

Nicol N. Schraudolph





Australian Government

Department of Communications, Information Technology and the Arts

Australian Research Council

NICTA Members

NICTA Partners







Department of State and



Modelling











Statistical Machine Learning Program



Why Stochastic Gradient?

Statistical Machine Learning Program



The Information Glut

The flood of information caused by

- Intiful, affordable sensors (such as webcams)
- ever-increasing networking of these sensors

overwhelms our processing ability in, e.g.,

- Science pulsar survey at Arecibo: | TB/day
- Solution by business Dell website: over 100 page requests/sec
- Security London: over 500'000 security cameras

We need intelligent, adaptive filters to cope!



The Challenge for ML

Coping with the info glut requires algorithms for

- Iarge adaptive models millions of degrees of freedom
- Iarge volumes of low-quality data noisy, correlated, non-stationary, outliers
- efficient real-time, online adaptation
 no fixed training set, life-long learning

Current optimization techniques are inadequate.



Online Learning Paradigm

classical optimization:



nested loops!

online learning:

online optimizer

training data stream

(*aka* adaptive filtering, stochastic approximation, ...)

Statistical Machine Learning Program



Stochastic Approximation

Classical formulation of optimization problem:

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} : E_{\boldsymbol{x}}[J(\boldsymbol{\theta}, \boldsymbol{x})] \approx \frac{1}{|X|} \sum_{\boldsymbol{x}_i \in X} J(\boldsymbol{\theta}, \boldsymbol{x}_i)$$

- $\ensuremath{{\ensuremath{{\rm G}}}}$ inefficient for large data sets X
- Inappropriate for never-ending, potentially non-stationary data streams
- \Rightarrow must resort to stochastic approximation:

$$\boldsymbol{\theta}_{t+1} \approx rg\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_t, \boldsymbol{x}_t) \quad (t = 0, 1, 2, \ldots)$$

Statistical Machine Learning Program



Stochastic Objective



Statistical Machine Learning Program



8

The Key Problem





Statistical Machine Learning Program

The Key Problem



9

Stochastic approximation breaks many optimizers:

- Se conjugate directions break down due to noise
- Ine minimizations (CG, quasi-Newton) inaccurate
- Wewton, Levenberg-Marquardt, Kalman filter too expensive for large-scale problems

This only leaves

- evolutionary alg.s very inefficient (don't use gradient)
- Simple gradient descent can be slow to converge



10

Stochastic Meta-Descent (SMD)

Statistical Machine Learning Program



Gain Vector Adaptation

Given stochastic gradient $g_t := \partial_{\theta} J(\theta_t, x_t)$, adapt θ by gradient descent with gain vector η : $\theta_{t+1} = \theta_t - \eta_t \cdot g_t$

Key idea:

simultaneously adapt **\eta** by exponentiated gradient:

$$\ln \eta_t = \ln \eta_{t-1} - \mu \partial_{\ln} \eta J(\boldsymbol{\theta}_t, \boldsymbol{x}_t)$$

$$\eta_t = \eta_{t-1} \cdot \exp(-\mu \partial_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_t, \boldsymbol{x}_t) \cdot \partial_{\ln} \eta \boldsymbol{\theta}_t)$$

$$\approx \eta_{t-1} \cdot \max(\frac{1}{2}, 1 - \mu \boldsymbol{g}_t \cdot \boldsymbol{v}_t)$$

Statistical Machine Learning Program



Single-Step Model

$$\begin{array}{ll} \text{Conventionally,} \quad \boldsymbol{v}_{t+1} \coloneqq \partial_{\ln} \boldsymbol{\eta}_t \ \boldsymbol{\theta}_{t+1} = - \ \boldsymbol{\eta}_t \cdot \boldsymbol{g}_t \\ & \quad \text{(recall that} \quad \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \boldsymbol{\eta}_t \cdot \boldsymbol{g}_t \text{)} \end{array}$$

giving
$$\boldsymbol{\eta}_t = \boldsymbol{\eta}_{t-1} \cdot \max(\frac{1}{2}, 1 + \mu \, \boldsymbol{\eta}_{t-1} \cdot \boldsymbol{g}_{t-1} \cdot \boldsymbol{g}_t)$$

\Rightarrow adaptation of η driven by autocorrelation of g:



Statistical Machine Learning Program



3

SMD's Multi-Step Model

To capture long-term dependence of θ on η :



define
$$v_{t+1} := \sum_{i=0}^{t} \lambda^i \frac{\partial \theta_{t+1}}{\partial \ln \eta_{t-i}}$$
 decay $0 \le \lambda \le 1$ (free parameter)

Statistical Machine Learning Program



14

SMD's v-update

$$\boldsymbol{v}_{t+1} = \lambda \boldsymbol{v}_t - \boldsymbol{\eta}_t \cdot (\boldsymbol{g}_t + \lambda \boldsymbol{H}_t \boldsymbol{v}_t)$$

- \mathbf{Q} we obtain a simple iterative update for \mathbf{v}
- \Im λ can smoothe over correlated input signals
- \bigcirc iteration similar to TD(λ) RL method (Sutton)
- generalizes Sutton's (1992) K1 algorithm linear to non-linear system, diagonal to full Hessian

Statistical Machine Learning Program



Fixpoint of
$$v_{t+1} = \lambda v_t - \eta_t \cdot (g_t + \lambda H_t v_t)$$

is a Levenberg-Marquardt style gradient step:

$$\boldsymbol{v} \rightarrow -[\lambda \boldsymbol{H} + (1-\lambda) \operatorname{diag}(\boldsymbol{\eta})^{-1}]^{-1} \boldsymbol{g}$$

✓ is too noisy to use directly; SMD achieves stability by means of the double integration v → η → θ
 ✓ g is well-behaved (self-normalizing property)
 ✓ non-convex fn.s: use Gauss-Newton approximation

Statistical Machine Learning Program

Fast Hv Product

Explicit computation of \mathbf{Hv} product would be $O(n^2)$. But: consider differential $\mathrm{d}\boldsymbol{g}(\boldsymbol{\theta}) = \boldsymbol{H}(\boldsymbol{\theta}) \,\mathrm{d}\boldsymbol{\theta}$

- Θ can set $d\theta := v$, forward-propagate through g()
- as efficient as 2-3 gradient evaluations (typ. O(n))
- \bigcirc matched approximations of **g** and $Hv \Rightarrow$ robust
- Sean even co-opt complex arithmetic:

$$g(\theta + i\epsilon d\theta) = g(\theta) + O(\epsilon^2) + i\epsilon dg(\theta) \quad (\epsilon \approx 10^{-150})$$

ΝΑΤΙΟΝΑ

16





SMD Benchmarks

Statistical Machine Learning Program



8

Four Regions Benchmark





Compare simple stochastic gradient (SGD), conventional gain vector adaptation (ALAP), stochastic meta-descent (SMD), and a global extended Kalman filter (GEKF).

Statistical Machine Learning Program



19

Benchmark: Convergence



Statistical Machine Learning Program



Computational Cost

Algorithm	<u>storage</u> weight	<u>flops</u> update	<u>CPU ms</u> pattern
SGD		6	0.5
SMD	3	18	I.0
ALAP	4	18	I.0
GEKF	>90	>1500	40

Statistical Machine Learning Program



21

Benchmark: CPU Usage



Statistical Machine Learning Program



22

Autocorrelated Data



i.i.d. uniform





Sobol





Brownian



Statistical Machine Learning Program



Conjugate Gradient deterministic stochastic (1000 pts) (1000 pts/iteration)

SMD stochastic (5 pts/iteration)

NATIONAL

ICT AUSTRALIA

23







overfits

diverges

converges

Statistical Machine Learning Program





24

SMD Applications

Statistical Machine Learning Program



25

Application: Turbulent Flow

(PhD thesis of M. Milano, Inst. of Computational Science, ETH Zürich)



Statistical Machine Learning Program



26

Hand Tracking with SMD

(PhD thesis of M. Bray, L. van Gool's Computer Vision Lab, ETH Zürich)

Annealed Particle Filter (114 sec/frame)

Stochastic Meta-Descent (3 sec/frame)



Statistical Machine Learning Program



Hand Tracking @ NICTA

with Desmond Chik (PhD) and Jochen Trumpf (SEACS)

- detailed hand model (26 dof, ~I0k vertices, skin blending, ...)
- In the second second
- Compare camera image(s) there, use resulting stochastic gradient to adjust model
- Search is a standard standa



28

Hand (et al.) Tracking



Statistical Machine Learning Program



29

SMD & Policy Gradient RL

with Jin Yu (PhD) and Doug Aberdeen (SML)

- SMD accelerates PG-RL
- complex interaction with temporal task structure
- had a poster at NIPS'05





Statistical Machine Learning Program



SMD for Online SVM

Online SVM aka NORMA (Kivinen, Smola, Williamson 2004):

- Online kernel method
- stochastic gradient in expansion coefficients
- 🥯 employs scalar gain η
- Applied SMD (with Vishy):
- \bigcirc **v** is function in RKHS



- <g,v> can be maintained incrementally in O(n)
- In presented at NIPS'05 workshop, submitted to JMLR

Statistical Machine Learning Program



Training CRFs (CoNLL-2000)

Sha, F. & Pereira, F. (2003). Shallow parsing with conditional random fields. In Proceedings of HLT-NAACL, 213–220. Association for Computational Linguistics.



(a) \mathcal{L}'_{λ} : CG (precond., mixed), L-BFGS

(b) \mathcal{L}'_{λ} : CG (precond., plain), GIS

www.nicta.com.au

Statistical Machine Learning Program



32

Training CRFs with SMD



I-D CRF chain for ConNLL-2000 Base NP chunking
 (predictable!) huge speed-up for online learning

Statistical Machine Learning Program



33

Training CRFs with SMD



BioNLP/NLPBA-2004 named-entity recognition task
 2-D CRF lattices for vision (M. Schmidt & K. Murphy, UBC) work well with loopy BP; more robust to overfitting

Statistical Machine Learning Program



SMD on 2-D CRF Lattices



ground truth

original

BFGS

logistic regression

SGD

SMD (loopy BP)



Statistical Machine Learning Program

log. regr.



35

Summary and Outlook

Summary:

data-rich ML problems need stochastic approximation

- Iclassical gradient methods are not up to the task
- SMD provides gain adaptation for stochastic gradient (Hv product gives cheap second-order information)

Wish List:

- Solution of the stability analysis for SMD (volunteers?)
- gain matrix version of SMD (rotation invariance)
- Online LBFGS; proof that online CG can't work

Statistical Machine Learning Program