

# Stochastic Approximation in Hilbert Spaces

Aymeric DIEULEVEUT

Supervised by Francis BACH

September 28, 2017



# Outline

1. Introduction:
  - ▶ Supervised Machine Learning
  - ▶ Stochastic Approximation
2. Finite dimensional results
3. Infinite dimensional results
4. Beyond quadratic loss: interpretation as a Markov chain.

# Supervised Machine Learning: definition & applications

**Goal:** predict a phenomenon from “explanatory variables”, given a set of observations.

# Supervised Machine Learning: definition & applications

**Goal:** predict a phenomenon from “explanatory variables”, given a set of observations.



Bio-informatics

Input: DNA/RNA sequence,  
Output: Disease predisposition /  
Drug responsiveness

```
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
```

Image classification

Input: Handwritten digits / Images,  
Output: Digit

# Supervised Machine Learning: definition & applications

**Goal:** predict a phenomenon from “explanatory variables”, given a set of observations.



Bio-informatics

Input: DNA/RNA sequence,  
Output: Disease predisposition /  
Drug responsiveness

```
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
```

Image classification

Input: Handwritten digits / Images,  
Output: Digit

“Large scale” learning framework: both the number of examples  $n$  and the number of explanatory variables  $d$  are large.

# Supervised Machine Learning: definition & applications

**Goal:** predict a phenomenon from “explanatory variables”, given a set of observations.



Bio-informatics

Input: DNA/RNA sequence,  
Output: Disease predisposition /  
Drug responsiveness

$n \rightarrow 10$  to  $10^4$

$d$  (e.g., number of basis)  $\rightarrow 10^6$

```
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
```

Image classification

Input: Handwritten digits / Images,  
Output: Digit

$n \rightarrow$  up to  $10^9$

$d$  (e.g., number of pixels)  $\rightarrow 10^6$

“Large scale” learning framework: both the number of examples  $n$  and the number of explanatory variables  $d$  are large.

## Supervised Machine Learning: mathematical framework

Consider an input/output pair  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ .  $(X, Y) \sim \rho$ , unknown distribution.

$\mathcal{Y} = \mathbb{R}$  (regression) or  $\{-1, 1\}$  (classification).

## Supervised Machine Learning: mathematical framework

Consider an input/output pair  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ .  $(X, Y) \sim \rho$ , unknown distribution.

$\mathcal{Y} = \mathbb{R}$  (regression) or  $\{-1, 1\}$  (classification).

**Goal:** find  $g : \mathcal{X} \rightarrow \mathbb{R}$ , such that  $g(X)$  is a *good* prediction for  $Y$ .



## Supervised Machine Learning: mathematical framework

Consider an input/output pair  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ .  $(X, Y) \sim \rho$ , unknown distribution.

$\mathcal{Y} = \mathbb{R}$  (regression) or  $\{-1, 1\}$  (classification).

**Goal:** find  $g : \mathcal{X} \rightarrow \mathbb{R}$ , such that  $g(X)$  is a *good* prediction for  $Y$ .

Measure accuracy with a loss function  $\ell : \mathcal{Y} \times \mathbb{R} \rightarrow \mathbb{R}_+$ : squared loss, logistic loss...

## Supervised Machine Learning: mathematical framework

Consider an input/output pair  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ .  $(X, Y) \sim \rho$ , unknown distribution.

$\mathcal{Y} = \mathbb{R}$  (regression) or  $\{-1, 1\}$  (classification).

**Goal:** find  $g : \mathcal{X} \rightarrow \mathbb{R}$ , such that  $g(X)$  is a *good* prediction for  $Y$ .

Measure accuracy with a loss function  $\ell : \mathcal{Y} \times \mathbb{R} \rightarrow \mathbb{R}_+$ : squared loss, logistic loss...

Risk (generalization error):

$$\mathcal{R}(g) := \mathbb{E}_\rho [\ell(Y, g(X))].$$

Parametric case: Prediction as a **linear function**  $g_\theta(X) = \langle \theta, \Phi(X) \rangle$  of features  $\Phi(X) \in \mathbb{R}^d$ . Notation:  $\mathcal{R}(\theta) := \mathcal{R}(g_\theta)$ .

## Supervised Machine Learning: mathematical framework

Consider an input/output pair  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ .  $(X, Y) \sim \rho$ , unknown distribution.

$\mathcal{Y} = \mathbb{R}$  (regression) or  $\{-1, 1\}$  (classification).

**Goal:** find  $g : \mathcal{X} \rightarrow \mathbb{R}$ , such that  $g(X)$  is a *good* prediction for  $Y$ .

Measure accuracy with a loss function  $\ell : \mathcal{Y} \times \mathbb{R} \rightarrow \mathbb{R}_+$ : squared loss, logistic loss...

Risk (generalization error):

$$\mathcal{R}(g) := \mathbb{E}_\rho [\ell(Y, g(X))].$$

**Parametric case:** Prediction as a **linear function**  $g_\theta(X) = \langle \theta, \Phi(X) \rangle$  of features  $\Phi(X) \in \mathbb{R}^d$ . Notation:  $\mathcal{R}(\theta) := \mathcal{R}(g_\theta)$ .

**Non-parametric case:** Prediction as a function  $g \in \mathcal{H}$ , for  $\mathcal{H}$  infinite-dimensional space.

# Empirical Risk minimization (I) - Parametric case

- ▶ **Data:**  $n$  observations  $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ ,  $i = 1, \dots, n$ , **i.i.d.**
- ▶ Empirical risk (or training error):

$$\hat{\mathcal{R}}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, \langle \theta, \Phi(x_i) \rangle).$$

- ▶ First approach: **Empirical risk minimization (regularized):**

$$\hat{\theta} := \operatorname{argmin}_{\theta \in \mathbb{R}^d} \hat{\mathcal{R}}(\theta) + \mu \Omega(\theta).$$

data fitting term + regularizer

## Empirical Risk minimization (II) - Parametric case

- ▶ For example, **least-squares regression**:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \theta, \Phi(x_i) \rangle)^2 + \mu \Omega(\theta),$$

## Empirical Risk minimization (II) - Parametric case

- ▶ For example, **least-squares regression**:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \theta, \Phi(x_i) \rangle)^2 + \mu \Omega(\theta),$$

- ▶ and **logistic regression**:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \log (1 + \exp(-y_i \langle \theta, \Phi(x_i) \rangle)) + \mu \Omega(\theta).$$

## Empirical Risk minimization (II) - Parametric case

- ▶ For example, **least-squares regression**:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \theta, \Phi(x_i) \rangle)^2 + \mu \Omega(\theta),$$

- ▶ and **logistic regression**:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \log (1 + \exp(-y_i \langle \theta, \Phi(x_i) \rangle)) + \mu \Omega(\theta).$$

- ▶ **Two fundamental questions:** (1) computing  $\hat{\theta}$  and (2) analyzing  $\hat{\theta}$ .

## Empirical Risk minimization (II) - Parametric case

- ▶ For example, **least-squares regression**:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{2n} \sum_{i=1}^n (y_i - \langle \theta, \Phi(x_i) \rangle)^2 + \mu \Omega(\theta),$$

- ▶ and **logistic regression**:

$$\min_{\theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \log (1 + \exp(-y_i \langle \theta, \Phi(x_i) \rangle)) + \mu \Omega(\theta).$$

- ▶ **Two fundamental questions**: (1) computing  $\hat{\theta}$  and (2) analyzing  $\hat{\theta}$ .

### 2 important insights for ML [Bottou and Bousquet, 2008]:

1. No need to optimize below *statistical error*,
2. True risk is more important than empirical risk.



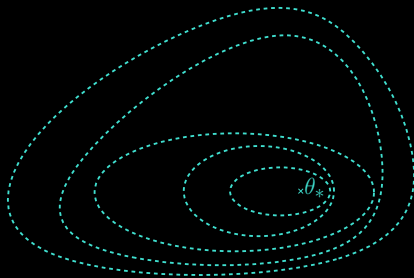
# Stochastic Approximation

► **Goal:**

$$\min_{\theta \in \mathbb{R}^d} f(\theta)$$

given unbiased gradient  
estimates  $f'_n$

►  $\theta_* := \operatorname{argmin}_{\mathbb{R}^d} f(\theta).$



# Stochastic Approximation

▶ **Goal:**

$$\min_{\theta \in \mathbb{R}^d} f(\theta)$$

given unbiased gradient  
estimates  $f'_n$

▶  $\theta_* := \operatorname{argmin}_{\mathbb{R}^d} f(\theta)$ .

▶ **Key algorithm: Stochastic Gradient Descent (SGD)** [Robbins and Monro, 1951]:

$$\theta_n = \theta_{n-1} - \gamma_n f'_n(\theta_{n-1})$$

▶  $\mathbb{E}[f'_n(\theta_{n-1}) | \mathcal{F}_{n-1}] = f'(\theta_{n-1})$  for a filtration  $(\mathcal{F}_n)_{n \geq 0}$ ,  $\theta_n$  is  $\mathcal{F}_n$  measurable.

# Stochastic Approximation

► **Goal:**

$$\min_{\theta \in \mathbb{R}^d} f(\theta)$$

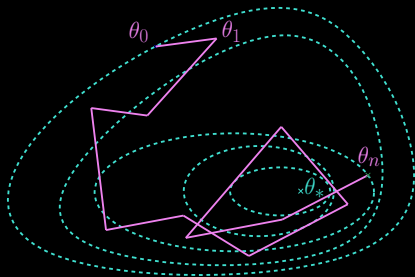
given unbiased gradient estimates  $f'_n$

►  $\theta_* := \operatorname{argmin}_{\mathbb{R}^d} f(\theta)$ .

► **Key algorithm: Stochastic Gradient Descent (SGD)** [Robbins and Monro, 1951]:

$$\theta_n = \theta_{n-1} - \gamma_n f'_n(\theta_{n-1})$$

►  $\mathbb{E}[f'_n(\theta_{n-1}) | \mathcal{F}_{n-1}] = f'(\theta_{n-1})$  for a filtration  $(\mathcal{F}_n)_{n \geq 0}$ ,  $\theta_n$  is  $\mathcal{F}_n$  measurable.



# Polyak Ruppert averaging

Introduced by Polyak and Juditsky [1992] and Ruppert [1988]:

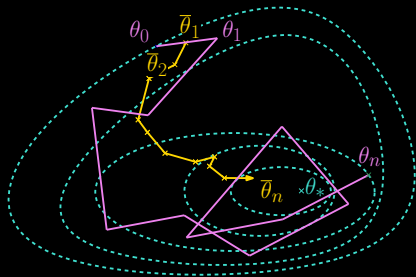
$$\bar{\theta}_n = \frac{1}{n+1} \sum_{k=0}^n \theta_k.$$

- ▶ off line averaging reduces the noise effect.

# Polyak Ruppert averaging

Introduced by Polyak and Juditsky [1992] and Ruppert [1988]:

$$\bar{\theta}_n = \frac{1}{n+1} \sum_{k=0}^n \theta_k.$$



- ▶ off line averaging reduces the noise effect.

# Stochastic Approximation (SA) in Machine Learning

Loss for a single pair of observations, for any  $k \leq n$ :

$$f_k(\theta) = \ell(y_k, \langle \theta, \Phi(x_k) \rangle).$$

SA for the **true risk** :

- ▶ For  $0 \leq k \leq n$ ,  $\mathcal{F}_k = \sigma((x_i, y_i)_{1 \leq i \leq k})$ .
- ▶ At step  $0 < k \leq n$ , use a **new point** independent of  $\theta_{k-1}$ :

$$\begin{aligned}\mathcal{R}(\theta) &= \mathbb{E} \ell(y_k, \langle \theta, \Phi(x_k) \rangle) \\ f'_k(\theta_{k-1}) &= \ell'(y_k, \langle \theta_{k-1}, \Phi(x_k) \rangle)\end{aligned}$$

# Stochastic Approximation (SA) in Machine Learning

Loss for a single pair of observations, for any  $k \leq n$ :

$$f_k(\theta) = \ell(y_k, \langle \theta, \Phi(x_k) \rangle).$$

SA for the **true risk** :

- ▶ For  $0 \leq k \leq n$ ,  $\mathcal{F}_k = \sigma((x_i, y_i)_{1 \leq i \leq k})$ .
- ▶ At step  $0 < k \leq n$ , use a **new point** independent of  $\theta_{k-1}$ :

$$\begin{aligned}\mathcal{R}(\theta) &= \mathbb{E} \ell(y_k, \langle \theta, \Phi(x_k) \rangle) \\ f'_k(\theta_{k-1}) &= \ell'(y_k, \langle \theta_{k-1}, \Phi(x_k) \rangle) \\ \mathbb{E}[f'_k(\theta_{k-1}) | \mathcal{F}_{k-1}] &= \mathcal{R}'(\theta_{k-1})\end{aligned}$$

# Stochastic Approximation (SA) in Machine Learning

Loss for a single pair of observations, for any  $k \leq n$ :

$$f_k(\theta) = \ell(y_k, \langle \theta, \Phi(x_k) \rangle).$$

SA for the **true risk** :

- ▶ For  $0 \leq k \leq n$ ,  $\mathcal{F}_k = \sigma((x_i, y_i)_{1 \leq i \leq k})$ .
- ▶ At step  $0 < k \leq n$ , use a **new point** independent of  $\theta_{k-1}$ :

$$\begin{aligned}\mathcal{R}(\theta) &= \mathbb{E} \ell(y_k, \langle \theta, \Phi(x_k) \rangle) \\ f'_k(\theta_{k-1}) &= \ell'(y_k, \langle \theta_{k-1}, \Phi(x_k) \rangle) \\ \mathbb{E}[f'_k(\theta_{k-1}) | \mathcal{F}_{k-1}] &= \mathcal{R}'(\theta_{k-1})\end{aligned}$$

Single pass through the data – “Automatic” regularization.

**Central algorithm in the thesis.**



## Outline: bibliography

- a) *Non-parametric Stochastic Approximation with Large Step-sizes*,  
A. Dieuleveut and F. Bach, in the *Annals of Statistics*
- b) *Harder, Better, Faster, Stronger Convergence Rates for Least-squares Regression*,  
A. Dieuleveut, N. Flammarion and F. Bach, in *Journal of Machine Learning Research*
- c) *Bridging the Gap between Constant Step Size Stochastic Gradient Descent and Markov Chains*,  
A. Dieuleveut, A. Durmus, F. Bach, under submission.

## Outline: bibliography

- a) *Non-parametric Stochastic Approximation with Large Step-sizes*,  
A. Dieuleveut and F. Bach, in the *Annals of Statistics*
- b) *Harder, Better, Faster, Stronger Convergence Rates for Least-squares Regression*,  
A. Dieuleveut, N. Flammarion and F. Bach, in *Journal of Machine Learning Research*
- c) *Bridging the Gap between Constant Step Size Stochastic Gradient Descent and Markov Chains*,  
A. Dieuleveut, A. Durmus, F. Bach, under submission.

	Quadratic loss	Smooth loss	FD	Non-parametric
a)	✓		✓	✓
b)	✓		✓	✓
c)	✓	✓	✓	

## Outline: bibliography

- a) *Non-parametric Stochastic Approximation with Large Step-sizes*,  
A. Dieuleveut and F. Bach, in the *Annals of Statistics*
- b) *Harder, Better, Faster, Stronger Convergence Rates for Least-squares Regression*,  
A. Dieuleveut, N. Flammarion and F. Bach, in *Journal of Machine Learning Research*
- c) *Bridging the Gap between Constant Step Size Stochastic Gradient Descent and Markov Chains*,  
A. Dieuleveut, A. Durmus, F. Bach, under submission.

	Quadratic loss	Smooth loss	FD	Non-parametric
a)	✓		✓	✓
b)	✓		✓	✓
c)	✓	✓	✓	

Part 1

## Outline: bibliography

- a) *Non-parametric Stochastic Approximation with Large Step-sizes*,  
A. Dieuleveut and F. Bach, in the *Annals of Statistics*
- b) *Harder, Better, Faster, Stronger Convergence Rates for Least-squares Regression*,  
A. Dieuleveut, N. Flammarion and F. Bach, in *Journal of Machine Learning Research*
- c) *Bridging the Gap between Constant Step Size Stochastic Gradient Descent and Markov Chains*,  
A. Dieuleveut, A. Durmus, F. Bach, under submission.

	Quadratic loss	Smooth loss	FD	Non-parametric
a)	✓		✓	✓
b)	✓		✓	✓
c)	✓	✓	✓	

Part 1 – Part 2

## Outline: bibliography

- a) *Non-parametric Stochastic Approximation with Large Step-sizes*,  
A. Dieuleveut and F. Bach, in the *Annals of Statistics*
- b) *Harder, Better, Faster, Stronger Convergence Rates for Least-squares Regression*,  
A. Dieuleveut, N. Flammarion and F. Bach, in *Journal of Machine Learning Research*
- c) *Bridging the Gap between Constant Step Size Stochastic Gradient Descent and Markov Chains*,  
A. Dieuleveut, A. Durmus, F. Bach, under submission.

	Quadratic loss	Smooth loss	FD	Non-parametric
a)	✓		✓	✓
b)	✓		✓	✓
c)	✓	✓	✓	

Part 1 – Part 2 – Part 3

# Outline

1. Introduction.
2. A warm up! Results in finite dimension, ( $d \gg n$ )
  - ▶ Averaged stochastic descent: adaptivity
  - ▶ Acceleration: two optimal rates
3. Non-parametric stochastic approximation
4. Stochastic approximation as a Markov chain: extension to non quadratic loss functions.

# Behavior of Stochastic Approximation in high dimension

Least-squares regression in **finite dimension**:

$$\mathcal{R}(\theta) = \mathbb{E}_\rho \left[ \left( \langle \theta, \Phi(X) \rangle - Y \right)^2 \right].$$

## Behavior of Stochastic Approximation in high dimension

Least-squares regression in **finite dimension**:

$$\mathcal{R}(\theta) = \mathbb{E}_\rho \left[ \left( \langle \theta, \Phi(X) \rangle - Y \right)^2 \right].$$

Let  $\Sigma = \mathbb{E} [\Phi(X)\Phi(X)^\top] \in \mathbb{R}^{d \times d}$ : for  $\theta_*$  the best linear predictor,

$$\mathcal{R}(\theta) - \mathcal{R}(\theta_*) = \left\| \Sigma^{1/2}(\theta - \theta_*) \right\|^2.$$

Let  $R^2 := \mathbb{E} [\|\Phi(X)\|^2]$ ,  $\sigma^2 := \mathbb{E} [(Y - \langle \theta_*, \Phi(X) \rangle)^2]$ .



# Behavior of Stochastic Approximation in high dimension

Least-squares regression in **finite dimension**:

$$\mathcal{R}(\theta) = \mathbb{E}_\rho \left[ \left( \langle \theta, \Phi(X) \rangle - Y \right)^2 \right].$$

Let  $\Sigma = \mathbb{E} [\Phi(X)\Phi(X)^\top] \in \mathbb{R}^{d \times d}$ : for  $\theta_*$  the best linear predictor,

$$\mathcal{R}(\theta) - \mathcal{R}(\theta_*) = \left\| \Sigma^{1/2}(\theta - \theta_*) \right\|^2.$$

Let  $R^2 := \mathbb{E} [\|\Phi(X)\|^2]$ ,  $\sigma^2 := \mathbb{E} [(Y - \langle \theta_*, \Phi(X) \rangle)^2]$ .

Consider stochastic gradient descent (*a.k.a.*, *Least-Mean-Squares*)

## Theorem

For any  $\gamma \leq \frac{1}{4R^2}$ , for any  $\alpha > 1$ , for any  $r \geq 0$ , for any  $n \in \mathbb{N}$ ,

$$\mathbb{E} \mathcal{R}(\bar{\theta}_n) - \mathcal{R}(\theta_*) \leq \frac{4\sigma^2 \gamma^{1/\alpha} \text{tr}(\Sigma^{1/\alpha})}{n^{1-1/\alpha}} + \frac{4 \left\| \Sigma^{1/2-r}(\theta_* - \theta_0) \right\|^2}{\gamma^{2r} n^{\min(2r, 2)}}.$$

# Theorem 1<sup>†</sup>, consequences

## Theorem

For any  $\gamma \leq \frac{1}{4R^2}$ , for any  $\alpha > 1$ , for any  $r \geq 0$ , for any  $n \in \mathbb{N}$ ,

$$\mathbb{E}\mathcal{R}(\bar{\theta}_n) - \mathcal{R}(\theta_*) \leq \underbrace{\frac{4\sigma^2\gamma^{1/\alpha} \text{tr}(\Sigma^{1/\alpha})}{n^{1-1/\alpha}}}_{\text{Variance}} + \underbrace{\frac{4\|\Sigma^{1/2-r}(\theta_* - \theta_0)\|^2}{\gamma^{2r}n^{\min(2r,2)}}}_{\text{Bias}}.$$

---

<sup>†</sup>Dieuleveut and Bach [2015].

# Theorem 1<sup>†</sup>, consequences

## Theorem

For any  $\gamma \leq \frac{1}{4R^2}$ , for any  $\alpha > 1$ , for any  $r \geq 0$ , for any  $n \in \mathbb{N}$ ,

$$\mathbb{E}\mathcal{R}(\bar{\theta}_n) - \mathcal{R}(\theta_*) \leq \underbrace{\frac{4\sigma^2\gamma^{1/\alpha} \text{tr}(\Sigma^{1/\alpha})}{n^{1-1/\alpha}}}_{\text{Variance}} + \underbrace{\frac{4\|\Sigma^{1/2-r}(\theta_* - \theta_0)\|^2}{\gamma^{2r} n^{\min(2r,2)}}}_{\text{Bias}}.$$

Variance term

$$\gamma\sigma^2 \text{tr}(\Sigma)$$

$$\alpha = 1$$

$$\frac{\sigma^2 d}{n}$$

$$\alpha \rightarrow \infty$$

Bias term

$$\frac{\|\theta_* - \theta_0\|^2}{\gamma n}$$

$$r = 1/2.$$

$$\frac{\|\Sigma^{-1/2}(\theta_* - \theta_0)\|^2}{\gamma^2 n^2}$$

$$r = 1.$$

<sup>†</sup>Dieuleveut and Bach [2015].

# Theorem 1<sup>†</sup>, consequences

## Theorem

For any  $\gamma \leq \frac{1}{4R^2}$ , for any  $\alpha > 1$ , for any  $r \geq 0$ , for any  $n \in \mathbb{N}$ ,

$$\mathbb{E}\mathcal{R}(\bar{\theta}_n) - \mathcal{R}(\theta_*) \leq \underbrace{\frac{4\sigma^2\gamma^{1/\alpha} \text{tr}(\Sigma^{1/\alpha})}{n^{1-1/\alpha}}}_{\text{Variance}} + \underbrace{\frac{4 \|\Sigma^{1/2-r}(\theta_* - \theta_0)\|^2}{\gamma^{2r} n^{\min(2r, 2)}}}_{\text{Bias}}.$$

Variance term

$$\gamma\sigma^2 \text{tr}(\Sigma)$$

$$\alpha = 1$$

$$\frac{\sigma^2 d}{n}$$

$$\alpha \rightarrow \infty$$

$$\frac{\|\theta_* - \theta_0\|^2}{\gamma n}$$

$$r = 1/2.$$

Recovers  
Bach and Moulines [2013]

Bias term

$$\frac{\|\Sigma^{-1/2}(\theta_* - \theta_0)\|^2}{\gamma^2 n^2}$$

$$r = 1.$$

Improves  
asymptotic Bias

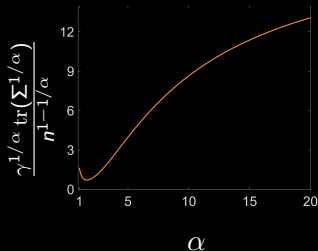
<sup>†</sup>Dieuleveut and Bach [2015].

# Theorem 1, consequences

## Theorem

For any  $\gamma \leq \frac{1}{4R^2}$ , for any  $n \in \mathbb{N}$ ,

$$\mathbb{E}\mathcal{R}(\bar{\theta}_n) - \mathcal{R}(\theta_*) \leq \inf_{\alpha > 1, r \geq 0} \left( \underbrace{\frac{4\sigma^2\gamma^{1/\alpha} \text{tr}(\Sigma^{1/\alpha})}{n^{1-1/\alpha}}}_{\text{Variance}} + \underbrace{\frac{4 \|\Sigma^{1/2-r}(\theta_* - \theta_0)\|^2}{\gamma^{2r} n^{\min(2r, 2)}}}_{\text{Bias}} \right).$$



## Adaptivity

Upper bound on the variance term as a function of  $\alpha$ .

$$d \gg n.$$

# Limits to SA performance: two lower bounds

Stochastic Approximation in Supervised ML

# Limits to SA performance: two lower bounds

## Stochastic Approximation in Supervised ML

Builds an estimator **given**  $n$   
**observations.**

↪ statistical **lower bound:**

$$\frac{\sigma^2 d}{n}$$

# Limits to SA performance: two lower bounds

## Stochastic Approximation in Supervised ML

Builds an estimator **given**  $n$  **observations**.

↪ statistical **lower bound**:

$$\frac{\sigma^2 d}{n}$$

Approximates the minimum of an ( $L$ -smooth) function **in**  $t$  **iterations**, using first order information.

↪ optimization **lower bound**:

$$\frac{L \|\theta_0 - \theta_*\|^2}{t^2}.$$



# Limits to SA performance: two lower bounds

## Stochastic Approximation in Supervised ML

Builds an estimator **given**  $n$  **observations**.

↪ statistical **lower bound**:

$$\frac{\sigma^2 d}{n}$$

Approximates the minimum of an ( $L$ -smooth) function **in**  $t$  **iterations**, using first order information.

↪ optimization **lower bound**:

$$\frac{L \|\theta_0 - \theta_*\|^2}{t^2}.$$

here,  $n = t$ .

# Limits to SA performance: two lower bounds

## Stochastic Approximation in Supervised ML

Approximates the minimum of an ( $L$ -smooth) function **in  $n$  iterations**, using first order information.

Builds an estimator **given  $n$  observations**.

↪ statistical **lower bound**:

$$\frac{\sigma^2 d}{n}$$

↪ optimization **lower bound**:

$$\frac{L \|\theta_0 - \theta_*\|^2}{n^2}.$$

here,  $n = t$ .

# Limits to SA performance: two lower bounds

## Stochastic Approximation in Supervised ML

Approximates the minimum of an ( $L$ -smooth) function **in  $n$  iterations**, using first order information.

Builds an estimator **given  $n$  observations**.

↪ statistical **lower bound**:

↪ optimization **lower bound**:

$$\frac{\sigma^2 d}{n}$$

$$\frac{L \|\theta_0 - \theta_*\|^2}{n^2}.$$

here,  $n = t$ .

Theorem 1, for Av-SGD, gives as upper bound:

$$\frac{\sigma^2 d}{n} + \min \left( \frac{L \|\theta_0 - \theta_*\|^2}{n}; \frac{L^2 \|\Sigma^{-1/2}(\theta_0 - \theta_*)\|^2}{n^2} \right).$$

## Acceleration<sup>†</sup>

Optimal rate (*for deterministic optimization*), is achieved by

**accelerated gradient descent:**

$$\begin{cases} \theta_n &= \eta_{n-1} - \gamma_n f'(\eta_{n-1}) \\ \eta_n &= \theta_n + \delta_n(\theta_n - \theta_{n-1}) . \end{cases}$$

---

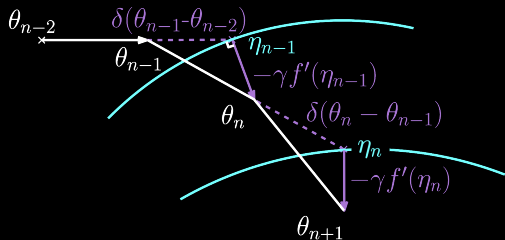
<sup>†</sup>Dieuleveut, Flammarion, Bach [2016]



## Acceleration<sup>†</sup>

Optimal rate (for deterministic optimization), is achieved by **accelerated gradient descent**:

$$\begin{cases} \theta_n = \eta_{n-1} - \gamma_n f'(\eta_{n-1}) \\ \eta_n = \theta_n + \delta_n(\theta_n - \theta_{n-1}) \end{cases}$$



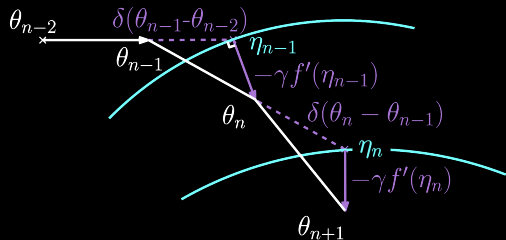
Problem: acceleration is sensitive to noise [d'Aspremont, 2008].

<sup>†</sup>Dieuleveut, Flammarion, Bach [2016]

## Acceleration<sup>†</sup>

Optimal rate (for deterministic optimization), is achieved by **accelerated gradient descent**:

$$\begin{cases} \theta_n &= \eta_{n-1} - \gamma_n f'(\eta_{n-1}) \\ \eta_n &= \theta_n + \delta_n(\theta_n - \theta_{n-1}). \end{cases}$$



Problem: acceleration is sensitive to noise [d'Aspremont, 2008].

Combining SGD, acceleration and averaging,

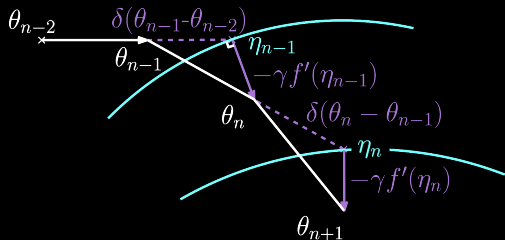
- ▶ using extra regularization,

<sup>†</sup>Dieuleveut, Flammarion, Bach [2016]

## Acceleration<sup>†</sup>

Optimal rate (for deterministic optimization), is achieved by **accelerated gradient descent**:

$$\begin{cases} \theta_n &= \eta_{n-1} - \gamma_n f'(\eta_{n-1}) \\ \eta_n &= \theta_n + \delta_n(\theta_n - \theta_{n-1}). \end{cases}$$



Problem: acceleration is sensitive to noise [d'Aspremont, 2008].

Combining SGD, acceleration and averaging,

- ▶ using extra regularization,
- ▶ and for “additive” noise model only,

<sup>†</sup>Dieuleveut, Flammarion, Bach [2016]

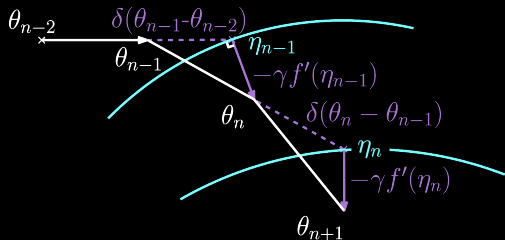




## Acceleration<sup>†</sup>

Optimal rate (for deterministic optimization), is achieved by **accelerated gradient descent**:

$$\begin{cases} \theta_n &= \eta_{n-1} - \gamma_n f'(\eta_{n-1}) \\ \eta_n &= \theta_n + \delta_n(\theta_n - \theta_{n-1}). \end{cases}$$



Problem: acceleration is sensitive to noise [d'Aspremont, 2008].

Combining SGD, acceleration and averaging,

- ▶ using extra regularization,
- ▶ and for “additive” noise model only,

we achieve **both of the optimal rates**.

*Caveat:* LMS recursion does not provide an additive noise oracle.

Different recursion with  $\Sigma$  known.

<sup>†</sup>Dieuleveut, Flammarion, Bach [2016]

## Acceleration and averaging

More precisely we consider:

$$\begin{aligned}\theta_n &= \nu_{n-1} - \gamma \mathcal{R}'_n(\nu_{n-1}) - \gamma \lambda (\nu_{n-1} - \theta_0) \\ \nu_n &= \theta_n + \delta (\theta_n - \theta_{n-1}),\end{aligned}$$

### Theorem

For any  $\gamma \leq 1/2R^2$ , for  $\delta = 1$ , and  $\lambda = 0$ ,

$$\mathbb{E} [\mathcal{R}(\bar{\theta}_n)] - \mathcal{R}(\theta_*) \leq 8 \frac{\sigma^2 d}{n+1} + 36 \frac{\|\theta_0 - \theta_*\|^2}{\gamma(n+1)^2}.$$

**Optimal rate from both statistical and optimization point of view.**

# Outline

1. Introduction.
2. A warm up! Results in finite dimension, ( $d \gg n$ )
3. Non-parametric stochastic approximation
  - ▶ Averaged stochastic descent: statistical rate of convergence
  - ▶ Acceleration: improving convergence in ill-conditioned regimes
4. Stochastic approximation as a Markov chain: extension to non quadratic loss functions.

# Non-parametric Random Design Least Squares Regression

Goal:

$$\min_g \mathcal{R}(g) = \mathbb{E}_\rho [(Y - g(X))^2]$$

# Non-parametric Random Design Least Squares Regression

Goal:

$$\min_g \mathcal{R}(g) = \mathbb{E}_\rho [(Y - g(X))^2]$$

- ▶  $\rho_X$  marginal distribution of  $X$  in  $\mathcal{X}$ ,
- ▶  $L^2_{\rho_X}$  set of squared integrable functions w.r.t.  $\rho_X$ .

# Non-parametric Random Design Least Squares Regression

Goal:

$$\min_g \mathcal{R}(g) = \mathbb{E}_\rho [(Y - g(X))^2]$$

- ▶  $\rho_X$  marginal distribution of  $X$  in  $\mathcal{X}$ ,
- ▶  $L_{\rho_X}^2$  set of squared integrable functions w.r.t.  $\rho_X$ .

Bayes predictor minimizes the quadratic risk over  $L_{\rho_X}^2$ :

$$g_\rho(X) = \mathbb{E}[Y|X].$$

# Non-parametric Random Design Least Squares Regression

Goal:

$$\min_g \mathcal{R}(g) = \mathbb{E}_\rho [(Y - g(X))^2]$$

- ▶  $\rho_X$  marginal distribution of  $X$  in  $\mathcal{X}$ ,
- ▶  $L_{\rho_X}^2$  set of squared integrable functions w.r.t.  $\rho_X$ .

Bayes predictor minimizes the quadratic risk over  $L_{\rho_X}^2$ :

$$g_\rho(X) = \mathbb{E}[Y|X].$$

Moreover, for any function  $g$  in  $L_{\rho_X}^2$ , the excess risk is:

$$\mathcal{R}(g) - \mathcal{R}(g_\rho) = \|g - g_\rho\|_{L_{\rho_X}^2}^2.$$



# Non-parametric Random Design Least Squares Regression

Goal:

$$\min_g \mathcal{R}(g) = \mathbb{E}_\rho [(Y - g(X))^2]$$

- ▶  $\rho_X$  marginal distribution of  $X$  in  $\mathcal{X}$ ,
- ▶  $L_{\rho_X}^2$  set of squared integrable functions w.r.t.  $\rho_X$ .

Bayes predictor minimizes the quadratic risk over  $L_{\rho_X}^2$ :

$$g_\rho(X) = \mathbb{E}[Y|X].$$

Moreover, for any function  $g$  in  $L_{\rho_X}^2$ , the excess risk is:

$$\mathcal{R}(g) - \mathcal{R}(g_\rho) = \|g - g_\rho\|_{L_{\rho_X}^2}^2.$$

$\mathcal{H}$  a space of functions: there exists  $g_{\mathcal{H}} \in \bar{\mathcal{H}}^{L_{\rho_X}^2}$  such that

$$\mathcal{R}(g_{\mathcal{H}}) = \inf_{g \in \mathcal{H}} \mathcal{R}(g).$$

# Reproducing Kernel Hilbert Space

## Definition

A Reproducing Kernel Hilbert Space (RKHS)  $\mathcal{H}$  is a space of functions from  $\mathcal{X}$  into  $\mathbb{R}$ , such that there exists a reproducing kernel  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , satisfying:

- ▶ For any  $x \in \mathcal{X}$ ,  $\mathcal{H}$  contains the function  $K_x$ , defined by:

$$\begin{aligned} K_x : \mathcal{X} &\rightarrow \mathbb{R} \\ z &\mapsto K(x, z). \end{aligned}$$

# Reproducing Kernel Hilbert Space

## Definition

A **Reproducing Kernel Hilbert Space (RKHS)**  $\mathcal{H}$  is a space of functions from  $\mathcal{X}$  into  $\mathbb{R}$ , such that there exists a reproducing kernel  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , satisfying:

- ▶ For any  $x \in \mathcal{X}$ ,  $\mathcal{H}$  contains the function  $K_x$ , defined by:

$$\begin{aligned} K_x : \mathcal{X} &\rightarrow \mathbb{R} \\ z &\mapsto K(x, z). \end{aligned}$$

- ▶ For any  $x \in \mathcal{X}$  and  $f \in \mathcal{H}$ , the *reproducing property* holds:

$$\langle K_x, f \rangle_{\mathcal{H}} = f(x).$$

# Why are RKHS so nice?

## ▶ Computation:

- ▶ Linear spaces of functions.
- ▶ Existence of gradients (Hilbert).
- ▶ Possible to compute inner products thanks to the reproducing property.
- ▶ Only deal with functions in the set  $\text{span}\{K_{x_i}, i = 1 \dots n\}$  (representer theorem).

↪ the algebraic framework is preserved !

# Why are RKHS so nice?

## ▶ Computation:

- ▶ Linear spaces of functions.
- ▶ Existence of gradients (Hilbert).
- ▶ Possible to compute inner products thanks to the reproducing property.
- ▶ Only deal with functions in the set  $\text{span}\{K_{x_i}, i = 1 \dots n\}$  (representer theorem).

↔ the algebraic framework is preserved !

- ▶ **Approximation:** many kernels satisfy  $\bar{\mathcal{H}}^{L^2_{\rho_X}} = L^2_{\rho_X}$ , there is no approximation error !

# Why are RKHS so nice?

## ► Computation:

- Linear spaces of functions.
- Existence of gradients (Hilbert).
- Possible to compute inner products thanks to the reproducing property.
- Only deal with functions in the set  $\text{span}\{K_{x_i}, i = 1 \dots n\}$  (representer theorem).

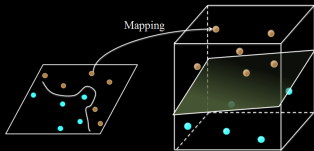
↪ the algebraic framework is preserved !

- **Approximation:** many kernels satisfy  $\bar{\mathcal{H}}^{L^2_{\rho_X}} = L^2_{\rho_X}$ , there is no approximation error !

- **Representation:** Feature map,

$$\begin{aligned} \mathcal{X} &\rightarrow \mathcal{H} \\ x &\mapsto K_x \end{aligned}$$

maps points from *any* set into a linear space to apply a linear method.



## Stochastic approximation in the RKHS.

As  $\mathcal{R}(g) = \mathbb{E} [(\langle g, K_X \rangle_{\mathcal{H}} - Y)^2]$ , for each pair of observations

$$(\langle g, K_{x_n} \rangle_{\mathcal{H}} - y_n)K_{x_n} = (g(x_n) - y_n)K_{x_n}$$

is an *unbiased stochastic gradient* of  $\mathcal{R}$  at  $g$ .

## Stochastic approximation in the RKHS.

As  $\mathcal{R}(g) = \mathbb{E} [(\langle g, K_X \rangle_{\mathcal{H}} - Y)^2]$ , for each pair of observations

$$(\langle g, K_{x_n} \rangle_{\mathcal{H}} - y_n) K_{x_n} = (g(x_n) - y_n) K_{x_n}$$

is an *unbiased stochastic gradient* of  $\mathcal{R}$  at  $g$ .

Consider the *stochastic gradient recursion*, starting from  $g_0 \in \mathcal{H}$ :

$$g_n = g_{n-1} - \gamma [\langle g_{n-1}, K_{x_n} \rangle_{\mathcal{H}} - y_n] K_{x_n},$$

where  $\gamma$  is the *step-size*.



## Stochastic approximation in the RKHS.

As  $\mathcal{R}(g) = \mathbb{E} [(\langle g, K_X \rangle_{\mathcal{H}} - Y)^2]$ , for each pair of observations

$$(\langle g, K_{x_n} \rangle_{\mathcal{H}} - y_n) K_{x_n} = (g(x_n) - y_n) K_{x_n}$$

is an *unbiased stochastic gradient* of  $\mathcal{R}$  at  $g$ .

Consider the *stochastic gradient recursion*, starting from  $g_0 \in \mathcal{H}$ :

$$g_n = g_{n-1} - \gamma [\langle g_{n-1}, K_{x_n} \rangle_{\mathcal{H}} - y_n] K_{x_n},$$

where  $\gamma$  is the *step-size*. Thus

$$g_n = \sum_{i=1}^n a_i K_{x_i},$$

with  $(a_n)_{n \geq 1}$ ,  $a_n = -\gamma_n (g_{n-1}(x_n) - y_n)$ . With averaging,

$$\bar{g}_n = \frac{1}{n+1} \sum_{k=0}^n g_k$$

**Total complexity:**  $O(n^2)$

# Kernel regression: Analysis

Assume  $\mathbb{E}[K(X, X)]$  and  $\mathbb{E}[Y^2]$  are finite. Define the *covariance operator*.

$$\Sigma = \mathbb{E} \left[ K_X K_X^\top \right].$$

We make two assumptions:

- ▶ Capacity condition: eigenvalue decay of  $\Sigma$ .
- ▶ Source condition: position of  $g_{\mathcal{H}}$  w.r.t. the kernel space  $\mathcal{H}$ .

# Kernel regression: Analysis

Assume  $\mathbb{E}[K(X, X)]$  and  $\mathbb{E}[Y^2]$  are finite. Define the *covariance operator*.

$$\Sigma = \mathbb{E} \left[ K_X K_X^\top \right].$$

We make two assumptions:

- ▶ Capacity condition: eigenvalue decay of  $\Sigma$ .
- ▶ Source condition: position of  $g_{\mathcal{H}}$  w.r.t. the kernel space  $\mathcal{H}$ .

$\Sigma$  is a **trace-class operator**, that can be decomposed over its eigen-spaces. Its power:  $\Sigma^\tau$ ,  $\tau > 0$ . are thus well defined.

## Capacity condition (CC)

**CC( $\alpha$ ):** for some  $\alpha > 1$ , we assume that  $\text{tr}(\Sigma^{1/\alpha}) < \infty$ .

## Capacity condition (CC)

**CC( $\alpha$ ):** for some  $\alpha > 1$ , we assume that  $\text{tr}(\Sigma^{1/\alpha}) < \infty$ .

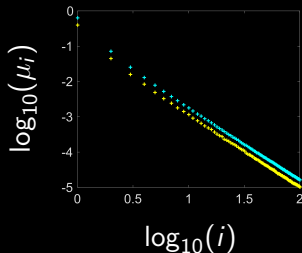
If we denote  $(\mu_i)_{i \in I}$  the sequence of non-zero eigenvalues of the operator  $\Sigma$ , in decreasing order, then  $\mu_i = O(i^{-\alpha})$ .

# Capacity condition (CC)

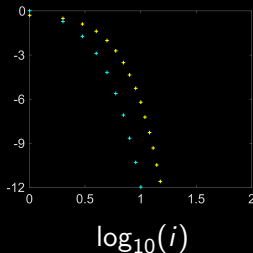
**CC( $\alpha$ ):** for some  $\alpha > 1$ , we assume that  $\text{tr}(\Sigma^{1/\alpha}) < \infty$ .

If we denote  $(\mu_i)_{i \in I}$  the sequence of non-zero eigenvalues of the operator  $\Sigma$ , in decreasing order, then  $\mu_i = O(i^{-\alpha})$ .

Sobolev first order kernel



Gaussian kernel



Eigenvalue decay of the covariance operator.

*Left:* min kernel,  $\rho_X = \mathcal{U}[0; 1]$ ,  $\rightarrow$  **CC( $\alpha = 2$ )**.

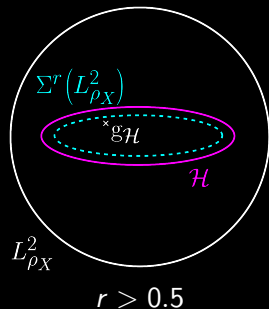
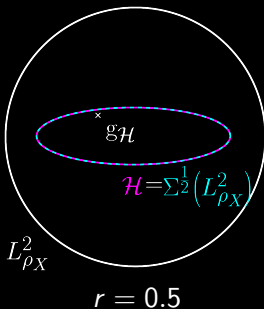
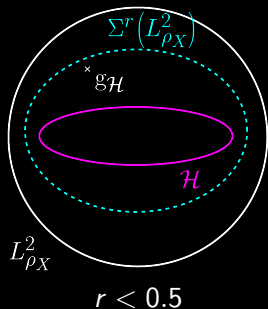
*Right:* Gaussian kernel,  $\rho_X = \mathcal{U}[-1; 1]$ .  $\rightarrow$  **CC( $\alpha$ )**,  $\forall \alpha \geq 1$ .

## Source condition (SC)

Concerning the optimal function  $g_{\mathcal{H}}$ , we assume:

**SC( $r$ ):** for some  $r \geq 0$ ,  $g_{\mathcal{H}} \in \Sigma^r(L^2_{\rho_X})$

Thus  $\|\Sigma^{-r}(g_{\mathcal{H}})\|_{L^2_{\rho_X}} < \infty$ .



## NPSA with large step sizes

### Theorem

Assume  $\text{CC}(\alpha)$  and  $\text{SC}(r)$ . Then for any  $\gamma \leq \frac{1}{4R^2}$ ,

$$\mathbb{E}\mathcal{R}(\bar{g}_n) - \mathcal{R}(g_{\mathcal{H}}) \leq \frac{4\sigma^2\gamma^{1/\alpha} \text{tr}(\Sigma^{1/\alpha})}{n^{1-1/\alpha}} + 4 \frac{\|\Sigma^{-r}(g_{\mathcal{H}} - g_0)\|_{L^2_{\rho_X}}^2}{\gamma^{2r} n^{\min(2r, 2)}}.$$

for  $\gamma = \gamma_0 n^{\frac{-2\alpha r - 1 + \alpha}{2\alpha r + 1}}$ , for  $\frac{\alpha - 1}{2\alpha} \leq r \leq 1$

$$\mathbb{E}\mathcal{R}(\bar{g}_n) - \mathcal{R}(g_{\mathcal{H}}) \leq n^{\frac{-2\alpha r}{2\alpha r + 1}} \left( 4\sigma^2 \text{tr}(\Sigma^{1/\alpha}) + 4 \|\Sigma^{-r}(g_{\mathcal{H}} - g_0)\|_{L^2_{\rho_X}}^2 \right).$$



# NPSA with large step sizes

## Theorem

Assume  $\text{CC}(\alpha)$  and  $\text{SC}(r)$ . Then for any  $\gamma \leq \frac{1}{4R^2}$ ,

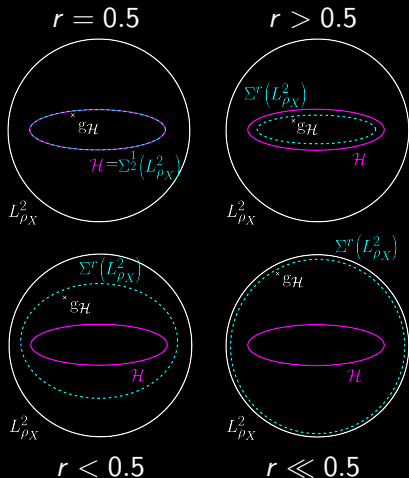
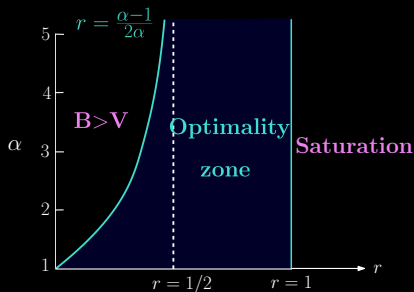
$$\mathbb{E}\mathcal{R}(\bar{g}_n) - \mathcal{R}(g_{\mathcal{H}}) \leq \frac{4\sigma^2\gamma^{1/\alpha} \text{tr}(\Sigma^{1/\alpha})}{n^{1-1/\alpha}} + 4 \frac{\|\Sigma^{-r}(g_{\mathcal{H}} - g_0)\|_{L^2_{\rho_X}}^2}{\gamma^{2r} n^{\min(2r, 2)}}.$$

for  $\gamma = \gamma_0 n^{\frac{-2\alpha r - 1 + \alpha}{2\alpha r + 1}}$ , for  $\frac{\alpha-1}{2\alpha} \leq r \leq 1$

$$\mathbb{E}\mathcal{R}(\bar{g}_n) - \mathcal{R}(g_{\mathcal{H}}) \leq n^{\frac{-2\alpha r}{2\alpha r + 1}} \left( 4\sigma^2 \text{tr}(\Sigma^{1/\alpha}) + 4 \|\Sigma^{-r}(g_{\mathcal{H}} - g_0)\|_{L^2_{\rho_X}}^2 \right).$$

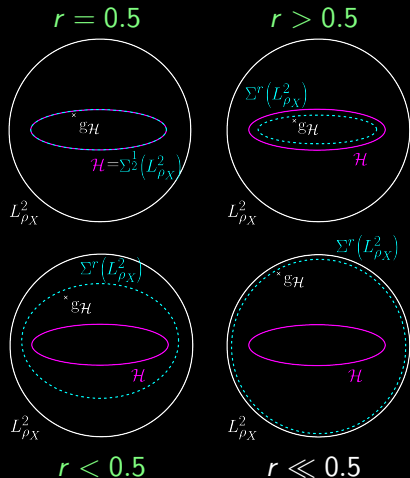
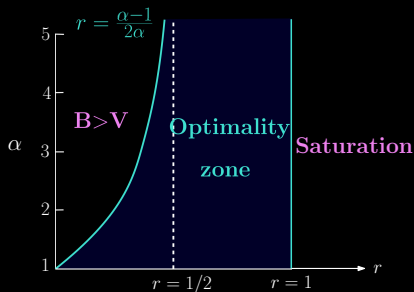
- ▶ **Statistically optimal rate.** [Caponnetto and De Vito, 2007].
- ▶ *Beyond:* online, minimal assumptions...

# Optimality regions



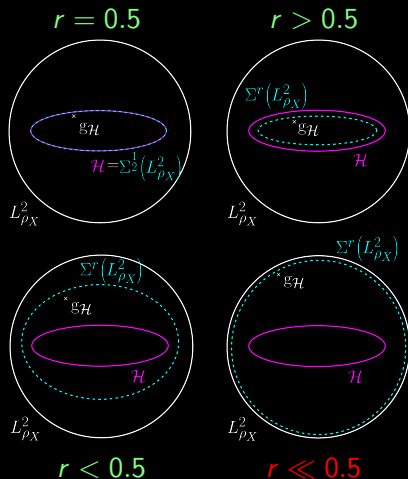
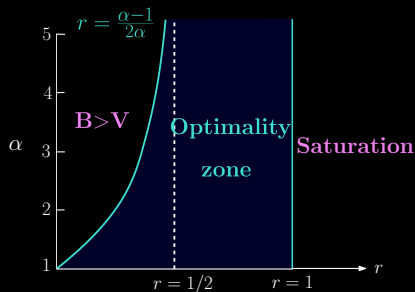
Optimal rate in RKHS can be achieved via large step size and averaging in many situations.

# Optimality regions



Optimal rate in RKHS can be achieved via large step size and averaging in many situations.

# Optimality regions



Optimal rate in RKHS can be achieved via large step size and averaging in many situations.

## Acceleration: Reproducing kernel Hilbert space setting

We consider the RKHS setting presented before.

### Theorem

Assume  $\text{CC}(\alpha)$  and  $\text{SC}(r)$ . Then for  $\gamma = \gamma_0 n^{-\frac{4r\alpha+2-\alpha}{2r\alpha+1}}$ , for  $\lambda = \frac{1}{\gamma n^2}$ , for  $r \geq \frac{\alpha-2}{2\alpha}$ ,

$$\mathbb{E}\mathcal{R}(\bar{g}_n) - \mathcal{R}(g_{\mathcal{H}}) \leq C_{\theta_0, \theta_*, \Sigma} n^{\frac{-2\alpha r}{2\alpha r+1}}.$$

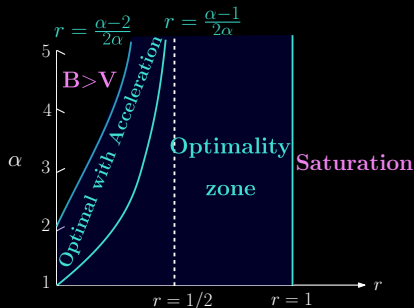
# Acceleration: Reproducing kernel Hilbert space setting

We consider the RKHS setting presented before.

## Theorem

Assume  $\text{CC}(\alpha)$  and  $\text{SC}(r)$ . Then for  $\gamma = \gamma_0 n^{-\frac{4r\alpha+2-\alpha}{2r\alpha+1}}$ , for  $\lambda = \frac{1}{\gamma n^2}$ , for  $r \geq \frac{\alpha-2}{2\alpha}$ ,

$$\mathbb{E}\mathcal{R}(\bar{g}_n) - \mathcal{R}(g_{\mathcal{H}}) \leq C_{\theta_0, \theta_*, \Sigma} n^{\frac{-2\alpha r}{2\alpha r+1}}.$$



## Least squares: some conclusions

- ▶ Provide **optimal rate of convergence** under two assumptions for non-parametric regression in Hilbert spaces: **large step sizes and averaging**.

## Least squares: some conclusions

- ▶ Provide **optimal rate of convergence** under two assumptions for non-parametric regression in Hilbert spaces: **large step sizes and averaging**.
- ▶ Sheds some light on FD case.



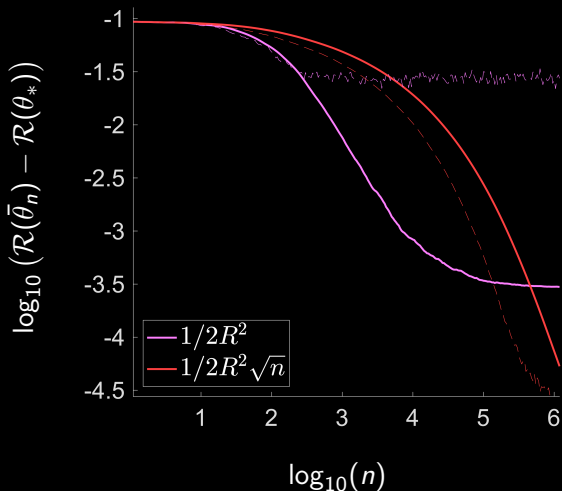
## Least squares: some conclusions

- ▶ Provide **optimal rate of convergence** under two assumptions for non-parametric regression in Hilbert spaces: **large step sizes and averaging**.
- ▶ Sheds some light on FD case.
- ▶ Possible to attain simultaneously optimal rate from the statistical and optimization point of view.

# Outline

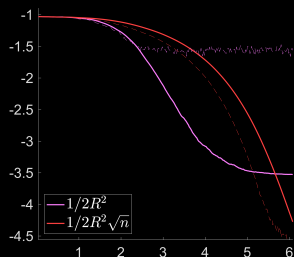
1. Introduction.
2. Non-parametric stochastic approximation
3. Faster rates with acceleration
4. Stochastic approximation as a Markov chain: extension to non quadratic loss functions.
  - ▶ Motivation
  - ▶ Assumptions
  - ▶ Convergence in Wasserstein distance.

## Motivation 1/ 2. Large step sizes!



Logistic regression. Final iterate (dashed), and averaged recursion (plain).

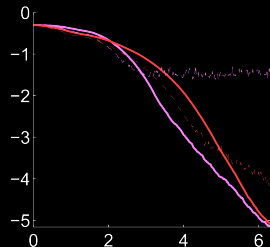
## Motivation 2/ 2. Difference between quadratic and logistic loss



Logistic Regression

$$\mathbb{E}\mathcal{R}(\bar{\theta}_n) - \mathcal{R}(\theta_*) = O(\gamma^2)$$

$$\text{with } \gamma = 1/(4R^2)$$



Least-Squares Regression

$$\mathbb{E}\mathcal{R}(\bar{\theta}_n) - \mathcal{R}(\theta_*) = O\left(\frac{1}{n}\right)$$

$$\text{with } \gamma = 1/(4R^2)$$

## SGD: an homogeneous Markov chain

Consider a  $L$ -smooth and  $\mu$ -strongly convex function  $\mathcal{R}$ .

## SGD: an homogeneous Markov chain

Consider a  $L$ -smooth and  $\mu$ -strongly convex function  $\mathcal{R}$ .

SGD with a step-size  $\gamma > 0$  is an homogeneous Markov chain:

$$\theta_{k+1}^\gamma = \theta_k^\gamma - \gamma [\mathcal{R}'(\theta_k^\gamma) + \varepsilon_{k+1}(\theta_k^\gamma)] ,$$

## SGD: an homogeneous Markov chain

Consider a  $L$ -smooth and  $\mu$ -strongly convex function  $\mathcal{R}$ .

SGD with a step-size  $\gamma > 0$  is an homogeneous Markov chain:

$$\theta_{k+1}^\gamma = \theta_k^\gamma - \gamma [\mathcal{R}'(\theta_k^\gamma) + \varepsilon_{k+1}(\theta_k^\gamma)] ,$$

- ▶ satisfies Markov property

# SGD: an homogeneous Markov chain

Consider a  $L$ -smooth and  $\mu$ -strongly convex function  $\mathcal{R}$ .

SGD with a step-size  $\gamma > 0$  is an homogeneous Markov chain:

$$\theta_{k+1}^\gamma = \theta_k^\gamma - \gamma [\mathcal{R}'(\theta_k^\gamma) + \varepsilon_{k+1}(\theta_k^\gamma)] ,$$

- ▶ satisfies Markov property
- ▶ is homogeneous, for  $\gamma$  constant,  $(\varepsilon_k)_{k \in \mathbb{N}}$  i.i.d.



## SGD: an homogeneous Markov chain

Consider a  $L$ -smooth and  $\mu$ -strongly convex function  $\mathcal{R}$ .

SGD with a step-size  $\gamma > 0$  is an homogeneous Markov chain:

$$\theta_{k+1}^\gamma = \theta_k^\gamma - \gamma [\mathcal{R}'(\theta_k^\gamma) + \varepsilon_{k+1}(\theta_k^\gamma)] ,$$

- ▶ satisfies Markov property
- ▶ is homogeneous, for  $\gamma$  constant,  $(\varepsilon_k)_{k \in \mathbb{N}}$  i.i.d.

Also assume:

- ▶  $\mathcal{R}'_k = \mathcal{R}' + \varepsilon_{k+1}$  is almost surely  $L$ -co-coercive.
- ▶ Bounded moments

$$\mathbb{E}[\|\varepsilon_k(\theta_*)\|^4] < \infty.$$

# Stochastic gradient descent as a Markov Chain: Analysis framework<sup>†</sup>

- ▶ Existence of a limit distribution  $\pi_\gamma$ , and linear convergence to this distribution:

$$\theta_n^\gamma \xrightarrow{d} \pi_\gamma.$$

---

<sup>†</sup>Dieuleveut, Durmus, Bach [2017].

# Stochastic gradient descent as a Markov Chain: Analysis framework<sup>†</sup>

- ▶ Existence of a limit distribution  $\pi_\gamma$ , and linear convergence to this distribution:

$$\theta_n^\gamma \xrightarrow{d} \pi_\gamma.$$

- ▶ Convergence of second order moments of the chain,

$$\bar{\theta}_{n,\gamma} \xrightarrow[n \rightarrow \infty]{L^2} \bar{\theta}_\gamma := \mathbb{E}_{\pi_\gamma} [\theta].$$

---

<sup>†</sup>Dieuleveut, Durmus, Bach [2017].

# Stochastic gradient descent as a Markov Chain: Analysis framework<sup>†</sup>

- ▶ Existence of a limit distribution  $\pi_\gamma$ , and linear convergence to this distribution:

$$\theta_n^\gamma \xrightarrow{d} \pi_\gamma.$$

- ▶ Convergence of second order moments of the chain,

$$\bar{\theta}_{n,\gamma} \xrightarrow[n \rightarrow \infty]{L^2} \bar{\theta}_\gamma := \mathbb{E}_{\pi_\gamma} [\theta].$$

- ▶ Behavior under the limit distribution ( $\gamma \rightarrow 0$ ):  $\bar{\theta}_\gamma = \theta_* + ?$ .

---

<sup>†</sup>Dieuleveut, Durmus, Bach [2017].

# Stochastic gradient descent as a Markov Chain: Analysis framework<sup>†</sup>

- ▶ Existence of a limit distribution  $\pi_\gamma$ , and linear convergence to this distribution:

$$\theta_n^\gamma \xrightarrow{d} \pi_\gamma.$$

- ▶ Convergence of second order moments of the chain,

$$\bar{\theta}_{n,\gamma} \xrightarrow[n \rightarrow \infty]{L^2} \bar{\theta}_\gamma := \mathbb{E}_{\pi_\gamma} [\theta].$$

- ▶ Behavior under the limit distribution ( $\gamma \rightarrow 0$ ):  $\bar{\theta}_\gamma = \theta_* + ?$ .

↷ Provable convergence improvement with extrapolation tricks.

---

<sup>†</sup>Dieuleveut, Durmus, Bach [2017].

# Existence of a limit distribution $\gamma \rightarrow 0$

**Goal:**  $(\theta_n^\gamma)_{n \geq 0} \xrightarrow{d} \pi_\gamma$ .

## Theorem

For any  $\gamma < L^{-1}$ , the chain  $(\theta_n^\gamma)_{n \geq 0}$  admits a unique stationary distribution  $\pi_\gamma$ . In addition for all  $\theta_0 \in \mathbb{R}^d$ ,  $n \in \mathbb{N}$ :

$$W_2^2(\theta_n^\gamma, \pi_\gamma) \leq (1 - 2\mu\gamma(1 - \gamma L))^n \int_{\mathbb{R}^d} \|\theta_0 - \vartheta\|^2 d\pi_\gamma(\vartheta).$$

# Existence of a limit distribution $\gamma \rightarrow 0$

**Goal:**  $(\theta_n^\gamma)_{n \geq 0} \xrightarrow{d} \pi_\gamma$ .

## Theorem

For any  $\gamma < L^{-1}$ , the chain  $(\theta_n^\gamma)_{n \geq 0}$  admits a unique stationary distribution  $\pi_\gamma$ . In addition for all  $\theta_0 \in \mathbb{R}^d$ ,  $n \in \mathbb{N}$ :

$$W_2^2(\theta_n^\gamma, \pi_\gamma) \leq (1 - 2\mu\gamma(1 - \gamma L))^n \int_{\mathbb{R}^d} \|\theta_0 - \vartheta\|^2 d\pi_\gamma(\vartheta).$$

**Wasserstein metric:** distance between probability measures.

## Behavior under limit distribution.

Ergodic theorem:  $\bar{\theta}_n \rightarrow \mathbb{E}_{\pi_\gamma}[\theta] =: \bar{\theta}_\gamma$ . Where is  $\bar{\theta}_\gamma$  ?



## Behavior under limit distribution.

Ergodic theorem:  $\bar{\theta}_n \rightarrow \mathbb{E}_{\pi_\gamma}[\theta] =: \bar{\theta}_\gamma$ . Where is  $\bar{\theta}_\gamma$  ?

If  $\theta_0 \sim \pi_\gamma$ , then  $\theta_1 \sim \pi_\gamma$ .

$$\theta_1^\gamma = \theta_0^\gamma - \gamma [\mathcal{R}'(\theta_0^\gamma) + \varepsilon_1(\theta_0^\gamma)] .$$

$$\mathbb{E}_{\pi_\gamma} [\mathcal{R}'(\theta)] = 0$$

## Behavior under limit distribution.

Ergodic theorem:  $\bar{\theta}_n \rightarrow \mathbb{E}_{\pi_\gamma}[\theta] =: \bar{\theta}_\gamma$ . Where is  $\bar{\theta}_\gamma$  ?

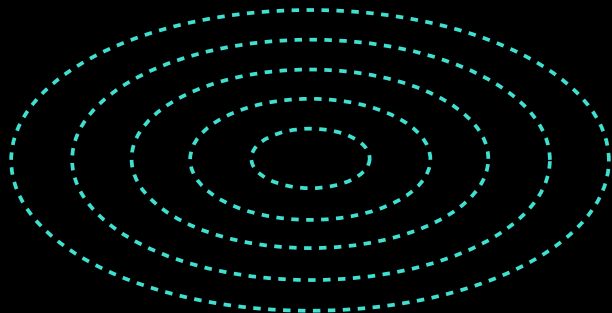
If  $\theta_0 \sim \pi_\gamma$ , then  $\theta_1 \sim \pi_\gamma$ .

$$\theta_1^\gamma = \theta_0^\gamma - \gamma [\mathcal{R}'(\theta_0^\gamma) + \varepsilon_1(\theta_0^\gamma)] .$$

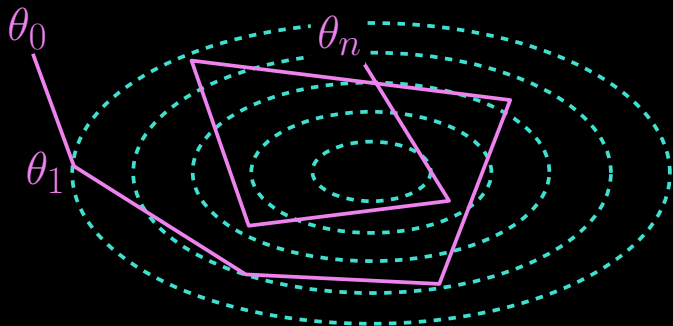
$$\mathbb{E}_{\pi_\gamma} [\mathcal{R}'(\theta)] = 0$$

In the **quadratic case** (linear gradients)  $\Sigma \mathbb{E}_{\pi_\gamma} [\theta - \theta_*] = 0$ :  $\bar{\theta}_\gamma = \theta_*$ !

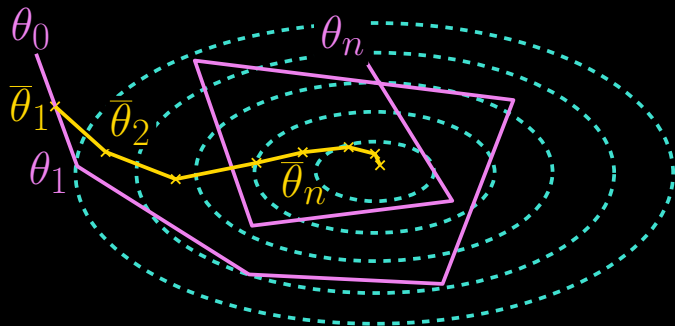
# Constant learning rate SGD: convergence in the quadratic case



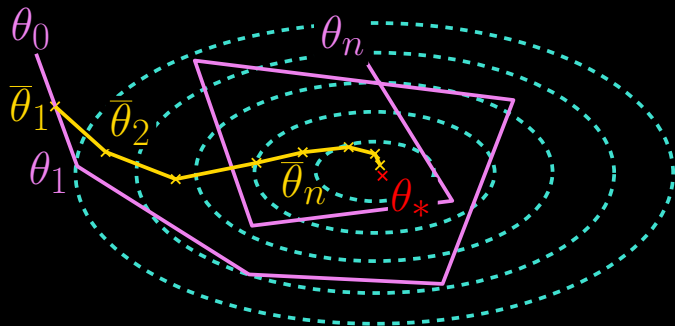
# Constant learning rate SGD: convergence in the quadratic case



# Constant learning rate SGD: convergence in the quadratic case



# Constant learning rate SGD: convergence in the quadratic case



## Behavior under limit distribution.

Ergodic theorem:  $\bar{\theta}_n \rightarrow \mathbb{E}_{\pi_\gamma}[\theta] =: \bar{\theta}_\gamma$ . Where is  $\bar{\theta}_\gamma$  ?

If  $\theta_0 \sim \pi_\gamma$ , then  $\theta_1 \sim \pi_\gamma$ .

$$\theta_1^\gamma = \theta_0^\gamma - \gamma [\mathcal{R}'(\theta_0^\gamma) + \varepsilon_1(\theta_0^\gamma)] .$$

$$\mathbb{E}_{\pi_\gamma} [\mathcal{R}'(\theta)] = 0$$

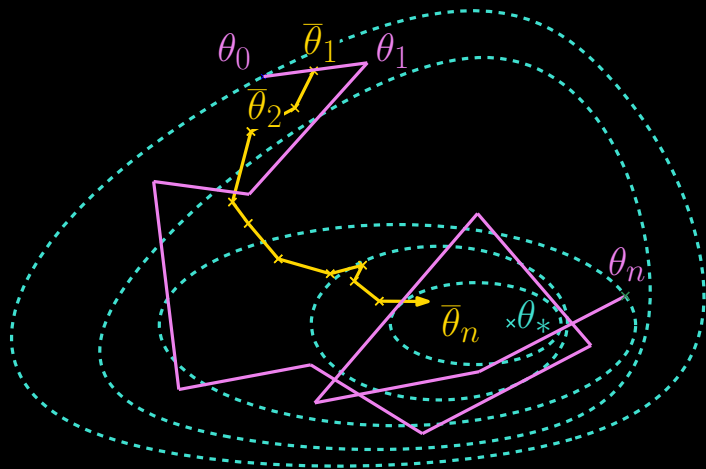
In the **quadratic case** (linear gradients)  $\Sigma \mathbb{E}_{\pi_\gamma} [\theta - \theta_*] = 0$ :  $\bar{\theta}_\gamma = \theta_*$ !

In the **general case**, Taylor expansion of  $\mathcal{R}$ , and same reasoning on higher moments of the chain leads to

$$\bar{\theta}_\gamma - \theta_* = \gamma \mathcal{R}''(\theta_*)^{-1} \mathcal{R}'''(\theta_*) \left( [\mathcal{R}''(\theta_*) \otimes I + I \otimes \mathcal{R}''(\theta_*)]^{-1} \mathbb{E}_{\pi_{\gamma, \varepsilon}} [\varepsilon(\theta)^{\otimes 2}] \right) + O(\gamma^2)$$

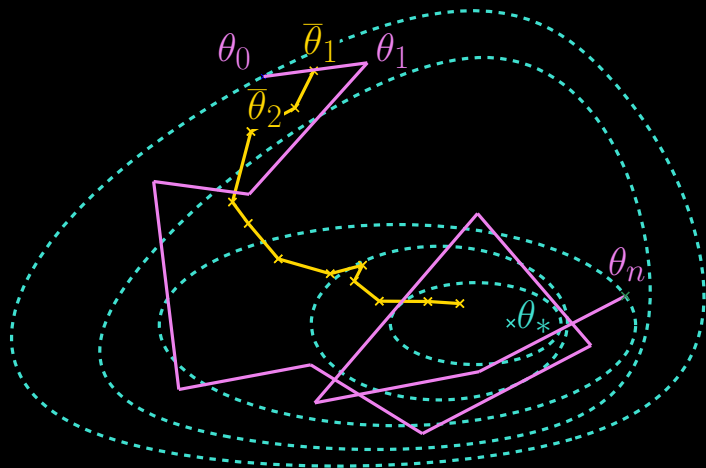
$$\text{Overall, } \bar{\theta}_\gamma - \theta_* = \gamma \Delta + O(\gamma^2).$$

# Constant learning rate SGD: convergence in the non-quadratic case

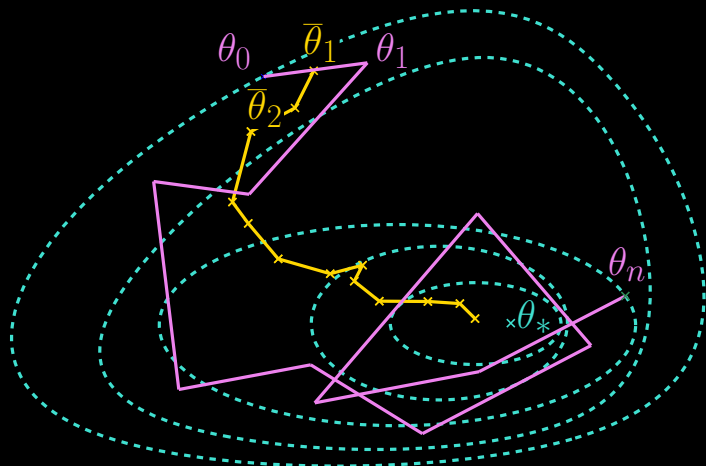




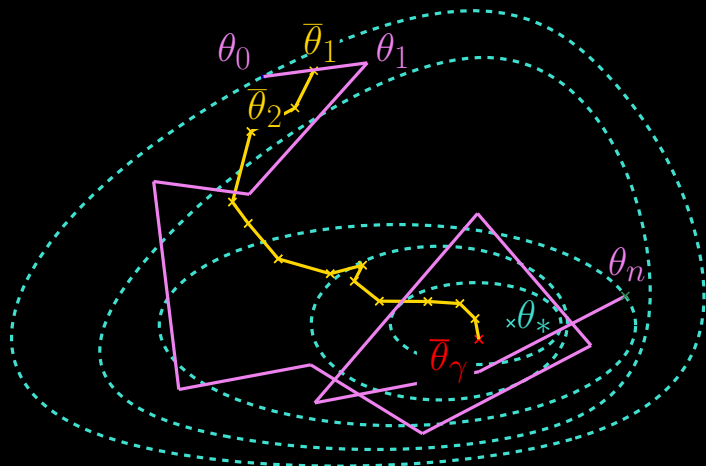
# Constant learning rate SGD: convergence in the non-quadratic case



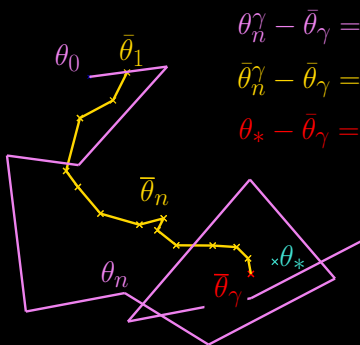
# Constant learning rate SGD: convergence in the non-quadratic case



# Constant learning rate SGD: convergence in the non-quadratic case



# Richardson extrapolation



$$\theta_n^\gamma - \bar{\theta}_\gamma = O_p(\gamma^{1/2})$$

$$\bar{\theta}_n^\gamma - \bar{\theta}_\gamma = O_p(n^{-1/2})$$

$$\theta_* - \bar{\theta}_\gamma = O(\gamma)$$

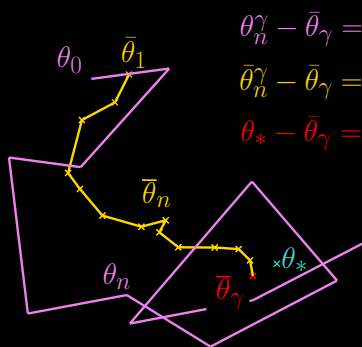
$\bullet \theta_*$

$\bullet \leftarrow \theta_* + \gamma\Delta$

Recovering convergence closer to  $\theta_*$  by **Richardson extrapolation**

$$2\bar{\theta}_{n,\gamma} - \bar{\theta}_{n,2\gamma}$$

# Richardson extrapolation



$$\theta_n^\gamma - \bar{\theta}_\gamma = O_p(\gamma^{1/2})$$

$$\bar{\theta}_n^\gamma - \bar{\theta}_\gamma = O_p(n^{-1/2})$$

$$\theta_* - \bar{\theta}_\gamma = O(\gamma)$$

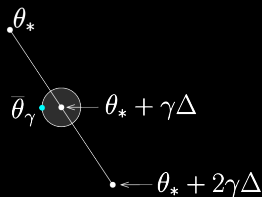
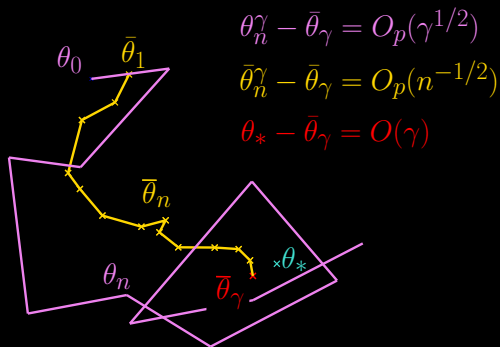
$\theta_*$

$$\bar{\theta}_\gamma \leftarrow \theta_* + \gamma \Delta$$

Recovering convergence closer to  $\theta_*$  by **Richardson extrapolation**

$$2\bar{\theta}_{n,\gamma} - \bar{\theta}_{n,2\gamma}$$

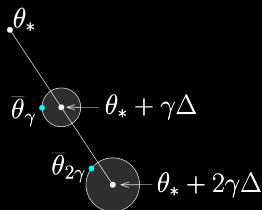
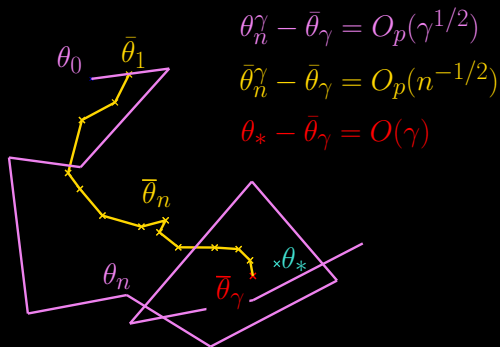
# Richardson extrapolation



Recovering convergence closer to  $\theta_*$  by **Richardson extrapolation**

$$2\bar{\theta}_{n,\gamma} - \bar{\theta}_{n,2\gamma}$$

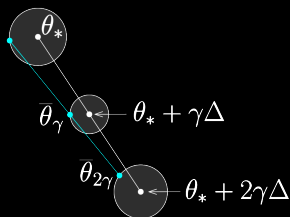
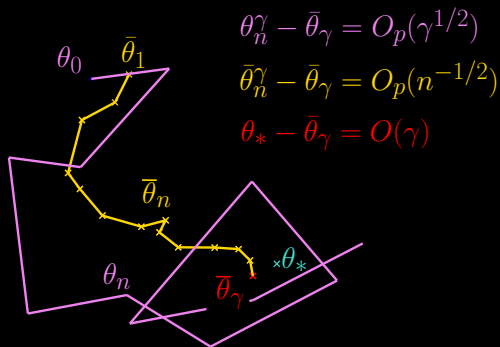
# Richardson extrapolation



Recovering convergence closer to  $\theta_*$  by **Richardson extrapolation**

$$2\bar{\theta}_{n,\gamma} - \bar{\theta}_{n,2\gamma}$$

# Richardson extrapolation

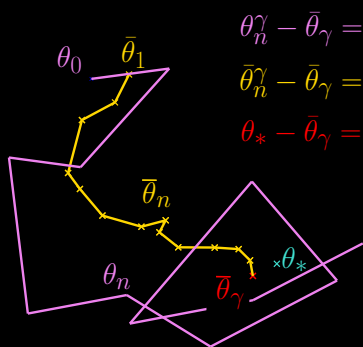


Recovering convergence closer to  $\theta_*$  by **Richardson extrapolation**

$$2\bar{\theta}_{n,\gamma} - \bar{\theta}_{n,2\gamma}$$



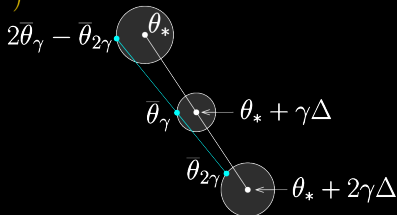
# Richardson extrapolation



$$\theta_n^\gamma - \bar{\theta}_\gamma = O_p(\gamma^{1/2})$$

$$\bar{\theta}_n^\gamma - \bar{\theta}_\gamma = O_p(n^{-1/2})$$

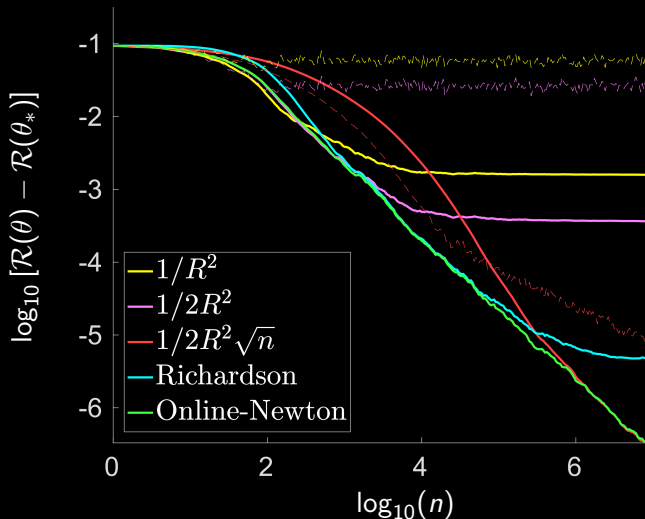
$$\theta_* - \bar{\theta}_\gamma = O(\gamma)$$



Recovering convergence closer to  $\theta_*$  by **Richardson extrapolation**

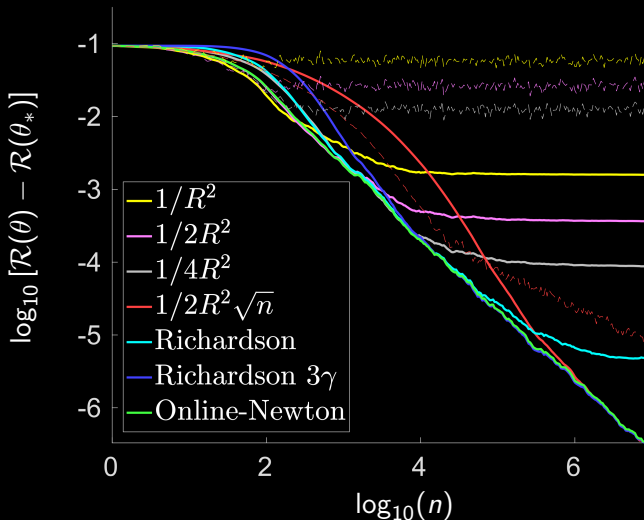
$$2\bar{\theta}_{n,\gamma} - \bar{\theta}_{n,2\gamma}$$

## Experiments: smaller dimension



Synthetic data, logistic regression,  $n = 8.10^6$

# Experiments: Double Richardson



Synthetic data, logistic regression,  $n = 8 \cdot 10^6$

“Richardson  $3\gamma$ ”: estimator built using *Richardson* on 3 different sequences:  $\tilde{\theta}_n^3 = \frac{8}{3}\bar{\theta}_{n,\gamma} - 2\bar{\theta}_{n,2\gamma} + \frac{1}{3}\bar{\theta}_{n,4\gamma}$

# Conclusion MC

Take home message:

- ▶ Precise description of the convergence in terms of Wasserstein distance.
- ▶ Decomposition as three sources of error: variance, initial conditions, and “drift”
- ▶ Detailed analysis of the position of the limit point: the direction does not depend on  $\gamma$  at first order.
- ▶ Extrapolation tricks can help.
- ▶ *Beyond*: new error decomposition (link with diffusions), ...

## Open directions

- ▶ *Markov chain, beyond strong convexity*

## Open directions

- ▶ *Markov chain, beyond strong convexity*
- ▶ *Adaptivity for non-parametric regression*

# Open directions

- ▶ *Markov chain, beyond strong convexity*
- ▶ *Adaptivity for non-parametric regression*
- ▶ *Complexity of non-parametric regression. Stochastic gradient descent and random features.*

# Open directions

- ▶ *Markov chain, beyond strong convexity*
- ▶ *Adaptivity for non-parametric regression*
- ▶ *Complexity of non-parametric regression. Stochastic gradient descent and random features.*
- ▶ *Density estimation.*



[noframenumbering]

- F. Bach and E. Moulines. Non-strongly-convex smooth stochastic approximation with convergence rate  $O(1/n)$ . *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- L. Bottou and O. Bousquet. The tradeoffs of large scale learning. In *Adv. NIPS*, 2008.
- A. Caponnetto and E. De Vito. Optimal rates for the regularized least-squares algorithm. *Foundations of Computational Mathematics*, 7(3):331–368, 2007.
- A. d’Aspremont. Smooth optimization with approximate gradient. *SIAM J. Optim.*, 19(3):1171–1183, 2008.
- A. Dieuleveut and F. Bach. Non-parametric stochastic approximation with large step sizes. *Annals of Statistics*, 2015.
- S. Lacoste-Julien, M. Schmidt, and F. Bach. A simpler approach to obtaining an  $O(1/t)$  rate for the stochastic projected subgradient method. ArXiv e-prints 1212.2002, 2012.
- B. T. Polyak and A. B. Juditsky. Acceleration of stochastic approximation by averaging. *SIAM J. Control Optim.*, 30(4):838–855, 1992.
- H. Robbins and S. Monro. A stochastic approximation method. *The Annals of mathematical Statistics*, 22(3):400–407, 1951.
- D. Ruppert. Efficient estimations from a slowly convergent Robbins-Monro process. Technical report, Cornell University Operations Research and Industrial Engineering, 1988.
- P. Tarrès and Y. Yao. Online learning as stochastic approximation of regularization paths. *IEEE Transactions in Information Theory*, (99):5716–5735, 2011.
- Y. Ying and M. Pontil. Online gradient descent learning algorithms. *Foundations of Computational Mathematics*, 2008.