

Lecture 16 - SGD and GA

Stochastic Gradient decent (SGD):

- In simulated annealing if your temperature is zero you can only go up, never go down.
- You still have to explore the neighborhood to find out which direction is up.
- If we know more about the landscape we can directly go up/down.

Steepest decent: use the derivative to take steps into the right direction.

$$\Rightarrow \underset{\theta}{\operatorname{argmin}} \quad J(\theta)$$

$$\Rightarrow \text{take the derivative } \nabla J(\theta)$$

$$\Rightarrow \theta^{(i+1)} = \theta^i - \alpha \nabla J(\theta)$$

↳ uses specified stepsize

• In many cases $J(\theta)$ has the form

$$J(\theta) = \sum_{j=1}^N J_j(\theta)$$

$$\Rightarrow \theta^{i+1} = \theta^i - \alpha \sum_{j=1}^N \nabla J_j(\theta) \quad \text{each data point adds a derivative to this sum.}$$

stochastic gradient decent: update after only one element of the sum at a time

$$\theta^{i+1} = \theta^i - \alpha \nabla J_j(\theta)$$

↳ we take $\nabla J_j(\theta)$ as an approximation for $\nabla J(\theta)$

$\Rightarrow \theta^{i+1}$ might go into the "wrong" direction

\Rightarrow In the long run we should still be going down.

SGD:

- initialize Θ
- while not done:
 - randomly reshuffle all data
 - For $j=1 \dots M$

$$\Theta^{\text{new}} = \Theta^{\text{old}} - \alpha \nabla J_j(\Theta)$$

\Rightarrow Faster trick: pick random j , do update, pick next random j .

\Rightarrow Can also do subset of data points (batch SGD)

Example: Linear Regression

$$J(\Theta) = \frac{1}{2} \sum_{i=1}^N (f_{\Theta}(x^{(i)}) - y^{(i)})^2$$

with $f_{\Theta}(x_i) = \Theta^T x_i$

$$\Rightarrow \Theta^{\text{new}} = \Theta^{\text{old}} + \alpha (f_{\Theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

α : How much do you believe the local approximation of the gradient leads to a minimum?

\Rightarrow You can keep memory of the best solution so far.

\Rightarrow last iteration might not be the best

\Rightarrow might have to do multiple random restarts.

\Rightarrow Can also use momentum:

$$V_{t+1} = \mu V_t - \alpha \nabla J(\Theta)$$

$$\Theta_{t+1} = \Theta_t + V_{t+1}$$

\Rightarrow keep memory of previous gradient

\Rightarrow helps with shallow ridges,

Genetic Algorithms:

- Motivated by Biology and evolution.

- We have:

- a population of individuals
- measure of fitness (optimality)
- operations for change: reproduction

crossover

mutation

- We change the population to gradually higher average fitness.

Reproduction:

- fitter individuals have more offspring
- population size remains constant
- less fit individuals die

Crossover: choose a random sized part of two individuals and swap:

$$\begin{bmatrix} A \\ B \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} \Rightarrow \begin{bmatrix} A \\ b \end{bmatrix} \begin{bmatrix} a \\ B \end{bmatrix}$$

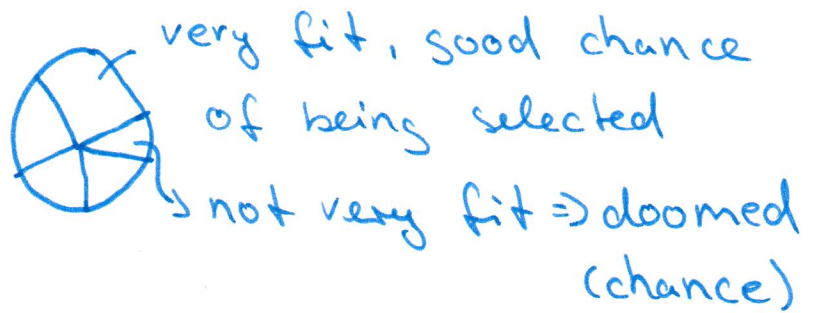
chance for crossover is typically high
($0.7 \leq P_c \leq 1$).

Mutation: randomly change part of an individual
 P_m is typically small ~ 0.001

- GA:
- compute fitness score of individuals
 - Select survivors/reproduction candidates according to fitness level.
 - apply crossover and mutations
 - repeat all steps for this new generation.

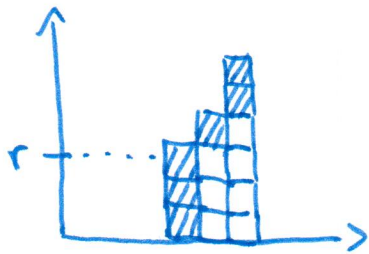
Selecting Survivors according to fitness level:

Idea: Fortune wheel



Example: 3 1 2 \Rightarrow fitness

3 4 6 \Rightarrow cumulative sum



$$F = \sum_i^n f_i$$

\hookrightarrow individual fitness

scaled fitness: $\hat{f}_i = \frac{f_i}{F}$

incremental scaled fitness: $\hat{g}_m = \sum_{i=1}^m \hat{f}_i$

$$\Rightarrow \hat{g}_N = 1$$