

Characteristic Functionals and the Innovations Approach to Stochastic Processes (With Applications to Random Neural Networks)

Rahul Parhi
UCSD ECE

UCSD Probability Seminar
6 February 2025

Outline

- Classical Formulation of Lévy Processes
- Innovation Model (Bode, Shannon, and Kailath, *ca.* 1950–1970)
- (Lévy) White Noises
- Generalized Stochastic Processes (Itô and Gelfand, *ca.* 1955)
- Generalized Lévy Processes and Stochastic Differential Equations
- Exact Characterization of the Law of Random Neural Networks

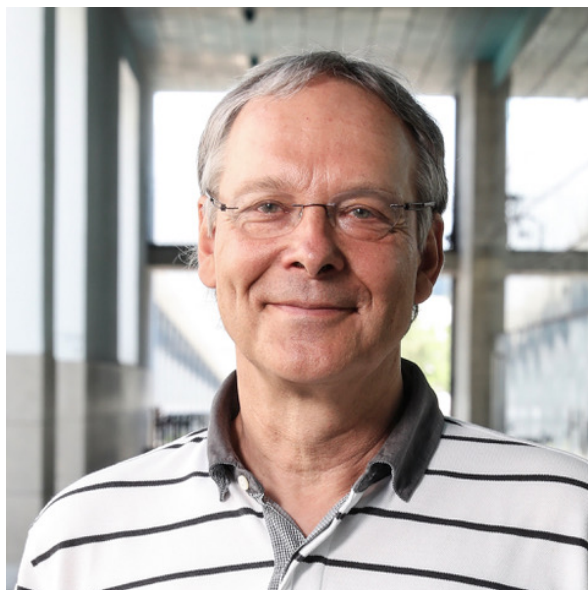
Joint Work With...



Pakshal Bohra



Mehrsa Pourya



Michael Unser



Ayoub El Biari

Classical Formulation of Lévy Processes

Definition (Lévy process)

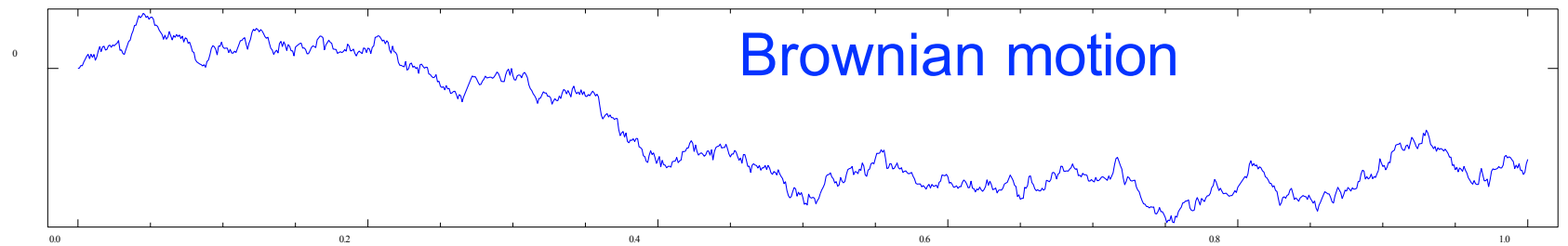
A stochastic process $(S_t)_{t \geq 0}$ is called a *Lévy process* if

1. $S_0 = 0$ almost surely;
2. Independence of increments: For any $0 \leq t_1 < \dots < t_n < \infty$, $X_{t_2} - X_{t_1}, \dots, X_{t_n} - X_{t_{n-1}}$ are mutually independent;
3. Stationary increments: For any $t_1 < t_2$, $S_{t_2} - S_{t_1} \stackrel{\mathcal{L}}{=} S_{t_2 - t_1}$;
4. Stochastic continuity: For any $\varepsilon > 0$ and $t \geq 0$

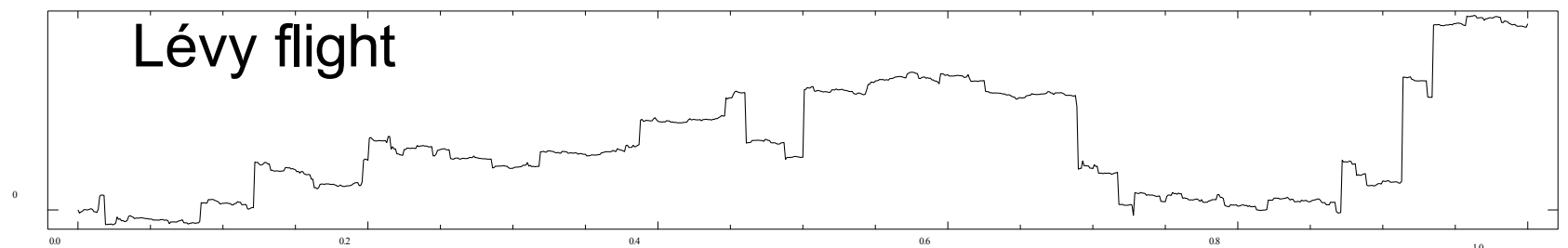
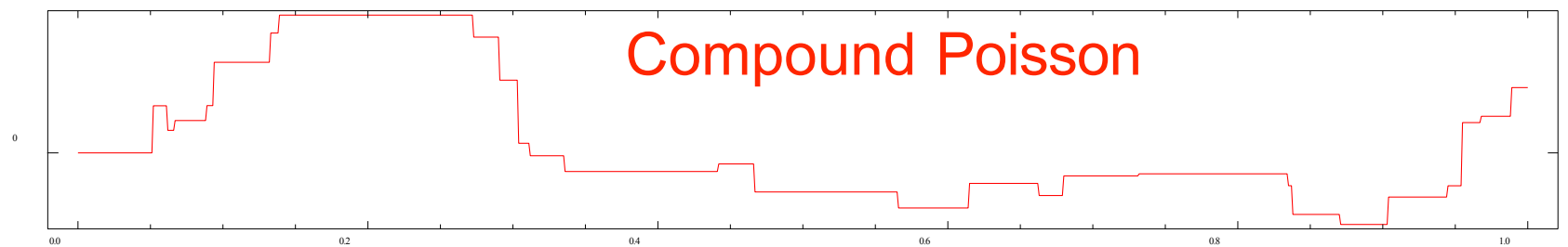
$$\lim_{h \rightarrow 0} \mathbf{P}(|S_{t+h} - S_t| > \varepsilon) = 0.$$

Lévy processes are càdlàg and tightly linked to infinite divisibility

Examples of Lévy Processes



Wiener 1923

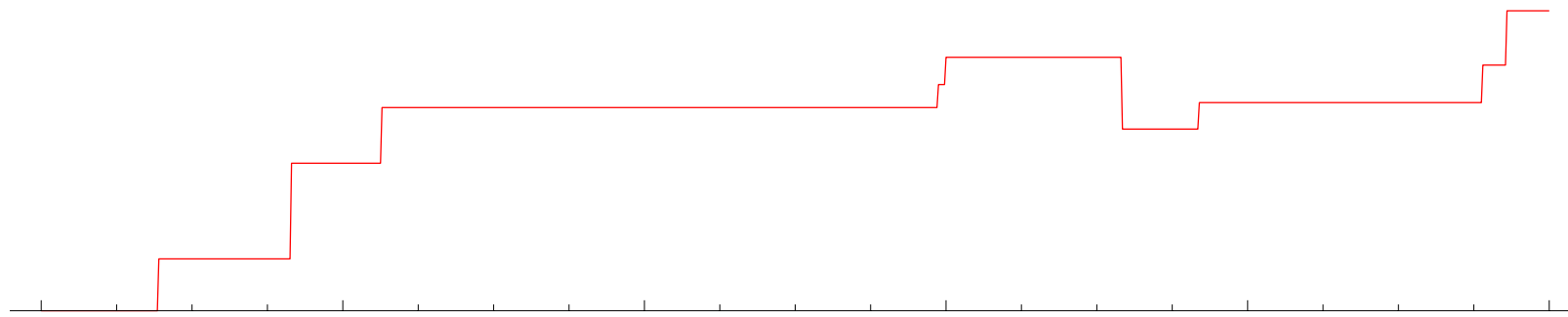


Lévy processes are often used to model various components of signal processing and communication systems

Decoupling the Underlying Randomness

Is there a way to decouple the correlation properties of a Lévy process from its underlying randomness?

Example: Compound Poisson process $s(t)$ (rate $\lambda > 0$ and jump law \mathbf{P}_V)



$$Ds(t) = \dot{s}(t) = \sum_k v_k \delta_{b_k}$$

innovation

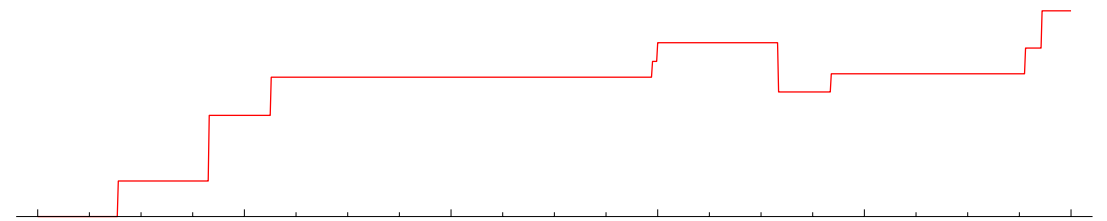
$v_k \sim \mathbf{P}_V$ i.i.d.

$\{b_k\}_k$ is a Poisson point process with rate λ

The underlying randomness of a compound Poisson process is completely determined by its innovation process

Innovation-Based Synthesis

$$\sum_k v_k \delta(\cdot - b_k) \xrightleftharpoons[\text{D}]{\text{???}}$$

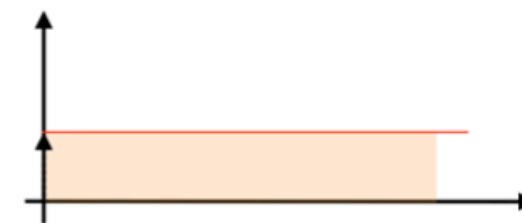


Integrate the innovation process?

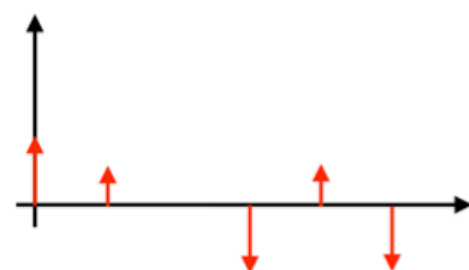
Recall: $s(0) = 0$ a.s.



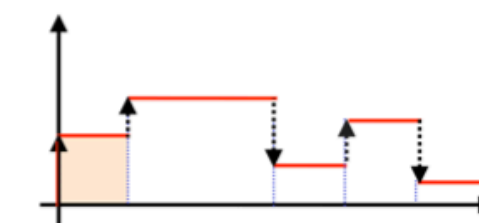
$$\xrightarrow{\text{D}^{-1}}$$



$$\xrightarrow{\text{D}^{-1}}$$



$$\xrightarrow{\text{D}^{-1}}$$



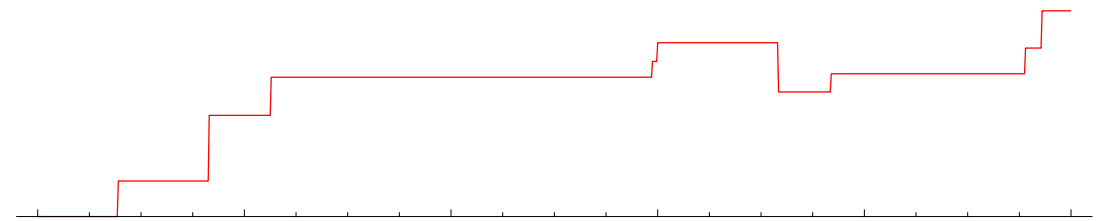
Compound Poisson Process

- SDE formulation

$$Ds \stackrel{\mathcal{L}}{=} w \quad \text{s.t.} \quad s(0) = 0$$

$$\text{Innovation: } w = \sum_k v_k \delta(\cdot - b_k)$$

Poisson innovation, or
Poisson white noise

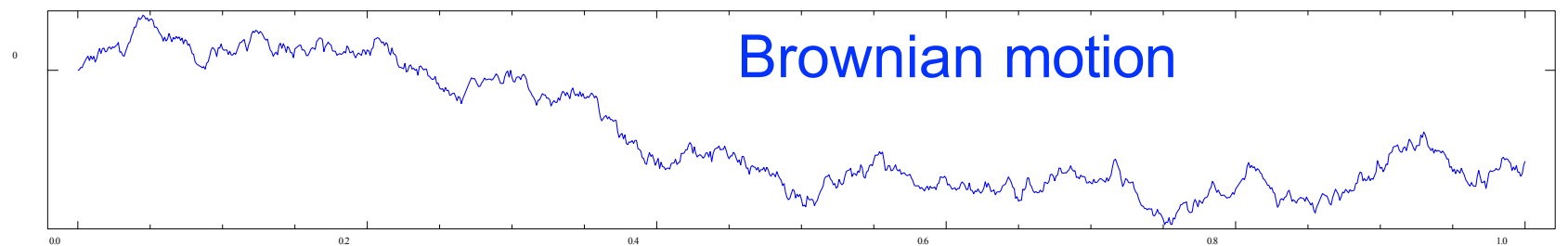


- “Formal” solution

$$\begin{aligned} s(t) &= D_0^{-1} w(t) = \sum_k v_k D_0^{-1} \{ \delta(\cdot - b_k) \}(t) \\ &= \sum_k v_k \left(H(t - b_k) - H(-b_k) \right) \end{aligned}$$

imposing boundary condition

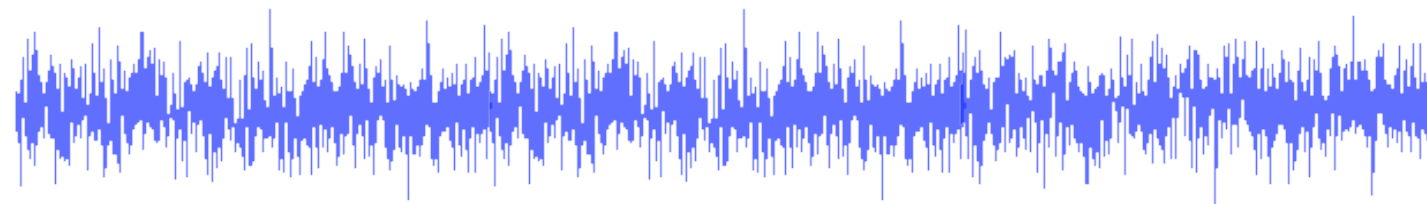
Brownian Motion and Gaussian White Noise



$b(t)$

D

$w(t)$



Gaussian innovation, or
Gaussian white noise

Gaussian white noise does not admit a pointwise representation

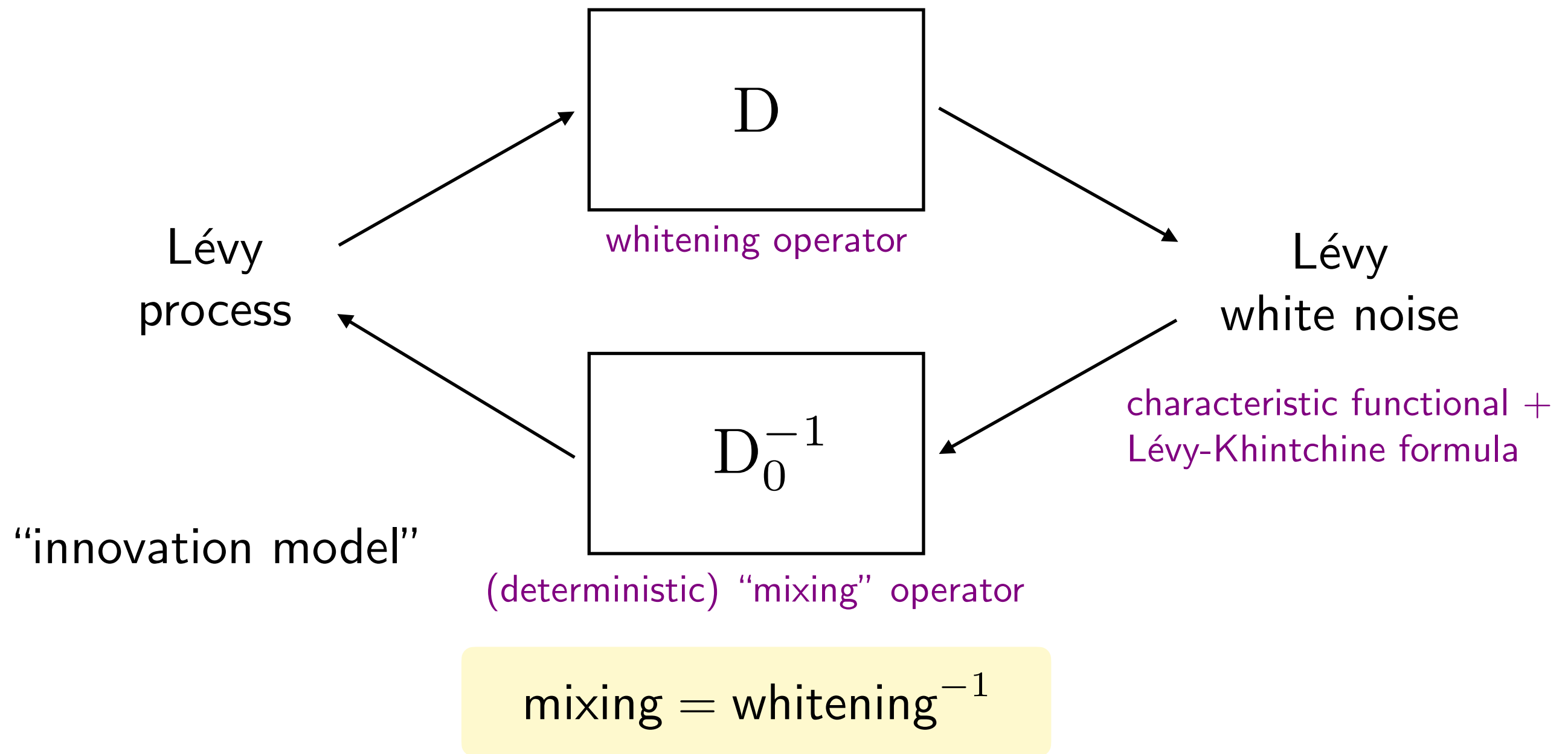
$w(t)$ is a random measure, or,
a random distribution (generalized function)

$$\mathbf{P}(w \in \mathcal{D}'(\mathbb{R})) = 1$$

$$\mathbf{P}(w \in \mathcal{S}'(\mathbb{R})) = 1$$

(Gelfand, 1955)

Innovations Approach to Lévy Processes



It is much easier to study white noise processes than it is to study general random processes.

Generalized Stochastic Processes



Kiyosi Itô



Israel Gelfand

MEMOIRS OF THE COLLEGE OF SCIENCE, UNIVERSITY OF KYOTO, SERIES A
Vol. XXVIII, Mathematics No. 3, 1953.

Stationary random distributions

By
Kiyosi Itô

(Received April 15, 1954)

In the same way as the concept of distributions by L. Schwartz [11]¹⁾ was introduced as an extended one of functions, we may define stationary random distributions as an extension of stationary random functions viz. stationary processes. Such consideration will enable us to establish a unified theory of stationary processes, Brownian motion processes, processes with stationary increments and other

Framework of generalized
stochastic processes

ca. 1950s

Stochastic counterpart to
Schwartz' theory of distributions

1.

Generalized random processes

Dokl. Akad. Nauk SSSR **100** (1955) 853–856. Zbl. **64**:111

1. Usually, a random process is defined by probability distributions of random variables $(x(t_1), \dots, x(t_n))$ for n arbitrary moments of time. However, one can give many examples of random processes which are important in practice for which such probability distributions do not exist. An example of such a process is the white noise that can be obtained, roughly speaking, as a superposition of all frequencies with random amplitudes; these amplitudes are independent, identically distributed Gaussian variables. In this note we introduce a description of random processes that appears to cover all practically important examples. If $x(t)$ is a random process, then any “linear apparatus” gives us a probability distribution

Warm-Up: Classical Probability Theory

$\mathbf{X} \in \mathbb{R}^d$ is a random vector with law $\mathbf{P}_{\mathbf{X}}$

This means that

$\mathbf{X} : (\Omega, \mathcal{F}, \mathbf{P}) \rightarrow (\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$ is measurable.
complete
probability space

The law of \mathbf{X} is given by the *pushforward measure*

$$\begin{aligned}\mathbf{P}_{\mathbf{X}}(A) &= (\mathbf{X}_{\#}\mathbf{P})(A) := \mathbf{P}(\mathbf{X}^{-1}(A)) = \mathbf{P}(\{\omega \in \Omega : \mathbf{X}(\omega) \in A\}) \\ &= \mathbf{P}(\mathbf{X} \in A).\end{aligned}$$

The *characteristic function* of \mathbf{X} is given by

$$\hat{\mathbf{P}}_{\mathbf{X}}(\boldsymbol{\xi}) = \mathbf{E}[e^{i\mathbf{X}^T \boldsymbol{\xi}}], \quad \boldsymbol{\xi} \in \mathbb{R}^d.$$

Once you have the characteristic function, you have everything.

Bochner's and Lévy's Theorem

Bochner's Theorem

A function $\hat{\mathbf{P}}$ is the characteristic function of a random variable $\mathbf{X} \in \mathbb{R}^d$ if and only if $\hat{\mathbf{P}}$ is continuous, positive-definite, and satisfies $\hat{\mathbf{P}}(\mathbf{0}) = 1$.

Lévy's Continuity Theorem

Let $(\mathbf{X}_n)_{n \in \mathbb{N}}$ and \mathbf{X} be random variables in \mathbb{R}^d . The sequence \mathbf{X}_n converges in law to \mathbf{X} if and only if, for all $\boldsymbol{\xi} \in \mathbb{R}^d$,

$$\hat{\mathbf{P}}_{\mathbf{X}_n}(\boldsymbol{\xi}) \xrightarrow{n \rightarrow \infty} \hat{\mathbf{P}}_{\mathbf{X}}(\boldsymbol{\xi}).$$

Generalized Stochastic Processes

A generalized stochastic process is a random variable that takes values in the dual of a nuclear space.

Let $(\mathcal{N}, \mathcal{N}')$ denote a nuclear space and its dual. A generalized stochastic process s is a random variable, i.e., a measurable map

$$s : (\Omega, \mathcal{F}, \mathbf{P}) \rightarrow (\mathcal{N}', \mathcal{B}_c(\mathcal{N}')).$$

The law of s is the probability measure $\mathbf{P}_s := s_{\#}\mathbf{P}$, which is defined on $\mathcal{B}_c(\mathcal{N}')$. The *characteristic functional* of s is given by

$$\hat{\mathbf{P}}_s(\varphi) = \mathbf{E}[e^{i\langle s, \varphi \rangle_{\mathcal{N}' \times \mathcal{N}}}], \quad \varphi \in \mathcal{N}.$$

Examples: $(\mathbb{R}^d, \mathbb{R}^d)$, $(\mathcal{D}(\mathbb{R}^d), \mathcal{D}'(\mathbb{R}^d))$, $(\mathcal{S}(\mathbb{R}^d), \mathcal{S}'(\mathbb{R}^d))$, etc.

Once you have the characteristic functional, you have everything.

Generalized Stochastic Processes

This formulation allows us to study classical stochastic processes $(s(\mathbf{x}))_{\mathbf{x} \in \mathbb{R}^d}$ as well as those that do not admit a pointwise representation such as white noise.

Example (Gaussian Processes)

A generalized stochastic process s that takes values in \mathcal{N}' is called *Gaussian* if its characteristic functional is of the form

$$\hat{\mathbf{P}}_s(\varphi) = \exp \left(i\mu_s(\varphi) - \frac{1}{2}\Sigma_s(\varphi, \varphi) \right), \quad \varphi \in \mathcal{N},$$

where $\mu_s(\varphi) = \mathbf{E}[\langle s, \varphi \rangle]$ denotes the mean functional and $\Sigma_s(\varphi_1, \varphi_2) = \mathbf{E}[(\langle s, \varphi_1 \rangle - \mu_s(\varphi_1))(\langle s, \varphi_2 \rangle - \mu_s(\varphi_2))]$ denotes the covariance functional of the process.

⇒ Backwards compatible with space-indexed Gaussian processes.

Is Nuclearity Necessary?

Bochner–Minlos Theorem

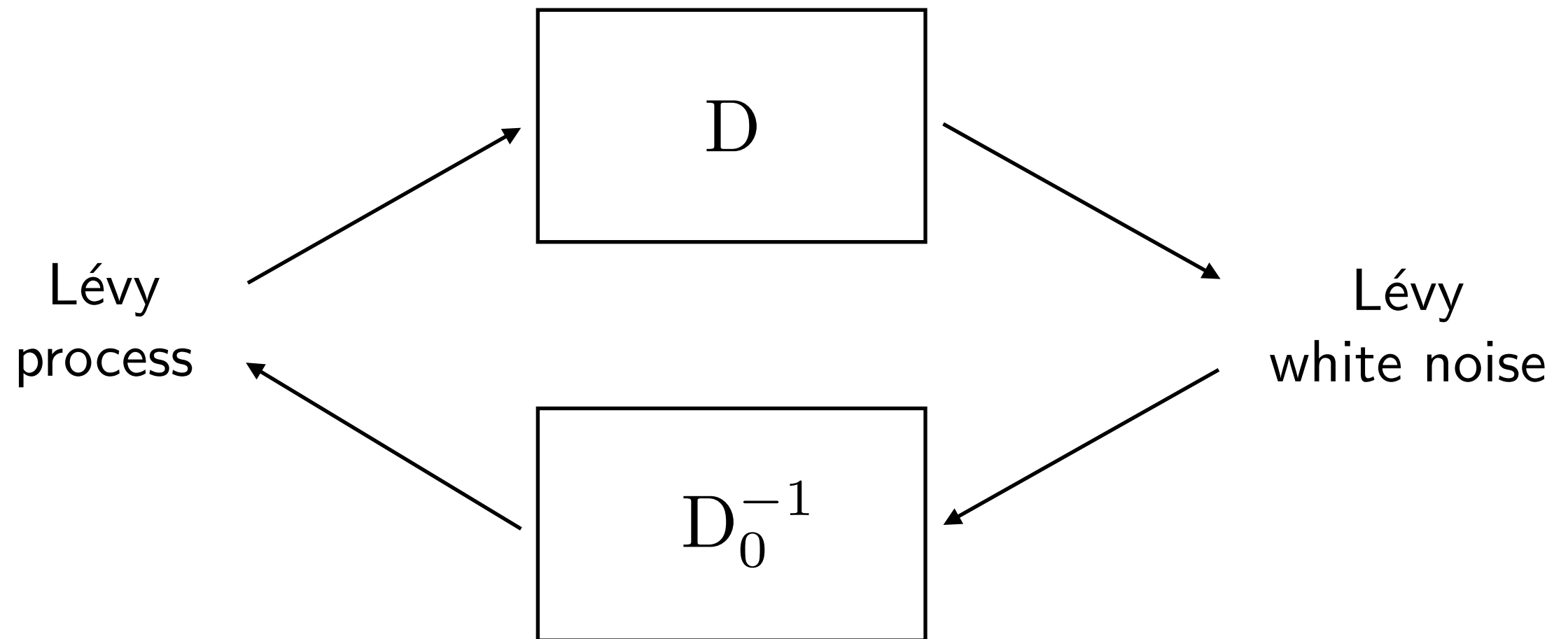
A function $\hat{\mathbf{P}}$ is the characteristic function of a stochastic process $s \in \mathcal{N}'$ if and only if $\hat{\mathbf{P}}$ is continuous, positive-definite, and satisfies $\hat{\mathbf{P}}(0) = 1$.

Lévy–Fernique Continuity Theorem

Let $(s_n)_{n \in \mathbb{N}}$ and s be stochastic processes in \mathcal{N}' . The sequence s_n converges in law to s if and only if, for all $\varphi \in \mathcal{N}$,

$$\hat{\mathbf{P}}_{s_n}(\varphi) \xrightarrow{n \rightarrow \infty} \hat{\mathbf{P}}_s(\varphi).$$

Classical Lévy Processes: A New Perspective



- SDE formulation

$$Ds \stackrel{\mathcal{L}}{=} w \quad \text{s.t.} \quad s(0) = 0$$

- “Formal” solution

$$s = D_0^{-1}w$$

What is the characteristic functional of a Lévy white noise?

Lévy White Noise Characteristic Functional

Theorem (Gelfand and Vilenkin, 1964)

Let $s(t)$ be a Lévy process on \mathbb{R} and let $w = Ds$. Then, $w \in \mathcal{D}'(\mathbb{R})$ almost surely and

$$\hat{\mathbf{P}}_w(\varphi) = \exp \left(\int_{\mathbb{R}} f(\varphi(x)) \, dx \right), \quad \varphi \in \mathcal{D}(\mathbb{R}),$$

where

$$f(\xi) = i\mu\xi - \frac{\sigma^2\xi^2}{2} + \int_{\mathbb{R}} e^{i\xi t} - 1 - i\xi \mathbf{1}_{[-1,1]}(t) \, d\mathbf{P}_V(t),$$

with $\mu \in \mathbb{R}$, $\sigma^2 \geq 0$, and \mathbf{P}_V is a Lévy measure, i.e.,

$$\int_{\mathbb{R}} \min\{1, t^2\} \, d\mathbf{P}_V(t) < \infty \quad \text{and} \quad \mathbf{P}_V(\{0\}) = 0.$$

$(\mu, \sigma^2, \mathbf{P}_V)$ is called a **Lévy triplet** and f is a **Lévy exponent**

Examples

- Poisson white noise w_{Poi} on \mathbb{R} with rate λ and jump law \mathbf{P}_V

$$\hat{\mathbf{P}}_{w_{\text{Poi}}}(\varphi) = \exp \left(\lambda \int_{\mathbb{R}} \int_{\mathbb{R}} \left(e^{iv\varphi(t)} - 1 \right) dt d\mathbf{P}_V(v) \right)$$

Duttweiler and Kailath (1973)

$$\text{cf.}, \xi \mapsto \exp \left(\lambda(e^{i\xi} - 1) \right)$$

- Gaussian white noise w_{Gauss} with unit variance

$$\hat{\mathbf{P}}_{w_{\text{Gauss}}}(\varphi) = \exp \left(-\frac{\|\varphi\|_{L^2}^2}{2} \right)$$

Gelfand and Vilenkin (1964)

$$\text{cf.}, \xi \mapsto \exp \left(-\frac{|\xi|^2}{2} \right)$$

Lévy Process Characteristic Functional

$$\hat{\mathbf{P}}_s(\varphi) = \mathbf{E}[e^{i\langle s, \varphi \rangle}], \quad \varphi \in \mathcal{D}(\mathbb{R})$$

$$s \stackrel{\mathcal{L}}{=} D_0^{-1}w$$

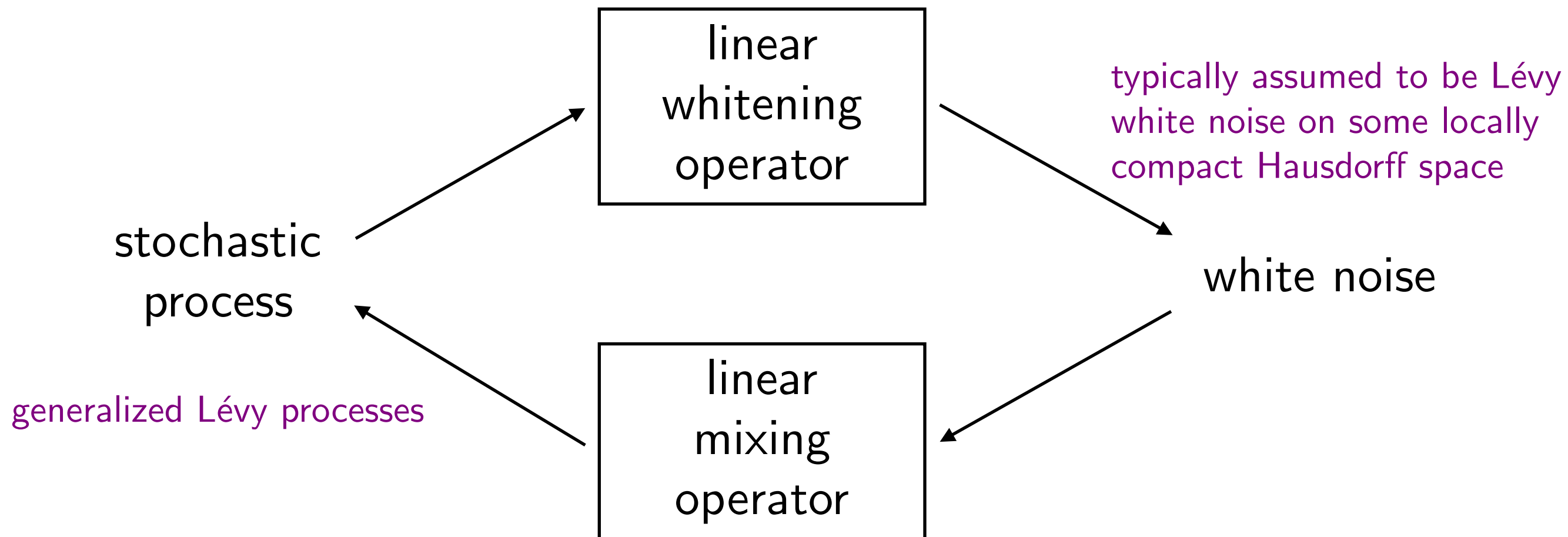
$$\begin{aligned} \langle s, \varphi \rangle &\stackrel{\mathcal{L}}{=} \langle D_0^{-1}w, \varphi \rangle \\ &= \langle w, D_0^{-1*}\varphi \rangle \end{aligned}$$

Fageot and Humeau (2021)

$$\begin{aligned} \hat{\mathbf{P}}_s(\varphi) &= \mathbf{E}[e^{i\langle s, \varphi \rangle}] = \mathbf{E}[e^{i\langle w, D_0^{-1*}\varphi \rangle}] \\ &= \hat{\mathbf{P}}_w(D_0^{-1*}\varphi) \end{aligned}$$

Given the characteristic functional of the innovation process, you automatically have the characteristic functional of the original process.

A More Abstract Innovation Model



$$\text{mixing} = \text{whitening}^{-1}$$

existence of suitable operators
is an active area of research

$$\hat{\mathbf{P}}_s(\varphi) = \hat{\mathbf{P}}_w(L^{-1*}\varphi)$$

Tempered Lévy White Noise

Theorem (Dalang and Humeau, 2017)

Let w be a Lévy white noise on \mathbb{R}^d , i.e., $w \in \mathcal{D}'(\mathbb{R}^d)$ and

$$\hat{\mathbf{P}}_w(\varphi) = \exp \left(\int_{\mathbb{R}^d} f(\varphi(\mathbf{x})) \, d\mathbf{x} \right), \quad \varphi \in \mathcal{D}(\mathbb{R}^d),$$

where

$$f(\xi) = i\mu\xi - \frac{\sigma^2\xi^2}{2} + \int_{\mathbb{R}} e^{i\xi t} - 1 - i\xi \mathbf{1}_{[-1,1]}(t) \, d\mathbf{P}_V(t),$$

with Lévy triplet $(\mu, \sigma^2, \mathbf{P}_V)$. Then, w is tempered (a.s.) if and only if there exists $\varepsilon > 0$ such that

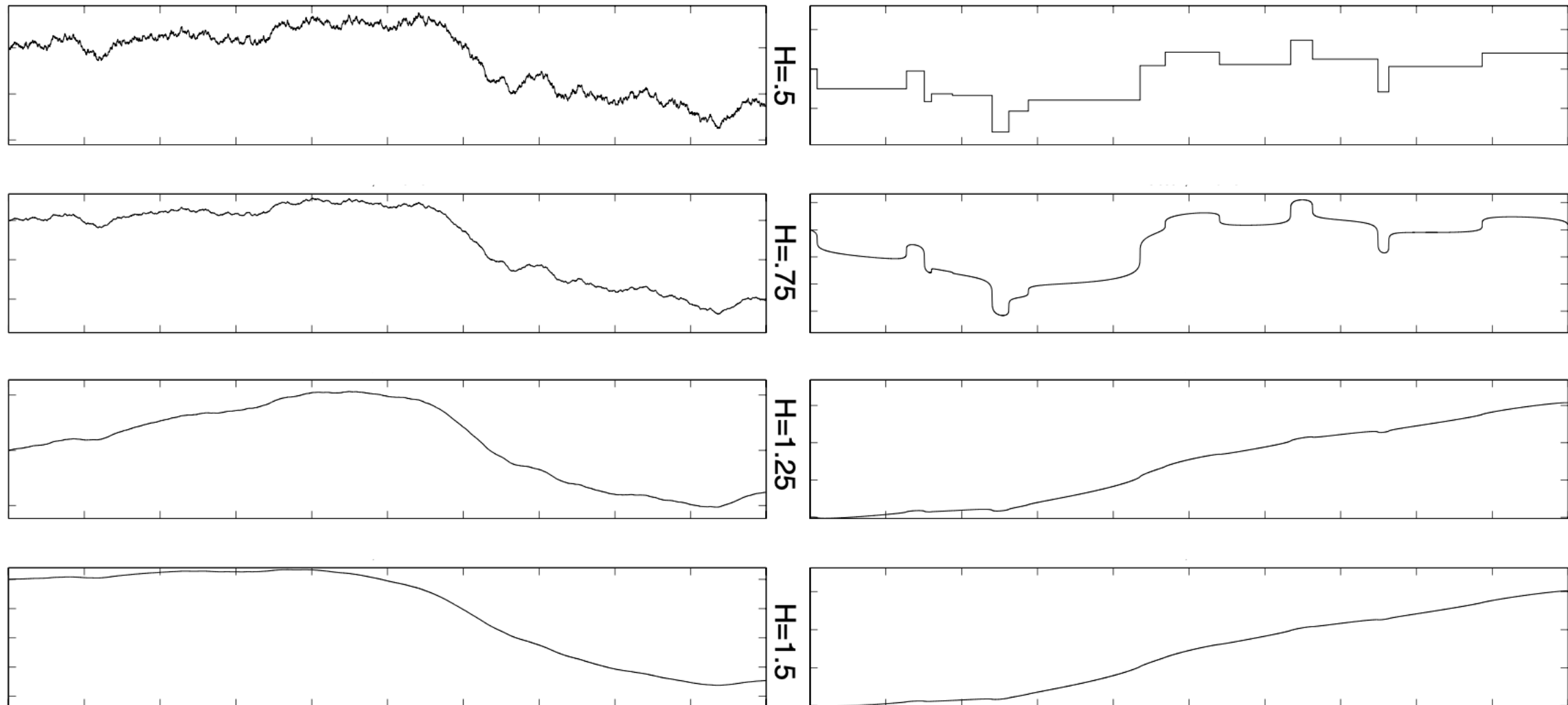
$$\mathbf{E}[|V|^\varepsilon] < \infty, \quad V \sim \mathbf{P}_V.$$

$$\mathbf{P}(w \in \mathcal{S}'(\mathbb{R}^d)) = 1 \Leftrightarrow \mathbf{P}(w \in \mathcal{D}'(\mathbb{R}^d) \setminus \mathcal{S}'(\mathbb{R}^d)) = 0 \Leftrightarrow \text{supp}(\mathbf{P}_w) \subset \mathcal{S}'(\mathbb{R}^d)$$

Fractional Order Processes

$$D^{H+\frac{1}{2}}s \stackrel{\mathcal{L}}{=} w$$

Blu and Unser (2007)



Gaussian

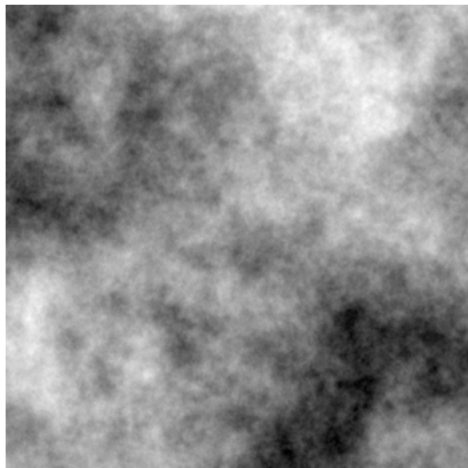
Sparse (generalized Poisson)

Mandelbrot and Van Ness (1968)

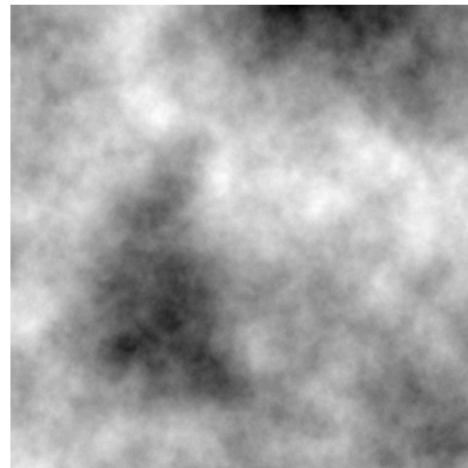
Fractional Order Processes

$$(-\Delta)^{\frac{H+1}{2}} s \stackrel{\mathcal{L}}{=} w$$

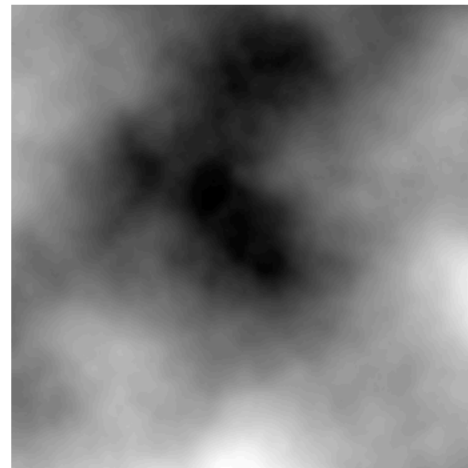
Gaussian



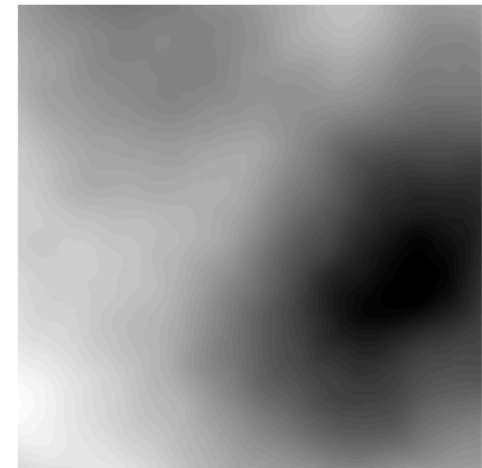
H=.5



H=.75

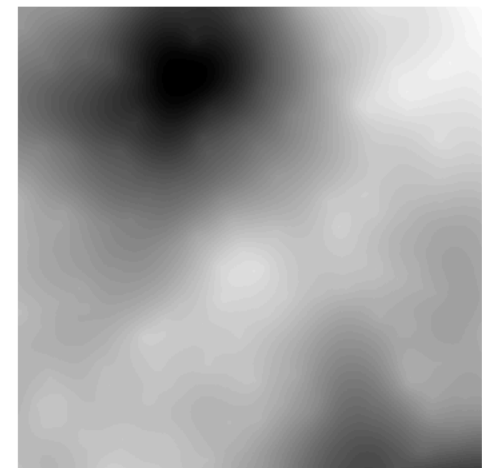
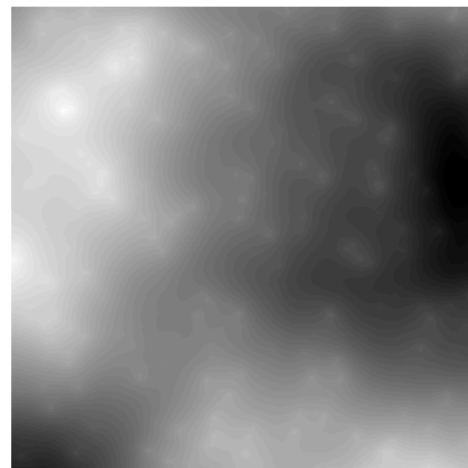
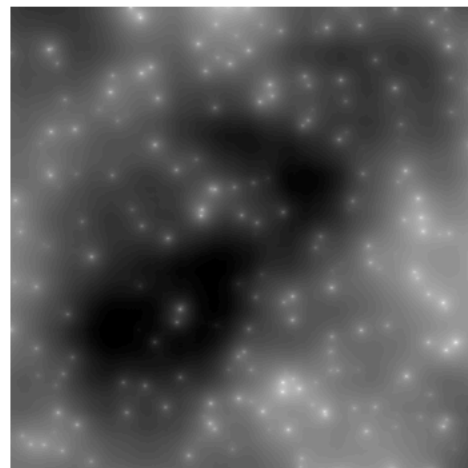
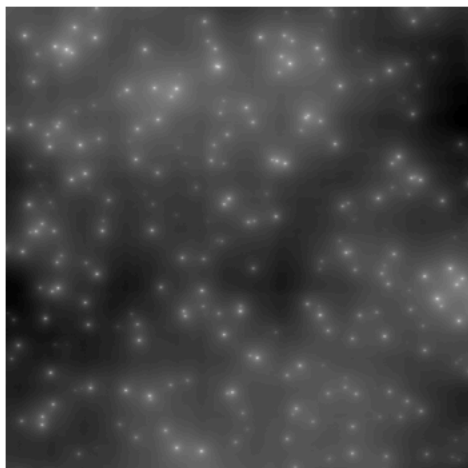


H=1.25

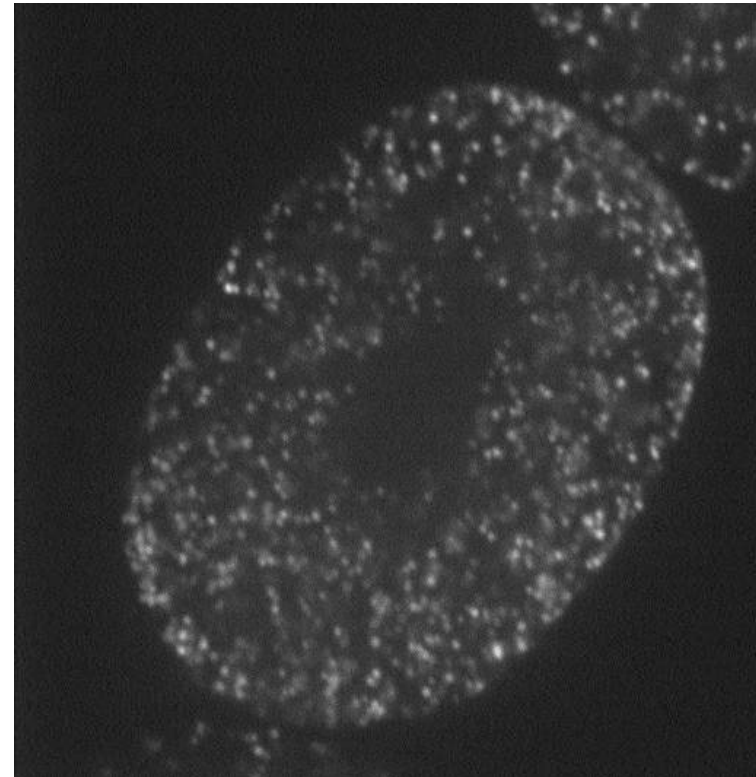


H=1.75

Sparse (generalized Poisson)



Sparse Stochastic Processes are Good Models

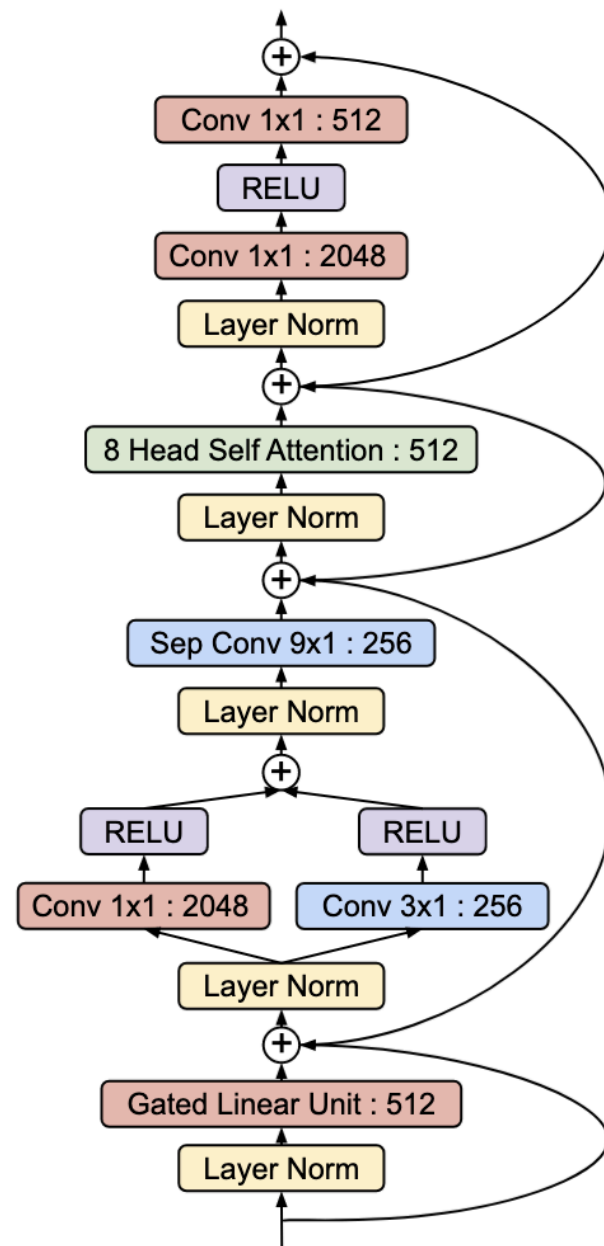


Applications to Random Neural Networks

Deep Neural Network Architectures

The Evolved Transformer

David So, Quoc Le, Chen Liang *Proceedings of the 36th International Conference on Machine Learning*, PMLR 97:5877-5886, 2019.



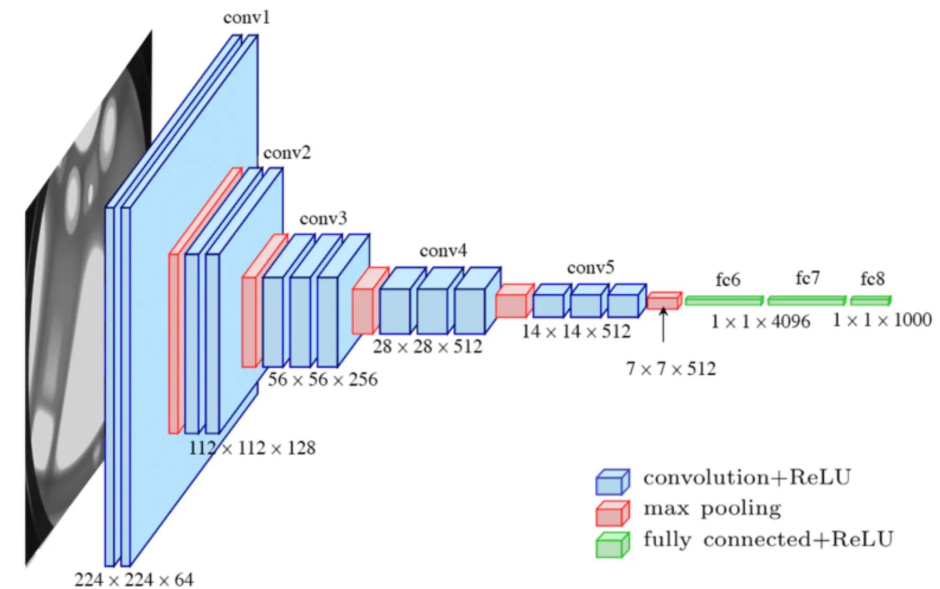
Google DeepMind

2023-10-26

ConvNets Match Vision Transformers at Scale

Samuel L Smith¹, Andrew Brock¹, Leonard Berrada¹ and Soham De¹

¹Google DeepMind

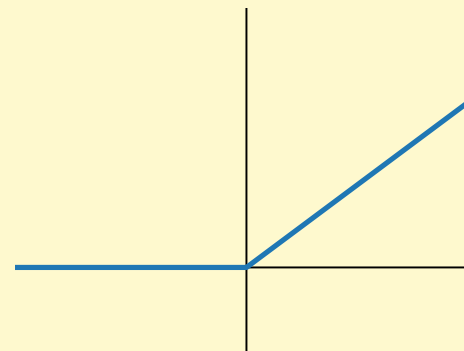


Very deep convolutional networks for large-scale image recognition

[K Simonyan, A Zisserman](#)

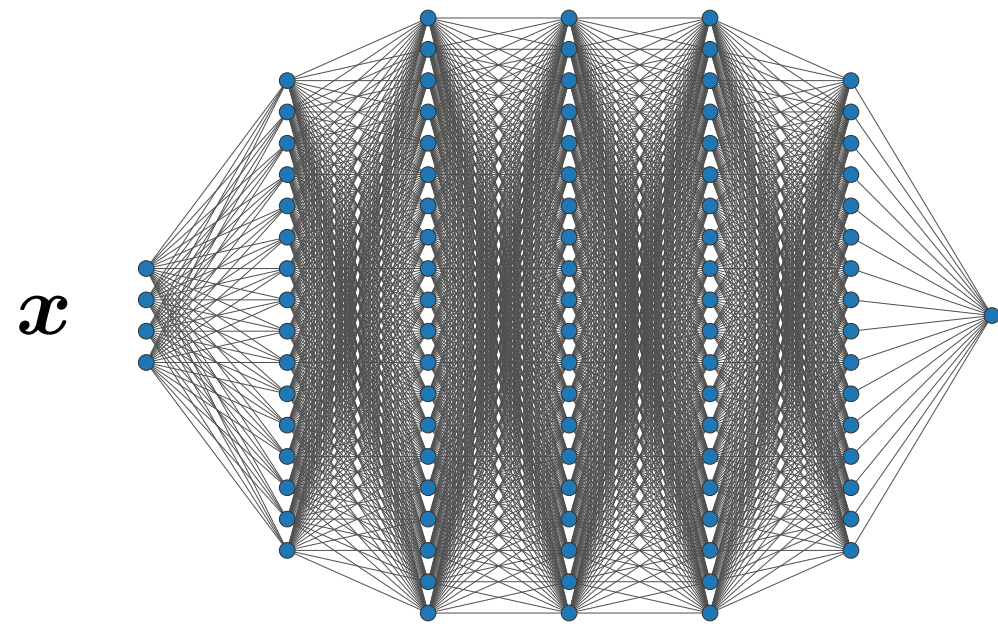
In this work we investigate the effect of the convolutional network depth on its accuracy in the large-scale image recognition setting. Our main contribution is a thorough evaluation of ...

☆ Save 📄 Cite Cited by 112161 Related articles 🔗



Rectified Linear Unit (ReLU)
 $\text{ReLU}(t) = \max\{0, t\} = t_+$

Neural Network Training



parameterized by a vector $\theta \in \mathbb{R}^P$
of neural network **weights**

$$f_{\theta}(x) = \sigma(\mathbf{W}_L \sigma(\mathbf{W}_{L-1} \sigma(\cdots \sigma(\mathbf{W}_1 x))))$$

Neural network training problem for the data $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$.

$$\min_{\theta \in \mathbb{R}^P} \underbrace{\sum_{n=1}^N \mathcal{L}(y_n, f_{\theta}(\mathbf{x}_n))}_{\text{data fidelity}} + \underbrace{\frac{\lambda}{2} \|\theta\|_2^2}_{\text{regularization}}$$

Neural Network Training

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^P} \underbrace{\sum_{n=1}^N \mathcal{L}(y_n, f_{\boldsymbol{\theta}}(\mathbf{x}_n))}_{\mathcal{L}(\boldsymbol{\theta})} + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

Gradient descent update on θ_i

$$\theta_i^{t+1} = \theta_i^t - \tau \left(\left. \frac{\partial \mathcal{L}}{\partial \theta_i} \right|_{\theta_i = \theta_i^t} + \lambda \theta_i^t \right)$$

step size
“learning rate”

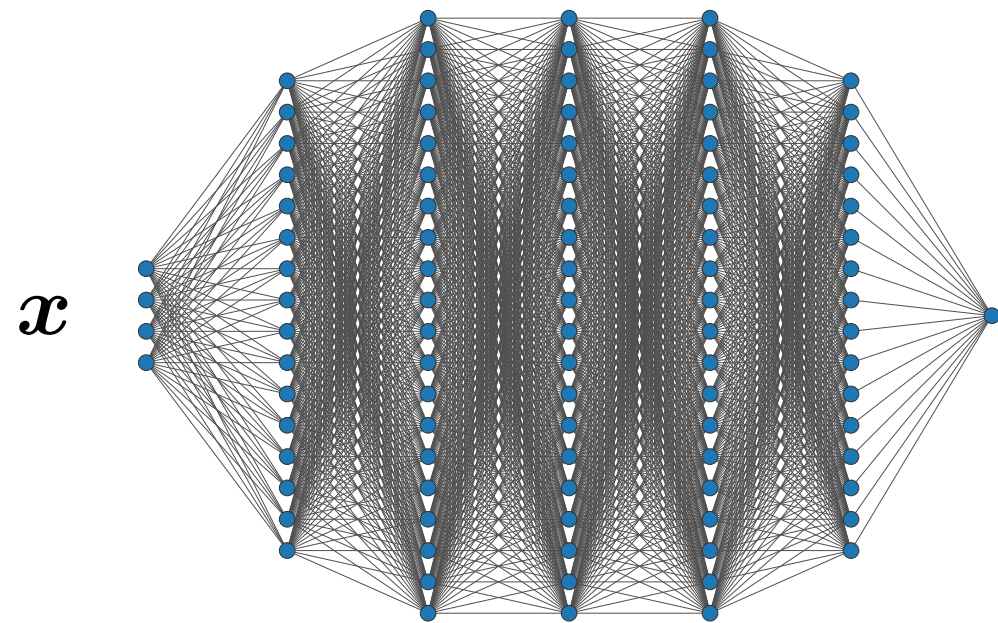
How do we choose θ_i^0 ?

randomly!

Hanson and Pratt (1988, NeurIPS)

Krogh and Hertz (1990, NeurIPS)

Random Neural Networks



parameterized by a vector $\theta \in \mathbb{R}^P$
of neural network **weights**

$$f_{\theta}(x) = \sigma(\mathbf{W}_L \sigma(\mathbf{W}_{L-1} \sigma(\cdots \sigma(\mathbf{W}_1 x))))$$

A random neural network is a random function.

f_{θ}

$$\theta \sim Q$$

A random neural network is a stochastic process.

$$(f_{\theta}(x))_{x \in \mathbb{R}^d}$$

Machine Learning Folklore

Folklore Theorem

The parameters of neural networks trained with GD do not move far from their initialization.

If we can understand neural networks at initialization,
then we can understand everything!

DEEP NEURAL NETWORKS AS GAUSSIAN PROCESSES

**Jaehoon Lee^{*†}, Yasaman Bahri^{*†}, Roman Novak, Samuel S. Schoenholz,
Jeffrey Pennington, Jascha Sohl-Dickstein**

Google Brain
{jaehlee, yasamanb, romann, schsam, jpennin, jaschasd}@google.com

ICLR 2018
1300+ citations

Neural Tangent Kernel: Convergence and Generalization in Neural Networks

Arthur Jacot
École Polytechnique Fédérale de Lausanne
arthur.jacot@netopera.net

Franck Gabriel
Imperial College London and École Polytechnique Fédérale de Lausanne
franckrgabriel@gmail.com

Clément Hongler
École Polytechnique Fédérale de Lausanne
clement.hongler@gmail.com

NeurIPS 2018
3800+ citations

Spawned a Burgeoning Industry of Research



Q Create account Log in ...

Neural network Gaussian process

3 languages

Article Talk

Read Edit View history Tools

From Wikipedia, the free encyclopedia

A **Neural Network Gaussian Process** (NNGP) is a [Gaussian process](#) (GP) obtained as the limit of a certain type of sequence of [neural networks](#). Specifically, a wide variety of network architectures converges to a GP in [the infinitely wide limit, in the sense of distribution](#).^{[1][2][3][4][5][6][7][8]} The concept constitutes an [intensional definition](#), i.e., a NNGP is just a GP, but distinguished by how it is obtained.



Q Create account Log in ...

Neural tangent kernel

2 languages

Article Talk

Read Edit View history Tools

From Wikipedia, the free encyclopedia

In the study of [artificial neural networks](#) (ANNs), the **neural tangent kernel** (NTK) is a [kernel](#) that describes the evolution of [deep artificial neural networks](#) during their training by [gradient descent](#). It allows ANNs to be studied using theoretical tools from [kernel methods](#).

In general, a kernel is a [positive-semidefinite symmetric](#) function of two inputs which represents some notion of similarity between the two inputs. The NTK is a specific kernel derived from a given neural network; in general, when the neural network parameters change during training, the NTK evolves as well. However, in the limit of large layer width the NTK becomes constant, revealing a duality between training the wide neural network and kernel methods: [gradient descent](#) in the [infinite-width limit](#) is fully equivalent to kernel gradient descent with the NTK. As a result, using gradient descent to minimize least-square loss for neural networks yields the same mean estimator as ridgeless kernel regression with the NTK. This duality enables simple [closed form](#) equations describing the training dynamics, [generalization](#), and predictions of wide neural networks.

The NTK was introduced in 2018 by Arthur Jacot, Franck Gabriel and Clément Hongler,^[1] who used it to study the convergence and generalization properties of fully connected neural networks. Later works^{[2][3]} extended the NTK results to other neural network architectures. In fact, the phenomenon behind NTK is not specific to neural networks and can be observed in generic nonlinear models, usually by a suitable scaling^[4].

Contradicting Machine Learning Folklore

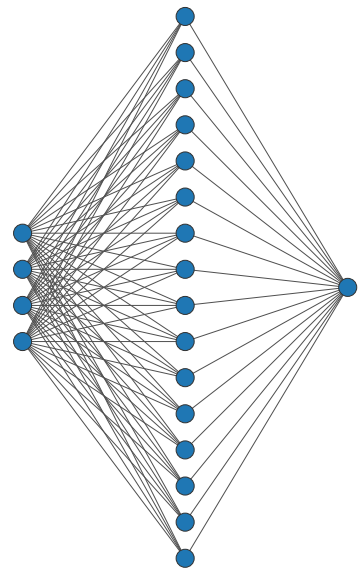
Folklore Theorem

Neural networks initialized with random parameters (with any law) converge to Gaussian processes in wide limits.

- Shallow neural networks with ReLU activation functions.
- Special kind of initialization of the network parameters.
- Finite-width networks are non-Gaussian processes.
- Gaussian and non-Gaussian processes in wide limits.

We provide a **complete characterization** of the statistical distribution of these neural network processes via their **characteristic functional**.

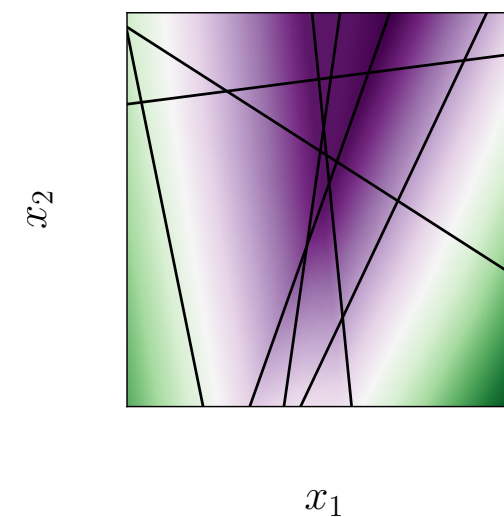
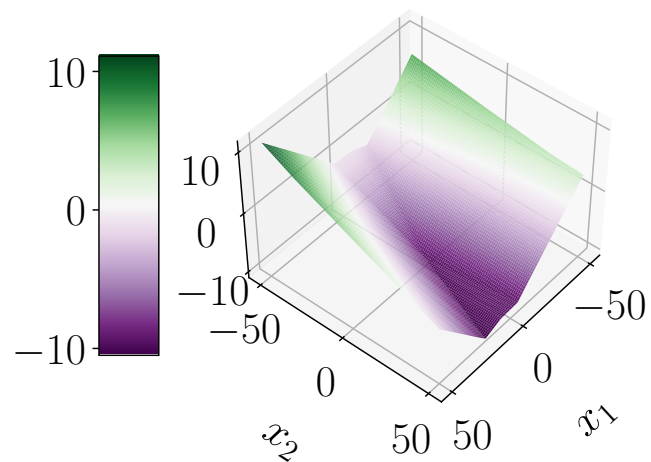
Shallow ReLU Neural Networks



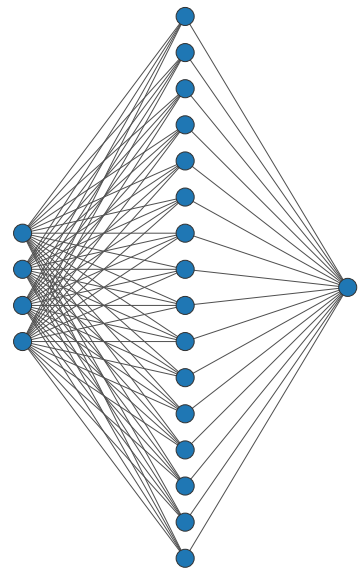
$$\mathbf{x} \mapsto \sum_{k=1}^K v_k \underbrace{(\mathbf{w}_k^T \mathbf{x} - b_k)_+}_{\text{ReLU neurons}}$$

K is the “width”
of the network

ReLU network



Shallow ReLU Neural Networks



$$\theta \sim Q$$

$$f_{\theta}(x) = \sum_{k=1}^K v_k (\mathbf{w}_k^{\top} x - b_k)_+$$

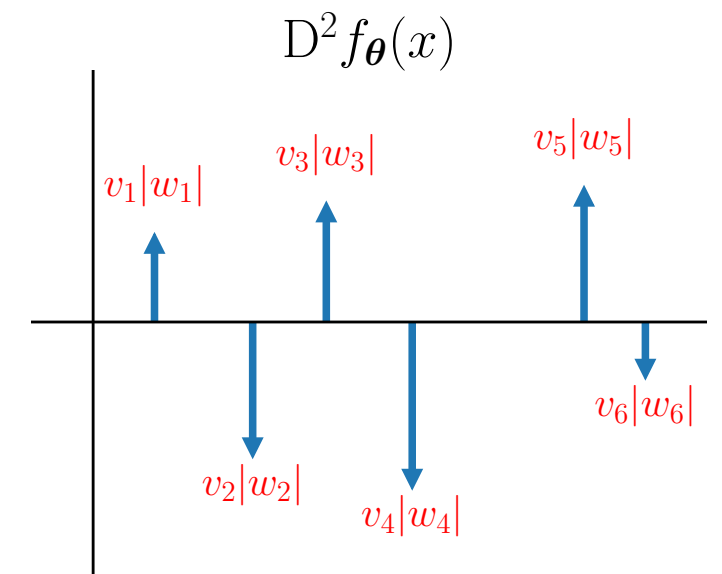
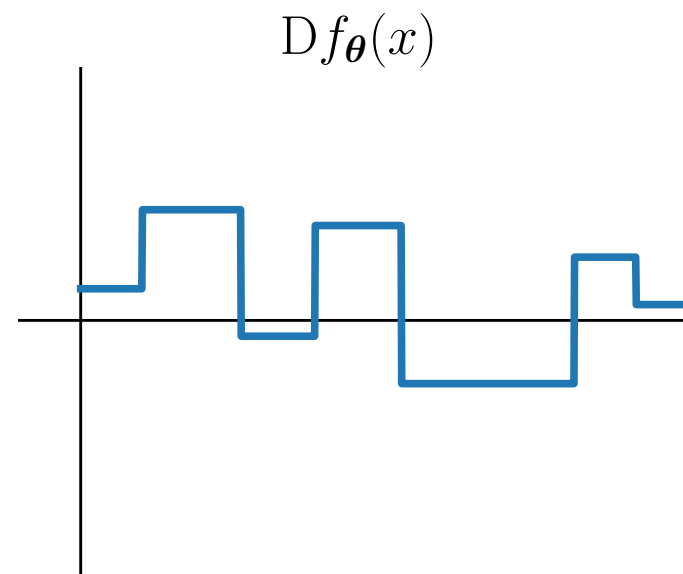
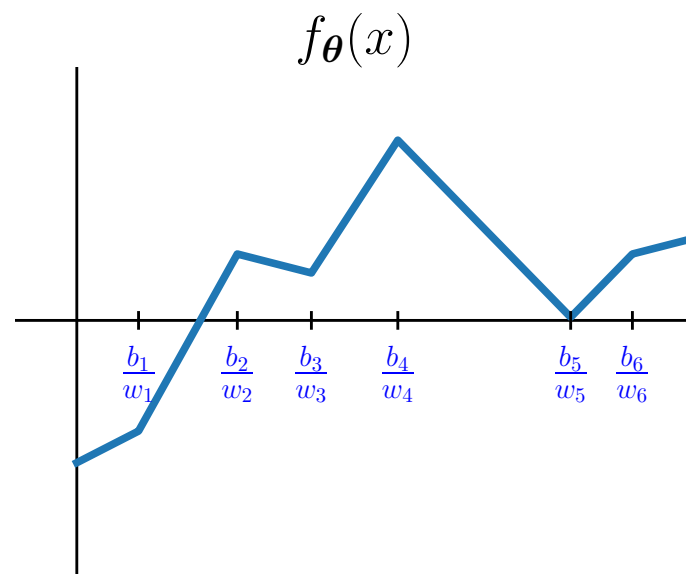
How do we study these random neural networks?

How do we derive a form of its **characteristic functional**?

We need to find a way to **whiten** the neural network.

Univariate Shallow ReLU Neural Networks

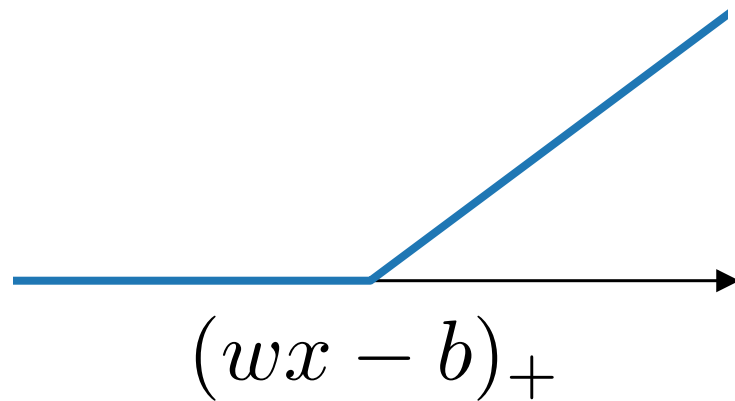
$$f_{\theta}(x) = \sum_{k=1}^K v_k (w_k x - b_k)_+$$



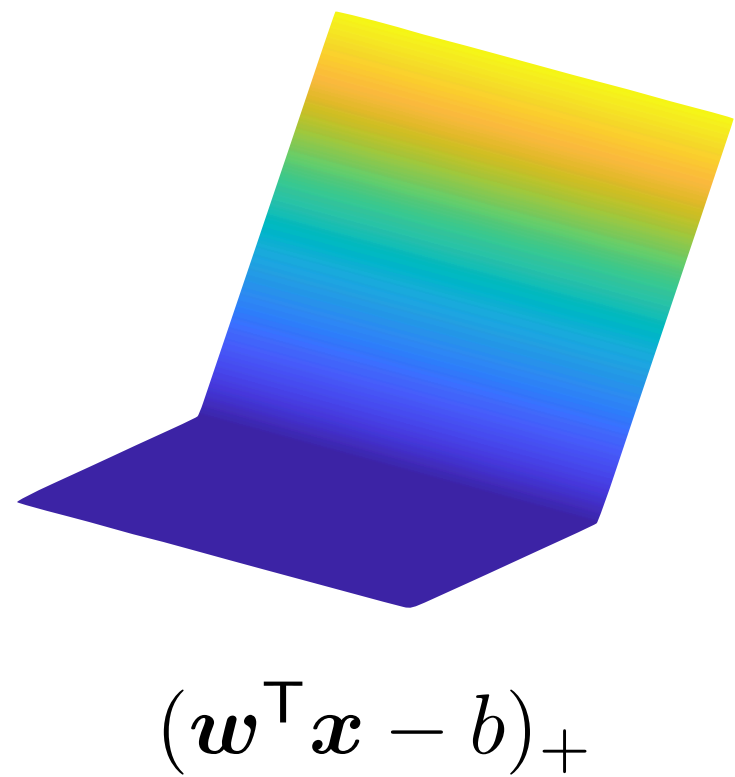
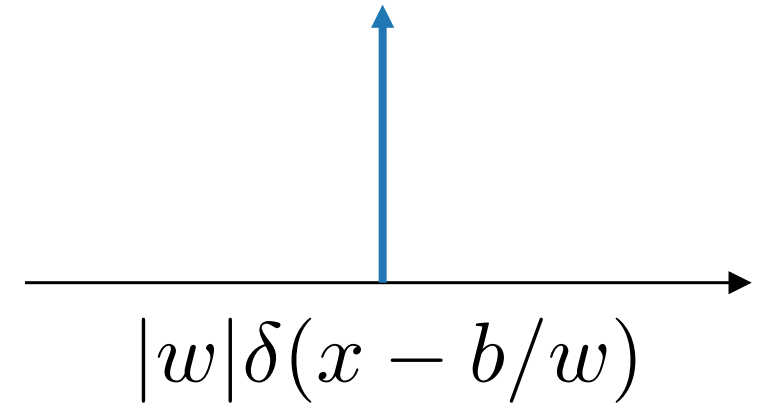
“looks like” a Poisson white noise

Second derivatives “whiten” univariate neural networks

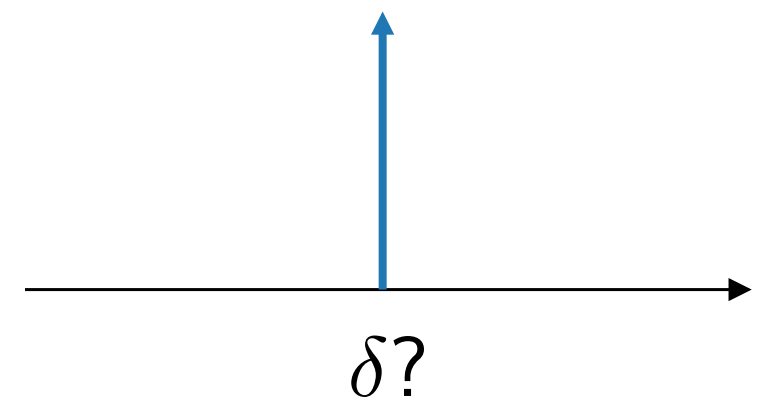
What About the Multivariate Case?



D^2



???



Whitening Operator of Shallow ReLU Networks

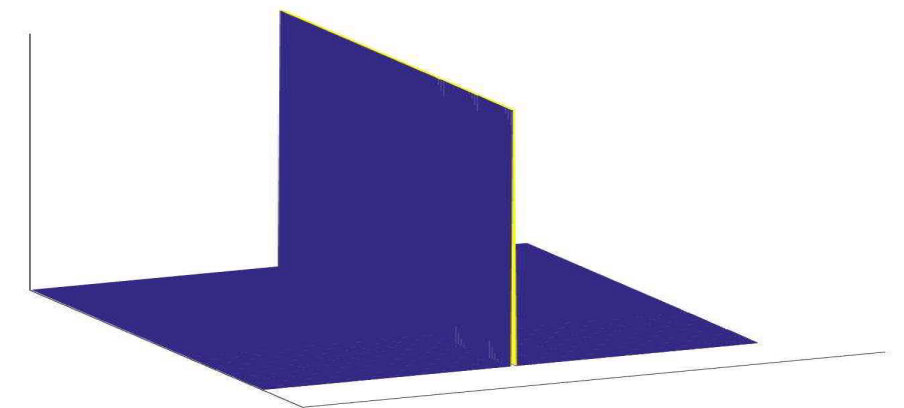
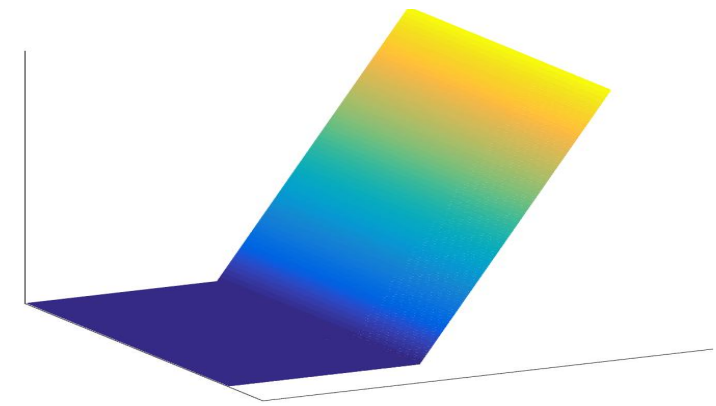
$$\begin{aligned} \text{ReLU Neuron: } & (\boldsymbol{w}^\top \boldsymbol{x} - b)_+ \\ = & \|\boldsymbol{w}\|_2 \left(\left[\frac{\boldsymbol{w}}{\|\boldsymbol{w}\|_2} \right]^\top \boldsymbol{x} - \frac{b}{\|\boldsymbol{w}\|_2} \right)_+ \end{aligned}$$

assume $\boldsymbol{w} \in \mathbb{S}^{d-1}$

$$\Delta\{(\boldsymbol{w}^\top \boldsymbol{x} - b)_+\} = \delta(\boldsymbol{w}^\top \boldsymbol{x} - b)$$

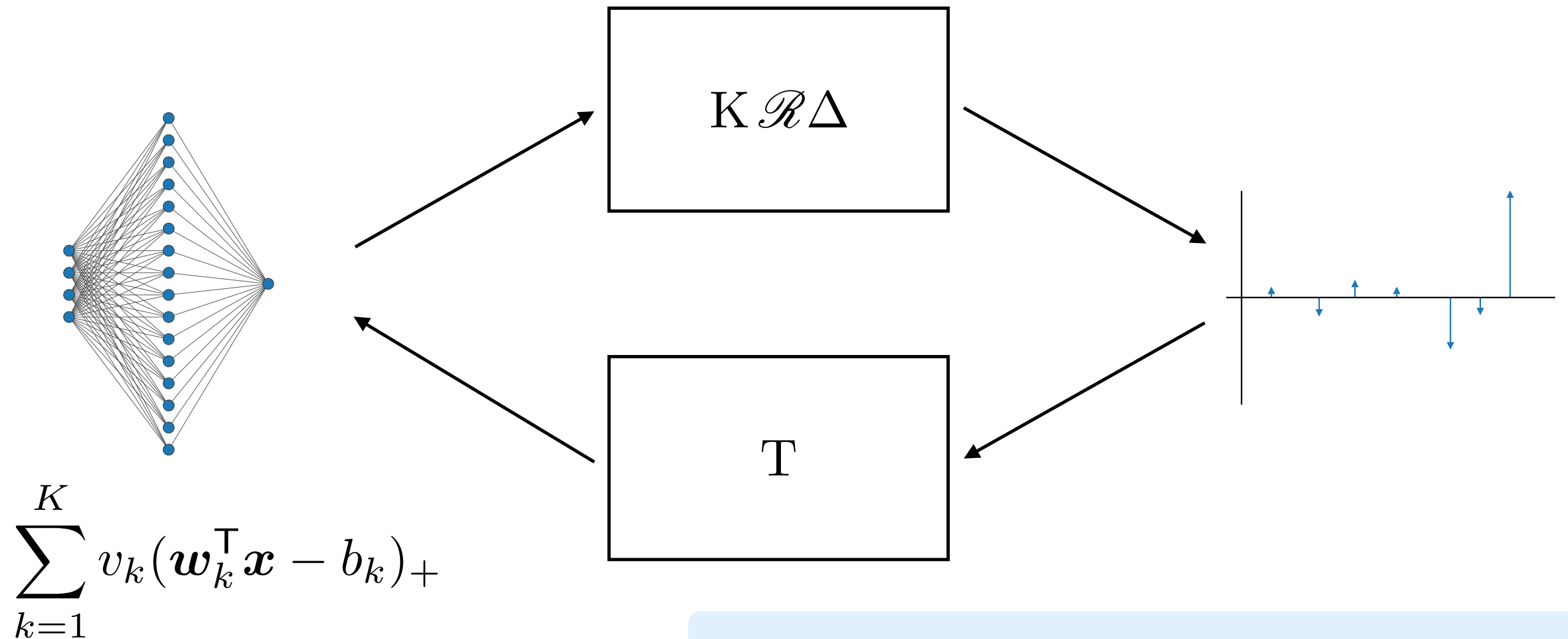
$$\mathbf{K} \mathcal{R} \Delta\{(\boldsymbol{w}^\top \boldsymbol{x} - b)_+\} = \delta_{(\boldsymbol{w}, b)}^e$$

$$\widehat{\mathbf{K}g}(\boldsymbol{\omega}) \propto |\boldsymbol{\omega}|^{d-1} \widehat{g}(\boldsymbol{\omega})$$



Radon-domain Dirac
centered at (\boldsymbol{w}, b)

Whitening Operator of Shallow ReLU Networks



Does such a well-behaved inverse exist?

$$T\{\psi\}(\mathbf{x}) = \int_{\mathbb{S}^{d-1} \times \mathbb{R}} k_{\mathbf{x}}(\mathbf{u}, t) \psi(\mathbf{u}, t) d\mathbf{u} dt$$

$$k_{\mathbf{x}}(\mathbf{u}, t) = (\mathbf{u}^\top \mathbf{x} - t)_+ - \frac{(\mathbf{u}^\top \mathbf{x} - t)}{2} - \frac{|t|}{2} + \mathbf{u}^\top \mathbf{x} \operatorname{sgn}(t)$$

Random Generation of Shallow ReLU Networks

$$\sum_k v_k (\mathbf{w}_k^\top \mathbf{x} - b_k)_+$$

1. Generate (v_k, \mathbf{w}_k, b_k) according to some point process on $\mathbb{S}^{d-1} \times \mathbb{R}$.

$$w = \sum_k v_k \delta_{(\mathbf{w}_k, b_k)}$$

2. Compute the characteristic functional $\hat{\mathbf{P}}_w$ of this point process
3. The characteristic functional of the random neural network is $\hat{\mathbf{P}}_w(\mathbf{T}^* \varphi)$

Random Generation of Shallow ReLU Networks

$$\sum_k v_k (\mathbf{w}_k^\top \mathbf{x} - b_k)_+$$

The v_k are drawn i.i.d. with respect to \mathbf{P}_V . The (\mathbf{w}_k, b_k) are drawn such that

finite absolute moment

1. The *activation thresholds* are mutually independent.
2. The expectation of the number of thresholds that intersect a finite volume in \mathbb{R}^d is a constant proportional to a rate $\lambda > 0$.
3. For every finite volume in \mathbb{R}^d , the thresholds are i.i.d. and “uniform” in the volume.

(\mathbf{w}_k, b_k) form a homogeneous Poisson point process on $\mathbb{S}^{d-1} \times \mathbb{R}$

Write $s \sim \mathcal{RP}(\lambda; \mathbf{P}_V)$ (ReLU Process) to denote this randomness.

Characteristic Functional of ReLU Processes

Theorem (PBEP, 2024)

The characteristic functional of random ReLU neural network

$$s \sim \mathcal{RP}(\lambda; \mathbf{P}_V)$$

is given by

$$\varphi \in \mathcal{S}(\mathbb{R}^d)$$

$$\hat{\mathbf{P}}_s(\varphi) = \exp \left(\lambda \int_{\mathbb{R}} \int_{\mathbb{R}} \int_{\mathbb{S}^{d-1}} \left(e^{i v^T \{ \varphi \}(\mathbf{u}, t)} - 1 \right) d\mathbf{u} dt d\mathbf{P}_V(v) \right).$$

$\hat{\mathbf{P}}_s$ gives us the full statistical distribution of s .

The law \mathbf{P}_s of s is the inverse Fourier transform of $\hat{\mathbf{P}}_s$.

Properties of ReLU Processes

Theorem (PBEPU, 2024)

Let $s \sim \mathcal{RP}(\lambda; \mathbf{P}_V)$ and $V \sim \mathbf{P}_V$. Then,

- The process is non-Gaussian.
- The autocovariance of the process is
$$C_s(\mathbf{x}, \mathbf{y}) \propto \mathbf{E}[V^2] \left(\|\mathbf{x} - \mathbf{y}\|_2^3 - \|\mathbf{x}\|_2^3 - \|\mathbf{y}\|_2^3 + 3\mathbf{x}^\top \mathbf{y} (\|\mathbf{x}\|_2 + \|\mathbf{y}\|_2) \right).$$
- The process is isotropic.
- The process wide-sense self-similar with Hurst exponent $H = 3/2$ (i.e., s and $a^H s(\cdot/a)$, $a > 0$ have the same second-order statistics).

Asymptotics of ReLU Processes

$$s_\lambda \sim \mathcal{RP}(\lambda; \mathbf{P}_V)$$

λ controls the width of the random ReLU neural network

$$s_\lambda|_\Omega(\mathbf{x}) = \mathbf{w}_0^\top \mathbf{x} + b_0 + \sum_{k=1}^{N_{\lambda,\Omega}} v_k (\mathbf{w}_k^\top \mathbf{x} - b_k)_+, \quad N_{\lambda,\Omega} \text{ Poisson with mean } \propto \lambda|\Omega|$$

As $\lambda \rightarrow \infty$ the network s_λ is an **infinite-width** network.

\mathbf{P}_V	fixed λ	$\lambda \rightarrow \infty$
Gaussian	s_λ is non-Gaussian	s_∞ is Gaussian
Laplacian	s_λ is non-Gaussian	s_∞ is Laplacian
α -stable, $\alpha < 2$	s_λ is non-Gaussian	s_∞ is α -stable

Random neural networks are non-Gaussian processes **even in wide limits.**

short proof via characteristic functional +
Lévy-Fernique theorem

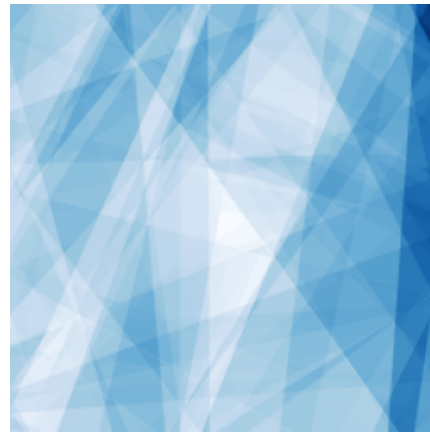
can check for non-Gaussianity by inspection

Pretty Pictures

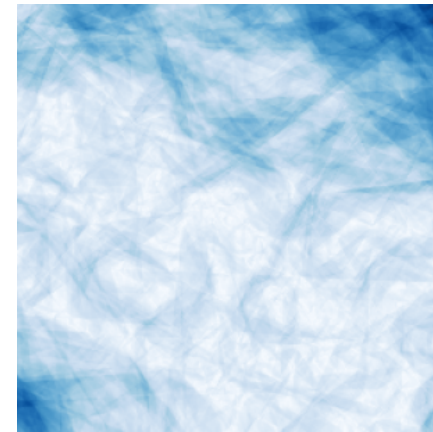
\mathbf{P}_V is Gaussian



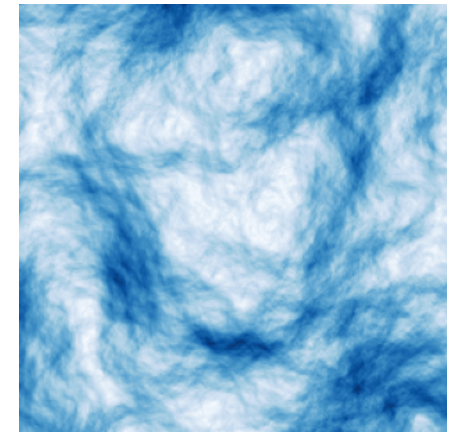
$\lambda = 1$



$\lambda = 10$



$\lambda = 100$

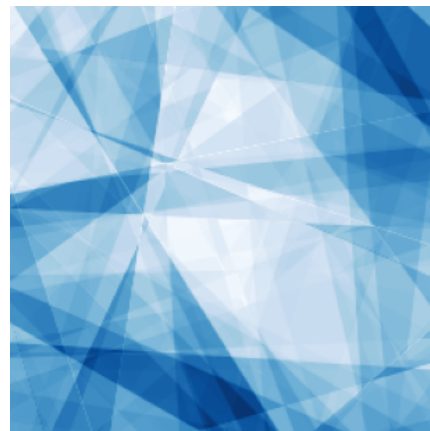


$\lambda = 1000$

\mathbf{P}_V is 1.25-stable



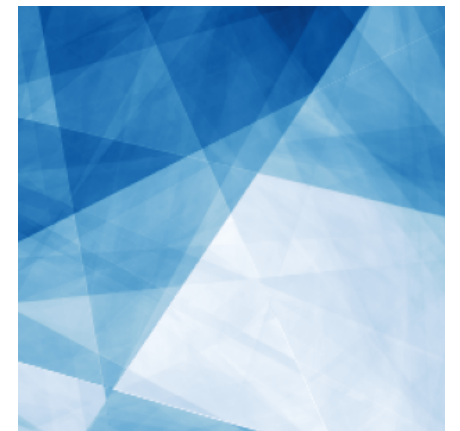
$\lambda = 1$



$\lambda = 10$



$\lambda = 100$



$\lambda = 1000$

Conclusion

Random shallow ReLU neural networks can be viewed as the deterministic “mixing” of a point process on the Radon domain

- Generic procedure to derive the characteristic functional of these random neural networks \Rightarrow law of the stochastic process
- The characteristic functional streamlines the derivation of properties (asymptotic and non-asymptotic) of stochastic processes
- We can prove things that “contradict” machine learning folklore

Next steps:

- Explore other Radon-domain point processes (and general Radon-domain innovation processes)
- Deep neural networks? Compositions of stochastic processes?

arXiv > stat > arXiv:2405.10229

Statistics > Machine Learning

[Submitted on 16 May 2024]

Random ReLU Neural Networks as Non-Gaussian Processes

Rahul Parhi, Pakshal Bohra, Ayoub El Biari, Mehrsa Pourya, Michael Unser