

$$l(w) = \sum_{i=1}^n (y_i - f_w(x_i))^2$$

"least squares"

$$\operatorname{argmin}_w l(w)$$

$$w^* \in \operatorname{argmin}_w \sum_{i=1}^n \left(y_i - \sum_{j=1}^m w_j x_{i,j} \right)^2$$

Black box optimization: ^(BBO) Optimization when we know little about l .

- Given w , can query $l(w)$.
- Many BBO algorithms.
 - Genetic algs.
 - CEM
 - CMA-ES
 - * - Hill-climbing / simulated annealing.
 - Finite difference methods.
 - SPSA
 - Bayesian optimization.

Lecture 4

- BBO

- Gradient Descent

Hill Climbing

1) Pick initial point w_0 arbitrarily

2) Randomly select a nearby point w' \nearrow eta

Example: $\forall j \in \{1, \dots, m\}$ $w'_j = w_j + \eta_j$

$\eta_j \sim N(0, \sigma^2)$

noise \leftarrow

is sampled from.

normal distribution

hyperparameters.

variance. (small)

hyperparameter

alternative: $\eta_j \sim U(c_1 - c)$
 \downarrow
uniform distribution.
 c is a hyperparameter.

3) IF $l(w') < l(w)$

$w \leftarrow w'$

4) Go to #2

duration is HP

Stop when tired of waiting, when decrease in l is very small (over last few improvements), or when no improvement has been found after trying many w' .

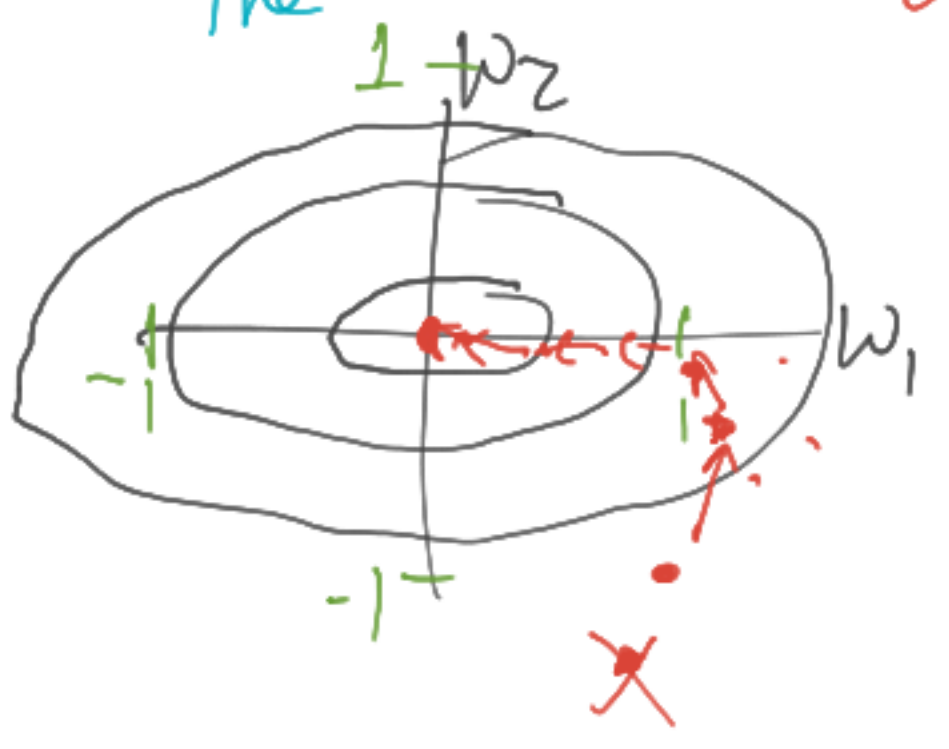
Example: $l(w) = w_1^2 + 2w_2^2$



- $w = (1, 0)$, $l(w) = 1$
- $w = (-1, 0)$, $l(w) = 1$
- $w = (0, 1)$, $l(w) = 2$
- $w = (0, -1)$, $l(w) = 2$
- $w = (0, 0)$, $l(w) = 0$

all points where $l(w)$ takes the same value.

Init $w = (1, -1)$
 ~~$w' = (.9, -1.2)$~~
 $w' = (1.1, -0.5)$
 ↳ new w



Pros

- Sometimes simple
- Don't require much knowledge about l .
- Some BBO algs use more knowledge.

Cons:

- Sophisticated BBO algs quite complex. (CMA-ES)
- Don't leverage additional knowledge about l .

Idea: What if we could use knowledge about l to pick a good direction?

$$l(w) = \sum_{i=1}^n (y_i - w_1 x_{i,1})^2$$

$$= \sum_{i=1}^n (y_i^2 - 2w_1 x_{i,1} y_i + w_1^2 x_{i,1}^2)$$

quadratic function of w_1

- Least squares loss, l , is always quadratic for linear models \Rightarrow one global minimum.

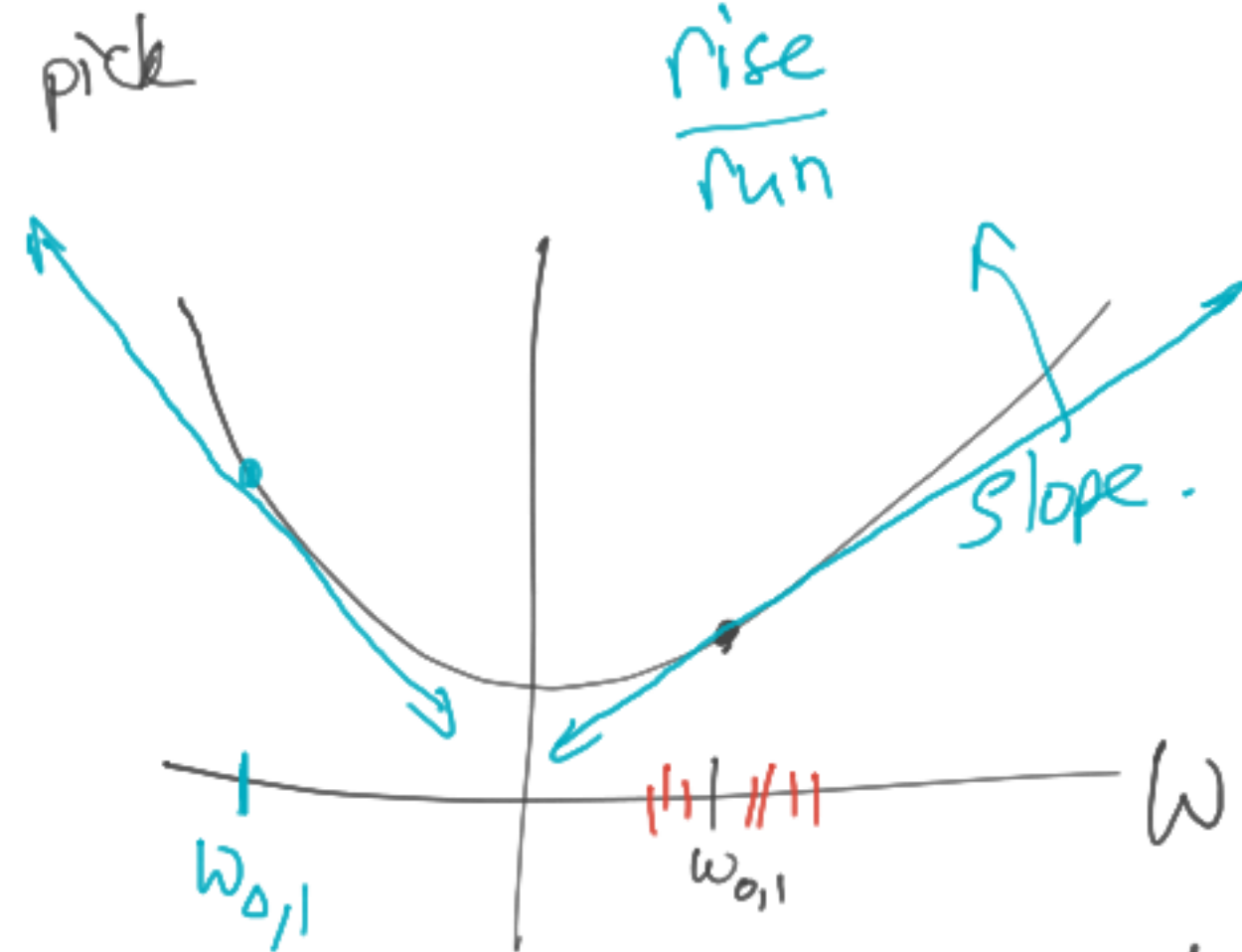
$x_{i,j} = j^{\text{th}}$ feature of i^{th} input

$w_j = j^{\text{th}}$ weight

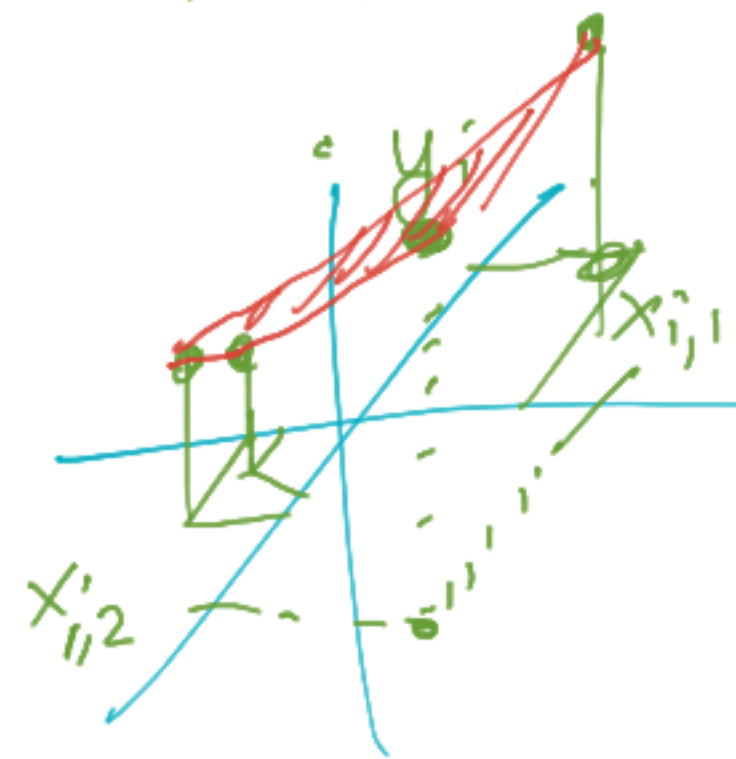
x_i is a vector.

$y_i \in \mathbb{R}$

$(x_i, y_i) = i^{\text{th}}$ data point



$w_{k,i} = j^{\text{th}}$ weight of k^{th} w.

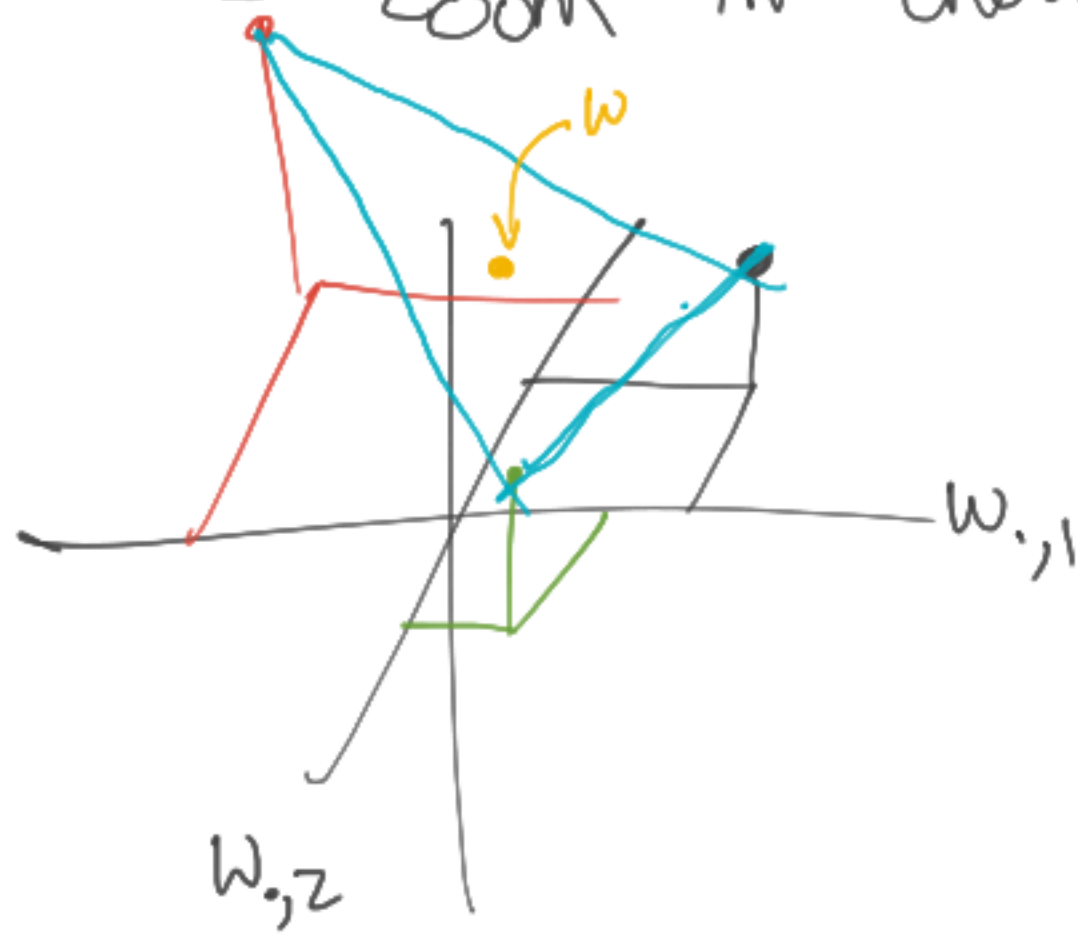


$M=2$

- Idea: use negative of the slope as our downwards direction

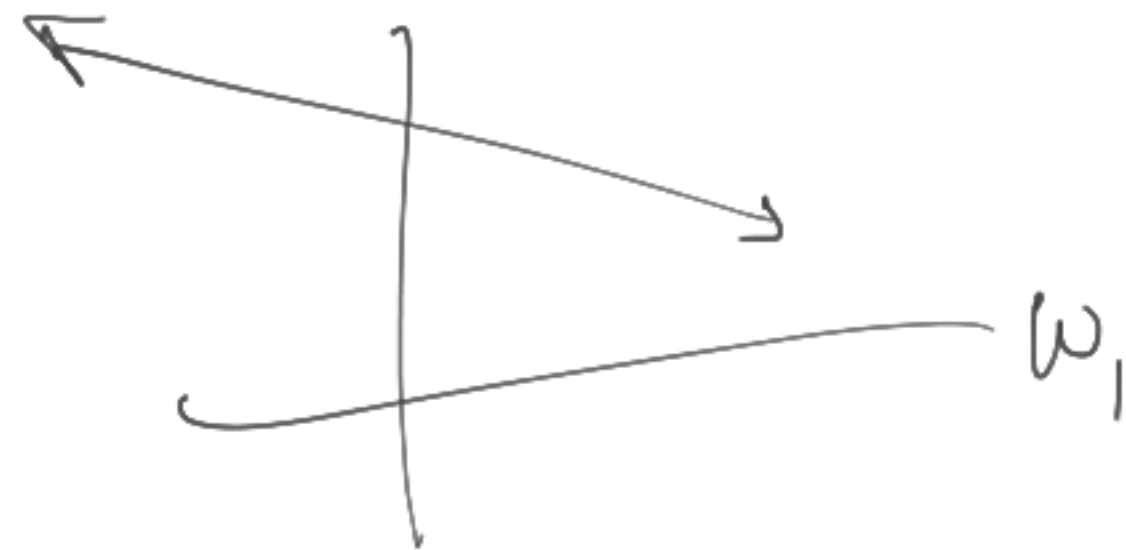
- Descent direction is a local property "descent direction"
↳ depends only on immediate surroundings, of current position.

- "Zoom in" enough, surface is flat.

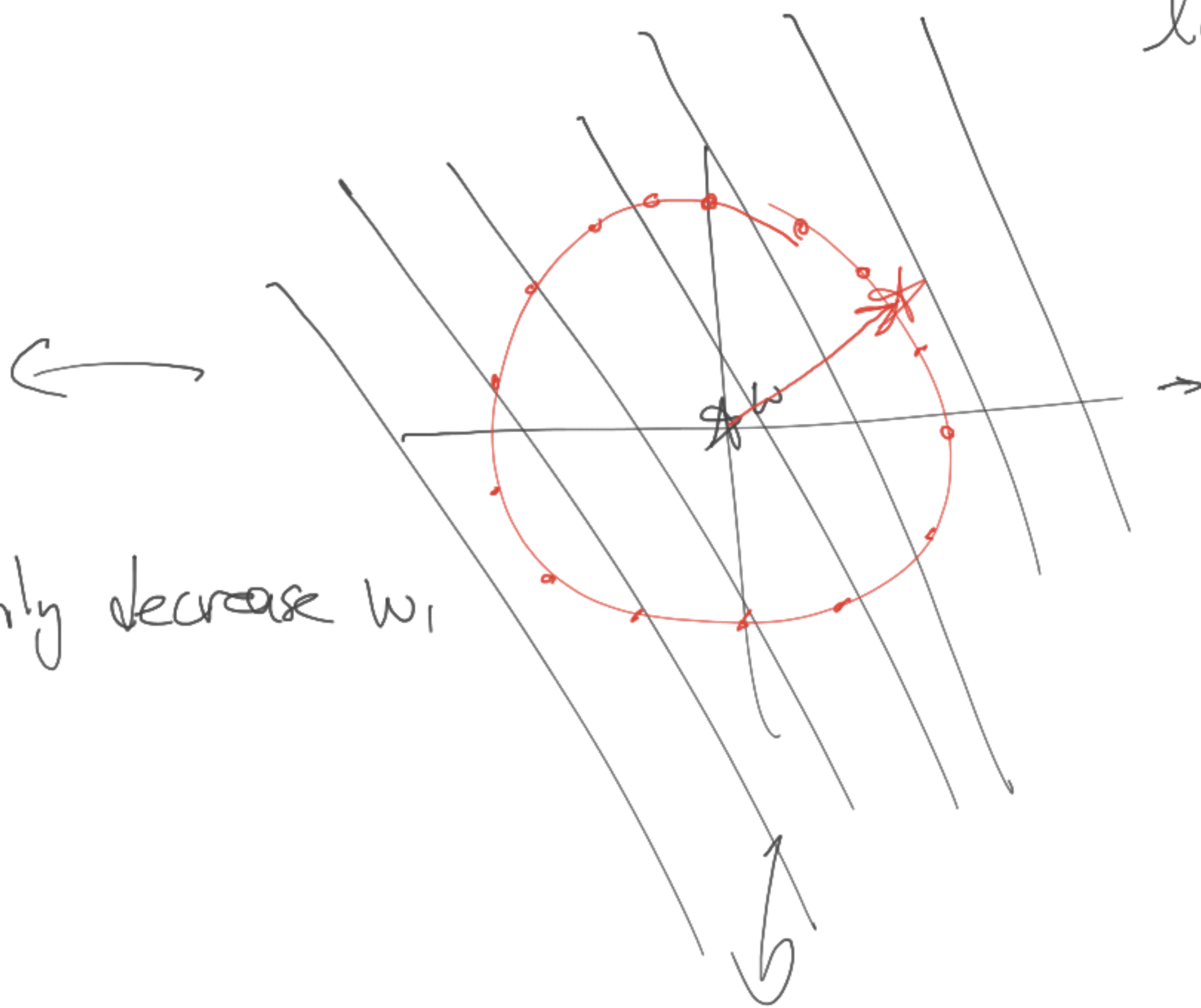


$$\frac{\partial l(w)}{\partial w_1} = \text{slope along } w_{1,1} \text{ axis}$$

$$g(w_1) = l(w_1, w_2)$$
$$\frac{\partial g(w_1)}{\partial w_1}$$



A4: Only decrease w_1



$l(w)$ level sets. (locally)

$$\frac{\partial l(w)}{\partial w_1} = -10$$

$$\frac{\partial l(w)}{\partial w_2} = -1$$

The direction of steepest ascent is the gradient

$$\nabla l(w) = \left(\frac{\partial l(w)}{\partial w_1}, \frac{\partial l(w)}{\partial w_2}, \dots, \frac{\partial l(w)}{\partial w_m} \right)$$

↓
nabla

The direction of steepest descent is negative
the gradient:

$$-\nabla l(w).$$

sequenz
von w 's

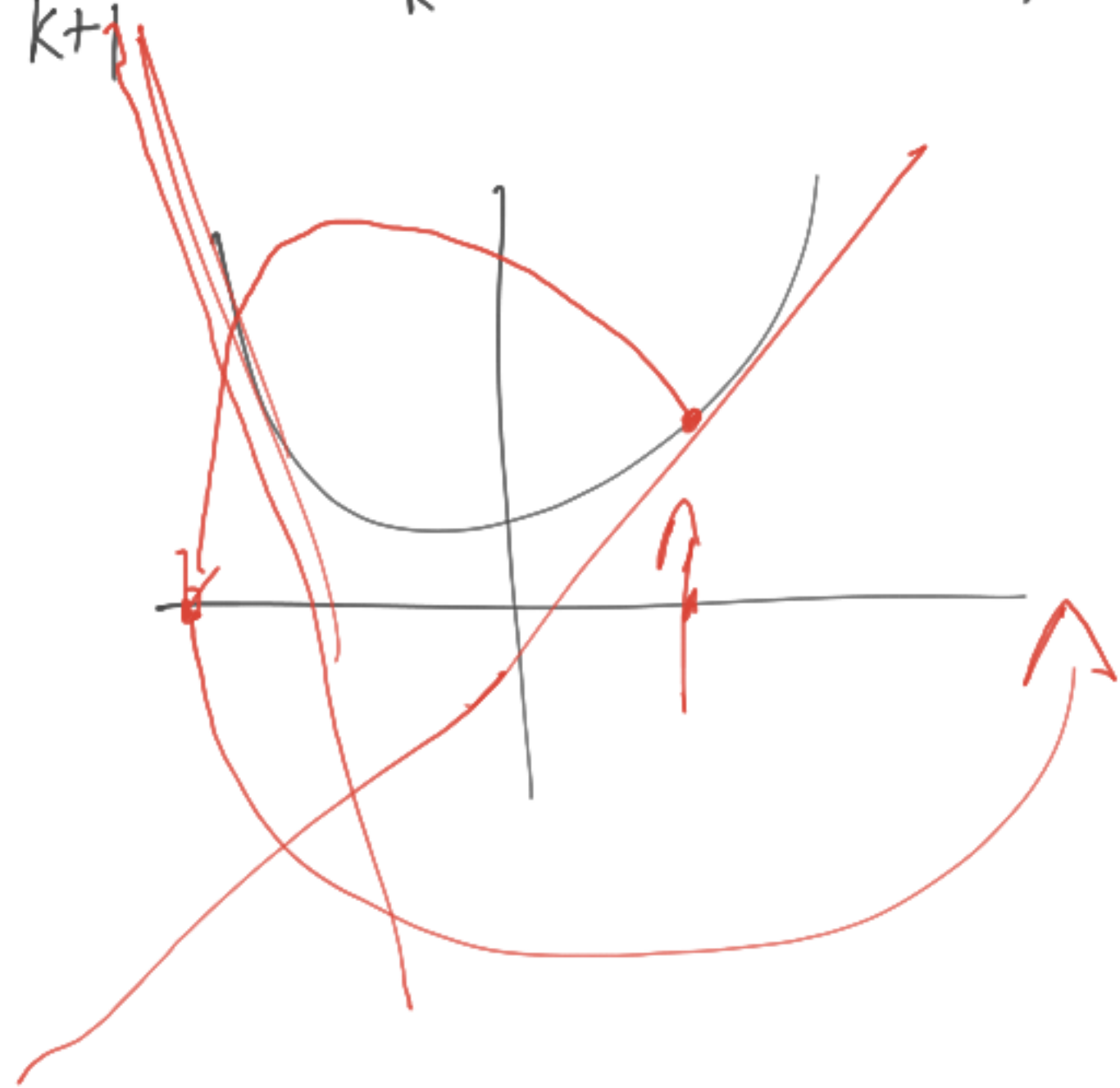
$$w' = w - \nabla l(w)$$

$$w_{k+1} = w_k - \nabla l(w_k)$$

Θ_j

$$w_{k+1,j} = w_{k,j}$$

$$- \frac{\partial l(w_k)}{\partial w_{k,j}}$$



$$w_{k+1} = w_k - \alpha \nabla l(w_k)$$

$$\forall j \in \{1, \dots, m\} \quad w_{k+1, j} = w_{k, j} - \alpha \frac{\partial l(w_k)}{\partial w_{k, j}}$$

Gradient descent.

"small"
Hyperparameter.

0.1

0.01

0.001