

Stochastic Approximations

Mauro Camara Escudero

School of Mathematics, University of Bristol

Table of contents

1. Robbins-Monro Procedure
2. Proof of Convergence
3. Application to EM Algorithm

Robbins-Monro Procedure

- Want to find **unique root** θ_* of a function $h(\theta)$ so that $h(\theta_*) = 0$.

- Want to find **unique root** θ_* of a function $h(\theta)$ so that $h(\theta_*) = 0$.
- Can only observe **noisy** measurements $H(\theta, X)$ with X RV and

$$h(\theta) = \mathbb{E}_X[H(\theta, X)] = \int H(\theta, x)p(x; \theta)dx$$

- Want to find **unique root** θ_* of a function $h(\theta)$ so that $h(\theta_*) = 0$.
- Can only observe **noisy** measurements $H(\theta, X)$ with X RV and

$$h(\theta) = \mathbb{E}_X[H(\theta, X)] = \int H(\theta, x)p(x; \theta)dx$$

- Under regularity and stability conditions the sequence

$$\theta_{k+1} = \theta_k + \gamma_{k+1}H(\theta_k, X_{k+1})$$

converges a.s. to θ_* . Here γ_{k+1} are *small* step-sizes.

Robbins-Monro Assumptions

- $(\gamma_k)_k$ big enough to explore, small enough for θ_k to converge.

$$\sum \gamma_k = \infty \quad \sum \gamma_k^2 < \infty$$

Robbins-Monro Assumptions

- $(\gamma_k)_k$ big enough to explore, small enough for θ_k to converge.

$$\sum \gamma_k = \infty \quad \sum \gamma_k^2 < \infty$$

- **Stability:** Recursion θ_k doesn't blow up. Related to convexity.

$$\sup_{\epsilon \leq |\theta - \theta_*| \leq \frac{1}{\epsilon}} (\theta - \theta_*)^\top h(\theta) < 0 \quad \forall \epsilon > 0$$

Robbins-Monro Assumptions

- $(\gamma_k)_k$ big enough to explore, small enough for θ_k to converge.

$$\sum \gamma_k = \infty \quad \sum \gamma_k^2 < \infty$$

- **Stability:** Recursion θ_k doesn't blow up. Related to convexity.

$$\sup_{\epsilon \leq |\theta - \theta_*| \leq \frac{1}{\epsilon}} (\theta - \theta_*)^\top h(\theta) < 0 \quad \forall \epsilon > 0$$

- **Existence:** $h(\theta)$ exists as \exists constant C bounding variance

$$\sigma^2(\theta) = \int |H(\theta, x)|^2 p(x; \theta) dx \leq C(1 + |\theta|^2)$$

Robbins-Monro Assumptions

- $(\gamma_k)_k$ big enough to explore, small enough for θ_k to converge.

$$\sum \gamma_k = \infty \quad \sum \gamma_k^2 < \infty$$

- **Stability:** Recursion θ_k doesn't blow up. Related to convexity.

$$\sup_{\epsilon \leq |\theta - \theta_*| \leq \frac{1}{\epsilon}} (\theta - \theta_*)^\top h(\theta) < 0 \quad \forall \epsilon > 0$$

- **Existence:** $h(\theta)$ exists as \exists constant C bounding variance

$$\sigma^2(\theta) = \int |H(\theta, x)|^2 p(x; \theta) dx \leq C(1 + |\theta|^2)$$

- **Memory:** Conditional expectation of X_{k+1} given past \mathcal{F}_k depends only on θ_k .

$$\mathbb{E}_X[g(\theta_k, X_{k+1}) \mid \mathcal{F}_k] = \int g(\theta_k, x)p(x; \theta_k)dx$$

Proof of Convergence

Robbins-Siegmund Lemma

If Z_k, B_k, C_k, D_k finite, non-negative RVs **known** given the **past** \mathcal{F}_k and satisfying

$$\mathbb{E}[Z_{k+1} \mid \mathcal{F}_k] \leq (1 + B_k)Z_k + C_k - D_k$$

on the set $\{\sum_k B_k < \infty, \sum_k C_k < \infty\}$ then

$$\sum_k D_k < \infty \quad \text{almost surely} \quad (1a)$$

$$Z_k \rightarrow Z < \infty \quad \text{almost surely} \quad (1b)$$

Robbins-Monro Convergence Proof (I)

- **Find expression for distance** $Z_{k+1} := |T_{k+1}|^2 := |\theta_{k+1} - \theta_*|^2$

1. *Expand the square and plug-in recursion formula*

$$Z_{k+1} = (\theta_k + \gamma_{k+1}H(\theta_k, X_{k+1}))^2 + \theta_*^2 - 2(\theta_k + \gamma_{k+1}H(\theta_k, X_{k+1}))\theta_*$$

Robbins-Monro Convergence Proof (I)

- **Find expression for distance** $Z_{k+1} := |T_{k+1}|^2 := |\boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_*|^2$

1. *Expand the square and plug-in recursion formula*

$$Z_{k+1} = (\boldsymbol{\theta}_k + \gamma_{k+1} H(\boldsymbol{\theta}_k, X_{k+1}))^2 + \boldsymbol{\theta}_*^2 - 2(\boldsymbol{\theta}_k + \gamma_{k+1} H(\boldsymbol{\theta}_k, X_{k+1})) \boldsymbol{\theta}_*$$

2. *Rearrange and rewrite in terms of T_k and Z_k*

$$Z_{k+1} = Z_k + 2\gamma_{k+1} T_k^\top H(\boldsymbol{\theta}_k, X_{k+1}) + \gamma_{k+1}^2 |H(\boldsymbol{\theta}_k, X_{k+1})|^2$$

- **Bound the expected distance** $\mathbb{E}[Z_{k+1} \mid \mathcal{F}_k]$

1. *Take conditional expectation given past \mathcal{F}_k . Use def of $\sigma^2(\boldsymbol{\theta})$ and $h(\boldsymbol{\theta})$.*

$$\mathbb{E}[Z_{k+1} \mid \mathcal{F}_k] = Z_k + 2\gamma_{k+1} T_k^\top h(\boldsymbol{\theta}_k) + \gamma_{k+1}^2 \sigma^2(\boldsymbol{\theta}_k)$$

2. *Use bound on $\sigma^2(\boldsymbol{\theta})$ and rearrange all constants.*

$$\mathbb{E}[Z_{k+1} \mid \mathcal{F}_k] \leq Z_k(1 + \gamma_{k+1}^2 C) + \gamma_{k+1}^2 \overline{C} + 2\gamma_{k+1} T_k^\top h(\boldsymbol{\theta}_k)$$

where $\overline{C} := C + \boldsymbol{\theta}_* C + 2T_k \boldsymbol{\theta}_* C$ is constant given \mathcal{F}_k .

Robbins-Monro Convergence Proof (II)

- **Prove convergence of Z_k to a finite RV Z**

1. Use stability condition $T_k^\top h(\theta_k) < 0$ to write it as Lemma, with

$$B_k := \gamma_{k+1}^2 C \quad C_k := \gamma_{k+1}^2 \bar{C} \quad D_k := -2\gamma_{k+1} T_k^\top h(\theta_k)$$

Robbins-Monro Convergence Proof (II)

- **Prove convergence of Z_k to a finite RV Z**

1. Use stability condition $T_k^\top h(\theta_k) < 0$ to write it as Lemma, with

$$B_k := \gamma_{k+1}^2 C \quad C_k := \gamma_{k+1}^2 \bar{C} \quad D_k := -2\gamma_{k+1} T_k^\top h(\theta_k)$$

2. Lemma gives convergence on $\{\sum \gamma_{k+1}^2 < \infty\}$

$$Z_k \xrightarrow{\text{a.s.}} Z < \infty \quad \text{and} \quad -\sum \gamma_{k+1} T_k^\top h(\theta_k) < \infty \text{ a.s.}$$

- **Prove Z is the zero random variable**

1. By contradiction assume $\exists \omega \in \Omega$ such that $Z(\omega) \neq 0$.
2. Using some analysis and the stability condition we finally find

$$\sum_k -\gamma_{k+1} T_k^\top(\omega) h(\theta_k(\omega)) \geq \alpha \sum_k \gamma_{k+1} = +\infty \quad \text{for some } \alpha > 0$$

which is a contradiction.

Application to EM Algorithm

Maximum Likelihood Estimation with Missing Data

- y observed data, x missing data

Maximum Likelihood Estimation with Missing Data

- y observed data, x missing data
- $p(\mathbf{x}, \mathbf{y} \mid \psi)$ complete-data likelihood

Maximum Likelihood Estimation with Missing Data

- \mathbf{y} observed data, \mathbf{x} missing data
- $p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi})$ complete-data likelihood
- Find $\boldsymbol{\psi}^*$ maximizing $p(\mathbf{y} \mid \boldsymbol{\psi})$ incomplete-data likelihood

Maximum Likelihood Estimation with Missing Data

- \mathbf{y} observed data, \mathbf{x} missing data
- $p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi})$ complete-data likelihood
- Find $\boldsymbol{\psi}^*$ maximizing $p(\mathbf{y} \mid \boldsymbol{\psi})$ incomplete-data likelihood
- Conditional distribution

$$p(\mathbf{x} \mid \mathbf{y}, \boldsymbol{\psi}) = \begin{cases} \frac{p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi})}{p(\mathbf{y} \mid \boldsymbol{\psi})} & \text{if } p(\mathbf{y} \mid \boldsymbol{\psi}) \neq 0 \\ 0 & \text{if } p(\mathbf{y} \mid \boldsymbol{\psi}) = 0 \end{cases}$$

- E-step

$$Q(\psi \mid \psi_k) := \mathbb{E}_{p(\mathbf{x} \mid \mathbf{y}, \psi_k)} [\log p(\mathbf{x}, \mathbf{y} \mid \psi)]$$

- M-step

$$\psi_{k+1} := \arg \max_{\psi \in \Psi} Q(\psi \mid \psi_k)$$

Intuition: Want to maximize $\log p(\mathbf{x}, \mathbf{y} \mid \psi)$ but it's unknown, so maximize current expectation given data and current estimate of ψ_k . EM useful when maximizing $Q(\psi \mid \psi_k)$ easier than maximizing $p(\mathbf{y} \mid \psi)$. $Q(\psi \mid \psi_k)$ can be **intractable**!

Replace $Q(\psi | \psi_k)$ with **Monte Carlo** estimate

$$Q(\psi | \psi_k) \approx \frac{1}{m_k} \sum_{i=1}^{m_k} \log p(\mathbf{x}_k^{(i)} | \psi)$$

where $\mathbf{x}_k^{(1)}, \dots, \mathbf{x}_k^{(m_k)}$ are drawn from $p(\mathbf{x} | \mathbf{y}, \psi_k)$.

Problem: At every new iteration, we discard previous samples.

- $p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi})$ follows a distribution from **curved exponential** family

$$p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi}) \propto \exp \left(-\xi(\boldsymbol{\psi}) + \langle \tilde{\boldsymbol{\theta}}(\mathbf{x}), \boldsymbol{\phi}(\boldsymbol{\psi}) \rangle \right) := \exp \left(L(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \boldsymbol{\psi}) \right)$$

- $p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi})$ follows a distribution from **curved exponential** family

$$p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi}) \propto \exp \left(-\xi(\boldsymbol{\psi}) + \langle \tilde{\boldsymbol{\theta}}(\mathbf{x}), \boldsymbol{\phi}(\boldsymbol{\psi}) \rangle \right) := \exp \left(L(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \boldsymbol{\psi}) \right)$$

- Define $\bar{\boldsymbol{\theta}}(\boldsymbol{\psi}) := \mathbb{E}_{p(\mathbf{x} \mid \mathbf{y}, \boldsymbol{\psi})} \left[\tilde{\boldsymbol{\theta}}(\mathbf{X}) \right]$. The E-step becomes

$$Q(\boldsymbol{\psi} \mid \boldsymbol{\psi}_k) = L(\bar{\boldsymbol{\theta}}(\boldsymbol{\psi}_k), \boldsymbol{\psi})$$

- $p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi})$ follows a distribution from **curved exponential** family

$$p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi}) \propto \exp \left(-\xi(\boldsymbol{\psi}) + \langle \tilde{\boldsymbol{\theta}}(\mathbf{x}), \boldsymbol{\phi}(\boldsymbol{\psi}) \rangle \right) := \exp \left(L(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \boldsymbol{\psi}) \right)$$

- Define $\bar{\boldsymbol{\theta}}(\boldsymbol{\psi}) := \mathbb{E}_{p(\mathbf{x} \mid \mathbf{y}, \boldsymbol{\psi})} \left[\tilde{\boldsymbol{\theta}}(\mathbf{X}) \right]$. The E-step becomes

$$Q(\boldsymbol{\psi} \mid \boldsymbol{\psi}_k) = L(\bar{\boldsymbol{\theta}}(\boldsymbol{\psi}_k), \boldsymbol{\psi})$$

- Assume there exists function $\hat{\boldsymbol{\psi}}$ that maximizes the E-step

$$L(\boldsymbol{\theta}, \hat{\boldsymbol{\psi}}(\boldsymbol{\theta})) \geq L(\boldsymbol{\theta}, \boldsymbol{\psi}) \quad \forall \boldsymbol{\theta} \quad \forall \boldsymbol{\psi}$$

- $p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi})$ follows a distribution from **curved exponential** family

$$p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi}) \propto \exp \left(-\xi(\boldsymbol{\psi}) + \langle \tilde{\boldsymbol{\theta}}(\mathbf{x}), \boldsymbol{\phi}(\boldsymbol{\psi}) \rangle \right) := \exp \left(L(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \boldsymbol{\psi}) \right)$$

- Define $\bar{\boldsymbol{\theta}}(\boldsymbol{\psi}) := \mathbb{E}_{p(\mathbf{x}|\mathbf{y}, \boldsymbol{\psi})} \left[\tilde{\boldsymbol{\theta}}(\mathbf{X}) \right]$. The E-step becomes

$$Q(\boldsymbol{\psi} \mid \boldsymbol{\psi}_k) = L(\bar{\boldsymbol{\theta}}(\boldsymbol{\psi}_k), \boldsymbol{\psi})$$

- Assume there exists function $\hat{\boldsymbol{\psi}}$ that maximizes the E-step

$$L(\boldsymbol{\theta}, \hat{\boldsymbol{\psi}}(\boldsymbol{\theta})) \geq L(\boldsymbol{\theta}, \boldsymbol{\psi}) \quad \forall \boldsymbol{\theta} \quad \forall \boldsymbol{\psi}$$

- Basically $\boldsymbol{\psi}_{k+1} = \hat{\boldsymbol{\psi}}(\bar{\boldsymbol{\theta}}(\boldsymbol{\psi}_k))$

- $p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi})$ follows a distribution from **curved exponential** family

$$p(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\psi}) \propto \exp \left(-\xi(\boldsymbol{\psi}) + \langle \tilde{\boldsymbol{\theta}}(\mathbf{x}), \boldsymbol{\phi}(\boldsymbol{\psi}) \rangle \right) := \exp \left(L(\tilde{\boldsymbol{\theta}}(\mathbf{x}), \boldsymbol{\psi}) \right)$$

- Define $\bar{\boldsymbol{\theta}}(\boldsymbol{\psi}) := \mathbb{E}_{p(\mathbf{x} \mid \mathbf{y}, \boldsymbol{\psi})} \left[\tilde{\boldsymbol{\theta}}(\mathbf{X}) \right]$. The E-step becomes

$$Q(\boldsymbol{\psi} \mid \boldsymbol{\psi}_k) = L(\bar{\boldsymbol{\theta}}(\boldsymbol{\psi}_k), \boldsymbol{\psi})$$

- Assume there exists function $\hat{\boldsymbol{\psi}}$ that maximizes the E-step

$$L(\boldsymbol{\theta}, \hat{\boldsymbol{\psi}}(\boldsymbol{\theta})) \geq L(\boldsymbol{\theta}, \boldsymbol{\psi}) \quad \forall \boldsymbol{\theta} \quad \forall \boldsymbol{\psi}$$

- Basically $\boldsymbol{\psi}_{k+1} = \hat{\boldsymbol{\psi}}(\bar{\boldsymbol{\theta}}(\boldsymbol{\psi}_k))$
- Define $\boldsymbol{\theta}_{k+1} = \bar{\boldsymbol{\theta}}(\boldsymbol{\psi}_k)$

Diagram of Function Compositions

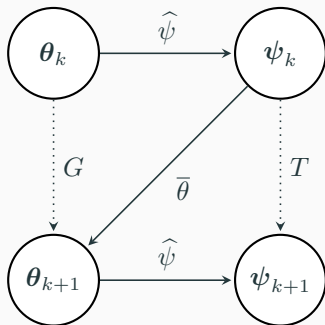
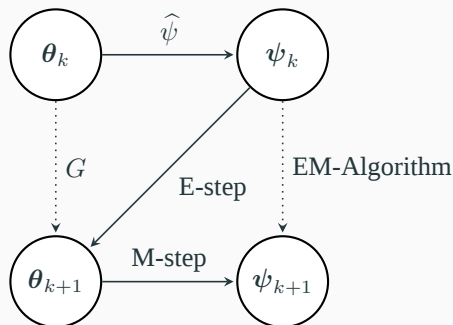


Diagram of Function Compositions



Can consider EM algorithm on **conditional expectations** θ_k 's instead of ψ_k 's.

$$\theta_{k+1} = G(\theta_k)$$

Aim: Find θ_* fixed point of G satisfying $G(\theta_*) = \theta_*$

Going back to Robbins-Monro

- Fixed point θ_* of G is unique root of $h(\theta) := G(\theta) - \theta$

Going back to Robbins-Monro

- Fixed point θ_* of G is unique root of $h(\theta) := G(\theta) - \theta$
- **Noisy** measurement of

$$h(\theta) = \mathbb{E}_{p(\mathbf{x}|\mathbf{y},\psi)} \left[\tilde{\theta}(\mathbf{X}) \right] - \theta$$

is a **Monte Carlo** estimate with $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$ drawn from $p(\mathbf{x} \mid \mathbf{y}, \psi)$

$$H(\theta, \mathbf{X}) = \frac{1}{m} \sum_{j=1}^m \tilde{\theta}(\mathbf{x}^{(j)}) - \theta$$

Going back to Robbins-Monro

- Fixed point θ_* of G is unique root of $h(\theta) := G(\theta) - \theta$
- Noisy measurement of

$$h(\theta) = \mathbb{E}_{p(\mathbf{x}|\mathbf{y},\psi)} \left[\tilde{\theta}(\mathbf{X}) \right] - \theta$$

is a Monte Carlo estimate with $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$ drawn from $p(\mathbf{x} \mid \mathbf{y}, \psi)$

$$H(\theta, \mathbf{X}) = \frac{1}{m} \sum_{j=1}^m \tilde{\theta}(\mathbf{x}^{(j)}) - \theta$$

- RM procedure converges to θ_* under regularity conditions.

$$\theta_k = \theta_{k-1} + \gamma_k \left(\frac{1}{m_k} \sum_{j=1}^{m_k} \tilde{\theta}(\mathbf{x}_k^{(j)}) - \theta_{k-1} \right)$$

where $\mathbf{x}_k^{(j)}$ are drawn from $p(\mathbf{x} \mid \mathbf{y}, \psi_{k-1})$

Going back to Robbins-Monro

- Fixed point θ_* of G is unique root of $h(\theta) := G(\theta) - \theta$
- Noisy measurement of

$$h(\theta) = \mathbb{E}_{p(\mathbf{x}|\mathbf{y},\psi)} \left[\tilde{\theta}(\mathbf{X}) \right] - \theta$$

is a Monte Carlo estimate with $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$ drawn from $p(\mathbf{x} | \mathbf{y}, \psi)$

$$H(\theta, \mathbf{X}) = \frac{1}{m} \sum_{j=1}^m \tilde{\theta}(\mathbf{x}^{(j)}) - \theta$$

- RM procedure converges to θ_* under regularity conditions.

$$\theta_k = \theta_{k-1} + \gamma_k \left(\frac{1}{m_k} \sum_{j=1}^{m_k} \tilde{\theta}(\mathbf{x}_k^{(j)}) - \theta_{k-1} \right)$$

where $\mathbf{x}_k^{(j)}$ are drawn from $p(\mathbf{x} | \mathbf{y}, \psi_{k-1})$

- Using linearity $Q(\psi | \psi_{k-1}) = -\xi(\psi) + \langle \theta_k, \phi(\psi) \rangle$ we get

$$\hat{Q}_k(\psi) = \hat{Q}_{k-1}(\psi) + \gamma_k \left(\frac{1}{m_k} \sum_{j=1}^{m_k} \log p(\mathbf{x}_k^{(j)}, \mathbf{y} | \psi) - \hat{Q}_{k-1}(\psi) \right)$$

Thank you



A. Benveniste, P. Priouret, and M. Métivier.

Adaptive Algorithms and Stochastic Approximations.

Springer-Verlag, Berlin, Heidelberg, 1990.



B. Delyon, M. Lavielle, and E. Moulines.

Convergence of a stochastic approximation version of the em algorithm.

Ann. Statist., 27(1):94–128, 03 1999.



A. P. Dempster, N. M. Laird, and D. B. Rubin.

Maximum likelihood from incomplete data via the em algorithm.

Journal of the Royal Statistical Society. Series B (Methodological), 39(1):1–38, 1977.



M. Duflo.

Random Iterative Models.

Springer-Verlag, Berlin, Heidelberg, 1st edition, 1997.



M. G. Gu and S. Li.

A stochastic approximation algorithm for maximum-likelihood estimation with incomplete data.

Canadian Journal of Statistics, 26(4):567–582, 1998.



H. Kushner and G. Yin.

Stochastic Approximation and Recursive Algorithms and Applications.

Stochastic Modelling and Applied Probability. Springer New York, 2003.



B. Polyak.

New stochastic approximation type procedures.

Avtomatica i Telemekhanika, 7:98–107, 01 1990.



H. Robbins and S. Monro.

A stochastic approximation method.

Ann. Math. Statist., 22(3):400–407, 09 1951.



H. Robbins and S. Monro.

A stochastic approximation method.

Ann. Math. Statist., 22(3):400–407, 09 1951.



C. F. J. Wu.

On the convergence properties of the em algorithm.

Ann. Statist., 11(1):95–103, 03 1983.