

Benchmarking Cross Entropy Method for partially Decomposable Problems

A Batch sampling method for large-scaled structured optimization



Siamak Ravanbakhsh, Russel Greiner
Computing Science Department, University of Alberta

Batch Sampling for Optimization

Monte Carlo methods for optimization includes both sequential and batch sampling methods.
Simulated annealing is a sequential method.
Cross Entropy and Probability Collectives are a batch sampling method

Cross Entropy Method

- Start from a prior
- Repeat until convergence
- Take samples from current distribution
- Calculate the loss for samples
- Select top (Elite) samples
- Find maximum Likelihood for elites

Advantages of batch to sequential methods

Highly parallelizable
Each instance for each iteration is independent of the other instances and may be parallelized.
Availability of Sensitivity Information
Sensitivity information (perturbation gives us extra information, value of the optimization function).
Our Alg. (CEED) uses Fisher Information to extract this information.

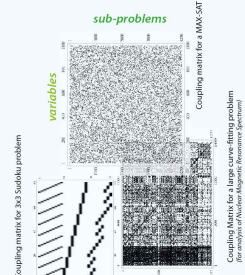
Partial Decomposability

The basic idea is to use sensitivity information that is available from batch sampling, to simultaneously optimize several (sub)-problems that share some variables.

- Consider problem of the form: $L(\beta) = \sum_{t=1}^n L_{p,t}(\beta^{p,t})$ when $X^p \subseteq X$
- We call such problems **partially decomposable**.
- CEED is designed to exploit such structures in the loss.

Coupling matrix representation

Each row represents a sub-problem and each column a variable.
Black squares show when a variable participates in a sub-problem



CEED Method

- Start from a Prior (with a product form)
- Repeat until convergence
- For each sub-problem p :
 - Take samples from this marginal
 - Calculate the loss for all samples
 - Select top (elite) samples
 - Calculate Maximum Likelihood dist. (with a product form) for these samples
- For each Variable V :
 - Consider all sub-problems that have V in their domain
 - Consider their posterior distribution over V
 - Combine these distributions to get a unified posterior
- Use joint distribution of unified posteriors as the next prior
- Return the mode of the last posterior as the optimal value found

The same sub-routine
Use Fisher Information to linearly combine the estimated max likelihood parameters. For discrete distributions this coincides with combination using inverse of variance-covariance matrix, asymptotically converge to this value for ML estimator (Cramer-Rao theorem). One might also combine distributions by finding a distribution with minimum sum of KL-divergence to all given dists.

Foveating for Continuous Domains

- We consider adaptive discretization of continuous domain.
- At each iteration, space is discretized into bins such that in the beginning of each iteration all bins have the same probability.
- Samples are taken uniformly from all bins and also uniformly within each bin
- Black lines show the boundary of bins for two variables.

Benchmark Problem

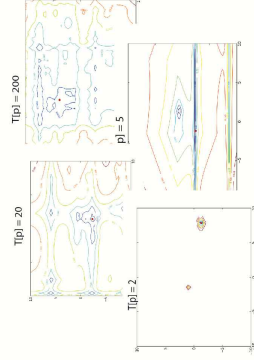
- We designed a partially decomposable multi-extremal loss, such that we can change
 - the number of variables
 - the number of sub-problems
 - the number of variables in each sub-problem
 - the difficulty of each sub-problem—i.e. number of terms ($T|p|$)

$$L(\beta) = \sum_{p=1}^P L_{p,t}(\beta^{p,t})$$

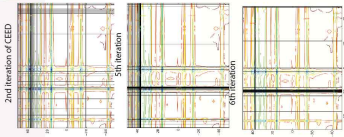
$$L_{p,t}(\beta^{p,t}) = \sum_{k=1}^{T|p|} c_k L_{p,t,k}(\beta^{p,t,k}) \quad X^{p,t,k} \subseteq X^p$$

$$L_{p,t,k}(\beta^{p,t,k}) = e^{-\sum_{j \in \mathcal{P}^{p,t,k}} (w_j x_j + c_j - s_{p,t,k})}$$

Optimization Landscape for various difficulties ($V=2$)

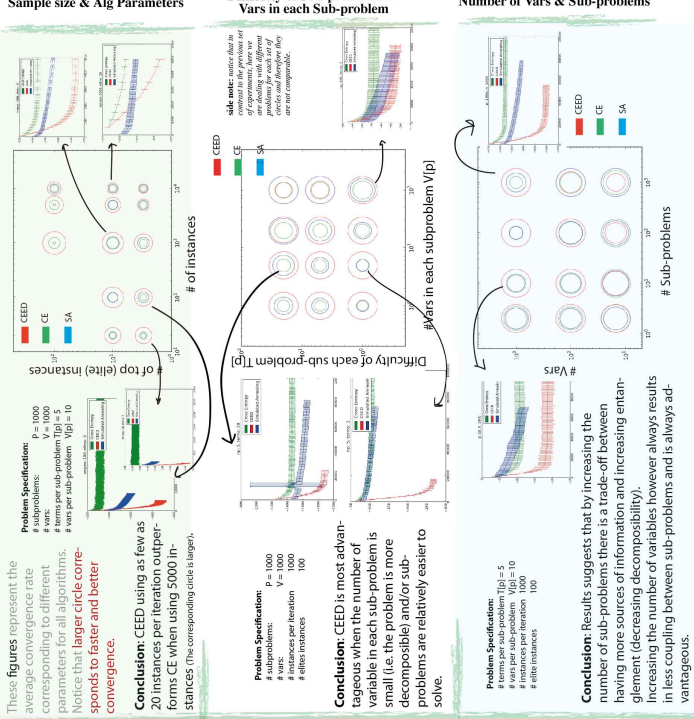


Contour plot of the Optimization Landscape and convergence of CEED by foveating



Side Note: The results shown in this paper are for a particular problem. The comparison between the parameters for all algorithms corresponds to number of elite instances corresponding to sub-problems. The ratio of elite instances to sub-problems corresponds to the algorithm's convergence.

CEED against CE and Simulated Annealing (SA)



These figures represent the average convergence rate corresponding to different parameters for all algorithms. Notice that larger circle corresponds to faster and better convergence.

Conclusion: CEED using as few as 20 instances per iteration outperforms CE when using 5000 instances. The corresponding circle is larger.

Conclusion: CEED is most advantageous when the number of variable in each sub-problem is small (i.e. the problem is more decomposable) and/or sub-problems are relatively easier to solve.

Conclusion: Results suggests that by increasing the number of sub-problems there is a trade-off between having more sources of information and increasing entanglement (decreasing decomposability). Increasing the number of variables however always results in less coupling between sub-problems and is always advantageous.

Future work

Energy minimization in molecular mechanics:

Any energy minimization task in molecular mechanics involves minimization of sum of many energy terms, therefore the partial decomposability structure. Each energy term usually involves only a subset of variables. For example energy terms corresponding to torsion angle in protein folding only depends on the configuration of four related residues. This motivates us to apply CEED to ab-initio protein folding problem or to similar but less explored problem of protein-docking.

Energy minimization in MRF:

In many problems of vision such as stereo-matching, denoising, image segmentation etc. the task is to find a labeling that minimizes an energy function. This energy function may be written as a sum of two terms for each individual label. One term penalizes deviation from the prior and the second penalizes distance from the observation. Again the problem demonstrates a decomposable structure and we intend to apply CEED to such problems.