

A Data-Driven Market Simulator for Small Data Environments

Blanka Horvath
King's College London &
The Alan Turing Institute

DataSig Workshop
16th March, 2021

The work is available for download on SSRN under the link
A Data-driven Market Simulator for Small Data Environments
SSRN:3632431

The work is available for download on SSRN under the link
A Data-driven Market Simulator for Small Data Environments
SSRN:3632431

Motivation for our Market Generators

The concept of “Model” (here, in form of a numerical program):

- ▶ Classical: (Program; Data) \Rightarrow Output
- ▶ Now: (**Architecture, ObjF; TrainData**) \Rightarrow Program
(Program, TestData) \Rightarrow Output

\Rightarrow This may redefine the concept of model governance too:

Model = $\underbrace{(\text{Architecture, ObjF; Dataset})}_{\text{Network}}$ “Quality of” training data shapes the DNN!

See: Deep Hedging, **Bühler, Gonon, Teichmann, Wood (2019)**

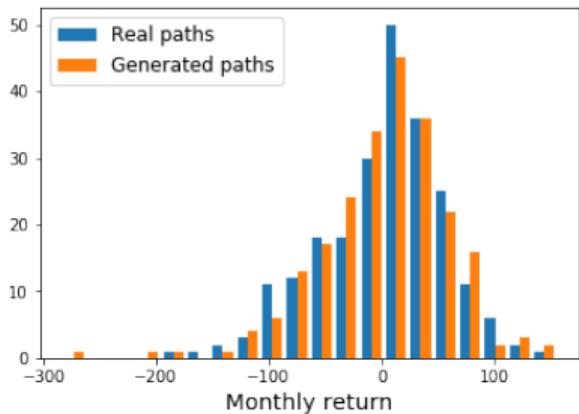
This talk is based on joint work with

Hans Bühler, Terry Lyons, Imanol Perez Arribas, Ben Wood



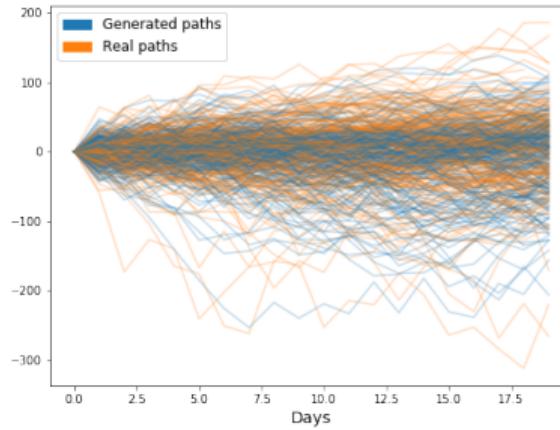
Market Simulation

Generation of Returns



vs.

Generation of paths



Market Simulation in Mathematical Finance

Transforming horizons of mathematical modelling towards models that can more and more accurately fit market data:

- (1) Classical and Neo-classical stochastic market models:

$$dS_t = rS_t dt + \sigma S_t dW_t, \quad S_0 = s.$$

Black-Scholes, diffusions, models with jumps, stochastic volatility, path-dependent models, LSV, multifactor models, Rough Volatility, ...

- (2) DNN-based Generative Modelling in other AI applications: Adapt to Finance?

Market Simulation in Mathematical Finance

Transforming horizons of mathematical modelling towards models that can more and more accurately fit market data:

- (1) Classical and Neo-classical stochastic market models:

$$dS_t = rS_t dt + \sigma S_t dW_t, \quad S_0 = s.$$

Black-Scholes, diffusions, models with jumps, stochastic volatility, path-dependent models, LSV, multifactor models, Rough Volatility, ...

- (2) DNN-based Generative Modelling in other AI applications: Adapt to Finance?
Data driven (no a-priori assumption on distribution of stochastic process)

Market Simulation in Mathematical Finance

Transforming horizons of mathematical modelling towards models that can more and more accurately fit market data:

- (1)** Classical and Neo-classical stochastic market models:

$$dS_t = rS_t dt + \sigma S_t dW_t, \quad S_0 = s.$$

Black-Scholes, diffusions, models with jumps, stochastic volatility, path-dependent models, LSV, multifactor models, Rough Volatility, ...

- (2)** DNN-based Generative Modelling in other AI applications: Adapt to Finance?
Data driven (no a-priori assumption on distribution of stochastic process) non-parametric?

$$f \in \mathcal{N}_r(I, d_1, \dots, d_{r-1}, O; \sigma_1, \dots, \sigma_r)$$

⇒ very flexible. Originally developed for static problems, adaptation to financial time-series modelling not straightforward (see later).

Market Simulation in Mathematical Finance

Transforming horizons of mathematical modelling towards models that can more and more accurately fit market data:

- (1)** Classical and Neo-classical stochastic market models:

$$dS_t = rS_t dt + \sigma S_t dW_t, \quad S_0 = s.$$

Black-Scholes, diffusions, models with jumps, stochastic volatility, path-dependent models, LSV, multifactor models, Rough Volatility, ...

- (2)** DNN-based Generative Modelling in other AI applications: Adapt to Finance?
Data driven (no a-priori assumption on distribution of stochastic process) non-parametric?

$$f \in \mathcal{N}_r(I, d_1, \dots, d_{r-1}, O; \sigma_1, \dots, \sigma_r)$$

⇒ very flexible. Originally developed for static problems, adaptation to financial time-series modelling not straightforward (see later). Interpretability, risk-management...

Market Simulation in Mathematical Finance

Transforming horizons of mathematical modelling towards models that can more and more accurately fit market data:

- (1)** Classical and Neo-classical stochastic market models:

$$dS_t = rS_t dt + \sigma S_t dW_t, \quad S_0 = s.$$

Black-Scholes, diffusions, models with jumps, stochastic volatility, path-dependent models, LSV, multifactor models, Rough Volatility, ...

- (2)** DNN-based Generative Modelling in other AI applications: Adapt to Finance?
Data driven (no a-priori assumption on distribution of stochastic process) non-parametric?

$$f \in \mathcal{N}_r(I, d_1, \dots, d_{r-1}, O; \sigma_1, \dots, \sigma_r)$$

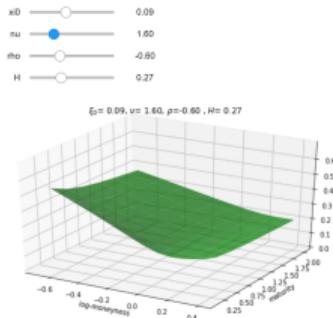
⇒ very flexible. Originally developed for static problems, adaptation to financial time-series modelling not straightforward (see later). Interpretability, risk-management...

- (1.5)** Extending (or augmenting) currently prevalent stochastic models: Models that are more adaptive to market environments by creating mixtures of (neo-)classical models

Market Simulation in Mathematical Finance

Transforming horizons of mathematical modelling towards models that can more and more accurately fit market data:

- (1.5)** Augmenting and extending currently prevalent stochastic models
models that are adaptive to market environments by mixing of (neo-)classical models
- ▶ **Mixture models:** A first step towards this was demonstrated in the Deep Learning Volatility framework: Take a mixture of two (or more) stochastic volatility models and calibrate it to data including the mixture parameter a .
$$a \times \text{Heston} + (1 - a) \times \text{rBergomi}$$



Generative Models in Machine Learning

(2) DNN-based Generative Modelling in other AI applications: Adapt to Finance?

Generative Models in Machine Learning

(2) DNN-based Generative Modelling in other AI applications: Adapt to Finance?

Generative models can be trained to transport **some source distribution** μ

Generative Models in Machine Learning

(2) DNN-based Generative Modelling in other AI applications: Adapt to Finance?

Generative models can be trained to transport **some source distribution** μ (say $\mu \sim \mathcal{U}[0, 1]$ or $\mu \sim \mathcal{N}(0, 1)$) to some **target distribution** observed in the data

Generative Models in Machine Learning

(2) DNN-based Generative Modelling in other AI applications: Adapt to Finance?

Generative models can be trained to transport **some source distribution** μ (say $\mu \sim \mathcal{U}[0, 1]$ or $\mu \sim \mathcal{N}(0, 1)$) to some **target distribution** observed in the data and generate more samples that are similar/indistinguishable from the ones observed (Data-driven as there are no assumptions made on the latter distribution).

Generative Models in Machine Learning

(2) DNN-based Generative Modelling in other AI applications: Adapt to Finance?

Generative models can be trained to transport **some source distribution** μ (say $\mu \sim \mathcal{U}[0, 1]$ or $\mu \sim \mathcal{N}(0, 1)$) to some **target distribution** observed in the data and generate more samples that are similar/indistinguishable from the ones observed (Data-driven as there are no assumptions made on the latter distribution).

Currently the most popular DNN-based generative models

- ▶ Restricted Boltzman Machine (RBM)

Generative Models in Machine Learning

(2) DNN-based Generative Modelling in other AI applications: Adapt to Finance?

Generative models can be trained to transport **some source distribution** μ (say $\mu \sim \mathcal{U}[0, 1]$ or $\mu \sim \mathcal{N}(0, 1)$) to some **target distribution** observed in the data and generate more samples that are similar/indistinguishable from the ones observed (Data-driven as there are no assumptions made on the latter distribution).

Currently the most popular DNN-based generative models

- ▶ Restricted Boltzman Machine (RBM) Kondratyev & Schwarz (2019)
- ▶ Generative Adverserial Networks (GAN)

Generative Models in Machine Learning

(2) DNN-based Generative Modelling in other AI applications: Adapt to Finance?

Generative models can be trained to transport **some source distribution** μ (say $\mu \sim \mathcal{U}[0, 1]$ or $\mu \sim \mathcal{N}(0, 1)$) to some **target distribution** observed in the data and generate more samples that are similar/indistinguishable from the ones observed (Data-driven as there are no assumptions made on the latter distribution).

Currently the most popular DNN-based generative models

- ▶ Restricted Boltzman Machine (RBM) Kondratyev & Schwarz (2019)
- ▶ Generative Adverserial Networks (GAN) several authors
- ▶ Variational Autoencoders (VAE)

Generative Models in Machine Learning

(2) DNN-based Generative Modelling in other AI applications: Adapt to Finance?

Generative models can be trained to transport **some source distribution** μ (say $\mu \sim \mathcal{U}[0, 1]$ or $\mu \sim \mathcal{N}(0, 1)$) to some **target distribution** observed in the data and generate more samples that are similar/indistinguishable from the ones observed (Data-driven as there are no assumptions made on the latter distribution).

Currently the most popular DNN-based generative models

- ▶ Restricted Boltzman Machine (RBM) Kondratyev & Schwarz (2019)
- ▶ Generative Adverserial Networks (GAN) several authors
- ▶ Variational Autoencoders (VAE) today's presentation

Bühler, H., Lyons, Perez-Arribas, Wood (2020)

Generative Models in Machine Learning

(2) DNN-based Generative Modelling in other AI applications: Adapt to Finance?

Generative models can be trained to transport **some source distribution** μ (say $\mu \sim \mathcal{U}[0, 1]$ or $\mu \sim \mathcal{N}(0, 1)$) to some **target distribution** observed in the data and generate more samples that are similar/indistinguishable from the ones observed (Data-driven as there are no assumptions made on the latter distribution).

Currently the most popular DNN-based generative models

- ▶ Restricted Boltzman Machine (RBM) Kondratyev & Schwarz (2019)
- ▶ Generative Adverserial Networks (GAN) several authors
- ▶ Variational Autoencoders (VAE) today's presentation

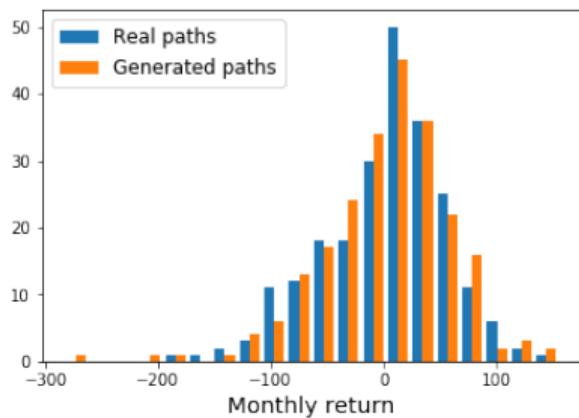
Bühler, H., Lyons, Perez-Arribas, Wood (2020)

Data driven, flexible, originally developed for static problems, adaptation to financial time-series modelling not straightforward.

- ▶ One approach to the incorporation of time-series aspect is by Causal Optimal Transport (see ongoing work by B. Acciaio, T. Xu and collaborators)
- ▶ In this work we take an approach via **Rough Paths** using **Signatures**.

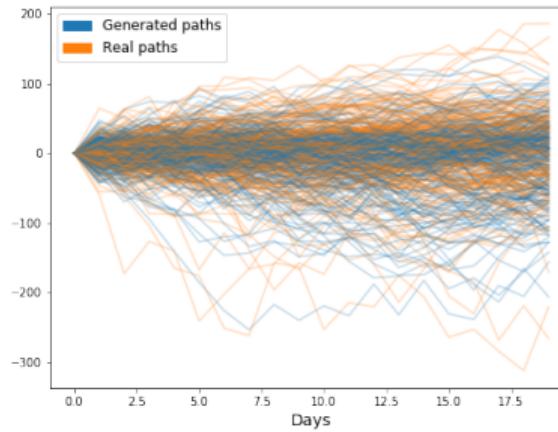
Market Simulation by Signatures

Returns-based generation (i)



vs.

Path-based generation (ii)



General Strategy

(Step 1) **Data extraction from time series** we subdivide original (say daily S&P index) data into partitions of: (1) daily data, (2) weekly path segments, i.e. 5 days, and (3) monthly path segments, i.e. 20 days.

General Strategy

- (Step 1) **Data extraction from time series** we subdivide original (say daily S&P index) data into partitions of: (1) daily data, (2) weekly path segments, i.e. 5 days, and (3) monthly path segments, i.e. 20 days.
- (Step 2) **Preprocessing the data** transforming data into (i) returns for all (1), (2), (3) and into (ii) log-signatures for (2), (3) of (leadlag) paths.

General Strategy

- (Step 1) **Data extraction from time series** we subdivide original (say daily S&P index) data into partitions of: (1) daily data, (2) weekly path segments, i.e. 5 days, and (3) monthly path segments, i.e. 20 days.
- (Step 2) **Preprocessing the data** transforming data into (i) returns for all (1), (2), (3) and into (ii) log-signatures for (2), (3) of (leadlag) paths.
- (Step 3) **Creating and training the VAE and the CVAE network:** VAE, a parsimonious generator model with “bottleneck structure”. Conditional Variational Autoencoder is learned to condition VAE on current market conditions (a) current level of the index (b) instantaneous volatility (c) signature of the previous path segment.

General Strategy

- (Step 1) **Data extraction from time series** we subdivide original (say daily S&P index) data into partitions of: (1) daily data, (2) weekly path segments, i.e. 5 days, and (3) monthly path segments, i.e. 20 days.
- (Step 2) **Preprocessing the data** transforming data into (i) returns for all (1), (2), (3) and into (ii) log-signatures for (2), (3) of (leadlag) paths.
- (Step 3) **Creating and training the VAE and the CVAE network:** VAE, a parsimonious generator model with “bottleneck structure”. Conditional Variational Autoencoder is learned to condition VAE on current market conditions (a) current level of the index (b) instantaneous volatility (c) signature of the previous path segment.
- (Step 4) **Postprocessing of the outputs of the VAEs**
transforming (i), (ii) back to paths; such as building paths of arbitrary length

General Strategy

- (Step 1) **Data extraction from time series** we subdivide original (say daily S&P index) data into partitions of: (1) daily data, (2) weekly path segments, i.e. 5 days, and (3) monthly path segments, i.e. 20 days.
- (Step 2) **Preprocessing the data** transforming data into (i) returns for all (1), (2), (3) and into (ii) log-signatures for (2), (3) of (leadlag) paths.
- (Step 3) **Creating and training the VAE and the CVAE network:** VAE, a parsimonious generator model with “bottleneck structure”. Conditional Variational Autoencoder is learned to condition VAE on current market conditions (a) current level of the index (b) instantaneous volatility (c) signature of the previous path segment.
- (Step 4) **Postprocessing of the outputs of the VAEs**
transforming (i), (ii) back to paths; such as building paths of arbitrary length
- (Step 5) **Performance evaluation** similar to the role of the discriminator in GANs, but here without feeding back to the generator.

General Strategy

- (Step 1) **Data extraction from time series** we subdivide original (say daily S&P index) data into partitions of: (1) daily data, (2) weekly path segments, i.e. 5 days, and (3) monthly path segments, i.e. 20 days.
- (Step 2) **Preprocessing the data** transforming data into (i) returns for all (1), (2), (3) and into (ii) log-signatures for (2), (3) of (leadlag) paths.
- (Step 3) **Creating and training the VAE and the CVAE network:** VAE, a parsimonious generator model with “bottleneck structure”. Conditional Variational Autoencoder is learned to condition VAE on current market conditions (a) current level of the index (b) instantaneous volatility (c) signature of the previous path segment.
- (Step 4) **Postprocessing of the outputs of the VAEs** transforming (i), (ii) back to paths; such as building paths of arbitrary length
- (Step 5) **Performance evaluation** similar to the role of the discriminator in GANs, but here without feeding back to the generator. **Good performance evaluation metrics?**

Rough Paths Approach to Generative Modelling of Markets

Levin, Lyons, & Ni. (2013) firstly proposed the signature of a path as a basis of functions for a functional on path space.

Definition (Signature of a path)

Let $X : [0, T] \rightarrow \mathbb{R}^d$ be a continuous path of bounded variation. The signature of X is then defined by the sequence of iterated integrals given by

$$\mathbb{X}_T^{<\infty} := (1, \mathbb{X}_t^1, \dots, \mathbb{X}_T^n, \dots), \quad \text{where}$$

$$\mathbb{X}_T^n := \int_{0 < u_1 < \dots < u_k < T} dX_{u_1} \otimes \dots \otimes dX_{u_k} \in (\mathbb{R}^d)^{\otimes n}$$

with \otimes the tensor product. Similarly, given $N \in \mathbb{N}$, the truncated signature of order N is defined by

$$\mathbb{X}_T^{< N} := (1, \mathbb{X}_T^1, \dots, \mathbb{X}_T^N).$$

If the path X has bounded variation – which is the case of discrete data – the integrals above can be defined using Riemann-Stieltjes integrals.

Rough Paths Approach to Generative Modelling

Why signatures?

Rough Paths Approach to Generative Modelling

Why signatures? Signatures provide a basis of functions for a functional on path space.

- ▶ While Fourier transforms and wavelets have a similar role approximating curves as a linear combination of basis functions, signatures do so in an un-parametrised way

Rough Paths Approach to Generative Modelling

Why signatures? Signatures provide a basis of functions for a functional on path space.

- ▶ While Fourier transforms and wavelets have a similar role approximating curves as a linear combination of basis functions, signatures do so in an un-parametrised way (model free, path by path characterisation possible).
- ▶ Robustness to missing data and irregular sampling and

Rough Paths Approach to Generative Modelling

Why signatures? Signatures provide a basis of functions for a functional on path space.

- ▶ While Fourier transforms and wavelets have a similar role approximating curves as a linear combination of basis functions, signatures do so in an un-parametrised way (model free, path by path characterisation possible).
- ▶ Robustness to missing data and irregular sampling and towards highly oscillatory data, and

Rough Paths Approach to Generative Modelling

Why signatures? Signatures provide a basis of functions for a functional on path space.

- ▶ While Fourier transforms and wavelets have a similar role approximating curves as a linear combination of basis functions, signatures do so in an un-parametrised way (model free, path by path characterisation possible).
- ▶ Robustness to missing data and irregular sampling and towards highly oscillatory data, and invariance under time re-parametrisation. **Liao, Lyons, Ni, Yang (2019)**
- ▶ Signatures provide the right framework for performance evaluation metrics on pathspace, one of the difficulties being that the pathspace $C([0, 1], \mathbb{R}^d)$ is infinite-dimensional and not locally compact. **Chevrev, Oberhauser (2018)**
- ▶ In fact, **Liao, Lyons, Ni, Yang (2019)** demonstrated the advantages of log-signatures (with all positive properties listed above, but lower dimensionality)

The Signature MMD Two-Sample Test

To assess whether a generative model is able to generate paths that are realistic with respect to a sample of real paths Y_1, \dots, Y_n , we sample from the generative model $n \in \mathbb{N}$, paths X_1, \dots, X_n and we apply the two-sample test proposed in [?]. More specifically, we compute the signature-based MMD test statistic
 $T(X_1, \dots, X_n; Y_1, \dots, Y_n)$

$$T(X_1, \dots, X_n; Y_1, \dots, Y_n) := \frac{1}{n(n-1)} \sum_{i,j;i \neq j} k(X_i, X_j) - \frac{2}{n^2} \sum_{i,j} k(X_i, Y_j) + \frac{1}{n(n-1)} \sum_{i,j;i \neq j} k(Y_i, Y_j), \quad (1)$$

where $k(\cdot, \cdot)$ is the so-called *signature kernel*. Then, given a fixed confidence level $\alpha \in (0, 1)$, we compute the threshold $c_\alpha := 4\sqrt{-n^{-1} \log \alpha}$. The generative model will be said to be realistic with a confidence α if $T_U^2 < c_\alpha$.

Rough Paths Approach to Generative Modelling

In what follows we work with the log-signatures (**Liao, Lyons, Ni, Yang (2019)**)

Definition (Log-signature)

Let $X : [0, T] \rightarrow \mathbb{R}^d$ be a path such that its signature $\mathbb{X}_{0,T}^{<\infty}$ is well-defined. The log-signature is then defined by

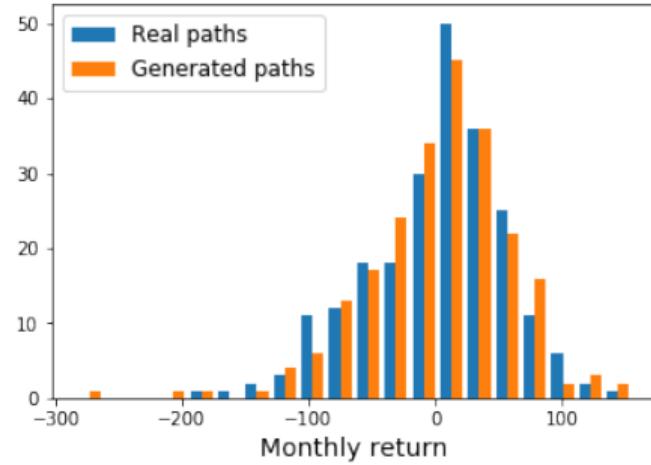
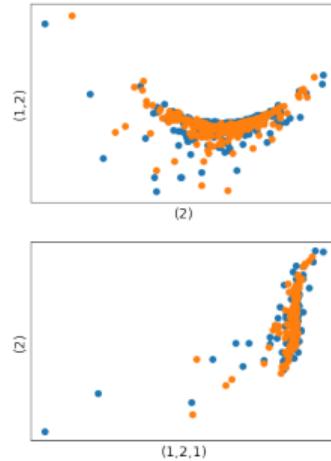
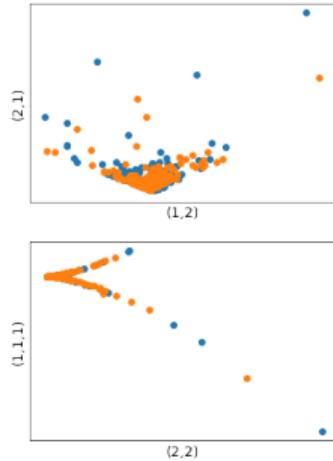
$$\log \mathbb{X}_T^{<\infty} := -\mathbb{X}_T^{<\infty} + \frac{1}{2}(\mathbb{X}_T^{<\infty})^{\otimes 2} - \frac{1}{3}(\mathbb{X}_T^{<\infty})^{\otimes 3} + \dots + (-1)^n \frac{1}{n}(\mathbb{X}_T^{<\infty})^{\otimes n} + \dots,$$

which can be shown to be well-defined.

- ▶ There is a one-to-one map between signatures and log-signatures.
- ▶ Log-signatures have all positive properties listed above.
- ▶ They allow for lower dimensional representation and are better suited to VAE.

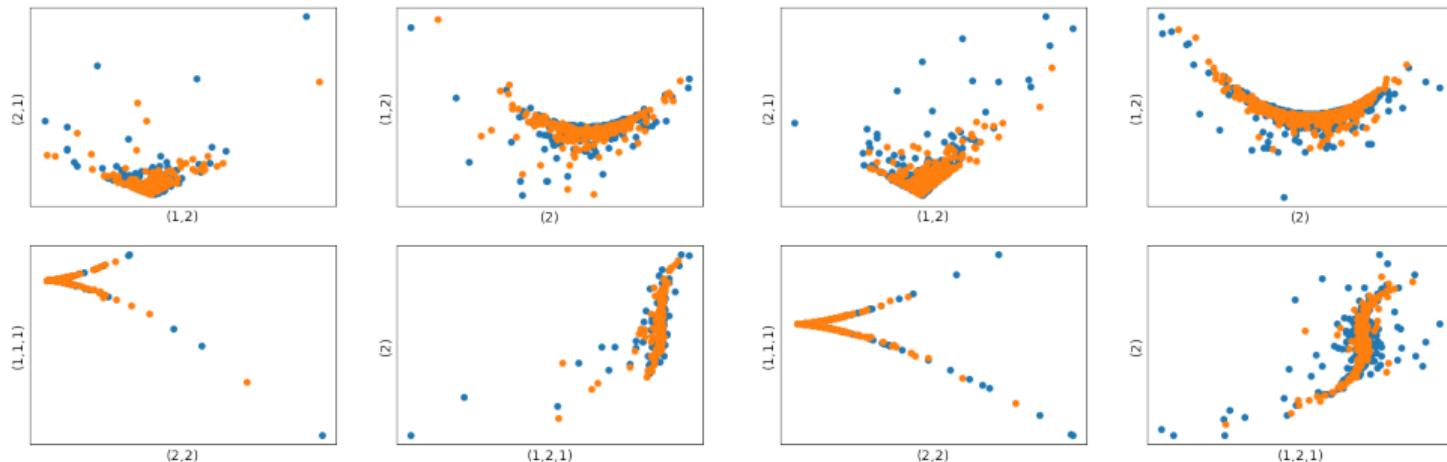
Numerical Results

- Signature based training (ii) at least as accurate as returns-based training (i).



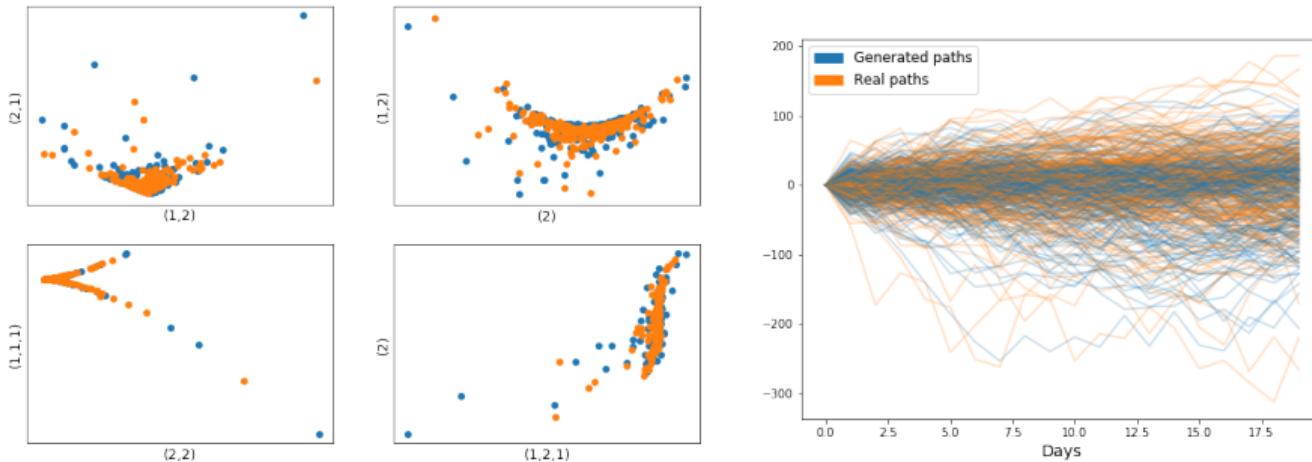
Numerical Results

- ▶ Signature-based generative model works well already with the low number of training samples available in data (1000 samples weekly-unconditioned **left image**, 250 samples monthly **right image** unconditioned).
- ▶ Increasing the number of samples in numerically generated data does not lead to better performance.



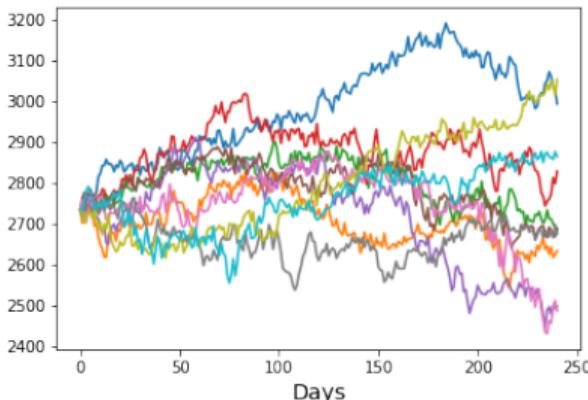
Numerical Results

- ▶ Signature-based generative model works well already with the low number of training samples available in data (1000 samples weekly-unconditioned, 250 samples monthly unconditioned).
- ▶ Increasing the number of samples in numerically generated data does not lead to better performance.



Conditional VAE and Postprocessing

- ▶ Conditional Variational Autoencoder is learned to condition VAE on current market conditions (a) current level of the index (b) instantaneous volatility (c) signature of the previous path segment.
- ▶ By conditioning, the number of available training samples is even smaller ⇒ powerful parsimonious generative model crucial.
- ▶ By conditioning on the signature of the previous path segment, one can generate + build paths far longer than the direct output of the generative model. Numerical results consistent between weekly vs. monthly paths. **Below: Yearly paths**



Thank you for your attention.