

Gradient Descent Methods

- basic ingredients of ML:

1) data: \bar{X}

2) model: $\theta \mapsto g_\theta$ θ : model parameters
e.g. weights, biases
(trainable)

3) cost function:

$$\bar{X}, g, \theta \mapsto E(\bar{X}, g_\theta)$$

evaluates how well model works

e.g. $\bar{X} = \{(x_i, y_i)\}_{i=1}^n$, n data points
(e.g. generated from
 $x \mapsto f(x) = y$ unknown)

$g_\theta: x_j \mapsto g_\theta(x_j)$ some parametrized fu.
e.g. $g_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
or $g_\theta(x) = \mathcal{NN}_\theta(x)$

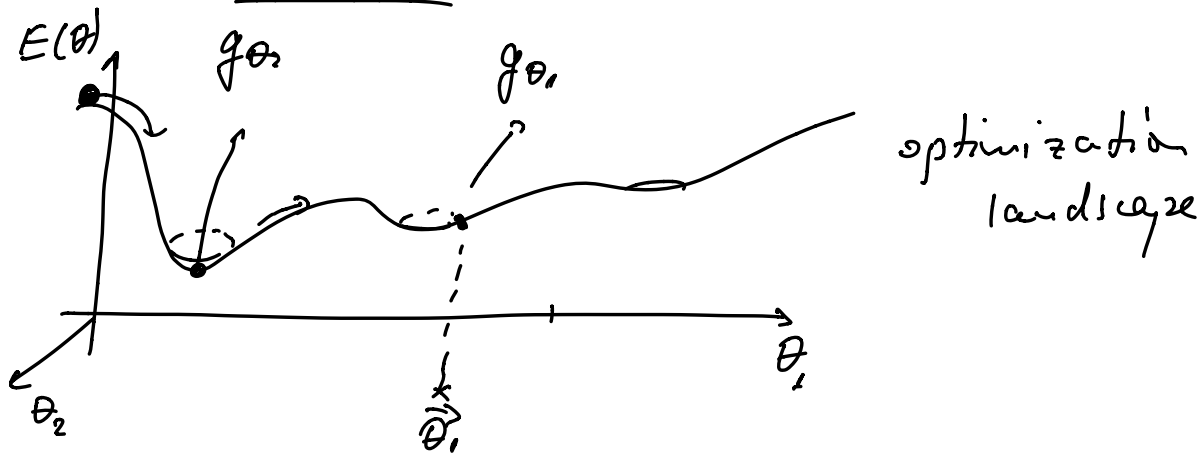
$$E: \bar{X}, g_\theta \mapsto \frac{1}{n} \sum_{j=1}^n [y_j - g_\theta(x_j)]^2$$

- how do we train the model?, i.e.

how do we find optimal values of param. θ ?

→ by minimizing cost function E wrt.
parameters θ .

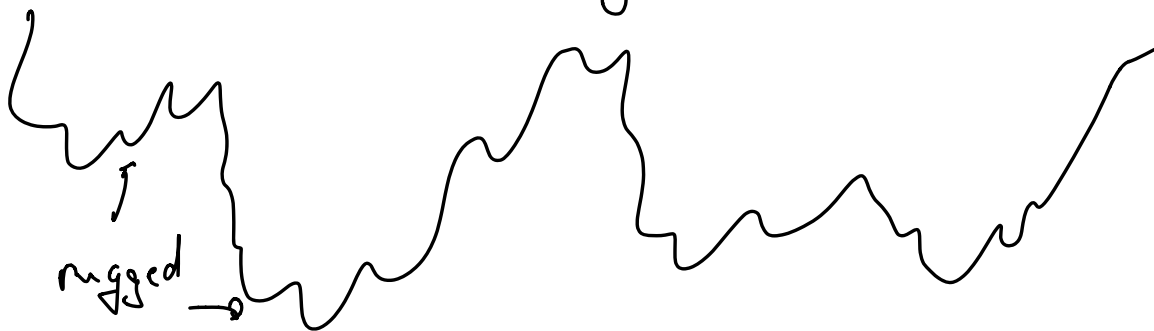
→ use Gradient descent methods.



basic idea: iteratively adjust params θ in directions where gradient of cost fn $\nabla_{\theta} E$ is large & negative

→ params θ "flow" towards a local minimum of cost fn / optimization landscape

— in ML: i) cost fn's are rugged, non-convex landscapes in high-dim. space with many local minima



ii) don't have access over true cost function we want to minimize
→ data is correct wise, uncertainty \mathbb{I}_2 -

→ gradient descent (GD)

initialize: θ_0 at random

want: new parameters θ_1 from old θ_0

$$\underbrace{\theta_{t+1}}_{\substack{\uparrow \\ \text{new} \\ \text{params}}} = \underbrace{\theta_t}_{\substack{\uparrow \\ \text{old} \\ \text{params}}} - \underbrace{\alpha_t}_{\substack{\uparrow \\ \text{learning} \\ \text{rate}}} \underbrace{\nabla_{\theta} E(\theta_t)}_{\substack{\uparrow \\ \text{gradient of} \\ \text{cost fn.}}}$$



→ re-write as

$$g_t := \alpha_t \nabla_{\theta} E(\theta_t)$$

$$\theta_{t+1} = \theta_t - g_t$$

Newton's Method

$$E(\theta - g) \approx \underbrace{E(\theta)}_{\substack{\uparrow \\ \text{expand} \\ \text{for small } g}} - \underbrace{g \cdot \nabla_{\theta} E}_{\substack{\uparrow \\ \text{gradient}}} + \frac{1}{2} g^t \underbrace{H(\theta)}_{\substack{\uparrow \\ \text{Hessian matrix}}} g \quad (*)$$

$$H_{mn} = \partial_{\theta_m} \partial_{\theta_n} E(\theta)$$

assume $\theta - g^*$ is an optimum of $E(\theta)$

$$\nabla_{\theta} E(\theta - g^*) = 0$$

take ∇_f of (*) at $f = \theta - f_*$

\Rightarrow

$$0 = \nabla_f E|_{\theta - f_*} \approx -\nabla_\theta E + H(\theta) f_*$$

$$\Rightarrow f_* = H^{-1}(\theta) \nabla_\theta E(\theta)$$

\Rightarrow Newton method update rule :

$$f_t = H^{-1}(\theta_t) \nabla_\theta E(\theta_t)$$

$$\theta_{t+1} = \theta_t - f_t$$

problems: 1) H may not be invertible

fix \rightarrow regularize: $H^{-1} \rightarrow (H + \epsilon \mathbb{I})^{-1}$
 \uparrow small regularizer

2) computing inverse H^{-1} extremely expensive when # params is large, e.g. $\theta \in \mathbb{R}^{10^4}$

- what do we learn from Newton's method about GD?

• in GD $\alpha_t = \text{const} \rightarrow$ in time t
 \rightarrow for all params θ

in Newton: $\alpha_t = H^{-1}$ is

$\rightarrow \theta$ -dependent, varies different θ

recall : $x \mapsto f(x) = ax^2$ ✓

$\Rightarrow f'' = a$; a^{-1} defines curvature at the extremum

H is generalization of second derivative

\rightarrow inverse eigenvalue define curvature in direction given by corresponding eigenvector

\rightarrow Newton's method adjusts the learning rate \propto proportional to local curvature of $E(\theta)$:

- take larger steps in flat directions (small curvature)

- take small steps in steep directions (high curvature)

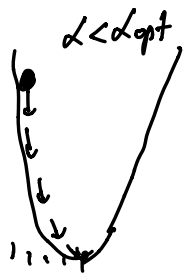
Example : optimal learning rate :

$E(\theta) = \theta^2$ is quadratic, $\theta \in \mathbb{R}$

$$E(\theta - \eta) = E(\theta) - \eta \partial_{\theta} E(\theta) + \frac{1}{2} \eta^2 \partial_{\theta}^2 E \text{ exact!}$$

\Rightarrow optimal learning rate
 $\alpha_{\eta,t} = (\partial_{\theta}^2 E)^{-1}$

→ use GD with const α :



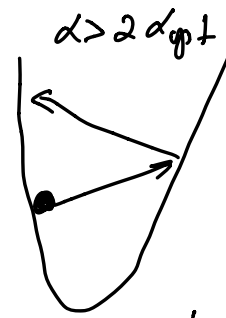
many small steps to reach minimum



single step



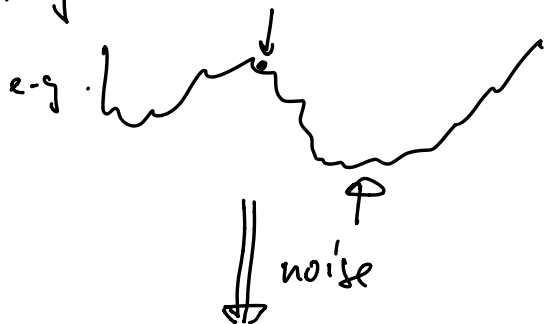
oscillate but reach min



shoot out, GD diverges

— Some disadvantages of GD

1) gets stuck in local minima:



suggestion: introduce small noise to smear out cost function barriers



→ upgrade GD to stochastic GD (SGD)

2) gradients are computationally expensive in large data sets X

→ do we need all data at every step t ?

⇒ use subsets of the data, called minibatches

e.g. $E(\theta) \sim \sum_{j=1}^n [y_j - g_\theta(x_j)]^2$



$$= \sum_{u=1}^M \sum_{j \in B_u} [y_j - g_\theta(x_j)]^2$$

M minibatches of size $|B_u|$

→ drop sum over k

$$E(\theta) \approx \sum_{j \in B_k} [y_j - g_\theta(x_j)]^2$$

→ B_k are determined randomly every step t

⇒ SGD

3) GD is sensitive to choice learning rate α_t

→ consider an adaptive schedule for α_t
(Runge-Kutta methods)

4) GD weighs all directions in parameter space uniformly:

→ make α_t θ -dependent (curvature of landscape)

(⇒ RMSProp, ADAM)

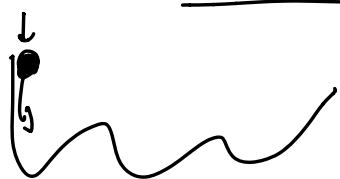
5) GD is sensitive to initial condition θ_0
(true for alternative methods)

6) GD can get stuck in saddle points



saddles are more common than
minima/maxima in higher-dim space

Gradient Descent with Momentum:



- add memory of the
direction (in param space)
we are coming from:

$$\beta_t = \gamma \beta_{t-1} + \alpha_t \nabla_{\theta} E(\theta_t)$$

$$\theta_{t+1} = \theta_t - \beta_t$$

$\gamma \in [0, 1]$: momentum parameter

• interpret β_t as a running average of
recently encountered gradients

- analogy with a viscous particle
in potential energy $E(z)$ described by $z(t)$:

$$m \frac{d^2 z}{dt^2} + \mu \frac{dz}{dt} = - \nabla_z E(z)$$

→ use finite differences

$$\frac{dz}{dt} \approx \frac{z_{t+\Delta t} - z_t}{\Delta t}$$

$$m \frac{z_{t+\Delta t} - z_t - z_t + z_{t-\Delta t}}{(\Delta t)^2} + \mu \frac{z_{t+\Delta t} - z_t}{\Delta t} = - \nabla_z E$$

$$m \frac{\Delta z_{t+\Delta t}}{(\Delta t)^2} - m \frac{\Delta z_t}{(\Delta t)^2} + \mu \frac{\Delta z_{t+\Delta t}}{\Delta t} = - \nabla_z E$$

$$\Rightarrow \Delta z_{t+\Delta t} = \frac{m}{m + \mu \Delta t} \Delta z_t - \frac{(\Delta t)^2}{m + \mu \Delta t} \nabla_z E$$

$$\text{identity: } \gamma = \frac{m}{m + \mu \Delta t} \quad ; \quad \alpha = \frac{(\Delta t)^2}{m + \mu \Delta t}$$

$\Rightarrow \gamma \propto m$

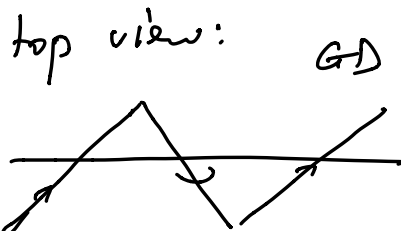
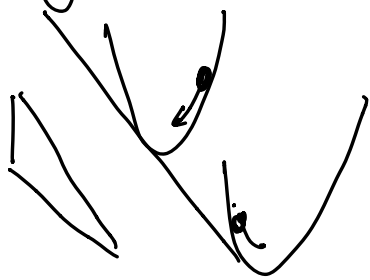
$m \hat{=}$ inertia \rightarrow momentum



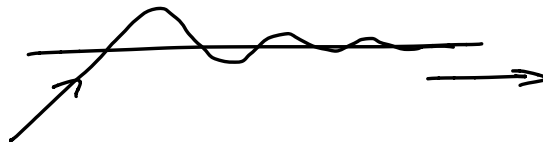
- why is momentum useful?

- gain speed in directions with persistent small gradients

- suppress oscillations in highly curved directions ($\gamma \rightarrow$ average over last few gradients)



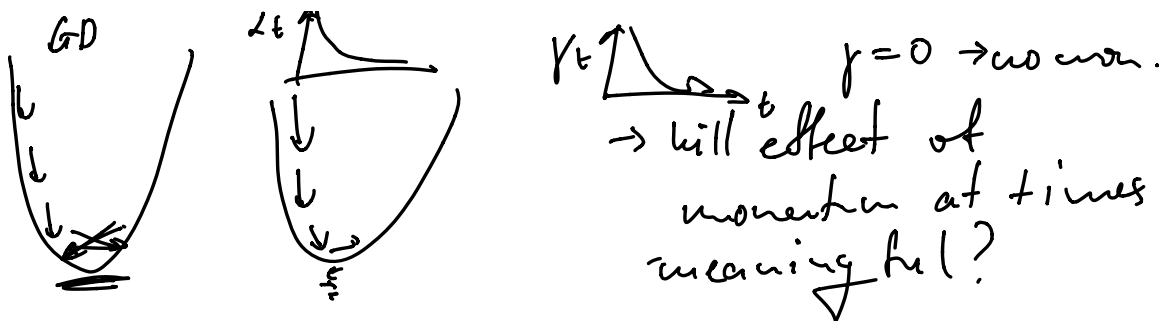
GD + momentum



- useful modification:

Nesterov Accelerated Gradient (NAG)

$$\begin{cases} g_t = \gamma g_{t-1} + \alpha_t \nabla_{\theta} E(\theta_t + \gamma g_{t-1}) \\ \theta_{t+1} = \theta_t - g_t \end{cases}$$



GD methods with second-order moments

-recall: learning rate limited by curvature of optimization landscape / cost fn.

idea: adjust learning rate using Hessian

issue: H^{-1} infeasible to compute if our model g_θ has $\sim 10^4$ params (i.e. $\theta \in \mathbb{R}^{10^4}$)

can we achieve a similar goal w/o H ?

Yes! \rightarrow keep a running estimate of second moment of the gradient

RMSprop:

$$\begin{cases} g_t = \nabla_\theta E(\theta_t) \\ v_t = \beta v_{t-1} + (1-\beta) g_t^2, \quad v_0 = 0 \\ \theta_{t+1} = \theta_t - \alpha_t \frac{g_t}{\sqrt{v_t + \epsilon}} \end{cases}$$

- $\beta \in [0, 1]$: controls averaging time of second moment of g_t
- ε : small regularizer ($\varepsilon \sim 10^{-7}$)
- α_t : learning rate schedule
- multiplication/division of two vectors is understood element-wise: $\frac{\vec{a}}{\vec{b}} = \left(\frac{a_i}{b_i} \right)$

→ why V_t is a running average?

$$\begin{aligned} V_t &= \beta V_{t-1} + (1-\beta) g_t^2 \\ &= \beta (\beta V_{t-2} + (1-\beta) g_{t-1}^2) + (1-\beta) g_t^2 \\ &= \dots \end{aligned}$$

$$= (1-\beta) \sum_{j=1}^t \beta^{t-j} g_j^2$$

↳ decay rate / (inv.) averaging time

- suppose $E(\theta)$ is noisy (e.g. data X is noisy)
→ can take an expectation over the noise \mathbb{E}

$$V_t = (1-\beta) \sum_{j=1}^t \beta^j g_j^2 / \mathbb{E}[\cdot]$$

$$\begin{aligned} \mathbb{E}[V_t] &= (1-\beta) \sum_{j=1}^t \beta^j \underbrace{\mathbb{E}[g_j^2]}_{\approx \mathbb{E}[g^2] \text{ indep. of } j} \\ &\quad - // - \end{aligned}$$

$$\approx \cancel{E[g^2] (1-\beta)} \underbrace{\sum_{j=1}^t \beta^j}_{\substack{= \frac{(1-\beta^t)}{(1-\beta)} \\ \nearrow \beta < 1}} = E[g^2] \underbrace{(1-\beta^t)}_{\substack{\uparrow \\ \text{creates} \\ \text{bias at} \\ \text{initial steps } t \\ \beta^t \xrightarrow{t \rightarrow \infty} 0}}$$

ADAM : (bias-corrected RMSProp)

$$g_t = \nabla_{\theta} E(\theta)$$

$$\left. \begin{aligned} m_t &= \beta_1 m_{t-1} + (1-\beta_1) g_t \\ v_t &= \beta_2 v_{t-1} + (1-\beta_2) g_t^2 \end{aligned} \right\} \text{biased}, \quad \begin{aligned} m_0 &= 0 \\ v_0 &= 0 \end{aligned}$$

$$\left. \begin{aligned} \hat{m}_t &= \frac{m_t}{1-\beta_1^t} \\ \hat{v}_t &= \frac{v_t}{1-\beta_2^t} \end{aligned} \right\} \text{correct bias}$$

$$\theta_t = \theta_{t-1} - \alpha_t \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}$$

β_1, β_2 : decay rates of first & second moment of the gradient estimate

$$\beta_j^t = \underbrace{\beta_j \cdots \beta_j}_{t \text{ times}}$$

ϵ : small regularizer

d_t : learning rate schedule

mult./division : element-wise operations