
The Single-Point Crossover Process

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The last decade has seen a burst of data optimisation techniques, mainly via
2 noisification mechanisms, in two separate contexts: (i) improve supervised learning
3 (as in dropout) and (ii) enhance individual privacy (as in differential privacy). In
4 this paper, we tackle both objectives for a *process* level of protection on data: we
5 want to keep the utility of data for supervised learning while being able to control —
6 or fool — *causal inference attacks* on subsets of its features. It is well known that
7 privacy protection and learnability are usually two conflicting objectives. Usually,
8 but not always: we show how to jointly control causal inference attacks *and*
9 learnability by a noise-free process that mixes observations in the training sample,
10 the Single-Point Crossover Process (SPCP). One key point is that the SPCP is
11 typically able to alter joint distributions *without* touching on marginals, nor altering
12 the sufficient statistics for the class — in other words, it saves (and sometimes
13 improves) generalization abilities for learning, but can fool causal inference attacks
14 into misleading *ad-hoc* conclusions. Extensive experiments validate the theory and
15 display the competitive advantages of SPCPs.

16 1 Introduction

17 — “Even when individuals are not ‘identifiable’, they may still be ‘reachable’, may still be com-
18 prehensibly represented in records that detail their attributes and activities, and may be subject to
19 *consequential inferences* and predictions taken on that basis” [5].

20 There are at least too good reasons to alter the training sample of a supervised learning problem.
21 One is privacy, a growing concern in the public sphere [5, 15, 17]. Two leading mechanisms for the
22 private release of data are differential privacy and k -anonymity [14, 15, 29, 38]. They guarantee the
23 *individual (identifiability) level* protection mentioned above [5] and rely on low-level modifications,
24 typically touching marginals, mostly with noise. The other good reason to alter a training sample
25 is to optimize learning: feature bagging [9, 25, 37], importance weighting in features, boosting
26 [23, 30], denoising autoencoders [40], independent noisification, dropout [3, 36, 39], all are methods
27 that optimise the presentation of observation variables to a learner with the objective to improve its
28 generalization abilities.

29 In 2014, a report to the US president by the Presidents Council of Advisors on Science and Technology
30 on big data and privacy judged, in its first recommendation, that “policies focused on the regulation of
31 data collection [...] and analysis (*absent identifiable actual uses of the data or products of analysis*) are
32 unlikely to yield effective strategies” [17]. True causality is deemed very sensitive but its imputation
33 in Big Data is judged “a research field in its infancy” [17]. It is however hard to exaggerate the recent
34 burst in causal inference techniques [8, 12, 13, 18, 19, 20, 21, 26, 27, 33], as well as the threats this
35 may pose on privacy issues [4, 5, 15, 24]. In this framework, rather than consequential inference [5],
36 we shall name them causal inference *attacks*, and we argue that they deserve an adequate *process*
37 *level* protection for data [4, 5], more than an all-purpose individual level protection. For example,
38 in a medical diagnosis data which gives a sickness state as a function of genetic, behavioral, habits

and infection history, we may want to make causal inference between specific traits and behaviour harder, or we may want to hide gender-prone infections [34], *while* making sure that the utility of the dataset for predicting the sickness state remains unaltered. Rather than making it harder, ultimately, we may just look to *reverse* the observable causal direction between specific observation variables, and thereby fool causal inference or causal rule mining. In short, the dataset’s utility for the black-box supervised prediction task remains within control, but it is surgically altered against fined-grained specific causal inference attacks among features.

Is such a task within technical reach ? It is well known that individual protection mechanisms do not get on well with optimizing learning [10, 29, 41]. Here comes an eventually surprising observation: coping with our protection level may require wrangling the complete data, but this may be done with a tight control of its utility for supervised learning, and it may even yield *better* models for prediction.

— This is our main contribution: to achieve these two goals, we introduce the *Single Point Crossover Process*, SPCP. An analogy may be done with the biological crossover: a population of DNA strands gets mixed with a crossover, but there is a single zone for chiasma (*i.e.* contact point) for the whole population. In the same way as DNA strands exchange genetic material during recombination, feature values get mixed between observations during a SPCP (in a more general way than the exchange in genetic recombination). The key to learning is that the SPCP may be done without changing the sufficient statistics for the class, nor touching marginals.

Organisation of the paper — Section §2 gives general definitions. Section §3 presents the Single-Point Crossover Process and its relationships with learnability. Section §5 presents the applications of the SPCP and details related experiments. A last Section discusses and concludes. A Supplementary Material (SM) [1] provides all proofs, additional results and the complete experiments performed. A movie is also included in SM, showing the effects of the SPCP on a popular domain for causal discovery [20].

2 General notations and definitions

Learning setting — We let $[m] \doteq \{1, 2, \dots, m\}$ and $\Sigma_m \doteq \{\sigma \in \{-1, 1\}^m\}$. $\mathcal{X} \subseteq \mathbb{R}^d$ is a domain of observations. Examples are couples (observation, label) $\in \mathcal{X} \times \Sigma_1$, sampled i.i.d. according to some unknown but fixed distribution \mathcal{D} . We denote $\mathcal{F} \doteq [d]$ the set of observation attributes (or features). $\mathcal{S} \doteq \{(\mathbf{x}_i, y_i), i \in [m]\} \sim \mathcal{D}_m$ is a training sample of $|\mathcal{S}| = m$ examples. For any vector $\mathbf{z} \in \mathbb{R}^d$, z_j denotes its coordinate j . Finally, notation $x \sim X$ for X a set denotes uniform sampling in X , and the mean operator is $\mu_{\mathcal{S}} \doteq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}}[y \cdot \mathbf{x}]$ [32].

In supervised learning, the task is to learn classifier $\mathcal{H} \ni h : \mathcal{X} \rightarrow \mathbb{R}$ from \mathcal{S} with good generalisation properties, that is, having a small *true risk* $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[L_{0/1}(y, h(\mathbf{x}))]$, with $L_{0/1}(z, z') \doteq 1_{zz' \leq 0}$ the 0/1 loss (1. is the indicator variable). In general, this is achieved by minimising over \mathcal{S} a φ -risk $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}}[\varphi(yh(\mathbf{x}))] = (1/m) \cdot \sum_i \varphi(y_i h(\mathbf{x}_i))$, where $\varphi(z) \geq 1_{z \leq 0}$ is a *surrogate* of the 0/1 loss. In this paper, φ is any differentiable proper symmetric (PS) loss [28] (symmetric meaning that there is no class-dependent misclassification cost). The logistic, square and Matsushita losses are examples of PS losses. Set \mathcal{H} is a predefined set of classifiers, such as linear separators, decision trees, etc. .

Matrix quantities — The set of unnormalised column stochastic (UCS) matrices, $\mathcal{M}_n \subset \mathbb{R}^{n \times n}$, is the superset of column stochastic matrices for which we drop the non-negativity constraint, thus keeping the sole constraint of unit per-column sums. We let $S_n \subset \mathcal{M}_n$ denote the symmetric group of order n . For any $A, B \in \mathbb{R}^{n \times n}$ and $M \in \mathcal{M}_n$, we let $\langle A, B \rangle_M \doteq \text{tr}((I_n - M)^\top A(I_n - M)B)$ denote the *centered inner product* of A and B with respect to M . It is a generalisation of the centered inner product used in kernel statistical tests of independence [18], for which $M = (1/n)\mathbf{1}\mathbf{1}^\top$.

Without loss of generality, we shall assume that indexes in \mathcal{S} cover first the positive class: $(y_i = +1 \wedge y_{i'} = -1) \Rightarrow i < i'$. A key subset of matrices of $\mathbb{R}^{m \times m}$ consists of block matrices whose coordinates on indexes corresponding to different classes in \mathcal{S} are zero: block-class matrices.

Definition 1 $A \in \mathbb{R}^{m \times m}$ is a **block-class matrix** iff $(y_i \cdot y_{i'} = -1) \Rightarrow A_{ii'} = 0, \forall i, i'$.

An asterisk exponent in a subset of matrices indicates the intersection of the set with block class matrices, such as for $\mathcal{M}_n^* \subset \mathcal{M}_n$ and $S_n^* \subset S_n$. Finally, matrix entries are noted with double indices like $M_{ii'}$, and replacing an index by a dot, “.”, indicates a sum over the index, like $M_{i.} \doteq \sum_{i'} M_{ii'}$.

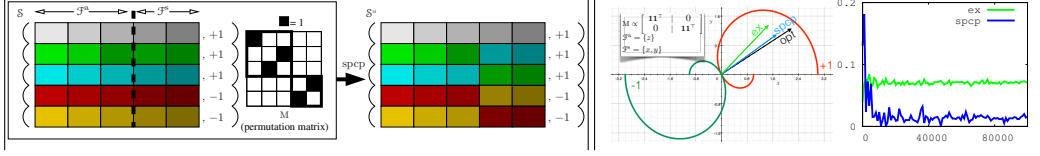


Figure 1: *Left*: example of SPCP with M a block-class permutation matrix (the two blocks are in bold). *Center+Right*: toy domain where $d = 3$, but all examples have zero z -coordinate (not shown). The SPCP uniformly mixes examples by class. The domain consists of two spirals (red for positive, green for negative examples) with \mathcal{D} = uniform distribution. Arrows depict respectively the optimal direction (black), and the directions learned by minimizing φ = square loss over \mathcal{S} (light green) and $\mathcal{S}^{\mathcal{S}}$ (cyan). The rightmost plot displays test errors (y -scale) on uniform sampling of datasets of different sizes (x -scale). The effect of the SPCP is to produce in $\mathcal{S}^{\mathcal{S}}$ two distinct examples that average the positive / negative examples, and yield a better approximation of the optimum.

90 3 The Single-Point Crossover Process

91 The single-point crossover process (SPCP) transforms \mathcal{S} in two steps: the split and the shuffle step. In
 92 the split step, a bi-partition of the features set \mathcal{F} is computed: $\mathcal{F} = \mathcal{F}_a \cup \mathcal{F}_s$. \mathcal{F}_a is the *anchor* set and
 93 \mathcal{F}_s is the *shuffle* set. To perform the shuffle step, we need some additional notations. Without loss
 94 of generality, we assume $\mathcal{F}_a \doteq [d_a]$ and $\mathcal{F}_s \doteq \{d_a + j, j \in [d_s]\}$, $d_a > 0, d_s > 0, d_a + d_s = d$. So,
 95 \mathcal{F}_a contains the first d_a features and \mathcal{F}_s contains the last d_s features. Let I_d be the identity matrix,
 96 and $[F^a | F^s] = I_d$ be a vertical block partition where $F^\pi \in \mathbb{R}^{d \times d_\pi}$ has columns representing the
 97 features of \mathcal{F}_π (for $\pi \in \{a, s\}$). Finally, we define the (row-wise) observation matrix $S \in \mathbb{R}^{m \times d}$ with
 98 $(S)_{ij} \doteq x_{ij}$. Let $\mathbf{1}_i$ be the i^{th} canonical basis vector.

99 **Definition 2** For any block partition $[F^a | F^s] = I_d$ and any *shuffle matrix* $M \in \mathcal{M}_n$, the single-point
 100 crossover process $\mathcal{T} \doteq \text{SPCP}(\mathcal{S}; F^a, F^s, M)$ returns m -sample $\mathcal{S}^{\mathcal{T}}$ such that its observation matrix is
 101 $S^M \doteq [SF^a | MSF^s]$, and each example $\mathcal{S}^{\mathcal{T}} \ni (x_i^M, y_i) \doteq ((S^M)^\top \mathbf{1}_i, y_i)$.

102 M can be fixed beforehand or learned with data. Figure 1 (left) presents the SPCP on a toy data with
 103 M a permutation matrix. Figure 1 (center + right) presents another example with M block-uniform.

104 **The Single-Point Crossover Process and learnability** — We now explore the effect of the SPCP on
 105 generalisation. We need two assumptions on \mathcal{H} and φ . The first is a weak linearity condition on \mathcal{H} :

106 (i) $\forall h \in \mathcal{H}, \exists$ classifiers h_π over features of \mathcal{F}_π ($\pi \in \{a, s\}$) s. t. $h(\mathbf{x}) = \sum_\pi h_\pi((F^\pi)^\top \mathbf{x})$.

107 Such an assumption is also made in the feature bagging model [37]. We let \mathcal{H}_s denote the set of all
 108 h_s . Any linear classifier satisfies (i), but also any linear combination of arbitrary classifiers, each
 109 learnt over one of \mathcal{F}_π for $\pi \in \{a, s\}$. The second postulates that key quantities are bounded [6]:

110 (ii) $0 \leq \varphi(z) \leq K_\varphi, \forall z$ and $|h_s((F^s)^\top \mathbf{x})| \leq K_s, \forall \mathbf{x} \in \mathcal{X}, \forall h_s \in \mathcal{H}_s$.

111 $(F^s)^\top \mathbf{x}$ picks the features of \mathbf{x} in \mathcal{F}_s . Let $R_S(\mathcal{H}) \doteq \mathbb{E}_{\sigma \sim \Sigma_m} [\sup_{h \in \mathcal{H}} |(1/m) \cdot \sum_i \sigma_i h(\mathbf{x}_i)|]$ be the
 112 empirical Rademacher complexity of \mathcal{H} .

113 **Definition 3** The Rademacher discrepancy (D) of \mathcal{H} with respect to $\mathcal{T} \doteq \text{SPCP}(\mathcal{S}; F^a, F^s, M)$ is:

$$D_{\mathcal{T}}(\mathcal{H}) \doteq \mathbb{E}_{\sigma \sim \Sigma_m} \left[\sup_{h \in \mathcal{H}_s} \left| \frac{1}{m} \sum_i \sigma_i (h((SF^s)^\top \mathbf{1}_i) - h((MSF^s)^\top \mathbf{1}_i)) \right| \right]. \quad (1)$$

114 Notice that the Rademacher discrepancy is computed over the shuffle set of features only, and
 115 $(SF^s)^\top \mathbf{1}_i = (F^s)^\top \mathbf{x}_i$. We investigate generalization under two different settings on $\mathcal{H} / \mathcal{T}$:

116 Setting (A) Classifiers h^s in (i) above are linear;

117 Setting (B) $M \in S_m^*$.

Theorem 4 Consider any \mathcal{H} , φ and split $\mathcal{F} = \mathcal{F}_a \cup \mathcal{F}_s$ such that (i) and (ii) hold. For any m and any $\delta > 0$, with probability $\geq 1 - \delta$ over i.i.d. m -sample \mathcal{S} , we have:

$$\mathbb{E}_{\mathcal{D}} [L_{0/1}(y, h(\mathbf{x}))] \leq \mathbb{E}_{\mathcal{S}, \mathcal{T}} [\varphi(yh(\mathbf{x}))] + D_{\mathcal{T}}(\mathcal{H}) + \frac{4}{b_{\varphi}} \cdot R_{\mathcal{S}}(\mathcal{H}) + (2K_{\varphi} + K_{\mathcal{S}}) \cdot \sqrt{\frac{2}{m} \log \frac{3}{\delta}},$$

for every classifier h and every $\mathcal{T} \doteq \text{SPCP}(\mathcal{S}; \mathcal{F}_a, \mathcal{F}_s, \mathbf{M})$ such that $\mathbf{M} \in \mathcal{M}_m^*$ and whichever of (A) or (B) holds. Here, $b_{\varphi} > 0$ is a constant depending on φ .

([1], Subsection 2.1) Notice that Theorem 4 requires that \mathbf{M} is a block-class matrix. A key to the proof is the invariance of the mean operator: $\boldsymbol{\mu}_{\mathcal{S}} = \boldsymbol{\mu}_{\mathcal{S}, \mathcal{T}}$. The roles of \mathcal{F}_a and \mathcal{F}_s may be interchanged for Setting (B). This allows to improve at no cost the result, which thus holds over the *min* over the permutation of the roles.

Theorem 4 gives a perhaps counterintuitive rationale for the SPCP that goes beyond our framework to machine learning at large: *learning over a SPCP'ed \mathcal{S} may improve generalization over \mathcal{D} as well*. By means of words, learning over transformed data may improve generalization over the initial domain. Figure 1 (right) gives a toy example for which this holds. It is also not hard to exhibit domains for which we even have:

$$\min_h \mathbb{E}_{\mathcal{S}, \mathcal{T}} [\varphi(yh(\mathbf{x}))] + D_{\mathcal{T}}(\mathcal{H}) < \min_h \mathbb{E}_{\mathcal{S}} [\varphi(yh(\mathbf{x}))]. \quad (2)$$

To save space, we present such an example in the Supplementary Material ([1], Subsection 2.2). Theorem 4 says that a key to good generalization is the “complexity” of the process encapsulated in $D_{\mathcal{T}}(\mathcal{H})$. We would typically want it to be small compared to the Rademacher complexity. The rest of this Section shows that (and when) this is indeed achievable.

Upperbounds on $D_{\mathcal{T}}(\mathcal{H})$ — The following Lemma establishes a first bound on $D_{\mathcal{T}}(\mathcal{H})$.

Lemma 5 if \mathcal{T} satisfies the conditions of Theorem 4, then $D_{\mathcal{T}}(\mathcal{H}) \leq 2 \cdot R_{\mathcal{S}'}(\mathcal{H}_{\mathcal{S}})$, for $\mathcal{S}' \doteq (\mathbf{I}_m - \mathbf{M})\mathbf{S}\mathcal{F}^{\mathcal{S}}$ in Setting (A), and $\mathcal{S}' \doteq \mathbf{S}\mathcal{F}^{\mathcal{S}}$ in Setting (B). \mathcal{S}' is the row-wise observation matrix of \mathcal{S}' .

Proof (Sketch) Consider for example Setting (B). In this case, recalling that $(\mathbf{S}\mathcal{F}^{\mathcal{S}})^{\top} \mathbf{1}_i = (\mathcal{F}^{\mathcal{S}})^{\top} \mathbf{x}_i$ and letting $\varsigma : [m] \rightarrow [m]$ denote the permutation that \mathbf{M} represents, we have:

$$\begin{aligned} D_{\mathcal{T}}(\mathcal{H}) &= \mathbb{E}_{\boldsymbol{\sigma} \sim \Sigma_m} \left[\sup_{h \in \mathcal{H}_{\mathcal{S}}} \left| \frac{1}{m} \sum_i \sigma_i (h((\mathcal{F}^{\mathcal{S}})^{\top} \mathbf{x}_i) - h((\mathcal{F}^{\mathcal{S}})^{\top} \mathbf{x}_{\varsigma(i)})) \right| \right] \\ &\leq \mathbb{E}_{\boldsymbol{\sigma} \sim \Sigma_m} \left[\sup_{h \in \mathcal{H}_{\mathcal{S}}} \left| \frac{1}{m} \sum_i \sigma_i h((\mathcal{F}^{\mathcal{S}})^{\top} \mathbf{x}_i) \right| \right] + \mathbb{E}_{\boldsymbol{\sigma} \sim \Sigma_m} \left[\sup_{h \in \mathcal{H}_{\mathcal{S}}} \left| \frac{1}{m} \sum_i \sigma_i h((\mathcal{F}^{\mathcal{S}})^{\top} \mathbf{x}_{\varsigma(i)}) \right| \right] \\ &= 2 \cdot R_{\mathcal{S}'}(\mathcal{H}_{\mathcal{S}}), \end{aligned} \quad (3)$$

as claimed. The case of Setting (A) follows the same path. ■

Lemma 5 says that $D_{\mathcal{T}}(\mathcal{H})$ is at most twice a Rademacher complexity over the *shuffle set*. This bound is however loose since many terms can cancel in the sum of eq. (3), and the inequality does not take this into account. We now exploit this to show bounds that can be much smaller.

Theorem 6 Under Setting (A), suppose any $h_{\mathcal{S}}$ is of the form $h_{\mathcal{S}}(\mathbf{x}) = \boldsymbol{\theta}^{\top} \mathbf{x}$ with $\|\boldsymbol{\theta}\|_2 \leq r_{\mathcal{S}}$, for some $r_{\mathcal{S}} > 0$. Let $\mathbf{K}^{\mathcal{S}} \doteq \mathbf{S}\mathcal{F}^{\mathcal{S}}(\mathbf{S}\mathcal{F}^{\mathcal{S}})^{\top}$. Then $\exists u \in (0, 1)$ depending only in \mathcal{S} such that for any $\mathbf{M} \in \mathcal{M}_m$,

$$D_{\mathcal{T}}(\mathcal{H}) \leq (ur_{\mathcal{S}}/m) \cdot \sqrt{\langle \mathbf{I}_m, \mathbf{K}^{\mathcal{S}} \rangle_{\mathbf{M}}}. \quad (4)$$

Notice that $\mathbf{K}^{\mathcal{S}}$ is a Gram matrix in the shuffle feature space. The proof technique [1] (Subsection 2.3) relies on a data-dependent expression for u which depends on the cosines of angles between the observations in \mathcal{S} . It can be used to refine and improve a popular bound on the empirical Rademacher complexity of linear classifiers [22] (we give the proof in [1] (Theorem 4)). We now investigate an upperbound on Setting (B) in which classifiers in $\mathcal{H}^{\mathcal{S}}$ are (rooted) directed acyclic graph (DAG), like decision trees, with bounded real valued predictions (say, $K_{\mathcal{S}} > 0$) at the leaves. Each classifier $h_{\mathcal{S}}$ defines a partition over \mathcal{X} . We let $\mathcal{H}_{+}^{\mathcal{S}}$ be the subset of $\mathcal{H}^{\mathcal{S}}$ in which all leaves have in absolute value the largest magnitude, i.e., $K_{\mathcal{S}}$. Remark that we may have $|\mathcal{H}_{+}^{\mathcal{S}}| \ll \infty$ while $|\mathcal{H}^{\mathcal{S}}| = \infty$ in general.

Theorem 7 Under Setting (B), suppose \mathcal{H}^s is DAG and assumption (ii) is satisfied. Suppose that $\log |\mathcal{H}_+^s| \geq (4\varepsilon/3) \cdot m$ for some $\varepsilon > 0$. Then, letting $\text{odd_cycle}(\mathbf{M})$ denote the set of odd cycles (excluding fixed points) of \mathbf{M} , we have:

$$D_{\mathcal{T}}(\mathcal{H}) \leq K_s \cdot \sqrt{\frac{2}{m} \cdot \log \frac{|\mathcal{H}_+^s|}{(1+\varepsilon)^{|\text{odd_cycle}(\mathbf{M})|}}} \quad (5)$$

([1], Subsection 2.4) The assumption on \mathcal{H}_+^s is not restrictive and would be met by decision trees, branching programs, etc. (and subsets). The Rademacher complexity of decision trees would roughly be the right-hand side of (5) *without* the denominator in the log. Hence, the Rademacher discrepancy may be much smaller than the Rademacher complexity for more “involved” SPCPs. The number of cycles is not the only relevant parameter of the SPCP on which relies non-trivial bounds on $D_{\mathcal{T}}(\mathcal{H})$: the Supplementary Material presents, for the interested reader, a proof that the number of fixed points is another parameter which can decrease significantly the *expected* Rademacher discrepancy (by a factor $\sqrt{1 - |\text{fixed_points}|/m}$), when SPCPs are picked at random.

4 The SPCP, learnability, inference and causality

The Single-Point Crossover Process and measures of independence — In this section, we assume that \mathcal{S} is subject to quantitative tests of independence, that is, assessing $\mathcal{U} \perp\!\!\!\perp \mathcal{V}$ for some $\mathcal{U}, \mathcal{V} \subset \mathcal{X}$. We therefore compute SPCPs such that $\mathcal{U} \subseteq \mathcal{F}_a$ and $\mathcal{V} \subseteq \mathcal{F}_s$, so that the SPCP alters the measure of independence. The problem is thus essentially the design of the shuffle matrix \mathbf{M} . One popular criterion to determine (conditional) (in)dependence is Hilbert-Schmidt Independence Criterion [13, 18, 19].

Definition 8 Let $\mathcal{U} \subset [d]$ and $\mathcal{V} \subset [d]$ be non-empty and disjoint. Let \mathbf{K}^u and \mathbf{K}^v be two kernel functions over \mathcal{U} and \mathcal{V} computed using \mathcal{S} . The (unnormalised) Hilbert-Schmidt Independence Criterion (HSIC) between \mathcal{U} and \mathcal{V} is defined as $\text{HSIC}(\mathbf{K}^u, \mathbf{K}^v) \doteq \langle \mathbf{K}^u, \mathbf{K}^v \rangle_{(1/m)\mathbf{1}\mathbf{1}^\top}$.

We choose not to normalise the HSIC: various exist but they mainly rely on a multiplicative factor depending on m only, so they do not affect the results to come. Our first result shows more than just how to pick \mathbf{M} in the SPCP so as to alter HSIC: (a) while storing kernels requires $O(m^2)$ space, controlling the evolution of the HSIC requires only *linear*-space information about kernels, and (b) this information can be computed beforehand, and can be efficiently approximated from low-rank approximations of the kernels [2]. Hereafter, we restrict ourselves to the scenario of the *decrease* of the HSIC, yet all our results would apply to the opposite polarity for the modification. We use for sake of readability the shorthand HSIC for $\text{HSIC}(\mathbf{K}^u, \mathbf{K}^v)$.

Theorem 9 Let HSIC and $\text{HSIC}_{\mathcal{T}}$ denote the HSIC before and after applying the SPCP \mathcal{T} to \mathcal{S} (\mathbf{K}^u and \mathbf{K}^v implicit). Let $\tilde{\mathbf{u}} \doteq (1/m) \sum_i \lambda_i (\mathbf{1}^\top \mathbf{u}_i) \mathbf{u}_i$, $\tilde{\mathbf{v}} \doteq (1/m) \sum_i \mu_i (\mathbf{1}^\top \mathbf{v}_i) \mathbf{v}_i$, where $\{\lambda_i, \mathbf{u}_i\}_{i \in [d]}$, $\{\mu_i, \mathbf{v}_i\}_{i \in [d]}$ are respective eigensystems of \mathbf{K}^u and \mathbf{K}^v . Then $\text{HSIC}_{\mathcal{T}} < \text{HSIC}$ iff $\tilde{\mathbf{u}}^\top (\mathbf{I}_m - \mathbf{M}) \tilde{\mathbf{v}} < 0$.

([1], Subsection 2.5) In the following Theorem, we compose SPCP processes with T different elementary permutation shuffling matrices. Notice that since the composition of permutation matrices is a permutation matrix, when the matrix of the final process $\text{HSIC}_{\mathcal{T}_T}$ is block class, Theorem 4 can be applied *directly* to $\text{HSIC}_{\mathcal{T}_T}$. We also let $\mathcal{R}^{u,v} \doteq m (1 - (\mathbf{K}_\cdot^u + \mathbf{K}_\cdot^v)/(2m^2))$.

Theorem 10 Suppose \mathcal{S} is shuffled by a sequence of $T = \epsilon m$ elementary permutation ($\epsilon > 0$) and the kernels \mathbf{K}^u and \mathbf{K}^v have unit diagonal. Suppose that the initial $\text{HSIC} > \mathcal{R}^{u,v}$. Then there exists such a sequence such that \mathcal{T}_T satisfies $\text{HSIC}_{\mathcal{T}_T} \leq \mathcal{R}^{u,v} + \alpha \cdot (\text{HSIC} - \mathcal{R}^{u,v})$, where $\alpha \doteq \exp(-8\epsilon)$.

The proof ([1], Subsection 2.6) states a more general result, not restricted to unit diagonal kernels. Theorem 10 is a worst-case result: some sequences of permutations may be much more efficient in decreasing HSIC. If we compare this bound to Theorem 3 in [19], then $\mathcal{R}^{u,v}$ may be *below* the expectation of the HSIC, and so Theorem 10 guarantees efficient jamming of dependence.

We finish with independence on a rather specific causality model, but that we feel is enlightening on one important potential of SPCP with respect to independence: trick statistical tests into keeping

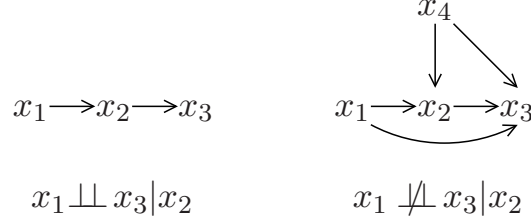


Figure 2: The Cornia-Mooij model [12]. Left: belief, right: true model.

Algorithm 1 SPCP-Train $\text{DT}(\mathcal{S}, T, \theta_0, M_0, F^a, F^s; \mathcal{G}_1, [\mathcal{G}_2])$

Input Sample \mathcal{S} , iterations T , classifier θ_0 , initial SPCP \mathcal{T} (matrices $M_0 \in S_m^*$, $[F^a | F^s] = I_m$);

Step 1 : **for** $t = 1, 2, \dots, T$

 Step 1.1 : let $M \leftarrow \arg \min_{M' \in S_m^{e*}} \mathcal{G}_1(M' \circ M_{t-1} | [\theta_{t-1}])$;

 Step 1.2 : let $M_t \leftarrow M \circ M_{t-1}$;

 [Step 1.3 : let $\theta_t \leftarrow \arg \min_{\theta \in \mathbb{R}^d} \mathcal{G}_2(\theta | M_t)$];

Return classifier θ_T and / or dataset $\mathcal{S}^{\mathcal{T}}(M_T)$

independence *and* then incur arbitrarily large errors in estimating causal effects. Also, its causal graph is so simple that it may be found as subgraph of real-world domains, thus for which the results we give would directly transfer. The model we refer to is the Cornia-Mooij (CM) model [1, 12]. In the CM model, there are $d = 3$ observation variables, and a true model which relies on a weak conditional dependence $x_1 \not\perp\!\!\!\perp x_3 | x_2$. [12] show that *if* one keeps the independence assumption H_0 that $x_1 \perp\!\!\!\perp x_3 | x_2$, this can lead to very high causal estimation errors, as measured by $|\mathbb{E}[x_3 | x_2] - \mathbb{E}[x_3 | \text{do}(x_2)]| / |x_2|$ [12]. We show that the SPCP is precisely able to trick statistics into keeping H_0 . The Supplementary Material (Subsection 2.7) presents in detail the model as well as how the SPCP may be chosen.

The Single-Point Crossover Process and causal calculus — Rather than the shuffle matrix M , the problem in this section is rather the design of \mathcal{F}_a and \mathcal{F}_s .

5 Experiments

Our applications use the same meta-level algorithm (Algorithm 1), which operates in Setting (A) \cap Setting (B) (M in S_m^* , linear classifiers), iteratively composing block-class elementary permutations. Here, $S_m^{e*} \subset S_m$ is the set of block-class elementary permutations. The iteration step minimises a criterion \mathcal{G}_1 over S_m^{e*} , and eventually after having updated the SPCP matrix, a criterion \mathcal{G}_2 over \mathcal{H} . The experimental setup and the results are provided *in extenso* in [1] (Section 3.4); Table 1 summarises them. Due to the lack of space, we comment experiments alongside the applications. The split step and the choice of \mathcal{F}_a are highly domain and task dependent: to keep experiments of reasonable length, we chose to put in \mathcal{F}_a the first half of features for each domain. In Abalone2D and Digoxin, this amounts to jam a particular ground truth (see below). Also, φ =logistic loss.

Data optimisation for efficient learning — We perform Algorithm 1 with $M_0 = I_m$, $\mathcal{G}_1(M | \theta) = \mathbb{E}_{\mathcal{S}^{\mathcal{T}}(M)} [\varphi(y\theta^\top x)]$ (Theorem 4) and $\mathcal{G}_2(\theta | M) = \mathbb{E}_{\mathcal{S}^{\mathcal{T}}} [\varphi(y\theta^\top x)] + \lambda \|\theta\|_2^2$ where λ is learnt through cross-validation. The algorithm returns classifier θ_T . The top row in Table 1, and [1] display that DT almost always find some permutations that reduce the test error compared to the initial data, validating this part of the theory, even when a specific data optimisation should care for a risk of over-fitting, which seems to occur on Ionosphere and Abalone2D [1] (Subsection 3.5) also compares DT as in Algorithm 1 to the one where we relax the constraint that permutations must be block-class (implying the invariance of the mean operator). The results are a clear advocacy for the constraint, as relaxing it brings poor results, from both the φ -risk and test error standpoints.

Disrupting dependence and causality without touching marginals — In our experiments, we run DT without step 1.3, and fixing \mathcal{G}_1 to be HSIC. We use two Gaussian kernels for L^u and L^v , each computed over its full subset of features [1]. Experimental results display that DT achieves in general

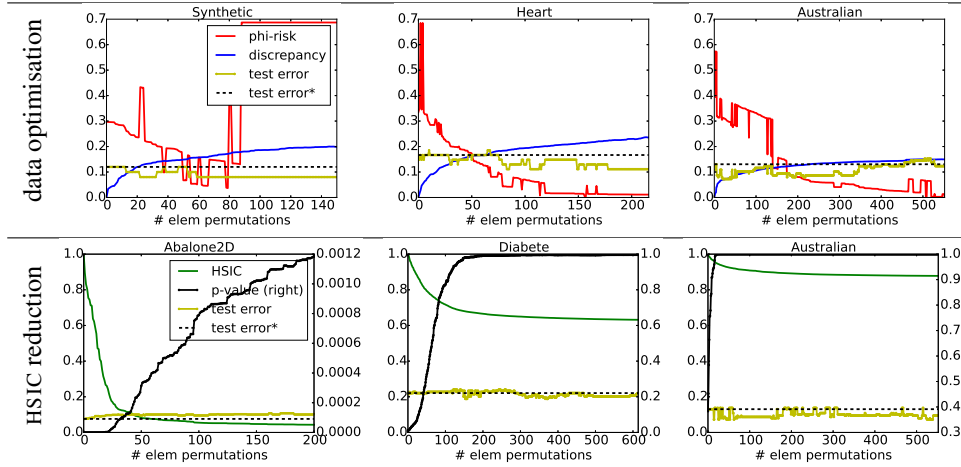


Table 1: Experiments performed with SPCP. From top row to bottom row: data optimisation, reduction in HSIC. References to domain names are provided in [1]. “test-error*” is test error over initial, non shuffled data.

a better control of the Rademacher discrepancy (than *e.g.* in data optimisation), and in several cases, manages to decrease the true error as well through the process. It also manages in general to relatively blow-up the p -value (computed as in [19]), even for domains for which the ground truth *implies the alternative hypothesis* H_1 . In Abalone2D for example, a gold standard for dependence [20], while the initial p -value is zero up to *thirteen* digits, after shuffling, the final p -value exceeds 1%. For the Digoxin domain, another popular domain [13] with ground truth, p is very small at the beginning (which corresponds to the ground truth $D \perp\!\!\!\perp U$; D = digoxin clearance, U = urin flow). After shuffling, we obtain $p > 0.4$, which easily brings $D \perp\!\!\!\perp U$, while the ground truth is $D \perp\!\!\!\perp U|C$ (C = creatinine clearance). In both cases, the effect on test error is minimal, considering that Abalone2D and Digoxin have $d = 2$ attributes only.

6 Discussion and conclusion

This paper introduces the Single-Point Crossover Process (SPCP), a process that cross-modifies data using a generalisation of stochastic matrices. This process can be used to cope with data optimisation for supervised learning, as well as for the problem of handling a process-level protection on data: causal inference attacks on a supervised learning dataset. In this case, the SPCP allows to release data with spotless low-level description (variable names, observed values), substantial utility (learnability), but disclosing dependences and causal effects under control, and thus that could even be crafted to be conflicting with a ground truth to protect¹. Note that causality in big data may still be in its infancy [17], it is however rapidly growing with a variety of techniques and lots of promises. We have chosen to focus here on two major components of the actual trends (stastical measures of causal inference and causal calculus), but our technique has more applications in the field of causal discovery: suppose for example that description features denote transactions. Since we jam joint distributions without touching on marginals, our technique has direct applications in causal rule mining, with the potential to fool any spawn of Apriori — that is, level-wise association rule mining algorithms [26].

The theory we develop fo SPCP introduces a new complexity measure of the process, the Rademacher discrepancy. We do believe that the SPCP is a good candidate in the pool of methods optimising data for learning, and may provide new metrics, algorithms and tools to devise improved solutions that fit to challenging domains not just restricted to optimizing learning or data privacy. One example is learning without entity resolution [31].

¹Note that the initial data may not be lost, as opposed to differential privacy: knowing the noise parameters does not allow to revert differential privacy protection, while a SPCP protection is reversible when M is invertible.

References

- [1] Anonymous. Supplementary material to this paper, 2016.
- [2] F. Bach. Sharp analysis of low-rank kernel matrix approximations. In *26th COLT*, pages 185–209, 2013.
- [3] P. Bachman, O. Alsharif, and D. Precup. Learning with pseudo-ensembles. In *NIPS*27*, pages 3365–3373, 2014.
- [4] S. Barocas and H. Nissenbaum. Big data’s end run around procedural privacy protection. *Communications of the ACM*, 57:31–33, 2014.
- [5] S. Barocas and H. Nissenbaum. Big data’s end run around anonymity and consent. In J. Lane, V. Stodden, S. Bender, and H. Nissenbaum, editors, *Privacy, Big Data, and the Public Good*, pages 44–75. Cambridge University Press, 2014.
- [6] P.-L. Bartlett and S. Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *JMLR*, 3:463–482, 2002.
- [7] O. Benjelloun, H. Garcia-Molina, D. Menestrina, Q. Su, S.-E. Whang, and J. Widom. Swoosh: a generic approach to entity resolution. *VLDB J.*, 18(1):255–276, 2009.
- [8] L. Bottou, J. Peters, J. Quiñonero Candela, D.-X. Charles, M. Chickering, E. Portugaly, D. Ray, P.-Y. Simard, and E. Snelson. Counterfactual reasoning and learning systems: the example of computational advertising. *JMLR*, 14:3207–3260, 2013.
- [9] R.-K. Bryll, R. Gutierrez-Osuna, and F.-K.-H. Quek. Attribute bagging: improving accuracy of classifier ensembles by using random feature subsets. *Pattern Recognition*, 36:1291–1302, 2003.
- [10] K. Chaudhuri, C. Monteleoni, and A.-D. Sarwate. Differentially private empirical risk minimization. *JMLR*, 12:1069–1109, 2011.
- [11] P. Christen. *Data Matching Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection*. Springer Data-Centric Systems and Applications, 2012.
- [12] N. Cornia and J.-M. Mooij. Type-II errors of independence tests can lead to arbitrarily large errors in estimated causal effects: An illustrative example. In *30th UAI Workshops*, pages 35–42, 2014.
- [13] G. Doran, K. Muandet, K. Zhang, and B. Schölkopf. A permutation-based kernel conditional independence test. In *30th UAI*, 2014.
- [14] C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9:211–407, 2014.
- [15] M. Enserink and G. Chin. The end of privacy. *Science*, 347:490–491, 2015.
- [16] L. Getoor and A. Machanavajjhala. Entity resolution: theory, practice & open challenges. *Proceedings of the VLDB Endowment*, 5(12):2018–2019, 2012.
- [17] S. Graham, W. Press, S.-J. Gates, M. Gorenberg, J.-P. Holden, E. Lander, C. Mundie, M. Savitz, and E. Schmidt. Big data and privacy: a technological perspective. CreateSpace Independent Publishing Platform, 2015. — President’s Council of Advisors on Science and Technology.
- [18] A. Gretton, O. Bousquet, A.-J. Smola, and B. Schölkopf. Measuring statistical dependence with Hilbert-Schmidt norms. In *16th ALT*, pages 63–77, 2005.
- [19] A. Gretton, K. Fukumizu, C.-H. Teo, L. Song, B. Schölkopf, and A. Smola. A kernel statistical test of independence. In *NIPS*20*, pages 585–592, 2007.
- [20] P.-O. Hoyer, D. Janzing, J.-M. Mooij, J. Peters, and B. Schölkopf. Nonlinear causal discovery with additive noise models. In *NIPS*21*, pages 689–696, 2008.
- [21] D. Janzing, J.-M. Mooij, K. Zhang, J. Lemeire, J. Zscheischler, P. Daniusis, B. Steudel, and B. Schölkopf. Information-geometric approach to inferring causal directions. *Artificial Intelligence*, 182-183:1–31, 2012.
- [22] S. Kakade, K. Sridharan, and A. Tewari. On the complexity of linear prediction: Risk bounds, margin bounds, and regularization. In *NIPS*21*, pages 793–800, 2008.
- [23] A. Kolcz and C.-H. Teo. Feature weighting for improved classifier robustness. In *6th CEAS*, 2009.

- 307 [24] M.-J. Kusner, Y. Sun, K. Sridharan, and K.-Q. Weinberger. Inferring the causal direction privately. In *19th*
308 *AISTATS*, 2016.
- 309 [25] A. Lazarevic and V. Kumar. Feature bagging for outlier detection. In *11th KDD*, pages 157–166, 2005.
- 310 [26] J. Li, T.-D. Le, L. Liu, J. Liu, Z. Jin, B. Sun, and S. Ma. From observational studies to causal rule mining.
311 *ACM Trans. IST*, 7:1–27, 2016.
- 312 [27] J.-M. Mooij, J. Peters, D. Janzing, J. Zscheischler, and B. Schölkopf. Distinguishing cause from effect
313 using observational data: methods and benchmarks. *JMLR*, 2016.
- 314 [28] R. Nock and F. Nielsen. On the efficient minimization of classification-calibrated surrogates. In *NIPS*21*,
315 pages 1201–1208, 2008.
- 316 [29] R. Nock, G. Patrini, and A. Friedman. Rademacher observations, private data, and boosting. *32nd ICML*,
317 2015.
- 318 [30] J. O’Sullivan, J. Langford, R. Caruana, and A. Blum. Featureboost: A meta-learning algorithm that
319 improves model robustness. In *17th ICML*, pages 703–710, 2000.
- 320 [31] G. Patrini, R. Nock, S. Hardy, and T. Caetano. Fast learning from distributed datasets without entity
321 matching. In *26th IJCAI*, 2016.
- 322 [32] G. Patrini, R. Nock, P. Rivera, and T. Caetano. (Almost) no label no cry. In *NIPS*27*, 2014.
- 323 [33] J. Pearl and E. Bareinboim. External validity: from do-calculus to transportability across populations.
324 *Statistical Science*, 29:579–595, 2014.
- 325 [34] S.-J. Rizvi and J.-R. Haritsa. Maintaining data privacy in association rule mining. In *Proc. of the 28th*
326 *VLDB*, pages 682–693, 2002.
- 327 [35] L. Song, A. Smola, A. Gretton, J. Bedo, and K.-M. Borgwardt. Feature selection via dependence
328 maximization. *JMLR*, 13:1393–1434, 2012.
- 329 [36] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to
330 prevent neural networks from overfitting. *JMLR*, 15:1929–1958, 2014.
- 331 [37] C.-A. Sutton, M. Sindelar, and A. McCallum. Reducing weight undertraining in structured discriminative
332 learning. In *10th HLT-NAACL*, 2006.
- 333 [38] L. Sweeney. Achieving k -anonymity privacy protection using generalization and suppression. *Int. J.*
334 *Uncertainty, Fuzziness and Knowledge-based systems*, 10:571–588, 2002.
- 335 [39] L. van der Maaten, M. Chen, S. Tyree, and K.-Q. Weinberger. Learning with marginalized corrupted
336 features. In *30th ICML*, pages 410–418, 2013.
- 337 [40] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol. Extracting and composing robust features with
338 denoising autoencoders. In *25th ICML*, pages 1096–1103, 2008.
- 339 [41] Y.-X. Wang, S.-E. Fienberg, and A.-J. Smola. Privacy for free: Posterior sampling and stochastic gradient
340 monte carlo. In *32nd ICML*, pages 2493–2502, 2015.