

## REMARKS AND OMITTED PROOFS\*

**Abstract.** We provide remarks and omitted proofs from “On the Connection Between Compression Learning and Scenario Based Single-Stage and Cascading Optimization Problems”, Margellos, Prandini & Lygeros, IEEE Trans. Autom. Control, 60(10), 2015.

**1. Proof of Theorem 1.** The proof of this theorem is based on the proof of Theorem 5 in [1], and Theorem 1 in [7]. Without loss of generality fix  $I_d = \{1, \dots, d\} \in \mathcal{I}_d$  and let  $(\delta_1, \dots, \delta_d)$  in  $\bar{\Delta}^d = \{(\delta_1, \dots, \delta_d) \in \Delta^d : d_{\mathbb{P}}(T, H_{I_d}) > \epsilon\}$ , where  $\epsilon \in (0, 1)$ . We have that

$$(1.1) \quad \begin{aligned} \mathbb{P}\left\{\delta \in \Delta : H_{I_d} \text{ is consistent with } (\delta, \mathbb{1}_T(\delta)) \text{ and } d_{\mathbb{P}}(T, H_{I_d}) > \epsilon\right\} \\ = \mathbb{P}\left\{\delta \in \Delta : H_{I_d} \text{ is consistent with } (\delta, \mathbb{1}_T(\delta))\right\} = 1 - d_{\mathbb{P}}(T, H_{I_d}) \leq 1 - \epsilon, \end{aligned}$$

where the first step follows from the fact that  $d_{\mathbb{P}}(T, H_{I_d})$  does not depend on  $\delta$  but only on  $(\delta_1, \dots, \delta_d) \in \bar{\Delta}^d$ , and the second step follows from the definition of a consistent hypothesis (Definition 2). In fact, the inequality in (1.1) is strict as  $d_{\mathbb{P}}(T, H_{I_d}) > \epsilon$ . Since the samples are extracted independently we have that

$$(1.2) \quad \begin{aligned} \mathbb{P}^{m-d}\left\{(\delta_{d+1}, \dots, \delta_m) \in \Delta^{m-d} : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=d+1}^m \text{ and } d_{\mathbb{P}}(T, H_{I_d}) > \epsilon\right\} \\ = \mathbb{P}^{m-d}\left\{(\delta_{d+1}, \dots, \delta_m) \in \Delta^{m-d} : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=d+1}^m\right\} \\ = \prod_{j=d+1}^m \mathbb{P}\left\{\delta_j \in \Delta : H_{I_d} \text{ is consistent with } (\delta_j, \mathbb{1}_T(\delta_j))\right\} \leq (1 - \epsilon)^{m-d}, \end{aligned}$$

where the first equality is due to the fact that  $d_{\mathbb{P}}(T, H_{I_d})$  does not depend on  $\delta_{d+1}, \dots, \delta_m$ , and since by the theorem's hypothesis,  $H_{I_d}$  is consistent with  $\{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in I_d}$ . However, (1.2) can be rewritten by the following conditional probability as  $\bar{\Delta}^d$  is a cylinder set with base the cartesian product of the domains of the first  $d$  multi-sample elements (since it is independent of  $\{\delta_i\}_{i=d+1}^m$ ), i.e.,

$$(1.3) \quad \mathbb{P}^m\left\{\{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m\} \mid \{(\delta_1, \dots, \delta_d) \in \bar{\Delta}^d\}\right\} \leq (1 - \epsilon)^{m-d}.$$

Deconditioning, we get

$$(1.4) \quad \begin{aligned} \mathbb{P}^m\left\{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m \text{ and } d_{\mathbb{P}}(T, H_{I_d}) > \epsilon\right\} \\ \leq \int_{\bar{\Delta}^d} (1 - \epsilon)^{m-d} \mathbb{1}_{\{(\delta_1, \dots, \delta_d) \in \bar{\Delta}^d\}} d\mathbb{P}^d(\{\delta_i\}_{i \in I_d}) \leq (1 - \epsilon)^{m-d}. \end{aligned}$$

The left-hand side of (2) in the statement of Theorem 1 can be then expressed as follows

$$(1.5) \quad \begin{aligned} \mathbb{P}^m\left\{\bigcup_{I_d \in \mathcal{I}_d} \{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m \text{ and } d_{\mathbb{P}}(T, H_{I_d}) > \epsilon\}\right\} \\ \leq \sum_{I_d \in \mathcal{I}_d} \mathbb{P}^m\left\{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m \text{ and } d_{\mathbb{P}}(T, H_{I_d}) > \epsilon\right\} \\ \leq \binom{m}{d} (1 - \epsilon)^{m-d}, \end{aligned}$$

concluding the proof (the first inequality is by the subadditivity of  $\mathbb{P}^m$  and the second one due to (1.4)). If  $\bar{\Delta}^d$  was empty, then the corresponding term in the summation would be zero, thus bounded by  $(1 - \epsilon)^{m-d}$ .

---

\*28.01.2022

**2. Proof of Theorem 3.** The proof of this theorem is based on the first part of the proof of Theorem 1 in [8]. Without loss of generality fix  $I_d = \{1, \dots, d\} \in \mathcal{I}_d$ . Let  $V : \Delta^d \rightarrow [0, 1]$  be a random variable such that  $V(\{\delta_i\}_{i=1}^d) = d_{\mathbb{P}}(T, H_{I_d})$  for all  $\{\delta_i\}_{i=1}^d \in \Delta^d$ . For  $v \in [0, 1]$  denote its probability distribution function by  $F(v) = \mathbb{P}^d\{\{\delta_i\}_{i=1}^d \in \Delta^d : d_{\mathbb{P}}(T, H_{I_d}) \leq v\} = \mathbb{P}^d\{V^{-1}([0, v])\} = \mathbb{P}_V^d\{V \leq v\}$ , where  $\mathbb{P}_V^d$  denotes the image probability of  $\mathbb{P}^d$  through  $V$ , defined over the Borel  $\sigma$ -algebra on  $[0, 1]$  (see Chapter 2 in [Campi, Selected Topics in Probability, Lecture Notes, 2008]). We used that, for  $a < b$ ,  $V^{-1}([a, b]) = \{\{\delta_i\}_{i=1}^d \in \Delta^d : V(\{\delta_i\}_{i=1}^d) \in [a, b]\}$ .

By the first part of Assumption 3 and due to the fact that for a fixed  $(\delta_1, \dots, \delta_d) \in \Delta^d$ ,  $d_{\mathbb{P}}(T, H_{I_d}) = V(\{\delta_i\}_{i=1}^d)$ , (1.3) becomes

$$(2.1) \quad \mathbb{P}^m\left\{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m\right\} \Big| \{(\delta_1, \dots, \delta_d) \in \Delta^d\} \\ = \left(1 - V(\{\delta_i\}_{i=1}^d)\right)^{m-d}.$$

Deconditioning and using the equivalences between a probability and its image, (2.1) yields then

$$(2.2) \quad \mathbb{P}^m\left\{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m\right\} \\ = \int_{\Delta^d} \left(1 - V(\{\delta_i\}_{i=1}^d)\right)^{m-d} d\mathbb{P}^d(\{\delta_i\}_{i \in I_d}) = \int_0^1 (1-v)^{m-d} dF(v).$$

By Assumption 3 the sets  $\left\{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m\right\}_{I_d \in \mathcal{I}_d}$  form a partition of  $\Delta^m$  up to a set of measure zero<sup>1</sup>. Since there are  $\binom{m}{d}$  sets in  $\mathcal{I}_d$ , we have that  $\binom{m}{d} \mathbb{P}^m\left\{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m\right\} = 1$ , which due to (2.2) leads to

$$(2.3) \quad \binom{m}{d} \int_0^1 (1-v)^{m-d} dF(v) = 1, \text{ for all } m \geq d \implies F(v) = v^d,$$

where the  $F(v)$  above is the unique solution satisfying the previous identities. This is due to [8], based on the fact that this is a moment problem for a distribution with finite support ( $[0, 1]$  here); see, e.g., Chapter 2, Section 12.9, Corollary 1 of [A. Shiryayev, Probability, Springer, 1996]. It can be verified using integration by parts. In place of (1.4) we thus have

$$(2.4) \quad \mathbb{P}^m\left\{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m \text{ and } d_{\mathbb{P}}(T, H_{I_d}) > \epsilon\right\} \\ = \int_{\Delta^d} \mathbb{P}^m\left\{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m\right\} \\ \Big| \{(\delta_1, \dots, \delta_d) \in \Delta^d : d_{\mathbb{P}}(T, H_{I_d}) > \epsilon\} \} \mathbb{1}_{\{V^{-1}((\epsilon, 1])\}} d\mathbb{P}^d(\{\delta_i\}_{i \in I_d}) \\ = \int_{\Delta^d} \left(1 - V(\{\delta_i\}_{i=1}^d)\right)^{m-d} \mathbb{1}_{\{V^{-1}((\epsilon, 1])\}} d\mathbb{P}^d(\{\delta_i\}_{i \in I_d}) \\ = \int_{\epsilon}^1 (1-v)^{m-d} dF(v) = \int_{\epsilon}^1 (1-v)^{m-d} dv^{d-1} dv,$$

where we used (2.1)<sup>2</sup>, since we still condition on some fixed  $\delta_1, \dots, \delta_d$ , and (2.3) as  $dF(v) = dv^{d-1} dv$ .

<sup>1</sup>To see this, for all  $I_d \in \mathcal{I}_d$  let  $S_{I_d} = \{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m\}$ . We first show that  $\mathbb{P}^m\{(\delta_1, \dots, \delta_m) \in \Delta^m : \Delta^m \setminus \cup_{I_d \in \mathcal{I}_d} S_{I_d}\} = 0$ . Equivalently,  $\cup_{I_d \in \mathcal{I}_d} S_{I_d} = \Delta^m$  up to a set of measure zero. Since  $\cup_{I_d \in \mathcal{I}_d} S_{I_d} \subseteq \Delta^m$ , we show that  $\cup_{I_d \in \mathcal{I}_d} S_{I_d} \supseteq \Delta^m$ , i.e. if  $(\delta_1, \dots, \delta_m) \in \Delta^m$  then there exists  $I_d \in \mathcal{I}_d$  such that  $(\delta_1, \dots, \delta_m) \in S_{I_d}$ . With  $\mathbb{P}^m$ -probability one, the last statement follows from Assumption 3 and the definition of  $S_{I_d}$ . We now show that  $S_{I_d^1} \cap S_{I_d^2} = \emptyset$  for all  $I_d^1, I_d^2 \in \mathcal{I}_d$  with  $I_d^1 \neq I_d^2$ . For the sake of contradiction assume that there exist  $I_d^1, I_d^2 \in \mathcal{I}_d$  with  $I_d^1 \neq I_d^2$  such that  $S_{I_d^1} \cap S_{I_d^2} \neq \emptyset$ . By the definition of  $S_{I_d^1}, S_{I_d^2}$ , this implies that there exists  $(\delta_1, \dots, \delta_m) \in \Delta^m$  such that both  $H_{I_d^1}$  and  $H_{I_d^2}$  are consistent with  $\{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m$ . However, by Assumption 3 a unique hypothesis consistent with respect to the  $m$ -multisample should exist almost surely, establishing a contradiction up to a measure zero set.

<sup>2</sup>Alternatively, we could see this by setting  $A = \{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m\}$ ,  $B = (\epsilon, 1]$ , and writing (2.4) as  $\mathbb{P}^m\{A \cap \{V \in B\}\} = \int_B \mathbb{P}^m\{A|V = v\} dF(v)$  (see [A. Shiryayev, Probability, Springer, 1996], Chapter 2, Section 7.5, eq. (17)). Noticing then that due to (2.1),  $\mathbb{P}^m\{A|V = v\} = (1-v)^{m-d}$  the last statement yields the result in (2.4).

Notice now that the first inequality in (1.5) holds with equality since the union in the left-hand side of (1.5) is disjoint as by Assumption 3 the associated sets form a partition of  $\Delta^m$  almost surely. We thus have

$$(2.5) \quad \mathbb{P}^m \left\{ \bigcup_{I_d \in \mathcal{I}_d} \{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{I_d} \text{ is consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m \text{ and } d_{\mathbb{P}}(T, H_{I_d}) > \epsilon\} \right\} \\ = \binom{m}{d} \int_{\epsilon}^1 (1-v)^{m-d} dv^{d-1} dv = \sum_{i=0}^{d-1} \binom{m}{i} \epsilon^i (1-\epsilon)^{m-i},$$

where the last equality follows by repeated integration by parts [8], and is the cumulative distribution of a binomial random variable. By the second part of Assumption 3, there exists a unique  $m_d(\delta_i, \mathbb{1}_T(\delta_i)) \in \mathcal{I}_d$  such that  $H_{m_d}$  is a consistent hypothesis. Eq. (2.5) is thus  $\mathbb{P}^m \{(\delta_1, \dots, \delta_m) \in \Delta^m : d_{\mathbb{P}}(T, H_{m_d}) > \epsilon\}$ , concluding the proof.

**3. Proof of Theorem 4.** The proof of this theorem is based on the proof of Theorem 2.1 in [20]. Fix any  $r \in \mathbb{N}$  and  $I_r \in \mathcal{I}_r$ , and let  $m_d^r \left( \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in \{1, \dots, m\} \setminus I_r} \right)$  be a set of  $d$  indices satisfying the first two parts of Assumption 4. Letting  $V(\{\delta_i\}_{i \in \{1, \dots, m\} \setminus I_r}) = d_{\mathbb{P}}(T, H_{m_d^r})$ , we have

$$(3.1) \quad \mathbb{P} \left\{ \delta \in \Delta : H_{m_d^r} \text{ is not consistent with } (\delta, \mathbb{1}_T(\delta)) \right\} = d_{\mathbb{P}}(T, H_{m_d^r}) = V(\{\delta_i\}_{i \in \{1, \dots, m\} \setminus I_r}).$$

Since the samples are extracted independently we have that

$$(3.2) \quad \mathbb{P}^r \left\{ \{\delta_j\}_{j \in I_r} \in \Delta^r : H_{m_d^r} \text{ is not consistent with } \{(\delta_j, \mathbb{1}_T(\delta_j))\}_{j \in I_r} \right\} \\ = \prod_{j \in I_r} \mathbb{P} \left\{ \delta_j \in \Delta : H_{m_d^r} \text{ is not consistent with } (\delta_j, \mathbb{1}_T(\delta_j)) \right\} = V(\{\delta_i\}_{i \in \{1, \dots, m\} \setminus I_r})^r.$$

Let  $\bar{F}(v) = \mathbb{P}^{m-r} \{ \{\delta_i\}_{i \in \{1, \dots, m\} \setminus I_r} \in \Delta^{m-r} : d_{\mathbb{P}}(T, H_{m_d^r}) \leq v \}$  denote the distribution function of  $V$ . Similarly to (2.4) in the proof of Theorem 3, we then have that

$$(3.3) \quad \mathbb{P}^m \left\{ (\delta_1, \dots, \delta_m) \in \Delta^m : H_{m_d^r} \text{ is not consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in I_r} \text{ and } d_{\mathbb{P}}(T, H_{m_d^r}) > \epsilon \right\} \\ = \int_{\epsilon}^1 \mathbb{P}^m \left\{ \{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{m_d^r} \text{ is not consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in I_r} \} \right. \\ \left. \mid V(\{\delta_i\}_{i \in \{1, \dots, m\} \setminus I_r}) = v \right\} d\bar{F}(v) \\ = \int_{\epsilon}^1 v^r d\bar{F}(v).$$

Construct the algorithm  $\{A_{m-r}\}_{m-r \geq d}$ , such that  $A_{m-r}(\{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in \{1, \dots, m\} \setminus I_r}) = H_{m-r} = H_{m_d^r}$ . By Theorem 1, since  $m_d^r$  satisfies the second part of Assumption 4, with  $m-r$  in place of  $m$  and  $v$  in place of  $\epsilon$ , we have that the constructed algorithm is PAC-T, hence for any  $v \in (0, 1)$ ,

$$(3.4) \quad \mathbb{P}^{m-r} \left\{ \{\delta_i\}_{i \in \{1, \dots, m\} \setminus I_r} \in \Delta^{m-r} : d_{\mathbb{P}}(T, H_{m_d^r}) > v \right\} \leq \binom{m-r}{d} (1-v)^{m-r-d} \\ \implies \bar{F}(v) \geq F(v) := 1 - \binom{m-r}{d} (1-v)^{m-r-d}.$$

We then have that

$$(3.5) \quad \int_{\epsilon}^1 v^r d\bar{F}(v) \leq \int_{\epsilon}^1 v^r dF(v) = \int_{\epsilon}^1 (m-r-d) \binom{m-r}{d} v^r (1-v)^{m-r-d-1} dv \\ = \binom{m-r}{d} \frac{1}{\binom{m-d}{r}} \sum_{i=0}^r \binom{m-d}{i} \epsilon^i (1-\epsilon)^{m-d-i},$$

where the first equality is due to the definition of  $F(v)$ , and the second one follows by repeated integration by parts. To see the inequality, notice that

$$(3.6) \quad \int_{\epsilon}^1 v^r d\bar{F}(v) = 1 - \epsilon^r \bar{F}(\epsilon) - \int_{\epsilon}^1 \bar{F}(v) r v^{r-1} dv \\ \leq 1 - \epsilon^r F(\epsilon) - \int_{\epsilon}^1 F(v) r v^{r-1} dv = \int_{\epsilon}^1 v^r dF(v),$$

where the equality follows from Theorem 11, Chapter 2, Section 6.11 of [A. Shiryayev, Probability, Springer, 1996], since  $v^r$  is an increasing function of  $v$ , hence it is treated as a generalized distribution (see also [20]). The inequality is since  $\bar{F}(v) \geq F(v)$  by (3.4). Intuitively,  $\bar{F}(v) \geq F(v)$  is a dominance condition for cumulative distributions, implying that  $F$  is concentrated to higher values compared to  $\bar{F}$  (if they admit a density, the density of  $F$  is also concentrated to higher values); hence, the associated expected value ( $\int_{\epsilon}^1 v^r dF(v)$ ) is related to the expected value of an increasing function of  $v$ ) would be higher compared to the one corresponding to  $\bar{F}$ .

Consider the third part of Assumption 4. Denote then by  $\bar{I}_r \in \mathcal{I}_r$  the associated set of indices, and let  $\bar{m}_d^r(\{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in \{1, \dots, m\} \setminus \bar{I}_r})$  be a set of  $d$  indices such that  $H_{\bar{m}_d^r} = G_d(\{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in \bar{m}_d^r})$  is not consistent with  $\{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in \bar{I}_r}$ . Since this happens with  $\mathbb{P}^m$ -probability one, we have the first inequality (this is not equality since  $\bar{I}_r$  is not necessarily unique) below, i.e.,

$$(3.7) \quad \mathbb{P}^m \left\{ (\delta_1, \dots, \delta_m) \in \Delta^m : d_{\mathbb{P}}(T, H_{\bar{m}_d^r}) > \epsilon \right\} \\ \leq \mathbb{P}^m \left\{ \bigcup_{I_r \in \mathcal{I}_r} \{(\delta_1, \dots, \delta_m) \in \Delta^m : H_{m_d^r} \text{ is not consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in I_r} \text{ and } d_{\mathbb{P}}(T, H_{m_d^r}) > \epsilon\} \right\} \\ \leq \binom{m}{r} \binom{m-r}{d} \frac{1}{\binom{m-d}{r}} \sum_{i=0}^r \binom{m-d}{i} \epsilon^i (1-\epsilon)^{m-d-i} = \binom{m}{d} \sum_{i=0}^r \binom{m-d}{i} \epsilon^i (1-\epsilon)^{m-d-i},$$

where the second inequality is by the subadditivity of  $\mathbb{P}^m$  since there are  $\binom{m}{r}$  sets in the union, and by (3.3) and (3.5). Set  $q(m, \epsilon) = \binom{m}{d} \sum_{i=0}^r \binom{m-d}{i} \epsilon^i (1-\epsilon)^{m-d-i}$ . Similarly to the last part of the proof of Theorem 2,  $\lim_{m \rightarrow \infty} q(m, \epsilon) = 0$ . Construct then algorithm  $\{A_m\}_{m \geq d+r}$ , where  $A_m : [\Delta \times \{0, 1\}]^m \rightarrow \mathcal{D}$  takes as input a labeled  $m$ -multisample and returns a hypothesis  $H_m = A_m(\{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i=1}^m)$  such that  $H_m = H_{\bar{m}_d^r}$ . By Definition 4, algorithm  $\{A_m\}_{m \geq d+r}$  is PAC-T, thus concluding the proof.

**4. Tightening Theorem 4.** If we strengthen Assumption 4 by further requiring that the set  $I_d \in \mathcal{I}_d^{m-r}$  that satisfies its conditions is unique, we could replace the right-hand side of (3.7) with  $\binom{r+d-1}{r} \sum_{i=0}^{r+d-1} \binom{m}{i} \epsilon^i (1-\epsilon)^{m-i}$ . The proof of this fact is identical to the one of Theorem 4, with the difference that (3.4) is replaced by  $\bar{F}(v) = F(v) := 1 - \sum_{i=0}^{d-1} \binom{m-r}{i} v^i (1-v)^{m-r-i}$ . Eq. (3.5) would then hold with equality, however, the final result would still hold with inequality due to the inequalities in (3.7).

Specializing this in a convex scenario optimization context, adopt an algorithm that removes the samples with indices in  $\bar{I}_r$ , and generates a hypothesis  $H_{\bar{m}_d^r}$  such that  $d_{\mathbb{P}}(T, H_{\bar{m}_d^r})$  is the probability of constraint violation. We further require that  $H_{\bar{m}_d^r}$  is inconsistent with the removed samples (i.e., the associated minimizer violates the constraints for the removed samples), and for any given  $(\delta_1, \dots, \delta_m)$ ,  $d_{\mathbb{P}}(T, H_{\bar{m}_d^r}) \leq d_{\mathbb{P}}(T, H_{m_d^r})$  for any set  $I_r$  in the family of sets of inconsistent hypotheses<sup>3</sup>. In words, we consider a conservative design, removing samples that lead to the lowest probability of constraint violation. We then have that,

$$(4.1) \quad \binom{r+d-1}{r} \mathbb{P}^m \left\{ (\delta_1, \dots, \delta_m) \in \Delta^m : d_{\mathbb{P}}(T, H_{\bar{m}_d^r}) > \epsilon \right\} \\ = \int_{\Delta^m} \mathbb{1}_{\{(\delta_1, \dots, \delta_m) : d_{\mathbb{P}}(T, H_{\bar{m}_d^r}) > \epsilon\}} \sum_{I_r \in \mathcal{I}_r} \mathbb{1}_{\{(\delta_1, \dots, \delta_m) : H_{m_d^r} \text{ is not consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in I_r}\}} d\mathbb{P}^m(\{\delta_i\}_{i=1}^m) \\ \leq \int_{\Delta^m} \sum_{I_r \in \mathcal{I}_r} \mathbb{1}_{\{(\delta_1, \dots, \delta_m) : H_{m_d^r} \text{ is not consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in I_r}\}} \mathbb{1}_{\{(\delta_1, \dots, \delta_m) : d_{\mathbb{P}}(T, H_{m_d^r}) > \epsilon\}} d\mathbb{P}^m(\{\delta_i\}_{i=1}^m) \\ = \sum_{I_r \in \mathcal{I}_r} \mathbb{P}^m \left\{ (\delta_1, \dots, \delta_m) \in \Delta^m : H_{m_d^r} \text{ is not consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in I_r} \text{ and } d_{\mathbb{P}}(T, H_{m_d^r}) > \epsilon \right\} \\ = \binom{r+d-1}{r} \sum_{i=0}^{r+d-1} \binom{m}{i} \epsilon^i (1-\epsilon)^{m-i},$$

<sup>3</sup>For the class of fully-supported scenario programs, for each  $(\delta_1, \dots, \delta_m)$ ,  $\binom{r+d-1}{r}$  denotes the number of sets of  $r$  samples, such that the minimizer obtained upon removing the samples in any of these sets violates the constraints corresponding to all  $r$  samples of that set; see Theorem 4.1 in [Gaertner & Welzl, "A Simple Sampling Lemma: Analysis and Applications in Geometric Optimization", Discrete Comput. Geom. vol. 35, pp. 569-590, 2001], and Theorem 2.3 in [Matusek, "On Geometric Optimization with Few Violated Constraints", Discrete Comput. Geom. vol. 14, pp. 365-384, 1995]. Figure 4.1 provides a pictorial illustration of this fact.

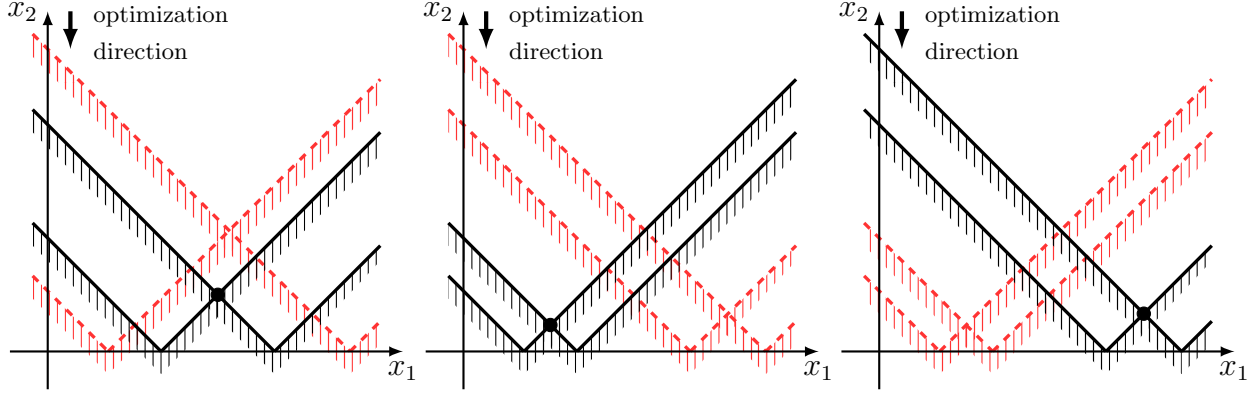


FIG. 4.1. Pictorial example of a convex scenario program with  $n_x = 2$  decision variables,  $x_1, x_2 \in \mathbb{R}$ , and  $m = 4$  samples giving rise to V-shaped constraints. The downwards pointing arrow indicates the optimization direction, i.e., minimizing  $x_2$ , while the feasibility region is outside the shaded area. Assume that each V-shaped constraint is a translation of the other with its vertex being sampled along the  $x_1$ -axis from a distribution that admits a density, thus rendering the scenario program non-degenerate and in fact fully-supported with  $\mathbb{P}^m$ -probability one, i.e.,  $d = n_x = 2$ . We discard  $r = 2$  samples/constraints (the ones indicated by “dashed-red”), and the resulting solution drifts to the point indicated by the “dot”. There are exactly  $\binom{r+d-1}{r} = 3$  sets of  $r = 2$  samples that can be discarded, such that the minimizer obtained upon the sample removal violates the constraints corresponding to the discarded samples, i.e., the associated hypothesis is not consistent with the two samples discarded each time. These three sets of discarded samples that exhibit this property are depicted in the three figure panels.

where the first equality is due to the definition of a probability as an integral, and since for each  $(\delta_1, \dots, \delta_m)$  the summation inside the integral is equal to  $\binom{r+d-1}{r}$  (by footnote 3, for each multisample this is exactly the number of nonzero terms in that summation), while the inequality is due to the fact that for any given  $(\delta_1, \dots, \delta_m)$ , if  $d_{\mathbb{P}}(T, H_{\bar{m}_d^r}) > \epsilon$  then  $d_{\mathbb{P}}(T, H_{\bar{m}_d^r}) > \epsilon$ . The second last equality is by the interplay between integral and probability, and the last one is due to (3.7) specialized to the case where Assumption 4 is strengthened. By (4.1),

$$(4.2) \quad \mathbb{P}^m \left\{ (\delta_1, \dots, \delta_m) \in \Delta^m : d_{\mathbb{P}}(T, H_{\bar{m}_d^r}) > \epsilon \right\} \leq \sum_{i=0}^{r+d-1} \binom{m}{i} \epsilon^i (1-\epsilon)^{m-i},$$

which is a sharper result compared to  $\binom{r+d-1}{r} \sum_{i=0}^{r+d-1} \binom{m}{i} \epsilon^i (1-\epsilon)^{m-i}$ . However, it holds only for the class of fully-supported, convex optimization programs, and for the constructed removal scheme that returns the solution with the lowest probability of constraint violation, while to compute the latter requires knowledge of  $\mathbb{P}$  which might be unavailable. A byproduct of this fact is that  $r + d$  constitutes the cardinality of a compression set for this problem, i.e., there exist  $r + d$  samples such that following the same procedure using only these samples yields the same solution had all the samples been employed. The result is tight, i.e., we have equality in (4.2), if the sets of multisamples for which the probability of constraint violation for each of the  $\binom{r+d-1}{r}$  solutions that violate the  $r$  removed samples exceeds  $\epsilon$ , are equally likely. This occurs when we have that  $\mathbb{P}^m \left\{ (\delta_1, \dots, \delta_m) \in \Delta^m : H_{\bar{m}_d^r} \text{ is not consistent with } \{(\delta_i, \mathbb{1}_T(\delta_i))\}_{i \in I_r} \text{ and } d_{\mathbb{P}}(T, H_{\bar{m}_d^r}) > \epsilon \right\}$  is the same for each of the  $\binom{r+d-1}{r}$  sets  $I_r$  giving rise to  $H_{\bar{m}_d^r}$ . In that case the inequality in (4.1) becomes an equality. An example where this is the case, is the scenario program corresponding to the minimum width interval when the samples are generated in an i.i.d. fashion from a uniform distribution on  $[0, 1]$ .

This derivation is inspired by a dual in some sense analysis in Section 5.2 of [20], where it is established that the right-hand side of (4.2) constitutes a lower bound for the case where the hypothesis is generated by an algorithm that returns the solution with the highest probability of constraint violation instead.

**5. Scenario approach and support constraints.** Consider the setting of Section III-B. Helly’s dimension is defined (see Definition 3.1 in [9]) as the least integer  $\zeta < \infty$  such that for any finite  $m \geq 1$ , with  $\mathbb{P}^m$ -probability one with respect to the choice of an  $m$ -multisample, the number of support samples is *at most*  $\zeta$ , i.e.,

$$(5.1) \quad \max_{m \geq 1} \text{ess sup}_{\{\delta_i\}_{i=1}^m \in \Delta^m} \left\{ \# \text{ of support constraints of } \mathcal{P}[\{\delta_i\}_{i=1}^m] \right\} \leq \zeta.$$

It follows then that we ought to have  $m \geq \zeta$ . Notice that  $\zeta$  is independent of the multi-sample; this implies that with non-zero probability there would exist a multisample of some length where the support constraints would be equal to  $\zeta$  (otherwise a tighter bound would exist), but this is not necessarily the case for all multisamples.

PROPOSITION 5.1. Consider Assumption 5 and fix  $d = \zeta$ . For any  $m \geq d$ , with  $\mathbb{P}^m$ -probability one with respect to the choice of an  $m$ -multisample  $\{\delta_1, \dots, \delta_m\}$ , there exists  $I_d \in \mathcal{I}_d$  such that  $x_d(\{\delta_i\}_{i \in I_d}) = x_m(\{\delta_i\}_{i=1}^m)$ , where  $x_m(\{\delta_i\}_{i=1}^m)$  is the minimizer of  $\mathcal{P}[\{\delta_i\}_{i=1}^m]$ , and  $x_d(\{\delta_i\}_{i \in I_d})$  denotes the optimizer of the same problem when fed only with the samples  $\{\delta_i\}_{i \in I_d}$ .

*Proof.* We show this by means of induction.

*Base case  $m = d$ :* The statement holds trivially for any  $d$ -multisample, as we have exactly  $d$  samples.

*Induction hypothesis for some  $m > d$ :* Consider an arbitrary  $m > d$ . Suppose that with  $\mathbb{P}^m$ -probability one with respect to an  $m$ -multisample  $\{\delta_1, \dots, \delta_m\}$ , there exists  $I_d \in \mathcal{I}_d$  such that  $x_d(\{\delta_i\}_{i \in I_d}) = x_m(\{\delta_i\}_{i=1}^m)$ .

*The  $(m+1)$ -th case:* Consider an  $(m+1)$ -multisample  $\{\delta_1, \dots, \delta_m, \delta_{m+1}\}$ . Since the number of support constraints is bounded by  $d = \zeta$  (independently of the multisample length) and by the definition of  $\zeta$ , with  $\mathbb{P}^{m+1}$ -probability one, at least  $m+1-d > 1$  of these samples will *not* be of support. Pick one of these, and without loss of generality assume this is the  $(m+1)$ -th sample, namely,  $\delta_{m+1}$ . By the definition of support constraints, removing the constraint associated to this sample will not change the (unique under Assumption 5) solution. We thus have that for the remaining samples  $\{\delta_1, \dots, \delta_m\}$  there exists  $I_d \in \mathcal{I}_d$  such that

$$(5.2) \quad x_{m+1}(\{\delta_i\}_{i=1}^{m+1}) = x_m(\{\delta_i\}_{i=1}^m) = x_d(\{\delta_i\}_{i \in I_d}),$$

where the first equality is due to the fact that  $\delta_{m+1}$  is not of support, and the second one follows from the induction hypothesis applied to the remaining  $m$ -multisample  $\{\delta_1, \dots, \delta_m\}$  (notice that the hypothesis refers to any multisample of length  $m$ ). Overall, there exists a subset of the  $(m+1)$ -multisample with length  $d$  that results in  $x_{m+1}(\{\delta_i\}_{i=1}^{m+1})$  up to a  $\mathbb{P}^{m+1}$ -measure zero set (follows from the fact that the induction hypothesis holds up to a  $\mathbb{P}^m$ -measure zero set), thus concluding the induction proof.  $\square$

A direct consequence of the construction in the proof, is that in the set of  $\zeta$  samples that are sufficient to return the same solution with the one that would have been obtained if all the samples were used, the support samples are always included. To see this, notice that if a support sample was not included in that set, i.e., it was removed from the  $m$ -multisample, then the solution would have to change by definition of support samples/constraints. A set of  $\zeta$  samples that satisfies the statement of Proposition 1 is not necessarily unique; in the particular case where the number of support constraints is *equal to*  $\zeta < \infty$ , i.e., (5.1) holds with equality, then such a set is unique and is equal to the set of support samples.

If the number of support constraints is *at most*  $\zeta = n_x$ , then Proposition 1 holds for any  $m$ -multisample as opposed to almost surely since the number of support constraints for convex problems cannot exceed  $n_x$  (e.g., see Theorem 3 in [7]), while if it is exactly *equal to*  $n_x$  we need to exclude measure zero cases that would prevent this from happening always, like selecting the same sample  $m$  times.