

Jinchao Xu

Deep Learning Algorithms and Analysis

Summer 2020

Contents

1	Monte Carlo Methods	5
1.1	Monte Carlo methods	5
1.1.1	A basic result	5
1.1.2	Application	8
1.2	Integral representations of functions	10
1.2.1	Fourier representation	10
1.2.2	Double Fourier representation	11
	References	13

Monte Carlo Methods

1.1 Monte Carlo methods

Let $\lambda \geq 0$ be a probability density function on $G \subset \mathbb{R}^d$ such that

$$(1.1) \quad \int_G \lambda(\omega) d\omega = 1.$$

For example:

$$(1.2) \quad \lambda(\omega) = \frac{1}{|G|},$$

if G is bounded. The expectation is defined:

$$(1.3) \quad \mathbb{E}g := \int_G g(\omega) \lambda(\omega) d\omega$$

and for any $h = h(\omega_1, \omega_2, \dots, \omega_n) : G \times G \cdots G \mapsto \mathbb{R}$

$$(1.4) \quad \mathbb{E}_n h := \int_{G \times G \times \dots \times G} h(\omega_1, \omega_2, \dots, \omega_n) \lambda(\omega_1) \lambda(\omega_2) \dots \lambda(\omega_n) d\omega_1 d\omega_2 \dots d\omega_n.$$

1.1.1 A basic result

Lemma 1. *For any $g \in L^\infty(G)$, we have*

$$(1.5) \quad \mathbb{E}_n \left(\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i) \right)^2 = \begin{cases} \frac{1}{n} \mathbb{E}((\mathbb{E}g - g)^2) \leq \frac{1}{n} \sup_{\omega, \omega' \in G} |g(\omega) - g(\omega')|^2 \\ \frac{1}{n} (\mathbb{E}(g^2) - (\mathbb{E}(g))^2) \leq \frac{1}{n} \mathbb{E}(g^2) \leq \frac{1}{n} \|g\|_{L^\infty}^2, \end{cases}$$

Proof. First note that

$$\begin{aligned}
(1.6) \quad \left(\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i) \right)^2 &= \frac{1}{n^2} \left(n\mathbb{E}g - \sum_{i=1}^n g(\omega_i) \right)^2 = \frac{1}{n^2} \left(\sum_{i=1}^n (\mathbb{E}g - g(\omega_i)) \right)^2 \\
&= \frac{1}{n^2} \sum_{i,j=1}^n (\mathbb{E}g - g(\omega_i))(\mathbb{E}g - g(\omega_j)) \\
&= \frac{I_1}{n^2} + \frac{I_2}{n^2}.
\end{aligned}$$

with

$$(1.7) \quad I_1 = \sum_{i=1}^n (\mathbb{E}g - g(\omega_i))^2, \quad I_2 = \sum_{i \neq j}^n ((\mathbb{E}g)^2 - \mathbb{E}(g)(g(\omega_i) + g(\omega_j)) + g(\omega_i)g(\omega_j)).$$

Consider I_1 , for any i ,

$$\mathbb{E}_n(\mathbb{E}g - g(\omega_i))^2 = \mathbb{E}_n(\mathbb{E}g - g)^2.$$

Thus,

$$\mathbb{E}_n I_1 = n\mathbb{E}((\mathbb{E}g - g)^2).$$

For I_2 , note that

$$\mathbb{E}_n g(\omega_i) = \mathbb{E}_n g(\omega_j) = \mathbb{E}(g)$$

and, for $i \neq j$,

$$\begin{aligned}
(1.8) \quad \mathbb{E}_n(g(\omega_i)g(\omega_j)) &= \int_{G \times G \times \dots \times G} g(\omega_j)g(\omega_i)\lambda(\omega_1)\lambda(\omega_2)\dots\lambda(\omega_n)d\omega_1 d\omega_2 \dots d\omega_n \\
&= \int_{G \times G} g(\omega_j)g(\omega_i)\lambda(\omega_1)\lambda(\omega_1)\lambda(\omega_2)d\omega_1 d\omega_2 \\
&= \mathbb{E}_n(g(\omega_i))\mathbb{E}_n(g(\omega_j)) = [\mathbb{E}(g)]^2.
\end{aligned}$$

Thus

$$(1.9) \quad \mathbb{E}_n(I_2) = \mathbb{E}_n \left(\sum_{i \neq j}^n ((\mathbb{E}g)^2 - \mathbb{E}(g)(\mathbb{E}(g(\omega_i)) + \mathbb{E}(g(\omega_j))) + \mathbb{E}(g(\omega_i)g(\omega_j))) \right) = 0.$$

Consequently, there exist the following two formulas for $\mathbb{E}_n \left(\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i) \right)^2$:

$$(1.10) \quad \mathbb{E}_n \left(\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i) \right)^2 = \frac{1}{n^2} \mathbb{E}_n I_1 = \begin{cases} \frac{1}{n} \mathbb{E}((\mathbb{E}g - g)^2) \\ \frac{1}{n} (\mathbb{E}(g^2) - (\mathbb{E}g)^2). \end{cases}$$

Based on the first formula above, since

$$|g(\omega) - \mathbb{E}g| = \left| \int_G (g(\omega) - g(\tilde{\omega}))\lambda(\tilde{\omega})d\tilde{\omega} \right| \leq \sup_{\omega, \omega' \in G} |g(\omega) - g(\omega')|,$$

it holds that

$$(1.11) \quad \mathbb{E}_n \left(\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i) \right)^2 \leq \frac{1}{n} \sup_{\omega, \omega' \in G} |g(\omega) - g(\omega')|^2.$$

Due to the second formula above,

$$(1.12) \quad \mathbb{E}_n \left(\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i) \right)^2 \leq \frac{1}{n} \mathbb{E}(g^2) \leq \frac{1}{n} \|g\|_{L^\infty}^2$$

which completes the proof. \square

Of course, we can use this to prove a high probability result.

Corollary 1. *Under the assumptions of the preceding lemma, we have*

$$(1.13) \quad \mathbb{P}[(\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i))^2 > \frac{k}{n} \|g\|_{L^\infty}^2] < \frac{1}{k}$$

Proof.

$$(1.14) \quad \mathbb{P}[(\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i))^2 > \epsilon] \leq \epsilon^{-1} \mathbb{E}[(\mathbb{E}g - \frac{1}{n} \sum_{i=1}^n g(\omega_i))^2] \leq \frac{1}{n\epsilon} \|g\|_{L^\infty}^2.$$

\square

This corollary implies that the set of ω_i where the estimate $n^{-1} \sum_{i=1}^n g(\omega_i)$ is far from the desired value $\mathbb{E}g$ is small.

The practical usefulness of this algorithm depends upon the existence of a *repeatable* process (for instance some physical process) which *generates* ω according to a desired distribution μ .

The precise meaning of this last statement is essentially that the strong law of large numbers holds. Specifically, if $\omega_1, \dots, \omega_n, \dots$ is a infinite sequence generated by the process, and $A \subset \Omega$ is any a measurable set, then

$$(1.15) \quad \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \chi_A(\omega_i) = \mu(A).$$

Generating n independent samples means generating $\omega_1, \dots, \omega_n$ from μ^n according to the above notion. The existence of a realizable process generating samples from a probability distribution, and the practical use of such processes is an interesting topic in the intersection of statistics, physics, and computer science. In addition, statistics/probability theory studies how to take samples from one probability distribution and transform them to samples from another distribution.

1.1.2 Application**Lemma 2.** *Let*

$$(1.16) \quad f(x) = \int_G g(x, \theta) \lambda(\theta) d\theta = \mathbb{E}(g)$$

with $\lambda(\theta) \geq 0$ and $\|\lambda(\theta)\|_{L^1(G)} = 1$. For any $n \geq 1$, there exist $\theta_i^* \in G$ such that

$$\|f - f_n\|_{L^2(\Omega)}^2 \leq \frac{1}{n} \int_G \|g(\cdot, \theta)\|_{L^2(\Omega)}^2 \lambda(\theta) d\theta = \frac{1}{n} \mathbb{E}(\|g(\cdot, \theta)\|_{L^2(\Omega)}^2)$$

where $\|g(\cdot, \theta)\|_{L^2(\Omega)}^2 = \int_\Omega [g(x, \theta)]^2 d\mu(x)$, and

$$f_n(x) = \frac{1}{n} \sum_{i=1}^n g(x, \theta_i^*).$$

Proof. Introducing a probability distribution $\lambda(\theta)$,

$$(1.17) \quad f(x) = \mathbb{E}(g).$$

By Lemma 1,

$$\mathbb{E}_n \left(\left(\mathbb{E}(g(x, \cdot)) - \frac{1}{n} \sum_{i=1}^n g(x, \theta_i) \right)^2 \right) \leq \frac{1}{n} \mathbb{E}(g^2)$$

and

$$\mathbb{E}_n(h(\theta_1, \theta_2, \dots, \theta_n)) \leq \frac{1}{n} \mathbb{E} \left(\int_\Omega g^2 d\mu(x) \right),$$

by taking integral where

$$h(\theta_1, \theta_2, \dots, \theta_n) = \int_\Omega \left(\mathbb{E}(g(x, \cdot)) - \frac{1}{n} \sum_{i=1}^n g(x, \theta_i) \right)^2 d\mu(x).$$

Sine $\mathbb{E}_n(1) = 1$ and $\mathbb{E}_n(h) \leq \frac{1}{n} \mathbb{E} \left(\int_\Omega g^2 d\mu(x) \right)$, there exists $(\theta_1^*, \theta_2^*, \dots, \theta_n^*) \in G \times G \times \dots \times G$ such that

$$h(\theta_1^*, \theta_2^*, \dots, \theta_n^*) \leq \frac{1}{n} \int_\Omega \mathbb{E}(g^2) d\mu(x).$$

Otherwise, $\mathbb{E}_n(h) > \frac{1}{n} \mathbb{E} \left(\int_\Omega g^2 d\mu(x) \right)$ if $h(\theta_1, \theta_2, \dots, \theta_n) > \frac{1}{n} \int_\Omega \mathbb{E}(g^2) d\mu(x)$. This implies that

$$\|f - f_n\|_{L^2(\Omega)}^2 \leq \frac{1}{n} \int_G \|g(\cdot, \theta)\|_{L^2(\Omega)}^2 \lambda(\theta) d\theta,$$

which completes the proof. \square

We also have a more general version of the above lemma.

Lemma 3. *Let*

$$(1.18) \quad f(x) = \int_G g(x, \theta) \lambda(\theta) d\theta = \mathbb{E}(g)$$

with $\|\lambda(\theta)\|_{L^1(\Theta)} = 1$. For any $n \geq 1$, there exist $\theta_i^* \in G$ such that

$$\|f - f_n\|_{H^m(\Omega)}^2 \leq \int_G \|g(\cdot, \theta)\|_{H^m(\Omega)}^2 \lambda(\theta) d\theta = \frac{1}{n} \mathbb{E}(\|g(\cdot, \theta)\|_{H^m(\Omega)}^2)$$

where

$$f_n(x) = \frac{1}{n} \sum_{i=1}^n g(x, \theta_i^*)$$

In particular, if

$$(1.19) \quad |D^\alpha g(x, \theta)| \leq C, \quad \forall x, \theta, |\alpha| \leq m$$

Then

$$\|f - f_n\|_{H^m(\Omega)} \leq \left(\frac{m+d}{m} \right)^{1/2} |\Omega|^{1/2} n^{-1/2}.$$

For any $f(x) = \int_G g(x, \theta) \rho(\theta) d\theta$ with $\|\rho\|_{L^1(\Theta)} \neq 1$. Let $\lambda(\theta) = \frac{\rho(\theta)}{\|\rho\|_{L^1(\Theta)}}$. Thus,

$$(1.20) \quad f(x) = \|\rho\|_{L^1(\Theta)} \int_G g(x, \theta) \lambda(\theta) d\theta$$

with $\|\lambda(\theta)\|_{L^1(\Theta)} = 1$. We can apply the above two lemmas to the given function $f(x)$.

1.2 Integral representations of functions

1.2.1 Fourier representation

Consider the Fourier transform:

$$(1.21) \quad \hat{f}(\omega) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} e^{-i\omega \cdot x} f(x) dx \quad \forall \omega \in \mathbb{R}^d.$$

We write $\hat{f}(\omega) = e^{i\theta(\omega)} |\hat{f}(\omega)|$. By Fourier inversion formula,

$$(1.22) \quad f(x) = \int_{\mathbb{R}^d} e^{i\omega \cdot x} \hat{f}(\omega) d\omega = \int_{\mathbb{R}^d} e^{i(\omega \cdot x + \beta(\omega))} |\hat{f}(\omega)| d\omega.$$

Since $f(x)$ is real-valued, it implies that, for x

$$(1.23) \quad \begin{aligned} f(x) &= \operatorname{Re} \int_{\mathbb{R}^d} e^{i\omega \cdot x} \hat{f}(\omega) d\omega \\ &= \operatorname{Re} \int_{\mathbb{R}^d} e^{i\omega \cdot x} e^{i\beta(\omega)} |\hat{f}(\omega)| d\omega \\ &= \int_{\mathbb{R}^d} \cos(\omega \cdot x + \beta(\omega)) |\hat{f}(\omega)| d\omega. \end{aligned}$$

Then we have

$$(1.24) \quad f(x) = \int_{\mathbb{R}^d} k(x, \omega) d\omega,$$

with

$$(1.25) \quad k(x, \omega) = \cos(\omega \cdot x + \beta(\omega)) |\hat{f}(\omega)|$$

and

$$(1.26) \quad |k(x, \omega)| \leq |\hat{f}(\omega)| = \rho(\omega).$$

Theorem 1. *There exist $\omega_i \in \mathbb{R}^d$, s.t., $G = \mathbb{R}$ and*

$$(1.27) \quad \int_{\Omega} (f(x) - f_n(x))^2 \leq \frac{1}{n} \int_{\mathbb{R}^d} |\hat{f}(\omega)| d\omega,$$

where

$$(1.28) \quad f_n(x) = \frac{\|f\|_{L^1}}{n} \sum_{i=1}^n \frac{\cos(\omega_i^* \cdot x + \beta_i^*)}{\rho(\omega_i^*)}.$$

Note that

$$(1.29) \quad f_n = \sum_{i=1}^n \frac{\cos(\omega_i^* \cdot x + \beta_i^*)}{\rho(\omega_i^*)} \in {}_n\mathbf{N}(\sigma, n),$$

with

$$(1.30) \quad \sigma(t) = \cos(t).$$

1.2.2 Double Fourier representation

Assume that σ is a locally Riemann integrable function and $\sigma \in L^1(\mathbb{R})$ and thus the Fourier transform of σ is well-defined and continuous. Since σ is non-zero and

$$(1.31) \quad \hat{\sigma}(\omega) = \frac{1}{2\pi} \int_{\mathbb{R}} \sigma(t) e^{-i\omega t} dt,$$

this implies that $\hat{\sigma}(a) \neq 0$ for some $a \neq 0$. Via a change of variables $t = w \cdot x + b$ and $dt = db$, this means that for all x and ω , we have

$$(1.32) \quad \begin{aligned} 0 \neq \hat{\sigma}(a) &= \frac{1}{2\pi} \int_{\mathbb{R}} \sigma(\omega \cdot x + b) e^{-ia(\omega \cdot x + b)} db \\ &= e^{-ia\omega \cdot x} \frac{1}{2\pi} \int_{\mathbb{R}} \sigma(\omega \cdot x + b) e^{-iab} db, \end{aligned}$$

and so

$$(1.33) \quad e^{ia\omega \cdot x} = \frac{1}{2\pi \hat{\sigma}(a)} \int_{\mathbb{R}} \sigma(\omega \cdot x + b) e^{-iab} db.$$

Likewise, since the growth condition also implies that $\sigma^{(k)} \in L^1$, we can differentiate the above expression under the integral with respect to x .

This allows us to write the Fourier mode $e^{ia\omega \cdot x}$ as an integral of neuron output functions. We substitute this into the Fourier representation of f (note that the assumption we make implies that $\hat{f} \in L^1$ so this is rigorously justified for a.e. x) to get

$$(1.34) \quad f(x) = \int_{\mathbb{R}^d} e^{i\omega \cdot x} \hat{f}(\omega) d\omega = \int_{\mathbb{R}^d} \int_{\mathbb{R}} \frac{1}{2\pi \hat{\sigma}(a)} \sigma(a^{-1}\omega \cdot x + b) \hat{f}(\omega) e^{-iab} db d\omega = \int_{\mathbb{R}^d \times \mathbb{R}} k(x, \theta) d\theta$$

where $\theta = (\omega, b)$ and

$$k(x, \theta) = \frac{1}{2\pi \hat{\sigma}(a)} \sigma(a^{-1}\omega \cdot x + b) \hat{f}(\omega) e^{-iab}.$$

Thus we have

$$(1.35) \quad \begin{aligned} |k(x, \theta)| &\leq \frac{1}{2\pi |\hat{\sigma}(a)|} \max_{x \in \Omega} |\sigma(a^{-1}\omega \cdot x + b)| |\hat{f}(\omega)|, \\ &\leq h(\omega, b) |\hat{f}(\omega)| = \rho(\theta) \end{aligned}$$

where

$$(1.36) \quad h(\omega, b) = \max_{x \in \Omega} |\sigma(a^{-1}\omega \cdot x + b)|.$$

if we ignore the coefficient. Thus, following the discussion in last section, the next step is to analyze $h(\omega, b)$ which we will discuss in the next section. Before that, let us introduce a special case of the above representation once the activation function is periodic.

References

- [1] L. Chen, P. Sun, and J. Xu. Optimal anisotropic meshes for minimizing interpolation errors in l^p -norm. *Mathematics of Computation*, 76(257):179–204, 2007.
- [2] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [3] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *European Conference on Computer Vision*, pages 630–645. Springer, 2016.
- [4] M. Leshno, V. Y. Lin, A. Pinkus, and S. Schocken. Multilayer feedforward networks with a nonpolynomial activation function can approximate any function. *Neural networks*, 6(6):861–867, 1993.
- [5] A. Pinkus. Approximation theory of the mlp model in neural networks. *Acta numerica*, 8:143–195, 1999.
- [6] L. R. Scott and S. Zhang. Finite element interpolation of nonsmooth functions satisfying boundary conditions. *Mathematics of Computation*, 54(190):483–493, 1990.
- [7] J. Xu. Finite element methods. Lecture notes, 2020.