

DM812
METAHEURISTICS

Lecture 12
Cross Entropy Method
Continuous Optimization

Marco Chiarandini

Department of Mathematics and Computer Science
University of Southern Denmark, Odense, Denmark
<marco@imada.sdu.dk>

Outline

1. Model Based Metaheuristics
Cross Entropy Method
2. Continuous Optimization
Numerical Analysis

Outline

1. Model Based Metaheuristics
Cross Entropy Method
2. Continuous Optimization
Numerical Analysis

Model Based Metaheuristics
Continuous Optimization CEM

Cross Entropy Method

Key idea: use **rare event-simulation** and **importance sampling** to proceed towards good solutions

- generate random solution samples according to a specified mechanism
- update the parameters of the random mechanism to produce a better “sample”

Notation:

- \mathcal{S} finite set of states
- f real valued performance functions on \mathcal{S}
- $\max_{s \in \mathcal{S}} f(s) = \gamma^* = f(s^*)$ (our problem)
- $\{p(s, \theta) \mid \theta \in \Theta\}$ family of discrete probability mass function on $s \in \mathcal{S}$
- $E_{\theta}[f(s)] = \sum_{s \in \mathcal{S}} f(s)p(s, \theta)$

We are interested in the probability that $f(s)$ is greater than some threshold γ under the probability $p(\cdot, \theta^*)$:

$$\ell = \Pr(f(s) \geq \gamma) = \sum_s I\{f(s) \geq \gamma\} p(s, \theta') = E_{\theta'} [I\{f(s) \geq \gamma\}]$$

if this probability is very small then we call $\{f(s) \geq \gamma\}$ a **rare event**

$$\ell = \sum_s I\{f(s) \geq \gamma\} p(s, \theta') = E_{\theta'} [I\{f(s) \geq \gamma\}]$$

Monte-Carlo simulation:

- draw a random sample
- compute unbiased estimator of ℓ : $\hat{\ell} = \frac{1}{N} \sum_{i=1}^N I\{f(s_i) \geq \gamma\}$
- if probability to sample $I\{f(s_i) \geq \gamma\}$ the estimation is not accurate

Importance sampling:

- use a **different** probability function g on \mathcal{S} to sample the solutions
- $\ell = \sum_s I\{f(s) \geq \gamma\} \frac{p(s, \theta')}{g(s)} g(s) = E_g [I\{f(s) \geq \gamma\} \frac{p(s, \theta')}{g(s)}]$
- compute unbiased estimator of ℓ :

$$\hat{\ell} = \frac{1}{N} \sum_{i=1}^N I\{f(s_i) \geq \gamma\} \frac{p(s, \theta')}{g(s)}$$

- How to determine g ?
Best choice would be:

$$g^*(s) := \frac{I\{f(s) \geq \gamma\} p(s, \theta')}{\ell},$$

as substituting $\hat{\ell} = \frac{1}{N} \sum_{i=1}^N I\{f(s_i) \geq \gamma\} \frac{p(s, \theta')}{g^*(s)} = \ell$.

But ℓ is unknown.

- It is convenient to choose g from $\{p(\cdot, \theta)\}$
- Choose the parameter θ such that the difference of $g = p(\cdot, \theta)$ to g^* is minimal
- **Cross entropy** or Kullback Leibler distance, measure of the distance between two probability distribution functions,

$$\mathcal{D}(g^*, g) = E_{g^*} \left[\ln \frac{g^*(s)}{g(s)} \right]$$

- Generalizing to probability density functions and Lebesgue integrals

$$\min \mathcal{D}(g^*, g) = \min_{\theta} \int g^*(s) \ln g^*(s) ds - \int g^*(s) \ln g(s, \theta) ds$$

- Minimizing the distance by means of sampling estimation leads to:

$$\hat{\theta} = \operatorname{argmax}_{\theta} E_{\theta'} [I\{f(s_i) \geq \gamma\} \frac{p(s, \theta')}{p(s, \theta)} \ln p(s, \theta)]$$

stochastic program (convex).

In some cases can be solved in closed form (eg, exponential, Bernoulli).

- Same result can be obtained by maximum likelihood estimation over the solutions s_i with performance $\geq \gamma$

$$L = \max_{\theta} \prod_{i=1}^N p(s_i, \theta)$$

- Estimation via stochastic counterpart:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \frac{1}{N} \sum_{i=1}^N I\{f(s_i) \geq \gamma\} \frac{p(s_i, \theta')}{p(s_i, \theta)} \ln p(s_i, \theta)$$

where s_1, \dots, s_N is a random sample from $p(\cdot, \theta')$.

- But still problems with sampling due to rare events.

Solution: Two-phase iterative approach:

- construct a sequence of levels $\hat{\gamma}_1, \hat{\gamma}_2, \dots, \hat{\gamma}_t$
- construct a sequence of parameters $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_t$

such that $\hat{\gamma}_t$ is close to optimal

and $\hat{\theta}_t$ assigns maximal probability to sample high quality solutions

- **Termination criterion:** if for some $t \geq d$ with, e.g., $d = 5$, $\hat{\gamma}_t = \hat{\gamma}_{t-1} = \dots = \hat{\gamma}_{t-d}$
- **Smoothed Updating:** $\hat{\theta}_t = \alpha \hat{\theta}'_t + (1 - \alpha) \hat{\theta}_{t-1}$ with $0.4 \leq \alpha \leq 0.9$
 θ'_t from the stochastic counterpart
- **Parameters:**
 - $N = cn$, n size of the problem (number of choices available for each solution component to decide)
 - $c > 1$ ($5 \leq c \leq 10$);
 - $\rho \approx 0.01$ for $n \geq 100$ and $\rho \approx \ln(n)/n$ for $n < 100$

Cross Entropy Method (CEM):

Define $\hat{\theta}_0$. Set $t = 1$

while *termination criterion* is not satisfied **do**

generate a sample (s_1, s_2, \dots, s_N) from the pdf $p(\cdot; \hat{\theta}_{t-1})$

set $\hat{\gamma}_t$ equal to the $(1 - \rho)$ -quantile with respect to f

$$\left(\hat{\gamma}_t = s^{(\lceil (1-\rho)N \rceil)} \right)$$

use the same sample (s_1, s_2, \dots, s_N) to solve the stochastic program

$$\hat{\theta}_t = \operatorname{argmax}_{\theta} \frac{1}{N} \sum_{i=1}^N I\{f(s_i) \leq \hat{\gamma}_t\} \ln p(s_i; \theta)$$

Example: TSP

- Solution representation: permutation representation
- Probabilistic model: matrix P where p_{ij} represents probability of vertex j after vertex i
- Tour construction: specific for tours
Define $P^{(1)} = P$ and $X_1 = 1$. Let $k = 1$
while $k < n - 1$ **do**
 - obtain $P^{(k+1)}$ from $P^{(k)}$ by setting the X_k -th column of $P^{(k)}$ to zero and normalizing the rows to sum up to 1.
Generate X_{k+1} from the distribution formed by the X_k -th row of $P^{(k)}$
 - set $k = k + 1$
- Update: take the fraction of times transition i to j occurred in those paths the cycles that have $f(s) \leq \gamma$

1. Model Based Metaheuristics
Cross Entropy Method
2. Continuous Optimization
Numerical Analysis

- We look at **unconstrained optimization** of **continuous, non-linear, non-convex, non-differentiable** functions
- Many applications above all in statistical estimation, (eg, likelihood estimation)
- Typically few variables (curse of dimensionality)

Standard Test Functions

- Rosenbrock's banana function

$$f(x, y) = (1 - x)^2 + 100(y - x^2)^2$$

Global minimum at $(x, y) = (1, 1)$ where $f(x, y) = 0$

Multidimensional extension is

$$f(x) = \sum_{i=1}^{N-1} [(1 - x_i)^2 + 100(x_{i+1} - x_i^2)^2] \quad \forall x \in \mathbb{R}^N.$$

Global minimum at $(x_1, \dots, x_N) = (1, \dots, 1)$

- Rastrigin's
- Schwefel's
- Sphere

Continue at: <http://www.cs.bham.ac.uk/research/projects/ecb/>

Smooth Functions Differentiable

Gradient Descent $f(x)$ decreases fastest moving in the direction of the negative gradient of f . Hence,

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \gamma \nabla f(\mathbf{x}_n)$$

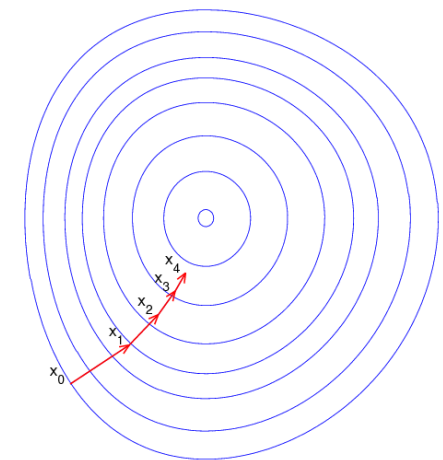
converges for appropriate x_0 and for $\gamma_n > 0$ small enough numbers.

- Problem is choosing γ

Secant Method

If only one-dimension and f hard to differentiate:

$$x_{n+1} = x_n - \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} f(x_n).$$



Newton's method in one dimension

Taylor expansion of $f(x)$:

$$f(x + \Delta x) = f(x) + f'(x)\Delta x + \frac{1}{2}f''(x)\Delta x^2,$$

attains its extremum when Δx solves the linear equation:

$$f'(x) + f''(x)\Delta x = 0 \quad \text{and } f''(x) > 0$$

Hence, if x_0 is chosen appropriately, the sequence below converges to x^*

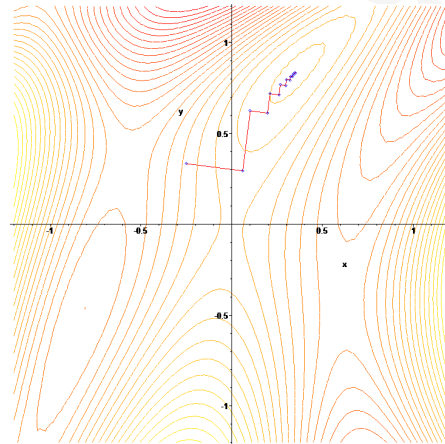
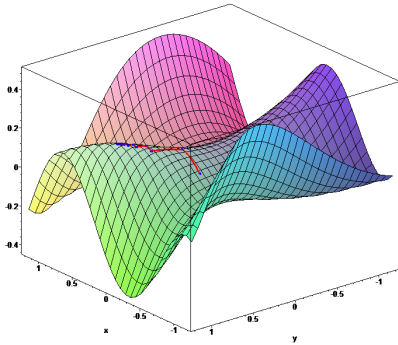
$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}, \quad n \geq 0$$

Newton's method generalized to several dimensions

first derivative \leftarrow gradient $\nabla f(\mathbf{x})$,

reciprocal of the second derivative \leftarrow inverse of Hessian matrix, $Hf(\mathbf{x})$

$$\mathbf{x}_{n+1} = \mathbf{x}_n - [Hf(\mathbf{x}_n)]^{-1}\nabla f(\mathbf{x}_n), \quad n \geq 0.$$



- Newton's method converges much faster towards a local maximum or minimum than gradient descent.
- However, finding the inverse of the Hessian may be an expensive operation, so approximations may be used instead
Quasi-Newton methods
 - Conjugate Gradient [Fletcher and Reeves (1964)]
 - BFGS (variable metric algorithm) [Broyden, Fletcher, Goldfarb and Shanno (1970)]