

Machine Intelligence 2

3 Stochastic Optimization

Prof. Dr. Klaus Obermayer

Fachgebiet Neuronale Informationsverarbeitung (NI)

SS 2018

Stochastic Optimization

Simulated Annealing

Mean-Field Annealing

Stochastic optimization

Supervised & unsupervised learning \rightarrow evaluation of cost function E^T

- real-valued arguments: gradient based techniques (e.g. ICA weights)
- discrete arguments: ?? (e.g. for cluster assignment)

\Rightarrow simulated annealing

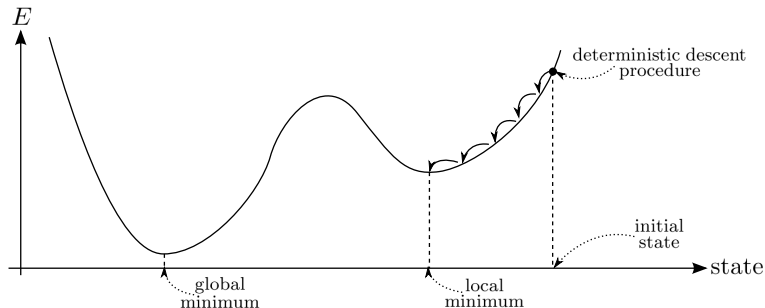
Setting

- discrete variables $s_i, i = 1, \dots, N$ (e.g. $s_i \in \{+1, -1\}$ or $s_i \in \mathbb{N}$)
- short-hand notation: \underline{s} ("state") – often $\{\underline{s}\}$ not a vector space (but called state space)
- **cost function:** $E : \underline{s} \mapsto E_{(\underline{s})} \in \mathbb{R}$ – not restricted to learning problems

Goal: find state \underline{s}^* , such that:

$$E \stackrel{!}{=} \min \quad (\text{desirable global minimum of } E)$$

Optimizing cost functions with local optima



- Deterministic descent may converge to local minima
- Grid-search, random search, multiple initializations
 \leadsto *Simulated Annealing*

Simulated Annealing

History: "Naturalistic" stochastic optimization

- ~> mimicking freezing and crystallization
(atom configurations in crystals often close to global minima of the energy)
- ~> slow cooling (glass, unordered vs. crystal, ordered) \Rightarrow annealing

\Rightarrow slowly lower temperature while maintaining thermal equilibrium

\Rightarrow computational temperature T or *noise parameter* $\beta = \frac{1}{T}$

Simulated Annealing

initialization: \underline{s}_0, β_0 small (\leadsto high temperature)

BEGIN Annealing loop ($t = 1, 2, \dots$)

$\underline{s}_t = \underline{s}_{t-1}$ (initialization of inner loop)

BEGIN State update loop (M iterations)

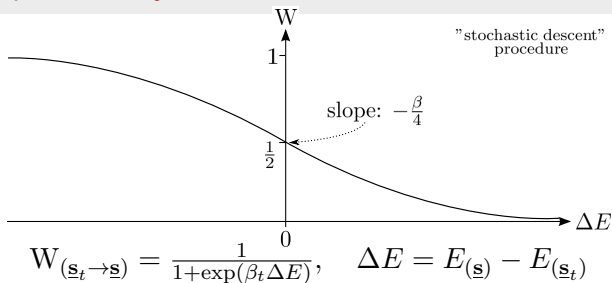
- choose a new candidate state \underline{s} randomly (local to \underline{s}_t – e.g. "bitflip")
- calculate difference in cost: $\Delta E = E(\underline{s}) - E(\underline{s}_t)$
- switch \underline{s}_t to \underline{s} with probability $W_{(\underline{s}_t \rightarrow \underline{s})} = \frac{1}{1 + \exp(\beta_t \Delta E)}$ otherwise keep the previous state \underline{s}_t

END State update loop

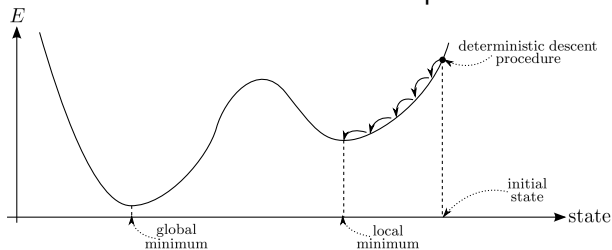
$\beta_t = \tau \beta_{t-1}$ ($\tau > 1 \implies$ increase of β)

END Annealing loop

Transition probability

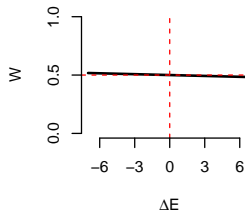


cost function with local optima:

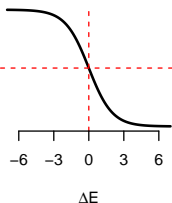


Annealing

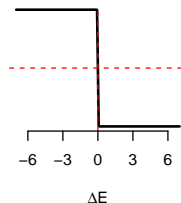
limiting cases for high vs. low temperature:



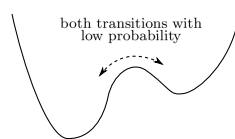
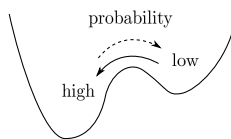
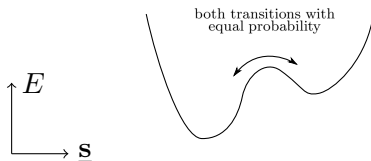
low β (high temperature)



intermediate β



high β



Annealing schedule & convergence

Convergence to the global optimum is guaranteed if: $\beta_t \sim \ln t$

- ⇒ robust optimization procedure
- ⇒ but: $\beta_t \sim \ln t$ is **too slow** for practical problems
- ⇒ therefore: $\beta_{t+1} = \tau \beta_t$, $\tau \in [1.01, 1.30]$ (exponential annealing)
- ⇒ additionally: the State Update loop has to be iterated often enough, e.g. $M = 500 - 2000$ (\leadsto thermal equilibrium)

Examples

1. Finding the global optimum of cost function (with continuous variables)

⇒ https://www.youtube.com/watch?v=iaq_Fpr4KZc

2. Solving Sudoku with Simulated Annealing

- initially fill columns randomly (without replacement)
- rows/3x3-boxes violate the Sudoku rules
- choose random column and two rows: switch the 2 numbers (stochastically)
- $s_i \in \{1, 2, \dots, 9\} \implies (9!)^9 \geq 10^{50}$ states
- cost function $E(\underline{s})$ total number of doubles in all rows/boxes (normalized)
- multiple global optima and also local optima
- 1000 steps per State Update loop

⇒ <https://www.youtube.com/watch?v=E8tkpzDne7I> (from 2:19)

The Gibbs distribution

- for constant β : noisy state change via Markov process $\underline{s}_{t'}$
- t' : iteration count of the State Update loop
- $\Pi_{(\underline{s}, t')}$: probability distribution across states

$$\Pi_{(\underline{s}, t')} \rightarrow \underbrace{P(\underline{s})}_{\substack{\text{stationary} \\ \text{distribution}}} \quad \text{for } t' \rightarrow \infty \text{ (and constant } \beta)$$

→ $P(\underline{s})$ can be calculated analytically!

Calculation of the stationary distribution

Assumption of *detailed balance*:

$$\underbrace{\text{probability of transition } \underline{s} \rightarrow \underline{s}'}_{P(\underline{s}) W(\underline{s} \rightarrow \underline{s}')} = \underbrace{\text{probability of transition } \underline{s}' \rightarrow \underline{s}}_{P(\underline{s}') W(\underline{s}' \rightarrow \underline{s})}$$

$$\begin{aligned} \frac{P(\underline{s})}{P(\underline{s}')} &= \frac{W(\underline{s}' \rightarrow \underline{s})}{W(\underline{s} \rightarrow \underline{s}')} = \frac{1 + \exp\left\{\beta(E(\underline{s}) - E(\underline{s}'))\right\}}{1 + \exp\left\{\beta(E(\underline{s}') - E(\underline{s}))\right\}} = \frac{1 + \exp(\beta\Delta E)}{1 + \exp(-\beta\Delta E)} \\ &= \exp(\beta\Delta E) \frac{1 + \exp(-\beta\Delta E)}{1 + \exp(-\beta\Delta E)} = \exp(\beta\Delta E) \end{aligned}$$

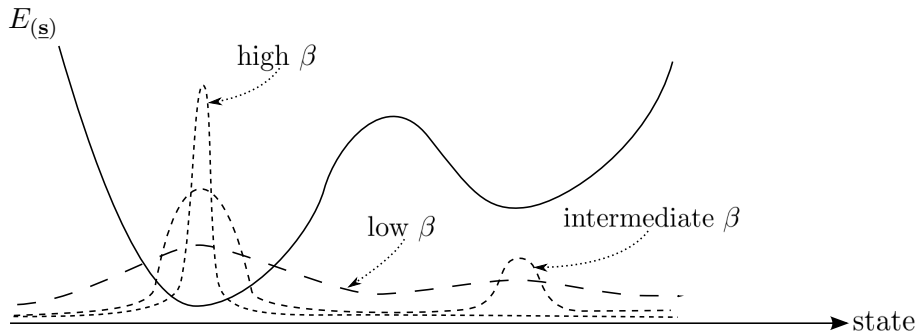
this condition is fulfilled for:

$$P(\underline{s}) = \frac{1}{Z} \exp(-\beta E) \quad (\text{Gibbs-Boltzmann-distribution})$$

normalization constant / partition function: $Z = \sum_{\underline{s}} \exp(-\beta E)$

Cost vs. probability distribution

$$P(\underline{s}) = \frac{1}{Z} \exp(-\beta E) \quad (\text{Gibbs-Boltzmann-distribution})$$



$\beta \downarrow$: broad, "delocalized" distribution

$\beta \uparrow$: distribution localized around (global) minima

Mean-field annealing

Simulated Annealing

- stochastic optimization: computationally expensive (sampling!)
- stationary distribution $P(\underline{s})$ known (for each β_t), why not evaluate?
- however: maxima of $P(\underline{s})$ equally hard to obtain as minima of $E(\underline{s})$
- moments? for $\beta \rightarrow \infty$: $\langle \underline{s} \rangle_P$ converges to \underline{s}^* of minimal cost ($P(\underline{s})$ singular)
- but: moments of $P(\underline{s})$ can – in general – not be calculated analytically

Approximation by Mean-Field Annealing

- ⇒ idea: approximate $P(\underline{s})$ by a computationally tractable distribution $Q(\underline{s})$
- ⇒ this distribution is then used to calculate the first moment $\langle \underline{s} \rangle_Q$
- ⇒ the first moment is tracked during the annealing schedule β_t
- ⇒ hope: $\langle \underline{s} \rangle_Q \rightarrow \underline{s}^*$ for $\beta_t \rightarrow \infty$

Factorizing distribution

Distribution $Q(\underline{s})$ to approximate $P(\underline{s})$

$$Q(\underline{s}) = \frac{1}{Z_Q} \exp \{-\beta E_Q\} = \frac{1}{Z_Q} \exp \left\{ -\beta \sum_k \underbrace{e_k}_{\text{parameters}} s_k \right\}$$

- Gibbs distribution with costs E_Q linear in the state variable \underline{s}_k

- factorizing distribution $Q(\underline{s}) = \prod_k Q_k(s_k)$ with
 $Q_k(s_k) = \frac{1}{Z_{Q_k}} \exp(-\beta e_k s_k)$

- $Q(\underline{s})$ factorizing $\iff s_k$ independent
 $\implies \langle \prod_k s_k \rangle_Q = \prod_k \langle s_k \rangle_Q$ (moments factorize)

- $\langle s_k \rangle_Q = \frac{\sum_{s_k} s_k \exp(-\beta e_k s_k)}{\sum_{s_k} \exp(-\beta e_k s_k)}$

→ family of distributions parametrized by the *mean fields* e_k

→ determine e_k such that this approximation is as good as possible

Mean-field approximation

Quantities

$$P(\underline{s}) = \frac{1}{Z_p} \exp(-\beta E_p) \quad \text{true distribution}$$

$$Q(\underline{s}) = \frac{1}{Z_Q} \exp \left(-\beta \overbrace{\sum_k e_k s_k}^{E_Q} \right) \quad \text{approximation: family of factorizing distributions}$$

$$e_k : \text{mean fields} \quad \text{parameters to be determined}$$

Good approximation of P by Q

→ minimization of the KL-divergence:

$$D_{\text{KL}}(Q||P) = \sum_{\underline{s}} Q(\underline{s}) \ln \frac{Q(\underline{s})}{P(\underline{s})} \stackrel{!}{=} \min_{\underline{e}}$$

Minimization of KL-divergence

$$D_{\text{KL}}(Q||P) = \sum_{\underline{s}} Q(\underline{s}) \ln \frac{Q(\underline{s})}{P(\underline{s})} \stackrel{!}{=} \min_{\underline{e}} \quad \begin{array}{l} P(\underline{s}) = \frac{1}{Z_p} \exp(-\beta E_p) \\ Q(\underline{s}) = \frac{1}{Z_Q} \exp\left(-\beta \sum_k e_k s_k\right) \end{array}$$

$$\begin{aligned} \frac{\partial}{\partial e_l} D_{\text{KL}} &= \frac{\partial}{\partial e_l} \left\{ \beta \sum_{\underline{s}} Q(\underline{s}) E_p - \beta \sum_{\underline{s}} Q(\underline{s}) E_Q + \ln Z_p - \ln Z_Q \right\} \\ &= \beta \frac{\partial}{\partial e_l} \langle E_p \rangle_Q - \underbrace{\beta \frac{\partial}{\partial e_l} \left(\sum_{\underline{s}} Q(\underline{s}) \sum_k e_k s_k \right)}_{-\beta \sum_k e_k \frac{\partial}{\partial e_l} \langle s_k \rangle_Q - \beta \langle s_l \rangle_Q} - \underbrace{\frac{1}{Z_Q} \sum_{\underline{s}} \frac{\partial}{\partial e_l} \exp(-\beta \sum_k e_k s_k)}_{+\beta \langle s_l \rangle_Q} \\ &= \beta \frac{\partial}{\partial e_l} \langle E_p \rangle_Q - \beta \sum_k e_k \frac{\partial}{\partial e_l} \langle s_k \rangle_Q \stackrel{!}{=} 0, \quad l = 1, \dots, N \end{aligned}$$

Result

$$\frac{\partial}{\partial e_l} \langle E_p \rangle_Q - \sum_k e_k \frac{\partial}{\partial e_l} \langle s_k \rangle_Q = 0$$

s_k are independent under Q :

$$\frac{\partial}{\partial e_l} \langle E_p \rangle_Q - e_l \frac{\partial}{\partial e_l} \langle s_l \rangle_Q = 0$$

$$\langle s_k \rangle_Q = \frac{\sum_{s_k} s_k \exp(-\beta e_k s_k)}{\sum_{s_k} \exp(-\beta e_k s_k)}$$

- coupled deterministic system of equations for $\{e_k\}$
- iterative solution procedure (usually no analytic result)

Mean-field annealing

Algorithm

initialization: $\langle \underline{s} \rangle_0, \beta_0$

BEGIN Annealing loop

Repeat

- calculate mean-fields: $e_k, \quad k = 1, \dots, N$
- calculate moments: $\langle s_k \rangle_Q, \quad k = 1, \dots, N$

Until $|e_k^{\text{old}} - e_k^{\text{new}}| < \varepsilon$

increase β

END Annealing loop

- \Rightarrow inner loop: fixed-point iteration for the mean-fields e_k (\leadsto EM-like)
- \Rightarrow deterministic (fast) rather than stochastic (slow) optimization method (given that mean-field equations can be easily evaluated, dep. on E_p)
- \Rightarrow moments $\langle s_k \rangle$ in general not from state space but $\langle s_k \rangle \rightarrow s_k^*$ for $\beta \rightarrow \infty$

Example (Ising model) – Setting and first Moments

Quadratic cost function $E(\underline{s})$ with binary variables $s_k \in \mathcal{S} = \{+1, -1\}$,

$$E_p(\underline{s}) = -\frac{1}{2} \sum_{\substack{i=1, j=1 \\ i \neq j}}^N W_{ij} s_i s_j,$$

real symmetric matrix \underline{W} , no self-coupling

→ Expressions required for the mean-field algorithm can be calculated:

$$\begin{aligned} \langle s_k \rangle_Q &= \frac{\sum_{s_k \in \mathcal{S}} s_k \exp(-\beta e_k s_k)}{\sum_{s_k \in \mathcal{S}} \exp(-\beta e_k s_k)} = \frac{(+1) \exp(-\beta e_k) + (-1) \exp(\beta e_k)}{\exp(-\beta e_k) + \exp(\beta e_k)} \\ &= \boxed{\tanh(-\beta e_k)} \end{aligned}$$

Example (Ising model) – Mean-fields

$$\begin{aligned}
 0 &= \frac{\partial}{\partial e_k} \langle E_p \rangle_Q - e_k \frac{\partial}{\partial e_k} \langle s_k \rangle_Q \\
 &= \frac{\partial}{\partial e_k} \left\langle -\frac{1}{2} \sum_{\substack{i=1, j=1 \\ i \neq j}}^N W_{ij} s_i s_j \right\rangle_Q - e_k \frac{\partial}{\partial e_k} \langle s_k \rangle_Q \\
 &= -\frac{1}{2} \frac{\partial}{\partial e_k} \sum_{\substack{i=1, j=1 \\ i \neq j}}^N W_{ij} \langle s_i \rangle_Q \langle s_j \rangle_Q - e_k \frac{\partial}{\partial e_k} \langle s_k \rangle_Q \\
 &= - \sum_{\substack{i=1 \\ i \neq k}}^N W_{ik} \langle s_i \rangle_Q \frac{\partial}{\partial e_k} \langle s_k \rangle_Q - e_k \frac{\partial}{\partial e_k} \langle s_k \rangle_Q
 \end{aligned}$$

$$\Rightarrow \boxed{e_k = - \sum_{\substack{i=1 \\ i \neq k}}^N W_{ik} \langle s_i \rangle_Q} \quad (\text{will be applied in exercise sheet 9})$$

Example (Ising model) – Fixed point iteration

Inner loop in mean-field annealing algorithm:

Repeat

- calculate mean-fields: $e_k = - \sum_{\substack{i=1 \\ i \neq k}}^N W_{ik} \langle s_i \rangle_Q, \quad k = 1, \dots, N$
- calculate moments: $\langle s_k \rangle_Q = \tanh(-\beta e_k), \quad k = 1, \dots, N$

Until $|e_k^{\text{old}} - e_k^{\text{new}}| < \varepsilon$

\leadsto fixed-point iteration for mean-fields will converge