

## Learning Optimization Algorithms with Average Case Convergence Rates



Peter Ochs  
Mathematical Optimization for Data Science  
Saarland University

— 07.05.2026 —



joint work: Michael Sucker

## Example:

$$\min_x f(x), \quad f(x) := \frac{1}{2} \|Ax - b\|^2.$$

**Example:**

$$\min_x f(x), \quad f(x) := \frac{1}{2} \|Ax - b\|^2.$$

**How do we solve the problem?**

## Example:

$$\min_x f(x), \quad f(x) := \frac{1}{2} \|Ax - b\|^2.$$

## How do we solve the problem?

- ◆ Inspect the properties of the problem.

Example: Smooth/Quadratic problem with  $L = \|A\|^2$ -Lipschitz gradient.

## Example:

$$\min_x f(x), \quad f(x) := \frac{1}{2} \|Ax - b\|^2.$$

## How do we solve the problem?

- ◆ Inspect the properties of the problem.

Example: Smooth/Quadratic problem with  $L = \|A\|^2$ -Lipschitz gradient.

- ◆ Embed the problem into a **class of problems** for which algorithms are available.

Example: Use Gradient Descent with step size  $\alpha = 1/L$

$$x^{(t+1)} = x^{(t)} - \alpha \nabla f(x^{(t)}).$$

*Worst case convergence guarantee:*

$$f(x^{(t)}) - \min f \leq O(1/t).$$

If we knew ...

## If we knew ...

- ◆ that for 99% of the problems that we have to solve the matrix  $A$  is diagonal

$$A = \begin{pmatrix} a_1 & 0 & \cdots & 0 \\ 0 & a_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & a_n \end{pmatrix} ?$$

## If we knew ...

- ◆ that for 99% of the problems that we have to solve the matrix  $A$  is diagonal

$$A = \begin{pmatrix} a_1 & 0 & \cdots & 0 \\ 0 & a_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & a_n \end{pmatrix} ?$$

- ◆ In that case, we could solve 99% of all problems directly

$$x_i^* = b_i/a_i, \quad i = 1, \dots, n.$$

## If we knew ...

- ◆ that for 99% of the problems that we have to solve the matrix  $A$  is diagonal

$$A = \begin{pmatrix} a_1 & 0 & \cdots & 0 \\ 0 & a_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & a_n \end{pmatrix} ?$$

- ◆ In that case, we could solve 99% of all problems directly

$$x_i^* = b_i/a_i, \quad i = 1, \dots, n.$$

- ◆ The worst-case does not explain what is usually observed in practice.

## If we knew ...

- ◆ that for 99% of the problems that we have to solve the matrix  $A$  is diagonal

$$A = \begin{pmatrix} a_1 & 0 & \cdots & 0 \\ 0 & a_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & a_n \end{pmatrix} ?$$

- ◆ In that case, we could solve 99% of all problems directly

$$x_i^* = b_i/a_i, \quad i = 1, \dots, n.$$

- ◆ The worst-case does not explain what is usually observed in practice.
- ◆ Sometimes the “**best**” class of problems is not obvious!

## If we knew ...

- ◆ that for 99% of the problems that we have to solve the matrix  $A$  is diagonal

$$A = \begin{pmatrix} a_1 & 0 & \cdots & 0 \\ 0 & a_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & a_n \end{pmatrix} ?$$

- ◆ In that case, we could solve 99% of all problems directly

$$x_i^* = b_i/a_i, \quad i = 1, \dots, n.$$

- ◆ The worst-case does not explain what is usually observed in practice.
- ◆ Sometimes the **“best” class of problems is not obvious!**

Can we construct an algorithm that adapts to hidden problem structures?

## Learning alleviates the bounds of analytic tractability by providing more:

- ◆ **Information:** Leverage more structure, *for example, statistical information.*
- ◆ **Automation:** Less “hand-crafting”.
- ◆ **Possibilities:** More building blocks.

## Learning alleviates the bounds of analytic tractability by providing more:

- ◆ **Information:** Leverage more structure, *for example, statistical information.*
- ◆ **Automation:** Less “hand-crafting”.
- ◆ **Possibilities:** More building blocks.

### **Our goal: Go beyond analytic tractability**

- ◆ Develop completely new algorithms.
- ◆ Define tight classes of problems.
- ◆ Break the barrier of worst-case estimates.
- ◆ We insist on having theoretical guarantees.

# Probabilistic Formulation of the Optimization Procedure

- ◆ Random parametric optimization problem with measurable  $\ell : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ :

$\min_{x \in \mathbb{R}^d} \ell(x, \theta)$ , random variable  $\theta : (\Omega, \mathfrak{A}, \mathbb{P}) \rightarrow \Theta$  defines the **class of problems**.

# Probabilistic Formulation of the Optimization Procedure

- ◆ Random parametric optimization problem with measurable  $\ell : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ :

$$\min_{x \in \mathbb{R}^d} \ell(x, \theta), \quad \text{random variable } \theta : (\Omega, \mathfrak{A}, \mathbb{P}) \rightarrow \Theta \text{ defines the class of problems.}$$

## Example:

- ◆ Regularized Inverse Problem:  $\ell(x, \lambda) = \frac{1}{2} \|Ax - b\|^2 + \lambda R(x)$ , i.e.  $\theta := \lambda$ ,  $\Theta = [0, 1]$ .

# Probabilistic Formulation of the Optimization Procedure

- ◆ Random parametric optimization problem with measurable  $\ell : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ :

$$\min_{x \in \mathbb{R}^d} \ell(x, \theta), \quad \text{random variable } \theta : (\Omega, \mathfrak{A}, \mathbb{P}) \rightarrow \Theta \text{ defines the class of problems.}$$

- ◆ Use a parametric (iterative) optimization algorithm with  $t = 0, 1, 2, \dots$ :

$$\mathcal{A} : \mathcal{H} \times \Theta \times \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad (\alpha, \theta, x^{(t)}) \mapsto \mathcal{A}(\alpha, \theta, x^{(t)}) =: x^{(t+1)}$$

## Example:

- ◆ Regularized Inverse Problem:  $\ell(x, \lambda) = \frac{1}{2} \|Ax - b\|^2 + \lambda R(x)$ , i.e.  $\theta := \lambda$ ,  $\Theta = [0, 1]$ .

# Probabilistic Formulation of the Optimization Procedure

- ◆ Random parametric optimization problem with measurable  $\ell : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ :

$$\min_{x \in \mathbb{R}^d} \ell(x, \theta), \quad \text{random variable } \theta : (\Omega, \mathfrak{A}, \mathbb{P}) \rightarrow \Theta \text{ defines the class of problems.}$$

- ◆ Use a parametric (iterative) optimization algorithm with  $t = 0, 1, 2, \dots$ :

$$\mathcal{A} : \mathcal{H} \times \Theta \times \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad (\alpha, \theta, x^{(t)}) \mapsto \mathcal{A}(\alpha, \theta, x^{(t)}) =: x^{(t+1)}$$

## Example:

- ◆ Regularized Inverse Problem:  $\ell(x, \lambda) = \frac{1}{2} \|Ax - b\|^2 + \lambda R(x)$ , i.e.  $\theta := \lambda$ ,  $\Theta = [0, 1]$ .
- ◆ Preconditioned GD:  $x^{(t+1)} = x^{(t)} - P \nabla \ell(x^{(t)}, \theta)$ , i.e.  $\alpha := P$ ,  $\mathcal{H} := \mathbb{R}^{d \times d}$ .

# Probabilistic Formulation of the Optimization Procedure

- ◆ Random parametric optimization problem with measurable  $\ell : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ :

$$\min_{x \in \mathbb{R}^d} \ell(x, \theta), \quad \text{random variable } \theta : (\Omega, \mathfrak{A}, \mathbb{P}) \rightarrow \Theta \text{ defines the class of problems.}$$

- ◆ Use a parametric (iterative) optimization algorithm with  $t = 0, 1, 2, \dots$ :

$$\mathcal{A} : \mathcal{H} \times \Theta \times \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad (\alpha, \theta, x^{(t)}) \mapsto \mathcal{A}(\alpha, \theta, x^{(t)}) =: x^{(t+1)}$$

- ◆ Stop the algorithm, if  $x^{(t)}$  satisfies a stopping criterion  $C_\varepsilon(\theta)$  and define stopping time

$$\tau_\alpha(\theta, x) := \inf\{t \in \mathbb{N} \mid x_{\alpha, \theta}^{(t)} \in C_\varepsilon(\theta)\}, \quad \text{where } x := (x_{\alpha, \theta}^{(t)})_{t \in \mathbb{N}}.$$

## Example:

- ◆ Regularized Inverse Problem:  $\ell(x, \lambda) = \frac{1}{2} \|Ax - b\|^2 + \lambda R(x)$ , i.e.  $\theta := \lambda$ ,  $\Theta = [0, 1]$ .
- ◆ Preconditioned GD:  $x^{(t+1)} = x^{(t)} - P \nabla \ell(x^{(t)}, \theta)$ , i.e.  $\alpha := P$ ,  $\mathcal{H} := \mathbb{R}^{d \times d}$ .

# Probabilistic Formulation of the Optimization Procedure

- ◆ Random parametric optimization problem with measurable  $\ell : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ :

$$\min_{x \in \mathbb{R}^d} \ell(x, \theta), \quad \text{random variable } \theta : (\Omega, \mathfrak{A}, \mathbb{P}) \rightarrow \Theta \text{ defines the class of problems.}$$

- ◆ Use a parametric (iterative) optimization algorithm with  $t = 0, 1, 2, \dots$ :

$$\mathcal{A} : \mathcal{H} \times \Theta \times \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad (\alpha, \theta, x^{(t)}) \mapsto \mathcal{A}(\alpha, \theta, x^{(t)}) =: x^{(t+1)}$$

- ◆ Stop the algorithm, if  $x^{(t)}$  satisfies a stopping criterion  $C_\varepsilon(\theta)$  and define stopping time

$$\tau_\alpha(\theta, x) := \inf\{t \in \mathbb{N} \mid x_{\alpha, \theta}^{(t)} \in C_\varepsilon(\theta)\}, \quad \text{where } x := (x_{\alpha, \theta}^{(t)})_{t \in \mathbb{N}}.$$

## Example:

- ◆ Regularized Inverse Problem:  $\ell(x, \lambda) = \frac{1}{2} \|Ax - b\|^2 + \lambda R(x)$ , i.e.  $\theta := \lambda$ ,  $\Theta = [0, 1]$ .
- ◆ Preconditioned GD:  $x^{(t+1)} = x^{(t)} - P \nabla \ell(x^{(t)}, \theta)$ , i.e.  $\alpha := P$ ,  $\mathcal{H} := \mathbb{R}^{d \times d}$ .
- ◆ Check if  $x^{(t)}$  satisfies  $\|\nabla \ell(x^{(t)}, \theta)\| \leq \varepsilon$ .

# Probabilistic Formulation of the Optimization Procedure

- ◆ **Random parametric optimization problem** with measurable  $\ell : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ :

$$\min_{x \in \mathbb{R}^d} \ell(x, \theta), \quad \text{random variable } \theta : (\Omega, \mathfrak{A}, \mathbb{P}) \rightarrow \Theta \text{ defines the **class of problems** .}$$

- ◆ Use a **parametric (iterative) optimization algorithm** with  $t = 0, 1, 2, \dots$ :

$$\mathcal{A} : \mathcal{H} \times \Theta \times \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad (\alpha, \theta, x^{(t)}) \mapsto \mathcal{A}(\alpha, \theta, x^{(t)}) =: x^{(t+1)}$$

- ◆ Stop the algorithm, if  $x^{(t)}$  satisfies a **stopping criterion**  $C_\varepsilon(\theta)$  and define **stopping time**

$$\tau_\alpha(\theta, x) := \inf\{t \in \mathbb{N} \mid x_{\alpha, \theta}^{(t)} \in C_\varepsilon(\theta)\}, \quad \text{where } \mathbf{x} := (x_{\alpha, \theta}^{(t)})_{t \in \mathbb{N}}.$$

- ◆ For best **expected case complexity** w.r.t. choice of algorithm  $\mathcal{A}$ , solve

$$\min_{\alpha \in \mathcal{H}} \left( \mathbb{E}_{\theta, x | \alpha} \{ \tau_\alpha(\theta, x) \} \right).$$

## Example:

- ◆ **Regularized Inverse Problem**:  $\ell(x, \lambda) = \frac{1}{2} \|Ax - b\|^2 + \lambda R(x)$ , i.e.  $\theta := \lambda$ ,  $\Theta = [0, 1]$ .
- ◆ **Preconditioned GD**:  $x^{(t+1)} = x^{(t)} - P \nabla \ell(x^{(t)}, \theta)$ , i.e.  $\alpha := P$ ,  $\mathcal{H} := \mathbb{R}^{d \times d}$ .
- ◆ Check if  $x^{(t)}$  satisfies  $\|\nabla \ell(x^{(t)}, \theta)\| \leq \varepsilon$ .

# Probabilistic Formulation of the Optimization Procedure

- ◆ **Random parametric optimization problem** with measurable  $\ell : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ :

$$\min_{x \in \mathbb{R}^d} \ell(x, \theta), \quad \text{random variable } \theta : (\Omega, \mathfrak{A}, \mathbb{P}) \rightarrow \Theta \text{ defines the **class of problems** .}$$

- ◆ Use a **parametric (iterative) optimization algorithm** with  $t = 0, 1, 2, \dots$ :

$$\mathcal{A} : \mathcal{H} \times \Theta \times \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad (\alpha, \theta, x^{(t)}) \mapsto \mathcal{A}(\alpha, \theta, x^{(t)}) =: x^{(t+1)}$$

- ◆ Stop the algorithm, if  $x^{(t)}$  satisfies a **stopping criterion**  $C_\varepsilon(\theta)$  and define **stopping time**

$$\tau_\alpha(\theta, \mathbf{x}) := \inf\{t \in \mathbb{N} \mid x_{\alpha, \theta}^{(t)} \in C_\varepsilon(\theta)\}, \quad \text{where } \mathbf{x} := (x_{\alpha, \theta}^{(t)})_{t \in \mathbb{N}}.$$

- ◆ For best **expected case complexity** w.r.t. choice of algorithm  $\mathcal{A}$ , solve

$$\min_{\alpha \in \mathcal{H}} \left( \mathbb{E}_{\theta, \mathbf{x} \mid \alpha} \{ \tau_\alpha(\theta, \mathbf{x}) \} \right).$$

Using **expectation is intractable** in general!

Routine in **Machine Learning**:  
**Empirical Risk Minimization**

## Routine in **Machine Learning**: **Empirical Risk Minimization**

- ◆ Want to minimize the **risk**  $\mathcal{R}(\alpha)$  defined as the **expected loss**:

$$\min_{\alpha \in \mathcal{H}} \mathcal{R}(\alpha), \quad \mathcal{R}(\alpha) := \mathbb{E}[f(\alpha, \theta)].$$

*Deep Learning notation:*

$\alpha$	$\sim$	<i>neural network parameters</i>
$f$	$\sim$	<i>network architecture including loss</i>
$\theta$	$\sim$	<i>data distribution</i>

## Routine in **Machine Learning**: **Empirical Risk Minimization**

- ◆ Want to minimize the **risk**  $\mathcal{R}(\alpha)$  defined as the **expected loss**:

$$\min_{\alpha \in \mathcal{H}} \mathcal{R}(\alpha), \quad \mathcal{R}(\alpha) := \mathbb{E}[f(\alpha, \theta)].$$

This is intractable since the distribution of  $\theta$  is **unknown**.

*Deep Learning notation:*

$\alpha$	$\sim$	neural network parameters
$f$	$\sim$	network architecture including loss
$\theta$	$\sim$	data distribution

## Routine in **Machine Learning**: **Empirical Risk Minimization**

- Want to minimize the **risk**  $\mathcal{R}(\alpha)$  defined as the **expected loss**:

$$\min_{\alpha \in \mathcal{H}} \mathcal{R}(\alpha), \quad \mathcal{R}(\alpha) := \mathbb{E}[f(\alpha, \theta)].$$

*Deep Learning notation:*

$\alpha$	$\sim$	neural network parameters
$f$	$\sim$	network architecture including loss
$\theta$	$\sim$	data distribution

This is intractable since the distribution of  $\theta$  is **unknown**.

- Hence, resort to minimizing the **empirical risk**  $\hat{\mathcal{R}}(\alpha, \theta_{[N]})$  over some dataset  $\theta_{[N]} := \{\theta_i\}_{i=1}^N$ :

$$\min_{\alpha \in \mathcal{H}} \hat{\mathcal{R}}(\alpha, \theta_{[N]}), \quad \hat{\mathcal{R}}(\alpha, \theta_{[N]}) := \frac{1}{N} \sum_{i=1}^N f(\alpha, \theta_i).$$

## Routine in **Machine Learning**: **Empirical Risk Minimization**

- Want to minimize the **risk**  $\mathcal{R}(\alpha)$  defined as the **expected loss**:

$$\min_{\alpha \in \mathcal{H}} \mathcal{R}(\alpha), \quad \mathcal{R}(\alpha) := \mathbb{E}[f(\alpha, \theta)].$$

*Deep Learning notation:*

$\alpha$	$\sim$	neural network parameters
$f$	$\sim$	network architecture including loss
$\theta$	$\sim$	data distribution

This is intractable since the distribution of  $\theta$  is **unknown**.

- Hence, resort to minimizing the **empirical risk**  $\hat{\mathcal{R}}(\alpha, \theta_{[N]})$  over some dataset  $\theta_{[N]} := \{\theta_i\}_{i=1}^N$ :

$$\min_{\alpha \in \mathcal{H}} \hat{\mathcal{R}}(\alpha, \theta_{[N]}), \quad \hat{\mathcal{R}}(\alpha, \theta_{[N]}) := \frac{1}{N} \sum_{i=1}^N f(\alpha, \theta_i).$$

Is the performance on  $\hat{\mathcal{R}}$  representative for the overall performance  $\mathcal{R}$ ?

# Is the performance on $\hat{\mathcal{R}}$ representative for $\mathcal{R}$ ?

## Is the performance on $\hat{\mathcal{R}}$ representative for $\mathcal{R}$ ?

- ◆ Yes, if we have **uniform generalization bounds**, i.e. bounds of the form:  $\forall \epsilon > 0$ :

$$\mathbb{P}\{\mathcal{R}(\alpha_{\theta_{[N]}^*}) \leq \inf_{\alpha \in \mathcal{H}} \hat{\mathcal{R}}(\alpha, \theta_{[N]}) + K(N, \alpha, \epsilon)\} \geq 1 - \epsilon.$$

## Is the performance on $\hat{\mathcal{R}}$ representative for $\mathcal{R}$ ?

- ◆ Yes, if we have **uniform generalization bounds**, i.e. bounds of the form:  $\forall \epsilon > 0$ :

$$\mathbb{P}\{\mathcal{R}(\alpha_{\theta_{[N]}^*}) \leq \inf_{\alpha \in \mathcal{H}} \hat{\mathcal{R}}(\alpha, \theta_{[N]}) + K(N, \alpha, \epsilon)\} \geq 1 - \epsilon.$$

## Is the performance on $\hat{\mathcal{R}}$ representative for $\mathcal{R}$ ?

- ◆ Yes, if we have **uniform generalization bounds**, i.e. bounds of the form:  $\forall \epsilon > 0$ :

$$\mathbb{P}\{\mathcal{R}(\alpha_{\theta_{[N]}^*}) \leq \inf_{\alpha \in \mathcal{H}} \hat{\mathcal{R}}(\alpha, \theta_{[N]}) + K(N, \alpha, \epsilon)\} \geq 1 - \epsilon.$$

## Is the performance on $\hat{\mathcal{R}}$ representative for $\mathcal{R}$ ?

- ◆ Yes, if we have **uniform generalization bounds**, i.e. bounds of the form:  $\forall \epsilon > 0$ :

$$\mathbb{P}\{\mathcal{R}(\alpha_{\theta_{[N]}^*}) \leq \inf_{\alpha \in \mathcal{H}} \hat{\mathcal{R}}(\alpha, \theta_{[N]}) + K(N, \alpha, \epsilon)\} \geq 1 - \epsilon.$$

## Is the performance on $\hat{\mathcal{R}}$ representative for $\mathcal{R}$ ?

- ◆ Yes, if we have **uniform generalization bounds**, i.e. bounds of the form:  $\forall \epsilon > 0$ :

$$\mathbb{P}\{\mathcal{R}(\alpha_{\theta_{[N]}^*}) \leq \inf_{\alpha \in \mathcal{H}} \hat{\mathcal{R}}(\alpha, \theta_{[N]}) + K(N, \alpha, \epsilon)\} \geq 1 - \epsilon.$$

## Is the performance on $\hat{\mathcal{R}}$ representative for $\mathcal{R}$ ?

- Yes, if we have **uniform generalization bounds**, i.e. bounds of the form:  $\forall \epsilon > 0$ :

$$\mathbb{P}\{\mathcal{R}(\alpha_{\theta_{[N]}^*}) \leq \inf_{\alpha \in \mathcal{H}} \hat{\mathcal{R}}(\alpha, \theta_{[N]}) + K(N, \alpha, \epsilon)\} \geq 1 - \epsilon.$$

- Such bounds are called **PAC-bounds**, which is an acronym for:

Probably      Approximately      Correct .

With high probability, the empirical risk is close to the true risk.

## Is the performance on $\hat{\mathcal{R}}$ representative for $\mathcal{R}$ ?

- Yes, if we have **uniform generalization bounds**, i.e. bounds of the form:  $\forall \epsilon > 0$ :

$$\mathbb{P}\{\mathcal{R}(\alpha_{\theta_{[N]}^*}) \leq \inf_{\alpha \in \mathcal{H}} \hat{\mathcal{R}}(\alpha, \theta_{[N]}) + K(N, \alpha, \epsilon)\} \geq 1 - \epsilon.$$

- Such bounds are called **PAC-bounds**, which is an acronym for:

Probably      Approximately      Correct .  
With high probability, the empirical risk is close to the true risk.

**PAC-Bayes** extends this to the Bayes-risk:

## Is the performance on $\hat{\mathcal{R}}$ representative for $\mathcal{R}$ ?

- Yes, if we have **uniform generalization bounds**, i.e. bounds of the form:  $\forall \epsilon > 0$ :

$$\mathbb{P}\{\mathcal{R}(\alpha_{\theta_{[N]}^*}) \leq \inf_{\alpha \in \mathcal{H}} \hat{\mathcal{R}}(\alpha, \theta_{[N]}) + K(N, \alpha, \epsilon)\} \geq 1 - \epsilon.$$

- Such bounds are called **PAC-bounds**, which is an acronym for:

Probably      Approximately      Correct .  
With high probability, the empirical risk is close to the true risk.

### PAC-Bayes extends this to the Bayes-risk:

- Such bounds hold for **distributions**  $\rho \in \mathcal{M}(\mathbb{P}_\alpha)$ :

$$\mathbb{P}\{\rho_{\theta_{[N]}^*}[\mathcal{R}(\alpha)] \leq \inf_{\rho \in \mathcal{M}(\mathbb{P}_\alpha)} \rho[\hat{\mathcal{R}}(\alpha, \theta_{[N]})] + K(\rho, N, \epsilon)\} \geq 1 - \epsilon,$$

where  $\mathcal{M}(\mathbb{P}_\alpha)$  denotes some class of (probability) measures on  $\mathcal{H}$  that satisfy a certain property w.r.t. the **distribution**  $\mathbb{P}_\alpha$ .

## Is the performance on $\hat{\mathcal{R}}$ representative for $\mathcal{R}$ ?

- Yes, if we have **uniform generalization bounds**, i.e. bounds of the form:  $\forall \epsilon > 0$ :

$$\mathbb{P}\{\mathcal{R}(\alpha_{\theta_{[N]}^*}) \leq \inf_{\alpha \in \mathcal{H}} \hat{\mathcal{R}}(\alpha, \theta_{[N]}) + K(N, \alpha, \epsilon)\} \geq 1 - \epsilon.$$

- Such bounds are called **PAC-bounds**, which is an acronym for:

Probably      Approximately      Correct .  
With high probability, the empirical risk is close to the true risk.

### PAC-Bayes extends this to the Bayes-risk:

- Such bounds hold for **posterior distributions**  $\rho \in \mathcal{M}(\mathbb{P}_\alpha)$ :

$$\mathbb{P}\{\rho_{\theta_{[N]}^*}^*[\mathcal{R}(\alpha)] \leq \inf_{\rho \in \mathcal{M}(\mathbb{P}_\alpha)} \rho[\hat{\mathcal{R}}(\alpha, \theta_{[N]})] + K(\rho, N, \epsilon)\} \geq 1 - \epsilon,$$

where  $\mathcal{M}(\mathbb{P}_\alpha)$  denotes some class of (probability) measures on  $\mathcal{H}$  that satisfy a certain property w.r.t. the **prior distribution**  $\mathbb{P}_\alpha$ .

*This is a naming convention! Not to be confused with prior and posterior in Bayesian analysis, which are linked by a likelihood.*

**PAC-Bayes** [Alquier '21]

**Learning-to-Optimize** [Chen et al. '21]



**PAC-Bayesian Learning of Optimization Algorithms**

[Sucker, O. '22], [Sucker, Fadili, O. '24]

**A Markovian Model for Learning-to-Optimize**

[Sucker, O. '24]

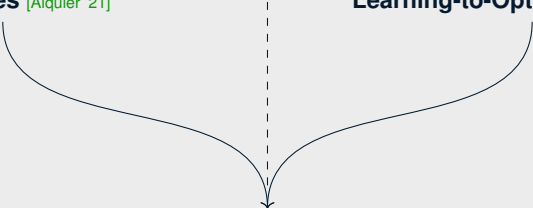
**A Generalization Result for Convergence in Learning-to-Optimize**

[Sucker, O. '24]

- 
- ◆ [Alquier '21]: “User-friendly introduction to PAC-Bayes bounds”, Foundations and Trends in Machine Learning 17(2):174–303, 2024 (arXiv 2021).
  - ◆ [Chen et al. '22]: “Learning to optimize: A primer and a benchmark”, Journal of Machine Learning Research 23(1):8562–8620, 2022 (arXiv 2021).

**PAC-Bayes** [Alquier '21]

**Learning-to-Optimize** [Chen et al. '21]



**PAC-Bayesian Learning of Optimization Algorithms**

[Sucker, O. '22], [Sucker, Fadili, O. '24]

**A Markovian Model for Learning-to-Optimize (next slides !)**

[Sucker, O. '24]

**A Generalization Result for Convergence in Learning-to-Optimize (later !)**

[Sucker, O. '24]

- 
- ◆ [Alquier '21]: “User-friendly introduction to PAC-Bayes bounds”, Foundations and Trends in Machine Learning 17(2):174–303, 2024 (arXiv 2021).
  - ◆ [Chen et al. '22]: “Learning to optimize: A primer and a benchmark”, Journal of Machine Learning Research 23(1):8562–8620, 2022 (arXiv 2021).

# Probabilistic Formulation of the Optimization Procedure

- ◆ **Random parametric optimization problem** with measurable  $\ell : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ :

$$\min_{x \in \mathbb{R}^d} \ell(x, \theta), \quad \text{random variable } \theta : (\Omega, \mathfrak{A}, \mathbb{P}) \rightarrow \Theta \text{ defines the **class of problems** .}$$

- ◆ Use a **parametric (iterative) optimization algorithm** with  $t = 0, 1, 2, \dots$ :

$$\mathcal{A} : \mathcal{H} \times \Theta \times \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad (\alpha, \theta, x^{(t)}) \mapsto \mathcal{A}(\alpha, \theta, x^{(t)}) =: x^{(t+1)}$$

- ◆ Stop the algorithm, if  $x^{(t)}$  satisfies a **stopping criterion**  $C_\varepsilon(\theta)$  and define **stopping time**

$$\tau_\alpha(\theta, \mathbf{x}) := \inf\{t \in \mathbb{N} \mid x_{\alpha, \theta}^{(t)} \in C_\varepsilon(\theta)\}, \quad \text{where } \mathbf{x} := (x_{\alpha, \theta}^{(t)})_{t \in \mathbb{N}}.$$

- ◆ For best **expected case complexity** w.r.t. choice of algorithm  $\mathcal{A}$ , solve

$$\min_{\alpha \in \mathcal{H}} \left( \mathbb{E}_{\theta, \mathbf{x} | \alpha} \{ \tau_\alpha(\theta, \mathbf{x}) \} \right).$$

Using the **expectation is intractable** in general!

↪ resolve using **PAC-Bayesian bounds**

# Probabilistic Formulation of the Optimization Procedure

- ◆ **Random parametric optimization problem** with measurable  $\ell : \mathbb{R}^d \times \Theta \rightarrow \mathbb{R}_{\geq 0}$ :

$$\min_{x \in \mathbb{R}^d} \ell(x, \theta), \quad \text{random variable } \theta : (\Omega, \mathfrak{A}, \mathbb{P}) \rightarrow \Theta \text{ defines the **class of problems** .}$$

- ◆ Use a **parametric (iterative) optimization algorithm** with  $t = 0, 1, 2, \dots$ :

$$\mathcal{A} : \mathcal{H} \times \Theta \times \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad (\alpha, \theta, x^{(t)}) \mapsto \mathcal{A}(\alpha, \theta, x^{(t)}) =: x^{(t+1)}$$

- ◆ Stop the algorithm, if  $x^{(t)}$  satisfies a **stopping criterion**  $C_\varepsilon(\theta)$  and define **stopping time**

$$\tau_\alpha(\theta, \mathbf{x}) := \inf\{t \in \mathbb{N} \mid x_{\alpha, \theta}^{(t)} \in C_\varepsilon(\theta)\}, \quad \text{where } \mathbf{x} := (x_{\alpha, \theta}^{(t)})_{t \in \mathbb{N}}.$$

- ◆ For best **expected case complexity** w.r.t. choice of algorithm  $\mathcal{A}$ , solve

$$\min_{\alpha \in \mathcal{H}} \left( \mathbb{E}_{\theta, \mathbf{x} | \alpha} \{ \tau_\alpha(\theta, \mathbf{x}) \} \right).$$

Using the **expectation is intractable** in general!

↪ resolve using PAC-Bayesian bounds

- **Stopping time needs access to the full trajectory**  $\mathbf{x} = (x_{\alpha, \theta}^{(t)})_{t \in \mathbb{N}}$
- **Expectation requires the conditional distribution of  $\mathbf{x}$  given  $(\alpha, \theta)$ !**

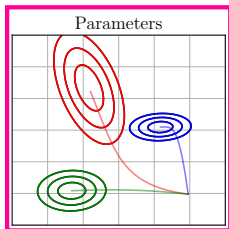
↪ model the distribution over algorithms/trajectories.

(Iterative) Optimization Algorithm:

$$x^{(t+1)} = \mathcal{A}(\alpha, \theta, x^{(t)})$$

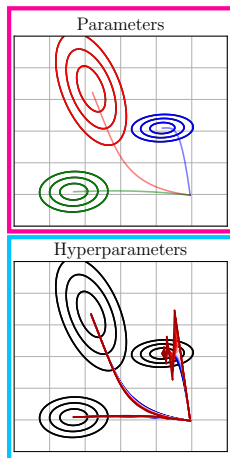
## (Iterative) Optimization Algorithm:

$$x^{(t+1)} = \mathcal{A}(\alpha, \theta, x^{(t)})$$



(Iterative) Optimization Algorithm:

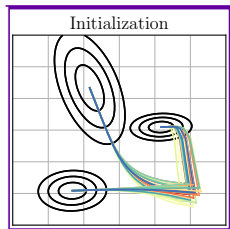
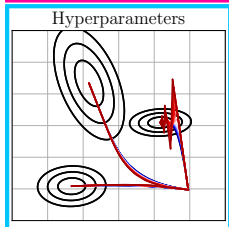
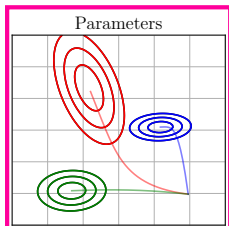
$$x^{(t+1)} = \mathcal{A}(\alpha, \theta, x^{(t)})$$



# From randomness of the algorithm to the distribution over trajectories

(Iterative) Optimization Algorithm:

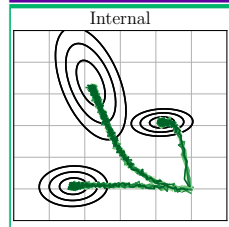
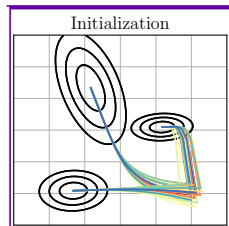
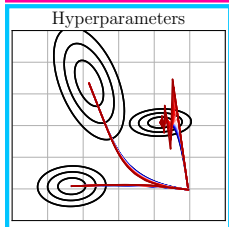
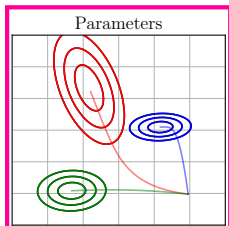
$$x^{(t+1)} = \mathcal{A}(\alpha, \theta, x^{(t)})$$



# From randomness of the algorithm to the distribution over trajectories

(Iterative) Optimization Algorithm:  $\eta^{(t+1)}$  used for stochastic algorithms

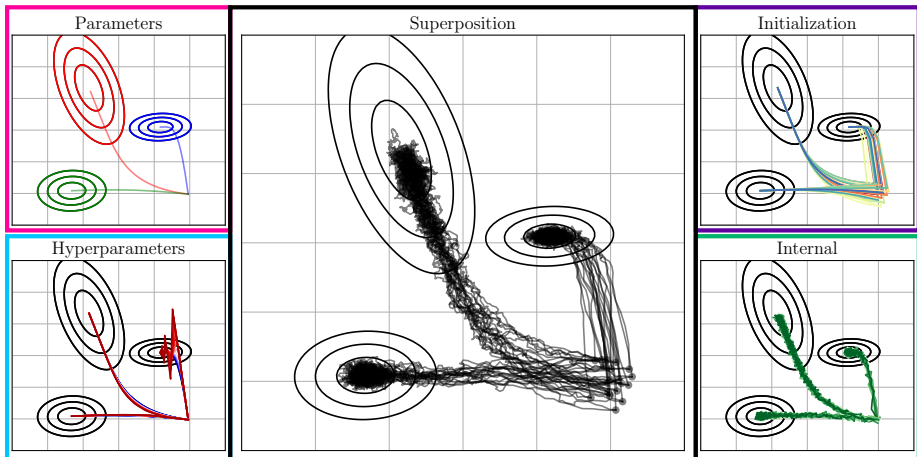
$$x^{(t+1)} = \mathcal{A}(\alpha, \theta, x^{(t)}, \eta^{(t+1)})$$



# From randomness of the algorithm to the distribution over trajectories

(Iterative) Optimization Algorithm:  $\mathcal{A}$  describes a discrete-time stochastic process

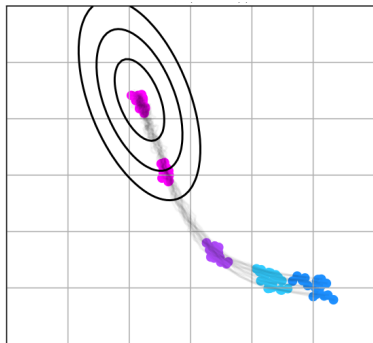
$$x^{(t+1)} = \mathcal{A}(\alpha, \theta, x^{(t)}, \eta^{(t+1)})$$



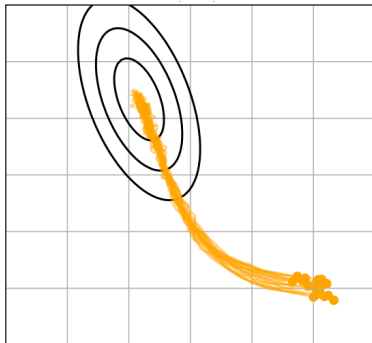
# From randomness of the algorithm to the distribution over trajectories

$$x^{(t+1)} = \mathcal{A}(\alpha, \theta, x^{(t)}, \eta^{(t+1)}) \xrightarrow{\text{Theorem}} \mathbb{P}_{\mathbf{x}|\alpha, \theta}\{\mathbf{B}\}, \mathbf{B} \subset (\mathbb{R}^d)^{\mathbb{N}}$$

exists regular conditional distribution over trajectories  $\mathbf{x}$  given  $\alpha, \theta$



distribution over iteration maps



distribution over trajectories

## PAC-Bayesian Bounds: Generalization of the Stopping Time

Use the measurable stopping time  $\tau_\alpha(\boldsymbol{\theta}, \boldsymbol{x}) := \inf\{t \leq t_{\max} \mid x_{\alpha, \boldsymbol{\theta}}^{(t)} \in C_\varepsilon(\boldsymbol{\theta})\}$   
and define

**average expected stopping time:**  $\bar{\tau}_\alpha := \mathbb{E}_{(\boldsymbol{\theta}, \boldsymbol{x}) | \alpha} \{\tau_\alpha(\boldsymbol{\theta}, \boldsymbol{x})\}$

**$n$ -th expected stopping time:**  $\bar{\tau}_{n, \alpha} := \mathbb{E}_{(\boldsymbol{\theta}, \boldsymbol{x}) | \alpha, \theta_n} \{\tau_\alpha(\boldsymbol{\theta}, \boldsymbol{x})\} .$

# PAC-Bayesian Bounds: Generalization of the Stopping Time

Use the measurable stopping time  $\tau_\alpha(\theta, \mathbf{x}) := \inf\{t \leq t_{\max} \mid x_{\alpha, \theta}^{(t)} \in C_\epsilon(\theta)\}$  and define

**average expected stopping time:**  $\bar{\tau}_\alpha := \mathbb{E}_{(\theta, \mathbf{x}) | \alpha} \{\tau_\alpha(\theta, \mathbf{x})\}$

**$n$ -th expected stopping time:**  $\bar{\tau}_{n, \alpha} := \mathbb{E}_{(\theta, \mathbf{x}) | \alpha, \theta_n} \{\tau_\alpha(\theta, \mathbf{x})\}$ .

**Theorem:** For every  $\lambda \in (0, \infty)$  and  $\epsilon > 0$  it holds that:

$$\mathbb{P}_{\theta_{[N]}} \left\{ \forall \rho \in \mathcal{M}(\mathbb{P}_\alpha) : \right. \\ \left. \rho[\bar{\tau}_\alpha] \leq \frac{1}{N} \sum_{n=1}^N \rho[\bar{\tau}_{n, \alpha}] + \frac{D_{\text{KL}}(\rho \parallel \mathbb{P}_\alpha) + \frac{\lambda^2}{2N} t_{\max}^2 - \log(\epsilon)}{\lambda} \right\} \geq 1 - \epsilon.$$

**Reminder:** (Abstract PAC-Bayesian bound from before)

$$\mathbb{P}_{\theta_{[N]}} \left\{ \forall \rho \in \mathcal{M}(\mathbb{P}_\alpha) : \rho_{\theta_{[N]}}^*[\mathcal{R}(\alpha)] \leq \rho[\hat{\mathcal{R}}(\alpha, \theta_{[N]})] + K(\rho, N, \epsilon) \right\} \geq 1 - \epsilon,$$

# PAC-Bayesian Bounds: Generalization of the Stopping Time

Use the measurable stopping time  $\tau_\alpha(\theta, \mathbf{x}) := \inf\{t \leq t_{\max} \mid x_{\alpha, \theta}^{(t)} \in C_\epsilon(\theta)\}$  and define

**average expected stopping time:**  $\bar{\tau}_\alpha := \mathbb{E}_{(\theta, \mathbf{x}) | \alpha} \{\tau_\alpha(\theta, \mathbf{x})\}$

**$n$ -th expected stopping time:**  $\bar{\tau}_{n, \alpha} := \mathbb{E}_{(\theta, \mathbf{x}) | \alpha, \theta_n} \{\tau_\alpha(\theta, \mathbf{x})\}$ .

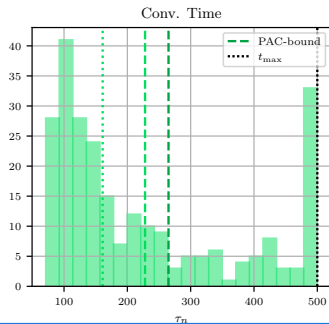
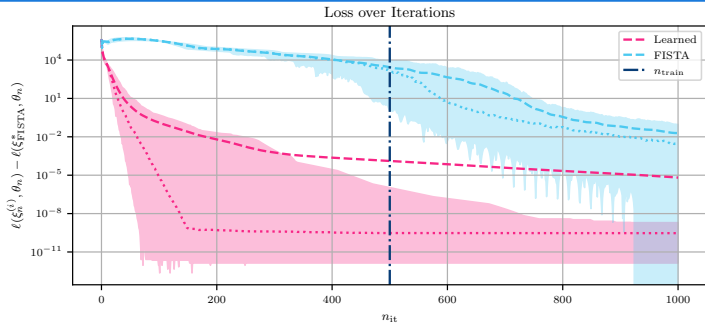
**Theorem:** For every  $\lambda \in (0, \infty)$  and  $\epsilon > 0$  it holds that:

$$\mathbb{P}_{\theta_{[N]}} \left\{ \forall \rho \in \mathcal{M}(\mathbb{P}_\alpha) : \right. \\ \left. \rho[\bar{\tau}_\alpha] \leq \frac{1}{N} \sum_{n=1}^N \rho[\bar{\tau}_{n, \alpha}] + \frac{D_{\text{KL}}(\rho \parallel \mathbb{P}_\alpha) + \frac{\lambda^2}{2N} t_{\max}^2 - \log(\epsilon)}{\lambda} \right\} \geq 1 - \epsilon.$$

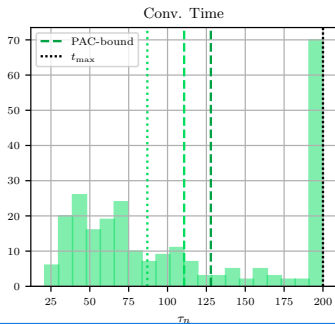
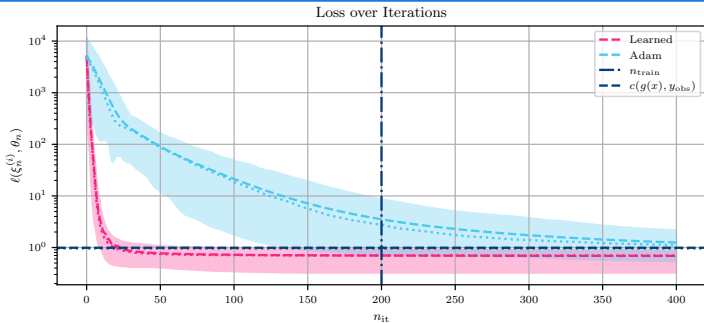
**Training / Learning the algorithm:** = **optimize the right hand side**

- ◆ Closed form solution of right hand side. Only optimize w.r.t.  $\lambda$ .
- ◆ Note that stopping times can be translated to **convergence rates**. (*easier to train with*)
- ◆ **Criticism:** Prior needs to be constructed. (*usually trained on independent data*)

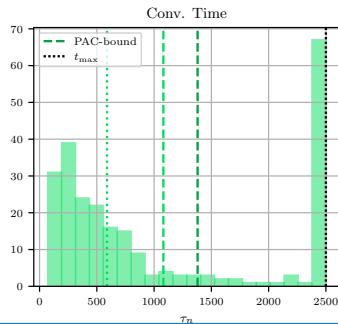
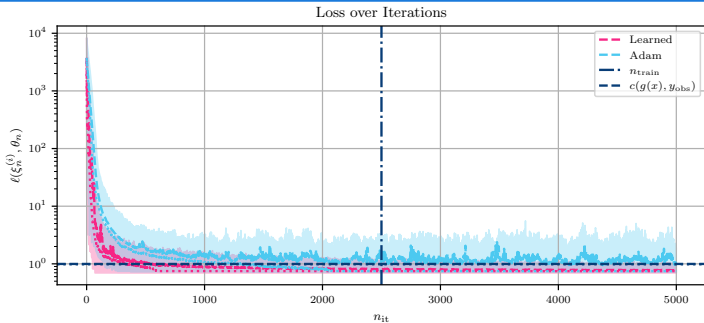
# Numerical Result: Lasso Problem



# Numerical Results: Training a neural network for regression



# Numerical Results: same problem but with mini-batch update



How to prove convergence of learned algorithms to a critical point ?

## How to prove convergence of learned algorithms to a critical point ?

### Classical strategy:

- ◆ Show that the algorithm (or the trajectory  $x$ ) has property  $a$  and  $b$ .

## How to prove convergence of learned algorithms to a critical point ?

### Classical strategy:

- ◆ Show that the algorithm (or the trajectory  $x$ ) has property  $a$  and  $b$ .
- ◆ Prove that  $a \wedge b$  imply convergence  $c$ :

$$x \text{ satisfies } a \wedge b \quad \implies \quad x \text{ satisfies } c$$

## How to prove convergence of learned algorithms to a critical point ?

### Classical strategy:

- ◆ Show that the algorithm (or the trajectory  $x$ ) has property  $a$  and  $b$ .
- ◆ Prove that  $a \wedge b$  imply convergence  $c$ :

$$x \text{ satisfies } a \wedge b \quad \implies \quad x \text{ satisfies } c$$

- ◆ **However**, we can only say that  $x$  satisfies a property with some probability.

## How to prove convergence of learned algorithms to a critical point ?

### Classical strategy:

- ◆ Show that the algorithm (or the trajectory  $x$ ) has property  $a$  and  $b$ .
- ◆ Prove that  $a \wedge b$  imply convergence  $c$ :

$$x \text{ satisfies } a \wedge b \quad \Longrightarrow \quad x \text{ satisfies } c$$

- ◆ **However**, we can only say that  $x$  satisfies a property with some probability.

### Probabilistic / measure theoretic strategy:

- ◆ Define  $A := \{x \mid x \text{ satisfies } a\}$ ,  $B := \{x \mid x \text{ satisfies } b\}$ ,  $C := \{x \mid x \text{ satisfies } c\}$ .
- ◆ Then the question of “implication” becomes an “inclusion”

$$A \cap B = \{x \mid x \text{ satisfies } a \wedge b\} \subset C$$

## How to prove convergence of learned algorithms to a critical point ?

### Classical strategy:

- ◆ Show that the algorithm (or the trajectory  $x$ ) has property  $a$  and  $b$ .
- ◆ Prove that  $a \wedge b$  imply convergence  $c$ :

$$x \text{ satisfies } a \wedge b \quad \Longrightarrow \quad x \text{ satisfies } c$$

- ◆ **However**, we can only say that  $x$  satisfies a property with some probability.

### Probabilistic / measure theoretic strategy:

- ◆ Define  $A := \{x \mid x \text{ satisfies } a\}$ ,  $B := \{x \mid x \text{ satisfies } b\}$ ,  $C := \{x \mid x \text{ satisfies } c\}$ .
- ◆ Then the question of “implication” becomes an “inclusion”

$$A \cap B = \{x \mid x \text{ satisfies } a \wedge b\} \subset C$$

- ◆ Given a probability measure  $\mu$  over trajectories  $x$ , we can conclude

$$\mu\{A \cap B\} \leq \mu\{C\}.$$

**Theorem:** [Attouch et al. 2013] Let  $f: \mathbb{R}^d \rightarrow \overline{\mathbb{R}}$  be a proper lsc Kurdyka-Łojasiewicz function that is continuous on  $\text{dom } f$  and  $\inf f > -\infty$ . Suppose there exists  $a, b > 0$  such that  $(x^{(t)})_{t \in \mathbb{N}}$  satisfies

◆ **Sufficient decrease condition:**

$$f(x^{(t+1)}) + a\|x^{(t+1)} - x^{(t)}\|^2 \leq f(x^{(t)}), \quad \forall t \in \mathbb{N};$$

◆ **Relative error condition:**

$$\exists v^{(t+1)} \in \partial f(x^{(t+1)}): \|v^{(t+1)}\| \leq b\|x^{(t+1)} - x^{(t)}\|, \quad \forall t \in \mathbb{N};$$

◆ **Boundedness condition:**

$$(x^{(t)})_{t \in \mathbb{N}} \text{ is bounded.}$$

**Then,  $(x^{(t)})_{t \in \mathbb{N}}$  converges to a critical point of  $f$ .**

- ◆ Collection of sequences that **converge to a stationary point**:

$$A_{\text{conv}} := \{(\boldsymbol{\theta}, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \rightarrow \tilde{x} \text{ as } t \rightarrow \infty \text{ and } 0 \in \partial_x \ell(\tilde{x}, \boldsymbol{\theta})\}$$

- ◆ Collection of sequences that **converge to a stationary point**:

$$A_{\text{conv}} := \{(\boldsymbol{\theta}, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \rightarrow \tilde{x} \text{ as } t \rightarrow \infty \text{ and } 0 \in \partial_x \ell(\tilde{x}, \boldsymbol{\theta})\}$$

- ◆ Collection of sequences that satisfy the **sufficient decrease condition**:

$$A_{\text{desc}} := \{(\boldsymbol{\theta}, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \in \text{dom } \ell \text{ for all } t \in \mathbb{N} \text{ and}$$

$$\exists a > 0: \forall t \in \mathbb{N}: \ell(x^{(t+1)}, \boldsymbol{\theta}) + a \|x^{(t+1)} - x^{(t)}\|^2 \leq \ell(x^{(t)}, \boldsymbol{\theta})\}$$

- ◆ Collection of sequences that **converge to a stationary point**:

$$A_{\text{conv}} := \{(\boldsymbol{\theta}, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \rightarrow \tilde{x} \text{ as } t \rightarrow \infty \text{ and } 0 \in \partial_x \ell(\tilde{x}, \boldsymbol{\theta})\}$$

- ◆ Collection of sequences that satisfy the **sufficient decrease condition**:

$$A_{\text{desc}} := \{(\boldsymbol{\theta}, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \in \text{dom } \ell \text{ for all } t \in \mathbb{N} \text{ and} \\ \exists a > 0: \forall t \in \mathbb{N}: \ell(x^{(t+1)}, \boldsymbol{\theta}) + a \|x^{(t+1)} - x^{(t)}\|^2 \leq \ell(x^{(t)}, \boldsymbol{\theta})\}$$

- ◆ Collection of sequences that satisfy the **relative error condition**:

$$A_{\text{rerr}} := \{(\boldsymbol{\theta}, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \in \text{dom } \ell \text{ for all } t \in \mathbb{N} \text{ and} \\ \exists b > 0: \forall t \in \mathbb{N}: \|v(x^{(t+1)}, \boldsymbol{\theta})\| \leq b \|x^{(t+1)} - x^{(t)}\|\}$$

- ◆ Collection of sequences that **converge to a stationary point**:

$$A_{\text{conv}} := \{(\boldsymbol{\theta}, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \rightarrow \tilde{x} \text{ as } t \rightarrow \infty \text{ and } 0 \in \partial_x \ell(\tilde{x}, \boldsymbol{\theta})\}$$

- ◆ Collection of sequences that satisfy the **sufficient decrease condition**:

$$A_{\text{desc}} := \{(\boldsymbol{\theta}, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \in \text{dom } \ell \text{ for all } t \in \mathbb{N} \text{ and} \\ \exists a > 0: \forall t \in \mathbb{N}: \ell(x^{(t+1)}, \boldsymbol{\theta}) + a \|x^{(t+1)} - x^{(t)}\|^2 \leq \ell(x^{(t)}, \boldsymbol{\theta})\}$$

- ◆ Collection of sequences that satisfy the **relative error condition**:

$$A_{\text{rerr}} := \{(\boldsymbol{\theta}, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \in \text{dom } \ell \text{ for all } t \in \mathbb{N} \text{ and} \\ \exists b > 0: \forall t \in \mathbb{N}: \|v(x^{(t+1)}, \boldsymbol{\theta})\| \leq b \|x^{(t+1)} - x^{(t)}\|\}$$

- ◆ Collection of sequences that are **bounded**:

$$A_{\text{bd}} := \{(\boldsymbol{\theta}, (x^{(t)})_{t \in \mathbb{N}}) \mid (x^{(t)})_{t \in \mathbb{N}} \text{ is bounded}\}$$

- Collection of sequences that **converge to a stationary point**:

$$A_{\text{conv}} := \{(\theta, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \rightarrow \tilde{x} \text{ as } t \rightarrow \infty \text{ and } 0 \in \partial_x \ell(\tilde{x}, \theta)\}$$

- Collection of sequences that satisfy the **sufficient decrease condition**:

$$A_{\text{desc}} := \{(\theta, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \in \text{dom } \ell \text{ for all } t \in \mathbb{N} \text{ and} \\ \exists a > 0: \forall t \in \mathbb{N}: \ell(x^{(t+1)}, \theta) + a \|x^{(t+1)} - x^{(t)}\|^2 \leq \ell(x^{(t)}, \theta)\}$$

- Collection of sequences that satisfy the **relative error condition**:

$((x, \theta) \mapsto v(x, \theta))$  is a measurable selection of  $\partial_x \ell(x, \theta)$ :

$$A_{\text{rerr}} := \{(\theta, (x^{(t)})_{t \in \mathbb{N}}) \mid x^{(t)} \in \text{dom } \ell \text{ for all } t \in \mathbb{N} \text{ and} \\ \exists b > 0: \forall t \in \mathbb{N}: \|v(x^{(t+1)}, \theta)\| \leq b \|x^{(t+1)} - x^{(t)}\|\}$$

- Collection of sequences that are **bounded**:

$$A_{\text{bd}} := \{(\theta, (x^{(t)})_{t \in \mathbb{N}}) \mid (x^{(t)})_{t \in \mathbb{N}} \text{ is bounded}\}$$

**Proposition:** All sets listed above are measurable and if  $\ell(\cdot, \theta)$  is a Kurdyka-Łojasiewicz function for every  $\theta$ , then we have that

$$A_{\text{desc}} \cap A_{\text{rerr}} \cap A_{\text{bd}} \subset A_{\text{conv}} .$$

- ◆ Recall:  $A := A_{\text{desc}} \cap A_{\text{rerr}} \cap A_{\text{bd}} \subset A_{\text{conv}}$
- ◆ Using monotonicity

$$\mathbb{P}_{(\theta, x)|\alpha} \{A_{\text{desc}} \cap A_{\text{rerr}} \cap A_{\text{bd}}\} \leq \mathbb{P}_{(\theta, x)|\alpha} \{A_{\text{conv}}\},$$

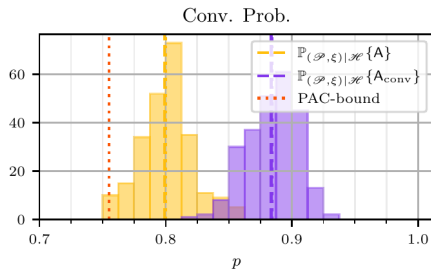
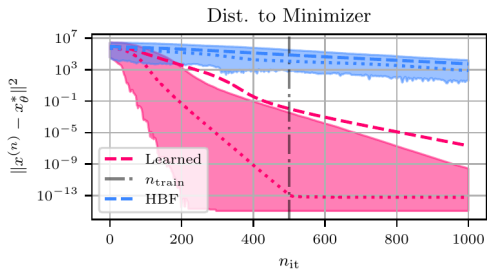
we obtain **generalization of the probability to converge**:

**Theorem: [Catoni-like PAC bound]** Let  $\ell(\cdot, \theta)$  is a Kurdyka-Łojasiewicz function for every  $\theta$ . Then, for any  $\lambda \in (0, \infty)$ , it holds that:

$$\mathbb{P}_{\theta_{[N]}} \left\{ \forall \rho \in \mathcal{M}(\mathbb{P}_\alpha) : \rho \left[ \mathbb{P}_{(\theta, x)|\alpha} \{A_{\text{conv}}\} \right] \geq 1 - \Phi_{\frac{\lambda}{N}}^{-1} \left( \frac{1}{N} \sum_{n=1}^N \rho \left[ \mathbb{P}_{(\theta_n, x_n)|\alpha, \theta_n} \{A^c\} \right] + \frac{D_{\text{KL}}(\rho \parallel \mathbb{P}_\alpha) - \log(\epsilon)}{\lambda} \right) \right\} \geq 1 - \epsilon.$$

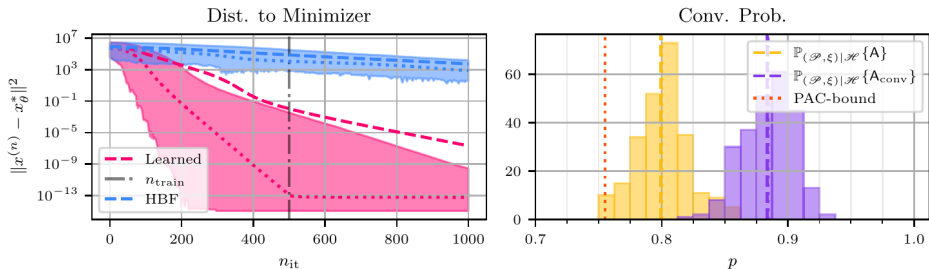
# Convergence of Learned Algorithms

## Smooth and strongly convex problem:

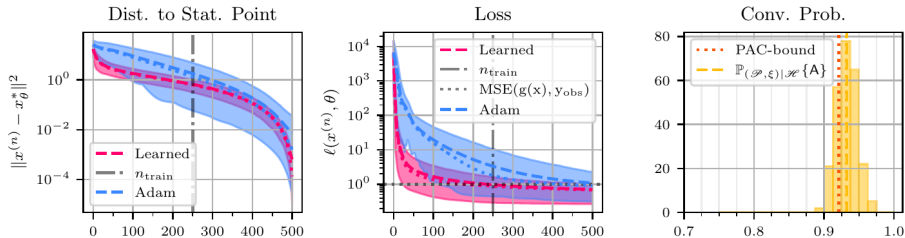


# Convergence of Learned Algorithms

## Smooth and strongly convex problem:



## Non-smooth and non-convex problem:



# Goal: Learning of optimization algorithms with theoretical guarantees

PAC-Bayes

Learning-to-Optimize

PAC-Bayesian Learning of Optimization Algorithms

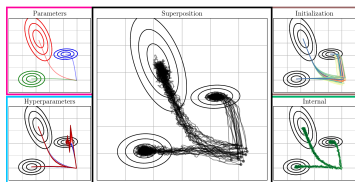
The diagram is enclosed in a blue rounded rectangle. It features two terms at the top: 'PAC-Bayes' on the left and 'Learning-to-Optimize' on the right. A vertical dashed line connects these two terms. From the ends of this dashed line, two solid curved lines extend downwards and outwards, meeting at a point above a downward-pointing arrow. This arrow points to the text 'PAC-Bayesian Learning of Optimization Algorithms' at the bottom of the diagram.

PAC-Bayes

Learning-to-Optimize

PAC-Bayesian Learning of Optimization Algorithms

- ◆ probabilistic framework for learning algorithms

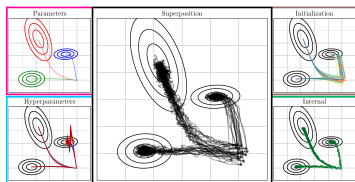


PAC-Bayes

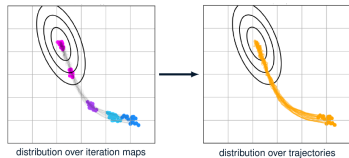
Learning-to-Optimize

PAC-Bayesian Learning of Optimization Algorithms

- ◆ probabilistic framework for learning algorithms



- ◆ analysis requires distributions over trajectories



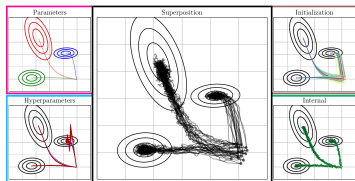
# Goal: Learning of optimization algorithms with theoretical guarantees

PAC-Bayes

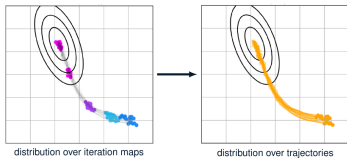
Learning-to-Optimize

PAC-Bayesian Learning of Optimization Algorithms

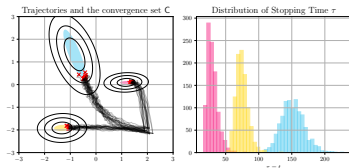
- probabilistic framework for learning algorithms



- analysis requires distributions over trajectories



- learning of **Stopping Times**

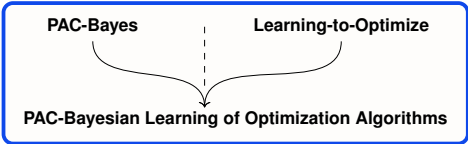


with PAC-Bayesian bounds

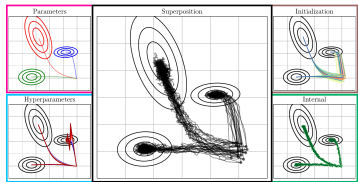
**Corollary:** For every  $\lambda \in (0, \infty)$  and  $\epsilon > 0$  it holds that:

$$\mathbb{P}_{\theta_0, N} \left\{ \forall \rho \in \mathcal{P}(\mathbb{P}_\alpha) : \right.$$

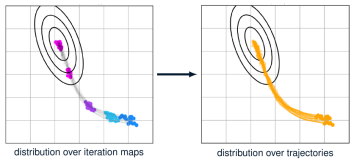
$$\left. \rho[\tau] \leq \frac{1}{N} \sum_{n=1}^N \rho[\tau_n | \alpha, \theta_n] + \frac{D_{\text{KL}}(\rho \| \mathbb{P}_\alpha) + \frac{\lambda^2}{3N} \rho_{\max} - \log(\epsilon)}{\lambda} \right\} \geq 1 - \epsilon.$$



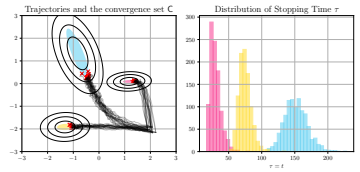
probabilistic framework for learning algorithms



analysis requires distributions over trajectories



learning of **Stopping Times**



with PAC-Bayesian bounds

**Corollary:** For every  $\lambda \in (0, \infty)$  and  $\epsilon > 0$  it holds that:

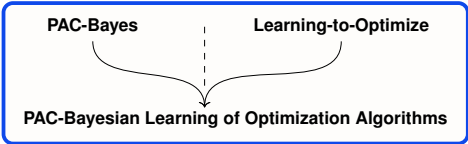
$$\mathbb{P}_{\theta, N} \left\{ \forall \rho \in \mathcal{P}(\mathbb{P}_\alpha) : \rho[\tau] \leq \frac{1}{N} \sum_{n=1}^N \rho[\mathbb{E}\{\tau_n \mid \alpha, \theta, \rho_n\}] + \frac{D_{\text{KL}}(\rho \parallel \mathbb{P}_\alpha)}{\lambda} + \frac{\lambda^2 \epsilon^2}{\lambda} \max - \log(\epsilon) \right\} \geq 1 - \epsilon.$$

learning of **algorithms that converge to critical points** by combining [Attouch et al. 2013] with PAC-Bayesian bounds

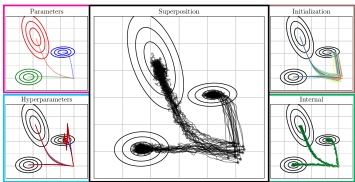
**Theorem:** Let  $\ell(\cdot, \theta)$  is a Kurdyka-Lojasiewicz function for every  $\theta$ . Then, for any  $\lambda \in (0, \infty)$ , it holds that:

$$\mathbb{P}_{\theta, N} \left\{ \forall \rho \in \mathcal{P}(\mathbb{P}_\alpha) : \rho[\mathbb{P}(\theta, \epsilon)]_{\text{in}}(\mathcal{A}_{\text{conv}}) \geq 1 - \Phi_\lambda^{-1} \left( \frac{1}{N} \sum_{n=1}^N \rho[\mathbb{P}(\theta_n, \epsilon_n)]_{\text{in}}(\mathcal{A}^c) + \frac{D_{\text{KL}}(\rho \parallel \mathbb{P}_\alpha)}{\lambda} - \log(\epsilon) \right) \right\} \geq 1 - \epsilon.$$

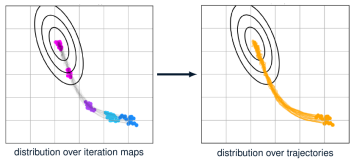
# Goal: Learning of optimization algorithms with theoretical guarantees



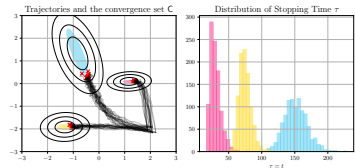
probabilistic framework for learning algorithms



analysis requires distributions over trajectories



learning of **Stopping Times**



with PAC-Bayesian bounds

**Corollary:** For every  $\lambda \in (0, \infty)$  and  $\epsilon > 0$  it holds that:

$$\mathbb{P}_{\theta, N} \left\{ \forall \rho \in \mathcal{P}(\mathbb{P}_\alpha) : \rho[\tau] \leq \frac{1}{N} \sum_{n=1}^N \rho[\mathbb{E}\{\tau_n \mid \alpha, \theta, \alpha_n\}] + \frac{D_{\text{KL}}(\rho \parallel \mathbb{P}_\alpha) + \frac{\lambda^2 \epsilon^2}{\lambda} \max - \log(\epsilon)}{\lambda} \right\} \geq 1 - \epsilon.$$

learning of **algorithms that converge to critical points** by combining [Atouch et al. 2013] with PAC-Bayesian bounds

**Theorem:** Let  $\ell(\cdot, \theta)$  is a Kurdyka-Lojasiewicz function for every  $\theta$ . Then, for any  $\lambda \in (0, \infty)$ , it holds that:

$$\mathbb{P}_{\theta, N} \left\{ \forall \rho \in \mathcal{P}(\mathbb{P}_\alpha) : \rho[\mathbb{P}(\theta, \epsilon)_{\text{in}}\{A^c\}] \geq 1 - \Phi_{\frac{1}{\lambda}}^{-1} \left( \frac{1}{N} \sum_{n=1}^N \rho[\mathbb{P}(\theta, \epsilon_n)_{\text{in}}\{A^c\}] + \frac{D_{\text{KL}}(\rho \parallel \mathbb{P}_\alpha) - \log(\epsilon)}{\lambda} \right) \right\} \geq 1 - \epsilon.$$

**What you observe is what you get !**

