# Lecture #3 — Gradient Descent

Readings: Murphy 8-8.2.1
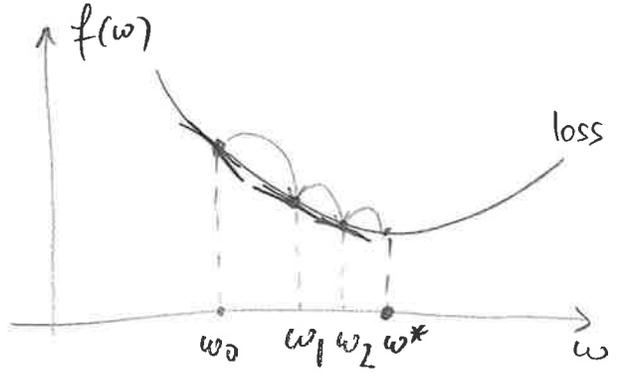         MIT Notes

> Standard ML paradigm is:

Define **Loss**, then optimize:

Example: $\boxed{\hat{w}_{LS} = \arg\min_{w} \| y - Xw \|_2^2}$

And we had derived: $\boxed{\hat{w}_{LS} = (X^T \cdot X)^{-1} X^T \cdot y}$ as a closed form solution.

> But, for most losses used in practice, there is no closed-form solution!
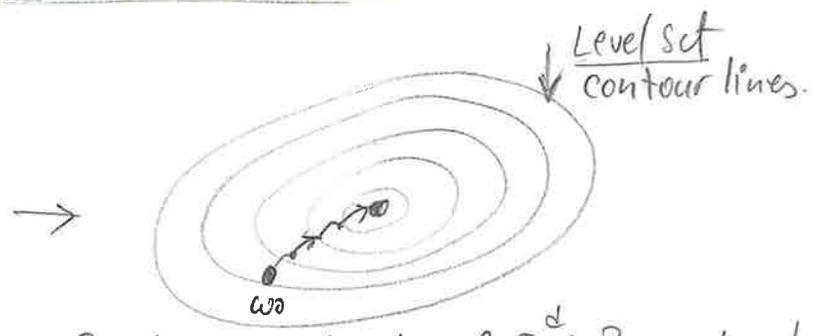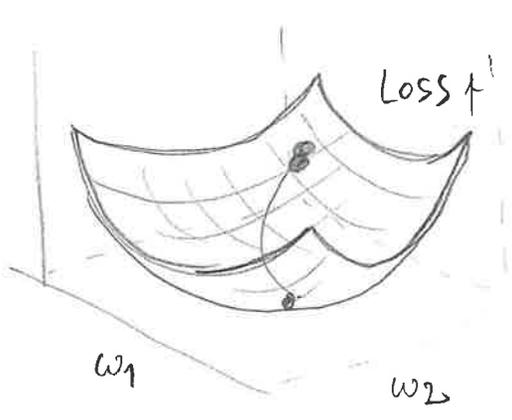
> Key idea: iterative methods (very popular)

## Gradient Descent in 1 dimension:



Step Direction: -gradient
Step Size: -proportional to |gradient|

use $\boxed{\eta} * |\text{gradient}|$

## Gradient Descent in multiple dimensions:



↓ Level Set
  contour lines.

Recall: For a function $f: \mathbb{R}^d \to \mathbb{R}$, a **Level Set** at value $\boxed{c}$ is the set of all points in the domain where $f$ takes the same value
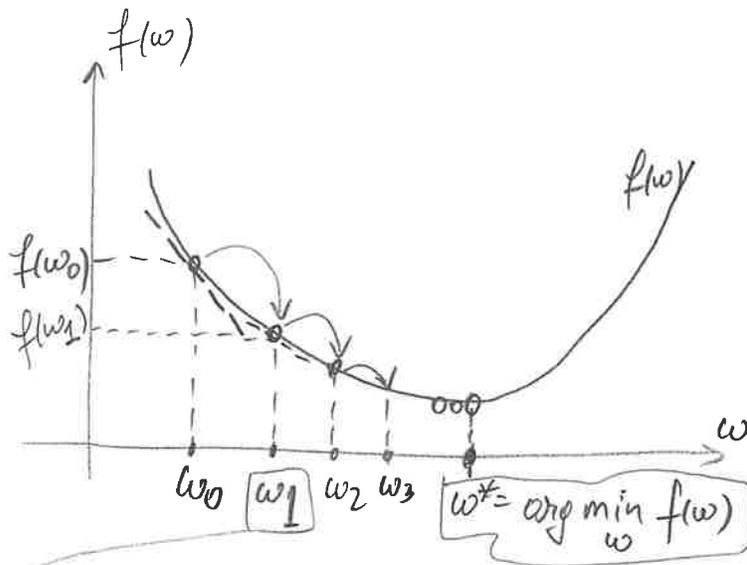
# Algorithm — Gradient Descent (GD)

for $t = 0, 1, 2, \ldots$

$$\omega_{t+1} = \omega_t - \eta \left. \frac{df(\omega)}{d\omega} \right|_{\omega = \omega_t}$$

Hyperparameters:
- initial point $\boxed{\omega_0}$
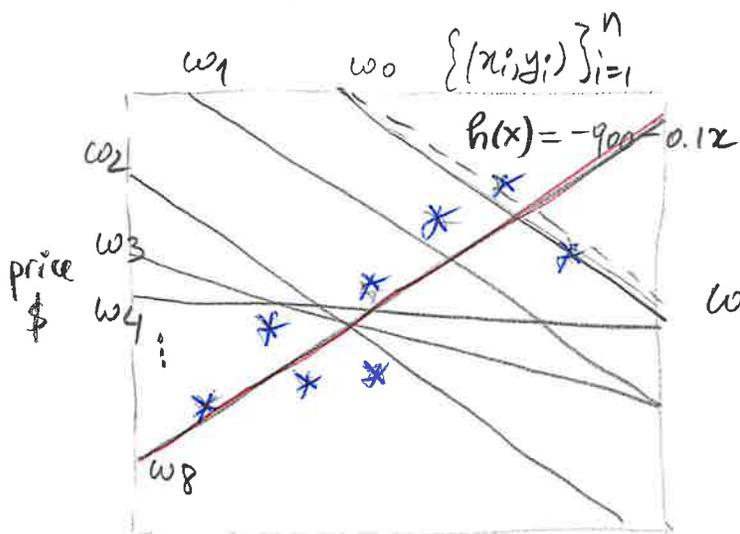- step size $\boxed{\eta}$



$$\omega_1 = \omega_0 - \eta \left. \frac{df(\omega)}{d\omega} \right|_{\omega = \omega_0}$$

$$\omega_2 = \omega_1 - \eta \left. \frac{df(\omega)}{d\omega} \right|_{\omega = \omega_1}$$

$\bullet\bullet\bullet$

$\boxed{\text{Note}}$: As $t \to \infty$, $\left. \frac{df(\omega)}{d\omega} \right|_{\omega = \omega_t} \to 0$

## 1-dimensional Linear Regression with 2 parameters

$$\omega_{t+1} \leftarrow \omega_t - \eta \nabla_\omega f(\omega_t)$$



$\{(x_i, y_i)\}_{i=1}^n$

$h(x) = -900 - 0.1x$

$\omega_{LS}^*$

$\omega[1]$

$\omega_0 \sim N(0, I_{d \times d} \sigma^2)$

$\omega[0]$

price $\$$

Size (feet$^2$)

Evolution of predictor
$$y = \omega[0] + \omega[1] x$$

Gradient Descent dynamics in parameter space $\omega$. Ovals show "level set" of the objective function.

# EXAMPLES:

## ① Quadratic Functions

$$\hat{w} = \arg\min_{w} \underbrace{aw^2 + bw + c}_{f(w) \text{ (loss)}}$$

$$w_0 \sim \mathcal{N}(0, I\sigma^2)$$

$$\frac{df(w)}{dw}\Big|_{w=w_0} = 2aw_0 + b$$

$$\boxed{w_1 = w_0 - \eta(2aw_0 + b)}$$

## ② Linear Regression

$$\hat{w} = \arg\min_{w} \underbrace{\frac{1}{2}\sum_i (y_i - x_i^T w)^2}_{f(w) \text{ (loss/cost)}}$$

$$\nabla_w f(w_0) = X^T(Xw_0 - y) \qquad \leftarrow \text{see Lecture 02 notes.}$$



$$\underbrace{\Box}_{d \times n} \left( \underbrace{\Box}_{n \times d} \underbrace{\Box}_{d \times 1} - \underbrace{\Box}_{n \times 1} \right)$$

$$\boxed{w_{t+1} = w_t - \eta \cdot X^T(Xw_0 - y)}$$

# ③ Lasso

$$\hat{w} = \arg\min_{w} \frac{1}{2}\|y - Xw\|_2^2 + \lambda \|w\|_1$$

$$\|w\|_1 \triangleq \sum_{i=1}^{n} |w_i| \quad \text{is not differentiable etc}$$

$f(w)$ (loss)

$$\boxed{\nabla_w f = X^T(Xw - y) + \lambda \sum_{i=1}^{n} \text{sign}(w_i)}$$

(see notes of Lecture 2)
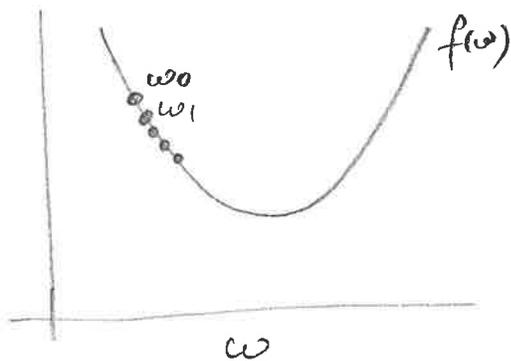
$$w_{t+1} = w_t - \eta \cdot \nabla_w f|_{w = w_t}$$



Derivative for each coordinate:

$$\frac{\partial |w_i|}{\partial w_i} = \begin{cases} +1, & w_i > 0 \\ (-1, 1), & w_i = 0 \\ -1, & w_i < 0 \end{cases}$$

$$f(w) = a w^2 + b w + c + |w|$$

$$\frac{df}{dw} = 2aw + b + \text{sign}(w) = 0$$

$\left\{\begin{array}{l} - \text{can find it } w \text{ is 1-dimensional} \\ - w \in \mathbb{R}^d \quad \text{sign}(w) \to 2^d \text{ possibilities} \end{array}\right.$

## How to choose a step size ?



Step size $\boxed{\eta}$ is too small

=> Slow convergence.



If $\boxed{\eta}$ is too large.

=> Diverging

> In practice: guess and check !?

Training
error
loss

validation
error

loss plateau

training iterative
steps

Train loss effectively minimized, but use |early stopping| to prevent overfitting.

Loss

Loss increasing is a sign that step size $\boxed{\eta}$ might be too large!

Based on **Géron's Book**: (pp. 145)

- Main problem with Batch Gradient Descent is:
  - uses the whole training set to compute the gradients at every step, which makes it very slow when training set is large!

- At the opposite extreme: Stochastic Gradient Descent — picks a random instance in the traing set at every step and computes the gradients based only on that single instance!

- As a compromise: mini-Batch Gradient Descent.

NOTES

- The actual implementations of Batch GD & Stochastic GD algorithms are provided at pp. 143 & 146 in Geron's Book.

- These algorithms are also available in Scikit-Learn as the SGDRegressor class.

- When using GD, you should ensure that all features have similar scale ( e.g, using Scikit-Learn's StandardScaler class), or else it will take longer to converge! [Geron, pp.141]