

# Introduction to Computational Statistics

First concepts in Stochastic Optimization – [juliette.chevallier@insa-toulouse.fr](mailto:juliette.chevallier@insa-toulouse.fr)

Spring School on Statistics and Machine Learning

---

## 1. Introduction & Motivation

## 2. Stochastic Gradient Descent

2.1 Intuition and First example

2.2 Practical Considerations

2.3 Stochastic Gradient Descent in High Dimension

## 3. EM Algorithm and Variants

3.1 Intuition and First example

3.2 Example: (Gaussian) Mixture Model

3.3 Convergence of the EM Algorithm

## 4. A Detour through Stochastic Approximation Theory

4.1 General Principle

4.2 Point-wise (Deterministic) Convergence

4.3 Robins-Monroe Algorithms

All materials for the course are available at

[plmlab.math.cnrs.fr/chevallier-teaching/hcmus-springschool-computationalstatistics](http://plmlab.math.cnrs.fr/chevallier-teaching/hcmus-springschool-computationalstatistics)



# Introduction to Computational Statistics

## Motivation: (Bayesian) inference

Given a **parametric model** and observations “from” this model,  
Estimate the parameters that **best fit the model to the data**

## Definition (**Estimator**)

Given a measurable space  $\mathcal{Y}$ , and the set of *admissible* parameters  $\Theta$

An estimator  $\hat{\theta}: \mathcal{Y} \rightarrow \Theta$  is a **function of the data**

**Example:** *Estimation of the mean and variance of  $\mathcal{N}(\mu, \sigma^2)$*

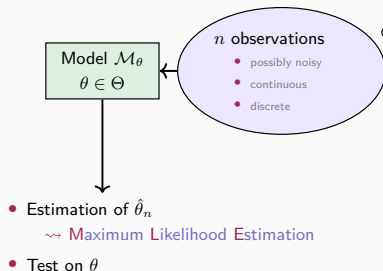
✓  $\bar{y}_n = \frac{1}{n} \sum_{i=1}^n y_i$  estimator of  $\mu$

✓  $\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}_n)^2$  estimator of  $\sigma^2$

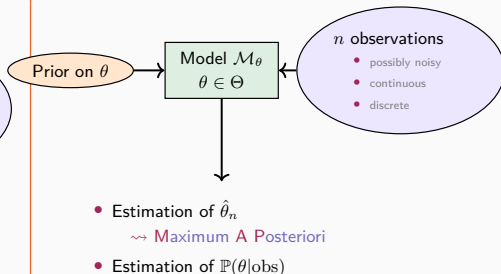
✗  $\frac{1}{n} \sum_{i=1}^n (y_i - \mu)^2$  *NOT* an estimator since it depends on  $\mu$

# Frequentist vs. Bayesian inference

## Frequentist



## Bayesian



*Combine both information, increase knowledge*

## When?

- When **additive information** is available, even weakly informative
- When  $\dim \theta \geq \# \text{observations}$   
 $\rightsquigarrow$  **Regularization**

**When?** No idea of possible prior  
 $\rightsquigarrow$  Better no prior than a wrong one

# Classical Estimators

## Some estimators:

1. Method of moments
2. Mean square error  $\rightsquigarrow$  Cf. linear model
3. Maximum likelihood or Maximum a posteriori (if Bayesian)

- 
- Let  $(f_\theta)_{\theta \in \Theta}$  be a family of pdf
  - Let  $(y_1, y_2, \dots, y_n)$  be an i.i.d sample with respect to unknown  $f_{\theta^*} \in (f_\theta)_{\theta \in \Theta}$

### Definition (Likelihood)

- Likelihood (of the observations)  $\mathcal{L}_n^{\text{MLE}} = \mathcal{L}_n \longleftrightarrow$  pdf of  $y_1, y_2, \dots, y_n$

$$\mathcal{L}_n: \theta \in \Theta \mapsto \mathcal{L}_n(\theta; y_1, y_2, \dots, y_n) = \prod_{i=1}^n f_\theta(y_i)$$

- Posterior likelihood  $\mathcal{L}_n^{\text{MAP}} \longleftrightarrow$  pdf of  $\theta | y_1, y_2, \dots, y_n$

$$\hat{\theta}^{\text{MLE}} \in \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_n^{\text{MLE}}(\theta)$$

$$\hat{\theta}^{\text{MAP}} \in \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_n^{\text{MAP}}(\theta)$$

**Remark:** MLE/MAP may not exist, may not be unique

- Several choice for the estimation of  $\theta$  given a  $n$ -sample

? *How can we choose the best one?*

↪ Quantify the quality/goodness-of-fit of the estimators and compare them

↪ **Loss function**

## Definition (Loss function)

Loss  $\ell(\theta^*, \hat{\theta}) \equiv$  function which quantifies the difference between  $\theta^*$  and  $\hat{\theta}$

- $\ell(a, b) \geq 0$
- $\ell(a, b) = \ell(b, a)$
- $\ell(a, a) = 0$
- $\ell(a, b) = 0 \implies a = b$
- $\mathcal{C}^0$ , differentiable, triangular inequality

**Remark:**  $\ell(\theta, \hat{\theta}(Y))$  is a random variable

↪ Even if the estimation is excellent,  $\ell(\theta, \hat{\theta}(Y))$  may be large

↪ Consider the **average loss**

## Definition (Risk function)

Risk  $R(\theta, \hat{\theta}) \equiv$  Average loss:

$$R(\theta, \hat{\theta}) = \mathbb{E}_{Y|\theta} [\ell(\theta, \hat{\theta}(Y))]$$

## Decision Theory in Practice

Let  $Y \sim \text{Bin}(100, \theta)$   $\theta \in [0, 1]$

Estimate quality measured by quadratic loss

- Naiv estimator:  $\hat{\theta}_1 = \frac{Y}{100}$

$\rightsquigarrow$  Quadratic loss  $\ell_1(\theta, \hat{\theta}_1) = \mathbb{E} \left[ \left( \theta - \frac{Y}{100} \right)^2 \right] = \frac{\theta(1-\theta)}{100}$

- Other choice:  $\hat{\theta}_2 = \frac{Y+3}{100}$

$\rightsquigarrow$  Quadratic loss  $\ell_2(\theta, \hat{\theta}_2) = \mathbb{E} \left[ \left( \theta - \frac{Y+3}{100} \right)^2 \right] = \frac{9}{100^2} + \frac{\theta(1-\theta)}{100}$

**Question:** Which estimator to choose?

**Problem:** The choice depends on the *unknown*  $\theta$

$\rightsquigarrow$  Consider the area under the curve

**Note:** It is equivalent to a Bayes uniform prior

$\rightsquigarrow$  If possible, choose a suitable prior!

# Contribution of the Bayesian Framework to Decision Theory

Note that [computations left to run]

$$\begin{aligned}\mathbb{E}_{Y|\theta} [\|\hat{\theta}(Y) - \theta\|^2] &= \mathbb{E}_{Y|\theta} \left[ \|\hat{\theta}(Y) - \mathbb{E}_{Y|\theta} [\hat{\theta}(Y)]\|^2 \right] + \|\mathbb{E}_{Y|\theta} [\hat{\theta}(Y)] - \theta\|^2 \\ &= \text{Var} (\hat{\theta}(Y)) + \text{Bias} (\hat{\theta}(Y))^2\end{aligned}$$

## Definition (Bayesian risk)

Given a prior  $\pi$  on  $\theta \in \Theta$ :

$$\hat{R}(\hat{\theta}) = \mathbb{E}_{\theta} [\mathbb{E}_{Y|\theta} [\ell(\theta, \hat{\theta}(Y))]]$$

Remark:

- $R(\theta, \hat{\theta}) = \int \ell(\theta, \hat{\theta}(Y)) \mathbb{P}_{\theta}(Y) dY$
- $\hat{R}(\hat{\theta}) = \int \int \ell(\theta, \hat{\theta}(Y)) \mathbb{P}_{\theta}(Y) \pi(\theta) dY d\theta$
- $\hat{R}$  does not depend on  $\theta$



## Quadratic Loss and Bayesian Prior

- Consider a discrete distribution, with discrete loss  $\ell(\theta, \hat{\theta}) = \mathbb{1}_{\theta \neq \hat{\theta}}$

Hence: *Bayesian risk*

$$\hat{R}(\hat{\theta}) = \sum_{\theta \in \Theta} \sum_{Y \in \mathcal{Y}} \ell(\theta, \hat{\theta}(Y)) \mathbb{P}(Y|\theta) \pi(\theta) = \sum_{\theta \in \Theta} \sum_{Y \in \mathcal{Y}} \ell(\theta, \hat{\theta}(Y)) \mathbb{P}(Y, \theta)$$

and [computations left to run]

$$\hat{\theta}(Y) \in \operatorname{argmin} \hat{R}(\hat{\theta}) \iff \hat{\theta}(Y) \in \operatorname{argmax} \mathbb{P}(\theta|Y) \rightsquigarrow \text{MAP}$$

- Likewise with  $L^2$  quadratic loss:  $\hat{\theta}(Y) = \mathbb{E}[\theta|Y] \rightsquigarrow \text{Mean A Posteriori}$

- 
- Previous result generalize well with improper prior
  - For estimate, one need to calculate either **expectations** or **maxima**, or **sample** to approximate these quantities (MLE, MAP, Mean a posteriori)

## Stochastic Gradient Descent

---

### 2.1 Intuition and First example

### 2.2 Practical Considerations

### 2.3 Stochastic Gradient Descent in High Dimension

## Reminder of Gradient Descent

Consider the problem  $\min_{\theta \in \Theta} J(\theta)$  with

- $\exists Y: (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow \mathcal{Y}$  random variable
- $\exists j: \Theta \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$

and  $J(\theta) = \mathbb{E}_Y[j(\theta, Y)]$



$$\min_{\theta \in \Theta} \mathbb{E}[j(\theta, Y)]$$

### Algorithm 1: Gradient Descent (GD)

**Data:**  $\theta_0 = \theta_{\text{init}} \in \Theta$

**Result:** Local minimizer  $\theta_{\text{end}} \in \Theta$

```
1 Set  $k = 0$ . Compute  $\nabla J(\theta_0)$ 
2 while  $\|\nabla J(\theta_k)\| \geq \varepsilon$  do
3    $\theta_{k+1} = \text{proj}_{\Theta} \{ \theta_k - \gamma \nabla J(\theta_k) \}$ 
4   Compute  $\nabla J(\theta_{k+1})$ 
5 return  $\theta_{k+1}$ 
```

- Typically,  $j(\theta, Y) := \ell(\theta, \hat{\theta}(Y))$
- Under regularity constraints for  $J$ , GD algorithm converges toward **local minimum** of  $J$
- Other deterministic minimization methods apply (e.g. Newton)  
**But:** All require the calculation  $\nabla J$ , involving an integral

↪ **Very long computation times**, especially for high-dimensional  $Y$

↪ **Solution:** Monte Carlo approach to computing  $\mathbb{E}$

---

**Algorithm 2:** Gradient Descent (GD)

---

**Data:**  $\theta_0 = \theta_{\text{init}} \in \Theta$ , sequence  $(\gamma_k)_k$

**Result:** Local minimizer  $\theta_{\text{end}} \in \Theta$

```
1 for  $k = 0 \rightarrow \text{maxIter}$  do  
2    $\theta_{k+1} = \text{proj}_{\Theta} \{ \theta_k - \gamma_k \nabla J(\theta_k) \}$   
3 return  $\theta_{k+1}$ 
```

---

---

**Algorithm 3:** Stochastic Gradient Descent (SGD)

---

**Data:**  $\theta_0 = \theta_{\text{init}} \in \Theta$ , sampler  $\mathbb{P}_Y$ , sequence  $(\gamma_k)_k$

**Result:** Local minimizer  $\theta_{\text{end}} \in \Theta$

```
1 for  $k = 0 \rightarrow \text{maxIter}$  do  
2   Sample  $y_k \sim \mathbb{P}_Y$   
3    $\theta_{k+1} = \text{proj}_{\Theta} \{ \theta_k - \gamma_k \frac{\partial j}{\partial \theta}(\theta_k, y_{k+1}) \}$   
4 return  $\theta_{k+1}$ 
```

---

# Intuition behind Stochastic Gradient Decent

1. Consider  $k_0 \in \mathbb{N}$
2. Sum the SGD formula  $k$  times from  $k_0$

$$\theta_{k+1} = \theta_k - \gamma_k \frac{\partial j}{\partial \theta}(\theta_k, y_{k+1}) \implies \theta_{k_0+k} = \theta_{k_0} - \sum_{\ell=0}^{k-1} \gamma_{k_0+\ell} \frac{\partial j}{\partial \theta}(\theta_{k_0+\ell}, y_{k_0+\ell+1})$$

Assume:

- $\theta \mapsto \frac{\partial j}{\partial \theta}(\theta, y)$  sufficiently regular
- $|\theta_{k_0+\ell} - \theta_{k_0+\ell'}|$  small (to be rigorously defined !)

Hence, according to **Cesàro mean lemma** (2nd line)

$$\begin{aligned} \theta_{k_0+k} &\simeq \theta_{k_0} - \sum_{\ell=0}^{k-1} \gamma_{k_0+\ell} \frac{\partial j}{\partial \theta}(\theta_{k_0}, y_{k_0+\ell+1}) \\ &\simeq \theta_{k_0} - \left( \sum_{\ell=0}^{k-1} \gamma_{k_0+\ell} \right) \nabla J(\theta_{k_0}) \end{aligned}$$

↪ One step of GD algorithm!

## Cesàro lemma

Let  $x_k$  be a sequence in an Hilbert space such that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n x_k = \mu$$

Let  $(\rho_k)_k$  a sequence of positive numbers  $\searrow 0$

Assume that

$$\varepsilon_k = (\rho_k - \rho_{k+1}) > 0 \text{ and } \sum_k \varepsilon_k \text{ diverge}$$

Then

$$\lim_{n \rightarrow \infty} \frac{\sum_{k=1}^n \rho_k x_k}{\sum_{k=1}^n \rho_k} = \mu$$

# Monte Carlo Sum and Stochastic Gradient

**Aim:** Compute the expectation  $\mathbb{E}[Y]$  of a random variable  $Y: \Omega \rightarrow \mathbb{R}$ ,  $Y \sim \mu$

- Approximation of  $\mathbb{E}[Y]$  by Monte Carlo sum:

- Let  $y_1, \dots, y_k \stackrel{i.i.d}{\sim} \mu$ . Denote  $\theta_k = \frac{1}{k} \sum_{\ell=1}^k y_\ell$
- According to law of large numbers,  $\theta_k \xrightarrow[k \rightarrow \infty]{p.s.} \mathbb{E}[Y]$

- Note that  $\frac{1}{k+1} = \frac{1}{k} - \frac{1}{k(k+1)}$ . Hence:

$$\begin{aligned}\theta_{k+1} &= \frac{1}{k+1} \sum_{\ell=1}^k y_\ell + \frac{y_{k+1}}{k+1} = \frac{1}{k} \sum_{\ell=1}^k y_\ell - \frac{1}{k+1} \left( \frac{1}{k} \sum_{\ell=1}^k y_\ell - y_{k+1} \right) \\ &= \theta_k - \frac{1}{k+1} (\theta_k - y_{k+1}) \\ &= \theta_k - \gamma_k \frac{\partial j}{\partial \theta}(\theta_k, y_{k+1}) \quad \text{where} \quad \gamma_k = \frac{1}{k+1} \quad \text{and} \quad j(\theta, y) = \frac{1}{2}(\theta - y)^2\end{aligned}$$

- However: Expectation of a random variable  $\equiv$  minimum value of the dispersion criterion of a point cloud  $\rightsquigarrow \mathbb{E}[Y] = \operatorname{argmin}_{\theta \in \mathbb{R}} \frac{1}{2}(\theta - Y)^2$

$\Rightarrow$  Monte Carlo approximation is a stochastic gradient on the dispersion function 13

### Remark:

- $\gamma_k$  converges to 0, but not too fast:  $\lim_{k \rightarrow +\infty} \gamma_k = 0$ , but  $\sum_k \gamma_k$  diverges
- Monte Carlo method converges p.s.  
Can we get the equivalent result for the stochastic gradient? Yes !
- There are central limit theorems (CLT) for Monte Carlo sums.  
What about stochastic gradient? Yes !

↪ Cf. *Stochastic Approximation theory*

## Stochastic Gradient Descent

---

2.1 Intuition and First example

**2.2 Practical Considerations**

2.3 Stochastic Gradient Descent in High Dimension



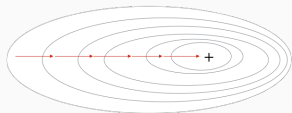
# Batch Size

- In all the above, we present a gradient step with **one simulation**  
But we could also run a Monte Carlo sum on **several simulations**  
     $\rightsquigarrow$  We refer to **batch size**

## Batch GD

At each iteration, gradient obtained by sampling  $n$  realizations  $y_{k+1,i}$

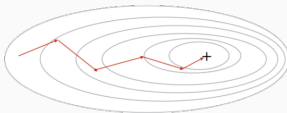
$$\theta_{k+1} = \theta_k - \frac{\gamma}{n} \sum_{i=1}^n \frac{\partial j}{\partial \theta}(\theta_k, y_{k+1,i})$$



## Mini-Batch GD

At each iteration, gradient obtained by sampling  $1 \leq m \leq n$  realizations  $y_{k+1,j}$

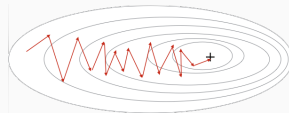
$$\theta_{k+1} = \theta_k - \frac{\gamma}{m} \sum_{j=1}^m \frac{\partial j}{\partial \theta}(\theta_k, y_{k+1,j})$$



## Stochastic GD

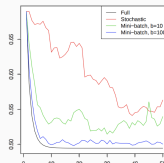
At each iteration, gradient obtained by sampling **one** realization  $y_{k+1}$

$$\theta_{k+1} = \theta_k - \gamma \frac{\partial j}{\partial \theta}(\theta_k, y_{k+1})$$



An **epoch** corresponds to the simulation of  $n$  observations

- Batch gradient descent: 1 iteration per epoch
- Mini-batch gradient descent:  $m$  iterations per epoch
- Stochastic gradient descent:  $n$  iterations per epoch



## Stop Criterion & Step Size

### Stop criterion:

- Cannot, as with gradient descent, rely on  $\|\theta_{k+1} - \theta_k\|$

As  $\gamma_k \rightarrow 0$  and  $|\frac{\partial j}{\partial \theta}| \leq n$ ,  $\|\theta_{k+1} - \theta_k\| \rightarrow 0$  by construction

- Nor on  $\frac{\partial j}{\partial \theta}$ , which is not informative about  $\nabla J$
  - An approximation of  $\mathbb{E} \left[ \frac{\partial j}{\partial \theta} \right]$  can be used
  - Most often a number is set by the user  $\rightsquigarrow$  `maxIter`
- 

### Step size:

- In theory, the highest asymptotic speed requires  $\gamma_k = \frac{1}{k}$
- But, in practice,  $\frac{1}{k^\alpha}$ ,  $\alpha \in ]\frac{1}{2}, 1[$  is sometimes better  
 $\rightsquigarrow$  It allows  $\frac{\partial j}{\partial \theta}$  to have a fairly high weight compared to the previous value.

Finally, we may also require a **burn-in time**, which consists in removing the first values (considered aberrant).

## Stochastic Gradient Descent

---

2.1 Intuition and First example

2.2 Practical Considerations

**2.3 Stochastic Gradient Descent in High Dimension**

# High-Dimensional Stochastic Gradient Descent

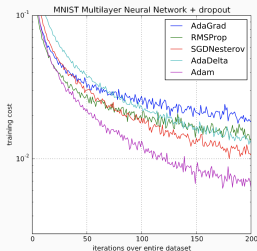
- Due to its **low numerical cost**, stochastic gradient is widely used in (very) large parametric models, such as **neural networks**

↪ Few possible/necessary adjustments to make SGD more practical in this context

- In high-dimensional parameter spaces, the topology of the objective function makes gradient descent difficult or even inefficient
  - Many local minima
  - Each “small” gradient  $\frac{\partial j}{\partial \theta}$  is uninformative compared to the “large” gradient  $\nabla J$
  - ...

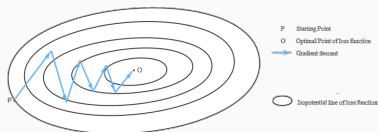
- Some suitable optimizers:

NAG, **AdaGrad**, Adadelata, **RMSProp**, **Adam**, AdaMax, Nadam, AMSGrad, AdamW, QHAdam, YellowFin, AggMo, QHM, Demon Adam, etc.



# Momentum

- The SGD updates the parameters after viewing only a subset of the training set
  - ↪ Add an inertia or **momentum** term:
    - Reduces variance
    - Limits oscillations along the convergence path
    - Avoids getting stuck too easily in a local minimum



**In practice:** we adapt the SGD algorithm to take account of previous gradients and smooth the update

$$\theta_{k+1} = \theta_k - \gamma \frac{\partial j}{\partial \theta}(\theta_k, y_{k+1}) \quad \rightsquigarrow \quad \begin{cases} v_{k+1} = \beta v_k - \gamma \frac{\partial j}{\partial \theta}(\theta_k, y_{k+1}) \\ \theta_{k+1} = \theta_k + v_k \end{cases}$$

- **Velocity**  $v$ : Direction in which parameters will be modified
- $\beta \in ]0, 1[$  quantifies the relative importance of previous gradients compared to the current one

Note:  $\beta \simeq 0.9$  in general

## SGD Enhancement Example: AdaGrad

- AdaGrad introduces a form of **learning rate adaptation** by accumulating the squares of the previous gradients
- Balance the power of gradients:
  - Smaller updates, *i.e.* low learning rates, for parameters associated with high gradients
  - And larger updates, *i.e.* high learning rates, for parameters associated with low gradients

### Implementation:

1. Compute the gradient  $g_{k+1} = \frac{\partial j}{\partial \theta}(\theta_k, y_{k+1})$
2. Gradient accumulation  $r_{k+1} = r_k + \|g_{k+1}\|$
3. Parameter update  $\theta_{k+1} = \theta_k - \frac{1}{\sqrt{r_{k+1}}} g_{k+1}$

## SGD Enhancement Example: RMSProp and Adam

- RMSProp is almost identical to AdaGrad, but the impact of older gradients is altered by a multiplicative coefficient  $\rho \in ]0, 1[$  (weight decay)

### Implementation:

1. Compute the gradient  $g_{k+1} = \frac{\partial j}{\partial \theta}(\theta_k, y_{k+1})$
2. Gradient accumulation  $r_{k+1} = \rho r_k + (1 - \rho) \|g_{k+1}\|$
3. Parameter update  $\theta_{k+1} = \theta_k - \frac{1}{\sqrt{r_{k+1}}} g_{k+1}$

- Adam is similar to RMSProp, but also adapts momentum

## EM Algorithm and Variants

---

### 3.1 Intuition and First example

### 3.2 Example: (Gaussian) Mixture Model

### 3.3 Convergence of the EM Algorithm



# Framework: Latent Variables Model

- **Aim:** Given a **latent variables model**, find the MLE (or MAP if Bayesian)

$$\theta^* \in \operatorname{argmax}_{\theta \in \Theta} q(y; \theta) = \operatorname{argmax}_{\theta \in \Theta} \mathbb{E}_Z [q(y, z; \theta)]$$

- **Model:** Latent or **Hierarchical** model

**Observations:**  $Y \leftrightarrow y = (y_i)_{i \in \llbracket 1, n \rrbracket} \in \mathcal{Y}$   
 $Y_i | Z_i; \theta \sim q(y_i | z_i; \theta)$

**Latent variables:**  $Z \leftrightarrow z = (z_i)_{i \in \llbracket 1, n \rrbracket} \in \mathcal{Z}$   
 $Z_i; \theta \sim q(z_i; \theta)$

$$\begin{cases} Y_i | Z_i; \theta \sim q(y_i | z_i; \theta) \\ Z_i; \theta \sim q(z_i; \theta) \\ \theta \sim q_{\text{prior}}(\theta) \end{cases}$$

**Parameter:**  $\theta \in \Theta$ ,  $\theta$  set of admissible parameters

If Bayesian framework,  $\theta \sim q_{\text{prior}}(\theta)$

- **Remarks**

- We observe only  $y$
- The law of  $Z$  may depend on  $\theta$
- Why EM instead of GDS?

To avoid the computation of the expectation of the gradient (*See in few minutes*)

# Intuitions behind the EM Algorithm

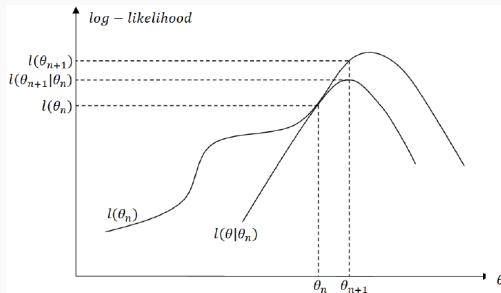
We seek  $\theta^* \in \operatorname{argmax}_{\theta \in \Theta} q(y; \theta)$ , equivalent to  $\theta^* \in \operatorname{argmax}_{\theta \in \Theta} \log q(y; \theta)$

**But:** Hard to derive due to the  $\mathbb{E}_Z \equiv \int_Z$

**An idea:** Construct a sequence  $(\theta_k)_k$  from an initial point  $\theta_0$  and

1. Find a function  $B$  such that  $B(\theta|\theta_k) \leq \log q(y; \theta)$  for all  $\theta$
2. Maximize  $B(\cdot|\theta_k)$  instead of  $\log q(y; \theta)$

↪ As soon as  $B$  is “optimal”, repeating step 1 & 2, we may hope a convergence toward **local minima** of  $q$



## Construction of the Minimizer $B$ (1/2)

- Let  $f_k$  be a **density function**, which may depend on the current value of  $\theta$

$$\log q(y; \theta) = \log \int_{\mathcal{Z}} q(y, z; \theta) dz = \log \int_{\mathcal{Z}} \frac{q(y, z; \theta)}{f_k(z)} f_k(z) dz$$

- According to the **Jensen inequality**,

$$\log \int_{\mathcal{Z}} \frac{q(y, z; \theta)}{f_k(z)} f_k(z) dz \geq \int_{\mathcal{Z}} \log \left( \frac{q(y, z; \theta)}{f_k(z)} \right) f_k(z) dz$$

- Since the Jensen convex inequality is optimal (we can reach the equality), we set

$$B(\theta | \theta_k) = \int_{\mathcal{Z}} \log \left( \frac{q(y, z; \theta)}{f_k(z)} \right) f_k(z) dz$$

- Aim:** Find the optimal density  $f_k$ , i.e. the one that “best” minimize  $\log q$ 
  - We seek  $f_k$  such that  $B(\theta_k | \theta_k)$  “touch”  $\log q(y; \theta_k)$
  - But:  $B(\cdot | \theta_k)$  minimizes  $\log q(y; \cdot)$ , so “touching” in  $\theta_k$  is equivalent for  $B(\cdot | \theta_k)$  to be maximal in  $\theta_k$   $\rightsquigarrow$  Find  $f_k$  such that  $B(\theta_k | \theta_k)$  is maximal

$$\implies f_k^* \in \operatorname{argmax}_{f_k} \int_{\mathcal{Z}} \log \left( \frac{q(y, z; \theta_k)}{f_k(z)} \right) f_k(z) dz \quad \text{s.t.} \quad f_k \geq 0 \quad \& \quad \int_{\mathcal{Z}} f_k(z) dz = 1$$

## “Reminder” on Functional differential

- **Lagrangian multiplier**  $G(f_k, \lambda) = \lambda(1 - G_1(f_k)) + G_2(f_k) - G_3(f_k)$ , with
  - $G_1(f_k) = \int_{\mathcal{Z}} f_k(z) \, dz$
  - $G_2(f_k) = \int_{\mathcal{Z}} \log q(y, z; \theta_k) f_k(z) \, dz$
  - $G_3(f_k) = \int_{\mathcal{Z}} \log f_x(z) f_k(z) \, dz$

### [Reminder ?] Functional differential/Directional derivative

- Let  $F(f)$  be a functional defined from a function  $f$
- Let  $\phi$  be a **test function** sufficiently regular

↪ *Differential of  $F$  in  $f$  in the direction  $\phi$ :*

$$dF[f, \phi] = \int_{\mathcal{X}} \frac{\partial F}{\partial f}(x) \phi(x) \, dx$$

$$dF[f, \phi] = \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} (F(f + \varepsilon \phi) - F(f)) = \left[ \frac{d}{d\varepsilon} F(f + \varepsilon \phi) \right]$$

## Construction of the Minimizer $B$ (2/2)

- **Lagrangian multiplier**  $G(f_k, \lambda) = \lambda(1 - G_1(f_k)) + G_2(f_k) - G_3(f_k)$ , with

$$G_1(f_k) = \int_{\mathcal{Z}} f_k(z) \, dz \quad G_2(f_k) = \int_{\mathcal{Z}} \log q(y, z; \theta_k) f_k(z) \, dz \quad G_3(f_k) = \int_{\mathcal{Z}} \log f_x(z) f_k(z) \, dz$$

$$\begin{cases} \frac{\partial G}{\partial f_t}(z) = -\lambda + \log q(y, z; \theta_k) - \log f_x(z) - 1 \end{cases} \quad (1)$$

$$\begin{cases} \frac{\partial G}{\partial \lambda}(z) = 1 - \int_{\mathcal{Z}} f_k(z) \, dz \end{cases} \quad (2)$$

$$(1) \iff \log f_k(z) = \log q(y, z; \theta_k) - \lambda - 1 \iff f_k(z) = e^{-\lambda-1} q(y, z; \theta_k)$$

$$(2) \iff 1 = \int_{\mathcal{Z}} e^{-\lambda-1} q(y, z; \theta_k) \, dz = e^{-\lambda-1} q(y; \theta_k) \iff e^{-\lambda-1} = \frac{1}{q(y; \theta_k)}$$

$$\implies f_k(z) = \frac{1}{q(y; \theta_k)} \times q(y, z; \theta_k) = q(z|y; \theta_k)$$

- We check that  $B(\theta_k|\theta_k) = q(y; \theta_k)$

$$B(\theta_k|\theta_k) = \int_{\mathcal{Z}} \log \left( \frac{q(y, z; \theta_k)}{f_k(z)} \right) q(z|y; \theta_k) \, dz = [\dots] = q(y; \theta_k)$$

## Maximization of $B \rightsquigarrow$ EM Algorithm

- **Maximization of  $B$ :**  $\implies \theta_{k+1} \in \underset{\theta \in \Theta}{\operatorname{argmax}} Q(\theta|\theta_k) + \log q(\theta)$

$$\begin{aligned} B(\theta|\theta_k) &= \mathbb{E}_{Z \sim q(\cdot|y;\theta_k)} [\log q(y, Z; \theta) - \log f_k(Z)] \\ &= \mathbb{E}_{Z \sim q(\cdot|y;\theta_k)} \left[ \log q(y, Z | \theta) + \underbrace{\log q(\theta)}_{\text{If Bayesian}} \right] + \underbrace{\mathcal{H}(f_k)}_{\text{Entropy of } f_k} \\ &= Q(\theta|\theta_k) + \log q(\theta) + \mathcal{H}(f_k) \end{aligned}$$

---

### Expectation-Maximization (EM) Algorithm

**E-step** Compute the conditional expectation

$$Q(\theta|\theta_k) = \mathbb{E}_{Z \sim q(\cdot|y;\theta_k)} [\log q(y, Z; \theta)]$$

**M-step** Maximize  $Q(\theta|\theta_k)$  in the feasible set:

$$\theta_{k+1} \in \underset{\theta \in \Theta}{\operatorname{argmax}} Q(\theta|\theta_k)$$

#### Remarks:

- ✓ We only need to compute **one** integral, no derivative (like in SGD)
- ✗ We have to compute an integral...

## EM Algorithm and Variants

---

3.1 Intuition and First example

**3.2 Example: (Gaussian) Mixture Model**

3.3 Convergence of the EM Algorithm

# Finite Mixture Model

**Mixture model of  $m$  components:** *Given:*

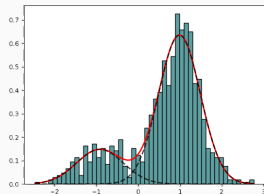
- $(\alpha_j)_{j \in \llbracket 1, m \rrbracket}$  the *proportions* of the mixture

$$\forall j \in \llbracket 1, m \rrbracket, \quad \alpha_j \in [0, 1] \quad \text{and} \quad \sum_{j=1}^m \alpha_j = 1;$$

- For all  $j$ ,  $f_j(\cdot; \omega_j)$  the *density* of the  $j$ -th sub-population, which (possibly) depends on a parameter  $\omega_j$ ; *and*
- $\theta = (\alpha_j, \omega_j)_{j \in \llbracket 1, m \rrbracket}$  the whole *parameters* of the mixture model.

We define:

$$q(\cdot; \theta) = \sum_{j=1}^m \alpha_j f_j(\cdot; \omega_j)$$





## Mixture Models vs. Latent Variables

**Mixture model:**  $q(\cdot; \theta) = \sum_{j=1}^m \alpha_j f_j(\cdot; \omega_j) \longleftrightarrow \begin{cases} \text{Proportions } \alpha_j \\ \text{Densities } f_j(\cdot; \omega_j) \end{cases}$

**Question:** How to generate data according to such a model?  $(y_i)_{i \in \llbracket 1, n \rrbracket}$

For all individual  $i \in \llbracket 1, n \rrbracket$ ,

- (i) Let  $\mathcal{P}_m = \{\mathcal{C}_1, \dots, \mathcal{C}_m\}$  be a partition of  $\llbracket 1, n \rrbracket$  into  $m$  classes, such that the individual  $i$  belongs to  $\mathcal{C}_j$  with probability  $\alpha_j$ :  $\mathbb{P}(i \in \mathcal{C}_j; \alpha_j) = \alpha_j$ ;
- (ii) Then,  $y_i$  is generated according to the density  $f_j(\cdot; \omega_j)$  if  $i \in \mathcal{C}_j$ .

$\rightsquigarrow$  **Latent variables**  $(z_i)_{i \in \llbracket 1, n \rrbracket}$  to encode the classes.

For all individual  $i \in \llbracket 1, n \rrbracket$ ,

$$\begin{cases} z_i \mid (\alpha_j)_{j \in \llbracket 1, m \rrbracket} \sim \sum_{k=1}^m \alpha_k \delta_k \\ y_i \mid z_i, \theta = (\alpha_j, \omega_j)_{j \in \llbracket 1, m \rrbracket} \sim f_{z_i}(\cdot; \omega_{z_i}) \end{cases}$$

# Hierarchical Writing of Mixture Models

- **Mixture Models:** For all  $i \in \llbracket 1, n \rrbracket$ , 
$$\begin{cases} z_i \mid (\alpha_j)_{j \in \llbracket 1, m \rrbracket} \sim \sum_{j=1}^m \alpha_j \delta_j, \\ y_i \mid z_i, (\alpha_j, \omega_j)_{j \in \llbracket 1, m \rrbracket} \sim f_{z_i}(\cdot; \omega_{z_i}). \end{cases}$$

- **Complete likelihood:** For all  $y = (y_i)_{i \in \llbracket 1, n \rrbracket}$  and  $z = (z_i)_{i \in \llbracket 1, n \rrbracket}$ ,

$$q(y, z; \theta) = \prod_{i=1}^n q(y_i, z_i; \theta) = \prod_{i=1}^n q(y_i | z_i; \theta) q(z_i; \theta) = \prod_{i=1}^n \alpha_{z_i} f_{z_i}(y_i; \omega_{z_i}).$$

- **Conditional likelihood:** For all  $i \in \llbracket 1, n \rrbracket$ ,

$$\begin{aligned} q(y_i; \theta) &= \sum_{j=1}^m q(y_i, \{z_i = j\}; \theta) \\ &= \sum_{j=1}^m q(y_i \mid \{z_i = j\}; \theta) q(\{z_i = j\}; \theta) = \sum_{j=1}^m \alpha_j f_j(y_i; \omega_j) \end{aligned}$$

## Parameters Estimation through the EM Algorithm

**E-step:** Compute the conditional expected log-likelihood

$$\begin{aligned} Q(\theta|\theta_k) &= \int_{\mathcal{Z}} \log q(y, z; \theta) q(z|y; \theta_k) dz \\ &= \mathbb{E}_{Z \sim q(\cdot | y; \theta_k)} [\log q(y, Z; \theta)] \end{aligned}$$

Here: 
$$Q(\theta|\theta_k) = \sum_{i=1}^n \sum_{j=1}^m [\log(\alpha_j) + \log(f_j(y_i; \omega_j))] \tau_{ij}^{(k)}$$

where  $\tau_{ij}^{(k)} = \mathbb{P}(Z_i = j | y_i; \theta_k) = \frac{\alpha_j^{(k)} f_j(y_i; \omega_j^{(k)})}{\sum_{\ell=1}^m \alpha_{\ell}^{(k)} f_{\ell}(y_i; \omega_{\ell}^{(k)})}$ .

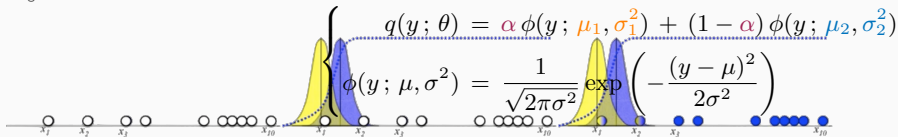
**M-step:** Maximize  $Q(\cdot | \theta_k)$  in the feasible set  $\theta$ :  $\theta_{k+1} \in \underset{\theta \in \Theta}{\operatorname{argmax}} Q(\theta | \theta_k)$

Here: 
$$\begin{cases} \forall j \in \llbracket 1, m \rrbracket, \quad \alpha_j^{(k+1)} = \frac{1}{n} \sum_{i=1}^n \tau_{ij}^{(k)} \\ (\omega_k^{(k+1)})_{j \in \llbracket 1, m \rrbracket} \in \underset{\omega = (\omega_j) \in \Omega}{\operatorname{argmax}} \sum_{i=1}^n \sum_{j=1}^m \tau_{ij}^{(k)} \log(f_j(y_i; \omega_j)) \end{cases}$$

# Example 1: One-dimensional Gaussian Mixture Model

**Likelihood:**  $\theta = (\alpha, \mu_1, \mu_2, \sigma_1^2, \sigma_2^2)$   $m = 2$

Images from Victor Lavrenko.



**E-step:**  $Q(\theta|\theta_k) \longleftrightarrow \begin{cases} \tau_{i1}^{(k)} \\ \tau_{i2}^{(k)} = 1 - \tau_{i1}^{(k)} \end{cases}$

$$\tau_{i1}^{(k)} = \frac{\alpha^{(k)} \phi(y; \mu_1^{(k)}, \sigma_1^{2(k)})}{\alpha^{(k)} \phi(y; \mu_1^{(k)}, \sigma_1^{2(k)}) + (1 - \alpha^{(k)}) \phi(y; \mu_2^{(k)}, \sigma_2^{2(k)})}$$

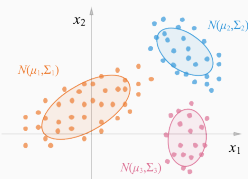
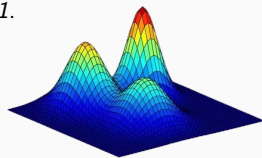
**M-step:**  $\alpha^{(k+1)} = \frac{1}{n} \sum_{i=1}^n \tau_{i1}^{(k)}$

$$\forall j \in \{1, 2\}, \quad \mu_j^{(k+1)} = \frac{\sum_{i=1}^n \tau_{ij}^{(k)} y_i}{\sum_{i=1}^n \tau_{ij}^{(k)}} \quad \& \quad \sigma_j^{2(k+1)} = \frac{\sum_{i=1}^n \tau_{ij}^{(k)} (y_i - \mu_j^{(k+1)})^2}{\sum_{i=1}^n \tau_{ij}^{(k)}}$$

# Multivariate Gaussian Mixtures

- **Quantitative data:**  $y_i \in \mathbb{R}^d$ , *Generalization of slide 31.*

- **Likelihood:**  $\theta = (\alpha_j, \mu_j, \Sigma_j)_{j \in \llbracket 1, m \rrbracket}$



$$\left\{ \begin{array}{l} q(y; \theta) = \sum_{j=1}^m \alpha_j \phi(y; \mu_j, \Sigma_j) . \\ \phi(y; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left( -\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right) . \end{array} \right.$$

- **Gaussian Mixture Model:**

$$\left\{ \begin{array}{l} \Theta = \left\{ (\alpha_j, \mu_j, \Sigma_j)_{j \in \llbracket 1, m \rrbracket} \in ([0, 1] \times \mathbb{R}^d \times \mathcal{S}_d \mathbb{R})^K \left| \sum_{j=1}^m \alpha_j = 1 \right. \right\} \\ \mathcal{M} = \{ \theta \in \Theta \mid y \in \mathbb{R}^d \mapsto q_K(x; \theta) \} \end{array} \right.$$

- **Estimation** through the EM algorithm (*See tutorials  $\rightsquigarrow$  Friday!*)

## EM Algorithm and Variants

---

3.1 Intuition and First example

3.2 Example: (Gaussian) Mixture Model

**3.3 Convergence of the EM Algorithm**

# Convergence of the EM Algorithm

- Several demonstrations of convergence, with assumptions that are more or less complicated to ensure in practice
- We state the version of [Delyon, Lavielle, Moulines], whose assumptions are often met

- (M1)
- $\Theta \subset \mathbb{R}^{n_\theta}$  open subset of  $\mathbb{R}^{n_\theta}$
  - $\exists S: \mathbb{R}^{n_z} \rightarrow \mathcal{S} \subset \mathbb{R}^{n_\theta}$  a Borelian function such that
    - Convex envelope  $\text{Conv}(S(\mathbb{R}^{n_z})) \subset \mathcal{S}$
    - $\forall \theta \in \Theta, \int_{\mathbb{R}^{n_z}} |S(y, z)| q(z|y; \theta) dz < +\infty$
    - $q(y, z; \theta) = \text{Exp}(-\psi(\theta) + \langle S(y, z) | \phi(\theta) \rangle)$

$\rightsquigarrow$  Exponential family,  $S \equiv$  Exhaustive statistics

(M2)  $\phi, \psi \in \mathcal{C}^2(\Theta)$

(M3)  $\bar{s}: \Theta \rightarrow \mathcal{S}$  s.t.  $\bar{s}(\theta) = \int_{\mathbb{R}^{n_z}} S(y, z) q(z|y; \theta) dz$  is  $\mathcal{C}^1(\Theta)$

(M4)  $\ell(\theta) = \log q(y; \theta)$  is  $\mathcal{C}^1(\Theta)$ . Moreover,  $\partial_\theta \int_{\mathbb{R}^{n_z}} \log q(y, z; \theta) dz = \int_{\mathbb{R}^{n_z}} \partial_\theta \log q(y, z; \theta) dz$

(M5) Let  $S: \mathcal{S} \times \Theta \rightarrow \mathbb{R}$  s.t.  $L(s, \theta) = \psi(\theta) + \langle S(y, z) | \phi(\theta) \rangle$

$\exists \hat{\theta}: \mathcal{S} \rightarrow \Theta$ , of class  $\mathcal{C}^1(\mathcal{S})$  and  $\forall s \in \mathcal{S}, \forall \theta \in \Theta, L(s, \hat{\theta}(s)) \geq L(s, \theta)$

## EM Convergence

Assume (M1–5).

Then for all  $\theta_0 \in \Theta$

- The sequence  $\ell(\theta_k)_k$  produced by the EM algorithm is increasing

- $\lim_{k \rightarrow \infty} d(\theta_k, \mathcal{L}) = 0$

$$\mathcal{L} = \{\theta \in \Theta \text{ s.t. } \partial_\theta \ell(\theta) = 0\}$$

# Convergence Considerations

## Remarks

- The sequence  $(\theta_k)_k$  is deterministic, and therefore depends on  $\theta_0$   
However, for a  $\theta_0$ , we may obtain a local minimum  
     $\rightsquigarrow$  Test several  $\theta_0$  and choose the one for which  $\ell(\theta_0)$  is maximum
- (M1) *Exponential family*: Many models fall into this category, including complex models  
     $\rightsquigarrow$  *Not a constraint*
- (M4) ensures regularity of the log-likelihood
- (M5) allows easy updating of  $\theta$  knowing  $S$
- Convergence toward a *stationary point of  $\ell$* : A saddle point, for example, may be encountered  
     $\rightsquigarrow$  *Convexity constraints*, at least local, are required to ensure that we reach a maximum



# Variants of the EM Algorithm

1. **Speeding-up** the EM Algorithm.

2. Limitations concerning the **M-step**.

**GEM:** Generalized EM Algorithm [?, ?]

*Idea:* Instead of maximize  $Q(\cdot|\theta_k)$  at each step, only find a point  $\theta_k$  “that makes it grow”

3. Limitations concerning the **E-step**.

**SEM:** Stochastic EM Algorithm [?]

**MCEM:** Monte-Carlo EM Algorithm [?]

**SAEM:** Stochastic-Approximation EM Algorithm [?]

*Idea:* Construct a **stochastic approximation** of  $Q(\cdot|\theta_k)$ , denoted  $Q_k$

# Stochastic Variants of the EM Algorithm

## SEM – Stochastic EM

**S-step:** Draw **one** observed sample

$$z_k \sim q(\cdot | y; \theta_k)$$

**“E”-step:** Estim. of  $Q(\cdot | \theta_k)$

$$Q_k(\theta) = \log q(y, z_k; \theta)$$

**M-step:** Maximize  $Q_{k+1}$ :

$$\theta_{k+1} \in \operatorname{argmax}_{\theta \in \Theta} Q_{k+1}(\theta)$$

- ✓ Very easy to install
- ✓ Randomness reduces the dependence on  $\theta_0$ , more general exploration of the modes of  $q(z|y; \theta)$
- ✗ Convergence proved on average only

## MCEM – Monte Carlo EM

**S-step:** Draw  **$m$**  samples

$$z_{k,j} \sim q(\cdot | y; \theta_k)$$

**“E”-step:** Monte-Carlo estim.

$$Q_k(\theta) = \frac{1}{m} \sum_{j=1}^m \log q(y, z_{k,j}; \theta)$$

**M-step:** Maximize  $Q_{k+1}$ :

$$\theta_{k+1} \in \operatorname{argmax}_{\theta \in \Theta} Q_{k+1}(\theta)$$

- ✗ Longer computation times
- ✗ No theoretical convergence results known (except in very specific cases)

## SAEM – Stochastic Approx.

**S-step:** Draw a sample

$$z_k \sim q(\cdot | y; \theta_k)$$

**SA-step:** Update  $Q_k(\theta)$  as

$$Q_{k+1}(\theta) = Q_k(\theta) + \gamma_k (\log q(y, z_k; \theta) - Q_k(\theta))$$

**M-step:** Maximize  $Q_{k+1}$ :

$$\theta_{k+1} \in \operatorname{argmax}_{\theta \in \Theta} Q_{k+1}(\theta)$$

- ✓ Convergence speed
- ✓ With a few more hypotheses (in particular about the  $\gamma_k$  sequence), theoretical convergence demonstrated

## A Detour through Stochastic Approximation Theory

---

### 4.1 General Principle

### 4.2 Point-wise (Deterministic) Convergence

### 4.3 Robins-Monroe Algorithms

# Stochastic Approximation

- A much broader framework than previously seen

Stochastic Gradient Descent, Stochastic EM

- General framework: We seek for  $\theta^*$  such that  $h(\theta^*) = \mathbb{E}_Y [H(\theta^*, Y)] = 0$ 
  - $H$  is known
  - The distribution  $\mathbb{P}_Y$  (which may depend on  $\theta^*$ ) is unknown

## Vocabulary

- The  $h$  function is called **mean field**
- The idea behind stochastic approximation is to determine  $\theta^*$  iteratively via a scheme of the form

$$\theta_n = \theta_{n-1} + \gamma_n h(\theta_{n-1}) + \gamma_n \eta_n$$

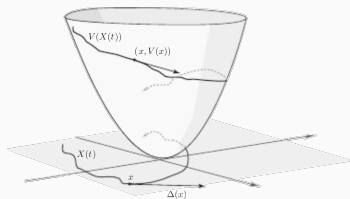
- $h(\theta) = \mathbb{E}_Y [H(\theta, Y)]$  is the function we're trying to cancel out
- $\eta_n = \mathbb{E}_Y [H(\theta_{n-1}, Y)] - h(\theta_{n-1})$  is a (**random**) perturbation

**Remark:** If  $\eta_n = 0$ , the schema is similar to that of the stochastic gradient

# Convergence Considerations

Two steps to show the convergence of the sequence  $(\theta_n)_n$

1. Find general conditions on the deterministic sequence of  $(\eta_n)_n$  and on  $h$  that ensure the deterministic convergence of  $(\theta_n)_n$
2. Show that these conditions are satisfied with probability 1, i.e. for almost any path of the noise process  $(\eta_1(\omega), \dots, \eta_n(\omega))$



$$\theta_n = \theta_{n-1} + \gamma_n h(\theta_{n-1}) + \gamma_n \eta_n$$

$$\longleftrightarrow \frac{\theta_n - \theta_{n-1}}{\gamma_n} = h(\theta_{n-1}) + \eta_n$$

↪ *Deterministic* convergence of scheme strongly linked to convergence of

mean-field equation

$$\frac{d\theta_t}{dt} = h(\theta_t)$$

## A Detour through Stochastic Approximation Theory

---

4.1 General Principle

4.2 Point-wise (Deterministic) Convergence

4.3 Robins-Monroe Algorithms

## Point-wise Convergence for Bounded $\theta$ : First Assumption

- Assume that  $\eta_n = e_n + r_n \rightsquigarrow$

$$\theta_n = \theta_{n-1} + \gamma_n h(\theta_n) + \gamma_n e_n + \gamma_n r_n$$

$$\gamma_n \geq 0$$

$$\lim_{n \rightarrow \infty} \sum_{i=1}^n \gamma_i = \infty$$

$$\lim_{n \rightarrow \infty} \gamma_n = 0$$

### Assumption (SA1) : Existence of a Lyapunov function

Let  $\mathcal{O} \subset \mathbb{R}^d$  an open set, with frontier denoted  $\partial\mathcal{O}$ .

- $h$  is a **continuous** vector field over  $\mathcal{O}$
- There exists a positive (or null)  $\mathcal{C}^1$  function  $V$  such that
  - $\forall x \in \mathcal{O}, \langle \nabla V(x) \mid h(x) \rangle \leq 0$
  - The set  $\mathcal{S} = \{x \mid \langle \nabla V(x) \mid h(x) \rangle = 0\}$  has an empty interior:  $\text{int}(V(\mathcal{S})) = \emptyset$

### Remark:

- We refer to the function  $V$  as a **Lyapunov function**
  - Not restrictive condition
  - For instance: If  $h = -\nabla Q$  with  $Q$  of class  $\mathcal{C}^1$ ,  $V = Q$  is an admissible choice
- Condition (ii) can be proved using the **Sard theorem** (*differential geometry*)  
 $\rightsquigarrow$  For all  $V \in \mathcal{C}^d$ ,  $V(\{\nabla V = 0\})$  has an empty interior

# Point-wise Convergence for Bounded $\theta$

## Definition ( $A$ -stable)

We say that algorithm  $\theta_n = \theta_{n-1} + \gamma_n h(\theta_n) + \gamma_n e_n + \gamma_n r_n$  is  $A$ -stable if:

- $\theta_n$  stays in a **compact** subset of  $\mathcal{O}$
- $\lim_{n \rightarrow \infty} \sum_{i=1}^n \gamma_i e_i < \infty$  and  $\lim_{n \rightarrow \infty} |r_n| = 0$

## Theorem (Point-wise Convergence)

Consider a  **$A$ -stable** algorithm and assume (SA1). **Then:**  $\lim_{n \rightarrow \infty} d(\theta_n, \mathcal{S}) = 0$

*Epecially, if  $|\mathcal{S}| < \infty$ ,  $\theta_n$  converges toward a point of  $\mathcal{S}$*

## Remark

- The  $A$ -stable condition is the baseline assumption for convergence  $\rightsquigarrow$  Ok !
- **But:** Assuming that  $\theta$  stays in a compact is **very restrictive** !  
 $\rightsquigarrow$  We can avoid this assumption by considering a sequence of growing compacts, and assuming that  $\theta_n$  does not diverge to  $\infty$

For more details: *Cf. the book/poly of Bernard Delyon*

Stochastic approximation with decreasing gain: Convergence and asymptotic theory (2000)



## Point-wise Convergence for **un**Bounded $\theta$

### Assumption (SA2) : Existence of a Lyapunov function

Let  $\mathcal{O} \subset \mathbb{R}^d$  an open set, with frontier denoted  $\partial\mathcal{O}$ .

- $h$  is a **continuous** vector field over  $\mathcal{O}$
- There exists  $K \subset \mathcal{O}$  compact
- There exists a positive (or null)  $\mathcal{C}^1$  function  $V$  such that
  - (i)  $V(x) \rightarrow +\infty$  if  $x \rightarrow \partial\mathcal{O}$  or  $|x| \rightarrow +\infty$
  - (ii)  $\langle \nabla V(x) \mid h(x) \rangle < 0$  if  $x \notin K$

### Theorem

**Assume:**

- Hypothesis (SA2)
- There exists a **compact**  $K_0 \subset \mathcal{O}$  such that  $\theta_n \in K_0$  an infinity of times
- For all  $N \in \mathbb{N}$ ,

$$\lim_{n \rightarrow \infty} \sum_{i=1}^n \gamma_i e_i \mathbb{1}_{\{V(\theta_{i-1}) \leq N\}} < \infty \quad \text{and} \quad \lim_{n \rightarrow \infty} |r_n| \mathbb{1}_{\{V(\theta_{i-1}) \leq N\}} = 0$$

**Then:** The algorithm is **A-stable**

## A Detour through Stochastic Approximation Theory

---

4.1 General Principle

4.2 Point-wise (Deterministic) Convergence

**4.3 Robins-Monroe Algorithms**

# Robins-Monroe Algorithms

*Let's go toward a random noise!*

## Robins-Monroe Stochastic Approximation

Let the scheme  $\theta_n = \theta_{n-1} + \gamma_n H(\theta_{n-1}, Y_n)$

- $Y_n \sim \mathbb{P}_{\theta_{n-1}}$  with  $\mathbb{P}(Y_n \in A | Y_{n-1}, Y_{n-2}, \dots, \theta_0) = \mathbb{P}_{\theta_{n-1}}(Y_n \in A)$
- The algorithm seeks to solve  $h(\theta) = \mathbb{E}_Y[H(\theta, Y)] = 0$

Rewrite the algorithm to make  $h$  explicit:

Set  $e_n = H(\theta_{n-1}, Y_n) - h(\theta_{n-1})$ ,

$$\theta_n = \theta_{n-1} + \gamma_n h(\theta_{n-1}) + \gamma_n e_n \quad (\text{RM})$$

**Example:** The **SAEM** and the **SGD** algorithms are of Robins-Monroe type

# Almost-Sure Convergence of Robins-Monroe Schemes

## Theorem

Assume:

(R1)  $\sum \gamma_n = +\infty$  and  $\sum \gamma_n^2 < \infty$

(R2) •  $h$  is *continuous*

•  $\mathcal{S} = \{\theta | h(\theta) = 0\}$  is finite

• There exists  $V$  of class  $\mathcal{C}^1$  such that

(i)  $\lim_{x \rightarrow \partial} V(x) = 0$  OR  $\theta_n$  remains in a *compact*

(ii)  $\langle \nabla V(\theta) | h(\theta) \rangle \leq 0$

(iii)  $\{ \langle \nabla V(\theta) | h(\theta) \rangle = 0 \} = \mathcal{S}$

(R3) For all compact  $K \subset \mathcal{O}$ ,  $\sup_{\theta \in K} \mathbb{E}_Y [\|H(\theta, Y)\|^2] < \infty$

Then: The (RM) algorithm converge toward  $\theta^*$  such that  $h(\theta^*) = 0$  with proba 1

# Convergence of Robins-Monroe Schemes

## *Proof (Compact case)*

- Compact case  $\implies$  For all  $n \in \mathbb{N}$ ,  $\theta_n \in K_0$
- **Aim:** Apply the previous theorem
  - (R1) and (R2) trivially induce the existence of a Lyapunov function (Made to that !)
  - $r_n = 0$  (and so  $\lim_{n \rightarrow \infty} |r_n| = 0$ )
  - What about  $\sum \gamma_n e_n$ ?

- Set  $X_n = \gamma_n e$  and  $S_n = \sum_{i=1}^n X_i$

$$\begin{aligned}\mathbb{E}_Y [|X_n|^2 \mid \mathcal{F}_{n-1}] &= \gamma_n^2 \mathbb{E}_Y [|e_n|^2 \mid \mathcal{F}_{n-1}] \\ &\leq \sup_{\theta \in K_0} \gamma_n^2 [|H(\theta, Y) - h(\theta)|^2 \mid \mathcal{F}_{n-1}]\end{aligned}$$

However,  $H$  and  $h$  continuous and therefore bounded on  $K_0 \oplus$  using (R3)

$$\mathbb{E}_Y [|X_n|^2 \mid \mathcal{F}_{n-1}] \leq cste \times \gamma_n^2 \implies \mathbb{E}_Y [|X_n|^2] \leq cste \times \gamma_n^2$$

Last, since  $\sum \gamma_n^2 < \infty$ ,  $S_n$  converges almost surely

□