

Machine Learning from Data

Lecture 9: Spring 2021

Today's Lecture

- Logistic Regression
- Gradient Descent

3 types

→ Probability

Previous Lecture

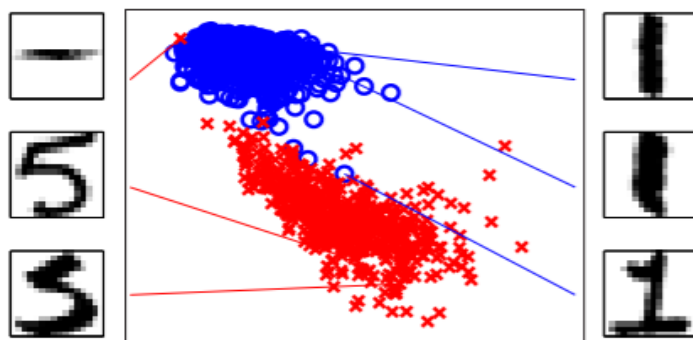
The linear signal:

$$s = \mathbf{w}^T \mathbf{X}$$

Ein 0

feature construction

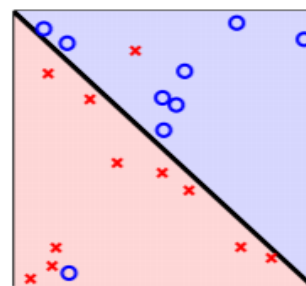
Good Features are Important



Before looking at the data, we can **reason** that symmetry and intensity should be good features based on our knowledge of the problem.

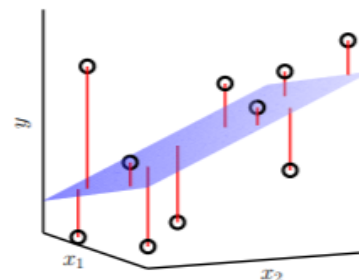
du

Algorithms



Linear Classification.

Pocket algorithm can tolerate errors
Simple and efficient



Linear Regression.

Single step learning:

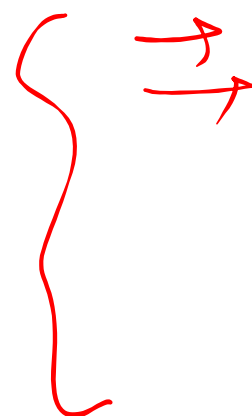
$$\mathbf{w} = \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Very efficient $O(Nd^2)$ *exact* algorithm.

yet R

Predicting a Probability (Logistic Regression)

→ Will someone have a heart attack over the next year?



age	62 years
gender	male
blood sugar	120 mg/dL40,000
HDL	50
LDL	120
Mass	190 lbs
Height	5' 10"
...	...

Classification: Yes/No

Logistic Regression: Likelihood of heart attack

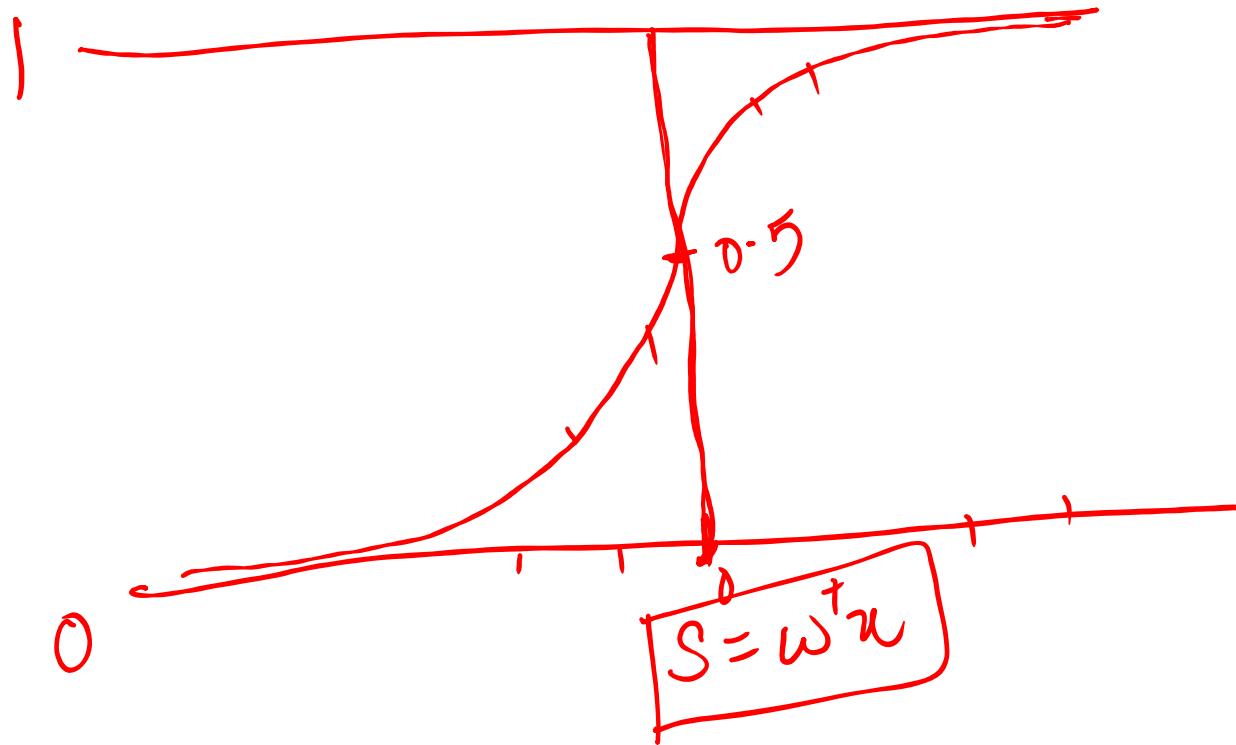
linear signal
↓ transform
sigmoid
Probability

± 1

$-1, +1$
 0.5
 0.3
 0.9

logistic regression $\equiv y \in [0, 1]$

Logistic Regression \longrightarrow Probability
 $s = w^T x \longrightarrow \sigma(w^T x)$ $y \in [0, 1]$



$$\sigma(s) = \frac{e^s}{1 + e^s}$$

(Mathematically convenient)

Properties of the Sigmoid

$$\bullet \quad \theta(s) = \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}}$$

$$\frac{1}{\theta(s)} = 1 + e^{-s}$$

$$\theta(-s) = \frac{e^{-s}}{1 + e^{-s}} = \frac{1}{1 + e^s} = 1 - \theta(s)$$

How to get the Probability?

- Patients : x_1, x_2, \dots, x_n
Observe : y_1, y_2, \dots, y_n (+1/-1)
(Yes/No)
 y_i \nearrow
Don't have the exact data
 \rightarrow Probability values

Target function

- Probabilistic target function (Ground Truth)

$$p[+1|x] = f(x)$$

$$p[-1|x] = 1 - f(x)$$

$$h(x) = \theta(w^T x) \approx \underset{\uparrow}{f(x)}$$

$$\frac{g(x)}{E_{in}(w)}$$

$$P[+1|x] \approx h(x)$$

$$P[-1|x] \approx 1 - h(x)$$

Observe $x_n \rightarrow y_n = +1$
 $x_n \rightarrow y_n = -1$

$$\therefore h(x_n) \approx \frac{1}{2} (1 + y_n)$$

then $h(x_n) \approx 1$
 then $h(x_n) \approx 0$

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N \left(h(x_n) - \frac{1}{2}(1+y_n) \right)^2$$

We cannot use this because:

i) Inconvenient to minimize.

ii) Where is the probability approx.?

What is a better error function?

CROSS ENTROPY ERROR

$$E_{in}(w) = \frac{1}{N} \sum_{i=1}^N \ln(1 + e^{-y_n(w^T x_n)})$$

Better behaved + Probabilistic interpretation.

Data: $(x_1, y_1) (x_2, y_2) \dots (x_N, y_N) \dots y_n \in \{-1, +1\}$

\downarrow
IID \longrightarrow One of the pillars of theory

Observing $y_1, y_2, \dots, y_N \longrightarrow x_1, x_2, \dots, x_N$

$$P[y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_n] \longrightarrow \text{LIKELIHOOD}$$

$$= P(y_1|x_1) \cdot P(y_2|x_2) \cdots P(y_n|x_n)$$

$$= \prod_{n=1}^N P(y_n|x_n)$$

$$\longrightarrow \begin{cases} h(x_n) & \text{if } y_n = +1 \\ 1 - h(x_n) & \text{if } y_n = -1 \end{cases}$$

$$h(x_n) = \underline{\underline{\theta(\omega^T x_n)}}$$

$$1 - h(x_n) = 1 - \theta(\omega^T x_n)$$

$$= \theta(-\omega^T x_n)$$

$$\boxed{\theta(-s) = 1 - \theta(s)}$$

If $y_n = +1$ then $h(x_n) \longrightarrow \theta(\omega^T x_n)$
 If $y_n = -1$ " $h(x_n) \longrightarrow \theta(-\omega^T x_n)$
 $\longrightarrow \theta(y_n(\omega^T x_n))$

$$\prod_{n=1}^N P(y_n | x_n) \rightarrow \theta(y_n(w^T x_n))$$

$$L = \prod_{n=1}^N \theta(y_n(w^T x_n)) \quad \text{Likelihood}$$

Means

Goal: Choose w 's \rightarrow maximize
 Maximum Likelihood (MLE)

$L(w) \rightarrow \text{maximize}$

$$\arg\max_w L(w) \iff \arg\max_w \ln L(w) \checkmark$$

$$\arg\max_w \frac{1}{N} \ln L(w) \iff \arg\min_w -\frac{1}{N} \ln L(w) \checkmark$$

Choose weights w to minimize $-\frac{1}{N} \ln L(w)$

$$\therefore \min_{\text{Cross entropy error}} -\frac{1}{N} \sum_{i=1}^N \ln \checkmark \theta(y_n w^T x_n) = \frac{1}{N} \sum_{i=1}^N \ln \frac{1}{\theta(y_n w^T x_n)}$$
$$\rightarrow J_{\text{in}}(w) = \frac{1}{N} \sum_{i=1}^N \ln(1 + e^{-y_n w^T x_n})$$

How to minimize E_{in} ?

•

$$\underline{\underline{E_{in}(\mathbf{w})}} = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n \cdot \mathbf{w}^T \mathbf{x}_n})$$

$$y_n = -1, \mathbf{w}^T \mathbf{x}_n \ll 0$$

$$y_n = +1, \mathbf{w}^T \mathbf{x}_n \gg 0$$

Classification: Pocket Algo.

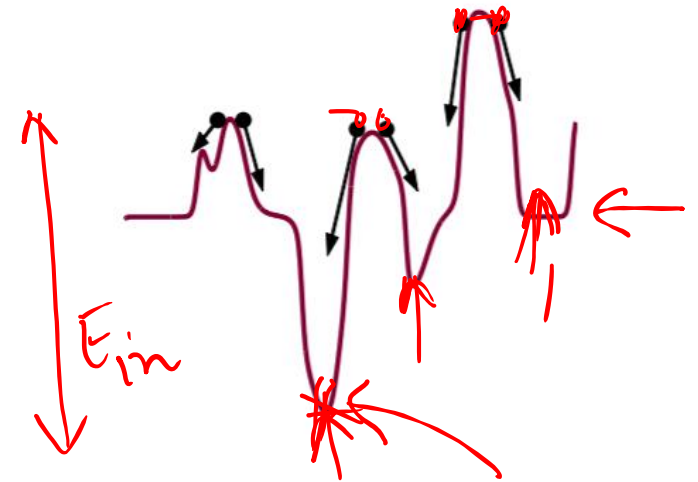
Regression: Pseudo Inverse

$$\underline{\underline{\nexists E_{in}(\mathbf{w}) = 0}}$$

Analytically not possible
Iteratively possible.

Ball Phenomena

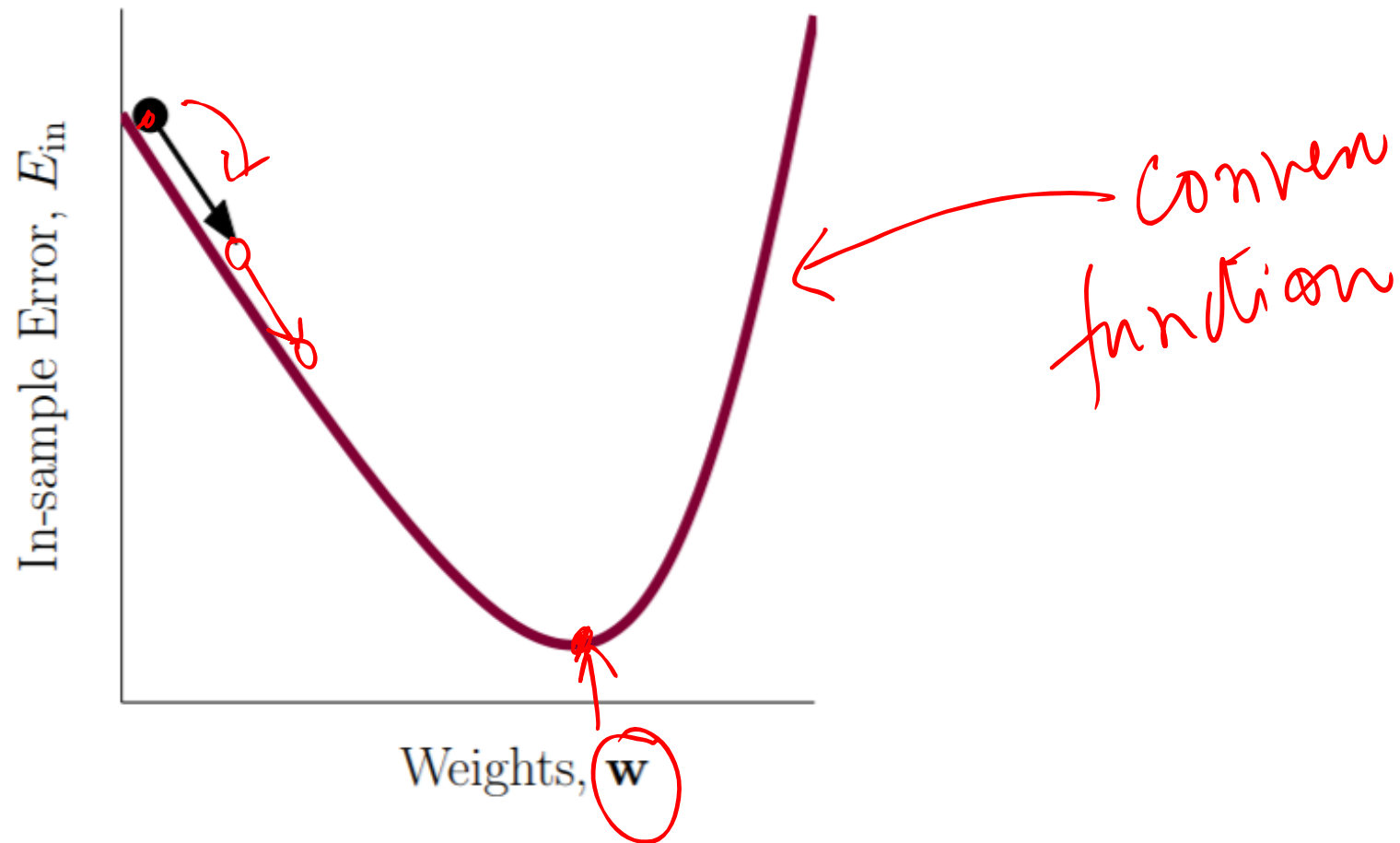
Ball on a complicated hilly terrain
— rolls down to a *local valley*
↑
this is called a *local minimum*



Questions:

How to get to the bottom of the deepest valey?

How to do this when we don't have gravity? ←



How do we roll down?

Iteration (t)

$\omega(t)$



Pick a direction i.e. a unit vector \hat{v} . Take a small step in that direction. η

Pick \hat{v} to make $E_{in}(\omega(t+1))$ as small as possible ΔE_{in}

$$\omega(t+1) = \omega(t) + \eta \hat{v}$$

$$E_{in}(\omega(t))$$

$$\Delta = \underbrace{E_{in}(\omega(t+1))}_{\text{By}} - \underbrace{E_{in}(\omega(t))}_{\text{Cauchy Schwartz inequality}}$$

$$\Delta E_{in} = \underbrace{\nabla E_{in}(\omega(t))(\eta \hat{v})}_{\text{Cauchy Schwartz inequality}} \geq -\eta \|\nabla E_{in}(\omega(t))\|$$

$$\hat{v} = - \frac{\nabla E_{in}(w(t))}{\|\nabla E_{in}(w(t))\|}$$

-ve of the gradient

Summarize

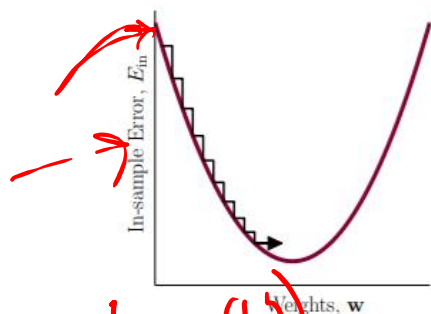
- i) Initialize $t=0$, $w(0)$
- ii) $\rightarrow t$: $w(t)$ Compute $\nabla E_{in}(w(t))$
- $$\underline{w(t+1)} = w(t) - \eta \frac{\nabla E_{in}(w(t))}{\|\nabla E_{in}(w(t))\|}$$
- $\rightarrow t+1$

Normalized gradient descent.

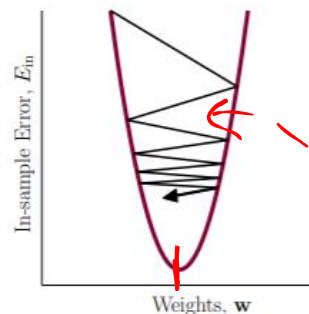
Step Size

$w(t+1) = w(t) - \eta_t \nabla E_{in}(w(t))$
 Batch gradient descent!

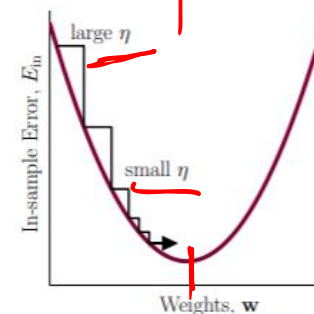
η too small



η too large

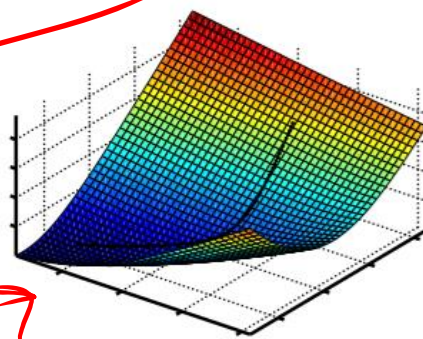


variable η_t - just right

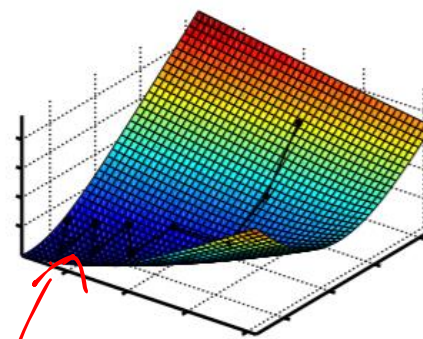


$\nabla E_{in} \approx 0$
 $\nabla E_{in} > 0$

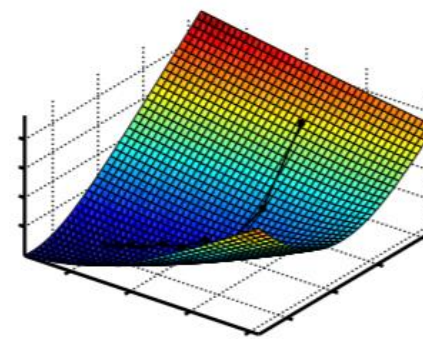
$\eta \propto \|\nabla E_{in}\|$
 variable η_t - just right



$\eta = 0.1$; 75 steps



$\eta = 2$; 10 steps



variable η_t ; 10 steps

- i) $E_{in}(w) \rightarrow$ differentiable
- ii) Inefficient. $\nabla E_{in} \rightarrow O(Nd)$

Updating weights $\rightarrow O(d+1)$

How to roll down efficiently?

Stochastic gradient descent

$$E_{in}(w) = \frac{1}{N} \sum_{i=1}^N e(w^T x_i, y_i)$$

Randomly

$x_n \rightarrow e(w^T x_n, y_n) \rightarrow$ minimize

$$w(t+1) \Rightarrow w(t) - \eta \nabla e(w^T x_n, y_n)$$

\uparrow
 ∇E_{in}

$$e(\omega^T x_n, y_n) = \ln(1 + e^{-y_n(\omega^T x_n)})$$

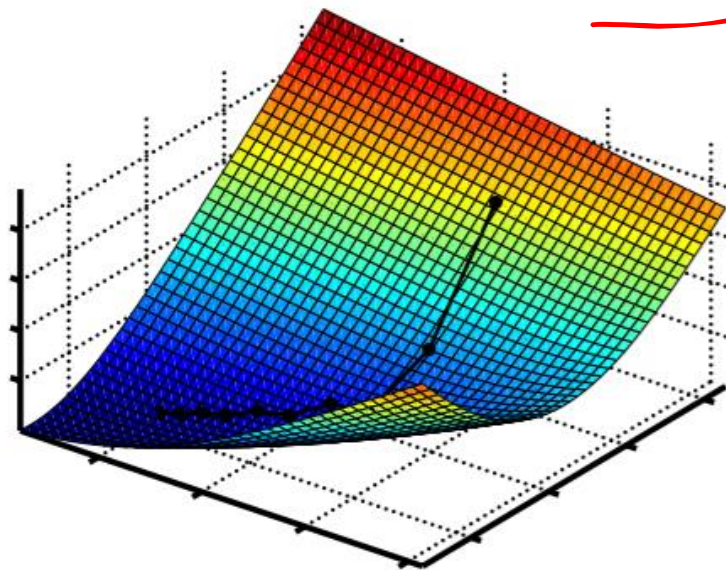
$$\nabla e(\omega^T x_n, y_n) = \frac{1}{1 + e^{-y_n \omega^T x_n}} - y_n x_n \cdot e^{-y_n(\omega^T x_n)}$$

$$= \frac{-y_n x_n}{1 + e^{y_n \omega^T x_n}}$$

N

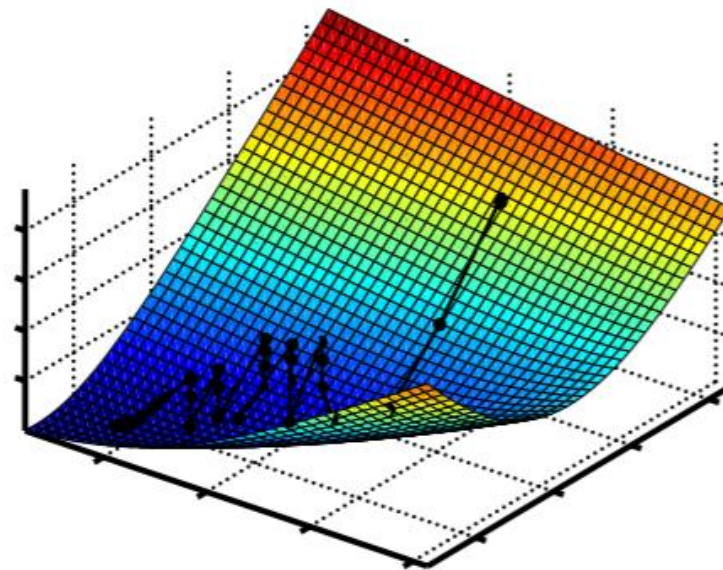
$$\omega(t+1) = \omega(t) + \frac{y_n x_n}{1 + e^{y_n \omega^T x_n}}$$

GD ✓ Bath



$\eta = 6$
10 steps
 $N = 10$

SGD ✓



$\eta = 2$
30 steps

Thanks!