

Project Type

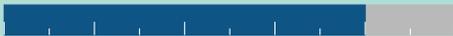
- Master Thesis
- Bachelor Thesis
- Research Project

Supervisors

-  Philipp Becker
-  philipp.becker@kit.edu
-  Maximilian Hüttenrauch
-  m.huettentrauch@kit.edu

Difficulty

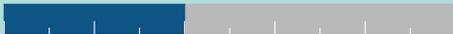
Algorithmic



Math



Application



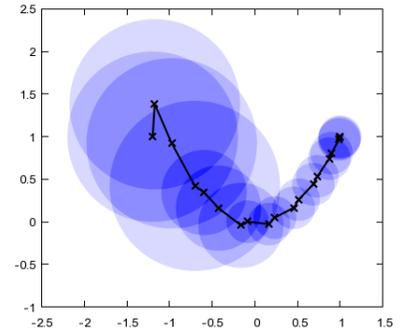
Information Theoretic Trust Regions for Gradient Descent

Description

Stochastic gradient-based optimization algorithms are important tools in machine learning, especially for optimizing neural networks. A prominent example is ADAM [3], a gradient descent based method with adaptive moment estimation. While it works well in practice, its update equations rely on heuristics and require step sizes which may not always be intuitive to choose. In this project, we want to develop more principled approaches for step-size control in gradient based methods using information theoretic principles.

MORE [1] is one such information-theoretic algorithm working with information theoretic trust regions. However, it has been designed black-box optimization, and therefore does not utilize any gradient information. More fits a local quadratic surrogate model for updating the policy. Modifying the fitting process of the surrogate such that it uses the additional gradient information would be a natural approach to extend MORE for scenarios where gradients are available, which should yield an information-theoretic alternative to ADAM as the covariance matrix of the search distribution is directly related to the step-size control of the gradient step.

Additionally, as MORE uses a Gaussian search distribution with a full covariance matrix, it scales poorly to higher dimensional problems. The naive choice to resolve this issue would be to assume a diagonal covariance which would completely neglect all correlations between the optimization variables. Factor Analyzers (FAs) provide an intuitive and theoretically justified way of interpolating between diagonal and full covariances, allowing to trade-off complexity with the amount of correlation modelled. Yet, FAs introduce latent variables and are thus harder to work with than simple Gaussians.



Tasks

- **Introducing Gradients into MORE.** The MORE algorithm is to be extended such that it can utilize gradient information. The natural approach for this is using the gradients to fit the squared surrogate. Alternative approaches to incorporate gradients can be explored.
- **Tackling Higher Dimensional Problems.** Using a Gaussian distribution prevents MORE from working with dimensions much greater than 50. FAs are to be employed to allow scaling the approach to higher dimensions. The decomposition introduced in [2] can be applied to handle the latent variables.
- **Evaluation** The approach is compared on several benchmark tasks to other common (stochastic) gradient descent approaches.
- **Theoretical Comparison** Ideally, we are able to formulate convergence guarantees similar to other common (stochastic) gradient descent approaches.

References

- [1] Abbas Abdolmaleki, Rudolf Lioutikov, Jan R Peters, Nuno Lau, Luis Pualo Reis, and Gerhard Neumann. Model-based relative entropy stochastic search. In *Advances in Neural Information Processing Systems*, pages 3537–3545, 2015.
- [2] Oleg Arenz, Mingjun Zhong, and Gerhard Neumann. Trust-region variational inference with gaussian mixture models. *arXiv preprint arXiv:1907.04710*, 2019.
- [3] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.