

SnapVX: A Network-Based Convex Optimization Solver

David Hallac*

Christopher Wong[†]

Steven Diamond[†]

Abhijit Sharang[†]

Rok Sosič[†]

Stephen Boyd*

Jure Leskovec[†]

*Department of Electrical Engineering**

Department of Computer Science[†]

Stanford University, Stanford, CA, 94305

HALLAC@STANFORD.EDU

CRWONG@CS.STANFORD.EDU

DIAMOND@CS.STANFORD.EDU

ABHISG@CS.STANFORD.EDU

ROK@CS.STANFORD.EDU

BOYD@STANFORD.EDU

JURE@CS.STANFORD.EDU

Editor: Cheng Soon Ong

Abstract

SnapVX is a high-performance solver for convex optimization problems defined on networks. For problems of this form, SnapVX provides a fast and scalable solution with guaranteed global convergence. It combines the capabilities of two open source software packages: Snap.py and CVXPY. Snap.py is a large scale graph processing library, and CVXPY provides a general modeling framework for small-scale subproblems. SnapVX offers a customizable yet easy-to-use Python interface with “out-of-the-box” functionality. Based on the Alternating Direction Method of Multipliers (ADMM), it is able to efficiently store, analyze, parallelize, and solve large optimization problems from a variety of different applications. Documentation, examples, and more can be found on the SnapVX website at <http://snap.stanford.edu/snapvx>.

Keywords: convex optimization, network analytics, graphs, data mining, ADMM

1. Introduction

Convex optimization is a widely used approach of modeling and solving problems in many different fields, as it offers well-established methods for finding globally optimal solutions. Numerous general-purpose optimization software packages exist (Sturm, 1999; Byrd et al., 2006; Mosek, 2010; Diamond and Boyd, 2016), but they typically rely on algorithms that are difficult to scale, so many modern machine learning problems cannot be solved by these common, yet general, approaches. Instead, solving large scale examples often requires developing problem-specific solution methods, which can be very fast but require significant optimization expertise to build. Furthermore, they are limited in scope, as they must be fine-tuned to only one particular type of problem and thus are hard to generalize.

In this paper, we build on the observation that many large convex optimization examples follow a common form in that they can often be split up into a series of subproblems using a network (graph) structure. Nodes are subproblems, representing anything from timestamps in a time-series data set to users in a social network. The edges then define the coupling, or relationships between the different nodes, and the combination of nodes and edges yields the

original convex optimization problem. This representation can refer to problems defined on actual networks, such as social or transportation systems, or questions classically modeled in other ways, such as control theory or time-series analysis.

Here, we present SnapVX, a solver that is both scalable and general on optimization problems defined over networks. It combines the graph capabilities of Snap.py (Leskovec and Sosič, 2016) with the general modeling framework from CVXPY (Diamond and Boyd, 2016). We show how SnapVX works, present syntax and supported features, scale it to large problems, and describe how it can be used to solve convex optimization problems from a variety of different fields. The full version of our solver, which is released under the BSD Open-Source License, can be found at the project website, <http://snap.stanford.edu/snapvx>. In addition to the source code, the download contains installation instructions, unit tests, documentation, and several examples to help users get started.

2. SnapVX General Form

Consider an optimization problem on an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, with vertex set \mathcal{V} and edge set \mathcal{E} , of the form:

$$\underset{x}{\text{minimize}} \quad \sum_{i \in \mathcal{V}} f_i(x_i) + \sum_{(j,k) \in \mathcal{E}} g_{jk}(x_j, x_k). \quad (1)$$

Variable x_i is associated with node i , for $i = 1, \dots, |\mathcal{V}|$ (the size of x_i can vary at each node). In problem (1), f_i is the cost function at node i , and g_{jk} is the cost associated with edge (j, k) . Constraints are represented by extended (infinite) values of f_i and g_{jk} . Note that nodes and edges can also have local, private optimization variables (which can also vary in size). These remain confined to a single node or edge, though, whereas the x_i 's are shared between different parts of the network. We consider only convex objectives and constraints for f_i and g_{jk} , so the entire problem is a convex optimization problem.

To solve problem (1), SnapVX uses a splitting algorithm based on the Alternating Direction Method of Multipliers (ADMM) (Boyd et al., 2011). With ADMM, each individual component of the graph solves its own subproblem, iteratively passing small messages over the network and eventually converging to an optimal solution. See the SnapVX developer documentation for more details on the underlying derivation. By splitting up the problem into a series of subproblems, SnapVX is able to solve much larger convex optimization examples than standard solvers, which solve the whole problem at once. ADMM also keeps the optimality benefits that general solvers enjoy: not only are we guaranteed to obtain an optimal solution, but we also are given a certificate of optimality in the form of primal and dual residuals; see (Parikh and Boyd, 2014).

SnapVX stores the network using a Snap.py graph structure. This allows fast traversal and easy manipulation, as well as efficient storage. On top of the standard graph, each node/edge is given convex objectives and constraints using CVXPY syntax. To find a global solution, SnapVX automatically splits up the problem and solves each subproblem using CVXPY, iteratively handling the ADMM message passing behind the scenes. Though wrapped in a Python layer, CVXPY uses ECOS and CVXOPT (two high-performance numerical optimization packages) as its primary underlying solvers (Andersen et al., 2014; Domahidi et al., 2013), so the Python overhead is not significant and allows for easier interpretability and improved user interface. To parallelize the solver, a worker pool coordinates

updates for the separate subproblems using Python’s multiprocessing library. SnapVX is built to run on a single machine, parallelizing across multiple cores and allowing “out-of-the-box” functionality on machines ranging from standard laptops to large-memory servers.

3. Syntax and Supported Features

We now present usage of SnapVX on a simple example. Complete documentation and more examples are available on the SnapVX website. Consider two nodes with an edge between them. We solve for a problem where each node has an unknown variable $x_i \in \mathbb{R}^1$. The first node’s objective is to minimize x_1^2 subject to $x_1 \leq 0$, the second’s is to minimize $|x_2 + 3|$, and the edge objective penalizes the square norm difference between the two variables, $\|x_1 - x_2\|_2^2$. The following code specifies the optimization problem and solves it:

```
from snapvx import *
gvx = TGraphVX() #Create a new graph
x1 = Variable(1, name='x1') #Create a variable for node 1
gvx.AddNode(1, Objective=square(x1), Constraints=[x1 <= 0]) #Add new node
x2 = Variable(1, name='x2') #Repeat for node 2
gvx.AddNode(2, abs(x2 + 3), [])
gvx.AddEdge(1, 2, Objective=square(norm(x1 - x2)), Constraints=[]) #Add edge between nodes
gvx.Solve() #Solve the problem
gvx.PrintSolution() #Print the solution on a node-by-node basis
```

As SnapVX is meant to be a general-purpose solver, it has many customizable options to help easily and efficiently solve a wide range of convex optimization problems. These include ADMM iteration limits, verbose mode (to list intermediate steps at each iteration), and defining customized convergence thresholds. Two key features are highlighted below:

- **Bulk Loading** - Often, the objectives and constraints at each node or edge will share a common form. For example, all the nodes could be trying to minimize $\|x - a_i\|_2$ for different values of a_i . Rather than requiring users to manually input each of these values, SnapVX allows for “bulk loading” of data. The functions *AddNodeObjectives* and *AddEdgeObjectives* allow the user to specify the general form of the objectives and an external file with separate data, and SnapVX fills in the details. This functionality makes practical and significantly speeds up the loading of very large data sets.
- **ADMM ρ -update** - The convergence time of ADMM depends on the value of the penalty parameter ρ (Nishihara et al., 2015), as it affects the tradeoff between primal and dual convergence, both of which need to be obtained for the overall problem to be solved. SnapVX users are not only able to select the value of ρ (it defaults to $\rho = 1$), but can also define a function to update ρ after each iteration based on the primal and dual residual values (He et al., 2000; Fougner and Boyd, 2015).

4. Scalability

One of the biggest benefits of SnapVX is that it allows us to solve large problems very efficiently. It does so by automatically parallelizing the ADMM updates across multiple cores of a single machine using Python’s multiprocessing class. Note that SnapVX is meant to be run on a single machine, rather than in a distributed computing environment. Convergence time depends on the problem complexity, but we empirically observe that it scales

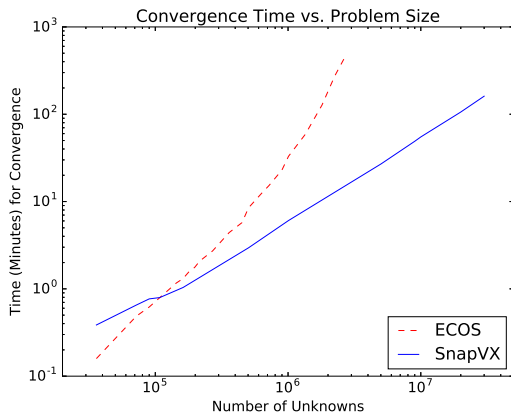


Figure 1: Timing comparison of SnapVX and a general-purpose solver.

# of Unknowns	SnapVX Solution Time (Minutes)
100,000	0.79
500,000	2.95
1 million	6.06
5 million	27.30
10 million	55.42
20 million	106.33
30 million	161.95

Table 1: SnapVX convergence time for several different problem sizes.

approximately linearly with problem size. We compare SnapVX and a general solver for a problem on a 3-regular graph. Each node solves for an unknown variable in \mathbb{R}^{9000} , where the node objectives are sum-of-Huber penalties (Huber, 1964) and edges have network lasso penalties (Hallac et al., 2015). By varying the number of nodes, we can span a wide range of problem sizes. The time required for SnapVX to converge, on a 40-core CPU where the entire problem can fit into memory, is shown in Table 1. Figure 1 displays a comparison to a general-purpose solver (ECOS, via CVXPY) in log-log scale.

5. Applications

SnapVX was created to allow users without significant optimization expertise to solve large network-based convex optimization problems by taking advantage of the powerful and scalable ADMM algorithm. Common examples from many different fields can be formulated in a SnapVX-friendly manner. Our software has been used to solve real-world machine learning problems ranging from housing price prediction (Hallac et al., 2015) to ride-sharing analysis (Ghosh et al., 2016) to high-dimensional regression (Yamada et al., 2016), often significantly outperforming other commonly used approaches. Examples in the SnapVX download package and on the project website — including financial modeling, network inference, image processing, PageRank, and time-series analysis — guide users towards applying the software to a variety of applications. Both user and developer documentation help new users get started and, if they are interested, contribute to the code base. A robust set of unit tests ensures that the software has been properly downloaded and that the code is functioning as expected. Overall, our open-source software provides an easy-to-use ADMM solver that can scale to large problems and apply to a wide variety of examples. With an active user base and rising interest from a range of scientific and engineering fields, we hope that SnapVX can become a useful tool for convex optimization problems, both in research and industry.

Acknowledgments

This work was supported by NSF IIS-1149837, NIH BD2K, DARPA XDATA, DARPA SIMPLEX, Stanford Data Science Initiative, Boeing, Bosch, Lightspeed, SAP, and Volkswagen.

References

- M. Andersen, J. Dahl, and L. Vandenberghe. CVXOPT: A Python package for convex optimization, version 1.1.7. <http://cvxopt.org/>, 2014.
- S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3:1–122, 2011.
- R. Byrd, J. Nocedal, and R. Waltz. KNITRO: An integrated package for nonlinear optimization. In *Large-scale nonlinear optimization*, pages 35–59. Springer, 2006.
- S. Diamond and S. Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- A. Domahidi, E. Chu, and S. Boyd. ECOS: An SOCP solver for embedded systems. In *European Control Conference*, pages 3071–3076, 2013.
- C. Fougner and S. Boyd. Parameter selection and pre-conditioning for a graph form solver. *arXiv preprint arXiv:1503.08366*, 2015.
- S. Ghosh, K. Page, and D. De Roure. An application of network lasso optimization for ride sharing prediction. *arXiv preprint arXiv:1606.03276*, 2016.
- D. Hallac, J. Leskovec, and S. Boyd. Network lasso: Clustering and optimization in large graphs. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 387–396, 2015.
- B. He, H. Yang, and S. Wang. Alternating direction method with self-adaptive penalty parameters for monotone variational inequalities. *Journal of Optimization Theory and Applications*, 106(2):337–356, 2000.
- P. Huber. Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, 35(1):73–101, 1964.
- J. Leskovec and R. Sosič. SNAP: A general purpose network analysis and graph mining library. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 8(1):1, 2016.
- ApS Mosek. The MOSEK optimization software. *Online at <http://www.mosek.com>*, 2010.
- R. Nishihara, L. Lessard, B. Recht, A. Packard, and M. Jordan. A general analysis of the convergence of ADMM. In *International Conference on Machine Learning*, pages 343–352, 2015.
- N. Parikh and S. Boyd. Proximal algorithms. *Foundations and Trends in Optimization*, 1: 123–231, 2014.
- J. Sturm. Using SeDuMi 1.02, a matlab toolbox for optimization over symmetric cones. *Optimization methods and software*, 11(1-4):625–653, 1999.
- M. Yamada, K. Takeuchi, T. Iwata, J. Shawe-Taylor, and S. Kaski. Localized lasso for high-dimensional regression. *arXiv preprint arXiv:1603.06743*, 2016.