth_l1_deriv_close(x):
        rderiv = 1. / (1 + np.exp(a + delta - x))
        lderiv = 1. / (1 + np.exp(x + delta - b))
        return 0.5 * (rderiv - lderiv)
    return smooth_l1_deriv_close
```

© 2020 YORAM SINGER

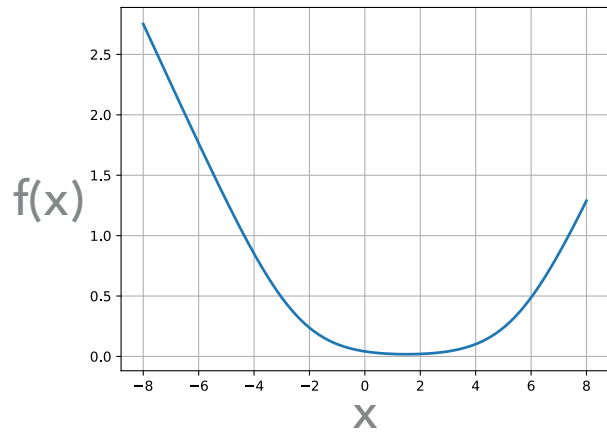
14

Derivative Descent

```
def derivative_descent(x0, deriv_func, eta):
    T = len(eta) + 1
    x = np.zeros(T)
    x[0] = x0
    for i in range(1, T):
        x[i] = x[i-1] - eta[i-1] deriv_func(x[i-1])
    return x
```

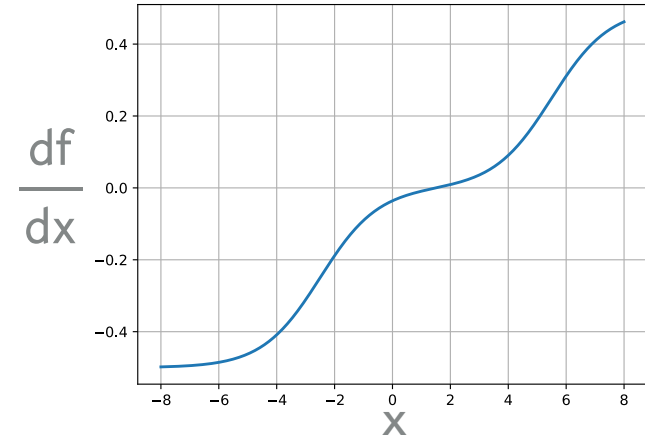
© 2020 YORAM SINGER

15



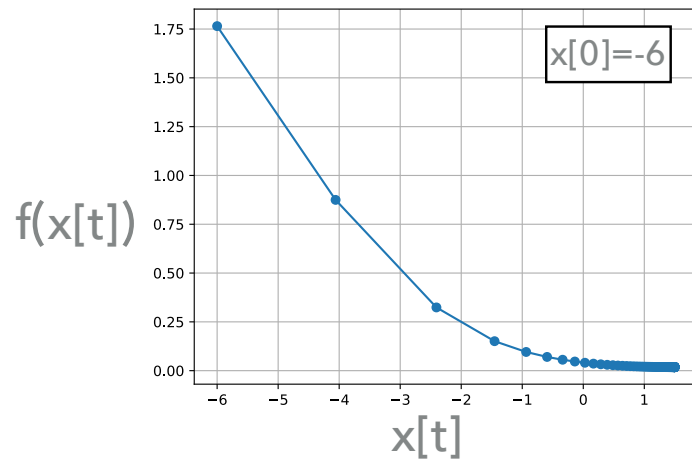
© 2020 YORAM SINGER

16



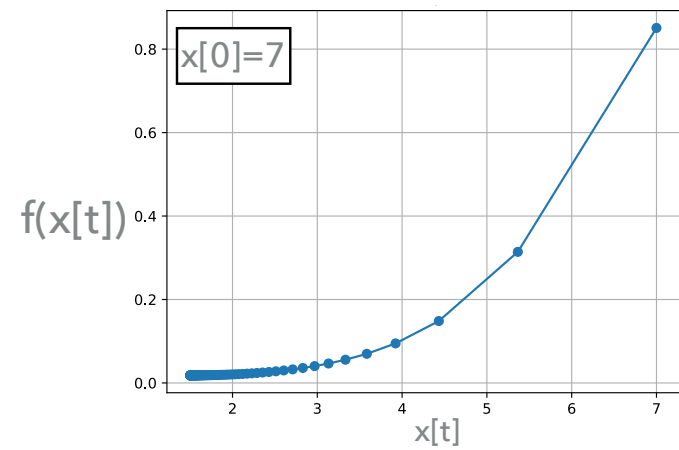
© 2020 YORAM SINGER

17



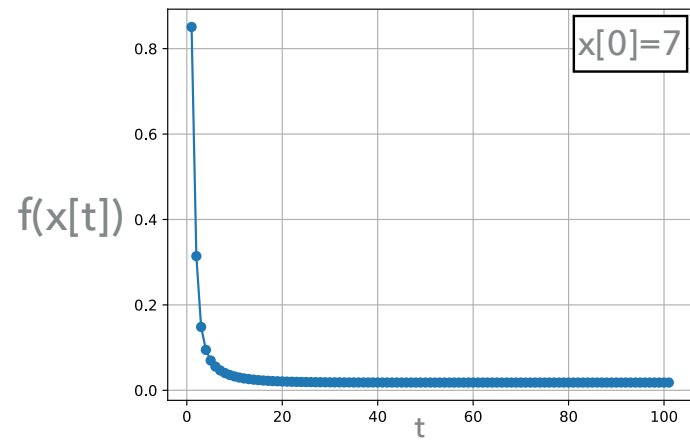
© 2020 YORAM SINGER

18



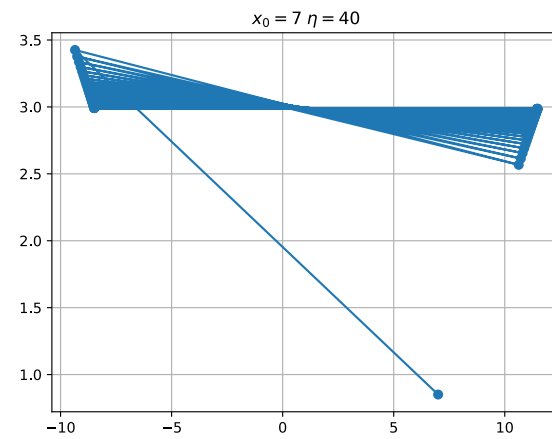
© 2020 YORAM SINGER

19



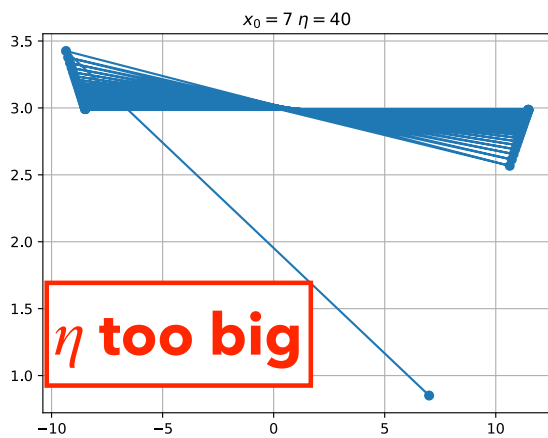
© 2020 YORAM SINGER

20



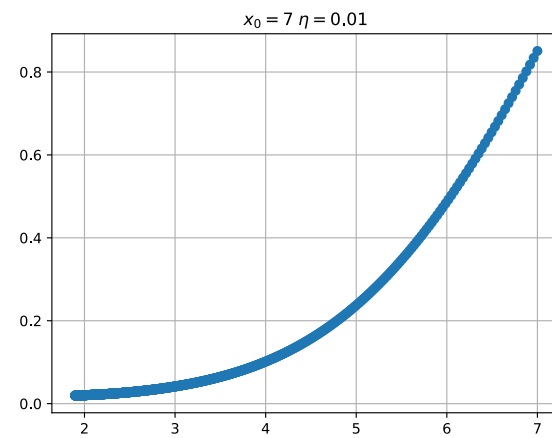
© 2020 YORAM SINGER

21



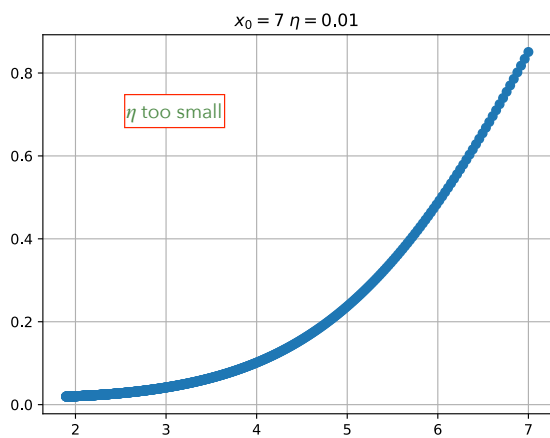
© 2020 YORAM SINGER

21



© 2020 YORAM SINGER

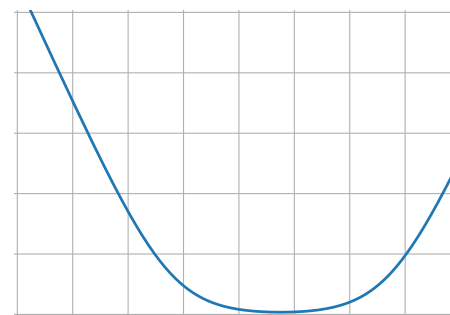
22



© 2020 YORAM SINGER

22

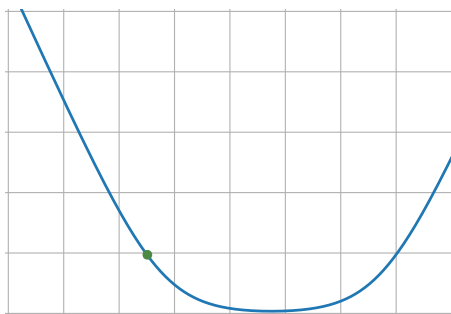
function f is called β -smooth: $f(x) \leq f(x_0) + f'(x_0)(x - x_0) + \frac{\beta}{2}(x - x_0)^2$



© 2020 YORAM SINGER

23

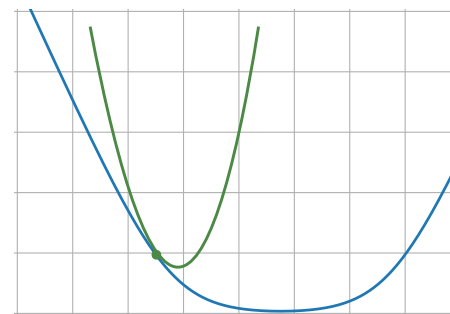
function f is called β -smooth: $f(x) \leq f(x_0) + f'(x_0)(x - x_0) + \frac{\beta}{2}(x - x_0)^2$



© 2020 YORAM SINGER

23

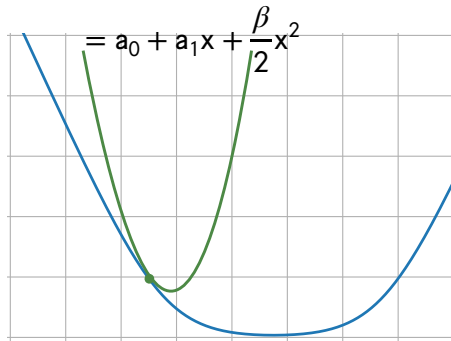
function f is called β -smooth: $f(x) \leq f(x_0) + f'(x_0)(x - x_0) + \frac{\beta}{2}(x - x_0)^2$



© 2020 YORAM SINGER

23

function f is called β -smooth: $f(x) \leq f(x_0) + f'(x_0)(x - x_0) + \frac{\beta}{2}(x - x_0)^2$

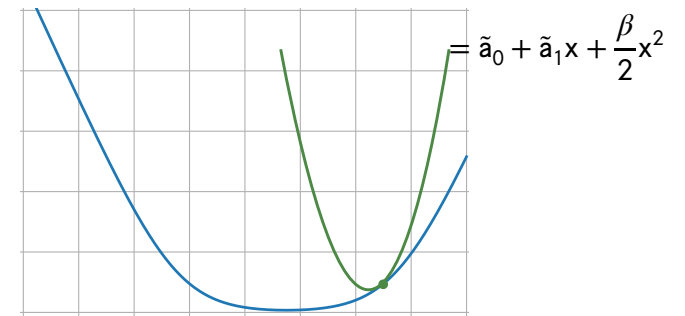


© 2020 YORAM SINGER

23

function f is called β -smooth: $f(x) \leq f(x_0) + f'(x_0)(x - x_0) + \frac{\beta}{2}(x - x_0)^2$

f is β -smooth $\Leftrightarrow f''(x) \leq \beta$



© 2020 YORAM SINGER

23

DD for Smooth Functions

- Set $\forall t : \eta_t = \frac{1}{\beta}$ which gives

$$x_{t+1} = x_t - \frac{1}{\beta} f'(x_t) = x_t - \frac{1}{\beta} g_t$$

© 2020 YORAM SINGER

24

DD for Smooth Functions

- Set $\forall t : \eta_t = \frac{1}{\beta}$ which gives

$$x_{t+1} = x_t - \frac{1}{\beta} f'(x_t) = x_t - \frac{1}{\beta} g_t$$

© 2020 YORAM SINGER

24

DD for Smooth Functions

- Set $\forall t : \eta_t = \frac{1}{\beta}$ which gives

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{1}{\beta} \mathbf{f}'(\mathbf{x}_t) = \mathbf{x}_t - \frac{1}{\beta} \mathbf{g}_t$$

- From Smoothness:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \mathbf{g}_t^T (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{\beta}{2} (\mathbf{x}_{t+1} - \mathbf{x}_t)^2$$

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{\beta} \mathbf{g}_t^2 + \frac{\beta}{2} \frac{\mathbf{g}_t^2}{\beta^2} = f(\mathbf{x}_t) - \frac{1}{2\beta} \mathbf{g}_t^2$$

© 2020 YORAM SINGER

24

DD for Smooth Functions

With $\eta_t = \frac{1}{\beta}$ and β -smoothness $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2\beta} \mathbf{g}_t^2$

© 2020 YORAM SINGER

25

DD for Smooth Functions

With $\eta_t = \frac{1}{\beta}$ and β -smoothness $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2\beta} \mathbf{g}_t^2$

However \mathbf{g}_t gets smaller as we approach the minimum

© 2020 YORAM SINGER

25

DD for Smooth Functions

With $\eta_t = \frac{1}{\beta}$ and β -smoothness $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2\beta} \mathbf{g}_t^2$

However \mathbf{g}_t gets smaller as we approach the minimum

It turns out that from convexity $\mathbf{g}_t \geq \frac{\Delta_t}{|\mathbf{x}_0 - \mathbf{x}^*|}$ where $\Delta_t = f(\mathbf{x}_t) - f(\mathbf{x}^*)$

© 2020 YORAM SINGER

25

DD for Smooth Functions

With $\eta_t = \frac{1}{\beta}$ and β -smoothness $f(x_{t+1}) \leq f(x_t) - \frac{1}{2\beta} g_t^2$

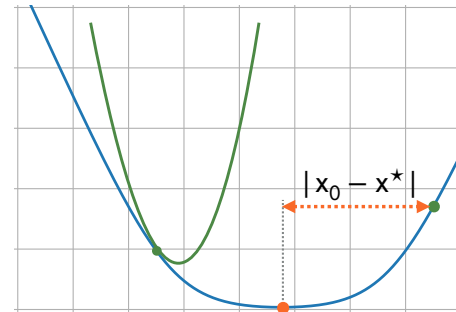
However g_t gets smaller as we approach the minimum

It turns out that from convexity $g_t \geq \frac{\Delta_t}{|x_0 - x^*|}$ where $\Delta_t = f(x_t) - f(x^*)$

This gives $f(x_{t+1}) - f(x^*) \leq \frac{2\beta |x_0 - x^*|}{t}$

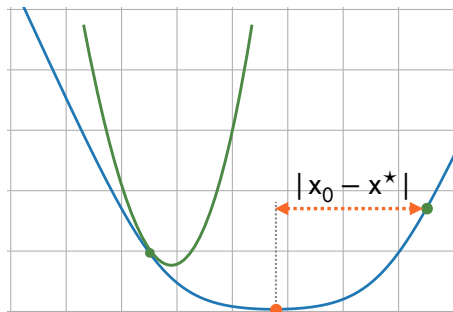
© 2020 YORAM SINGER

25



© 2020 YORAM SINGER

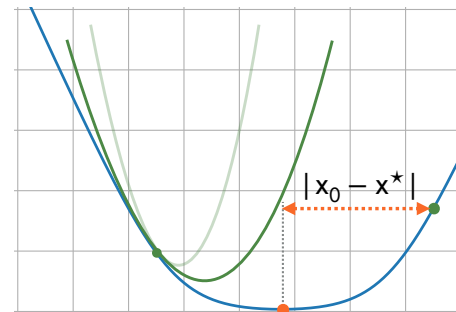
26



The closer to the optimum we start the faster we converge

© 2020 YORAM SINGER

26

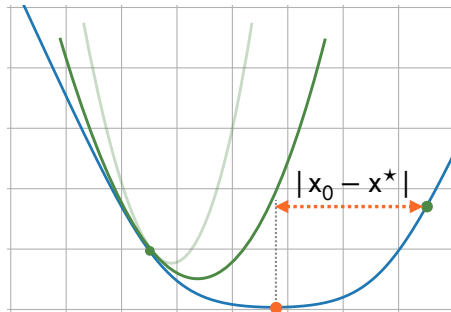


The closer to the optimum we start the faster we converge

© 2020 YORAM SINGER

26

The smaller β is the faster we convergence

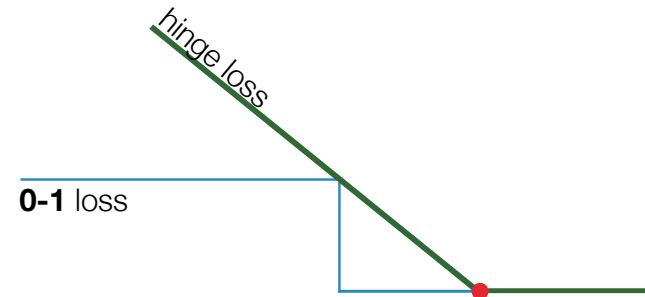


The closer to the optimum we start the faster we converge

© 2020 YORAM SINGER

26

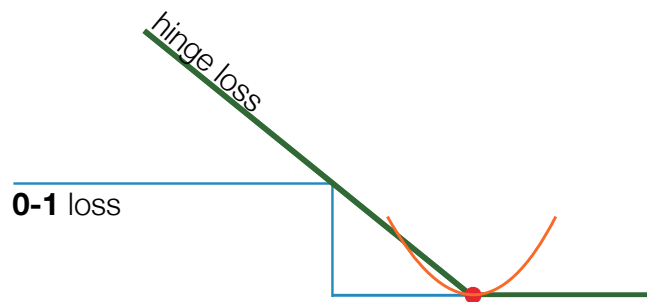
Non-smooth Functions



© 2020 YORAM SINGER

27

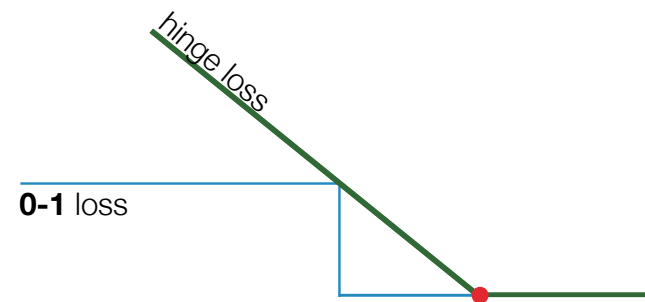
Non-smooth Functions



© 2020 YORAM SINGER

27

Non-smooth Functions

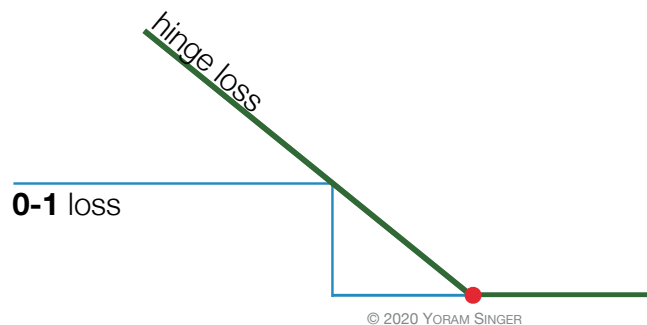


© 2020 YORAM SINGER

27

Non-smooth Functions

When loss function is not smooth:

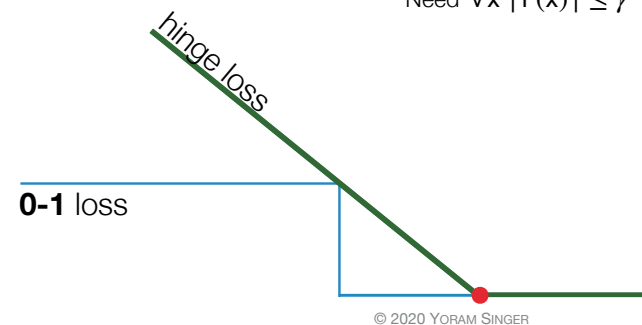


27

Non-smooth Functions

When loss function is not smooth:

Need $\forall x \ |f'(x)| \leq \gamma$



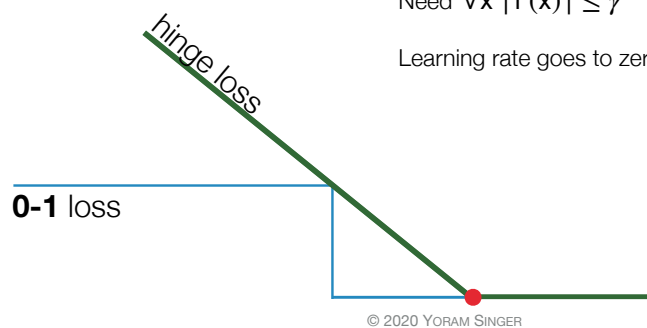
27

Non-smooth Functions

When loss function is not smooth:

Need $\forall x \ |f'(x)| \leq \gamma$

Learning rate goes to zero $\eta_0 t^{-1/2}$



27

Non-smooth Functions

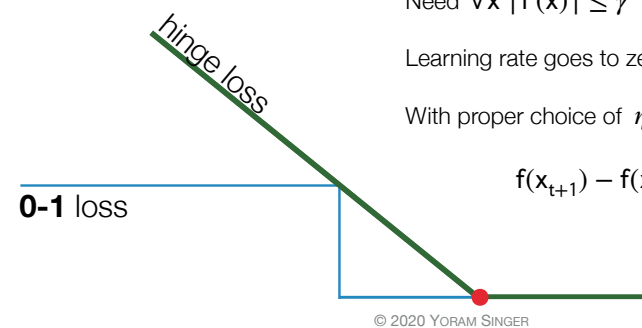
When loss function is not smooth:

Need $\forall x \ |f'(x)| \leq \gamma$

Learning rate goes to zero $\eta_0 t^{-1/2}$

With proper choice of η_0

$$f(x_{t+1}) - f(x^*) \leq \frac{\gamma |x_0 - x^*|}{\sqrt{t}}$$



27

Non-smooth Functions

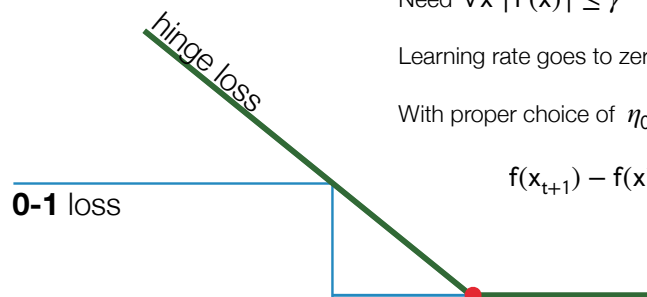
When loss function is not smooth:

Need $\forall \mathbf{x} \quad |\mathbf{f}'(\mathbf{x})| \leq \gamma$

Learning rate goes to zero $\eta_0 t^{-1/2}$

With proper choice of η_0

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \leq \frac{\gamma \|\mathbf{x}_0 - \mathbf{x}^*\|}{\sqrt{t}}$$



© 2020 YORAM SINGER

27

Non-smooth Functions

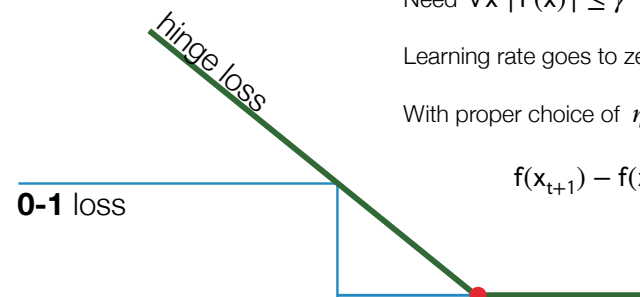
When loss function is not smooth:

Need $\forall \mathbf{x} \quad |\mathbf{f}'(\mathbf{x})| \leq \gamma$

Learning rate goes to zero $\eta_0 t^{-1/2}$

With proper choice of η_0

$$f(\mathbf{x}_{t+1}) - f(\mathbf{x}^*) \leq \frac{\gamma \|\mathbf{x}_0 - \mathbf{x}^*\|}{\sqrt{t}}$$



© 2020 YORAM SINGER

27

Next...

From one variable to a vector of variables

Gradients and their properties

Convexity of multivariate functions

Smoothness of multivariate functions

Gradient Descent

Stochastic Gradient Descent (SGD)

SGD for generalized linear models

SGD for non-linear models

© 2020 YORAM SINGER

28