

# Scenario Trees – A Process Distance Approach

Raimund M. Kovacevic<sup>a,b</sup>, Alois Pichler<sup>a,\*</sup>

<sup>a</sup>University of Vienna, Austria. Department of Statistics and Operations Research .

<sup>b</sup>Funded by WWTF.

---

## Abstract

The approximation of stochastic processes by trees is an important topic in multistage stochastic programming. In this paper we focus on improving the approximation of large trees by smaller (tractable) trees. The quality of the approximation is measured by the nested distance, recently introduced in [Pflug, Pf09]. The nested distance is derived from the Wasserstein distance. It additionally takes into account the effect of information, which is increasing over time.

After discussing the basic relations between processes and trees and reviewing the nested distance we introduce and analyze an algorithm for finding good approximations. The algorithm, step by step, improves the probabilities on a tree, and also improves the paths. For the important case of quadratic nested distances the algorithm, generalizing multistage, k-means clustering, finds locally best approximating trees in finitely many iterations.

*Keywords:* Stochastic processes and trees, Wasserstein and Kantorovich Distance, Tree approximation, Optimal transport, Facility location

*2010 MSC:* 90C15, 60B05, 90-08

---

## 1. Introduction

Stochastic programming is an important methodology for decision making under uncertainty. Typical formulations may involve minimizing expected loss or maximizing expected profit and can be extended by risk averse optimization (e.g. using expected utility or coherent risk measures). In addition it is possible to deal with stochastic constraints e.g. by recourse formulations or by using probabilistic constraints. See [SDR09] for a recent and comprehensive treatment of the main techniques.

An important subset of stochastic programs is given by multistage stochastic programs, where decisions have to be taken at multiple points in time. As a simple example we sketch the problem

$$\begin{aligned} & \text{maximize} && \mathbb{E} H(\xi, x) \\ & \text{(in } x) && \\ & \text{subject to} && x_t \in \mathbb{X}_t \quad t \in \{0, \dots, T\}, \\ & && x_t \text{ is measurable with respect to } \mathcal{F}_t, \end{aligned} \tag{1}$$

which uses a real valued profit functions  $H: \Xi \times \mathcal{X} \rightarrow \mathbb{R}$ , a stochastic process  $(\xi_t)_t$ , and expectation utility in the objective. The decision vector  $x = (x_0, x_1, \dots, x_T)$  models decisions at points  $t \in \{0, \dots, T\}$  in time. Typical constraints are expressed by equations and inequalities and e.g. model budgeting, bookkeeping or storage constraints and any lower and upper bounds on the decisions ( $x_t \in \mathbb{X}_t$ ). The measurability constraints refer to a filtration, modeling the increase in information over time, and express the fact that decisions have to be taken without knowing the future (*nonanticipativity*). Because the decisions  $x$  also form a stochastic processes, optimization must be done in function spaces. Unfortunately, only in rare cases an analytic solution can be given.

---

\*Corresponding Author

URL: <http://isor.univie.ac.at/> (Alois Pichler )

In stochastic programming the underlying processes are typically replaced by approximations: The process  $\xi$  is replaced by a finitely valued stochastic scenario process  $\xi'$  and the decisions  $x_t$  are replaced by (often high dimensional) vectors  $x'_t$ . In order to model nonanticipativity, it is assumed that the decisions  $x'$  are adapted to some suitable filtration  $\mathcal{F}'$ , related to the discretized process  $\xi'$ . This filtration usually is modeled by a tree structure, such that all the relevant information (values, probabilities, decisions) is related to nodes in the tree.

In a tree-based framework it is possible to rewrite problem (1) in straightforward manner

$$\begin{aligned} & \text{maximize} && \sum_{i>0} p_i H(x_{i-}, \xi_i) \\ & \text{(in } x) && \\ & \text{subject to} && x_i \in \mathbb{X}_i. \end{aligned}$$

The notation will be clarified later on (Section 4), but it is important to note that the random process  $\xi$  and the random decision process  $x$  is now modeled just by real values, sitting on the nodes  $i$  of the tree. This means that optimization can be done numerically with  $x_i \in \mathbb{R}^n$  for some  $n$ .

From this sketch it should be clear that the construction of trees – in fact the approximation of a process by a tree – is an important topic in multistage stochastic programming. Different approaches have been used in literature. Besides just simulating trees by Monte Carlo simulation, the most popular approach consists in constructing trees such that the conditional moments (up to some order) of the tree are close to the conditional moments of the real process (moment method). This approach was generalized in [HW01] (cf. also [Kla02]) by minimizing the Euclidean distance between whole collections of statistical properties. Important work also was done on using probability metrics, e.g. based on the Wasserstein/ Kantorovich-distances [DGKR03]. The Wasserstein distance is a transportation distance intending to minimize total costs that have to be taken into account when passing from a given distribution to a desired one. Other concepts of distances for multistage stochastic programming emphasize the role of filtrations and use distances between filtrations, as introduced in [Boy71], see also [Kud74], [HR09]. Recently, [Pfl09] proposed a new type of distance – the nested, or process distance – that extends in a natural way the Wasserstein distance between probability distributions to a distance between processes. Both aspects, the distributional and the filtration, are accounted for by this measure.

The present paper aims at using the process distance for tree construction. We will focus on improving the distance between a given, big scenario tree (constructed e.g. by simulation or any other means) and a smaller scenario tree, suitable for solving a stochastic multistage optimization problem. While this can be done in a relative simple way when using the Wasserstein distance, approximations and iterative approaches for finding local optima have to be used in the tree-case. We introduce suitable algorithms, discuss their properties and provide numerical examples.

In order to clarify the notation the concept of trees and their link to stochastic processes and filtrations are reviewed in Section 2. The Wasserstein distance – as the most similar “ordinary” probability metric – and its key properties as regards the approximation quality are discussed in Section 3. This is the basis to introduce the nested distance in Section 4 and elaborate in Section 5 how to improve the values and the probabilities within a given tree structures in order to improve the approximation quality.

## 2. Trees and filtrations

When stating stochastic optimization problems it is often advantageous to use the notion of stochastic processes and filtered probability spaces to describe the objects being studied. This is, however, not adequate when implementing concrete realizations in computer models. Typically the models are reformulated in terms of finite state spaces. For multistage stochastic decision problems the basic data structure is given by stochastic trees.

The most important aspect of the nested distance, defined later in Section 4, consists in accounting for the increase of information over time. Usually this is modeled by a filtration. Let  $(\Omega, \mathcal{F}_T, P)$  be a

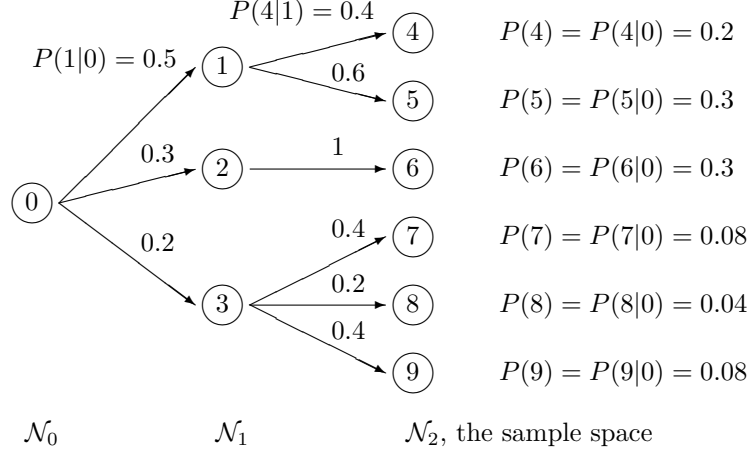


Figure 1: An exemplary finite tree process  $\nu = (\nu_0, \nu_1, \nu_2)$  with nodes  $\mathcal{N} = \{0, \dots, 9\}$  and leaves  $\mathcal{N}_2 = \{4, \dots, 9\}$  at  $T = 2$  stages. The filtrations, generated by the respective atoms, are  $\mathcal{F}_2 = \sigma(\{4\}, \{5\}, \dots, \{9\})$ ,  $\mathcal{F}_1 = \sigma(\{4, 5\}, \{6\}, \{7, 8, 9\})$  and  $\mathcal{F}_0 = \sigma(\{4, 5, \dots, 9\})$  (cf. [PR07, Section 3.1.1]).

probability space and  $\mathcal{F} = (\mathcal{F}_t)_{t \in \{0, 1, \dots, T\}}$  a family of sigma algebras. Then  $\mathcal{F} = (\mathcal{F}_t)_{t \in \{0, 1, \dots, T\}}$  is a *filtration* provided that  $\mathcal{F}_t \subset \mathcal{F}_\tau$  whenever  $t \leq \tau$ . A triple  $(\Omega, \mathcal{F}, P)$ , where  $\mathcal{F} = (\mathcal{F}_t)_{t \in \{0, 1, \dots, T\}}$  is a filtration, is called *filtered probability space* or *stochastic basis*.

Because the nested distance is able to compare processes and their approximating scenario trees, it is necessary to review the main properties of scenario trees and to carefully discuss the relations between processes and trees.

### Trees

A tree is a directed graph  $(\mathcal{N}, A)$  without circles. The vertices  $\mathcal{N}$  will be called *nodes* (following [PR07, p. 216]) in the following. A node  $m \in \mathcal{N}$  is a *direct predecessor* or *parent* of the node  $n \in \mathcal{N}$  if  $(m, n) \in A$ . The parental relation between  $m$  and  $n$  is denoted by  $m = n_-$ . The set of *direct successors* (or *children*) of a vertex  $m$  is  $m_+$ , such that  $m = n_-$  iff  $n \in m_+$ . A node  $m \in \mathcal{N}$  is said to be a *predecessor* of  $n \in \mathcal{N}$  – in symbols  $m \supset n$  – if  $n_- = n_1$ ,  $n_{1-} = n_2$ , and finally  $n_{k-} = m$  for some sequence  $n_k \in \mathcal{N}$ . It holds in particular that  $n_- \supset n$ .

In addition we assume that any node  $n \in \mathcal{N}$  is at a certain *stage* within the tree such that the following properties are fulfilled:

- Nodes at the same stage  $t$  are collected in  $\mathcal{N}_t$ , such that  $\mathcal{N}$  is the union of the disjoint subsets  $\mathcal{N}_0, \mathcal{N}_1, \dots, \mathcal{N}_T$ ,  $\mathcal{N} = \dot{\bigcup}_{t=0}^T \mathcal{N}_t$ ;
- $r \in \mathcal{N}$  is a *root* node if  $r$  is a predecessor of all nodes,  $r \supset n$  ( $n \in \mathcal{N}$ ), and  $\mathcal{N}_0$  (stage 0) contains the root node  $r$ . By convention, the unique node of a rooted tree (the root node) is denoted by 0, hence  $\mathcal{N}_0 = \{0\}$ ;
- $i \in \mathcal{N}$  is a *leaf* node if  $i_+ = \emptyset$ .  $\mathcal{N}_T$  collects all *leaf nodes* of the tree, and  $T$  is the height of the tree;
- $n \in \mathcal{N}_t$  iff  $n_+ \subset \mathcal{N}_{t+1}$ .

In concrete implementations all vertices can be numbered consecutively, starting with 0 for the root node (cf. Figure 1).

Any tree induces a filtration

Any tree with height  $T$  and finitely many nodes  $\mathcal{N}$  naturally induces a filtration  $\mathcal{F}$ : First use  $\mathcal{N}_T$  as sample space. For any  $n \in \mathcal{N}$  define the atom<sup>1</sup>  $a(n) \subset \mathcal{N}_T$  in a backward recursive way by

$$a(n) := \begin{cases} \{n\} & \text{if } n \in \mathcal{N}_T \\ \bigcup_{j \in n_+} a(j) & \text{else.} \end{cases}$$

Employing these atoms, the related sigma algebra is defined by

$$\mathcal{F}_t := \sigma(a(n) : n \in \mathcal{N}_t).$$

From the construction of the atoms it is evident that  $\mathcal{F}_0 = \{\emptyset, \mathcal{N}_T\}$  for a rooted tree and that  $\mathcal{F} = (\mathcal{F}_0, \dots, \mathcal{F}_T)$  is a filtration on the sample space  $\mathcal{N}_T$ , i.e. it holds that  $\mathcal{F}_t \subset \mathcal{F}_{t+1}$ . Notice that node  $m$  is a predecessor of  $n$ , i.e.  $m \supset n$ , if and only if

$$a(m) \supset a(n).$$

This observation suggests the symbol  $m \supset n$  introduced in the previous section for the predecessor relation in a tree structure.

Employing the atoms  $a(n)$  a *tree process* can be defined by

$$\begin{aligned} \nu : \{0, \dots, T\} \times \mathcal{N}_T &\rightarrow \mathcal{N} \\ (t, i) &\mapsto n \text{ if } i \in a(n) \text{ and } n \in \mathcal{N}_t \text{ (i.e. } n \supset i), \end{aligned}$$

such that each

$$\begin{aligned} \nu_t : \mathcal{N}_T &\rightarrow \mathcal{N}_t \\ i &\mapsto \nu(t, i) \end{aligned}$$

is  $\mathcal{F}_t$ -measurable. Moreover, the process  $\nu$  is *adapted* to its *natural filtration*, i.e.

$$\mathcal{F}_t = \sigma(\nu_0, \dots, \nu_t) = \sigma(\nu_t).$$

It is natural to introduce the notation  $i_t := \nu_t(i)$  which denotes the state of the tree process for any final outcome  $i \in \mathcal{N}_T$  at stage  $t$ . It then holds that  $i_T = i$ , and moreover that  $i_t \supset i_\tau$  whenever  $t \leq \tau$ , and finally – for a rooted tree –  $i_0 = 0$ . The *sample path* from the root node 0 to a final node  $i \in \mathcal{N}_T$  is

$$(\nu_t(i))_{t=0}^T = (i_t)_{t=0}^T.$$

Any filtration induces a tree

On the other hand, given a filtration  $\mathcal{F} = (\mathcal{F}_0, \dots, \mathcal{F}_T)$  on a finite sample space  $\Omega$  it is possible to define a tree, representing the filtration: Just consider the set  $A_t$  collecting all atoms generating  $\mathcal{F}_t$  ( $\mathcal{F}_t = \sigma(A_t)$ ), and define the nodes

$$\mathcal{N} := \{(a, t) : a \in A_t\}$$

and the arcs

$$A = \{((a, t), (b, t+1)) : a \in A_t, a \supset b \in A_{t+1}\}.$$

$(\mathcal{N}, A)$  then is a directed tree respecting the filtration  $\mathcal{F}$ .

Hence filtrations on a finite sample space and finite trees are equivalent structures up to possibly different labels, and in the following we will not distinguish between them.

---

<sup>1</sup>A  $\mathcal{F}$ -measurable set  $a \in \mathcal{F}$  is an atom if  $b \subset a$  implies that  $P(a) = 0$ .

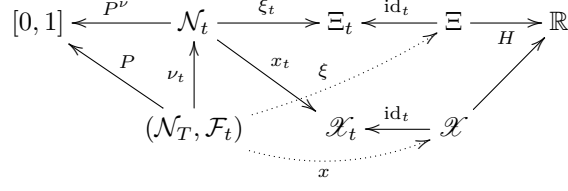


Figure 2: The probabilistic setup: Diagram for the filtered probability space  $(\mathcal{N}_T, (\mathcal{F}_t)_t, P)$  (left) , and the value process  $\xi$  (including the decision process  $x$ ) (right).

### Measures on trees

Let  $P$  be a probability measure on  $\mathcal{F}_T$ , such that  $(\mathcal{N}_T, \mathcal{F}_T, P)$  is a probability space. The notions introduced above allow to *extend* the probability measure to the entire tree via the definition (cf. Figure 1)

$$P^\nu(A) := P \left( \bigcup_{t \in \{0, \dots, T\}} \nu_t^{-1}(A \cap \mathcal{N}_t) \right) \quad (A \subset \mathcal{N}).$$

In particular this definition includes the unconditional probabilities

$$P(\{n\}) =: P(n)$$

for each node. Furthermore it can be used to define conditional probabilities

$$P(\{n\} | \{m\}) =: P(n | m),$$

representing the probability of transition from  $n$  to  $m$ , if  $m \supset n$ .

### Value and decision processes

In a multi-period, discrete time setup the *outcomes* or *realizations* of a stochastic process are of interest, not the concrete model (the sample space): in focus is the sample space

$$\Xi := \Xi_0 \times \dots \times \Xi_T$$

of the stochastic process

$$\xi : \{0, \dots, T\} \times \mathcal{N}_T \rightarrow \Xi.$$

The process is measurable with respect to each  $\mathcal{F}_t = \sigma(\nu_t)$ , from which follows (cf. [Shi96, Theorem II.4.3]) that  $\xi$  can be decomposed as

$$\xi_t = \xi_t \circ \nu_t,$$

(i.e.  $\text{id}_t \circ \xi = \xi_t \circ \nu_t$ , where  $\text{id}_t : \Xi \rightarrow \Xi_t$  is the natural projection) as depicted in Figure 1. Notice that  $\xi_t \in \Xi_t$  is an observation of the stochastic process at stage  $t$  and measurable with respect to  $\mathcal{F}_t$  (in symbols  $\xi_t \triangleleft \mathcal{F}_t$ ), and at this stage  $t$  all prior observations

$$\xi_{0:t} := (\xi_0, \dots, \xi_t)$$

are  $\mathcal{F}_t$ -measurable as well.

In *multistage* stochastic programming, a decision maker has the possibility to influence the results to be expected at the very end of the process by making a decision  $x_t$  at any stage  $t$  of time, having available the information which occurred up to the time when the decision is made, that is  $\xi_{0:t}$ . The decision has to be taken prior to the next observation  $\xi_{t+1}$  (e.g., a decision about a new portfolio allocation has to be made *before* knowing next days security prices).

This *nonanticipativity* property of the decisions is modeled by the assumption that any  $x_t$  is measurable with respect to  $\mathcal{F}_t$  ( $x_t \triangleleft \mathcal{F}_t$ ), such that again

$$x_t = x_t \circ \nu_t$$

(i.e.  $\text{id}_t \circ x = x_t \circ \nu_t$ ).

### 3. The Wasserstein distance for probability measures – definition and computation

Probability metrics are functionals that quantify distances between random objects like random variables, random vectors or even random processes. See e.g. [Rac91] for an encyclopedic treatment or [GS02] for a comprehensive overview of relations between some classical probability metrics. An important group of probability metrics is given by the Wasserstein, or Kantorovich distances.

**Definition 1.** Given two probability spaces  $\mathbb{P} := (\Xi, \Sigma, P)$  and  $\mathbb{P}' := (\Xi', \Sigma', P')$  and a convex function  $d: \Xi \times \Xi' \rightarrow \mathbb{R}$ , the *Wasserstein distance of order  $r \geq 1$*  – denoted  $d_r(P, P')$  – is the optimal value of the optimization problem

$$\begin{aligned} & \text{minimize} && (\int d(\xi, \xi')^r \pi(d\xi, d\xi'))^{\frac{1}{r}} \\ & \text{(in } \pi) && \\ & \text{subject to} && \pi(M \times \Xi') = P(M) \quad (M \in \Sigma), \\ & && \pi(\Xi \times N) = P'(N) \quad (N \in \Sigma'), \end{aligned} \tag{2}$$

where the infimum in (2) is among all bivariate probability measures  $\pi \in \mathcal{P}(\Xi \times \Xi')$  which are measures on the product sigma algebra  $\Sigma \otimes \Sigma'$ . Often  $\Xi = \Xi'$  and in typical applications of interest  $d: \Xi \times \Xi' \rightarrow \mathbb{R}$  is a distance function.

Of particular interest is the Wasserstein distance of order  $r = 2$  with a Euclidean norm  $d(\xi, \xi') = \|\xi - \xi'\|_2$ . We shall refer to this combination as the *quadratic Wasserstein distance*.

The Wasserstein distance was treated first in [Mon81] in an entirely different context. A very comprehensive summary can be found in [Vil03]. In the Russian literature (cf. [Ver06]) the Wasserstein distance is rather known under the name Kantorovich distance. As a matter of fact the Wasserstein distance depends on the sigma algebras  $\Sigma$  and  $\Sigma'$ . This fact is neglected by writing  $d_r(P, P')$ .

Basically, (2) can be interpreted as a transportation problem. The resulting functional  $d_r(\cdot, \cdot)$  can be shown to be a full distance. Furthermore, convergence in  $d_r(\cdot, \cdot)$  is equivalent to weak convergence plus convergence of the  $r$ -th moment (cf. [Vil03]). It has been shown (see e.g. [DGKR03]) that single stage expected loss minimization problems with objective function  $\mathbb{E}_\xi H(\xi, x)$  are (under some regularity conditions on the loss function) Lipschitz continuous with respect to the Wasserstein distance.

*Remark 1.* It should be noted that the Wasserstein distance is a well-defined distance of probability measures, even if the sample spaces  $\Xi$  and  $\Xi'$  are entirely different. The link between different spaces is provided by the distance function – or cost function –  $d$ .

If  $P = \sum_i p_i \delta_{\xi_i}$  and  $P' = \sum_j p'_j \delta_{\xi'_j}$  are discrete measures on a space  $\Xi$  ( $\Xi'$ , respectively), then the Wasserstein distance can be computed by the linear program (LP)

$$\begin{aligned} & \text{minimize} && \sum_{i,j} d_{i,j}^r \pi_{i,j} \\ & \text{(in } \pi) && \\ & \text{subject to} && \sum_j \pi_{i,j} = p_i, \\ & && \sum_i \pi_{i,j} = p'_j, \\ & && \pi_{i,j} \geq 0, \end{aligned} \tag{3}$$

where  $d_{i,j}$  is the matrix with entries  $d_{i,j} = d(\xi_i, \xi'_j)$ .

*Remark 2.* It can be derived from the complementary slackness conditions for linear programs that the optimizing transport plan  $\pi_{i,j}$  in (3) is sparse, i.e. it has at most  $|\Xi| + |\Xi'| - 1$  non-zero entries. This corresponds to the number of entries in one row plus one column of the matrix  $\pi$  or  $d$ .

#### 3.1. Scenario approximation with Wasserstein distances

Given a probability measure  $P$  one might ask for the best approximating probability measure, with support  $Q$ . The following Lemma 1 reveals that the probability measure  $P_Q^*$ , which is the best approximation of  $P$  located just on  $Q$ , i.e.

$$d_r(P, P_Q^*) \leq d_r(P, P') \quad (P'(Q) = 1), \tag{4}$$

can be computed in a direct way.

**Lemma 1** (Lower bounds and best approximation). *Let  $P$  and  $P'$  be probability measures.*

(i) *The Wasserstein distance has the lower bound*

$$\mathbf{d}_r(P, P')^r \geq \int \min_{\xi \in \Xi'} d(\xi, \xi')^r P(d\xi). \quad (5)$$

(ii) *The lower bound in (5) is attained if the transport map  $\mathbf{T} : \Xi \rightarrow \Xi'$  with  $\mathbf{T}(\xi) \in \operatorname{argmin}_{\xi'} d(\xi, \xi')$  is measurable. The pushforward  $P^* := P \circ \mathbf{T}^{-1}$  satisfies<sup>2</sup>*

$$\mathbf{d}_r(P, P^*)^r = \int \min_{\xi' \in \Xi'} d(\xi, \xi')^r P(d\xi). \quad (6)$$

(iii) *If  $\Xi = \Xi'$  is vector space, then*

$$\mathbf{d}_r(P, P^{\tilde{\mathbf{T}}}) \leq \mathbf{d}_r(P, P^{\mathbf{T}}),$$

where  $\tilde{\mathbf{T}}(\xi) := \mathbb{E}_P(\tilde{\xi} | \mathbf{T}(\tilde{\xi}) = \mathbf{T}(\xi))$ .

*Proof of Lemma 1.* Let  $\pi$  have the marginals of  $P$  and  $P'$ . Then

$$\begin{aligned} \int_{\Xi \times \Xi'} d(\xi, \xi')^r \pi(d\xi, d\xi') &\geq \int_{\Xi} \int_{\Xi'} \min_{q \in \Xi'} d(\xi, q)^r \pi(d\xi, d\xi') \\ &= \int_{\Xi} \min_{q \in \Xi'} d(\xi, q)^r P(d\xi). \end{aligned}$$

Taking the infimum leads to (5).

Employing the transport map  $\mathbf{T}$ , define the transport plan  $\pi := P \circ (\operatorname{id}_{\Xi} \times \mathbf{T})^{-1}$  where  $\operatorname{id}_{\Xi}$  is the identity on  $\Xi$ , i.e.

$$\pi(A \times B) = P(\{\xi : (\xi, \mathbf{T}(\xi)) \in A \times B\}) = P(\{\xi : \xi \in A, \mathbf{T}(\xi) \in B\}).$$

$\pi$  is feasible, hence it has the marginals  $\pi(A \times \Xi') = P(\{\xi : \xi \in A, \mathbf{T}(\xi) \in \Xi'\}) = P(A)$  and  $\pi(\Xi \times B) = P(\{\xi : \mathbf{T}(\xi) \in B\}) = P^{\mathbf{T}}(B)$ . Thus

$$\int \int_{\Xi \times \Xi'} d(\xi, \xi')^r \pi(d\xi, d\xi') = \int_{\Xi} d(\xi, \mathbf{T}(\xi))^r P(d\xi) = \int_{\Xi} \min_{\xi' \in \Xi'} d(\xi, \xi')^r P(d\xi),$$

which is (6).

For the last assertion apply the conditional Jensen's inequality  $\varphi \circ \mathbb{E}(X | \mathbf{T}) \leq \mathbb{E}(\varphi(X) | \mathbf{T})$  to  $\varphi(y) := d(x, y)$  and obtain

$$d(x, \mathbb{E}(\operatorname{id} | \mathbf{T}) \circ \mathbf{T}) \leq \mathbb{E}(d(x, \operatorname{id}) | \mathbf{T}) \circ \mathbf{T}.$$

The measure  $\tilde{\pi}(A \times B) := P(A \cap \tilde{\mathbf{T}}^{-1}(B))$  has marginals  $P$  and  $P^{\tilde{\mathbf{T}}}$ , from which follows that

$$\begin{aligned} \mathbf{d}_r(P, P^{\tilde{\mathbf{T}}})^r &\leq \int d(\xi, \tilde{\mathbf{T}}(\xi))^r P(d\xi) = \int d(\xi, \mathbb{E}(\operatorname{id} | \mathbf{T}) \circ \mathbf{T}(\xi))^r P(d\xi) \\ &\leq \int \mathbb{E}(d(\xi, \operatorname{id})^r | \mathbf{T})(\mathbf{T}(\xi)) P(d\xi) = \int d(\xi, \mathbf{T}(\xi))^r P(d\xi) = \mathbf{d}_r(P, P^{\mathbf{T}})^r. \end{aligned}$$

□

It should be noted that the measure  $P_Q^*$  does *not depend* on the order  $r$ . Moreover, given a probability measure  $P$ , Lemma 1 allows to find the best approximation, which is located just on finitely many points  $Q = \{q_1 \dots q_n\}$ . For this consider  $\Xi' = Q$ , define  $p_j^* := P(\mathbf{T} = q_j)$  (the collection of distinct sets  $\{\mathbf{T} = q_j\}$  is a Voronoi tessellation; for a comprehensive treatment see [GL00] and the oeuvre of Gilles Pagès, e.g. [BPP05]) and  $P_Q^* := \sum_j p_j^* \delta_{q_j}$ , as above. Then  $\mathbf{d}_r(P, P_Q^*)^r = \int \min_{q \in Q} d(\xi, q)^r P(d\xi)$ , and no better approximation is possible by Lemma 1. Usually the  $q_j$  are called *quantizers*, which we will adopt in the following.

<sup>2</sup>see also [DGKR03, Theorem 2]

*Optimal probabilities*

According to Lemma 1 the best approximating measure for  $P = \sum_i p_i \delta_{\xi_i}$ , which is located on  $Q$ , is  $P_Q^* = \sum_j p_j^* \delta_{q_j}$ . The respective linear program is

$$\begin{aligned} & \text{minimize} && \sum_{i,j} d_{i,j}^r \pi_{i,j} \\ & \text{(in } \pi) && \\ & \text{subject to} && \sum_j \pi_{i,j} = p_i, \\ & && \pi_{i,j} \geq 0, \end{aligned}$$

which is solved by the optimal transport plan

$$\pi_{i,j}^* := \begin{cases} p_i & \text{if } d(\xi_i, q_j) = \min_{q \in Q} d(\xi_i, q) \\ 0 & \text{else} \end{cases} \quad (7)$$

such that

$$p_j^* = \sum_i \pi_{i,j}^* \quad \text{and} \quad \mathbf{d}_r(P, P_Q^*)^r = \mathbb{E}_{\pi^*} d^r. \quad (8)$$

Observe as well that the matrix  $\pi^*$  in (7) has just  $|\Xi|$  non-zero entries, which is less than in Remark 3: In every row  $i$  of  $\pi^*$  there is just one non-zero entry  $\pi_{i,j}^*$ .

Given the support points  $Q$ , it is hence an easy exercise to look up the closest points according to (7), and sum up their probabilities according (8), such that the solution of (4) is immediately obtained by  $P_Q^* = \sum_j p_j^* \delta_{q_j}$ .

*Optimal supporting points – facility location*

Given the previous results on optimal probabilities the problem of finding a sufficiently good approximation of  $P$  in the Wasserstein distance reduces to the problem of looking up good locations  $Q$ , that is to minimize the function

$$\{q_1, \dots, q_n\} \mapsto \mathbf{d}_r(P, P_{\{q_1, \dots, q_n\}}^*)^r = \int \min_{q \in \{q_1, \dots, q_n\}} d(\xi, q)^r P(d\xi). \quad (9)$$

This problem is often referred to as *facility location* [DH02]. It is not convex, and no closed form solution exists in general, it hence has to be handled with adequate numerical algorithms. Moreover the facility location problem is NP-hard.

For the important case of the quadratic Wasserstein distance Lemma 1 and its proof give rise for an adaption of the k-means clustering algorithm (also referred to as Lloyd's algorithm, [Llo82]). The approach is described in Algorithm 1.

**Theorem 1.** *The measures  $P^k$  generated by Algorithm 1 are improved approximations for  $P$ , they satisfy*

$$\mathbf{d}_r(P, P^{k+1}) \leq \mathbf{d}_r(P, P^k).$$

*Algorithm 1 terminates after finitely many iterations.*

*In the case of the quadratic Wasserstein distance Algorithm 1 terminates at a local minimum  $\{q_1, \dots, q_n\}$  of (9).*

*Proof.* Algorithm 1 is an iterative refinement technique, which finds the measure

$$P^k = \sum_{j=1}^n P(T_j^k) \delta_{q_j^k}$$

after  $k$  iterations. By construction of (10) it is an improvement due to Lemma 1, (ii) and (iii), and hence

$$\mathbf{d}_r(P, P^{k+1}) \leq \mathbf{d}_r(P, P^k).$$



**Algorithm 1**

Facility location for  $P = \sum_i p_i \delta_{\xi_i}$  in the special case of the Euclidean distance and quadratic Wasserstein distance (order  $r = 2$ ).

**Initialization** ( $k = 0$ ):

Choose  $n$  points  $\{q_i^0 : i = 1, \dots, n\}$ , for example by randomly picking  $n$  distinct points from  $\{\xi_i : i\}$ .

**Assignment Step:**

In each step  $k$  assign each  $\xi_i$  to the cluster with the closest mean,

$$T_j^k := \{\xi_i : \|\xi_i - q_j^k\| \leq \|\xi_i - q_{j'}^k\| \text{ for all } j' = 1, \dots, n\}$$

and set

$$P^k := \sum_{j=1}^n P(T_j^k) \delta_{q_j^k}.$$

**Update Step:**

Set

$$q_j^{k+1} := \sum_{\xi_i \in T_j^k} \frac{P(\xi_i)}{P(T_j^k)} \xi_i. \quad (10)$$

**Iteration:**

Set  $k \leftarrow k + 1$  and continue with an assignment step until  $\{q_j^{(k)} : j = 1, \dots, n\}$  is met again.

The algorithm terminates after finitely many iterations because there are just finitely many Voronoi-combinations  $T_j$ .

For the Euclidean distance and  $r = 2$  the expectation  $\mathbb{E}X = \sum_i p_i x_i$  minimizes the function

$$q \mapsto \sum_i p_i \cdot \|q - \xi_i\|_2^2 = \mathbb{E} \|q - \xi\|_2^2.$$

In this case  $P^k$  thus is a local minimum of (9).  $\square$

For other distances than the quadratic Wasserstein distance,  $P^k$  is possibly a good starting point, but in general *not* a local (global) minimum of (9).

**4. A nested distance for stochastic processes**

The Wasserstein distance basically is a distance for random variables. However, it can be used for stochastic processes as well:

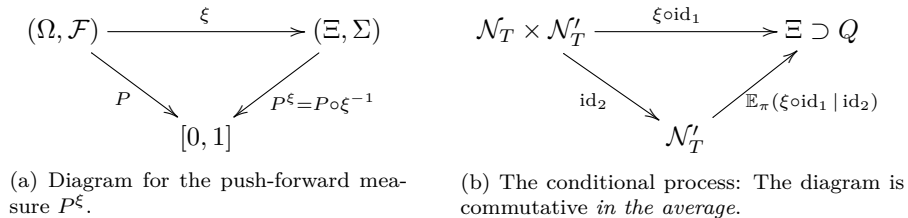


Figure 3: Pushforward measure and the projected process.

A random variable

$$\xi: (\Omega, \mathcal{F}) \rightarrow (\Xi, \Sigma)$$

on a probability space  $(\Omega, \mathcal{F}, P)$  naturally induce the push-forward measure

$$P^\xi := P \circ \xi^{-1}: \Sigma \rightarrow [0, 1]$$

on the state space  $(\Xi, \Sigma)$  (cf. Figure 3a and Figure 1), such that in particular the distance  $d_r(P^\xi, P^{\xi'})$  is available by employing a distance function

$$d: \Xi \times \Xi' \rightarrow \mathbb{R}$$

defined on the *state spaces*  $\Xi \times \Xi'$ . If the spaces  $\Xi$  and  $\Xi'$  are interpreted as containing the sample paths of Stochastic processes, it is possible to consider a process as a random variables, and to apply the Wasserstein distance to the sample paths and their distribution. However, the gradually increasing information, which is the essential ingredient of stochastic processes, is simply ignored by just having a look at the final sigma algebra  $\sigma(\xi) \subset \Sigma_T$  instead of the entire filtration  $(\Sigma_0, \dots, \Sigma_T)$ . Figure 4 shows a situation where similar paths (small  $\varepsilon$ ) lead to a small value of the Wasserstein distance between the first and the second case, which neglects the fact that in the second case perfect information about the final outcome is available already at the intermediary step.

These considerations led to the proposal of a nested distance in [Pff09, PP12].

#### 4.1. Definition of the nested distance

In the following we will use a multistage (nested) distance concept that shares many properties of the Wasserstein distances but accounts for the effects of filtrations. It was introduced first in [Pff09] and analyzed in [PP12]. In order to introduce nested distances we have to generalize the distributional concepts used so far from random variables to stochastic processes. For this consider the process  $(\xi_t)_{t \in \{0, \dots, T\}}$ , where  $\xi_t: (\Omega, \mathcal{F}) \rightarrow (\Xi_t, \Sigma_t)$  are random variables with possibly different state spaces  $(\Xi_t, \Sigma_t)$ . Define the product space  $\Xi := \Xi_1 \times \dots \times \Xi_T$ , which can be equipped itself with the product sigma algebra  $\Sigma := \sigma(\Sigma_1 \otimes \dots \otimes \Sigma_T)$ . Then

$$\begin{aligned} \xi: (\Omega, \mathcal{F}) &\rightarrow (\Xi, \Sigma) \\ \omega &\mapsto (\xi_t(\omega))_{t \in \{0, \dots, T\}} \end{aligned}$$

is a random variable, mapping any outcome  $\omega \in \Omega$  to its path  $(\xi_t(\omega))_{t=0}^T$ . The law of of the process  $\xi$ ,

$$P^\xi := P \circ \xi^{-1}: \Sigma \rightarrow [0, 1],$$

is the push-forward measure on  $\Xi = \Xi_1 \times \dots \times \Xi_T$ . The situation for processes is completely analogous to random variables, provided that a distance function

$$d: \Xi \times \Xi' \rightarrow \mathbb{R}$$

on  $\Xi = \Xi_0 \times \dots \times \Xi_T$  and  $\Xi' = \Xi'_0 \times \dots \times \Xi'_T$  is available. For metric spaces  $(\Xi_t, d_t)$  the functions  $d(\xi, \xi') = \sum_{t=0}^T d_t(\xi_t, \xi'_t)$  or  $d(\xi, \xi') = \max_{t \in \{0, \dots, T\}} d_t(\xi_t, \xi'_t)$  are immediate candidates. We shall

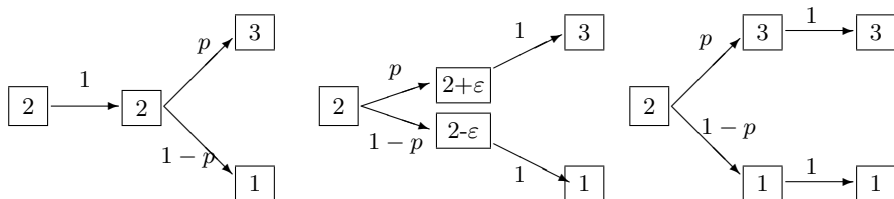


Figure 4: Three tree processes to illustrate three different flows of information: If  $\varepsilon$  is small the Wasserstein distance will be small too (cf. [HRS06]).

call a metric (*weighted*) *Euclidean*, if  $d(\xi, \xi') = \left( \sum_{t=0}^T w_t \|\xi_t - \xi'_t\|_2^2 \right)^{1/2}$ , where  $w_t > 0$  are positive weights and each norm  $\|\cdot\|_2$  satisfies the parallelogram law.

With these preparations the nested distances can be defined as follows:

**Definition 2.** For two filtered probability spaces  $\mathbb{P} := (\Xi, \Sigma, P)$ ,  $\mathbb{P}' := (\Xi', \Sigma', P')$  and a real-valued, convex function  $d: \Xi \times \Xi' \rightarrow \mathbb{R}$  the *nested distance of order  $r \geq 1$*  – denoted  $\mathbf{d}_r(\mathbb{P}, \mathbb{Q})$  – is the optimal value of the optimization problem

$$\begin{aligned} & \text{minimize} && \left( \int d(\xi, \xi')^r \pi(d\xi, d\xi') \right)^{\frac{1}{r}} \\ & \text{(in } \pi) && \\ \text{subject to} &&& \pi(M \times \Xi' \mid \Sigma_t \otimes \Sigma'_t) = P(M \mid \Sigma_t) && (M \in \Sigma_T, t \in \{0, \dots, T\}), \\ &&& \pi(\Xi \times N \mid \Sigma_t \otimes \Sigma'_t) = P'(N \mid \Sigma'_t) && (N \in \Sigma'_T, t \in \{0, \dots, T\}), \end{aligned} \quad (11)$$

where the infimum in (11) is among all bivariate probability measures  $\pi \in \mathcal{P}(\Xi \times \Xi')$ , which are measures on the product sigma algebra  $\Sigma_T \otimes \Sigma'_T$ . We will refer to the nested distance also as *process distance*, or *multistage distance*. The nested distance  $\mathbf{d}_2$  (order  $r = 2$ ), with  $d$  a weighted Euclidean distance is referred to as *quadratic nested distance*.

Note that the minimization (2) for the Wasserstein distance  $\mathbf{d}_r(P, P')$  is a relaxation of (11). Hence the Wasserstein distance is always less or equal to the nested distance,

$$\mathbf{d}_r(P, P') \leq \mathbf{d}_r(\mathbb{P}, \mathbb{P}').$$

It is possible therefore to decompose the nested distance into the Wasserstein and the effect  $\mathbf{d}_r(\mathbb{P}, \mathbb{P}')$  –  $\mathbf{d}_r(P, P')$  caused by the filtration related to the additional constraints in (11).

The multistage distance  $\mathbf{d}_r(\cdot, \cdot)$  also preserves important regularity properties (Lipschitz and Hölder continuity) of the objective function of multistage stochastic programs (see [PP12, Section 6]).

#### 4.2. The nested distance for trees.

The Wasserstein distance between discrete probability measures can be calculated by solving the linear program (3). To establish a similar linear program for the nested distance we use trees that model the whole filtration. Then problem (11) reads

$$\begin{aligned} & \text{minimize} && \sum_{i,j} \pi_{i,j} \cdot d_{i,j}^r \\ & \text{(in } \pi) && \\ \text{subject to} &&& \sum_{j \subset n} \pi(i, j \mid m, n) = P(i \mid m) && (m \supset i, n), \\ &&& \sum_{i \subset m} \pi(i, j \mid m, n) = P'(j \mid n) && (n \supset j, m), \\ &&& \pi_{i,j} \geq 0 \text{ and } \sum_{i,j} \pi_{i,j} = 1, \end{aligned} \quad (12)$$

where again  $\pi_{i,j}$  is a matrix defined on the samples ( $i \in \mathcal{N}_T, j \in \mathcal{N}'_T$ ) and  $m \in \mathcal{N}_t, n \in \mathcal{N}'_t$  are arbitrary nodes. The conditional probabilities  $\pi(i, j \mid m, n)$  are given by

$$\pi(i, j \mid m, n) = \frac{\pi_{i,j}}{\sum_{i' \subset m, j' \subset n} \pi_{i',j'}}.$$

The nested structure of the transportation plan  $\pi$ , which is induced by the trees, is schematically depicted in Figure 5.

The constraints in (12) can be written in more detail

$$\begin{aligned} P(i) \cdot \sum_{i' \subset m, j' \subset n} \pi_{i',j'} &= P(m) \cdot \sum_{j' \subset n} \pi_{i,j'} && (m \supset i, n) && \text{and} \\ P'(j) \cdot \sum_{i' \subset m, j' \subset n} \pi_{i',j'} &= P'(n) \cdot \sum_{i' \subset m} \pi_{i',j} && (m, n \supset j). \end{aligned}$$

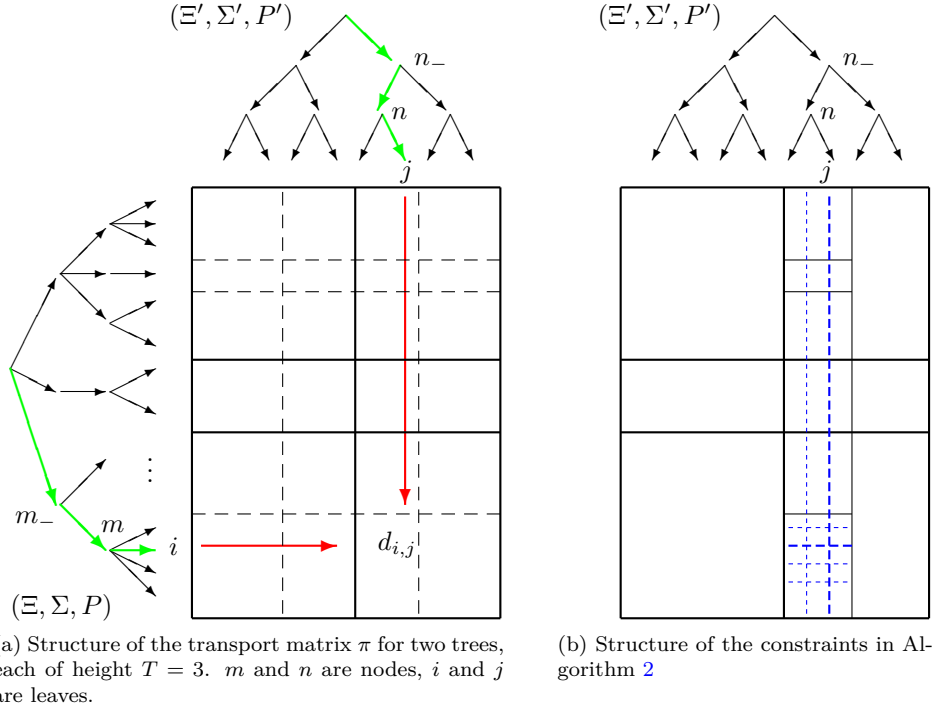


Figure 5: Schematic structure of the distance matrix  $d$  and the transport matrix  $\pi$ , as it is imposed by the structures of the trees and the respective constraints.

As  $P$  and  $P'$  are given, this shows that (12) is equivalent to the linear program

$$\begin{aligned}
& \text{minimize} && \sum_{i,j} \pi_{i,j} \cdot d_{i,j}^r \\
& \text{(in } \pi) && \\
& \text{subject to} && P(i) \cdot \sum_{i' \subset m, j' \subset n} \pi_{i',j'} = P(m) \cdot \sum_{j' \subset n} \pi_{i,j'} \quad (m \supset i), \\
& && P'(j) \cdot \sum_{i' \subset m, j' \subset n} \pi_{i',j'} = P'(n) \cdot \sum_{i' \subset m} \pi_{i',j} \quad (n \supset j), \\
& && \pi_{i,j} \geq 0 \text{ and } \sum_{i,j} \pi_{i,j} = 1.
\end{aligned}$$

*Remark 3.* As a matter of fact many constraints in (12) are linearly dependent. For computational reasons (loss of significance during numerical evaluations, which can impact linear dependencies and the feasibility) it is advisable to remove linear dependencies. This is partially accomplished by the simpler program

$$\begin{aligned}
& \text{minimize} && \sum_{i,j} \pi_{i,j} \cdot d_{i,j}^r \\
& \text{(in } \pi) && \\
& \text{subject to} && \sum_{j \in (n_-)_+} \pi(m, j | m_-, n_-) = P(m | m_-) \quad (m \in \mathcal{N} \setminus \mathcal{N}_0), \\
& && \sum_{i \in (m_-)_+} \pi(i, n | m_-, n_-) = P'(n | n_-) \quad (n \in \mathcal{N}' \setminus \mathcal{N}'_0), \\
& && \pi_{i,j} \geq 0 \text{ and } \sum_{i,j} \pi_{i,j} = 1,
\end{aligned} \tag{13}$$

which by [PP12, Lemma 10] is equivalent to (12) and can be reformulated as an LP as well. Further constraints can be removed from (13) by taking into account that  $\sum_{i=m_-} \frac{P(i)}{P(m_-)} = 1$ . Hence, for each node  $m$  it is possible to drop one constraint out of all  $|(m_-)_+|$  related equations.

It should be noted that instead of solving (13) the nested distance can be calculated in a recursive way. First define

$$dl_r(i, j) := d(\xi_i, \xi'_j) \tag{14}$$

for  $i \in \mathcal{N}_T$ ,  $j \in \mathcal{N}'_T$ . Given  $\mathbf{d}_r(i, j)$  for  $i \in \mathcal{N}_{t+1}$  and  $j \in \mathcal{N}'_{t+1}$  set

$$\mathbf{d}_r(m, n)^r := \sum_{i \in m_+, j \in n_+} \pi(i, j | m, n) \cdot \mathbf{d}_r(i, j)^r \quad (m \in \mathcal{N}_t, n \in \mathcal{N}'_t) \quad (15)$$

for  $m \in \mathcal{N}_t$ ,  $n \in \mathcal{N}'_t$ , where the conditional probabilities  $\pi(\cdot, \cdot | m, n)$  solve

$$\begin{aligned} & \text{minimize} && \sum_{i \in m_+, j \in n_+} \pi(i, j | m, n) \cdot \mathbf{d}_r(i, j)^r \\ & \text{in } \pi(\cdot, \cdot | m, n) && \\ & \text{subject to} && \sum_{j \in n_+} \pi(i, j | m, n) = P(i | m) \quad (i \in m_+), \\ & && \sum_{i \in m_+} \pi(i, j | m, n) = P'(j | n) \quad (j \in n_+), \\ & && \pi(i, j | m, n) \geq 0. \end{aligned}$$

The values  $\mathbf{d}_r(i, j)$  can be interpreted as conditional nested distances for the trees starting in nodes  $i$  ( $j$ , resp.). Finally the transport plan  $\pi$  on the leaves is recomposed by

$$\pi(i, j) = \pi(i, j | i_{T-1}, j_{T-1}) \cdot \pi(i_{T-1}, j_{T-1} | i_{T-2}, j_{T-2}) \cdot \dots \cdot \pi(i_1, j_1 | 0, 0)$$

and the nested distance is given by  $\mathbf{d}_r(\mathbb{P}, \mathbb{P}')^r = \mathbf{d}_r(0, 0)$ .

## 5. Improving an approximating tree

Lemma 1 and the succeeding remark explain how to approximate a probability measure  $P$  by a measure  $P_Q^*$ , which is located just on the points  $Q = \{q_1 \dots q_n\}$ : the measure  $P_Q^* = \sum_j P(\mathbf{T} = q_j) \cdot \delta_{q_j}$  (cf. (7)) was found to be the best choice with respect to the Wasserstein distance, irrespective of the order  $r \geq 1$ . In this section we address the question of looking up processes, i.e. trees, which are close in nested distance to a given tree.

As in Section 3 we split the problem in two parts:

- (i) Find *probabilities* on a given tree structure, which improve the nested distance to a given tree;
- (ii) facility location: Improve the *locations*, which are the scenarios of the tree, to again improve the approximation overall.

### 5.1. Optimal probabilities

In a multistage context we have to answer the question which probability measure  $P_Q^*$  is *best* to approximate  $\mathbb{P} = (\Xi, \Sigma, P)$ , provided that the *states*  $Q \subset \Xi'$  and *filtration*  $\Sigma'$  of the stochastic processes are given: Knowing the branching structure of the tree, we seek for the best probabilities such that the multistage distance to  $\mathbb{P}$  is as small as possible. The best approximation,  $P_Q^*$ , satisfies

$$\mathbf{d}_r(\mathbb{P}, (\Xi', \Sigma', P_Q^*)) \leq \mathbf{d}_r(\mathbb{P}, (\Xi', \Sigma', P')) \quad (P'(Q) = 1),$$

where  $Q = \{q_1, \dots, q_n\} \subset \Xi'$ .

Compared to the Wasserstein distances it is considerably more difficult to find those optimal probabilities. From (13) it follows that the corresponding transport plan  $\pi^*$  necessarily satisfies

$$\begin{aligned} & \text{minimize} && \sum_{i,j} \pi_{i,j}^* \cdot d_{i,j}^r \\ & \text{(in } \pi^*) && \\ & \text{subject to} && \sum_j \pi^*(m, j | m_-, n_-) = P(m | m_-), \\ & && \sum_i \pi^*(i, n | m_-, n_-) = \sum_i \pi^*(i, n | \tilde{m}_-, n_-), \quad m, \tilde{m} \in \mathcal{N}_t \\ & && \pi_{i,j}^* \geq 0 \text{ and } \sum_{i,j} \pi_{i,j}^* = 1. \end{aligned} \quad (16)$$

The constraint

$$\sum_i \pi^*(i, n | m_-, n_-) = \sum_i \pi^*(i, n | \tilde{m}_-, n_-) \quad (m, \tilde{m} \in \mathcal{N}_t) \quad (17)$$

for nodes  $m$  and  $\tilde{m}$  at the same stage  $t$  in (16) ensures that

$$P_Q^*(n|n_-) := \sum_i \pi^*(i, n|m_-, n_-)$$

is well defined (as it is independent of  $m$ ), allowing thus to reconstruct a measure  $P_Q^*$  by  $P_Q^* = \sum_j \delta_{q_j} \cdot \sum_i \pi_{i,j}^*$ .

Unfortunately, problem (16) does not allow an immediate solution in general. Moreover the constraints (17), for

$$\frac{\sum_{i \in m_-} \pi_{i,n}^*}{\sum_{i \in m_-, j \in n_-} \pi_{i,j}^*} = \pi^*(n|m_-, n_-) = \pi^*(n|\tilde{m}_-, n_-) = \frac{\sum_{i \in \tilde{m}_-} \pi_{i,n}^*}{\sum_{i \in \tilde{m}_-, j \in n_-} \pi_{i,j}^*},$$

are not linear any more – in fact they are multilinear in  $\pi^*$ .

#### Recursive computation of the nested distance

Formulation (16) and the fact that the nested distance can be calculated in a recursive way (see (14) and (15)) leads to the idea of calculating improved probabilities in a recursive way too:

Assume that  $\pi$  is feasible for given quantizers  $Q$ . Define

$$\mathbf{dl}_r(i, j) := d(\xi_i, q_j) \quad (18)$$

for  $i \in \mathcal{N}_T$ ,  $j \in \mathcal{N}'_T$  and, given  $\mathbf{dl}_r(i, j)$  for  $i \in \mathcal{N}_{t+1}$  and  $j \in \mathcal{N}'_{t+1}$ , recursively compute

$$\mathbf{dl}_r(m, n)^r := \sum_{i \in m_+, j \in n_+} \pi^*(i, j|m, n) \cdot \mathbf{dl}_r(i, j)^r \quad (m \in \mathcal{N}_t) \quad (19)$$

for  $m \in \mathcal{N}_t$ ,  $n \in \mathcal{N}'_t$ , where the conditional probabilities  $\pi^*(\cdot, \cdot|m, n)$  solve

$$\begin{aligned} & \text{minimize} && \sum_{m \in \mathcal{N}_t} \pi(m, n) \cdot \sum_{i \in m_+, j \in n_+} \tilde{\pi}(i, j|m, n) \cdot \mathbf{dl}_r(i, j)^r \\ & \text{in } \tilde{\pi}(\cdot, \cdot|m, n) && \\ & \text{subject to} && \sum_{j \in n_+} \tilde{\pi}(i, j|m, n) = P(i|m) \quad (i \in m_+), \\ & && \sum_{i \in m_+} \tilde{\pi}(i, j|m, n) = \sum_{i \in \tilde{m}_+} \tilde{\pi}(i, j|\tilde{m}, n) \quad (j \in n_+), \\ & && \tilde{\pi}(i, j|m, n) \geq 0. \end{aligned} \quad (20)$$

Recomposing the transport plan  $\pi^*$  on the leaves  $i \in \mathcal{N}_T$  and  $j \in \mathcal{N}'_T$  by

$$\pi^*(i, j) = \pi^*(i_T, j_T|i_{T-1}, j_{T-1}) \cdot \pi^*(i_{T-1}, j_{T-1}|i_{T-2}, j_{T-2}) \cdot \dots \cdot \pi^*(i_1, j_1|0, 0) \quad (21)$$

leads to improved probabilities:

**Theorem 2.** *Let  $P'$  be the measure related to the feasible transport probabilities  $\pi$  and  $P'^*$  be related to the probabilities  $\pi^*$  by*

$$P'^* := \sum_j \delta_{q_j} \cdot \sum_i \pi^*(i, j).$$

*Then  $\mathbf{dl}_r(\mathbb{P}, \mathbb{P}^*) \leq \mathbf{dl}_r(\mathbb{P}, \mathbb{P}')$  and the improved distance is given by*

$$\mathbf{dl}_r(\mathbb{P}, \mathbb{P}^*) = \mathbf{dl}_r(0, 0).$$

*Proof.* Observe that the measures  $\pi$  and  $\pi^*$  have the iterative decomposition

$$\begin{aligned} \pi(i, j) &= \pi(i_T, j_T) \\ &= \pi(i_T, j_T|i_{T-1}, j_{T-1}) \cdot \pi(i_{T-1}, j_{T-1}|i_{T-2}, j_{T-2}) \cdot \dots \cdot \pi(i_1, j_1|0, 0), \end{aligned}$$

for all leaves  $i \in \mathcal{N}_T$  and  $i \in \mathcal{N}'_T$  (cf. [Dur04, Chapter 4, Theorem 1.6]). The terminal distance ( $t = T$ ), given the entire history up to  $(i, j)$ , is  $\mathbf{dl}_{T,r}(i, j) := d(i, j)$ , which serves as a starting value for the

iterative procedure. To improve a given transport plan  $\pi$  the algorithm in (20) fixes the conditional probabilities  $\pi(m, n)$  in an iterative step at stage  $t$ .

The constraints in (20) ensure, for

$$\sum_{i \in m_+} \sum_{j \in n_+} \pi^*(i, j | m, n) = \sum_{i \in m_+} P(i | m) = 1,$$

that  $\pi^*$  again is a probability measure for each  $m \in \mathcal{N}'_t$ , and hence, by (21),  $\pi^*$  is a probability measure on  $\mathcal{N}'_T \times \mathcal{N}'_T$ . Furthermore the constraints ensure that  $\pi^*$  respects the tree structures of both trees:  $\pi^*$  is feasible for (7). Finally it holds that

$$\sum_{i, j} \pi_{i, j}^* d(i, j)^r = \mathbf{dl}_r(0, 0)^r$$

due to the recursive construction.

As the initial  $\pi$  is feasible as well for all equations in (20) it follows from the construction that

$$\mathbf{dl}_r(\mathbb{P}, \mathbb{P}'^*)^r = \mathbf{dl}_r(0, 0)^r = \mathbb{E}_{\pi^*} d^r \leq \mathbb{E}_{\pi} d^r.$$

As  $\pi$  was chosen arbitrarily it follows that

$$\mathbf{dl}_r(\mathbb{P}, \mathbb{P}'^*) \leq \mathbf{dl}_r(\mathbb{P}, \mathbb{P}'),$$

which shows that  $\mathbb{P}'^*$  is an improvement.  $\square$

## 5.2. Optimal scenarios – facility location

Consider quantizers

$$Q = \{q_1, \dots, q_n\}$$

where each  $q_j = (q_{j,0}, \dots, q_{j,T})$  is a path in the tree. Given a fixed, feasible measure  $\pi$  define

$$D_{\pi}(\{q_1, \dots, q_n\})^r := \mathbb{E}_{\pi} d^r = \sum_{i, j} \pi_{i, j} d(\xi_i, q_j)^r. \quad (22)$$

The problem of finding optimal quantizers consists in solving the minimization problem

$$\min_{q_1, \dots, q_n} D_{\pi}(\{q_1, \dots, q_n\}). \quad (23)$$

Again it is difficult to solve (23), which can be considered as a facility location problem. However, in an iterative procedure as proposed in the following, a few steps of significant descent in each iteration will be sufficient to considerably improve the overall approximation.

In many applications the gradient of function (22) is available as an analytic expression, for example if  $d(\xi_i, \xi'_j) = (\sum_t d_t(\xi_i, \xi'_j)^p)^{1/p}$ . In this situation the derivative of  $D_{\pi}(\{q_1, \dots, q_n\})^r$  is given by

$$\nabla_{\xi'_{j,t}} D(\xi') = D_{\pi}(\xi')^{1-r} \cdot \sum_i \pi_{i,j} d(\xi_i, \xi'_j)^{r-p} \cdot d_t(\xi_{i,t}, \xi'_{j,t})^{p-1} \cdot \nabla_{\xi'_j} d_t(\xi_{i,t}, \xi'_{j,t}) \quad (j \in \mathcal{N}'_t).$$

If in addition the metric at stage  $t$  is a norm,  $d_t(\xi_{i,t}, \xi'_{j,t}) = \|\xi_{i,t} - \xi'_{j,t}\|_s$ , then it holds that

$$\nabla_{\xi'_{j,t}} d_t(\xi_{i,t}, \xi'_{j,t}) = d_t(\xi_{i,t}, \xi'_{j,t})^{1-s} \cdot \|\xi_{i,t} - \xi'_{j,t}\|_s^{s-2} \cdot (\xi_{i,t} - \xi'_{j,t})$$

which can be obtained by direct computation.

To compute the minimum in (23) a few steps by the steepest descent method will ensure some successive improvements. Another possible method is the limited memory BFGS (Broyden-Fletcher-Goldfarb-Shanno) method, cf [Noc80].

In the special case of the quadratic nested distance the facility location problem can be accomplished by explicit evaluations. This is by far the fastest procedure, and summarized in Algorithm 2, Step 3.

**Theorem 3.** *For a quadratic nested distance the scenarios*

$$q_t(n_t) := \sum_{m_t \in \mathcal{N}_t} \frac{\pi(m_t, n_t)}{\sum_{m_t \in \mathcal{N}_t} \pi(m_t, n_t)} \cdot \xi_t(m_t)$$

(cf. (24)) are the best possible choice to solve the facility location problem (23).

*Proof.* The explicit decomposition of the nested distance allows for the re-arrangement

$$\begin{aligned} \text{dl}_2(\mathbb{P}, \mathbb{P}')^2 &= \sum_{i,j} \pi_{i,j} d(\xi_i, q_j)^2 \\ &= \sum_{i,j} \pi_{i,j} \sum_{t=0}^T w_t \cdot \|\xi_{i_t} - q_{j_t}\|_2^2 \\ &= \sum_{t=0}^T w_t \cdot \sum_{n_t \in \mathcal{N}'_t} \left( \sum_{m_t \in \mathcal{N}_t} \pi(m_t, n_t) \|\xi(m_t) - q_t(n_t)\|_2^2 \right). \end{aligned}$$

By the same reasoning as in the proof of Theorem 1 the assertion follows for every  $n_t \in \mathcal{N}'_t$  by considering and minimizing every map

$$q \mapsto \sum_{m_t \in \mathcal{N}_t} \pi(m_t, n_t) \cdot \|\xi(m_t) - q\|_2^2$$

separately. □

### 5.3. The overall algorithm

As it is not possible to improve the probabilities and solve the facility location problem in one single step Algorithm 2 describes the course of action. Starting with an initial guess for the quantizers (resp. the scenario paths) and using the related transport probabilities  $\pi^0$  the algorithm iterates between improving the quantizers (Step 2) and improving the transport probabilities (Step 3). Step 2 goes backward in time and uses conditional versions  $\text{dl}_r^{k+1}(m, n)$  of the nested distance, which are related to nodes  $m$  and  $n$ , in order to resemble an approximation of the full nested distance. To improve the locations  $q$ , Step 3 either uses classical optimization algorithms for the general case, or a version of the k-means algorithm in the important case of the quadratic nested distance.

The algorithm leads to an improvement in each iteration step (Theorem 2 and Theorem 3) and converges in finitely many steps.

**Theorem 4.** *Provided that the minimization (23) can be done exactly – as is the case for the quadratic nested distance – Algorithm 2 terminates at a stationary  $\text{dl}_r(P, P^{k^*})$  after finitely many iterations  $k^*$ .*

*Proof.* It is possible – although very inadvisable for computational purposes – to rewrite the computation of  $\text{dl}_r^{k+1}(0, 0)$  in Algorithm 2 as a single linear program of the form

$$\begin{aligned} &\text{minimize} && c(\pi^{k+1} | \pi^k) \\ &\text{in } \pi^{k+1} && \\ &\text{subject to} && A\pi^{k+1} = b, \\ & && \pi^{k+1} \geq 0, \end{aligned}$$

where the matrix  $A$  and the vector  $b$  collect all linear conditions from (20), and  $\pi \mapsto c(\pi | \tilde{\pi})$  is multilinear. Note that the constraints  $\Pi := \{\pi : A\pi = b, \pi \geq 0\}$  form a convex polytope, which is independent of the iterate  $\pi^k$ . Without loss of generality one may assume that  $\pi^k$  is an edge of the polytope  $\Pi$ . Because  $\Pi$  has finitely many edges and each edge  $\pi \in \Pi$  can be associated with a unique quantization scenario  $q(\pi)$ , by assumption it is clear that the decreasing sequence

$$\text{dl}_r^{k+2}(\mathbb{P}, \mathbb{P}^{k+2}) = c(\pi^{k+2} | \pi^{k+1}) \leq c(\pi^{k+1} | \pi^k) = \text{dl}_r^{k+1}(\mathbb{P}, \mathbb{P}^{k+1})$$

cannot improve further whenever the stationary point is met. □



**Algorithm 2**

Sequential improvement of the measure  $P^k$  to approximate  $P = \sum_i p_i \delta_{\xi_i}$  in the nested distance on the trees  $(\mathcal{F}_t)_{t \in \{0, \dots, T\}}$  ( $(\mathcal{F}'_t)_{t \in \{0, \dots, T\}}$ , resp.).

**Step 1 – Initialization**

Set  $k \leftarrow 0$ , and let  $q^0$  be process quantizers with related transport probabilities  $\pi^0(i, j)$  between scenario  $i$  of the original  $\mathbb{P}$ -tree and scenario  $q_j^0$  of the approximating  $\mathbb{P}'$ -tree;  $\mathbb{P}^0 := \mathbb{P}'$ .

**Step 2 – Improve the quantizers**

Find improved quantizers  $q_j^{k+1}$ :

- In case of the quadratic Wasserstein distance (Euclidean distance and Wasserstein of order  $r = 2$ ) set

$$q^{k+1}(n_t) := \sum_{m_t \in \mathcal{N}_t} \frac{\pi^k(m_t, n_t)}{\sum_{m_t \in \mathcal{N}_t} \pi^k(m_t, n_t)} \cdot \xi_t(m_t), \quad (24)$$

- or solve (23), for example by applying the steepest descent method, or the limited memory BFGS method.

**Step 3 – Improve the probabilities**

Setting  $\pi \leftarrow \pi^k$  and  $q \leftarrow q^{k+1}$  use (18), (19), (20) and (21) to calculate all conditional probabilities  $\pi^{k+1}(\cdot, \cdot | m, n) = \pi^*(\cdot, \cdot | m, n)$ , the unconditional transport probabilities  $\pi^{k+1}(\cdot, \cdot)$  and the distance  $\mathbf{d}_r^{k+1}(0, 0) = \mathbf{d}_r(0, 0)$ .

**Step 4**

Set  $k \leftarrow k + 1$  and continue with **Step 2** if

$$\mathbf{d}_r^{k+1}(0, 0) < \mathbf{d}_r^k(0, 0) - \varepsilon,$$

where  $\varepsilon > 0$  is the desired improvement in each cycle  $k$ .

Otherwise, set  $q^* \leftarrow q^k$ , define the measure

$$P^{k+1} := \sum_j \delta_{q_j^{k+1}} \cdot \sum_i \pi^{k+1}(i, j),$$

for which  $\mathbf{d}_r(\mathbb{P}, \mathbb{P}^{k+1}) = \mathbf{d}_r^{k+1}(0, 0)$  and stop.

*Remark.* In case of the quadratic nested distance ( $r = 2$ ) and the Euclidean distance the choice  $\varepsilon = 0$  is possible.

## 5.4. A numerical example and derived applications

To illustrate the results we have implemented all steps of the discussed algorithms in MATLAB<sup>®</sup>. All LPs were solved using the function `linprog`. It is a central observation that optimization for Euclidean norms and the quadratic Wasserstein distance is fastest. This is because the facility location problem can be avoided and replaced by computing the conditional expectation in a direct way. Moreover, when applying the methods, it was a repeated pattern that the first few iteration steps improve the distance significantly, whereas following steps just give minor improvements of the objective.

The computation times collected in Table 1 have been noted for an iteration step in Algorithm 2 on a customary, standard laptop.

Stages	4	5	5	6	* 7	7
Nodes of the initial tree	53	309	188	1,365	1,093	2,426
Nodes of the approximating tree	15	15	31	63	127	127
Time/ sec.	1	10	4	160	157	1,044

Table 1: Time to perform an iteration in Algorithm 2.  
The example indicated by the asterisk (\*) corresponds to Figure 6.

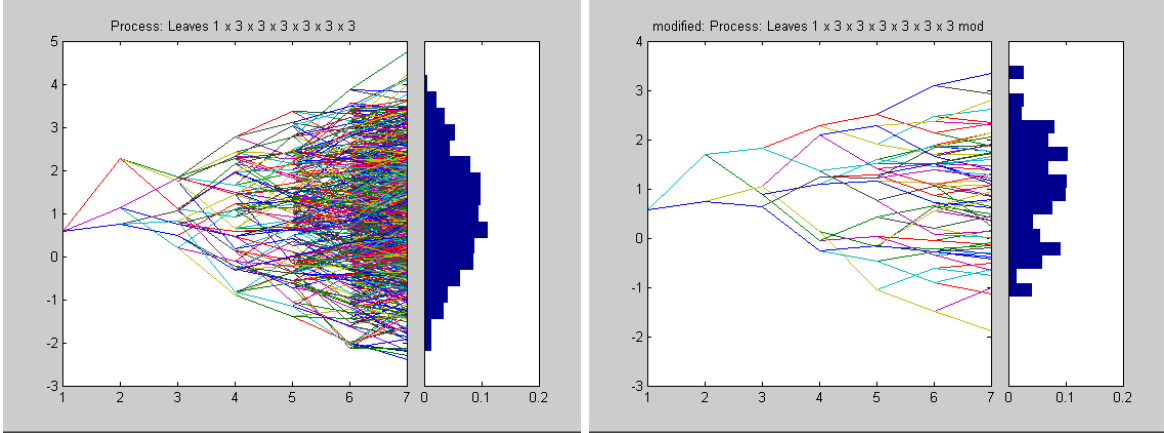


Figure 6: The initial tree with 1093 nodes at 7 stages (left) and a binary, approximating tree, which has 127 nodes (right). Their nested distance is 2.32. The tree structure is depicted, annotated is the histogram of the paths.

Figure 6 exemplary depicts the situation of the latter example with 7 stages. The computed distance of 2.32 allows the rough interpretation, that the scenarios of the initial tree – on average – can be squeezed into a “pipe” of radius  $2.32/7 = 0.3$  along a branch of the approximating tree.

## 6. Summary and outlook

In this paper we address the problem of approximating stochastic processes in discrete time by trees, which are discrete stochastic processes. For this purpose we build on the recently introduced nested distances, generalizations of the well known Wasserstein or Kantorovich distances. In addition to their properties as classical probability metrics they are able to account for the effects of filtrations related to stochastic processes.

In particular we use the nested distance to compare trees, which are important tools for discretizing stochastic optimization problems. The aim is to reduce the distance between a given – usually large – tree, and a smaller tree supposed to approximate the given tree. This problem is of fundamental interest in stochastic programming, as the number of variables of the initial process can be reduced significantly by the techniques and algorithms proposed.

The paper analyzes the relations between processes and trees, reviews the main properties of Wasserstein distances and nested distances and finally proposes and analyzes an iterative algorithm for improving the nested distance between trees. For the important special case of nested distances of order 2 based on Euclidean distances the algorithm can be enhanced by using k-means clustering in order to improve calculation speed.

While first numerical experiences are encouraging, some interesting issues have to be approached in future research: As an example the speed of the algorithm could be further increased by parallelization, as in its Step 3 many conditional distances can be calculated independently and in parallel for each

stage. Furthermore, we will aim at extending the algorithm to improve distances directly between stochastic processes and an approximating tree.

## 7. Acknowledgment

We wish to express our gratitude to Prof. Georg Ch. Pflug for his continual advice. We thank the referees for their constructive criticism.

## 8. Bibliography

- [Boy71] Edward S. Boylan. Epiconvergence of martingales. *Ann. Math. Statist.*, 42:552–559, 1971. [2](#)
- [BPP05] V. Bally, G. Pagès, and J. Printems. A quantization tree method for pricing and hedging multidimensional american options. *Mathematical Finance*, 15(1):119–168, 2005. [7](#)
- [DGKR03] Jitka Dupačová, Nicole Gröwe-Kuska, and Werner Römisch. Scenario reduction in stochastic programming. *Mathematical Programming, Ser. A*, 95(3):493–511, 2003. [2](#), [6](#), [7](#)
- [DH02] Z. Drezner and H. W. Hamacher. *Facility Location: Applications and Theory*. Springer, New York, NY, 2002. [8](#)
- [Dur04] Richard A. Durrett. *Probability. Theory and Examples*. Duxbury Press, Belmont, CA, second edition, 2004. [14](#)
- [GL00] Siegfried Graf and Harald Luschgy. *Foundations of Quantization for Probability Distributions*, volume 1730 of *Lecture Notes in Mathematics*. Springer-Verlag Berlin Heidelberg, 2000. [7](#)
- [GS02] A. L. Gibbs and F. E. Su. On choosing and bounding probability metrics. *International Statistical Review*, 70(3):419–435, 2002. [6](#)
- [HR09] Holger Heitsch and Werner Römisch. Scenario tree modeling for multistage stochastic programs. *Math. Program. Ser. A*, 118:371–406, 2009. [2](#)
- [HRS06] Holger Heitsch, Werner Römisch, and Cyrille Strugarek. Stability of multistage stochastic programs. *SIAM J. Optimization*, 17(2):511–525, 2006. [10](#)
- [HW01] Kjetil Høyland and Stein W. Wallace. Generating scenario trees for multistage decision problems. *Management Science*, 47:295–307, 2001. [2](#)
- [Kla02] Pieter Klaassen. Comment on "generating scenario trees for multistage decision problems". *Management Science*, 45(11):1512–1516, Nov. 2002. [2](#)
- [Kud74] Hirokichi Kudō. A note on the strong convergence of  $\sigma$ -algebras. *Ann. Probability*, 2:76–83, 1974. [2](#)
- [Llo82] Stuart P. Lloyd. Least square quantization in PCM. *IEEE Transactions of Information Theory*, 28(2):129–137, 1982. [8](#)
- [Mon81] Gaspard Monge. Mémoire sue la théorie des déblais et de remblais. *Histoire de l'Académie Royale des Sciences de Paris, avec les Mémoires de Mathématique et de Physique pour la même année*, pages 666–704, 1781. [6](#)
- [Noc80] Jorge Nocedal. Updating quasi-Newton matrices with limited storage. *Mathematics of Computation*, 35(151):773–782, 1980. [15](#)

- [Pfl09] Georg Ch. Pflug. Version-independence and nested distribution in multistage stochastic optimization. *SIAM Journal on Optimization*, 20:1406–1420, 2009. [1](#), [2](#), [10](#)
- [PP12] Georg Ch. Pflug and Alois Pichler. A distance for multistage stochastic optimization models. *SIAM Journal on Optimization*, 22(1):1–23, 2012. [10](#), [11](#), [12](#)
- [PR07] Georg Ch. Pflug and Werner Römisch. *Modeling, Measuring and Managing Risk*. World Scientific, River Edge, NJ, 2007. [3](#)
- [Rac91] Svetlozar T. Rachev. *Probability metrics and the stability of stochastic models*. John Wiley and Sons Ltd., West Sussex PO19, 1UD, England, 1991. [6](#)
- [SDR09] Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on Stochastic Programming*. MPS-SIAM Series on Optimization 9, 2009. [1](#)
- [Shi96] Albert Nikolayevich Shiryaev. *Probability*. Springer, New York, 1996. [5](#)
- [Ver06] Anatoly M. Vershik. Kantorovich metric: Initial history and little-known applications. *Journal of Mathematical Sciences*, 133(4):1410–1417, 2006. [6](#)
- [Vil03] Cédric Villani. *Topics in Optimal Transportation*, volume 58 of *Graduate Studies in Mathematics*. American Mathematical Society, 2003. [6](#)