

# Gradient Descent Method in Learning

*online vs. batch*<sup>a</sup>

Yuan Yao

Department of Mathematics,  
University of California, Berkeley

---

<sup>a</sup>Some joint work with Andrea Caponnetto, Lorenzo Rosasco, Steve Smale, Pierre Tarrès, with help from Yiming Ying and D.-X. Zhou, etc.

# Outline of the talk

---

- Batch vs Online learning

# Outline of the talk

---

- Batch vs Online learning
- Gradient Descent Method in both settings

# Outline of the talk

---

- Batch vs Online learning
- Gradient Descent Method in both settings
- How to do regularizations

# Outline of the talk

---

- Batch vs Online learning
- Gradient Descent Method in both settings
- How to do regularizations
- Lower and Upper Bounds

# Outline of the talk

---

- Batch vs Online learning
- Gradient Descent Method in both settings
- How to do regularizations
- Lower and Upper Bounds
- Exponential Rates for Plug-in Classifiers

# Outline of the talk

---

- Batch vs Online learning
- Gradient Descent Method in both settings
- How to do regularizations
- Lower and Upper Bounds
- Exponential Rates for Plug-in Classifiers
- “Random Connections” with Random Projections

# Outline of the talk

---

- Batch vs Online learning
- Gradient Descent Method in both settings
- How to do regularizations
- Lower and Upper Bounds
- Exponential Rates for Plug-in Classifiers
- “Random Connections” with Random Projections
- Future Directions



# Batch vs. Online

Given a sequence of examples  $(z_i)_{i \in \mathbb{N}} \in (\mathcal{X} \times \mathcal{Y})^\infty$

- Batch Learning: truncation set  $\mathbf{z}_T = (z_i)_{i=1}^T$ , find a mapping

$$\mathbf{z}_T \mapsto f_{\mathbf{z}_T}$$

- Online Learning: a Markov Decision Process

$$f_{t+1} = T_t(f_t, z_{t+1})$$

where  $f_t$  only depends on  $z_1, \dots, z_t$ .

# Why Online?

---

- Low computational cost:
  - online needs  $\geq O(t)$  steps
  - batch typically needs  $\geq O(T^3)$  (inverting a matrix)
- Fast convergence: order optimality
- Temporal dependence of samples:
  - Markov Chain sampling*: large-scale networks
  - Mixing processes*: exponential-mixing and polynomial-mixing
  - Games*: competitive (non-statistical) analysis

# Where we start: Penalizations

---

$$\min_{f \in \mathcal{H}} \frac{1}{T} \sum_{i=1}^T V(y_i, f(x_i)) + \lambda \|f\|_{\mathcal{H}}^2$$

where we choose the loss  $V(y, f(x))$ :

- $L_2$  loss: for order optimality analysis
- $L_1$  loss (soft margin): for sparsity, e.g. Basis Pursuit and SVM regression

# Where we start: Penalizations

$$\min_{f \in \mathcal{H}} \frac{1}{T} \sum_{i=1}^T V(y_i, f(x_i)) + \lambda \|f\|_{\mathcal{H}}^2$$

where we choose the loss  $V(y, f(x))$ :

- $L_2$  loss: for order optimality analysis
- $L_1$  loss (soft margin): for sparsity, e.g. Basis Pursuit and SVM regression

and  $\mathcal{H} = \mathcal{H}_K$  a reproducing kernel Hilbert space (RKHS).

# Why RKHS

---

- $\mathcal{H}_K$  can be dense in continuous function space  $\mathcal{C}(X)$

# Why RKHS

---

- $\mathcal{H}_K$  can be dense in continuous function space  $\mathcal{C}(X)$
- the gradient takes a simple form in  $\mathcal{H}_K$ :

$$\text{grad}_f V = V'_f(y, f(x)) K_x$$

# Why RKHS

---

- $\mathcal{H}_K$  can be dense in continuous function space  $\mathcal{C}(X)$
- the gradient takes a simple form in  $\mathcal{H}_K$ :

$$\text{grad}_f V = V'_f(y, f(x)) K_x$$

- super-fast convergence rate for plug-in classifiers...

# Exponential Rates in Classifications

Assume that

- The regression function  $f_\rho \in \mathcal{H}_K$
- $f_\rho$  has margin  $\gamma > 0$ :

$$\mathbf{Prob}\{x \in X : |f_\rho(x)| \leq \gamma\} = 0$$

Then there is a  $f_t : (z_i)_1^t \rightarrow \mathcal{H}_K$ , s.t.

$$\mathbb{E}R(f_t) - R(f_\rho) \leq O(\exp(-c\gamma t^{1/4}))$$

where  $R(f)$  is the misclassification error of the *plug-in classifier*  $\text{sign}(f)$ .



# What is RKHS

- $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is a *Mercer* kernel, i.e. a *continuous*, symmetric and positive definite function
- $\mathcal{H}_K = \overline{\text{span}\{K_x : x \in \mathcal{X}\}}$  where the closure is w.r.t. the inner product  $\langle K_x, K_{x'} \rangle_K = K(x, x')$
- *Reproducing* property:  $f(x) = \langle f, K_x \rangle_K$
- $\mathcal{H}_K$  is a subspace (closed iff finite dimension) in  $\mathcal{L}^2_{\rho_X} \cap \mathcal{C}(\mathcal{X})$

# Gradient Descent Algorithms

For  $L_2$  loss and  $\mathcal{H} = \mathcal{H}_K$ ,

- Batch ( $L_2$ Boost):

$$\hat{f}_{t+1} = \hat{f}_t - \eta_t \left[ \frac{1}{T} \sum_{i=1}^T (\hat{f}_t(x_i) - y_i) K_{x_i} + \lambda_T \hat{f}_t \right]$$

- Online:

$$f_{t+1} = f_t - \eta_t [(f_t(x_{t+1}) - y_{t+1}) K_{x_{t+1}} + \lambda_t f_t]$$

# Our Theoretical Goal

---

Convergence of  $(\hat{f}_t) \in \mathcal{H}_K$  and  $(f_t) \in \mathcal{H}_K$  to the regression function

$$f_\rho(x) := \mathbb{E}[y|x] \in \mathcal{L}_{\rho_X}^2$$

and its rates when  $f_\rho$  takes some sparse form.

# Our Theoretical Goal

---

Convergence of  $(\hat{f}_t) \in \mathcal{H}_K$  and  $(f_t) \in \mathcal{H}_K$  to the regression function

$$f_\rho(x) := \mathbb{E}[y|x] \in \mathcal{L}_{\rho_X}^2$$

and its rates when  $f_\rho$  takes some sparse form.

But,  $\mathcal{L}_{\rho_X}^2$  is too large a space to search, so we need *regularizations*.

# Regularization

Two parameters:  $\lambda_t$  (or  $\lambda_T$ ) and  $\eta_t$ :

- $\lambda_T = 0$  and  $\eta_t = c$ : Landwebter iterations/ $L_2$ Boost
- $\lambda_T = 0$  and  $\eta_t \downarrow 0$ : Yao et al. (2005)
- $\lambda_t = \lambda > 0$  and  $\eta_t \downarrow 0$ :  $f_t \rightarrow f_\lambda \neq f_\rho$ , Smale and Yao (2005) etc.
- $\lambda_t \downarrow 0$  and  $\eta_t \downarrow 0$ :  $f_t \rightarrow f_\rho$ , Yao and Tarrès (2005)
- $\lambda_t = 0$  and  $\eta_t \downarrow 0$ :  $f_t \rightarrow f_\rho$ , Ying et al. (2006)

# Sparsity of Regression Function

---

We are going to assume that the regression function is sparse/smooth w.r.t. the following *basis*

- roughly speaking, kernel principle components,
- or more precisely, the eigenfunctions of the *covariance operator* of  $\rho_{\mathcal{X}}$  on  $\mathcal{H}_K$ .

# Covariance operator

---

- Define an integral operator

$$\begin{aligned} L_K : \mathcal{L}_{\rho_X}^2 &\rightarrow \mathcal{H}_K \\ f &\mapsto \int_X f(x') K(x', \cdot) d\rho_X \end{aligned}$$

# Covariance operator

- Define an integral operator

$$\begin{aligned} L_K : \mathcal{L}_{\rho_X}^2 &\rightarrow \mathcal{H}_K \\ f &\mapsto \int_X f(x') K(x', \cdot) d\rho_X \end{aligned}$$

- The *covariance operator*, is the restriction  $L_K|_{\mathcal{H}_K} : \mathcal{H}_K \rightarrow \mathcal{H}_K$ , i.e.  $\mathbb{E}_x[\langle \cdot, K_x \rangle K_x]$



# Covariance operator

- Define an integral operator

$$\begin{aligned} L_K : \mathcal{L}_{\rho_X}^2 &\rightarrow \mathcal{H}_K \\ f &\mapsto \int_X f(x') K(x', \cdot) d\rho_X \end{aligned}$$

- The *covariance operator*, is the restriction  $L_K|_{\mathcal{H}_K} : \mathcal{H}_K \rightarrow \mathcal{H}_K$ , i.e.  $\mathbb{E}_x[\langle \cdot, K_x \rangle K_x]$
- $L_K : \mathcal{L}_{\rho_X}^2 \rightarrow \mathcal{L}_{\rho_X}^2$  compact  $\Rightarrow$  orthonormal eigen-system  $(\lambda_i, \phi_i)_{i \in \mathbb{N}}$ ,  $\phi_i \in \mathcal{L}_{\rho_X}^2 \cap \mathcal{H}_K$  bi-orthogonal and

$$\sum \lambda_i \leq \sup_{x \in \mathcal{X}} K(x, x) =: \kappa < \infty$$

# Sparsity Assumption

Assume that

$$f_\rho = L_K^r g, \quad g \in \mathcal{L}_{\rho_X}^2, r > 0$$

i.e.  $f_\rho$  has at least *power-law decay* coordinates w.r.t. the basis of eigenfunctions of  $L_K : \mathcal{L}_{\rho_X}^2 \rightarrow \mathcal{L}_{\rho_X}^2$ :

$$f_\rho = \sum_i \lambda_i^r g_i \phi_i, \quad \sum \lambda_i < \infty, \sum g_i^2 < \infty$$

# Sparsity Assumption

Assume that

$$f_\rho = L_K^r g, \quad g \in \mathcal{L}_{\rho_X}^2, r > 0$$

i.e.  $f_\rho$  has at least *power-law decay* coordinates w.r.t. the basis of eigenfunctions of  $L_K : \mathcal{L}_{\rho_X}^2 \rightarrow \mathcal{L}_{\rho_X}^2$ :

$$f_\rho = \sum_i \lambda_i^r g_i \phi_i, \quad \sum \lambda_i < \infty, \sum g_i^2 < \infty$$

Note: if  $K$  is a stochastic density kernel,  $L_K^r$  is used to construct diffusion wavelets by Coifman et al.

# Lower Rates in Learning

Let  $\mathbb{P}(b, r)$  ( $b > 1$  and  $r \in (1/2, 1]$ ) be the set of probability measure  $\rho$  on  $\mathcal{X} \times \mathcal{Y}$ , such that:

- the eigenvalues  $\lambda_i$ , arranged in a nonincreasing order, decay at  $O(i^{-b})$
- $f_\rho = L_K^r g$  for some  $g \in \mathcal{L}_{\rho_X}^2$
- almost surely  $|y| \leq M_\rho$

# Minimax Lower Rates

[Caponnetto-DeVito'05] The minimax lower rate:

$$\liminf_{t \rightarrow \infty} \inf_{\mathbf{z}_t \mapsto f_t} \sup_{\rho \in \mathbb{P}(b,r)} \mathbf{Prob} \left\{ \mathbf{z}_t \in \mathbb{Z}^t : \frac{\|f_t - f_\rho\|_2}{t^{-\frac{rb}{2rb+1}}} > C \right\} = 1$$

where the inf is taken over all algorithms mapping  $(z_i)_1^t \mapsto f_t$ .

# Minimax Lower Rates

[Caponnetto-DeVito'05] The minimax lower rate:

$$\liminf_{t \rightarrow \infty} \inf_{\mathbf{z}_t \mapsto f_t} \sup_{\rho \in \mathbb{P}(b,r)} \mathbf{Prob} \left\{ \mathbf{z}_t \in \mathbb{Z}^t : \frac{\|f_t - f_\rho\|_2}{t^{-\frac{rb}{2rb+1}}} > C \right\} = 1$$

where the inf is taken over all algorithms mapping  $(z_i)_1^t \mapsto f_t$ .

- The  $\rho$  in  $\sup_{\rho \in \mathbb{P}(b,r)}$  depends on sample size  $t$ !

# Minimax Lower Rates

[Caponnetto-DeVito'05] The minimax lower rate:

$$\liminf_{t \rightarrow \infty} \inf_{\mathbf{z}_t \mapsto f_t} \sup_{\rho \in \mathbb{P}(b,r)} \mathbf{Prob} \left\{ \mathbf{z}_t \in \mathbb{Z}^t : \frac{\|f_t - f_\rho\|_2}{t^{-\frac{rb}{2rb+1}}} > C \right\} = 1$$

where the inf is taken over all algorithms mapping  $(z_i)_1^t \mapsto f_t$ .

- The  $\rho$  in  $\sup_{\rho \in \mathbb{P}(b,r)}$  depends on sample size  $t$ !
- Suitable for online learning, instead of batch learning.

# Individual Lower Rates

[Caponnetto-DeVito'05] The individual lower rate: for each  $B > b$ ,

$$\inf_{((z_i)_1^t \mapsto f_t)_{t \in \mathbb{N}}} \sup_{\rho \in \mathbb{P}(b, r)} \limsup_{t \rightarrow \infty} \frac{\|f_t - f_\rho\|_2}{t^{-\frac{rB}{2rB+1}}} > 0.$$



# Individual Lower Rates

[Caponnetto-DeVito'05] The individual lower rate: for each  $B > b$ ,

$$\inf_{((z_i)_1^t \mapsto f_t)_{t \in \mathbb{N}}} \sup_{\rho \in \mathbb{P}(b, r)} \limsup_{t \rightarrow \infty} \frac{\|f_t - f_\rho\|_2}{t^{-\frac{rB}{2rB+1}}} > 0.$$

Note: taking  $b = 1$  and  $B = 1$ , it suggests *eigenvalue independent* minimax and individual lower rates:

$$t^{-\frac{r}{2r+1}}$$

# Upper Bounds for Batch Learning

Theorem (Yao-Rosasco-Caponnetto'05). Assume that  $f_\rho = L_K^r g$  ( $r > 0$ ). There exist  $\lambda_T, \eta_t$  and an early stopping rule  $t^* : \mathbb{N} \rightarrow \mathbb{N}$ , such that

- if  $r > 0$ ,  $\|\hat{f}_{t^*(T)} - f_\rho\|_2 \leq O(T^{-\frac{r}{2r+2}})$
- if  $r > 1/2$ ,  $\|\hat{f}_{t^*(T)} - f_\rho\|_K \leq O(T^{-\frac{r-1/2}{2r+2}})$

In fact, one may choose  $\lambda_T = 0$ ,  $\eta_t = \frac{1}{\kappa^2(t+1)^\theta}$  and  $t^*(T) = \lceil T^{-\frac{1}{(2r+2)(1-\theta)}} \rceil$ .

# Improvements

---

[Bauer-Pereverzev-Rosasco'06] For  $\theta = 0$  and  $r > 1/2$ ,

$$\|\hat{f}_{t^*(T)} - f_\rho\|_2 \leq O(T^{-\frac{r}{2r+1}})$$

which meets the lower rates.

# Upper Bounds for Online Learning

Theorem (Tarrès-Yao'06). Assume that  $f_\rho = L_K^r g$  ( $r > 0$ ). There exist  $\lambda_t$  and  $\eta_t$  such that

- if  $r > 0$ ,  $\|f_t - f_\rho\|_2 \leq O(t^{-\max\{\frac{r}{2r+1}, 1/3\}})$
- if  $r > 1/2$ ,  $\|f_t - f_\rho\|_K \leq O(t^{-\max\{\frac{r-1/2}{2r+1}, 1/4\}})$

In fact,  $\lambda_t \sim O(t^{-1/(2r+1)})$  and  $\eta_t \sim O(t^{-2r/(2r+1)})$ .

# Upper Bounds for Online Learning

Theorem (Tarrès-Yao'06). Assume that  $f_\rho = L_K^r g$  ( $r > 0$ ). There exist  $\lambda_t$  and  $\eta_t$  such that

- if  $r > 0$ ,  $\|f_t - f_\rho\|_2 \leq O(t^{-\max\{\frac{r}{2r+1}, 1/3\}})$
- if  $r > 1/2$ ,  $\|f_t - f_\rho\|_K \leq O(t^{-\max\{\frac{r-1/2}{2r+1}, 1/4\}})$

In fact,  $\lambda_t \sim O(t^{-1/(2r+1)})$  and  $\eta_t \sim O(t^{-2r/(2r+1)})$ .

*Note:* the upper rates *saturate* when  $r \geq 1$  and  $r \geq 3/2$ !

# Breaking Saturation

---

It is expected that with  $\lambda_t = 0$  and suitable choices  $\eta_t \rightarrow 0$  and  $\sum_t \eta_t = \infty$ , one has

$$\|f_t - f_\rho\|_2 \leq O(t^{-\frac{r}{2r+1}})$$

for *all*  $r > 0$ .

# Breaking Saturation

It is expected that with  $\lambda_t = 0$  and suitable choices  $\eta_t \rightarrow 0$  and  $\sum_t \eta_t = \infty$ , one has

$$\|f_t - f_\rho\|_2 \leq O(t^{-\frac{r}{2r+1}})$$

for all  $r > 0$ .

*A Positive Answer:* Ying et al. (2006) give results suggesting its truth.

# Plug-in Classifiers: Exponential Rates

Taking  $r = 1/2$  or equivalently assuming that  $f_\rho \in \mathcal{H}_K$ , combining a recent result by Audibert & Tsybakov'05,

$$\mathbb{E}R(f) - R(f_\rho) \leq \rho_X(x : |f(x) - f_\rho(x)| \geq \gamma)$$

we obtain that for online learning

$$\mathbb{E}R(f_t) - R(f_\rho) \leq O(\exp(-c\gamma t^{1/4}))$$

where  $f_\rho$  has margin  $\gamma$ . Similar holds for batch learning.



# Random Projection Perspective

---

Given  $\mathbf{x}_T \in \mathcal{X}^T$ , define a sampling operator on  $\mathcal{H}_K$

$$\begin{aligned} S_{\mathbf{x}_T} : \mathcal{H}_K &\rightarrow l_2(\mathbf{x}_T) \\ f &\mapsto (f(x_i))_1^T = (\langle f, K_{x_i} \rangle_K)_1^T \end{aligned}$$

# Random Projection Perspective

Given  $\mathbf{x}_T \in \mathcal{X}^T$ , define a sampling operator on  $\mathcal{H}_K$

$$\begin{aligned} S_{\mathbf{x}_T} : \mathcal{H}_K &\rightarrow l_2(\mathbf{x}_T) \\ f &\mapsto (f(x_i))_1^T = (\langle f, K_{x_i} \rangle_K)_1^T \end{aligned}$$

- $S_{\mathbf{x}_T} f$  takes  $T$  random measurements/projections of  $f$ .

# Random Projection Perspective

Given  $\mathbf{x}_T \in \mathcal{X}^T$ , define a sampling operator on  $\mathcal{H}_K$

$$\begin{aligned} S_{\mathbf{x}_T} : \mathcal{H}_K &\rightarrow l_2(\mathbf{x}_T) \\ f &\mapsto (f(x_i))_1^T = (\langle f, K_{x_i} \rangle_K)_1^T \end{aligned}$$

- $S_{\mathbf{x}_T} f$  takes  $T$  random measurements/projections of  $f$ .
- Adjoint operator  $S_{\mathbf{x}_T}^* \mathbf{y} = \frac{1}{T} \sum_{i=1}^T y_i K_{x_i}$ .

# Random Projection Perspective

Given  $\mathbf{x}_T \in \mathcal{X}^T$ , define a sampling operator on  $\mathcal{H}_K$

$$\begin{aligned} S_{\mathbf{x}_T} : \mathcal{H}_K &\rightarrow l_2(\mathbf{x}_T) \\ f &\mapsto (f(x_i))_1^T = (\langle f, K_{x_i} \rangle_K)_1^T \end{aligned}$$

- $S_{\mathbf{x}_T} f$  takes  $T$  random measurements/projections of  $f$ .
- Adjoint operator  $S_{\mathbf{x}_T}^* \mathbf{y} = \frac{1}{T} \sum_{i=1}^T y_i K_{x_i}$ .
- $S_{\mathbf{x}_T} S_{\mathbf{x}_T}^*$  is the Gram matrix  $(K(x_i, x_j))^{T \times T}$ .

# Compressed Sensing

---

- $f$  is sparse w.r.t. certain basis/frames (unknown)
- $S_{\mathbf{x}_T}$  takes some random measurements of  $f$  such that the Uniform Uncertainty Principle holds, or equivalently, for small enough  $T_0$  and all  $T \leq T_0$ ,  $S_{\mathbf{x}_T} S_{\mathbf{x}_T}^*$  has a *uniform lower bound* (depending on the sparsity of  $f$ ) on the smallest eigenvalue.

# Compressed Sensing

- $f$  is sparse w.r.t. certain basis/frames (unknown)
- $S_{\mathbf{x}_T}$  takes some random measurements of  $f$  such that the Uniform Uncertainty Principle holds, or equivalently, for small enough  $T_0$  and all  $T \leq T_0$ ,  $S_{\mathbf{x}_T} S_{\mathbf{x}_T}^*$  has a *uniform lower bound* (depending on the sparsity of  $f$ ) on the smallest eigenvalue.

However in Learning, since

$$\mathbb{E}[S_{\mathbf{x}_T}^* S_{\mathbf{x}_T}] = L_K|_{\mathcal{H}_K}$$

where  $L_K$  is a compact operator with eigenvalues convergent to 0, NO lower bound!

# Learning vs. Compressed Sensing

---

To control the *condition number* (or smallest eigenvalue) of the Gram matrix  $S_{\mathbf{x}_T} S_{\mathbf{x}_T}^*$ :

- Learning uses *regularization*
- Scott & Nowak'05 uses Vapnik's *structural risk minimization*
- Donoho, Candès & Tao, use *Random Matrix Theory*

# Learning vs. Compressed Sensing

---

To control the *condition number* (or smallest eigenvalue) of the Gram matrix  $S_{\mathbf{x}_T} S_{\mathbf{x}_T}^*$ :

- Learning uses *regularization*
- Scott & Nowak'05 uses Vapnik's *structural risk minimization*
- Donoho, Candès & Tao, use *Random Matrix Theory*

Moreover, there is another kind of “condition number” in machine learning:

*margin*



# Margin (normalized)

---

Definitions.

- $f \in \mathcal{H}_K$  has *margin*  $\gamma > 0$ , if

$$\rho_X\{x \in X : |f(x)| \geq \gamma \|f\| \|K_x\|\} = 1$$

- $f \in \mathcal{H}_K$  has *margin*  $\gamma > 0$  with error  $\epsilon \in [0, 1]$ , if

$$\rho_X\{x \in X : |f(x)| \geq \gamma \|f\| \|K_x\|\} \geq 1 - \epsilon$$

# Margin and Random Projections

[Balcan-Blum-Vempala'05] If  $f \in \mathcal{H}_K$  has margin  $\gamma$ , then with i.i.d. examples of number

$$t \geq \frac{8}{\epsilon} \max \left\{ \frac{1}{\gamma^2}, \ln \frac{1}{\delta} \right\}$$

there is a  $f_t$  such that with confidence  $1 - \delta$ ,  $f_t$  has margin  $\gamma/2$  with error  $\epsilon$ .

# Margin and Random Projections

[Balcan-Blum-Vempala'05] If  $f \in \mathcal{H}_K$  has margin  $\gamma$ , then with i.i.d. examples of number

$$t \geq \frac{8}{\epsilon} \max \left\{ \frac{1}{\gamma^2}, \ln \frac{1}{\delta} \right\}$$

there is a  $f_t$  such that with confidence  $1 - \delta$ ,  $f_t$  has margin  $\gamma/2$  with error  $\epsilon$ .

- In fact,  $f_t$  can be realized by the Orthogonal Projections on to  $\text{span}\{K_i : 1 \leq i \leq t\}$ .

# Future Directions

---

- Step-Size Adaptation
  - Cross-Validation
  - Averaging process acceleration
  - Stochastic Meta-Descent (SMD)
- Dependent Sampling
  - Markov Chain sampling
  - Mixing process
- Various aspects of Random Projections