# Chapter 3: Stochastic approximations

or stochastic iterative algorithms

## 3.1. Introduction

- $x^*$ = solution to some problem
- Examples:
  - Solve $f(x) = a$

    $\min f(x)$ — optimization

    $A v_i = \lambda_i v_i$ — Eigenvalues/vectors
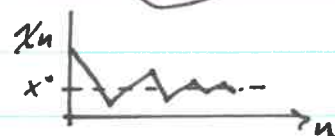
    RL: find policy optimizing reward

- Goal: Calculate/find/estimate $x^*$
- Approach 1: Deterministic recursion ← iteration ~ iteration

  $$x_{n+1} = T(x_n) = x_n + F(x_n)$$

  ↳ map, transfer fct

  - $x_n \xrightarrow{n \to \infty} x^*$ from $x_1$
  - $x^*$ attracting fixed pt of $T$



- Approach 2: Stochastic recursion (iteration)

  $$X_{n+1} = X_n + a_n Y_n$$

  $Y_n$ function of $X_n$ related to $X_n$

  $$\text{"} = X_n + a_n \underbrace{F(X_n, C_n)}_{Y_n}$$

  $P(Y_n | X_n)$

  $$\text{"} = X_n + a_n \underset{\substack{\text{deterministic} \\ \text{map}}}{H(X_n)} + b_n \underset{\text{noise}}{M_n}$$

  - Find $a_n, b_n \downarrow 0$ as $n \to \infty$ such that

    $$\underset{RV}{X_n} \xrightarrow{n \to \infty} \underset{\text{value/constant}}{x^*} \quad \text{in probability}$$

  - $\lim\limits_{n \to \infty} P(|X_n - x^*| > \varepsilon) = 0$

  - Stochastic recursion = non-homogeneous Markov chain

    Why?          why?
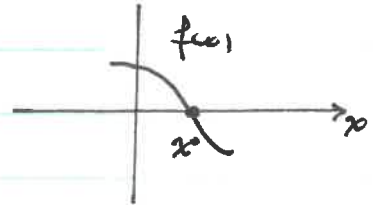
- Applications:
  - Stochastic gradient descent
  - Stochastic annealing
  - Reinforcement learning
  - etc.

## 3.2. Finding zero

- $f : \mathbb{R}^d \to \mathbb{R}^d$
- Solve $f(x) = a$
  - Take $a = 0$ w/o loss generality
  - Can have one, no, or multiple solutions
  - Solution: $x^* \ni f(x^*) = 0$



### 3.2.1 Deterministic recursion

$$\begin{cases} x_{n+1} = T(x_n) & \text{iteration / map / transfer function} \\ x_1 = x & \text{initial value} \end{cases}$$

- Fixed-point condition: $T(x^*) = x^*$   ✓ num
- Contraction condition: $\| DT(x^*) \| < 1$
  
  $\underbrace{\qquad}_{\text{Jacobian}}$

- In $\mathbb{R}$: $|T'(x^*)| < 1$
- Convergence: $X_n \to x^*$ as $n \to \infty$

Example: Newton-Raphson map: $T(x) = x - \dfrac{f(x)}{f'(x)}$

$$T(x^*) = x^* - \frac{\overset{=0}{f(x^*)}}{f'(x^*)} = x^*$$
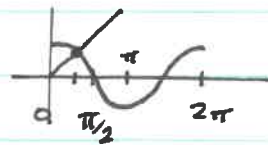
Example: $\cos(x) = x$

$\Rightarrow f(x) = \cos(x) - x = 0$



$T_1(x) = \cos(x)$

$T_2(x) = x - \dfrac{f(x)}{f'(x)} = x + \dfrac{\cos(x) - x}{\sin(x) + 1}$   N-R   See demo

- Iteration: $X_1 \to X_2 \to \cdots \to X_n$   deterministic from $X_1 = x$



  - Stop when $|X_n - x^*| < \varepsilon$
  - In practice, stop when $|X_n - X_{n-1}| < \varepsilon$   no improvement
  - If many solutions: try different initial point

- Rem: $X_{n+1} = T(X_n) = X_n + F(X_n)$, $F(x^*) = 0$

## 3.2.2 Stochastic recursion
Ref: Robbins - Monro 1951

- Solve $f(x) = d$
- $f(x)$ not known / given exactly
- Estimate: $y$

$$E[Y|x] = \sum_y y \, P(y|x) = f(x)$$

$x \to \boxed{\pm} \to f(x)$

$x \to \square \to Y$

random variable

- Unbiased estimate of $f(x)$
- Unbiased function call

- Example: Additive noise model $\quad Y = f(x) + \xi \qquad E[\xi] = 0$

$$E[Y|x] = f(x) + E[\xi] = f(x)$$

- Recursion: $\quad X_{n+1} = X_n - a_n (Y_n - \alpha) \qquad X_1 = x_1$

initial value

- $Y_n$: fct call for $X_n$ = estimate of $f(X_n)$
- $a_n$: annealing sequence

$$a_n \searrow 0 \quad as \quad n \to \infty$$

$$\sum_{n=1}^{\infty} a_n = \infty \qquad \sum_{n=1}^{\infty} a_n^2 < \infty$$

- Convergence: $X_n \xrightarrow{n \to \infty} x^*$ in probability

$$\lim_{n \to \infty} P(|X_n - x^*| > \varepsilon) = 0$$

- In practice: $a_n = \frac{1}{n}$

- Iteration: $\quad X_1 \to X_2 \to \cdots \to X_n$

- Markov chain (non-homogeneous)
- Probably almost correct (PAC): Stop at $X_n \ni$

$$P(|X_n - x^*| \geq \varepsilon) < \delta \qquad \forall n \geq n_0$$

- In practice: Stop when $|X_n - X_{n-1}| < \varepsilon$ no improvement in solution

Example: $f(x) = \cos(x) + x \quad \rightarrow \quad Y = \cos(x) + x + \mathcal{U}[-1,1]$
See demo
additive noise

Example: Mean estimator

$$S_n = \frac{1}{n} \sum_{i=1}^{n} X_i \qquad X_i \sim P$$

$$S_{n+1} = \frac{n}{n+1} S_n + \frac{X_{n+1}}{n+1}$$

$$= S_n - \frac{1}{n+1} \left( S_n - X_{n+1} \right)$$

$$= S_n - a_n Y_n$$

$$E[Y_n | S_n = s] = E[S_n - X_n | S_n = s]$$
$$= s - E[X_{n+1}]$$
$$= s - \mu \quad = f(s)$$

$$f(s) = 0 \quad \Rightarrow \quad s^* = \mu \quad \Rightarrow \quad S_n \rightarrow \mu \quad \text{in probability as}$$
$$n \rightarrow \infty \quad (LLN)$$

Rem: Linear + additive noise models

$$X_{n+1} = X_n - a_n \left( \underbrace{X_n + Z_n}_{Y_n} \right)$$
noise
noisy version of $f(x) = x$

$$= \underbrace{(1 - a_n) X_n}_{\nearrow 1} - \underbrace{a_n Z_n}_{\searrow 0}$$

reinforce update
exploitation

random update/change
exploration

Rem: | Before | Now |
|---|---|
| $Y$ = noisy obs. of $f(x)$ | Put noise in $f(x)$ |
| inherent noise | explicit noise |
| noise is bad | noise is good $\longleftarrow$ exploration |

## 3.3. Optimization

- Potential / cost / loss : $V : \mathcal{X} \to \mathbb{R}$
- State / solution space : $\mathcal{X} = \mathbb{R}^d$ or discrete space    ✓ See Sec. 3.4
- Minimization problem :

$$\min_{x \in D} V(x)$$



- $D \subseteq \mathcal{X}$  Constraint set
- Assume unique solution (fn now) : $x^* = \underset{x \in D}{\arg\min}\ V(x)$

- Examples :   · MLE
                · Neural net training

- Rem : $\mathcal{X} = \mathbb{R}^d$
  - Gradient : $\nabla V(x) = \begin{pmatrix} \frac{\partial}{\partial x_1} V(x_1 \dots x_n) \\ \vdots \\ \frac{\partial}{\partial x_n} V(x_1 \dots x_n) \end{pmatrix} = \begin{pmatrix} \partial_1 V \\ \vdots \\ \partial_n V \end{pmatrix}$

  - Hessian matrix : $\operatorname{Hess} V(x) = \dfrac{\partial^2 V(x)}{\partial x_i\, \partial x_j}$

  - Critical point : $\nabla V(x^*) = 0$    n equations of n variables
    - Min (local or global) :  $\operatorname{Hess} V(x)$  positive definite
    - Max (  "    "    "  ) :  $\operatorname{Hess} V(x)$  negative def
    - Saddle point :  positive and negative eigenvalues

  - In $\mathbb{R}^1$ :   $V'(x^*) = 0$
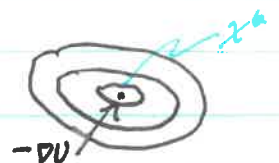    - Min :  $V''(x^*) > 0$
    - Max :  $V''(x^*) < 0$

  - $\nabla V(x) = $ direction of greatest ~~descent~~ ascent at $x$
    $$V_{\vec{d}}'(x) = \nabla V(x) \cdot \vec{d}$$
    $$\| V_{\vec{d}}'(x) \| = \| \nabla V(x) \| \, \| \vec{d} \| \cos\theta$$
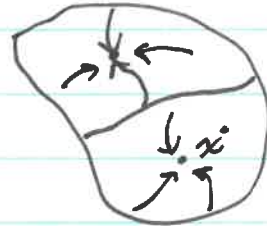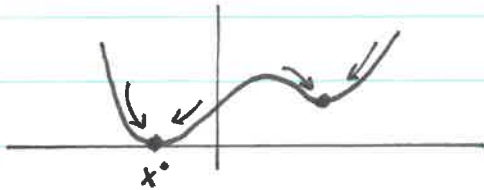    max at $\theta = 0$   i.e.  $\vec{d} \parallel \nabla V$

## 3.3.1 Gradient descent
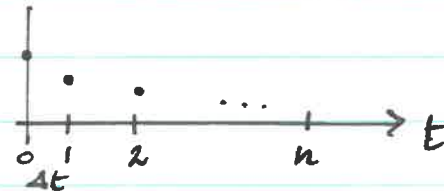
- Iteration: $x_{n+1} = x_n - \gamma \nabla V(x_n)$

  $x_{n=1} = x$     initial value/seed



- Fixed point: $\nabla V(x^*) = 0$
- Convergence: $x_n \to x^*$ if $x_1$ in basin of attraction of $x^*$
- $V(x)$ can have multiple global/local min
- Convergence to global min not guaranteed



- Continuous-time limit:



$$x_{n+1} = x_n - \gamma \nabla V(x_n) \qquad x_0 \text{ initial value}$$

$$\downarrow$$

$$x(t+\Delta t) = x(t) - \tilde{\gamma} \nabla V(x(t)) \Delta t \qquad \gamma = \tilde{\gamma}\Delta t$$

$$\frac{x(t+\Delta t) - x(t)}{\Delta t} = -\tilde{\gamma} \nabla V(x(t))$$

$$\Rightarrow \quad \dot{x}(t) = \frac{d}{dt} x(t) = -\tilde{\gamma} \nabla V(x(t)) \qquad \text{Gradient ODE}$$

$$x(t) \to x^* \quad \text{as } t \to \infty \quad \text{if } x(0) \text{ in basin of attraction of } x^*$$

- Example: $V(x) = \dfrac{x^2}{2} \qquad x^* = 0$

$$\dot{x}(t) = -\gamma V'(x(t)) \qquad \to \quad x(t) = x(0) e^{-\gamma t}$$
$$= -\gamma x(t)$$

Converges exponentially to $x^* = 0$ from any $x(0)$

· Other gradient dynamics

1- Newton - Raphson :

$$V'(x^*) = f(x^*) = 0$$

$$\begin{aligned} x_{n+1} &= T(x_n) \\ &= x_n - \frac{V'(x_n)}{V''(x_n)} \\ &= x_n - \gamma_n V'(x_n) \end{aligned}$$

$$T(x) = x - \frac{f(x)}{f'(x)}$$

$$\gamma_n = V''(x_n)^{-1}$$

adjusted learning rate

2- Gradient with momentum :

$$p_{i+1} = -\varepsilon \nabla V(x_i) + (1-\varepsilon \gamma) p_i$$
$$x_{i+1} = x_i + \varepsilon p_{i+1}$$



· Heavy ball descent with "oscillations"

· Continuous limit :

friction force        potential force

$$\dot{p}(t) = -a\, p(t) - b\, \nabla V(x(t))$$
$$\dot{x}(t) = c\, p(t)$$

momentum

· Comes from Newton's law :   $F = ma = m\ddot{x}$

$$\Rightarrow \ddot{x}(t) = F/m \rightarrow \begin{cases} \dot{x} = v \\ \dot{v} = \ddot{x} = F/m \end{cases}$$

2nd order        2 first-order eqs

· Adam (adagrad) :

$$M_{i+1} = \beta_1 m_i + (1-\beta_1) \nabla V(\theta_i)$$
$$V_{i+1} = \beta_2 v_i + (1-\beta_2) \nabla V(\theta_i)^2$$

$$\hat{m}_{i+1} = \frac{M_{i+1}}{1-\beta_1^{i+1}}$$

$$\hat{v}_{i+1} = \frac{v_{i+1}}{1-\beta_2^{i+1}}$$

$$\theta_{i+1} = \theta_i - \alpha \frac{\hat{m}_{i+1}}{\sqrt{\hat{v}_{i+1}} + \varepsilon}$$

2 momenta
reinforced

### 3.3.2 Stochastic gradient descent (SGD)

$$X_{n+1} = X_n - a_n G_n \qquad\qquad X_1 = x$$

$a_n$ = learning rate, $x$ = initial value

· Gradient estimate:

$$E[G|x] = \sum_g g\, P(g|x) = \nabla V(x)$$

$$E[G_n|x] = E[G_n|X_n = x] = \nabla V(x)$$

· Example: Additive noise: $\quad G_n = \nabla V(X_n) + \xi_n \quad$ (noise) $\quad E[\xi_n] = 0$

$$E[G_n|X_n = x] = \nabla V(x) + E[\xi_n]$$
$$= \nabla V(x)$$

· Convergence conditions (sufficient):

$$\sum_{n=1}^{\infty} a_n = \infty \qquad \sum_{n=1}^{\infty} a_n^2 < \infty$$

· $X_n \to x^*$ as $n \to \infty$ in probability
· $\lim_{n\to\infty} P(|X_n - x^*| > \varepsilon) = 0$ 
  
  $X_n$ : Markov chain non-homogeneous

· Example: $a_n = 1/n$

· Example: Kiefer-Wolfowitz 1952

$$X_{n+1} = X_n - a_n \left( \frac{Y_n^+ - Y_n^-}{C_n} \right) \leftarrow \text{discrete derivative}$$

$$Y_n^+ \sim P(\cdot \mid X_n + C_n)$$
$$Y_n^- \sim P(\cdot \mid X_n - C_n)$$

$E[Y|x] = V(x)$

estimate of $V(x + \Delta x)$
estimate of $V(x - \Delta x)$

$$\sum a_n = \infty \qquad \sum a_n C_n < \infty \qquad \sum a_n^2 C_n^{-2} < \infty$$

· In practice: $C_n = \Delta x$ fixed
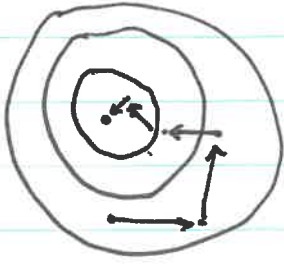$$Y_n^+ \sim P(\cdot \mid X_n + \Delta x)$$
$$Y_n^- \sim P(\cdot \mid X_n - \Delta x)$$
$$\sum a_n = \infty \qquad \sum a_n^2 < \infty$$

· Noisy estimation of $V(x)$
· $\nabla V$ by discrete derivative

· Example: Random descent



$$\nabla V = \begin{pmatrix} \partial_1 V \\ \partial_2 V \\ \vdots \\ \partial_n V \end{pmatrix}$$

direction of largest ~~descent~~ ascent

· Random directional derivative: $G = \nabla_{\vec{d}} V = \nabla V \cdot \vec{d}$

$$E[G|x] = \nabla V \qquad \text{unbiased gradient}$$

· In $\mathbb{R}^2$: $\quad \nabla_x V = \partial_1 V \qquad \vec{d} \sim \mathcal{U}\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\}$
$\qquad\qquad\qquad \nabla_y V = \partial_2 V$

$$E[G|x] = \sum_i p_i \, \partial_i V = \frac{1}{2} \nabla V(x) \propto \nabla V(x)$$

· In $\mathbb{R}^d$: · Uniform directions
· Any uniform subsets of $m \leq d$ directions
· Unbiased over directions = no preferred direction

· Example: Drop out
· Neural network (model) with $d$ parameters
· Loss: $L(\theta)$
· Choose subset of parameters randomly $\{\theta\}_{chosen}$
· Gradient estimate:
$$\nabla_\theta L \approx \nabla_{\{\theta\}_{chosen}} L$$
↳ same as dropping parameters

· Unbiased: $E[\nabla_{\{\theta\}_{chosen}} L] = \nabla_\theta L$

· Rem: Mini batch optimization
$$L(\theta) = \frac{1}{n} \sum_{i=1}^n C_i(\theta) \stackrel{e.g.}{=} \frac{1}{n} \sum_{i=1}^n |y_i - f(x_i, \theta)|^2$$

$$\nabla_{\theta_i} L(\theta) \approx \frac{1}{m} \sum_{j=1}^m \nabla_{\theta_i} C_j(\theta)$$

estimate loss on random subset of data
$m \leq n$

### 3.3.3 Langevin dynamics

- Gradient dynamics:

$$\dot{x}(t) = -\gamma \nabla V(x(t)), \qquad x(0) = x \quad \text{initial value}$$

- SGD: $\dot{X}(t) = -\gamma \nabla V(X(t)) + \sigma \underbrace{\xi(t)}_{\substack{\text{noise} \\ \text{noise amplitude}}}$

$$\underbrace{\hspace{4cm}}_{\text{noisy gradient}}$$

- Stochastic differential equation (SDE):

$$dX(t) = -\gamma \nabla V(X(t)) \, dt + \sigma \, dW(t)$$

$$\dot{X}(t) = \frac{X(t+dt) - X(t)}{dt} = -\gamma \nabla V(X(t)) + \sigma \xi(t)$$

$$X(t+dt) - X(t) = -\gamma \nabla V(X(t)) dt + \underbrace{\sigma \xi(t) dt}_{dW(t)}$$

$$\Rightarrow \quad X(t+dt) = X(t) - \gamma \nabla V(X(t)) dt + \sigma \, \Delta W(t)$$

- Gaussian white noise: $\Delta W(t) \sim \mathcal{N}(0, \Delta t)$

$$= \sqrt{\Delta t} \, Z, \qquad Z \sim \mathcal{N}(0,1)$$

$$\Rightarrow \quad X_{n+1} = X_n - \gamma \nabla V(X_n) \Delta t + \sigma \sqrt{\Delta t} \, Z \qquad \text{Euler–Maruyama Scheme}$$

See CW3

- Stationary distribution:

$$P(x,t) \longrightarrow p^+(x) = \frac{e^{-2\gamma V(x)/\sigma^2}}{Z} \qquad \text{Gibbs density}$$

$$= \frac{e^{-\beta V(x)}}{Z} \qquad \beta = \frac{2\gamma}{\sigma^2}$$

- Annealing: Decrease $\sigma$ in time:

$$dX(t) = -\nabla V(X(t)) dt + \sigma_t \, dW(t)$$

$$\sigma_t \searrow 0 \quad \text{as} \quad t \to \infty$$

See simulated annealing
See CW3

3.4 Simulated annealing

· Potential / cost / loss : $V : \mathcal{X} \to \mathbb{R}$

· Minimization problem : $\min_{x \in D} V(x)$

  · $\nabla V$ exists $\to$ use GD or SGD
  · $\nabla V$ doesn't exist (e.g. $\mathcal{X}$ discrete) ?

· Rem : · Exhaustive search : $O(|\mathcal{X}|)$
  · Random search : Needs structure / guiding
  · Use $V(x)$ in search : Explore : $x \to x'$
    Reinforce : Accept if $\Delta V < 0$

· Gibbs distribution : $P_T(x) = \dfrac{e^{-V(x)/T}}{Z_T}$

  · Partition function / normalization : $Z_T = \sum_x e^{-V(x)/T}$

    $n = \int dx \, e^{-V(x)/T}$

  · $T$ = temperature (usually $> 0$)
  · Inverse temperature : $\beta = T^{-1}$    $T \to 0^+$    $\beta \to \infty$

  $$P_\beta(x) = \dfrac{e^{-\beta V(x)}}{Z_\beta} \qquad Z_\beta = \sum_x e^{-\beta V(x)}$$

· Laplace principle : $P_T(\cdot)$ concentrates on $\min V(x)$ as $T \to 0$.

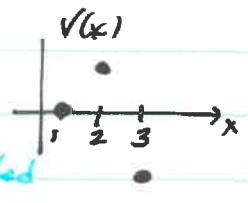· Example : $\mathcal{X} = \{1, 2, 3\}$    $V(1) = 0$, $V(2) = 1$, $V(3) = -1$
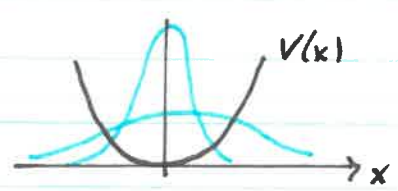


  $P_T \to (0 \quad 0 \quad 1)$    $T \to 0$    $(\beta \to \infty)$ peaked

  $P_T \to (\frac{1}{3} \quad \frac{1}{3} \quad \frac{1}{3})$    $T \to \infty$    $(\beta \to 0)$ uniform

· Example : $\mathcal{X} = \mathbb{R}$    $V(x) = \dfrac{x^2}{2}$    $P_T(x) = \dfrac{e^{-x^2/2T}}{\sqrt{2\pi T}}$



  $\mathrm{Var}(X) = T \downarrow 0$ as $T \to 0$

- Simulated annealing (SA) algorithm:
  - Sample $P_T$ with low $T \to$ samples around $x^*$
  - Metropolis - Hastings with $P_T$
  - Time-dependent $T \to T_n$
  - Steps:
  1 - $X_1 = x$, $T_1 = T$ initial value
  2 - Proposal: $X \to x'$    Metropolis or MH   symmetric   non-symmetric
  3 - Accept with prob

$$\rho = \min \left\{ 1, \frac{P_{T_1}(x')}{P_{T_1}(x)} \right\} = e^{-\Delta V / T_1}$$

$$\text{if } u[0,1] < \rho:$$
$$X_2 = x'$$
$$\text{else}$$
$$X_2 = x$$

  4 - Annealing: $T_1 \to T_2$ decrease
  5 - Repeat

- Rem:
  - Time-dependent Markov chain
  - Semi-greedy: $\Delta V \leq 0$   $x'$ accepted for sure   exploit
  -    $\Delta V > 0$   $x'$   "   with prob. $\rho$   explore

- Annealing schedule: $T_n \downarrow 0$ as $n \to \infty$
  - Decrease too fast: $X_n \not\to x^*$   get stuck in local min
  - "   "   slow: $X_n$ too noisy
  - Log schedule: $T_n = \dfrac{T_1}{\log n + 1}$   can be slow
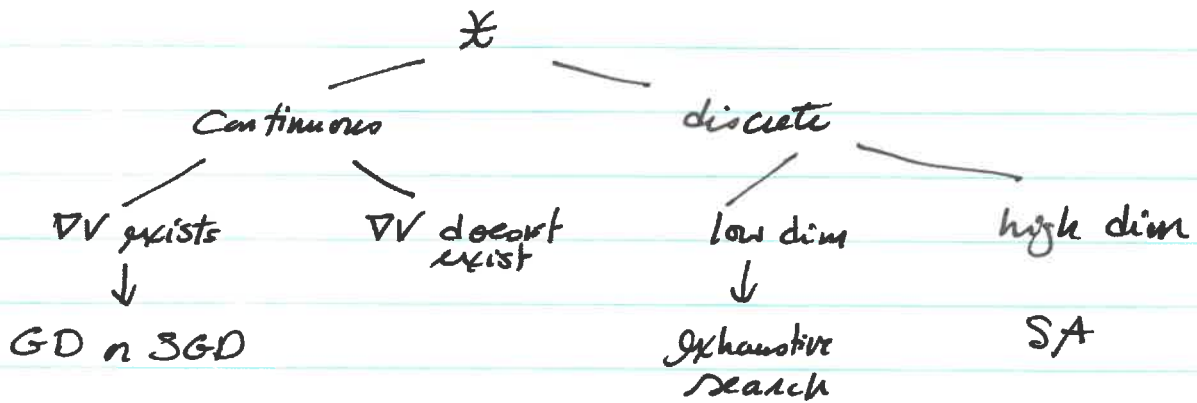  - Geometric schedule: $T_n = \dfrac{T_1}{k^n}$   $k > 1$   no convergence guarantee Simulated quench

$\beta_n = \beta_1 (\log n + 1)$
$\beta_n = \beta_1 \, k^n$
See CW3

## 3.5 Remarks on optimization

- Potential: $V: \mathcal{X} \to \mathbb{R}$      $V(x)$
- Minimization: $\min\limits_{x \in D} V(x)$

```
                        𝒳
              ╱                    ╲
         Continuous            discrete
          ╱      ╲              ╱        ╲
   ∇V exists   ∇V doesn't    low dim    high dim
       ↓        exist          ↓
    GD ∩ SGD              Exhaustive      SA
                            search
```

- Ultimate goal: $\text{cost}(\text{optimization}) \sim \text{cost}(\text{simulation})$
- Why use randomness/noise in optimization?
  - $V(x)$ rugged, many mins
  - $V(x)$ non convex



- Neural network training: $L(D, \theta)$
             data   parameters

  - $D$ only a sample of "true" underlying distribution
  - Compute $\nabla_\theta L$ over subset of parameters    drop out
  - Estimate $L$ on subset of $D$    mini batch

  → Estimation of "true" $L, \nabla_\theta L$
  "Noisy" estimate of $L, \nabla_\theta L$

- Rem: Why not training by solving $\nabla_\theta L = 0$ ?
      "    "    "   using MCMC/SA ?