



Nonconvex Optimization for Latent Data Models: An Incremental and An Online Point of View

Samsung-HSE Laboratory, National Research University Higher School of Economics

Belhal Karimi

CMAP (Ecole Polytechnique) & XPOP (INRIA Saclay)

July 24, 2019

Overview

1. Statistical Learning in Latent Data Models

2. Nonconvex Risk Minimization

Incremental Method for Non-smooth Nonconvex Objective: with Applications to Logistic Regression and Bayesian Deep Learning

Online Optimization of Nonconvex Expected Risk: with Applications to Online and Reinforcement Learning

3. Conclusion

1. Statistical Learning in Latent Data Models

Supervised Learning

- Given input-output pair of random variables (X, Y) taking values in arbitrary input set $X \subset \mathbb{R}^p$ and arbitrary output set $Y \subset \mathbb{R}^q$ from unknown distribution \mathcal{P} .
- *Modeling* phase: $M_\theta : X \mapsto Y$ of parameter $\theta \in \mathbb{R}^d$, called the *predictor*.
- Performance measured using a *loss function* $\ell : Y \times Y \mapsto \mathbb{R}$ where $\ell(y, y')$ is the loss incurred when the true output is y whereas y' is predicted.
- *Training* phase boils down to computing the following quantity:

$$\arg \min_{\theta \in \mathbb{R}^d} \bar{\mathcal{L}}(\theta) = \arg \min_{\theta \in \mathbb{R}^d} \{ \mathcal{L}(\theta) + R(\theta) \} \quad (1)$$

with

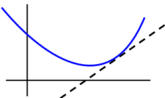
$$\mathcal{L}(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{P}} [\ell(y, M_\theta(x))] \quad \text{or} \quad \mathcal{L}(\theta) = n^{-1} \sum_{i=1}^n \ell(y_i, M_\theta(x_i)) . \quad (2)$$

Nonconvex Optimization

$$\mathcal{L} : \mathbb{R}^d \rightarrow \mathbb{R}$$

Convex function

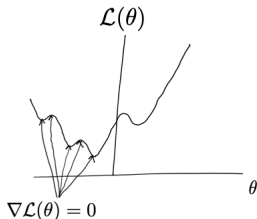
$$f(\lambda w_1 + (1-\lambda)w_2) \leq \lambda f(w_1) + (1-\lambda)f(w_2), \\ 0 \leq \lambda \leq 1$$



$$\mathcal{C} \subseteq \mathbb{R}^d$$

Convex set

$$\forall w_1, w_2 \in \mathcal{C}, \lambda w_1 + (1-\lambda)w_2 \in \mathcal{C} \\ 0 \leq \lambda \leq 1$$



- For convex problems, we generally use $|\mathcal{L}(\theta) - \mathcal{L}(\theta^*)|$ (or $\|\theta - \theta^*\|^2$) as stability condition. where θ^* is the optimal solution that can efficiently be found in the convex case.
- Yet, in the nonconvex case we use $\|\nabla \mathcal{L}(\theta)\|^2$, as advocated in (Nesterov, 2004) and (Ghadimi and Lan, 2013).
- A point θ^* is said to be ε -stationary if $\|\nabla \mathcal{L}(\theta^*)\|^2 \leq \varepsilon$. A stochastic iterative algorithm is said to achieve ε -stationarity in $\Gamma > 0$ iterations if $\mathbb{E}[\|\nabla \mathcal{L}(\theta^{(R)})\|^2] \leq \varepsilon$.

Latent Data Models

- Models where the input-output relationship is not completely characterized by the observed $(x, y) \in X \times Y$ pairs in the training set
- Dependence on a set of unobserved latent variables $z \in Z \subset \mathbb{R}^m$.
- Mandatory: Simulation step to complete the observed data with realizations of the latent variables.
- Formally, this specificity in our setting implies extending the loss function ℓ to accept a third argument as follows:

$$\ell(y, M_{\theta}(x)) = \int_Z \ell(z, y, M_{\theta}(x)) dz . \quad (3)$$

Some Examples of Latent Data Models

- Include the incomplete data framework, *i.e.*, some observations are missing, but are far broader than that: for example, the latent structure could stem from the unknown labels in mixture models or hidden states in Hidden Markov Models.
- **Missing Data:** y stands for the observed data and the latent variables z are the missing data.
- **Mixed Effects Models:** the latent variables z are the random effects and identifying the structure of the latent data mainly corresponds to the inter-individual variability among the individuals of the dataset.
- **Mixture Models:** the latent variables correspond to the unknown mixture labels taking values in a discrete finite set.

Quick Overview

- *L-Smoothness*: A function $f : \mathbb{R}^d \mapsto \mathbb{R}$ is L -smooth if and only if it is differentiable and its gradient is L -Lipschitz-continuous, i.e., for all $(\theta, \vartheta) \in \mathbb{R}^d \times \mathbb{R}^d$:

$$\|\nabla f(\theta) - \nabla f(\vartheta)\| \leq L \|\theta - \vartheta\| . \quad (4)$$

- We will deal with constrained ($\arg \min_{\theta \in \Theta} \bar{\mathcal{L}}(\theta)$) and unconstrained problems.

Algorithm	Gradient	Non-gradient	MC	Step.
SGD	$\mathcal{O}(1/\varepsilon^2)$ (Ghadimi and Lan, 2013)	?	x	γ_k
GD	$\mathcal{O}(n/\varepsilon)$ (Nesterov, 2004)	?	x	γ
SVRG/SAGA	$\mathcal{O}(n^{2/3}/\varepsilon)$ (Reddi et al., 2016)	$\mathcal{O}(n^{2/3}/\varepsilon)$ (Karimi et al., 2019c)	x	γ_k
MISO	$\mathcal{O}(n/\varepsilon)$ (Karimi et al., 2019b)	$\mathcal{O}(n/\varepsilon)$ (Karimi et al., 2019b)	x	—
MISSO	$\mathcal{O}(n/\varepsilon)$ (Karimi et al., 2019b)	$\mathcal{O}(n/\varepsilon)$ (Karimi et al., 2019b)	✓	—
Biased SA	$\mathcal{O}(c_0 + \frac{\log(n)}{\varepsilon\sqrt{n}})$ (Karimi et al., 2019a)	$\mathcal{O}(c_0 + \frac{\log(n)}{\varepsilon\sqrt{n}})$ (Karimi et al., 2019a)	✓	γ_k

Table 1: ERM methods: Table comparing the complexity, measured in terms of iterations, of different algorithms for non-convex optimization. MC stands for Monte Carlo integration of the drift term and Step. for stepsize.

2. Nonconvex Risk Minimization

2. Nonconvex Risk Minimization

2.1 Incremental Method for Non-smooth Nonconvex Objective: with Applications to Logistic Regression and Bayesian Deep Learning

Large-scale machine learning

Constrained Minimization of large sum of functions

We are interested in the minimization of a large finite-sum of functions:

$$\min_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta}) := \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\boldsymbol{\theta}), \quad (5)$$

where Θ is a convex, compact, and closed subset of \mathbb{R}^p , and for any $i \in \llbracket 1, n \rrbracket$, the function $\mathcal{L}_i : \mathbb{R}^p \rightarrow \mathbb{R}$ is bounded from below and is (possibly) nonconvex and non-smooth.

Some examples Given data points $(x_i, i \in \llbracket 1, n \rrbracket)$ and observations $(y_i, i \in \llbracket 1, n \rrbracket)$

- Maximum likelihood estimation: $\mathcal{L}(\boldsymbol{\theta}) \triangleq - \sum_{i=1}^N \log p_i(y_i, \boldsymbol{\theta})$
- Variational inference: $\mathcal{L}(\boldsymbol{\theta}) \triangleq \sum_{i=1}^N \text{KL} (q_i(w; \boldsymbol{\theta}) || p_i(w|y_i, x_i))$
- Logistic regression: $\mathcal{L}(\boldsymbol{\theta}) \triangleq \sum_{i=1}^N \log(1 + e^{-y_i \langle \boldsymbol{\theta}, x_i \rangle})$

Majorization-Minimization principle

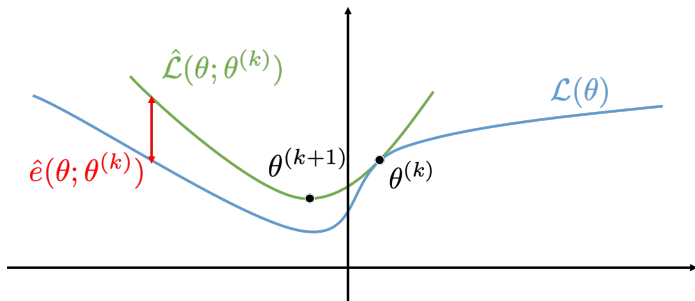


Figure 1: Majorization-Minimization principle

- Iteratively minimize locally tight upper bounds on the objective
- Drives the objective function downwards
- Examples: the proximal gradient algorithm ([Beck and Teboulle, 2009](#)), the EM algorithm ([McLachlan and Krishnan, 2007](#)) and variational inference ([Wainwright and Jordan, 2008](#)).

Notations and Assumptions

- Constrained optimization: Θ convex subset of \mathbb{R}^p .
- For all $i \in \llbracket 1, n \rrbracket$, \mathcal{L}_i is continuously differentiable on Θ .
- For all $i \in \llbracket 1, n \rrbracket$, \mathcal{L}_i is bounded from below, i.e. there exist a constant $M_i \in \mathbb{R}$ such as for all $\theta \in \Theta$, $\mathcal{L}_i(\theta) \geq M_i$.
- $\widehat{\mathcal{L}}_i(\theta; \cdot) : \mathbb{R}^p \rightarrow \mathbb{R}$ is a surrogate of \mathcal{L}_i at θ if the following properties are satisfied:
 1. the function $\vartheta \rightarrow \widehat{\mathcal{L}}_i(\theta; \vartheta)$ is continuously differentiable on Θ
 2. for all $\vartheta \in \Theta$, $\widehat{\mathcal{L}}_i(\theta; \vartheta) \geq \mathcal{L}_i(\vartheta)$, $\widehat{\mathcal{L}}_i(\theta; \theta) = \mathcal{L}_i(\theta)$ and
$$\nabla \widehat{\mathcal{L}}_i(\theta; \vartheta) \Big|_{\vartheta=\theta} = \nabla \mathcal{L}_i(\vartheta) \Big|_{\vartheta=\theta}.$$

Incremental Surrogate Minimization

The incremental scheme of (Mairal, 2015) computes surrogate functions, at each iteration of the algorithm, for a mini-batch of components:

Algorithm 1 MISO algorithm

Initialization: given an initial parameter estimate $\hat{\theta}^{(0)}$, for all $i \in \llbracket 1, n \rrbracket$ compute a surrogate function $\vartheta \rightarrow \hat{\mathcal{L}}_i(\hat{\theta}^{(0)}; \vartheta)$.

Iteration k: given the current estimate $\hat{\theta}^{(k)}$:

1. Pick i_k uniformly from $\llbracket 1, n \rrbracket$.
2. Update $\mathcal{A}_i^{k+1}(\theta)$ as:

$$\mathcal{A}_i^{k+1}(\theta) = \begin{cases} \hat{\mathcal{L}}_i(\theta; \hat{\theta}^{(k)}), & \text{if } i = i_k \\ \mathcal{A}_i^k(\theta), & \text{otherwise.} \end{cases}$$

3. Set $\hat{\theta}^{(k+1)} \in \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \mathcal{A}_i^{k+1}(\theta)$.
-

Intractable surrogate functions

- In many cases of interest those surrogates are intractable.
- Denote by $z = (z_i \in Z, i \in \llbracket 1, n \rrbracket)$ where Z is a subset of \mathbb{R}^{m_i} as set of latent variables.
- For all $i \in \llbracket 1, n \rrbracket$, let μ_i be a σ -finite measure on the Borel σ -algebra $\mathcal{Z} = \mathcal{B}(Z)$.
- $\mathcal{P}_i = \{p_i(z_i, \theta); \theta \in \Theta\}$ be a family of probability densities with respect to μ_i , and $r_{i,\theta} : Z \times \Theta \rightarrow \mathbb{R}$ be functions such that:

$$\widehat{\mathcal{L}}_i(\theta; \bar{\theta}) := \int_Z r_i(\theta; \bar{\theta}, z_i) p_i(z_i; \bar{\theta}) \mu_i(dz_i) \quad \forall (\theta, \bar{\theta}) \in \Theta \times \Theta. \quad (6)$$

The surrogate function denoted $\widehat{\mathcal{L}}_i(\theta; \vartheta)$ is fully defined by the pair $(r_i(\theta; \bar{\theta}, z_i), p_i(z_i, \bar{\theta}))$.

Examples of intractable surrogates

Incremental EM algorithm

- In the missing data context, let $c_i(z_i, \theta)$ be the joint likelihood of the observations and the latent data referred to as the complete likelihood.
- $g_i(\theta) \triangleq \int_{\mathcal{Z}} c_i(z_i, \theta) \mu_i(z_i)$ is the likelihood of the observations (in which the latent variables are marginalized).

The incremental EM algorithm falls into the incremental MM framework:

- For $i \in \llbracket 1, n \rrbracket$ and $\theta \in \Theta$ the loss function $\ell_i(\theta) \triangleq -\log g_i(\theta)$
- for $\vartheta \in \Theta$ the surrogate function $\widehat{\mathcal{L}}_i(\theta; \vartheta)$, introduced in the pioneering paper ([Neal and Hinton, 1998](#)), is defined by

$$\widehat{\mathcal{L}}_i(\theta; \vartheta) \triangleq \int_{\mathcal{Z}} \log \frac{p_i(z_i, \theta)}{c_i(z_i, \vartheta)} p_i(z_i, \theta) \mu_i(z_i) = \text{KL}(p_i(z_i, \theta) \parallel p_i(z_i, \vartheta)) + \ell_i(\vartheta) \quad (7)$$

In most cases, this surrogate is intractable.

Examples of intractable surrogates

Variational Inference Let $x = (x_i, i \in \llbracket 1, n \rrbracket)$ and $y = (y_i, i \in \llbracket 1, n \rrbracket)$ be i.i.d. input-output pairs and w be a global latent variable taking values in W a subset of \mathbb{R}^J .

A natural decomposition of the joint distribution is:

$$p(y, w|x) = \pi(w) \prod_{i=1}^n p(y_i|x_i, w). \quad (8)$$

The variational inference problem boils down to minimizing the following KL divergence:

$$\begin{aligned} \min_{\theta \in \Theta} \mathcal{L}(\theta) &:= \text{KL}(q(w; \theta) \parallel p(w|y, x)) \\ &:= \mathbb{E}_{q(w; \theta)} [\log(q(w; \theta)/p(w|y, x))] . \end{aligned} \quad (9)$$

Using (8), we decompose $\mathcal{L}(\theta) = n^{-1} \sum_{i=1}^n \mathcal{L}_i(\theta) + \text{const.}$ where:

$$\begin{aligned} \mathcal{L}_i(\theta) &:= -\mathbb{E}_{q(w; \theta)} [\log p(y_i|x_i, w)] \\ &+ \frac{1}{n} \mathbb{E}_{q(w; \theta)} [\log q(w; \theta)/\pi(w)] = r_i(\theta) + R(\theta) . \end{aligned} \quad (10)$$

Examples of intractable surrogates

Variational Inference

- MISSO method with a quadratic surrogate function defined as:

$$\widehat{\mathcal{L}}_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}) := \mathcal{L}_i(\bar{\boldsymbol{\theta}}) + \langle \nabla_{\boldsymbol{\theta}} \mathcal{L}_i(\bar{\boldsymbol{\theta}}) | \boldsymbol{\theta} - \bar{\boldsymbol{\theta}} \rangle + \frac{L}{2} \|\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2. \quad (11)$$

- Let $t : \mathbb{R}^d \times \Theta \mapsto \mathbb{R}^d$ be a differentiable function *w.r.t.* $\boldsymbol{\theta} \in \Theta$ s.t. $w = t(z, \bar{\boldsymbol{\theta}})$, where $z \sim \mathcal{N}_d(0, I)$, is distributed according to $q(\cdot, \bar{\boldsymbol{\theta}})$.
- By (Blundell et al., 2015, Proposition 1), the gradient of $-r_i(\cdot)$ in (10) is (re-parametrization trick):

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{q(w; \bar{\boldsymbol{\theta}})} [\log p(y_i | x_i, w)] = \mathbb{E}_{z \sim \mathcal{N}_d(0, I)} [J_{\boldsymbol{\theta}}^t(z, \bar{\boldsymbol{\theta}}) \nabla_w \log p(y_i | x_i, w) |_{w=t(z, \bar{\boldsymbol{\theta}})}].$$

- For most cases, the term $\nabla R(\bar{\boldsymbol{\theta}})$ can be evaluated in closed form.

$$r_i(\boldsymbol{\theta}; \bar{\boldsymbol{\theta}}, z) := \left\langle \nabla_{\boldsymbol{\theta}} d(\bar{\boldsymbol{\theta}}) - J_{\boldsymbol{\theta}}^t(z, \bar{\boldsymbol{\theta}}) \nabla_w \log p(y_i | x_i, w) \Big|_{w=t(z, \bar{\boldsymbol{\theta}})} \mid \boldsymbol{\theta} - \bar{\boldsymbol{\theta}} \right\rangle + \frac{L}{2} \|\boldsymbol{\theta} - \bar{\boldsymbol{\theta}}\|^2.$$

Proposed Method: MISSO

- *Minimization by Incremental Stochastic Surrogate Optimization* (MISSO) method: expectation approximated by *Monte Carlo* integration.
- Denote by $M \in \mathbb{N}$ the Monte Carlo batch size and let $z_m \in \mathcal{Z}$, $m = 1, \dots, M$ be a set of samples. These samples can be drawn (Case 1) i.i.d. from the distribution $p_i(\cdot; \bar{\theta})$ or (Case 2) from a Markov chain with the stationary distribution $p_i(\cdot; \bar{\theta})$.
- We define the following stochastic surrogate:

$$\tilde{\mathcal{L}}_i(\theta; \bar{\theta}, \{z_m\}_{m=1}^M) := \frac{1}{M} \sum_{m=1}^M r_i(\theta; \bar{\theta}, z_m) \quad (12)$$

MISSO Algorithm

Algorithm 2 MISSO algorithm

Initialization: $\hat{\theta}^{(0)}$; a sequence of non-negative numbers $\{M_{(k)}\}_{k=0}^{\infty}$.

For all $i \in \llbracket 1, n \rrbracket$, draw $M_{(0)}$ samples from $p_i(\cdot; \hat{\theta}^{(0)})$ and $\tilde{\mathcal{A}}_i^0(\theta) := \tilde{\mathcal{L}}_i(\theta; \hat{\theta}^{(0)}, \{z_{i,m}^{(0)}\}_{m=1}^{M_{(0)}})$.

Iteration k: given the current estimate $\hat{\theta}^{(k)}$:

1. Pick a function index i_k uniformly on $\llbracket 1, n \rrbracket$.
2. Draw $M_{(k)}$ Monte-Carlo samples from $p_i(\cdot; \hat{\theta}^{(k)})$.
3. Update the individual surrogate functions recursively as:

$$\tilde{\mathcal{A}}_i^{k+1}(\theta) = \begin{cases} \tilde{\mathcal{L}}_i(\theta; \hat{\theta}^{(k)}, \{z_{i,m}^{(k)}\}_{m=1}^{M_{(k)}}), & \text{if } i = i_k \\ \tilde{\mathcal{A}}_i^k(\theta), & \text{otherwise.} \end{cases} \quad (13)$$

4. Set $\hat{\theta}^{(k+1)} \in \arg \min_{\theta \in \Theta} \tilde{\mathcal{L}}^{(k+1)}(\theta) := \frac{1}{n} \sum_{i=1}^n \tilde{\mathcal{A}}_i^{k+1}(\theta)$.
-

Convergence analysis: Assumptions

(S1) For all $i \in \llbracket 1, n \rrbracket$ and $\bar{\theta} \in \Theta$, the function $\widehat{\mathcal{L}}_i(\theta; \bar{\theta})$ is convex w.r.t. θ , and it holds

$$\widehat{\mathcal{L}}_i(\theta; \bar{\theta}) \geq \mathcal{L}_i(\theta), \quad \forall \theta \in \Theta, \quad (14)$$

where the equality holds when $\theta = \bar{\theta}$.

(S2) For any $\bar{\theta}_i \in \Theta$, $i \in \llbracket 1, n \rrbracket$ and some $\epsilon > 0$, the difference function $\widehat{e}(\theta; \{\bar{\theta}_i\}_{i=1}^n) := \frac{1}{n} \sum_{i=1}^n \widehat{\mathcal{L}}_i(\theta; \bar{\theta}_i) - \mathcal{L}(\theta)$ is defined for all $\theta \in \Theta_\epsilon$ and differentiable for all $\theta \in \Theta$, where $\Theta_\epsilon = \{\theta \in \mathbb{R}^d, \inf_{\theta' \in \Theta} \|\theta - \theta'\| < \epsilon\}$ is an ϵ -neighborhood set of Θ . Moreover, for some constant L , the gradient satisfies

$$\|\nabla \widehat{e}(\theta; \{\bar{\theta}_i\}_{i=1}^n)\|^2 \leq 2L \widehat{e}(\theta; \{\bar{\theta}_i\}_{i=1}^n), \quad \forall \theta \in \Theta. \quad (15)$$

Convergence analysis: Assumptions

(H1) For all $i \in \llbracket 1, n \rrbracket$, $\bar{\theta} \in \Theta$, $z_i \in Z$, the measurable function $r_i(\theta; \bar{\theta}, z_i)$ is convex in θ and is lower bounded.

(H2) For the samples $\{z_{i,m}\}_{m=1}^M$, there exists finite constants C_r and C_{gr} such that

$$C_r := \sup_{\bar{\theta} \in \Theta} \sup_{M > 0} \frac{1}{\sqrt{M}} \mathbb{E}_{\bar{\theta}} \left[\sup_{\theta \in \Theta} \left| \sum_{m=1}^M \left\{ r_i(\theta; \bar{\theta}, z_{i,m}) - \hat{\mathcal{L}}_i(\theta; \bar{\theta}) \right\} \right| \right] \quad (16)$$

$$C_{gr} := \sup_{\bar{\theta} \in \Theta} \sup_{M > 0} \sqrt{M} \mathbb{E}_{\bar{\theta}} \left[\sup_{\theta \in \Theta} \left| \frac{1}{M} \sum_{m=1}^M \frac{\hat{\mathcal{L}}'_i(\theta, \theta - \bar{\theta}; \bar{\theta}) - r'_i(\theta, \theta - \bar{\theta}; \bar{\theta}, z_{i,m})}{\|\bar{\theta} - \theta\|} \right|^2 \right] \quad (17)$$

for all $i \in \llbracket 1, n \rrbracket$, and we denoted by $\mathbb{E}_{\bar{\theta}}[\cdot]$ the expectation *w.r.t.* a Markov chain $\{z_{i,m}\}_{m=1}^M$ with initial distribution $\xi_i(\cdot; \bar{\theta})$, transition kernel $P_{i,\bar{\theta}}$, and stationary distribution $p_i(\cdot; \bar{\theta})$.

Constrained Optimization

- As problem (5) is a constrained optimization, we consider the following stationarity measure:

$$g(\bar{\theta}) := \inf_{\theta \in \Theta} \frac{\mathcal{L}'(\bar{\theta}, \theta - \bar{\theta})}{\|\bar{\theta} - \theta\|} \quad \text{and} \quad g(\bar{\theta}) = g_+(\bar{\theta}) - g_-(\bar{\theta}), \quad (18)$$

where $g_+(\bar{\theta}) := \max\{0, g(\bar{\theta})\}$, $g_-(\bar{\theta}) := -\min\{0, g(\bar{\theta})\}$ denote the positive and negative part of $g(\bar{\theta})$, respectively.

- Note that $\bar{\theta}$ is a stationary point if and only if $g_-(\bar{\theta}) = 0$ (Fletcher et al., 2002).
- Furthermore, suppose that the sequence $\{\hat{\theta}^{(k)}\}_{k \geq 0}$ has a limit point $\bar{\theta}$ that is a stationary point, then one has $\lim_{k \rightarrow \infty} g_-(\hat{\theta}^{(k)}) = 0$.

Non-asymptotic analysis

Theorem 1

Under (S1), (S2), (H1), (H2). For any $K_{\max} \in \mathbb{N}$, let K be an independent discrete r.v. drawn uniformly from $\{0, \dots, K_{\max} - 1\}$ and define the following quantity:

$$\Delta_{(K_{\max})} := 2nL\mathbb{E}[\tilde{\mathcal{L}}^{(0)}(\hat{\boldsymbol{\theta}}^{(0)}) - \tilde{\mathcal{L}}^{(K_{\max})}(\hat{\boldsymbol{\theta}}^{(K_{\max})})] + \sum_{k=0}^{K_{\max}-1} \frac{4LC_r}{\sqrt{M_{(k)}}}, \quad (19)$$

Then we have following non-asymptotic bounds:

$$\mathbb{E}[\|\nabla \hat{e}^{(K)}(\hat{\boldsymbol{\theta}}^{(K)})\|^2] \leq \frac{\Delta_{(K_{\max})}}{K_{\max}} \quad (20)$$

$$\mathbb{E}[g_-(\hat{\boldsymbol{\theta}}^{(K)})] \leq \sqrt{\frac{\Delta_{(K_{\max})}}{K_{\max}}} + \frac{C_{gr}}{K_{\max}} \sum_{k=0}^{K_{\max}-1} M_{(k)}^{-1/2}. \quad (21)$$

Non-asymptotic analysis

$$\Delta_{(K_{\max})} := 2nL\mathbb{E}[\tilde{\mathcal{L}}^{(0)}(\hat{\theta}^{(0)}) - \tilde{\mathcal{L}}^{(K_{\max})}(\hat{\theta}^{(K_{\max})})] + \sum_{k=0}^{K_{\max}-1} \frac{4LC_r}{\sqrt{M_{(k)}}}, \quad (22)$$

Then we have following non-asymptotic bounds:

$$\mathbb{E}[\|\nabla \hat{e}^{(K)}(\hat{\theta}^{(K)})\|^2] \leq \frac{\Delta_{(K_{\max})}}{K_{\max}} \quad (23)$$

$$\mathbb{E}[g_{-}(\hat{\theta}^{(K)})] \leq \sqrt{\frac{\Delta_{(K_{\max})}}{K_{\max}}} + \frac{C_{\text{gr}}}{K_{\max}} \sum_{k=0}^{K_{\max}-1} M_{(k)}^{-1/2}. \quad (24)$$

- $\Delta_{(K_{\max})}$ is finite for any $K_{\max} \in \mathbb{N}$
- MISO method can be analyzed as a special case of the MISSO method satisfying $C_r = C_{\text{gr}} = 0$
- Then, Eq. (24) gives a non-asymptotic rate of $\mathbb{E}[g_{-}^{(K)}] \leq \mathcal{O}(\sqrt{nL/K_{\max}})$.

Asymptotic analysis

Under an additional assumption on the sequence of batch size $M_{(k)}$:

Theorem 2

Under (S1), (S2), (H1), (H2). In addition, assume that $\{M_{(k)}\}_{k \geq 0}$ is a non-decreasing sequence of integers which satisfies

$\sum_{k=0}^{\infty} M_{(k)}^{-1/2} < \infty$. Then:

- 1. the negative part of the stationarity measure converges almost surely to zero, i.e., $\lim_{k \rightarrow \infty} g_{-}(\hat{\theta}^{(k)}) = 0$ a.s..*
- 2. the objective value $\mathcal{L}(\hat{\theta}^{(k)})$ converges almost surely to a finite number $\underline{\mathcal{L}}$, i.e., $\lim_{k \rightarrow \infty} \mathcal{L}(\hat{\theta}^{(k)}) = \underline{\mathcal{L}}$ a.s..*

In particular, the first result above shows that the sequence $\{\hat{\theta}^{(k)}\}_{k \geq 0}$ produced by the MISSO method satisfies an *asymptotic stationary point condition*.

Logistic Regression with Missing Covariates

- $y = (y_i, i \in \llbracket 1, n \rrbracket)$ vector of binary responses and $z_i = (z_{i,p}) \in \mathbb{R}^p$ covariates.
- z_i is not fully observed. $z_{i,mis}$: missing values and $z_{i,obs}$: observed.
- $(z_i, i \in \llbracket 1, n \rrbracket) \sim \mathcal{N}(\beta, \Omega)$ where $\beta \in \mathbb{R}^p$ (i.i.d.)
- Model defined by

$$\text{logit}(\mathbb{P}(y_{ij} = 0 | z_i)) = d_{ij}^\top z_i$$

- Exponential family: Sufficient statistics are $\tilde{S}_i(z_i) \triangleq (z_i, z_i^\top z_i)$.
- MISSO algorithm consists in picking a set I_k , sampling a Monte Carlo batch $\{z_i^{k,m}\}_{m=0}^{M_k-1}$ for $i \in I_k$ and computing the quantities $(s_i^{1,k}, s_i^{2,k})$ as follows:

$$(s_i^{1,k}, s_i^{2,k}) = \begin{cases} \left(\frac{1}{M_k} \sum_{m=0}^{M_k-1} z_i^{k,m}, \frac{1}{M_k} \sum_{m=0}^{M_k-1} (z_i^{k,m})^\top z_i^{k,m} \right) & \text{if } i \in I_k \\ (s_i^{1,k-1}, s_i^{2,k-1}) & \text{otherwise} \end{cases} \quad (25)$$

Then $\beta^k = \frac{1}{N} \sum_{i=1}^N s_i^{1,k}$ and $\Omega^k = \frac{1}{N} \sum_{i=1}^N s_i^{2,k} - (\beta^k)^\top \beta^k$

Logistic Regression on TraumaBase

- TraumaBase (<http://traumabase.eu>) dataset: 15 trauma centers in France, covering measurements on patients from the initial to last stage of trauma. (6384 and $p = 16$ quantitative variables of influence)
- Predict binary response: severe trauma or not.
- We apply the MISSO method to fit a logistic regression model

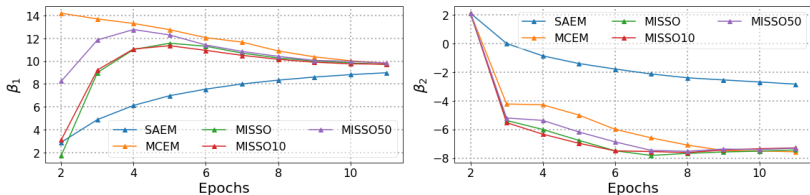


Figure 2: Convergence of two components of the vector of parameters β for the SAEM, the MCEM and the MISSO methods. The convergence is plotted against the number of passes over the data.

Bayesian LeNet-5 on MNIST

layer type	width	stride	padding	input shape	nonlinearity
convolution (5×5)	6	1	0	$1 \times 32 \times 32$	ReLU
max-pooling (2×2)		2	0	$6 \times 28 \times 28$	
convolution (5×5)	6	1	0	$1 \times 14 \times 14$	ReLU
max-pooling (2×2)		2	0	$16 \times 10 \times 10$	
fully-connected	120			400	ReLU
fully-connected	84			120	ReLU
fully-connected	10			84	

Table 1: LeNet-5 architecture

- $N = 60\,000$ handwritten digits, 28×28 images, $d = 784$
- $p(w) = \mathcal{N}(0, I)$, $p(y_i|x_i, w) = \text{Softmax}(f(x_i, w))$ where f is a NN.
- Variational distribution for layer j : $q(w_\ell, \theta_\ell)$ is a Gaussian distribution $\mathcal{N}(\mu_\ell, \sigma^2 I)$

Bayesian LeNet-5 on MNIST

- MISSO Updates:

$$\mu_\ell^{(k)} = \frac{1}{n} \sum_{i=1}^n \mu_\ell^{(\tau_i^k)} - \frac{\gamma}{n} \sum_{i=1}^n \hat{\delta}_{\mu_\ell, i}^{(k)} \quad \text{and} \quad \sigma^{(k)} = \frac{1}{n} \sum_{i=1}^n \sigma^{(\tau_i^k)} - \frac{\gamma}{n} \sum_{i=1}^n \hat{\delta}_{\sigma, i}^{(k)}, \quad (26)$$

where $\hat{\delta}_{\mu_\ell, i}^{(k)} = \hat{\delta}_{\mu_\ell, i}^{(k-1)}$ and $\hat{\delta}_{\sigma, i}^{(k)} = \hat{\delta}_{\sigma, i}^{(k-1)}$ for $i \neq i_k$ and:

$$\hat{\delta}_{\mu_\ell, i_k}^{(k)} = -\frac{1}{M^{(k)}} \sum_{m=1}^{M^{(k)}} \nabla_w \log p(y_{i_k} | x_{i_k}, w) \Big|_{w=t(\hat{\theta}^{(k-1)}, z_m^{(k)})} + \nabla_{\mu_\ell} R(\hat{\theta}^{(k-1)}),$$

$$\hat{\delta}_{\sigma, i_k}^{(k)} = -\frac{1}{M^{(k)}} \sum_{m=1}^{M^{(k)}} z_m^{(k)} \nabla_w \log p(y_{i_k} | x_{i_k}, w) \Big|_{w=t(\hat{\theta}^{(k-1)}, z_m^{(k)})} + \nabla_{\sigma} R(\hat{\theta}^{(k-1)})$$

with $R(\theta) = n^{-1} \sum_{\ell=1}^d (-\log(\sigma) + (\sigma^2 + \mu_\ell^2)/2 - 1/2)$.

Bayesian LeNet-5 on MNIST

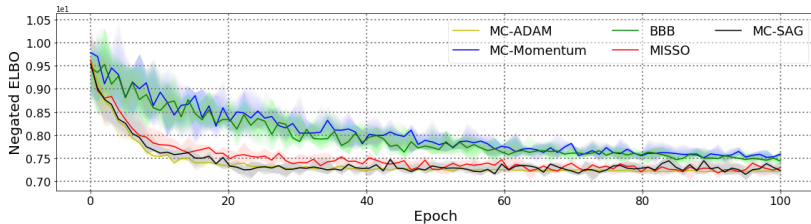


Figure 3: (Incremental Variational Inference) Negated ELBO versus epochs elapsed for fitting the Bayesian LeNet-5 on MNIST using different algorithms. The solid curve is obtained from averaging over 5 independent runs of the methods, and the shaded area represents the standard deviation.

2. Nonconvex Risk Minimization

2.2 Online Optimization of Nonconvex Expected Risk: with Applications to Online and Reinforcement Learning

Stochastic Approximation (SA) Scheme

- Consider a smooth Lyapunov function $V : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ (possibly nonconvex) that we wish to find its *stationary point*.
- SA scheme (Robbins and Monro, 1951) is a stochastic process:

$$\boldsymbol{\eta}_{n+1} = \boldsymbol{\eta}_n - \gamma_{n+1} H_{\boldsymbol{\eta}_n}(X_{n+1}), \quad n \in \mathbb{N}$$

where $\boldsymbol{\eta}_n \in \mathcal{H} \subseteq \mathbb{R}^d$ is the n th state, $\gamma_n > 0$ is the step size.

- The *drift term* $H_{\boldsymbol{\eta}_n}(X_{n+1})$ depends on an **i.i.d. random element** X_{n+1} and the mean-field satisfies

$$h(\boldsymbol{\eta}_n) = \mathbb{E}[H_{\boldsymbol{\eta}_n}(X_{n+1}) | \mathcal{F}_n] = \nabla V(\boldsymbol{\eta}_n),$$

where \mathcal{F}_n is the filtration generated by $\{\boldsymbol{\eta}_0, \{X_m\}_{m \leq n}\}$.

- In this case, the SA scheme is better known as the SGD method.

Biased SA Scheme

In this work, we relax a few restrictions of the classical SA. Consider:

$$\boldsymbol{\eta}_{n+1} = \boldsymbol{\eta}_n - \gamma_{n+1} H_{\boldsymbol{\eta}_n}(X_{n+1}), \quad n \in \mathbb{N}. \quad (27)$$

- The **mean field** $h(\boldsymbol{\eta}) \neq \nabla V(\boldsymbol{\eta})$
 \implies relevant to *non-gradient* method where the gradient is hard to compute, e.g., online EM.
- $\{X_n\}_{n \geq 1}$ is not i.i.d. and form a **state-dependent Markov chain**
 \implies relevant to *SGD with non-iid noise* and *policy gradient*. E.g., $\boldsymbol{\eta}_n$ controls the policy in a Markov decision process, and the gradient estimate $H_{\boldsymbol{\eta}_n}(x)$ is computed from the intermediate reward.

Biased SA Scheme

In this work, we relax a few restrictions of the classical SA. Consider:

$$\eta_{n+1} = \eta_n - \gamma_{n+1} H_{\eta_n}(X_{n+1}), \quad n \in \mathbb{N}. \quad (27)$$

- The **mean field** $h(\eta) \neq \nabla V(\eta)$ but satisfies for some $c_0 \geq 0, c_1 > 0$,

$$c_0 + c_1 \langle \nabla V(\eta) | h(\eta) \rangle \geq \|h(\eta)\|^2$$

- $\{X_n\}_{n \geq 1}$ is not i.i.d. and form a **state-dependent Markov chain**:

$$\mathbb{E}[H_{\eta_n}(X_{n+1}) | \mathcal{F}_n] = P_{\eta_n} H_{\eta_n}(X_n) = \int H_{\eta_n}(x) P_{\eta_n}(X_n, dx),$$

where $P_{\eta_n} : X \times \mathcal{X} \rightarrow \mathbb{R}_+$ is Markov kernel with a unique stationary distribution π_{η_n} , and the mean field $h(\eta) = \int H_{\eta}(x) \pi_{\eta}(dx)$.

Prior Work & Biased SA Scheme

We consider two cases depending on the noise sequence

$$\mathbf{e}_{n+1} = H_{\eta_n}(X_{n+1}) - h(\eta_n)$$

Case 1: When $\{\mathbf{e}_n\}_{n \geq 1}$ is Martingale difference —

$$\mathbb{E}[\mathbf{e}_{n+1} | \mathcal{F}_n] = 0 \text{ and other conditions...}$$

- *Asymptotic Analysis*: with smooth $h(\cdot)$ (Robbins and Monro, 1951), (Benveniste et al., 1990), (Borkar, 2009).
- *Non-asymptotic Analysis*: focus on $h(\boldsymbol{\eta}) = \nabla V(\boldsymbol{\eta})$,
 - Convex case: rate of $\mathcal{O}(1/n)$ in (Moulines and Bach, 2011), biased SA also studied in TD learning (Dalal et al., 2018).
 - Nonconvex case: (Ghadimi and Lan, 2013), (Bottou et al., 2018) studied convergence with martingale noise.

Prior Work & Biased SA Scheme

We consider two cases depending on the noise sequence

$$\mathbf{e}_{n+1} = H_{\eta_n}(X_{n+1}) - h(\eta_n)$$

Case 2: When $\{\mathbf{e}_n\}_{n \geq 1}$ is state-controlled Markov noise —

$\mathbb{E}[\mathbf{e}_{n+1} | \mathcal{F}_n] = P_{\eta_n} H_{\eta_n}(X_n) - h(\eta_n) \neq 0$ and other conditions....

- *Asymptotic Analysis*: studied with $h(\eta) = \nabla V(\eta)$ in (Kushner and Yin, 2003), similar biased SA setting in (Tadić and Doucet, 2017).
- *Non-asymptotic Analysis*: not many work here...
 - Sun et al. (2018) and Duchi et al. (2012) assumed $h(\eta) = \nabla V(\eta)$ & **state-independent** Markov chain.
 - Bhandari et al. (2018) studied a similar setting but focuses on linear SA with convex Lyapunov function.

Our Contributions

- First *non-asymptotic analysis* of biased SA scheme under the relaxed settings for *nonconvex* Lyapunov function.
- For both cases, with N being a r.v. drawn from $\{1, \dots, n\}$, we show

$$\mathbb{E}[\|h(\boldsymbol{\eta}_N)\|^2] = \mathcal{O}\left(c_0 + \frac{\log n}{\sqrt{n}}\right)$$

where c_0 is the *bias* of the mean field. If unbiased, then we find a stationary point.

- Analysis of two stochastic algorithms:
 - Online expectation maximization in (Cappé and Moulines, 2009)
 - Online policy gradient for infinite horizon reward maximization (Baxter and Bartlett, 2001).
- We provide the first *non-asymptotic* rates for the above algorithms.

General Assumptions

(A1) For all $\boldsymbol{\eta} \in \mathcal{H}$, there exists $c_0 \geq 0, c_1 > 0$ such that

$$c_0 + c_1 \langle \nabla V(\boldsymbol{\eta}) | h(\boldsymbol{\eta}) \rangle \geq \|h(\boldsymbol{\eta})\|^2$$

(A2) For all $\boldsymbol{\eta} \in \mathcal{H}$, there exists $d_0 \geq 0, d_1 > 0$ such that

$$d_0 + d_1 \|h(\boldsymbol{\eta})\| \geq \|\nabla V(\boldsymbol{\eta})\|$$

(A3) Lyapunov function V is L -smooth. For all $(\boldsymbol{\eta}, \boldsymbol{\eta}') \in \mathcal{H}^2$,

$$\|\nabla V(\boldsymbol{\eta}) - \nabla V(\boldsymbol{\eta}')\| \leq L\|\boldsymbol{\eta} - \boldsymbol{\eta}'\|$$

- **(A1), (A2)** assume that the mean field $h(\boldsymbol{\eta})$ is indirectly related to the gradient of a Lyapunov function $V(\boldsymbol{\eta})$ ($h(\boldsymbol{\eta}) \neq \nabla V(\boldsymbol{\eta})$).
- If $c_0 = d_0 = 0$, then the SA scheme is *un-biased*.
- **(A3)** is the standard smoothness assumption.

Stopping Criterion

- We adopt a stopping rule similar to (Ghadimi and Lan, 2013) that is typical for *nonconvex* problems.
- Fix any $n \geq 1$ and let $N \in \{0, \dots, n\}$ be a discrete random variable (independent of $\{\mathcal{F}_n, n \in \mathbb{N}\}$) with

$$\mathbb{P}(N = \ell) = \left(\sum_{k=0}^n \gamma_{k+1}\right)^{-1} \gamma_{\ell+1}, \quad (28)$$

where N serves as the terminating iteration for (27).

- Throughout this talk, we assume N is distributed as (28) and study the estimator η_N .

Case 1: Martingale Difference Noise

(A4) $\{\mathbf{e}_n\}_{n \geq 1}$ is a Martingale difference sequence such that $\mathbb{E}[\mathbf{e}_{n+1} | \mathcal{F}_n] = \mathbf{0}$, $\mathbb{E}[\|\mathbf{e}_{n+1}\|^2 | \mathcal{F}_n] \leq \sigma_0^2 + \sigma_1^2 \|h(\boldsymbol{\eta}_n)\|^2$ for any $n \in \mathbb{N}$.

Theorem 3

Let A1, A3 and A4 hold and $\gamma_{n+1} \leq (2c_1 L(1 + \sigma_1^2))^{-1}$ for all $n \geq 0$. Let $V_{0,n} := \mathbb{E}[V(\boldsymbol{\eta}_0) - V(\boldsymbol{\eta}_{n+1})]$, we have

$$\mathbb{E}[\|h(\boldsymbol{\eta}_N)\|^2] \leq \frac{2c_1(V_{0,n} + \sigma_0^2 L \sum_{k=0}^n \gamma_{k+1}^2)}{\sum_{k=0}^n \gamma_{k+1}} + 2c_0,$$

If we set $\gamma_k = (2c_1 L(1 + \sigma_1^2)\sqrt{k})^{-1}$, then the SA scheme (27) finds an $\mathcal{O}(c_0 + \log n/\sqrt{n})$ quasi-stationary point within n iterations.

Case 1: Martingale Difference Noise

(A4) $\{\mathbf{e}_n\}_{n \geq 1}$ is a Martingale difference sequence such that $\mathbb{E}[\mathbf{e}_{n+1} | \mathcal{F}_n] = \mathbf{0}$, $\mathbb{E}[\|\mathbf{e}_{n+1}\|^2 | \mathcal{F}_n] \leq \sigma_0^2 + \sigma_1^2 \|h(\boldsymbol{\eta}_n)\|^2$ for any $n \in \mathbb{N}$.
 \implies can be satisfied when X_n is *i.i.d.* similar to the SGD setting.

Theorem 3

Let A1, A3 and A4 hold and $\gamma_{n+1} \leq (2c_1 L(1 + \sigma_1^2))^{-1}$ for all $n \geq 0$.
Let $V_{0,n} := \mathbb{E}[V(\boldsymbol{\eta}_0) - V(\boldsymbol{\eta}_{n+1})]$, we have

$$\mathbb{E}[\|h(\boldsymbol{\eta}_N)\|^2] \leq \frac{2c_1 (V_{0,n} + \sigma_0^2 L \sum_{k=0}^n \gamma_{k+1}^2)}{\sum_{k=0}^n \gamma_{k+1}} + 2c_0,$$

If we set $\gamma_k = (2c_1 L(1 + \sigma_1^2) \sqrt{k})^{-1}$, then the SA scheme (27) finds an $\mathcal{O}(c_0 + \log n / \sqrt{n})$ quasi-stationary point within n iterations.

\implies if $h(\boldsymbol{\eta}) = \nabla V(\boldsymbol{\eta})$ it recovers (*Ghadimi and Lan, 2013, Theorem 2.1*).

Case 2: State-dependent Markov Noise

In this case, $\{\mathbf{e}_n\}_{n \geq 1}$ is not a Martingale sequence. Instead,

(A5) There exists a Borel measurable function $\hat{H} : \mathcal{H} \times X \rightarrow \mathcal{H}$,

$$\hat{H}_\eta(x) - P_\eta \hat{H}_\eta(x) = H_\eta(x) - h(\eta), \quad \forall \eta \in \mathcal{H}, x \in X.$$

(A6) For all $\eta \in \mathcal{H}$ and $x \in X$, $\|\hat{H}_\eta(x)\| \leq L_{PH}^{(0)}$, $\|P_\eta \hat{H}_\eta(x)\| \leq L_{PH}^{(0)}$, and

$$\sup_{x \in X} \|P_\eta \hat{H}_\eta(x) - P_{\eta'} \hat{H}_{\eta'}(x)\| \leq L_{PH}^{(1)} \|\eta - \eta'\|, \quad \forall (\eta, \eta') \in \mathcal{H}^2.$$

(A7) It holds that $\sup_{\eta \in \mathcal{H}, x \in X} \|H_\eta(x) - h(\eta)\| \leq \sigma$.

- **(A5)** refers to the existence of solution to Poisson equation.
- **(A6)** requires smoothness of $\hat{H}_\eta(x) \Leftarrow$ satisfied if the $P_\eta, H_\eta(X)$ are smooth *w.r.t.* η + the Markov chain is geometrically ergodic.
- *Remark:* **(A7)** requires the update is *uniformly bounded* for all $x \in X$. In fact, **(A5)–(A7)** can not imply **(A4)**, nor vice versa.

Case 2: State-dependent Markov Noise

Theorem 4

Let A1–A3, A5–A7 hold. Suppose that the step sizes satisfy

$$\gamma_{n+1} \leq \gamma_n, \quad \gamma_n \leq a\gamma_{n+1}, \quad \gamma_n - \gamma_{n+1} \leq a'\gamma_n^2, \quad \gamma_1 \leq 0.5(c_1(L + C_h))^{-1},$$

for $a, a' > 0$ and all $n \geq 0$. Let $V_{0,n} := \mathbb{E}[V(\boldsymbol{\eta}_0) - V(\boldsymbol{\eta}_{n+1})]$,

$$\mathbb{E}[\|h(\boldsymbol{\eta}_N)\|^2] \leq \frac{2c_1(V_{0,n} + C_{0,n} + (\sigma^2 L + C_\gamma) \sum_{k=0}^n \gamma_{k+1}^2)}{\sum_{k=0}^n \gamma_{k+1}} + 2c_0,$$

- $C_h := (L_{PH}^{(1)}(d_0 + \frac{d_1}{2}(a+1) + ad_1\sigma) + L_{PH}^{(0)}(L + d_1\{1+a'\}))$.
- $C_\gamma := L_{PH}^{(1)}(d_0 + d_0\sigma + d_1\sigma) + LL_{PH}^{(0)}(1 + \sigma)$.
- $C_{0,n} := L_{PH}^{(0)}((1 + d_0)(\gamma_1 - \gamma_{n+1}) + d_0(\gamma_1 + \gamma_{n+1}) + 2d_1)$.

Case 2: State-dependent Markov Noise

Theorem 4

Let A1–A3, A5–A7 hold. Suppose that the step sizes satisfy

$$\gamma_{n+1} \leq \gamma_n, \quad \gamma_n \leq a\gamma_{n+1}, \quad \gamma_n - \gamma_{n+1} \leq a'\gamma_n^2, \quad \gamma_1 \leq 0.5(c_1(L + C_h))^{-1},$$

for $a, a' > 0$ and all $n \geq 0$. Let $V_{0,n} := \mathbb{E}[V(\eta_0) - V(\eta_{n+1})]$,

$$\mathbb{E}[\|h(\eta_N)\|^2] \leq \frac{2c_1(V_{0,n} + C_{0,n} + (\sigma^2 L + C_\gamma) \sum_{k=0}^n \gamma_{k+1}^2)}{\sum_{k=0}^n \gamma_{k+1}} + 2c_0,$$

- If $\gamma_k = (2c_1 L(1 + C_h)\sqrt{k})^{-1}$, then $\mathbb{E}[\|h(\eta_N)\|^2] = \mathcal{O}(c_0 + \log n/\sqrt{n})$ as in our case 1 with Martingale noise.
- Key idea to the proof is to use the Poisson equation assumption **(A5)** [see [Lemma 3](#)], which is new to the SA analysis.

Latent Data Model

- **Goal:** Consider a stream of i.i.d. observations $\{Y_n\}_{n \geq 1}$, $Y_n \sim \pi$ and fit a parametric family $\{g(y; \theta) : \theta \in \Theta\}$ with $\theta \in \Theta \subset \mathbb{R}^d$.
- Augment Y with *latent variable* $Z \Rightarrow$ complete data: $X = (Y, Z)$.
- **Exponential Family Distribution:** the complete data distribution:

$$f(x; \theta) = h(x) \exp(\langle S(x) | \phi(\theta) \rangle - \psi(\theta)) ,$$

where $S : X \rightarrow S$ is the sufficient statistics and $S \subset \mathbb{R}^m$. We have $g(y; \theta) = \int_Z f(x; \theta) \mu(dz)$.

Two Important Operations

- **Expectation:** Given $\theta \in \Theta$ and an observation y , the conditional expected sufficient statistics is

$$\bar{\mathbf{s}}(y; \theta) = \mathbb{E}_{\theta} [S(X) | Y = y]$$

- **Maximization:** Given a sufficient statistics $\mathbf{s} \in S$, we can estimate θ by maximizing the regularized log-likelihood

$$\bar{\theta}(\mathbf{s}) := \arg \max_{\theta \in \Theta} \{ \langle \mathbf{s} | \phi(\theta) \rangle - \psi(\theta) - R(\theta) \}$$

where $R(\theta)$ is a (strongly convex) regularization function.

Regularized Online EM (ro-EM)

- As $\{Y_n\}_{n \geq 1}$ arrive in a streaming fashion, the ro-EM method (modified from (Cappé and Moulines, 2009)) does:

$$\begin{aligned} \text{E-step: } \hat{\mathbf{s}}_{n+1} &= \hat{\mathbf{s}}_n + \gamma_{n+1} \{ \bar{\mathbf{s}}(Y_{n+1}; \hat{\boldsymbol{\theta}}_n) - \hat{\mathbf{s}}_n \}, \\ \text{M-step: } \hat{\boldsymbol{\theta}}_{n+1} &= \bar{\boldsymbol{\theta}}(\hat{\mathbf{s}}_{n+1}). \end{aligned}$$

- Let us interpret **E-step** as an SA update (27) with drift term

$$H_{\hat{\mathbf{s}}_n}(Y_{n+1}) = \hat{\mathbf{s}}_n - \bar{\mathbf{s}}(Y_{n+1}; \bar{\boldsymbol{\theta}}(\hat{\mathbf{s}}_n)),$$

whose mean field is given by $h(\hat{\mathbf{s}}_n) = \hat{\mathbf{s}}_n - \mathbb{E}_\pi[\bar{\mathbf{s}}(Y_{n+1}; \bar{\boldsymbol{\theta}}(\hat{\mathbf{s}}_n))]$

- What should be the Lyapunov function? We use the KL divergence

$$V(\mathbf{s}) := \mathbb{E}_\pi \left[\log \frac{\pi(Y)}{g(Y; \bar{\boldsymbol{\theta}}(\mathbf{s}))} \right] + R(\bar{\boldsymbol{\theta}}(\mathbf{s})).$$

Special Case: Gaussian Mixture Model

Goal: fit $\{Y_n\}_{n \geq 1}$ in an GMM with $\theta = (\{\omega_m\}_{m=1}^{M-1}, \{\mu_m\}_{m=1}^M)$. Consider:

$$g(y; \theta) \propto \left(1 - \sum_{m=1}^{M-1} \omega_m\right) \exp\left(-\frac{(y - \mu_M)^2}{2}\right) + \sum_{m=1}^{M-1} \omega_m \exp\left(-\frac{(y - \mu_m)^2}{2}\right),$$

and the regularizer (with $\epsilon > 0$)

$$R(\theta) = \epsilon \sum_{m=1}^M \{\mu_m^2/2 - \log(\omega_m)\} - \epsilon \log\left(1 - \sum_{m=1}^{M-1} \omega_m\right).$$

Consider the assumption:

- **(A9)** The samples Y_n are i.i.d. and $|Y_n| \leq \bar{Y}$ for any $n \geq 0$.
- It can be verified that if **(A9)** holds, then **(A1)**, **(A3)**, **(A4)** are satisfied [cf. [Proposition 3, 4](#)].

Convergence Analysis

Corollary 1

Under A9 and set $\gamma_k = (2c_1L(1 + \sigma_1^2)\sqrt{k})^{-1}$. The ro-EM method for GMM finds $\hat{\mathbf{s}}_N$ such that

$$\mathbb{E}[\|\nabla V(\hat{\mathbf{s}}_N)\|^2] = \mathcal{O}(\log n/\sqrt{n})$$

The expectation is taken w.r.t. N and the observation law π .

- First *explicit non-asymptotic* rate given for online EM method.
- Note that rigorous convergence proof for *global convergence* of (online) EM methods are rare.

Policy Gradient: Markov Decision Process

- Consider a Markov Decision Process (MDP) (S, A, R, P) :
 - S is a finite set of spaces (state-space)
 - A is a finite set of action (action-space)
 - $R : S \times A \rightarrow [0, R_{\max}]$ is a reward function
 - P is the transition model, *i.e.*, given an action $a \in A$, $P^a = \{P_{s,s'}^a\}$ is a matrix, $P_{s,s'}^a$ is the probability of transiting from the s th state to the s' th state upon taking action a .
- $\{(S_t, A_t)\}_{t \geq 1}$ forms a Markov chain with the transition probability from (s, a) to (s', a') as:

$$Q_{\eta}((s, a); (s', a')) := \Pi_{\eta}(a'; s') P_{s,s'}^a .$$

- The parameter η controls the *conditional probability of taking action a' given the state s'* .

Policy Optimization Problem

- **Goal:** Find a policy η to maximize the average reward:

$$J(\eta) := \sum_{s \in \mathcal{S}, a \in \mathcal{A}} v(s, a) R(s, a) .$$

where $v(s, a)$ is the invariant distribution of $\{(S_t, A_t)\}_{t \geq 1}$.

- What is the gradient of $J(\eta)$ w.r.t. η ?

$$\nabla J(\eta) = \lim_{T \rightarrow \infty} \mathbb{E}_{\eta} \left[R(S_T, A_T) \sum_{i=0}^{T-1} \nabla \log \Pi_{\eta}(A_{T-i}; S_{T-i}) \right].$$

- REINFORCE algorithm ([Williams, 1992](#)) uses the sample average approximation. Let $M \gg 1, T \gg 1$,

$$\nabla J(\eta) \approx (1/M) \sum_{m=1}^M \left\{ R(S_T^m, A_T^m) \sum_{i=0}^{T-1} \nabla \log \Pi_{\eta}(A_{T-i}^m; S_{T-i}^m) \right\}$$

where $(S_1^m, A_1^m, \dots, S_T^m, A_T^m) \sim \Pi_{\eta}$ are drawn from a *roll-out* for each $m \implies$ *needs many samples and η to be static.*

Policy Optimization Problem

- **Goal:** Find a policy η to maximize the average reward:

$$J(\eta) := \sum_{s \in \mathcal{S}, a \in \mathcal{A}} v(s, a) R(s, a).$$

where $v(s, a)$ is the invariant distribution of $\{(S_t, A_t)\}_{t \geq 1}$.

- What is the gradient of $J(\eta)$ w.r.t. η ?

$$\nabla J(\eta) = \lim_{T \rightarrow \infty} \mathbb{E}_{\eta} \left[R(S_T, A_T) \sum_{i=0}^{T-1} \nabla \log \Pi_{\eta}(A_{T-i}; S_{T-i}) \right].$$

- We use a *biased* estimate of $\nabla J(\eta)$. Let $\lambda \in [0, 1)$ and $T \gg 1$, we have ([Baxter and Bartlett, 2001](#))

$$\nabla J(\eta) \approx \widehat{\nabla}_T J(\eta) := R(S_T, A_T) \sum_{i=0}^{T-1} \lambda^i \nabla \log \Pi_{\eta}(A_{T-i}; S_{T-i}),$$

where $(S_1, A_1, \dots, S_T, A_T) \sim \Pi_{\eta}$.

Online Policy Gradient (PG)

- We update the policy on-the-fly with an online policy gradient update (Baxter and Bartlett, 2001; Tadić and Doucet, 2017):

$$G_{n+1} = \lambda G_n + \nabla \log \Pi_{\eta_n}(A_{n+1}; S_{n+1}) , \quad (29a)$$

$$\eta_{n+1} = \eta_n + \gamma_{n+1} G_{n+1} R(S_{n+1}, A_{n+1}) . \quad (29b)$$

- We can interpret (29b) as an SA step with the drift term:

$$H_{\eta_n}(X_{n+1}) = G_{n+1} R(S_{n+1}, A_{n+1})$$

- Let the joint state be $X_n = (S_n, A_n, G_n) \in \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d$. We observe that $\{X_n\}_{n \geq 1}$ also forms a Markov chain. In particular,

$$h(\eta) = \lim_{T \rightarrow \infty} \mathbb{E}_{\tau_T \sim \Pi_\eta, S_1 \sim \bar{\Pi}_\eta} [\widehat{\nabla}_T J(\eta)] .$$

Convergence Analysis

- Focus on an exponential family policy (or soft-max):

$$\Pi_{\eta}(a; s) = \left\{ \sum_{a' \in \mathcal{A}} \exp \left(\langle \eta \mid \mathbf{x}(s, a') - \mathbf{x}(s, a) \rangle \right) \right\}^{-1}.$$

- **(A10)** For any $s \in \mathcal{S}$, we have $\|\mathbf{x}(s, a)\| \leq \bar{b}$.
- **(A11)** For any $\eta \in \mathcal{H}$, $\{S_t, A_t\}_{t \geq 1}$ is geometrically ergodic with $\|\mathbf{Q}_{\eta}^n - \mathbf{1}(\mathbf{v}_{\eta})^{\top}\| \leq \rho^n K_R$. The invariant distribution \mathbf{v}_{η} and its Jacobian $J_{\mathbf{v}_{\eta}}^{\eta}(\eta)$ are Lipschitz continuous

$$\|\mathbf{v}_{\eta} - \mathbf{v}_{\eta'}\| \leq L_Q \|\eta - \eta'\|, \quad \|J_{\mathbf{v}_{\eta}}^{\eta}(\eta) - J_{\mathbf{v}_{\eta}}^{\eta}(\eta')\| \leq L_v \|\eta - \eta'\|.$$

- Under **(A10)**, **(A11)**, the function $J(\eta)$ is $R_{\max} |\mathcal{S}| |\mathcal{A}|$ -smooth,

$$(1 - \lambda)^2 \Gamma^2 + 2 \langle \nabla J(\eta) \mid h(\eta) \rangle \geq \|h(\eta)\|^2,$$

where $\Gamma := 2\bar{b} R_{\max} K_R \frac{1}{(1-\rho)^2}$. Other required assumptions are satisfied too [cf. [Proposition 5-7](#)].

Convergence Analysis (cont'd)

Corollary 2

Under A10, A11 and set $\gamma_k = (2c_1L(1 + C_h)\sqrt{k})^{-1}$. For any $n \in \mathbb{N}$, the policy gradient algorithm (29) finds a policy that

$$\mathbb{E}[\|\nabla J(\boldsymbol{\eta}_N)\|^2] = \mathcal{O}\left((1 - \lambda)^2\Gamma^2 + c(\lambda)\log n/\sqrt{n}\right), \quad (30)$$

where $c(\lambda) = \mathcal{O}(\frac{1}{1-\lambda})$. Expectation is taken w.r.t. N and (A_n, S_n) .

- It shows the *first convergence rate* for the online PG method.
- Our result shows the *variance-bias trade-off* with $\lambda \in (0, 1)$.
- While setting $\lambda \rightarrow 1$ reduces the bias, but it decreases the convergence rate with $c(\lambda)$.

Sketches of Proofs

To simplify the notations next, we assume that $c_0 = 0, c_1 = 1$ and apply the following bound

Lemma 1. Let A1, A3 hold. With appropriate step size, it holds:

$$\begin{aligned} & \sum_{k=0}^n (\gamma_{k+1}/2) \|h(\boldsymbol{\eta}_k)\|^2 \\ & \leq V(\boldsymbol{\eta}_0) - V(\boldsymbol{\eta}_{n+1}) + L \sum_{k=1}^{n+1} \gamma_k^2 \|\mathbf{e}_k\|^2 - \sum_{k=1}^{n+1} \gamma_k \langle \nabla V(\boldsymbol{\eta}_{k+1}) | \mathbf{e}_k \rangle. \end{aligned}$$

- **Case 1:** $\{\mathbf{e}_n\}_{n \geq 1}$ is Martingale $\implies \mathbb{E}[\langle \nabla V(\boldsymbol{\eta}_{k+1}) | \mathbf{e}_k \rangle | \mathcal{F}_k] = 0$, taking the *total expectation* on both sides of the lemma suffices.
- **Case 2:** $\{\mathbf{e}_n\}_{n \geq 1}$ isn't Martingale $\implies \mathbb{E}[\langle \nabla V(\boldsymbol{\eta}_{k+1}) | \mathbf{e}_k \rangle | \mathcal{F}_k] \neq 0$.
- In the latter case, we control the sum $\mathbb{E}[\sum_{k=1}^{n+1} \gamma_k \langle \nabla V(\boldsymbol{\eta}_{k+1}) | \mathbf{e}_k \rangle]$ using the Poisson's equation.

Consider the following bound

Lemma 2. Let A1–A3, A5–A7 hold. With appropriate step size, it holds:

$$\begin{aligned} & \mathbb{E} \left[- \sum_{k=0}^n \gamma_{k+1} \langle \nabla V(\boldsymbol{\eta}_k) | \mathbf{e}_{k+1} \rangle \right] \\ & \leq C_h \sum_{k=0}^n \gamma_{k+1}^2 \mathbb{E}[\|h(\boldsymbol{\eta}_k)\|^2] + C_\gamma \sum_{k=0}^n \gamma_{k+1}^2 + C_{0,n}, \end{aligned}$$

- **Proof Idea:** note that $\mathbf{e}_{k+1} = H_{\boldsymbol{\eta}_k}(X_{k+1}) - h(\boldsymbol{\eta}_k)$, using **(A5)**, we have the decomposition

$$\sum_{k=0}^n \left\langle \nabla V(\boldsymbol{\eta}_k) | \hat{H}_{\boldsymbol{\eta}_k}(X_{k+1}) - P_{\boldsymbol{\eta}_k} \hat{H}_{\boldsymbol{\eta}_k}(X_{k+1}) \right\rangle \equiv A_1 + A_2 + A_3 + A_4 + A_5$$

$$A_1 = \sum_{k=1}^n \left\langle \nabla V(\boldsymbol{\eta}_k) | \hat{H}_{\boldsymbol{\eta}_k}(X_{k+1}) - P_{\boldsymbol{\eta}_k} \hat{H}_{\boldsymbol{\eta}_k}(X_k) \right\rangle$$

$$A_2 = \sum_{k=1}^n \left\langle \nabla V(\boldsymbol{\eta}_k) | P_{\boldsymbol{\eta}_k} \hat{H}_{\boldsymbol{\eta}_k}(X_k) - P_{\boldsymbol{\eta}_{k-1}} \hat{H}_{\boldsymbol{\eta}_{k-1}}(X_k) \right\rangle$$

$$A_3 = \sum_{k=1}^n \left\langle \nabla V(\boldsymbol{\eta}_k) - \nabla V(\boldsymbol{\eta}_{k-1}) | P_{\boldsymbol{\eta}_{k-1}} \hat{H}_{\boldsymbol{\eta}_{k-1}}(X_k) \right\rangle$$

$$A_4 = \sum_{k=1}^n \left\langle \nabla V(\boldsymbol{\eta}_k) | P_{\boldsymbol{\eta}_{k-1}} \hat{H}_{\boldsymbol{\eta}_{k-1}}(X_k) \right\rangle$$

$$A_5 = \gamma_1 \left\langle \nabla V(\boldsymbol{\eta}_0) | \hat{H}_{\boldsymbol{\eta}_0}(X_1) \right\rangle - \gamma_{n+1} \left\langle \nabla V(\boldsymbol{\eta}_n) | P_{\boldsymbol{\eta}_n} \hat{H}_{\boldsymbol{\eta}_n}(X_{n+1}) \right\rangle$$

- We can bound A_1 by observing that the inner product is a Martingale, bound A_2 using **(A6)**, **(A7)**, etc..

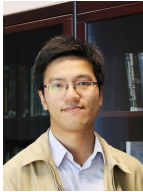
3. Conclusion

Take-aways

- We derived incremental and online methods for the optimization problem in machine learning.
 - For either ERM or Expected risk problems
 - When the objective function is a likelihood or not
 - For latent data models
- We conducted finite-time analysis of these methods for nonconvex loss functions and non necessarily gradient methods.
- Applications to several models of interest in machine learning with rigorous verification of the assumptions.

Perspectives

- Incremental algorithms: choice of the indices at each iteration. Optimal sampling strategies: (Roux et al., 2012) or (Horváth and Richtárik, 2018).
- Optimal mini-batch size of stochastic and incremental algorithms. See (Gower et al., 2019) (*variance-cost* trade off).
- Interplay between the Monte Carlo batch and the mini-batch of indices drawn at each iteration (*bias-variance* trade off).
- Complexity of $\mathcal{O}(n/\epsilon)$ was found for the MISO method. $\mathcal{O}(n^{2/3}/\epsilon)$ for quadratic surrogates in (Qian et al., 2019).
- Storage and Computation:
 - When data is big, need to store a lot for SAG, SAGA (less for SVRG).
 - Distributed first-order optimization procedures. Parallel or asynchronous method, or even centralized or decentralized architecture.



Thank you! Questions?

4. References

- Baxter, J. and Bartlett, P. L. (2001). Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15:319–350.
- Beck, A. and Teboulle, M. (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci*, 2:pp. 183.
- Benveniste, A., Priouret, P., and Métivier, M. (1990). *Adaptive Algorithms and Stochastic Approximation*.
- Bhandari, J., Russo, D., and Singal, R. (2018). A finite time analysis of temporal difference learning with linear function approximation. In *Conference On Learning Theory*, pages 1691–1692.
- Blundell, C., Cornebise, J., Kavukcuoglu, K., and Wierstra, D. (2015). Weight uncertainty in neural network. In *International Conference on Machine Learning*, pages 1613–1622.
- Borkar, V. S. (2009). *Stochastic approximation: a dynamical systems viewpoint*, volume 48. Springer.
- Bottou, L., Curtis, F. E., and Nocedal, J. (2018). Optimization methods for large-scale machine learning. *SIAM Review*, 60(2):223–311.
- Cappé, O. and Moulines, E. (2009). On-line Expectation Maximization algorithm for latent data models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(3):593–613.
- Dalal, G., Szorenyi, B., Thoppe, G., and Mannor, S. (2018). Finite sample analysis of two-timescale stochastic approximation with applications to reinforcement learning. In *Conference On Learning Theory*.
- Duchi, J. C., Agarwal, A., Johansson, M., and Jordan, M. I. (2012). Ergodic mirror descent. *SIAM Journal on Optimization*, 22(4):1549–1578.
- Fletcher, R., Gould, N. I., Leyffer, S., Toint, P. L., and Wächter, A. (2002). Global convergence of a trust-region sqp-filter algorithm for general nonlinear programming. *SIAM Journal on Optimization*, 13(3):635–659.
- Ghadimi, S. and Lan, G. (2013). Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368.
- Gower, R. M., Loizou, N., Qian, X., Sailanbayev, A., Shulgin, E., and Richtárik, P. (2019). SGD: general analysis and improved rates. *CoRR*, abs/1901.09401.
- Horváth, S. and Richtárik, P. (2018). Nonconvex variance reduced optimization with arbitrary sampling. *arXiv preprint arXiv:1809.04146*.

- Karimi, B., Miasojedow, B., Moulines, E., and Wai, H.-T. (2019a). Non-asymptotic analysis of biased stochastic approximation scheme. In Beygelzimer, A. and Hsu, D., editors, *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 1944–1974, Phoenix, USA. PMLR.
- Karimi, B., Wai, H.-T., and Moulines, E. (2019b). A doubly stochastic surrogate optimization scheme for non-convex finite-sum problems. *Submitted paper*.
- Karimi, B., Wai, H.-T., Moulines, E., and Lavielle, M. (2019c). On the global convergence of (fast) incremental expectation maximization methods. *Submitted paper*.
- Kingma, D. P. and Welling, M. (2014). Auto-encoding variational bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*.
- Kushner, H. and Yin, G. G. (2003). *Stochastic approximation and recursive algorithms and applications*, volume 35. Springer Science & Business Media.
- Mairal, J. (2015). Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM Journal on Optimization*, 25(2):829–855.
- McLachlan, G. and Krishnan, T. (2007). *The EM Algorithm and Extensions*.
- Moulines, E. and Bach, F. R. (2011). Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Advances in Neural Information Processing Systems*, pages 451–459.
- Neal, R. and Hinton, G. (1998). A view of the em algorithm that justifies incremental, sparse, and other variants. In *Learning in Graphical Models*, pages 355–368. Kluwer Academic Publishers.
- Nesterov, Y. (2004). *Introductory Lectures on Convex Optimization: A Basic Course*. Springer Publishing Company, Incorporated, 1 edition.
- Paisley, J., Blei, D., and Jordan, M. (2012). Variational bayesian inference with stochastic search. In *ICML*. icml.cc / Omnipress.
- Qian, X., Sailanbayev, A., Mishchenko, K., and Richtárik, P. (2019). Miso is making a comeback with better proofs and rates. *arXiv preprint arXiv:1906.01474*.
- Reddi, S. J., Sra, S., Póczos, B., and Smola, A. (2016). Fast incremental method for nonconvex optimization. *arXiv preprint arXiv:1603.06159*.
- Robbins, H. and Monro, S. (1951). A stochastic approximation method. *The Annals of Mathematical Statistics*, 22(3):400–407.
- Roux, N. L., Schmidt, M., and Bach, F. R. (2012). A stochastic gradient method with an exponential convergence rate for finite training sets. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems 25*, pages 2663–2671. Curran Associates, Inc.

- Sun, T., Sun, Y., and Yin, W. (2018). On Markov chain gradient descent. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R., editors, *Advances in Neural Information Processing Systems 31*, pages 9918–9927. Curran Associates, Inc.
- Tadić, V. B. and Doucet, A. (2017). Asymptotic bias of stochastic gradient search. *The Annals of Applied Probability*, 27(6):3255–3304.
- Wainwright, M. J. and Jordan, M. I. (2008). Graphical models, exponential families, and variational inference. *Found. Trends Mach. Learn.*, 1:pp. 1.
- Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3-4):229–256.