

456 **Supplementary Material for “A General Framework for Learning**  
 457 **under Corruption: Label Noise, Attribute Noise, and Beyond”**

458 **S1 Additional discussions on related work**

459 Here, we detail discussions on the relations with existing paradigms as shown in Tab. 1. As a reminder,  
 460 we review the commonly used notations. Let  $E : Y \rightsquigarrow X$  be an experiment and  $F : X \rightsquigarrow Y$  be a  
 461 posterior kernel. The clean distribution  $P$  can be represented either in a discriminative manner as  
 462  $\pi_X \times F$  or in a generative manner as  $\pi_Y \times E$ . However, we cannot observe samples drawn from  
 463 the clean distribution  $P$ , but observe samples from some corrupted distribution  $\tilde{P}$ . The corruption is  
 464 generally represented as  $\kappa_{Z\tilde{Z}}$ , where the variables  $z = (x, y) \in Z$  are referred to as parameters and  
 465 the differentials  $d\tilde{z} = d\tilde{x}d\tilde{y}$  are referred to as corrupted variables.  $\delta_{Z\tilde{Z}}$  denotes a kernel induced by  
 466 the Dirac delta measure from  $(Z, \mathcal{Z})$  to  $(Z, \tilde{\mathcal{Z}})$ .

467 **S1.1 Simple corruptions**

468 The most well-known and widely studied corruptions in the literature are the simple cases, where  
 469 the corruption solely acts on the feature space  $X$  or the label space  $Y$ . We discuss examples of the  
 470 simple corruptions S- $\tilde{X}$  and S- $\tilde{Y}$ , as illustrated in Fig. 1a, in the following.

471 **Attribute noise** The problem of attribute noise concerns errors that are introduced into the observa-  
 472 tions of attribute  $X$ , leaving the labels untouched [30, 31, 4, 19]. Widely studied examples of such  
 473 errors include erroneous attribute values and missing attribute values. Instead of observing  $(X, Y)$ , in  
 474 the first case, one can only observe a distorted version of  $X$ , e.g.  $(X + N, Y)$  with some independent  
 475 noise random variable  $N \perp X$ ; in the second case, one’s observation of  $X$  contains missing values.

476 Let  $\mathbf{X} = (x_{ij})_{1 \leq i \leq n, 1 \leq j \leq d}$  be the complete input matrix, with  $|X| = n$ , and  $\mathbf{M} =$   
 477  $(m_{ij})_{1 \leq i \leq n, 1 \leq j \leq d}$  be the associated missingness indicator matrix such that  $m_{ij} = 1$  if  $x_{ij}$  is observed  
 478 and  $m_{ij} = 0$  if  $x_{ij}$  is missing. Then the corresponding observed input matrix is  $\mathbf{X}_o = \mathbf{X} \odot \mathbf{M}$  and  
 479 its missing counterpart is  $\mathbf{X}_m = \mathbf{X} - \mathbf{X}_o$ , where  $\odot$  denotes Hadamard product. The missing value  
 480 mechanisms are further categorized into three types based on their dependencies [5, 6]:<sup>8</sup>

- 481 • Missing completely at random (MCAR): the cause of missingness is entirely random, i.e.,  $p(\mathbf{M} |$   
 482  $\mathbf{X}) = p(\mathbf{M})$  does not depend on  $\mathbf{X}_o$  or  $\mathbf{X}_m$ . This corresponds to having a trivial Markov kernel  
 483 acting on the clean distribution,  $\kappa_{X\tilde{X}} \equiv \mu \in \mathcal{P}(X)$ .
- 484 • Missing not at random (MNAR): the cause of missingness depends on both observed variables and  
 485 missing variables, i.e.,  $p(\mathbf{M} | \mathbf{X}) = p(\mathbf{M} | \mathbf{X}_o, \mathbf{X}_m)$ . This case corresponds to our non-trivial  $\kappa_{X\tilde{X}}$ .
- 486 • Missing at random (MAR): the cause of missingness depends on observed variables but not on  
 487 missing variables, i.e.,  $p(\mathbf{M} | \mathbf{X}) = p(\mathbf{M} | \mathbf{X}_o)$ . This case is a sub-case of the non-trivial  $\kappa_{X\tilde{X}}$ ,  
 488 which is not specifiable by our taxonomy because of the different premises it is built on.

489 We underline that the conditional distributions of  $\mathbf{M}$  described above are *not* an equivalent description  
 490 of our Markov kernels. The missing data case is also known as finite Selection Bias, as discussed in  
 491 § S2.3, we know there exists a Markov kernel describing this corruption but the definition per se is a  
 492 *non-stochastic corruption*.

493 Hence, attribute noise is an example of S- $\tilde{X}$  corruption that can be generally formulated as the  
 494 corrupted experiment illustrated in the transition diagram  $Y \rightsquigarrow^E X \rightsquigarrow^{\kappa_{X\tilde{X}}} \tilde{X}$ , and the corrupted  
 495 distribution is given by  $\tilde{P} = (\kappa_{X\tilde{X}} \delta_{Y\tilde{Y}}) \circ (\pi_Y \times E)$ .

496 **Class-conditional noise (CCN)** The problem of CCN arises in situations where, instead of ob-  
 497 serving the clean labels, one can only observe corrupted labels that have been flipped with a  
 498 label-dependent probability, while the marginal distribution of the instance remains unchanged  
 499 [5, 34, 7, 19]. CCN is an example of S- $\tilde{Y}$  corruption that can be formulated as a corrupted pos-  
 500 terior illustrated in the transition diagram  $X \rightsquigarrow^F Y \rightsquigarrow^{\kappa_{Y\tilde{Y}}} \tilde{Y}$ , and the corrupted distribution is  
 501 given by  $\tilde{P} = (\delta_{X\tilde{X}} \kappa_{Y\tilde{Y}}) \circ (\pi_X \times F)$ . For classification tasks,  $Y$  and  $\tilde{Y}$  are assumed to be  
 502 finite spaces. Therefore the corruption  $\kappa_{Y\tilde{Y}}$  can be represented by a column-stochastic matrix

<sup>8</sup>Assume the rows  $x_i, m_i$  are assigned a joint distribution. and  $X$  and  $M$  are treated as random variables.

503  $\mathbf{T} = (\rho_{ij})_{1 \leq i \leq |\tilde{Y}|, 1 \leq j \leq |Y|}$  which specifies the probability of the clean label  $Y = j$  being flipped to  
504 the corrupted label  $\tilde{Y} = i$ , i.e.,  $\forall i, j, \rho_{ij} = p(\tilde{Y} = i | Y = j)$ . The corrupted joint distribution can be  
505 rewritten as  $\tilde{P} = \sum_Y p(\tilde{Y} | Y)p(Y | X)p(X)$ . In the literature,  $\mathbf{T}$  is known as the noise transition  
506 matrix with its elements  $\rho_{ij}$  referred to as the noise rates, and is useful for designing loss correction  
507 approaches (our results in § 5 significantly generalize existing loss correction results in CCN to our  
508 broad class of simple, dependent and combined corruptions) [34]. Prior to the proposal of the CCN  
509 model, early studies primarily focused on a symmetric subcase of  $\mathbf{T}$  in binary classification, known  
510 as random classification noise (RCN) [32, 33, 12]. Note that in RCN, the output of the corruption  $\kappa_{\tilde{Y}}$   
511 remains constant w.r.t. its parameters. Recently, some variants of CCN have been further developed,  
512 for example, in Ishida et al. [13, 14], complementary labels that can be modeled via a symmetric  $\mathbf{T}$   
513 whose diagonal elements are all equal to zero are studied.

## 514 S1.2 Dependent corruptions

515 Although simple corruptions have been well studied and understood, more complexities arise in  
516 dependent cases, yet they receive relatively less attention and understanding. We discuss examples of  
517 the dependent corruptions 1D- $\tilde{X}$ , 1D- $\tilde{Y}$ , 2D- $\tilde{X}$  and 2D- $\tilde{Y}$ , as illustrated in Fig. 1a, in the following.

518 **Style transfer** Style transfer refers to the process of migrating the artistic style of a given image to  
519 the content of another image [35, 36]. The primary objective is to recreate the second image with the  
520 designated style of the first image. In recent developments, it has also been applied to audio signals  
521 [37]. If we represent the style of the first image by  $Y$ , and the second image and the reconstructed  
522 image as  $X$  and  $\tilde{X}$  respectively, style transfer serves as an illustrative example of 1D- $\tilde{X}$  “corruption”.  
523 Note that the aim here is to *learn how to corrupt* instead of learning in the presence of corruption.  
524 We mention this connection because our framework can also be used also with different purposes,  
525 but underline that our BR results are not applicable to this case. The process of style transfer can be  
526 formulated as a corrupted posterior illustrated in the transition diagram  $X \xrightarrow{F} Y \xrightarrow{\kappa_{Y\tilde{X}}} \tilde{X}$ , and the  
527 corrupted distribution is given by  $\tilde{P} = (\kappa_{Y\tilde{X}}\delta_{Y\tilde{Y}}) \circ (\pi_X \times F)$ .

528 **Adversarial noise** In contrast to additive random attribute noise, adversarial noise is specifically  
529 crafted by adversaries for each instance with the intent of changing the model’s prediction of the  
530 correct label [38, 39, 40, 41, 42]. Such adversarial examples raise significant security concerns as  
531 they can be utilized to attack machine learning systems, even in scenarios where the adversary has no  
532 access to the underlying model. The adversarial noise is an example of 2D- $\tilde{X}$  corruption that can be  
533 formulated as a corrupted experiment illustrated in the transition diagram  $Y \xrightarrow{E} X \xrightarrow{\kappa_{XY\tilde{X}}} \tilde{X}$ , and  
534 the corrupted distribution is given by  $\tilde{P} = (\kappa_{XY\tilde{X}}\delta_{Y\tilde{Y}}) \circ (\pi_Y \times E)$ .

535 **Instance-dependent noise (IDN)** As a counterpart to CCN, the problem of IDN arises in situations  
536 where, instead of observing the clean labels, one can only observe corrupted labels that have been  
537 flipped with an instance-dependent (but not label-dependent) probability [18, 8]. It is a special case  
538 of the ILN noise model, which we will describe later. IDN is an example of 1D- $\tilde{Y}$  corruption that can  
539 be formulated as a corrupted experiment illustrated in the transition diagram  $Y \xrightarrow{E} X \xrightarrow{\kappa_{X\tilde{Y}}} \tilde{Y}$ ,  
540 and the corrupted distribution is given by  $\tilde{P} = (\delta_{X\tilde{X}}\kappa_{X\tilde{Y}}) \circ (\pi_Y \times E)$ .

541 **Instance- and label-dependent noise (ILN)** ILN is the most general label noise model, which  
542 arises in situations where, instead of observing clean labels, one can only observe corrupted labels that  
543 have been flipped with an instance- and label-dependent probability [8, 43, 44, 45]. ILN is an example  
544 of 2D- $\tilde{Y}$  corruption that can be formulated as a corrupted posterior illustrated in the transition diagram  
545  $X \xrightarrow{F} Y \xrightarrow{\kappa_{XY\tilde{Y}}} \tilde{Y}$ , and the corrupted distribution is given by  $\tilde{P} = (\delta_{X\tilde{X}}\kappa_{XY\tilde{Y}}) \circ (\pi_X \times F)$ .

546 Compared to the matrix representation  $\mathbf{T}$  of the CCN corruption  $\kappa_{Y\tilde{Y}}$ , the ILN corruption  $\kappa_{XY\tilde{Y}}$  can  
547 be represented by a matrix-valued function of the instance  $\mathbf{T}(x) = (\rho_{ij}(x))_{1 \leq i \leq |\tilde{Y}|, 1 \leq j \leq |Y|}$  which  
548 specifies the probability that the instance  $X = x$  with the clean label  $Y = j$  being flipped to the  
549 corrupted label  $\tilde{Y} = i$ , i.e.,  $\forall i, j, \rho_{ij}(x) = p(\tilde{Y} = i | Y = j, X = x)$ . Some subcases of ILN have

550 also been studied in the literature, for example, the boundary-consistent noise, which considers a  
 551 label flip probability based on a score function of the instance and label. The score aligns with the  
 552 underlying class-posterior probability function, resulting in instances closer to the optimal decision  
 553 boundary having a higher chance of its label being flipped [23].

### 554 S1.3 Combined corruptions

555 Given the simple and dependent corruptions, we can combine them to generate 2-parameter joint  
 556 corruptions, i.e.,  $\kappa_{ZZ} : X \times Y \rightsquigarrow \tilde{X} \times \tilde{Y}$ . Below, we discuss some examples of combined noise  
 557 models illustrated in Fig. 1b.

558 **Combined simple noise** The simplest combined corruption is the combined simple noise, where  
 559 the observations of attribute  $X$  are subject to some errors and the observed labels  $Y$  are flipped  
 560 with a label-dependent probability [19]. Combined simple noise is an example of (S- $\tilde{X}$ , S- $\tilde{Y}$ )  
 561 corruption that can be formulated as a corrupted experiment illustrated in the transition diagram

$$562 \tilde{Y} \xrightarrow{\kappa_{Y\tilde{Y}}} Y \xrightarrow{E} X \xrightarrow{\kappa_{X\tilde{X}}} \tilde{X}, \text{ and the corrupted distribution is given by } \tilde{P} = (\kappa_{X\tilde{X}}\kappa_{Y\tilde{Y}}) \circ (\pi_Y \times E).$$

563 **Target shift** In the literature, target shift refers to the situation where the prior probability  $p(Y)$  is  
 564 changed while the conditional distribution  $p(X | Y)$  remains invariant across training and test domains  
 565 [46, 47, 48, 49]. The definition is established by assuming certain invariance from a generative  
 566 perspective of the learning problem, that is, considering it as a corruption of the experiment according  
 567 to  $P = \pi_Y \times E$ . However, when examining the learning problem from a discriminative perspective,  
 568 the change in  $p(Y)$  may cause changes in both  $p(X)$  and  $p(Y | X)$  due to the Bayes rule. Existing  
 569 frameworks for the categorization of target shift do not capture these implications, as they are based  
 570 on the notion of invariance from a single perspective of the  $E$  direction. In contrast, our framework  
 571 categorizes corruptions based on their dependencies and therefore is advantageous by offering dual  
 572 perspectives from both the  $E$  and  $F$  directions. Specifically, target shift is an example of (1D- $\tilde{X}$ ,  
 573 2D- $\tilde{Y}$ ) corruption and can be formulated either as a corrupted experiment illustrated in the transition

$$574 \text{ diagram } \tilde{X} \xrightarrow{\kappa_{Y\tilde{X}}} Y \xrightarrow{E} X \xrightarrow{\kappa_{XY\tilde{Y}}} \tilde{Y}, \text{ or as a corrupted posterior illustrated in the transition diagram}$$

$$575 \tilde{Y} \xrightarrow{\kappa_{XY\tilde{Y}}} X \xrightarrow{F} Y \xrightarrow{\kappa_{Y\tilde{X}}} \tilde{X}. \text{ The corrupted distribution is given by } \tilde{P} = (\kappa_{Y\tilde{X}}\kappa_{XY\tilde{Y}}) \circ (\pi_Y \times E)$$

$$576 \text{ or } \tilde{P} = (\kappa_{Y\tilde{X}}\kappa_{XY\tilde{Y}}) \circ (\pi_X \times F).$$

**Mutually contaminated distributions (MCD)** The problem of MCD arises in binary classification  
 situations where, instead of observing samples from the clean class-conditional distributions  $p(X |$   
 $Y = \pm 1)$ , one can only observe samples from corrupted class-conditional distributions  $\tilde{p}(X | Y =$   
 $\pm 1)$ , with

$$\begin{pmatrix} \tilde{p}(X | Y = +1) \\ \tilde{p}(X | Y = -1) \end{pmatrix} = \begin{pmatrix} 1 - \alpha & \alpha \\ \beta & 1 - \beta \end{pmatrix} \begin{pmatrix} p(X | Y = +1) \\ p(X | Y = -1) \end{pmatrix}$$

577 as described in [6, 50, 51]. The coefficients  $\alpha$  and  $\beta$  are defined as the fraction of data points having  
 578 a flipped label, given that the true one was respectively +1 or -1.

579 In comparison, CCN corrupts the class-posterior probability  $p(Y | X)$  while MCD corrupts the  
 580 class-conditional distribution  $p(X | Y)$ ; consequently, the marginal distribution of  $p(X)$  remains  
 581 unchanged in CCN but may be changed in MCD. Therefore  $\alpha$  and  $\beta$  in MCD are not the noise rates  
 582  $\rho_{12}$  and  $\rho_{21}$  in CCN, however, they are shown to be related by an invertible transformation [6]. In  
 583 other words, CCN is shown to be a subclass of the MCD, but what else is included in the MCD model  
 584 is not explored. Therefore, here we model MCD as (2D- $\tilde{X}$ , S- $\tilde{Y}$ ) corruption, which can be formulated

$$585 \text{ as a corrupted experiment illustrated in the transition diagram } \tilde{Y} \xrightarrow{\kappa_{Y\tilde{Y}}} Y \xrightarrow{E} X \xrightarrow{\kappa_{XY\tilde{X}}} \tilde{X}, \text{ and the}$$

$$586 \text{ corrupted distribution is given by } \tilde{P} = (\kappa_{XY\tilde{X}}\kappa_{Y\tilde{Y}}) \circ (\pi_Y \times E).$$

587 **Covariate shift** In the literature, covariate shift refers to the situation where the marginal distribution  
 588  $p(X)$  is changed while the class-posterior probability  $p(Y | X)$  remains invariant across training and  
 589 test domains [3, 52, 53, 54]. Similarly to target shift, the definition is established by assuming certain

590 invariance from a discriminative perspective of the learning problem. However, when examining the  
 591 learning problem from a generative perspective, the change in  $p(X)$  may cause changes in  $p(Y)$  and  
 592  $p(X | Y)$  due to the Bayes rule. Covariate shift in its general definition is an example of (2D- $\tilde{X}$ , 1D- $\tilde{Y}$ )  
 593 corruption and can be formulated either as a corrupted posterior illustrated in the transition diagram

594  $\tilde{Y} \xrightarrow{\kappa_{X\tilde{Y}}} X \xrightarrow{F} Y \xrightarrow{\kappa_{XY\tilde{X}}} \tilde{X}$ , or as a corrupted experiment illustrated in the transition diagram

595  $\tilde{X} \xrightarrow{\kappa_{XY\tilde{X}}} Y \xrightarrow{E} X \xrightarrow{\kappa_{X\tilde{Y}}} \tilde{Y}$ . The corrupted distribution is given by  $\tilde{P} = (\kappa_{XY\tilde{X}} \kappa_{X\tilde{Y}}) \circ (\pi_X \times F)$

596 or  $\tilde{P} = (\kappa_{XY\tilde{X}} \kappa_{X\tilde{Y}}) \circ (\pi_Y \times E)$ .

597 **Generalized target shift** In the literature, generalized target shift refers to the situation where  
 598 the prior probability  $p(Y)$  and the conditional distribution  $p(X | Y)$  both change across training  
 599 and test domains, however, with some invariance assumptions in the latent space [55, 56, 57].

600 Generalized target shift is an example of (2D- $\tilde{X}$ , 2D- $\tilde{Y}$ ) corruption that can be formulated as a

601 corrupted experiment illustrated in the transition diagram  $\tilde{Y} \xrightarrow{\kappa_{XY\tilde{X}}} Y \xrightarrow{E} X \xrightarrow{\kappa_{XY\tilde{X}}} \tilde{X}$ , and the

602 corrupted distribution is given by  $\tilde{P} = (\kappa_{XY\tilde{X}} \kappa_{X\tilde{Y}}) \circ (\pi_Y \times E)$ . Note that simpler scenarios can  
 603 also result in a generalized target shift, however, it is important to avoid degenerating to the simple  
 604 S- $\tilde{X}$  corruption, as it would violate the requirement of corrupting the label distribution.

605 **Concept drift** Concept drift refers to the situation where  $\tilde{p}(Y | X) \neq p(Y | X)$  [17]. As in the case  
 606 of generalized target shift, this case can be associated with every corruption in our framework, so the  
 607 most general correspondence is the (2D- $\tilde{X}$ , 2D- $\tilde{Y}$ ) joint Markov kernel.

## 608 S2 Appendix for “A general framework for corruption”, Section 3

### 609 S2.1 The superposition operation

610 We further describe the superposition operation between kernels, also known as “parallel combination”  
 611 [27] or Kronecker product when in finite spaces, by specifying the action of the resulting kernel on  
 612 functions and measures.

613 **Definition S1.** Let  $\kappa_1$  be a Markov kernel from  $(X, \mathcal{X})$  to  $(Y, \mathcal{Y})$  and  $\kappa_2$  be a Markov kernel from  
 614  $(Z, \mathcal{Z})$  to  $(W, \mathcal{W})$ . Hence, the **superposition** of the two is a kernel  $\kappa_1 \kappa_2$  from  $(X \times Z, \mathcal{X} \times \mathcal{Z})$  to  
 615  $(Y \times W, \mathcal{Y} \times \mathcal{W})$  such that:

$$\begin{aligned} (\kappa_1 \kappa_2) f(x, z) &= \int_{Y \times W} (\kappa_1 \kappa_2)(x, z, dydw) f(y, w) \\ &= \int_Y \kappa_1(x, dy) \int_W \kappa_2(z, dw) f(y, w), \end{aligned}$$

616 for every  $f$  positive  $\mathcal{Y} \times \mathcal{W}$ -measurable, or equivalently

$$\begin{aligned} \mu(\kappa_1 \kappa_2)(B) &= \int_B \int_{X \times Z} (\kappa_1 \kappa_2)(x, z, dydw) \mu(dx dz) \\ &= \int_B \int_X \kappa_1(x, dy) \int_Z \kappa_2(z, dw) \mu(dx dz), \end{aligned}$$

617 for every measure  $\mu$  on  $(X \times Z, \mathcal{X} \times \mathcal{Z})$ ,  $B \in \mathcal{Y} \times \mathcal{W}$ .

618 Both the operators are well defined, as we can rewrite them

$$\begin{aligned} (\kappa_1 \kappa_2) f(x, z) &= \kappa_1(\kappa_2 f(y, w)) = \kappa_1 \hat{f}(y, z), \\ \mu(\kappa_1 \kappa_2)(B_Y \times B_W) &= \kappa_1(\mu \kappa_2)(B_Y \times B_W) = \hat{\mu} \kappa_1(B_Y \times B_W). \end{aligned}$$

619 Hence we are just iteratively applying the standard kernel-induced operators to a parameterized  
 620 function or partially to a joint measure.

When dealing with finite spaces, Markov kernels are column-stochastic matrices. The superposition operation is then equal to the *Kronecker product*, between two matrices,

$$\kappa_1 \kappa_2 := \kappa_1 \otimes \kappa_2 = \begin{pmatrix} (\kappa_1)_{11} \kappa_2 & \dots & (\kappa_1)_{1n} \kappa_2 \\ \vdots & \ddots & \vdots \\ (\kappa_1)_{m1} \kappa_2 & \dots & (\kappa_1)_{mn} \kappa_2 \end{pmatrix}$$

621 with  $\kappa_1$  being a  $|Y| \times |X|$  matrix and  $\kappa_2$  a  $|W| \times |Z|$  matrix.

## 622 S2.2 The Bayesian inversion theorem

623 In this section, we present some existing results coming from category theory applied to Bayesian  
624 learning, which allows us to define and use the inverse kernel as introduced in the main text. The  
625 background knowledge required for following this section is rather different from that of the other  
626 sections. However, even if readers choose not to delve into the specific details, they can still  
627 comprehend our results by only referring to the notions in [Def. S4](#).

628 In Dahlqvist et al. [26], they address the question of Markov kernel inversion through the lens of  
629 category theory.<sup>9</sup> They investigate how and when the (weak) inversion is defined, both directly on the  
630 category of measurable spaces and indirectly by considering the associated Markov linear operator  
631 (Markov transition [42]). We only focus on illustrating the first result, given the focus on Markov  
632 kernels we had in the paper. We will use category theory terminology, and then connect it to our  
633 probabilistic vocabulary.

634 The first step is the construction of the  $\text{Krn}$  category, similar to our notion of space of Markov kernels  
635  $\mathcal{M}(X, Y)$  but with an equivalence relation acting on it. They start by considering Polish metric  
636 spaces, the category  $\text{Pol}$  with continuous mappings, a subcase of which are the closed sets  $X \subseteq \mathbb{R}^d$   
637 equipped with the usual topology. The category of measurable spaces considered for defining kernels,  
638  $\text{Mes}$ , is the one induced by a functor  $\mathcal{B} : \text{Pol} \rightarrow \text{Mes}$ , i.e. all the measurable spaces with the same  
639 underlying set of a Polish space but equipped by the Borel  $\sigma$ -algebra and interpreting continuous  
640 mapping as measurable ones. We call these spaces *standard Borel spaces*, and use them as the  
641 building block of the  $\text{Krn}$  category.

642 The category  $\text{Mes}$  is embedded by a functor  $F$  into the Kleisli category of  $\mathbb{G}$ , a monad over  $\text{Mes}$   
643 representing probability distributions over some set. The functor  $F$  acts *identically* on sets and  
644 *maps* measurable functions  $f : X \rightarrow Y$  to Kleisli arrows  $F(f) = \delta_Y \circ f$ . This means, in more  
645 familiar terms, that we build a trivial kernel  $\delta_Y(f^{-1}(dy))$ , i.e. the image measure of the dirac delta  
646 through  $f$ . It further induces the category  $1 \downarrow F$  of probabilities  $p : 1 \rightarrow \mathbb{G}X$  and trivial morphisms  
647  $f : (X, p) \rightarrow_\delta (Y, q)$  as degenerate arrows  $F(f)$  s.t.  $q = F(f) \circ_{\mathbb{G}} p = \mathbb{G}(f)(p)$ , where  $\circ_{\mathbb{G}}$  corresponds  
648 our  $\circ$  combination between a kernel and a distribution. In other words,  $F(f)$  induces a measure-  
649 preserving map, so  $1 \downarrow F$  includes all measure-preserving maps induced by *degenerate* arrows. When  
650 the arrows used are not degenerate, we obtain the supercategory  $1 \downarrow \mathcal{Kl}$ , with the same objects. We  
651 denote arrows in this category as  $f : (X, p) \rightarrow (Y, q)$ . Notice that the kernels included in the category  
652  $1 \downarrow \mathcal{Kl}$  are what we would call  $\mathcal{M}(X, Y)$ , where  $X$  has marginal  $p$  and  $Y$  has marginal  $q$ .

653 This last category includes Markov kernels as we have defined them in this paper. They are considered  
654 as *typed kernels*, i.e. their definition is tied to a fixed input and a fixed output (probabilities) instead  
655 of being characterized for every input probability and every reachable output. This remark is crucial  
656 for understanding our notion of exhaustiveness – we will later underline why.

657 Markov kernels cannot be inverted as they are, because of their non-singularity. They characterize it  
658 with their Lemma 3, proving that for a kernel  $f : (X, p) \rightarrow (Y, q)$  there are *p-negligibly many points*  
659 *jumping to q-negligible sets*. Once the non-singularity is understood, we can define an equivalence  
660 relation that, when acting on  $1 \downarrow \mathcal{Kl}$ , allows a well-posed definition of the inverse kernel.

**Definition S2.** For all objects  $(X, p), (Y, q)$ ,  $R_{(X,p),(Y,q)}$  is the smallest equivalence relation on  $\text{Hom}_{1 \downarrow \mathcal{Kl}}(X, Y)$  such that

$$(f, f') \in R_{(X,p),(Y,q)} \Leftrightarrow f = f' \cdot p - a.s.$$

661 In Lemma 4, they prove  $R$  to be a congruence relation on  $1 \downarrow \mathcal{Kl}$ . This congruence relation allows us  
662 to define the quotient category, with the proper morphisms.

<sup>9</sup>For a general overview, see Mac Lane [41].

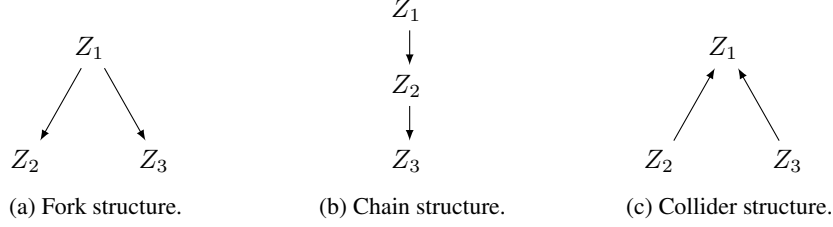


Figure S1: Possible non-degenerate relations among three probability spaces.

663 **Definition S3.** *The category  $\text{Krn}$  is the quotient category  $(1 \downarrow \mathcal{Kl})/R$ .*

Having defined the category, we have to build the functions that are going to constitute the weak inversion operator, i.e. a bijection between  $\text{Hom}_{\text{Krn}}((X, p), (Y, q))$  and  $\text{Hom}_{\text{Krn}}((Y, q), (X, p))$ . They are two mapping between the  $\text{Krn}$  category and the space of couplings associated to  $(X, p), (Y, q)$ . The first is equivalent to our  $\times$  kernel operation, applied to a kernel (i.e. conditional probability) and a probability, and is formally written as

$$\alpha_Y^X : \text{Hom}_{\text{Krn}}((X, p), (Y, q)) \rightarrow \Gamma((X, p), (Y, q)) \quad \text{s.t.} \quad \alpha_Y^X(f)(B_X \times B_Y) := \int_{x \in B_X} f(x)(B_Y) dp,$$

664 with  $\Gamma((X, p), (Y, q)) \subset \mathcal{G}(X, Y)$  the typed couplings associated to the marginals  $(X, p), (Y, q)$ .

The second is defined as its inverse operation, and it is decomposing a joint probability along a fixed marginal distribution, i.e.,

$$D_Y^X : \Gamma((X, p), (Y, q)) \rightarrow \text{Hom}_{\text{Krn}}((X, p), (Y, q)) \quad \text{s.t.} \quad D_Y^X(\gamma) := G(\pi_Y) \circ \pi_X^\dagger,$$

such that

$$\gamma(B_X \times B_Y) := \int_{x \in B_X} D_Y^X(\gamma)(x)(B_Y) dp,$$

665 with  $(\cdot)^\dagger$ : adjoint operator. Being one the inverse of the other, they are both obviously bijective and  
666 proving the one-to-one correspondence between typed kernels and typed couplings.

667 Hence, we formally define the pseudo-inverse as in the following:

668 **Definition S4.** *The inverse of a typed kernel  $\kappa$  from  $(X_1, p_1)$  to  $(X_2, p_2)$ , given by  $\kappa^\dagger \circ \kappa :=$   
669  $D_X^Y \circ \mathcal{G}(\pi_2 \times \pi_1) \circ \alpha_Y^X(\kappa)$  with  $\mathcal{G}(\pi_2 \times \pi_1) : \Gamma((X_1, p_1), (X_2, p_2)) \rightarrow \Gamma((X_2, p_2), (X_1, p_1))$  being  
670 the permutation map, is defined as*

- 671 1.  $\kappa^\dagger : (X_2, p_2) \rightarrow (X_1, p_1) \in \text{Krn}$  when  $\kappa$  is seen as element of  $\text{Krn}$ , such that  $\kappa^\dagger \circ \kappa = \delta_{X_1}$  and  
672  $\kappa \circ \kappa^\dagger = \delta_{X_2}$ ;
- 673 2.  $\kappa^\dagger : X_2 \rightsquigarrow X_1 \in \mathcal{M}(X_2, X_1)$  when  $\kappa$  is seen as element of  $\mathcal{M}(X_1, X_2)$ , such that  $\kappa^\dagger \circ \kappa =_R \delta_{X_1}$   
674 and  $\kappa \circ \kappa^\dagger =_R \delta_{X_2}$ .

### 675 S2.3 Exhaustiveness of the taxonomy

676 In the previous section, we define the operations  $\alpha$  and  $D$  for typed kernels, which are one the inverse  
677 of the other by construction. They are the operations representing the bijection between the space  
678 of Markov kernels typed for  $p, q$  and the space of couplings with marginals  $p, q$ . Hence, they are  
679 proving that for each couple of probability spaces, there exists a Markov kernel sending one into the  
680 other corresponding to a possible associated coupling.

681 This means that every pairwise stochastic corruption in the supervised learning setting is described  
682 by our taxonomy. Other possibilities are, having more than two spaces involved in the corruption  
683 process and having a deterministic mapping describing the corruption process as it has been defined.  
684 We discuss them in the following, providing examples.

685 **Stochastic corruption for more than two spaces** When in the presence of three probability spaces,  
686 we have only two possible corruption configurations. We represent them in Fig. S1, where arrows  
687 represent non-trivial Markov kernels. We remark that we do not consider the triangular structures

688 as in Fig. S1a and c with the spaces  $Z_2, Z_3$  coupled in some way, otherwise they would just be  
 689 considered as a single (product) probability space, i.e. a pairwise corruption.

690 The most simple case is Fig. S1b, in which the spaces influence each other in a chain fashion. This is  
 691 a clear subcase of our framework as we can integrate  $Z_2$  by considering  $\kappa_{Z_1 Z_3} := \kappa_3 \circ \kappa_2 \circ \kappa_1$ . We  
 692 then obtain a pairwise corruption  $\kappa_{Z_1 Z_3}$ , but we would pay the price of losing information about the  
 693 role of the ‘latent’ corruption process. To have a complete idea of how the chained corruption works,  
 694 we can additionally study it as an iterative process and analyze its single components. This entire  
 695 reasoning is true for a number of spaces  $Z_i, i \in [n]$  with  $n > 3$ , and well models several settings for  
 696 dynamical learning, e.g. online corrupted learning or concept drift over time [58, 61, 59].

697 The second option is, they act as per the diagrams in Fig. S1a and Fig. S1c, i.e. a triangular structure.  
 698 In particular, case (a) reflects assumptions made in settings combining data from different domains  
 699 [29, 7, 62], where we get to observe different data distributions obtained from the same clean one.  
 700 They can be seen, in our framework, as a pairwise dependence between  $Z_1 \times Z_2$  and  $Z_3$ , or  $Z_1$   
 701 and  $Z_2 \times Z_3$ . However, this formulation assumes some coupling on  $Z_2, Z_3$ , more complex than  
 702 our originally assumed corruption. For now, we do not investigate the consequences of this gap as  
 703 changes of the corruption effect, leaving it for future investigation. A similar idea can be stated  
 704 for  $n > 2$  spaces in the Cartesian product space, and for combinations of fork structures with fork  
 705 structures via superposition.

706 **Corruptions via deterministic mappings** We now want to give examples of how corruption  
 707 processes can be not stochastic. From the previous section, we know that even if there is no direct  
 708 way of modeling the specific corruption process with a Markov kernel, there exists a Markov kernel  
 709 representing some coupling between two distributions. We do not define any method to find the *best*  
 710 Markov corruption corresponding to a deterministic one, since it depends on the specific task one is  
 711 considering.

712 The first relevant example is the one of **Selection Bias**. Even if being a widely studied, common case  
 713 of corruption, we show here that when considered with its classical formulation we cannot directly  
 714 find a Markov kernel corresponding to it.

715 We start by introducing the Selection Bias type corruption, as done in [14]. It is characterized here as  
 716 a distributional corruption, unlike other cases in which only the selection variable is modeled. We  
 717 consider a target, clean distribution and a source, corrupted one from which we aim to learn. They  
 718 are defined on the same set  $Z \subseteq \mathbb{R}^d$  and Borel  $\sigma$ -algebra  $\mathcal{Z}$ . We define it as:

719 1. Support condition:

$$\tilde{P} \ll P \iff \exists! \alpha = \frac{d\tilde{P}}{dP} \text{ a.s., } \alpha \in L^1,$$

720 where we can equivalently say  $\mu - \text{a.s.}, \mu := 0.5 * (\tilde{P} + P)$ , or  $P - \text{a.s.}$ ;

721 2. Selection condition:

$$\|\alpha\|_\infty < +\infty.$$

722 The Support condition is equivalent to:

$$\tilde{P}(A) = \int_A \alpha(z)P(dz) \forall A.$$

723 Comparing it with a Markov kernel action on the input probability  $P$ , we get the condition

$$\int_A \int_Z \kappa(z, d\tilde{z})P(dz) = \int_A \alpha(z)P(dz) \quad \forall A.$$

724 A guess that satisfies our requirement is  $\kappa(z, d\tilde{z}) := \delta_z(d\tilde{z})\alpha(z)$ , which is a transition kernel, i.e. a  
 725 family of positive measures parameterized by  $z$ , *but not* a Markov kernel unless  $\alpha \equiv 1$ . This kernel is  
 726 defined such that  $P$  is corrupted into  $\tilde{P}$ , but it does not preserve mass for every probability measure.  
 727 *Is this the only possible guess?*

728 Assuming the existence of a transition kernel  $\hat{\kappa}(z, d\tilde{z}) \neq \delta_z(d\tilde{z})\alpha(z)$ , possibly Markov, implies

$$\int_A \int_Z \hat{\kappa}(z, d\tilde{z})P(dz) = \int_A \tilde{P}(dz) = \int_A \alpha(z)P(dz) \forall A.$$

729 We then can define a measure  $\hat{\mu}(dz) := \int_{\mathcal{Z}} \hat{\kappa}(z, d\tilde{z})P(dz)$ . This measure  $\hat{\mu}$  is almost surely equal to  
730  $\tilde{P}$  by definition, w.r.t. a reference measure  $\mu_1$ . The same is true for  $\tilde{P}$  and  $\alpha P$  w.r.t.  $\mu_2$ . Hence  $\hat{\mu}(dz)$   
731 is equal to  $\alpha(z)P(dz)$  w.r.t. to  $\mu$ ,  $\mu_1 \ll \mu$  and  $\mu_2 \ll \mu$ .

732 Since the same argument can be repeated for the kernel  $\kappa$ , the two kernels are forced to be equal  
733  $\mu$ -almost surely. That because, for two measures with the same value on every set, their Radon-  
734 Nikodym derivatives are the same almost everywhere w.r.t. a finite<sup>10</sup> reference measure. Hence,  
735 Selection Bias cannot be directly represented as a Markov kernel if we impose it to be acting *exactly*  
736 as the weak derivative  $\alpha$ .

737 Another relevant example is the one of **Markov kernel reconstruction**  $R$ , as introduced in [7]. They  
738 are considered in finite space settings and are defined as the *left inverse of the stochastic matrix*  
739 *representing the Markov kernel*. It is underlined by the authors that the  $R$  of the Markov kernel is  
740 not necessarily a Markov kernel; in fact, it is not even ensured to be a matrix with positive entries.  
741 A reconstruction  $R$  is then sending a corrupted probability  $\tilde{P}$  into the original clean probability  $P$   
742 without being a stochastic mapping.

### 743 S3 Appendix for “Consequences of corruption in supervised learning”, 744 Section 4: Proofs

745 We restate for clarity all the assumptions underlying the proofs.

746 **A1** We assume the loss function to be bounded in order to avoid problems when applying Fubini-  
747 Tonelli’s theorem.

748 **A2** We define the set  $\ell \circ \mathcal{H} := \{ (x, y) \mapsto \ell(h(x), y) \mid h \in \mathcal{H} \}$ .

749 **A3** When minimizing the risk for the corrupted distribution  $\tilde{P}$ , we assume that  $f^* \in$   
750  $\arg \min_f \mathbb{E}_{\tilde{P}}[f(X, Y)]$  belongs to the minimization space  $\ell \circ \mathcal{H}$ .

751 Theorems 3, 4 are here proved by means of two Lemmas on the dependent noise combined with  
752 identical simple noise.

753 **Lemma S5** (BR under X corruption). *Let  $(\ell, \mathcal{H}, P)$  be a learning problem with the input space  $X$*   
754 *and output space  $Y$ . Let  $E : Y \rightsquigarrow X$  be an experiment,  $\kappa^{\tilde{X}} \in \{ \kappa_{X\tilde{X}}, \kappa_{YX\tilde{X}} \}$  be the corruption*  
755 *on  $X$  with at most 2 parameters, then we can form the corrupted experiment as per the transition*

756 *diagram*  $Y \xrightarrow{\tilde{E}} X \xrightarrow{\kappa^{\tilde{X}}} \tilde{X}$  and obtain

$$\mathbb{E}_{Y \sim \pi_Y} CBR_{\ell \circ \mathcal{H}}(\kappa^{\tilde{X}} E_Y) = \mathbb{E}_{Y \sim \pi_Y} CBR_{\kappa^{\tilde{X}}(\ell \circ \mathcal{H})}(E_Y) .$$

757 Moreover, if  $\kappa^{\tilde{X}} = \kappa_{X\tilde{X}}$ , we have

$$BR_{\ell \circ \mathcal{H}}(\pi_Y \times \kappa^{\tilde{X}} E) = BR_{\kappa^{\tilde{X}}(\ell \circ \mathcal{H})}(\pi_Y \times E) .$$

758 *Proof.* Assume the full corruption  $\kappa$  has an associated kernel

$$\kappa(x, y, d\tilde{x}d\tilde{y}) := (\kappa^{\tilde{X}} \delta)(x, y, d\tilde{x}d\tilde{y}) = \kappa_y^{\tilde{X}}(x, d\tilde{x})\delta_y(d\tilde{y}), \quad (\text{S1})$$

759 Let  $E_y(dx) := E(y, dx)$  and  $A \in \tilde{\mathcal{X}} \times \tilde{\mathcal{Y}}$ , we have

$$\begin{aligned} \tilde{P}(A) &= \sum_Y \int_A \int_X \kappa(x, y, d\tilde{x}d\tilde{y}) P(dx dy) \\ &= \sum_Y \int_A \int_X \kappa_y^{\tilde{X}}(x, d\tilde{x})\delta_y(d\tilde{y}) E_y(dx)\pi_y \\ &= \sum_Y \int_A (\kappa^{\tilde{X}} E)_y(d\tilde{x}) \delta_y(d\tilde{y})\pi_y \\ &= \int_A \tilde{E}_{\tilde{y}}(d\tilde{x})\pi_{\tilde{y}}, \end{aligned}$$

<sup>10</sup>It is enough to ask “finite on all balls”, see [63], Theorem 5.8.8.



760 then we can write

$$\begin{aligned}
\mathbb{E}_{\tilde{Y} \sim \pi_{\tilde{Y}}} CBR_{\ell \circ \mathcal{H}}(\kappa^{\tilde{X}} E_{\tilde{Y}}) &= \sum_{\tilde{Y}} \pi_{\tilde{y}} \inf_{f \in \ell \circ \mathcal{H}} \int_{\tilde{X}} f(\tilde{x}, \tilde{y}) \tilde{E}_{\tilde{y}}(d\tilde{x}) \\
&= \sum_{Y, \tilde{Y}} \delta_y(d\tilde{y}) \pi_y \inf_{f \in \ell \circ \mathcal{H}} \int_{\tilde{X} \times X} f(\tilde{x}, \tilde{y}) \kappa_y^{\tilde{X}}(x, d\tilde{x}) E_y(dx) \\
&= \sum_Y \pi_y \inf_{f \in \ell \circ \mathcal{H}} \int_{\tilde{X}} \delta f(\tilde{x}, y) \int_X \kappa_y^{\tilde{X}}(x, d\tilde{x}) E_y(dx) \\
&= \sum_Y \pi_y \inf_{f \in \ell \circ \mathcal{H}} \int_X E_y(dx) (\kappa_y^{\tilde{X}} \delta f)(x, y) \tag{S2} \\
&= \sum_Y \pi_y \inf_{f \in \kappa^{\tilde{X}}(\ell \circ \mathcal{H})} \int_X E_y(dx) f(x, y) \\
&= \mathbb{E}_{Y \sim \pi_Y} CBR_{\kappa^{\tilde{X}}(\ell \circ \mathcal{H})}(E_Y). \tag{S3}
\end{aligned}$$

Since the X corruption  $\kappa^{\tilde{X}}$  has an identity mapping on Y,  $\mathbb{E}_{\tilde{Y} \sim \pi_{\tilde{Y}}}[\cdot] = \mathbb{E}_{Y \sim \pi_Y}[\cdot]$  and we obtain

$$\mathbb{E}_{Y \sim \pi_Y} CBR_{\ell \circ \mathcal{H}}(\kappa^{\tilde{X}} E_Y) = \mathbb{E}_{Y \sim \pi_Y} CBR_{\kappa^{\tilde{X}}(\ell \circ \mathcal{H})}(E_Y).$$

761 If  $\kappa^{\tilde{X}} = \kappa_{X\tilde{X}}$ , then the associated kernel (S1) takes the simple form  $\kappa^{\tilde{X}}(x, d\tilde{x})$  and the above  
762 equations from (S2) become

$$\begin{aligned}
BR_{\ell \circ \mathcal{H}}(\pi_Y \times \kappa^{\tilde{X}} E) &= \inf_{f \in \ell \circ \mathcal{H}} \sum_Y \pi_y \int_X E_y(dx) (\kappa^{\tilde{X}} f)(x, y) \\
&= \inf_{f \in \kappa^{\tilde{X}}(\ell \circ \mathcal{H})} \sum_Y \pi_y \int_X E_y(dx) f(x, y) \\
&= BR_{\kappa^{\tilde{X}}(\ell \circ \mathcal{H})}(\pi_Y \times E).
\end{aligned}$$

763 □

764 **Theorem** (BR under (S- $\tilde{X}$ , S- $\tilde{Y}$ ), (2D- $\tilde{X}$ , S- $\tilde{Y}$ ) joint corruption, **Theorem 3**). *Let  $(\ell, \mathcal{H}, P)$  be*  
765 *a learning problem,  $E : Y \rightsquigarrow X$  an experiment and  $\kappa^{\tilde{X}} \in \{\kappa_{X\tilde{X}}, \kappa_{YX\tilde{X}}\}$  a corruption as in*  
766 ***Lemma S5**. Let  $\kappa_{Y\tilde{Y}}$  be a simple corruption on Y. Then we can form the corrupted experiment as*  
767 *per the transition diagram  $\tilde{Y} \xrightarrow{\kappa_{Y\tilde{Y}}} Y \xrightarrow{E} X \xrightarrow{\kappa^{\tilde{X}}} \tilde{X}$  and obtain*

$$\mathbb{E}_{\tilde{Y} \sim \kappa_{Y\tilde{Y}} \pi_Y} CBR_{\ell \circ \mathcal{H}}(\kappa^{\tilde{X}} E_{\tilde{Y}}) = \mathbb{E}_{Y \sim \pi_Y} CBR_{\kappa^{\tilde{X}}(\kappa_{Y\tilde{Y}} \ell \circ \mathcal{H})}(E_Y)$$

*Proof.* We assume the full corruption  $\kappa$  has an associated kernel

$$\kappa(x, y, d\tilde{x}d\tilde{y}) := (\kappa^{\tilde{X}} \kappa_{Y\tilde{Y}})(x, y, d\tilde{x}d\tilde{y}) = \kappa_y^{\tilde{X}}(x, d\tilde{x}) \kappa_y^{Y\tilde{Y}}(d\tilde{y}).$$

768 With this corruption formulation, we can replicate the proof of **Lemma S5** up to (S2) by simply  
769 plugging in  $\kappa_y^{Y\tilde{Y}}(d\tilde{y})$  instead of  $\delta_y(d\tilde{y})$ . Therefore, we obtain the thesis. □

770 **Lemma S6** (BR under Y corruption). *Let  $E : Y \rightsquigarrow X$  and  $F : X \rightsquigarrow Y$  be an experiment and a*  
771 *posterior on it, and  $\kappa^{\tilde{Y}} \in \{\kappa_{Y\tilde{Y}}, \kappa_{XY\tilde{Y}}\}$  be the corruption on Y with at most 2 parameters, then we*  
772 *can form the corrupted posterior as per the transition diagram  $X \xrightarrow{F} Y \xrightarrow{\kappa^{\tilde{Y}}} \tilde{Y}$  and obtain*

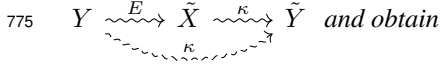
$$\mathbb{E}_{X \sim \pi_X} CBR_{\ell \circ \mathcal{H}}(\kappa^{\tilde{Y}} F_X) = \mathbb{E}_{X \sim \pi_X} CBR_{\kappa^{\tilde{Y}} \ell \circ \mathcal{H}}(F_X).$$

773 Moreover, if  $\kappa^{\tilde{Y}} = \kappa_{Y\tilde{Y}}$ , the equation simplifies as

$$BR_{\ell \circ \mathcal{H}}(\pi_X \times \kappa^{\tilde{Y}} F) = BR_{\kappa^{\tilde{Y}}(\ell \circ \mathcal{H})}(\pi_X \times F).$$

774 Equivalently, we can form the corrupted experiment as per the transition diagram

775  $Y \xrightarrow{E} \tilde{X} \xrightarrow{\kappa} \tilde{Y}$  and obtain



$$BR_{\ell \circ \mathcal{H}}(\kappa(\pi_Y \times E)) = BR_{\kappa(\ell \circ \mathcal{H})}(\pi_Y \times E),$$

776 where  $\kappa = \kappa^{\tilde{Y}} \delta_x$ , the combination of  $\kappa^{\tilde{Y}}$  with the identity kernel on  $X$ .

777 *Proof.* Assume the full corruption  $\kappa$  has an associated kernel

$$\kappa(x, y, d\tilde{x}d\tilde{y}) := (\kappa^{\tilde{Y}} \delta)(x, y, d\tilde{x}d\tilde{y}) = \kappa_x^{\tilde{Y}}(y, d\tilde{y})\delta_x(d\tilde{x}). \quad (\text{S4})$$

778 Let  $F_x(dy) := F(x, dy)$  and  $A \in \tilde{\mathcal{X}} \times \tilde{\mathcal{Y}}$ , we have

$$\begin{aligned} \tilde{P}(A) &= \sum_Y \int_A \int_X \kappa(x, y, d\tilde{x}d\tilde{y}) P(dx dy) \\ &= \sum_Y \int_A \int_X \kappa_x^{\tilde{Y}}(y, d\tilde{y})\delta_x(d\tilde{x}) F_x(dy)\pi_x \\ &= \int_A \int_X (\kappa^{\tilde{Y}} F)_x(d\tilde{y}) \delta_x(d\tilde{x})\pi_x \\ &= \int_A \tilde{F}_{\tilde{x}}(d\tilde{y})\pi_{\tilde{x}}, \end{aligned}$$

779 then we can write

$$\begin{aligned} \mathbb{E}_{\tilde{X} \sim \pi_{\tilde{X}}} CBR_{\ell \circ \mathcal{H}}(\kappa^{\tilde{Y}} F_{\tilde{X}}) &= \int_{\tilde{X}} \pi_{\tilde{x}} \inf_{f \in \ell \circ \mathcal{H}} \sum_{\tilde{Y}} f(\tilde{x}, \tilde{y}) \tilde{F}_{\tilde{x}}(d\tilde{y}) \\ &= \int_{\tilde{X} \times X} \delta_x(d\tilde{x})\pi_x \inf_{f \in \ell \circ \mathcal{H}} \sum_{Y, \tilde{Y}} f(\tilde{x}, \tilde{y}) \kappa_x^{\tilde{Y}}(y, d\tilde{y}) F_x(dy) \\ &= \int_X \pi_x \inf_{f \in \ell \circ \mathcal{H}} \sum_{Y, \tilde{Y}} \delta f(x, \tilde{y}) \kappa_x^{\tilde{Y}}(y, d\tilde{y}) F_x(dy) \\ &= \int_X \pi_x \inf_{f \in \ell \circ \mathcal{H}} \sum_Y F_x(dy) (\kappa_x^{\tilde{Y}} \delta f)(x, y) \end{aligned} \quad (\text{S5})$$

$$\begin{aligned} &= \int_X \pi_x \inf_{f \in \kappa^{\tilde{Y}}(\ell \circ \mathcal{H})} \sum_Y F_x(dy) f(x, y) \\ &= \mathbb{E}_{X \sim \pi_X} CBR_{\kappa^{\tilde{Y}}(\ell \circ \mathcal{H})}(F_X). \end{aligned} \quad (\text{S6})$$

Since the  $Y$  corruption  $\kappa^{\tilde{Y}}$  has an identity mapping on  $X$ ,  $\mathbb{E}_{\tilde{X} \sim \pi_{\tilde{X}}}[\cdot] = \mathbb{E}_{X \sim \pi_X}[\cdot]$  and we obtain

$$\mathbb{E}_{X \sim \pi_X} CBR_{\ell \circ \mathcal{H}}(\kappa^{\tilde{Y}} F_X) = \mathbb{E}_{X \sim \pi_X} CBR_{\kappa^{\tilde{Y}}(\ell \circ \mathcal{H})}(F_X).$$

780 If  $\kappa^{\tilde{Y}} = \kappa_{Y\tilde{Y}}$ , then the associated kernel (S4) takes the simple form  $\kappa^{\tilde{Y}}(y, d\tilde{y})$  and the above  
781 equations from (S5) become

$$\begin{aligned} BR_{\ell \circ \mathcal{H}}(\pi_X \times \kappa^{\tilde{Y}} F) &= \inf_{f \in \ell \circ \mathcal{H}} \int_X \pi_x \sum_Y F_x(dy) (\kappa^{\tilde{Y}} f)(x, y) \\ &= \inf_{f \in \kappa^{\tilde{Y}}(\ell \circ \mathcal{H})} \int_X \pi_x \sum_Y F_x(dy) f(x, y) \\ &= BR_{\kappa^{\tilde{Y}}(\ell \circ \mathcal{H})}(\pi_X \times F). \end{aligned}$$

In this case, reminding that  $h(x, y) = (h(x), id(y))$ , we have

$$BR_{\kappa^{\tilde{Y}}(\ell \circ \mathcal{H})}(\pi_X \times F) = BR_{\kappa^{\tilde{Y}}(\ell \circ \mathcal{H})}(\pi_X \times F).$$

782 Similarly, the results can also be expressed in terms of  $E$  using the generic corruption formulation:

$$\begin{aligned}
BR_{\ell \circ \mathcal{H}}(\kappa \pi_Y \times E) &= BR_{\ell \circ \mathcal{H}}(\kappa P) = \inf_{f \in \ell \circ \mathcal{H}} \int_{\tilde{X}} \sum_{\tilde{Y}} f(\tilde{x}, \tilde{y}) \kappa P(d\tilde{x}d\tilde{y}) \\
&= \inf_{f \in \ell \circ \mathcal{H}} \int_{\tilde{X} \times X} \sum_{Y, \tilde{Y}} f(\tilde{x}, \tilde{y}) \kappa(x, y, d\tilde{x}d\tilde{y}) P(dxdy) \\
&= \inf_{f \in \ell \circ \mathcal{H}} \int_X \sum_Y \kappa f(x, y) P(dxdy) \\
&= BR_{\kappa(\ell \circ \mathcal{H})}(\pi_Y \times E) \tag{S7}
\end{aligned}$$

783 Note that the last result in (S7), even if fitting in the comparison of experiments literature [20], does  
784 not give us any new insights since it is not based on the corruption decomposition formula in (S4).  
785 We provide (S7) here for completeness.  $\square$

786 **Theorem** (BR under (S- $\tilde{X}$ , S- $\tilde{Y}$ ), (S- $\tilde{X}$ , 2D- $\tilde{Y}$ ) joint corruption, **Theorem 4**). *Let  $(\ell, \mathcal{H}, P)$  be a  
787 learning problem,  $F : X \rightsquigarrow Y$  a posterior and  $\kappa^{\tilde{Y}} \in \{\kappa_{Y\tilde{Y}}, \kappa_{XY\tilde{Y}}\}$  a corruption as in **Lemma S6**.  
788 Let  $\kappa_{X\tilde{X}}$  be a simple corruption on  $X$ . Then we can form the corrupted experiment as per the*

789 transition diagram  $\tilde{X} \xrightarrow{\kappa_{X\tilde{X}}} X \xrightarrow{F} Y \xrightarrow{\kappa^{\tilde{Y}}} \tilde{Y}$  and obtain

$$\mathbb{E}_{\tilde{X} \sim \kappa_{X\tilde{X}} \pi_X} CBR_{\ell \circ \mathcal{H}}(\kappa^{\tilde{Y}} F_{\tilde{X}}) = \mathbb{E}_{X \sim \pi_X} CBR_{\kappa_{X\tilde{X}}(\kappa^{\tilde{Y}} \ell \circ \mathcal{H})}(F_X). \tag{S8}$$

*Proof.* We assume the full corruption  $\kappa$  has an associated kernel

$$\kappa(x, y, d\tilde{x}d\tilde{y}) := (\kappa^{\tilde{Y}} \kappa^{X\tilde{X}})(x, y, d\tilde{x}d\tilde{y}) = \kappa_x^{\tilde{Y}}(y, d\tilde{y}) \kappa_x^{X\tilde{X}}(d\tilde{x}).$$

790 With this corruption formulation, we can replicate the proof of **Lemma S6** up to (S5) by simply  
791 plugging in  $\kappa_x^{X\tilde{X}}(d\tilde{x})$  instead of  $\delta_x(d\tilde{x})$ . Therefore, we obtain the thesis.  $\square$

792 **Remark S7.** *When using the continuous notation for  $Y$  we do so for simplicity and homogeneity.  
793 Notice that all its associated kernel are actually (parameterized) squared matrices, hence transposable.  
794 So in **Theorem 4** and the associated Lemma, the operator acting on the function is actually the  
795 transpose of the corruption matrix.*

$$\sum_{\tilde{y}} C_{\tilde{y}y}(x) \ell_{\tilde{y}}(h(x)) = \sum_{\tilde{y}} C_{y\tilde{y}}^T(x) \ell_{\tilde{y}}(h(x)) = (\ell_y \circ h)_x^*(x).$$

796 **Theorem** (BR under (1D, 2D) joint corruption, **Theorem 5**). *Let  $(\ell, \mathcal{H}, P)$  be a learning problem,  
797  $E : Y \rightsquigarrow X$  and  $F : X \rightsquigarrow Y$  be an experiment and a posterior on it.*

798 1. *Let  $\kappa_{Y\tilde{X}}$  be a corruption on  $X$  and  $\kappa_{XY\tilde{Y}}$  be a corruption on  $Y$ , then we can form the jointly  
799 corrupted experiment as per the transition diagram  $\tilde{X} \xrightarrow{\kappa_{Y\tilde{X}}} Y \xrightarrow{E} X \xrightarrow{\kappa_{XY\tilde{Y}}} \tilde{Y}$  and obtain*

$$BR_{\ell \circ \mathcal{H}}[\kappa_{Y\tilde{X}}(\pi_Y \times \kappa_{XY\tilde{Y}} E)] = \mathbb{E}_{Y \sim \pi_Y} CBR_{\kappa_{Y\tilde{X}}(\kappa_{XY\tilde{Y}} \ell \circ \mathcal{H})}(E_Y). \tag{S9}$$

800 2. *Let  $\kappa_{X\tilde{Y}}$  be a corruption on  $Y$  and  $\kappa_{XY\tilde{X}}$  be a corruption on  $X$ , then we can form the jointly  
801 corrupted posterior as per the transition diagram  $\tilde{Y} \xrightarrow{\kappa_{X\tilde{Y}}} X \xrightarrow{F} Y \xrightarrow{\kappa_{XY\tilde{X}}} \tilde{X}$  and obtain*

$$BR_{\ell \circ \mathcal{H}}[\kappa_{X\tilde{Y}}(\pi_X \times \kappa_{XY\tilde{X}} F)] = \mathbb{E}_{X \sim \pi_X} CBR_{\kappa_{X\tilde{Y}}(\kappa_{XY\tilde{X}} \ell \circ \mathcal{H})}(F_X). \tag{S10}$$

802 *Proof.* For proving point (1), assume the full corruption  $\kappa$  has an associated kernel

$$\kappa(x, y, d\tilde{x}d\tilde{y}) := (\kappa_{Y\tilde{X}} \kappa_{XY\tilde{Y}})(x, y, d\tilde{x}d\tilde{y}) = \kappa^{Y\tilde{X}}(y, d\tilde{x}) \kappa_y^{XY\tilde{Y}}(x, d\tilde{y}).$$

803 Let  $E_y(dx) := E(y, dx)$  and  $A \in \tilde{\mathcal{X}} \times \tilde{\mathcal{Y}}$ , we have

$$\begin{aligned}
\tilde{P}(A) &= \sum_Y \int_A \int_X \kappa(x, y, d\tilde{x}d\tilde{y}) P(dx dy) \\
&= \sum_Y \int_A \int_X \kappa^{Y\tilde{X}}(y, d\tilde{x}) \kappa_y^{XY\tilde{Y}}(x, d\tilde{y}) E_y(dx) \pi(dy) \\
&= \int_A \sum_Y \kappa^{Y\tilde{X}}(y, d\tilde{x}) (\pi_Y \times (\kappa^{XY\tilde{Y}} E))(dy, d\tilde{y}) \\
&= \int_A \kappa^{Y\tilde{X}}(\pi_Y \times \kappa^{XY\tilde{Y}} E)(d\tilde{x}, d\tilde{y}) ,
\end{aligned}$$

804 which is less interpretable as a corruption action if compared to the previous theorems, since the  
805 effect on  $E$  and  $\pi_Y$  cannot be totally distinguished. However, we can still write

$$\begin{aligned}
\mathbb{E}_{\tilde{Y} \sim \tilde{\pi}_{\tilde{Y}}} CBR_{\ell \circ \mathcal{H}}(\kappa_{Y\tilde{X}} \kappa_{XY\tilde{Y}} E) &= \sum_Y \inf_{f \in \ell \circ \mathcal{H}} \int_{\tilde{X}} f(\tilde{x}, \tilde{y}) \tilde{P}(d\tilde{x}d\tilde{y}) \\
&= \sum_Y \int_{\tilde{X}} \kappa^{Y\tilde{X}}(y, d\tilde{x}) \pi_y \inf_{f \in \ell \circ \mathcal{H}} \sum_{\tilde{Y}} \int_X f(\tilde{x}, \tilde{y}) \kappa_x^{XY\tilde{Y}}(y, d\tilde{y}) E_y(dx) \\
&= \sum_Y \int_{\tilde{X}} \kappa^{Y\tilde{X}}(y, d\tilde{x}) \pi_y \inf_{\ell \circ h \in \ell \circ \mathcal{H}} \int_X (\kappa_x^{XY\tilde{Y}} \ell)(h(\tilde{x}), y) E_y(dx) \\
&= \sum_Y \pi_y \inf_{f \in \ell \circ \mathcal{H}} \int_X \kappa^{Y\tilde{X}}[(\kappa_x^{XY\tilde{Y}} \ell)_y \circ h](y) E_y(dx) \\
&= \sum_Y \pi_y \inf_{f \in \kappa_{Y\tilde{X}}(\kappa_{XY\tilde{Y}} \ell \circ \mathcal{H})} \int_X E_y(dx) f(x, y, h) \\
&= \mathbb{E}_{Y \sim \pi_Y} CBR_{\kappa_{Y\tilde{X}}(\kappa_{XY\tilde{Y}} \ell \circ \mathcal{H})}(E) ,
\end{aligned}$$

806 with  $f(x, y, h) := \kappa^{Y\tilde{X}}[(\kappa_x^{XY\tilde{Y}} \ell)_y \circ h](y)$ . In particular, notice that  $\kappa^{XY\tilde{Y}}$  acts only on  $\ell$ , while  
807  $\kappa_{Y\tilde{X}}$  acts on both  $\ell$  and  $h$ , which forces us to use  $f(x, y, h)$  instead of  $f(x, y)$ .

808 For point (2), assume the full corruption  $\kappa$  has an associated kernel

$$\kappa(x, y, d\tilde{x}d\tilde{y}) := (\kappa_{X\tilde{Y}} \kappa_{XY\tilde{X}})(x, y, d\tilde{x}d\tilde{y}) = \kappa^{X\tilde{Y}}(x, d\tilde{y}) \kappa_x^{XY\tilde{X}}(y, d\tilde{x}) .$$

809 Let  $F_x(dy) := F(x, dy)$  and  $A \in \tilde{\mathcal{X}} \times \tilde{\mathcal{Y}}$ , we have

$$\begin{aligned}
\tilde{P}(A) &= \sum_Y \int_A \int_X \kappa(x, y, d\tilde{x}d\tilde{y}) P(dx dy) \\
&= \sum_Y \int_A \int_X \kappa^{X\tilde{Y}}(x, d\tilde{y}) \kappa_x^{XY\tilde{X}}(y, d\tilde{x}) F_x(dy) \pi(dx) \\
&= \int_A \int_X \kappa^{X\tilde{Y}}(x, d\tilde{y}) (\pi_X \times \kappa^{XY\tilde{X}} F)(d\tilde{x}, dx) \\
&= \int_A \kappa^{X\tilde{Y}}(\pi_X \times \kappa^{XY\tilde{X}} F)(d\tilde{x}, d\tilde{y}) ,
\end{aligned}$$

810 Hence we can repeat a similar argument for the  $F$  case and find a minimization space of functions  
811  $f(x, y, h) := \kappa^{XY\tilde{X}}[(\kappa^{X\tilde{Y}} \ell)_x \circ h](x)$ . Thus, we obtain the thesis.  $\square$

812 **Corollary** (BR under (1D, 1D) joint corruption, **Corollary 6**). *Let  $(\ell, \mathcal{H}, P)$  be a learning problem,*  
813  *$E : Y \rightsquigarrow X$  and  $F : X \rightsquigarrow Y$  be an experiment and a posterior on it. Let  $\kappa_{Y\tilde{X}}$  be a corruption on  $X$*   
814 *and  $\kappa_{X\tilde{Y}}$  be a corruption on  $Y$ , then we can form the jointly corrupted experiment as per the transition*  
815 *diagram  $\tilde{X} \xrightarrow{\kappa_{Y\tilde{X}}} Y \xrightarrow{E} X \xrightarrow{\kappa_{X\tilde{Y}}} \tilde{Y}$  or equivalently  $\tilde{Y} \xrightarrow{\kappa_{X\tilde{Y}}} X \xrightarrow{F} Y \xrightarrow{\kappa_{Y\tilde{X}}} \tilde{X}$ . We*  
816 *obtain*

$$BR_{\ell \circ \mathcal{H}}[\kappa_{Y\tilde{X}}(\pi_Y \times \kappa_{X\tilde{Y}} E)] = BR_{\kappa_{Y\tilde{X}}(\kappa_{X\tilde{Y}} \ell \circ \mathcal{H})}(\pi_Y \times E) ,$$

817 or equivalently

$$BR_{\ell \circ \mathcal{H}}[\kappa_{X\tilde{Y}}(\pi_X \times \kappa_{Y\tilde{X}} F)] = BR_{\kappa_{X\tilde{Y}}(\kappa_{Y\tilde{X}} \ell \circ \mathcal{H})}(\pi_X \times F) .$$

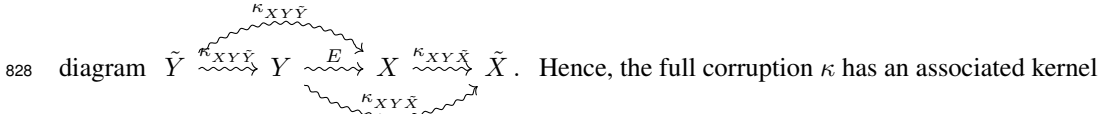
*Proof.* We assume the full corruption  $\kappa$  has an associated kernel

$$\kappa(x, y, d\tilde{x}d\tilde{y}) := (\kappa^{Y\tilde{X}} \kappa^{X\tilde{Y}})(x, y, d\tilde{x}d\tilde{y}) = \kappa^{Y\tilde{X}}(y, d\tilde{x}) \kappa^{X\tilde{Y}}(x, d\tilde{y}) .$$

818 With this corruption formulation, we can replicate the proof of **Theorem 5** by simply plugging in  
 819  $\kappa^{X\tilde{Y}}(x, d\tilde{y})$  instead of  $\kappa^{XY\tilde{Y}}(x, y, d\tilde{y})$  in the first point, and by simply plugging in  $\kappa^{Y\tilde{X}}(y, d\tilde{x})$   
 820 instead of  $\kappa^{XY\tilde{X}}(x, y, d\tilde{x})$  in the second point. We then in both cases obtain functions  $f(x, y, h) :=$   
 821  $\kappa^{Y\tilde{X}}[(\kappa^{X\tilde{Y}}\ell)_x \circ h](y)$ , i.e. comparing a point  $x$  with a kernel on  $\mathcal{P}(X)$  parameterized by  $y$ .

822 After getting the identities w.r.t. CBRs, we can further take the inf operator out of the outside  
 823 expectations and obtain identities w.r.t. BRs, as the kernels  $\kappa^{Y\tilde{X}}$  and  $\kappa^{X\tilde{Y}}$  are not parameterized by  
 824  $x$  or  $y$  anymore. Therefore, we obtain the thesis.  $\square$

825 **Analysis of Bayes Risk under (2D- $\tilde{X}$ , 2D- $\tilde{Y}$ )** Let  $(\ell, \mathcal{H}, P)$  be a learning problem,  $E : Y \rightsquigarrow X$   
 826 and  $F : X \rightsquigarrow Y$  be an experiment and a posterior on it. Let  $\kappa_{XY\tilde{X}}$  be the corruption on  $X$  and  $\kappa_{XY\tilde{Y}}$   
 827 be the corruption on  $Y$ . Then we can form the jointly corrupted experiment as per the transition



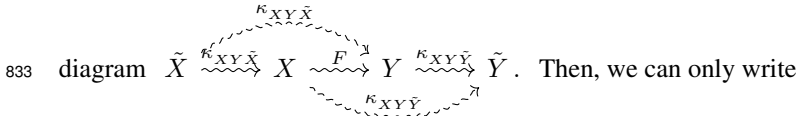
$$\kappa(x, y, d\tilde{x}d\tilde{y}) := (\kappa_{XY\tilde{X}} \kappa_{XY\tilde{Y}})(x, y, d\tilde{x}d\tilde{y}) = \kappa_y^{\tilde{X}}(x, d\tilde{x}) \kappa_x^{\tilde{Y}}(y, d\tilde{y}) .$$

829 Let  $\tilde{P} = \kappa P$  and  $A \in \tilde{\mathcal{X}} \times \tilde{\mathcal{Y}}$ , we have

$$\begin{aligned} \tilde{P}(A) &= \sum_Y \int_A \int_X \kappa(x, y, d\tilde{x}d\tilde{y}) P(dxdy) \\ &= \sum_Y \int_A \int_X \kappa_y^{\tilde{X}}(x, d\tilde{x}) \kappa_x^{\tilde{Y}}(y, d\tilde{y}) E_y(dx) \pi_y . \end{aligned}$$

830 This is not further decomposable as an action on the experiment and an action on the prior, given the  
 831 double dependence of both factors in the kernel.

832 A similar observation can be done for the posterior kernel, i.e. considering the equivalent transition



$$\begin{aligned} BR_{\ell \circ \mathcal{H}}(\kappa(\pi_Y \times E)) &= \inf_{f \in \ell \circ \mathcal{H}} \mathbb{E}_{(\tilde{X}, \tilde{Y}) \sim \tilde{P}} f(\tilde{X}, \tilde{Y}) \\ &= \inf_{f \in \ell \circ \mathcal{H}} \sum_{\tilde{Y}} \tilde{\pi}_{\tilde{y}} \int_{\tilde{X}} f(\tilde{x}, \tilde{y}) \tilde{E}_{\tilde{y}}(d\tilde{x}) \\ &= \inf_{f \in \ell \circ \mathcal{H}} \sum_{Y, \tilde{Y}} \pi_y \int_{\tilde{X} \times X} \kappa_x^{\tilde{Y}}(y, d\tilde{y}) \kappa_y^{\tilde{X}}(x, d\tilde{x}) f(\tilde{x}, \tilde{y}) E_y(dx) , \end{aligned}$$

834 which is, from a Bayesian Risk point of view, equivalent to the non-decomposed joint corruption  
 835 effect.

## 836 S4 Appendix for ‘‘Corruption-corrected learning’’, Section 5

837 In this section, we give the proofs of the results used in § 5.

### 838 S4.1 Proof of Lemma 7: Factorization of the pseudo-inverse

839 **Lemma.** *The feasible factorization of a Markov transition  $\kappa$  is also a valid factorization for its*  
 840 *pseudo-inverse  $\kappa^\dagger$ , both for the full transition or considering their parameterized versions.*

841 *Proof.* Let's consider a Markov kernel  $\kappa : X_1 \times Y_1 \rightarrow X_2 \times Y_2$ . Also assume that  $\kappa = \kappa_X \kappa_Y$ ,  
 842 i.e. factorizes by superposition with  $\kappa_{(\cdot)} : X_1 \times Y_1 \rightarrow (\cdot)_2, (\cdot)_2 \in \{X_2, Y_2\}$ .

843 Supposing that  $\kappa_X^\dagger \kappa_Y^\dagger$  is a pseudo-inverse, we can write by using the definition of pseudo-inverse  
 844 **Def. S4:**

$$\begin{aligned} \delta_{x'_1, y'_1}(dx_1 dy_1) &= \int_{X_2 \times Y_2} \kappa(dx_2 dy_2, x'_1, y'_1) \kappa^\dagger(dx_1 dy_1, x_2, y_2) \\ &= \int_{X_2 \times Y_2} (\kappa_X)_{y'_1}(dx_2, x'_1) (\kappa_Y)_{x'_1}(dy_2, y'_1) (\kappa_X^\dagger)_{y'_1}(dx_1, x_2) (\kappa_Y^\dagger)_{x'_1}(dy_1, y_2) \end{aligned}$$

845 which shows that  $\kappa^\dagger =_R \kappa_X^\dagger \kappa_Y^\dagger$ , and proves the Lemma for the parameterized case. Being the regular  
 846 case a subcase, we obtain the thesis.  $\square$

## 847 S4.2 Proof for Theorem 8

848 In addition to the assumptions **A1** – **A3** stated for proving the BR theorems (§ S3), we assume here:

849 **A4** We will assume the existence of an invertible function  $\tilde{h}^* \in \mathcal{H}$ ;

850 **A5** We ask the corrupted optimum to satisfy the equality  $\kappa^\dagger(\ell \circ \tilde{h}^*) = \tilde{\ell} \circ \tilde{h}^*$ .

851 **Theorem.** Let  $(\ell, \mathcal{H}, P)$  be a clean learning problem and  $(\kappa^\dagger(\ell \circ \mathcal{H}), \kappa P)$  its associated corrupted  
 852 one, not necessarily with a  $\circ$ -factorized structure. Let  $\kappa^\dagger$  be the joint cleaning kernel reversing  $\kappa$ ,  
 853 such that assumptions **A4** and **A5** hold for the said problems. The factorization of  $\kappa^\dagger$  is assumed to be  
 854 feasible and to have an equality result of the form Eq. (5). We write  $\kappa^\dagger(dz, \tilde{z}) = \kappa^X(dx, \cdot) \kappa^Y(dy, \cdot)$ ,  
 855 with  $(\cdot)$  some feasible parameters. Hence, we can prove the following points:

856 1. When  $\kappa^\dagger$  is either  $(id_X, S-Y)$  or  $(id_X, 2D-Y)$ , we can write the corrected loss as

$$\tilde{\ell}(h(\tilde{x}), \tilde{y}) = (\kappa^Y \ell)(h(\tilde{x}), \tilde{y}) \quad \forall (\tilde{x}, \tilde{y}) \in \tilde{X} \times \tilde{Y},$$

857 with  $\kappa^Y \ell = \kappa_{\tilde{x}}^Y \ell$  for the second case.

858 2. When  $\kappa^\dagger$  is  $(S-X, S-Y)$ ,  $(2D-X, S-Y)$  or  $(S-X, 2D-Y)$ , we have

$$\tilde{\ell}(\tilde{x}, \tilde{y}, h) = \mathbb{E}_{u \sim \kappa^X h(\tilde{x})} [\kappa^Y \ell(u, \tilde{y})] \quad \forall (\tilde{x}, \tilde{y}) \in \tilde{X} \times \tilde{Y},$$

859 with  $\kappa_{\tilde{x}}^X h(\tilde{x})(A) := \kappa^X(h^{-1}(A), \tilde{x})$ ,  $A \subset \mathcal{P}(Y)$  being the push-forward probability measure of  
 860  $\kappa^X(\cdot, \tilde{x})$  through  $h$ ,  $h$  seen as a function. For the cases that involve a 2D corruption, we have  
 861  $\kappa^Y \ell = \kappa_{\tilde{x}}^Y \ell$  for the former  $\kappa^\dagger$  factorization,  $\kappa^X h(\tilde{x}) = \kappa_{\tilde{y}}^X h(\tilde{x})$  for the latter.

862 3. When  $\kappa^\dagger$  is a  $(1D-X, 1D-Y)$  corruption, we can write the corrected loss as

$$\tilde{\ell}(\tilde{x}, \tilde{y}, h) = \mathbb{E}_{u \sim \kappa^X h(\tilde{y})} [\kappa^Y \ell(u, \tilde{x})] \quad \forall (\tilde{x}, \tilde{y}) \in \tilde{X} \times \tilde{Y},$$

863 with  $\kappa_{\tilde{x}}^X h(\tilde{y})(B) := \kappa^X(h^{-1}(B), \tilde{y})$ ,  $B \subset \mathcal{P}(X)$ .

864 4. When  $\kappa^\dagger$  is a  $(2D, 1D)$  corruption, we can write the corrected loss as

$$\tilde{\ell}(\tilde{x}, \tilde{y}, h) = \mathbb{E}_{u \sim \kappa^X h(\tilde{y})} [\kappa_{\tilde{x}}^Y \ell(u, \tilde{y})], \quad \tilde{\ell}(\tilde{x}, \tilde{y}, h) = \mathbb{E}_{u \sim \kappa_{\tilde{y}}^X h(\tilde{x})} [\kappa^Y \ell(u, \tilde{x})] \quad \forall (\tilde{x}, \tilde{y}) \in \tilde{X} \times \tilde{Y}.$$

865 for the  $(1D-X, 2D-Y)$ ,  $(2D-X, 1D-Y)$  respectively.

*Proof.* Given the assumptions **A4**, **A5**, we can write:

$$\tilde{\ell}(\tilde{h}^*(\tilde{x}), \tilde{y}) = \sum_Y \int_X \ell(\tilde{h}^*(x), y) \kappa^\dagger(dx dy, \tilde{x}, \tilde{y}).$$

866 We now look at all the feasible corruption combinations in Fig. 1b; given Lemma S4.1, are sure that  
 867 there factorizations on  $\kappa$  are also valid for  $\kappa^\dagger$ . Hence, we can consider the single point of the theorem  
 868 being sure that they cover every possible  $\kappa$  case having an Bayes Risk equality result.

869 Consider the  $\kappa^\dagger$  from point (1), i.e.  $\kappa^\dagger$  is either  $(id_X, S\text{-}Y)$  or  $(id_X, 2D\text{-}Y)$ . They act on  $\ell \circ h$  as

$$\begin{aligned}\tilde{\ell}(\tilde{h}^*(\tilde{x}), \tilde{y}) &= \sum_Y \int_X \ell(\tilde{h}^*(x), y) \delta(dx, \tilde{x}) \kappa^Y(dy, \tilde{x}, \tilde{y}) \\ &= \int_X (\kappa^Y \ell)_{\tilde{x}}(h(x), \tilde{y}) \delta(dx, \tilde{x}) \\ &= (\kappa^Y \ell)_{\tilde{x}}(h(\tilde{x}), \tilde{y}).\end{aligned}$$

870 Hence, the case  $\kappa^Y(dy, \tilde{x}, \tilde{y}) = \kappa_{\tilde{x}}^Y(dy, \tilde{x}, \tilde{y})$  and its subcase  $\kappa^Y(dy, \tilde{y})$  combined with an identity  
871 kernel on  $X$  do not change the hypothesis function.

872 For the more complex cases in point (2),  $\kappa^X(dx, \tilde{x}) \neq \delta_x(dx)$ , we have:

$$\begin{aligned}\tilde{\ell}(\tilde{h}^*(\tilde{x}), \tilde{y}) &= \sum_Y \int_X \ell(\tilde{h}^*(x), y) \kappa^X(dx, \tilde{x}) \kappa^Y(dy, \tilde{x}, \tilde{y}) \\ &= \sum_Y \int_{\tilde{h}^*(X)} \ell(u, y) \kappa^X((\tilde{h}^*)^{-1}(du), \tilde{x}) \kappa^Y(dy, \tilde{x}, \tilde{y}) \\ &= \int_{\tilde{h}^*(X)} (\kappa^Y \ell)_{\tilde{x}}(u, \tilde{y}) \kappa^X((\tilde{h}^*)^{-1}(du), \tilde{x}),\end{aligned}\tag{S11}$$

where  $u = u(dy) \in \mathcal{P}(Y)$ . The following equality for the expectation of  $u$ , the image measure of  $\kappa^\dagger$  through  $\tilde{h}^*$ , and the kernel chain composition holds:

$$\mathbb{E}_{\kappa^X((\tilde{h}^*)^{-1}(\cdot), \tilde{x})}[u] = \int_{\tilde{h}^*(X)} u \kappa^X((\tilde{h}^*)^{-1}(du), \tilde{x}) = \kappa^X \circ \tilde{h}^*(\tilde{x}) \in \mathcal{P}(Y),$$

873 that can be verified easily by recalling the alternative definition of  $\mathcal{H}$  as a subset of  $\mathcal{M}(X, Y)$  and  
874 using the definition of  $\kappa^\dagger \circ \tilde{h}^*$ . We remark that  $\kappa^X((\tilde{h}^*)^{-1}(du), \tilde{x})$  is then a probability in  $\mathcal{P}(\mathcal{P}(Y))$ .  
875 Hence we can rewrite Eq. (S11) as

$$\begin{aligned}\tilde{\ell}(\tilde{h}^*(\tilde{x}), \tilde{y}) &= \int_{\mathcal{P}(Y)} (\kappa^Y \ell)_{\tilde{x}}(u, \tilde{y}) \kappa^X((\tilde{h}^*)^{-1}(du), \tilde{x}) \\ &= \mathbb{E}_{\kappa^X((\tilde{h}^*)^{-1}(\cdot), \tilde{x})}[(\kappa^Y \ell)_{\tilde{x}}(u, \tilde{y})],\end{aligned}$$

876 with  $\kappa^X$  having support included in  $\tilde{h}^*(X)$ .

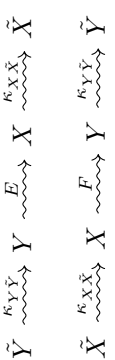
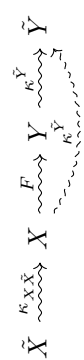
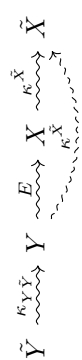
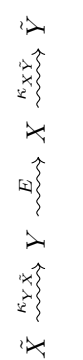
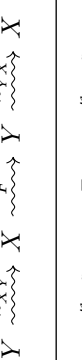
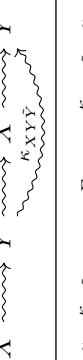

As for more dependent corruptions of  $X$ , i.e.  $\kappa^X(dx, \tilde{x}, \tilde{y})$ , the action on the hypothesis will be dependent from  $\tilde{y}$ , so

$$\tilde{\ell}(\tilde{h}^*(\tilde{x}), \tilde{y}) = \mathbb{E}_{\kappa_{\tilde{y}}^X((\tilde{h}^*)^{-1}(\cdot), \tilde{x})}[(\kappa^Y \ell)_{\tilde{x}}(u, \tilde{y})].$$

877 where only the simple  $Y$  noise can be considered, given the missing result for the BR equality in the  
878 (D2, D2) joint corruption case.

879 As for the cases involving the 1D with 1D or 2D, i.e. points (3) and (4), we follow the same procedure  
880 by using the action formula of dependent corruptions as described in the proof of **Theorem 5**, and  
881 obtain the thesis.  $\square$

Table S1: Action of corruption on joint distribution and minimization sets, written using **P1,P2,P3**. We keep the continuous notation also for  $Y$  for ease of notation. Here we refer to the Bayes Risk results as the informal equivalence  $(\kappa(\tilde{\ell} \circ \tilde{H}), P) \equiv_{BR} (\tilde{\ell} \circ \tilde{H}, \kappa P)$ , and call  $\tilde{\ell} \circ H$  the minimization set associated to  $\kappa P$  (differently from the § S3 Theorems, where we used  $\ell \circ H$ ).

	Integral representation	Graphical model	Bayes Risk results
S- $\tilde{X}$ , S- $\tilde{Y}$	$\begin{aligned} \tilde{P} &= (\kappa_{Y\tilde{Y}} * \kappa_{X\tilde{X}}) \circ (\pi_Y \times E) \\ &= \sum_Y \kappa(d\tilde{y}, y) \pi_Y(dy) \int_X \kappa(d\tilde{x}, x) E(dx, y) \\ \tilde{P} &= (\kappa_{Y\tilde{Y}} * \kappa_{X\tilde{X}}) \circ (\pi_X \times F) \\ &= \int_X \kappa(d\tilde{x}, x) \pi_X(dx) \sum_Y \kappa(d\tilde{y}, y) F(dy, x) \end{aligned}$		$\begin{aligned} &(\kappa_{X\tilde{X}}(\kappa_{Y\tilde{Y}} \tilde{\ell} \circ \tilde{H}), P) \equiv_{BR} (\tilde{\ell} \circ \tilde{H}, \kappa_{Y\tilde{Y}} \pi_Y \times \kappa_{X\tilde{X}} E_{\tilde{Y}}) \\ &(\kappa_{X\tilde{X}}(\kappa_{Y\tilde{Y}} \tilde{\ell} \circ \tilde{H}), P) \equiv_{BR} (\tilde{\ell} \circ \tilde{H}, \kappa_{X\tilde{X}} \pi_X \times \kappa_{Y\tilde{Y}} F_{\tilde{X}}) \end{aligned}$
S- $\tilde{X}$ , 2D- $\tilde{Y}$	$\tilde{P} = (\kappa_{X\tilde{X}} * \kappa_{XY\tilde{Y}}) \circ (\pi_X \times F)$ $= \int_X \kappa(d\tilde{x}, x) \pi_X(dx) \sum_Y \kappa(d\tilde{y}, y) F(dy, x)$		$(\kappa_{X\tilde{X}}(\kappa_{XY\tilde{Y}} \tilde{\ell} \circ \tilde{H}), P) \equiv_{BR} (\tilde{\ell} \circ \tilde{H}, \kappa_{X\tilde{X}} \pi_X \times \kappa_{XY\tilde{Y}} F_{\tilde{X}})$
2D- $\tilde{X}$ , S- $\tilde{Y}$	$\tilde{P} = (\kappa_{Y\tilde{Y}} * \kappa_{XY\tilde{X}}) \circ (\pi_Y \times E)$ $= \sum_Y \kappa(d\tilde{y}, y) \pi_Y(dy) \int_X \kappa(d\tilde{x}, x, y) E(dx, y)$		$(\kappa_{XY\tilde{X}}(\kappa_{Y\tilde{Y}} \tilde{\ell} \circ \tilde{H}), P) \equiv_{BR} (\tilde{\ell} \circ \tilde{H}, \kappa_{Y\tilde{Y}} \pi_Y \times \kappa_{XY\tilde{X}} E_{\tilde{Y}})$
1D- $\tilde{X}$ , 1D- $\tilde{Y}$	$\tilde{P} = (\kappa_{X\tilde{Y}} * \kappa_{Y\tilde{X}}) \circ (\pi_Y \times E)$ $= \sum_Y \kappa(d\tilde{x}, y) \pi_Y(dy) \int_X \kappa(d\tilde{y}, x) E(dx, y)$		$(\kappa_{Y\tilde{X}}(\kappa_{X\tilde{Y}} \tilde{\ell} \circ \tilde{H}), P) \equiv_{BR} (\tilde{\ell} \circ \tilde{H}, \kappa_{Y\tilde{X}}(\pi_Y \times \kappa_{X\tilde{Y}} E))$
1D- $\tilde{X}$ , 2D- $\tilde{Y}$	$\tilde{P} = (\kappa_{X\tilde{Y}} * \kappa_{Y\tilde{X}}) \circ (\pi_X \times F)$ $= \int_X \kappa(d\tilde{y}, x) \pi_X(dx) \sum_Y \kappa(d\tilde{x}, y) F(dy, x)$		$(\kappa_{Y\tilde{X}}(\kappa_{X\tilde{Y}} \tilde{\ell} \circ \tilde{H}), P) \equiv_{BR} (\tilde{\ell} \circ \tilde{H}, \kappa_{X\tilde{Y}}(\pi_X \times \kappa_{Y\tilde{X}} F))$
2D- $\tilde{X}$ , 1D- $\tilde{Y}$	$\tilde{P} = (\kappa_{Y\tilde{X}} * \kappa_{XY\tilde{Y}}) \circ (\pi_Y \times E)$ $= \sum_Y \kappa(d\tilde{x}, y) \pi_Y(dy) \int_X \kappa(d\tilde{y}, x, y) E(dx, y)$		$(\kappa_{XY\tilde{Y}}(\kappa_{Y\tilde{X}} \tilde{\ell} \circ \tilde{H}), P) \equiv_{BR} (\tilde{\ell} \circ \tilde{H}, \kappa_{XY\tilde{Y}}(\pi_Y \times \kappa_{XY\tilde{Y}} E))$
2D- $\tilde{X}$ , 2D- $\tilde{Y}$	<p>No integral can be isolated, all priors and posteriors are affected by both factors.</p>		<p>No result using the factorization.</p>



884 **Supplementary References**

- 885 [30] George Shackelford and Dennis Volper. Learning k-dnf with noise in the attributes. In *Proceedings of the*  
886 *first annual workshop on Computational learning theory*, pages 97–103, 1988.
- 887 [31] Sally A. Goldman and Robert H. Sloan. Can pac learning algorithms tolerate random attribute noise?  
888 *Algorithmica*, 14(1):70–84, 1995.
- 889 [4] Xingquan Zhu and Xindong Wu. Class noise vs. attribute noise: A quantitative study. *The Artificial*  
890 *Intelligence Review*, 22(3):177, 2004.
- 891 [19] Robert C Williamson and Zac Cranko. Information processing equalities and the information-risk bridge.  
892 *arXiv preprint arXiv:2207.11987*, 2022.
- 893 [5] Donald B Rubin. Inference and missing data. *Biometrika*, 63(3):581–592, 1976.
- 894 [6] Roderick JA Little and Donald B Rubin. *Statistical analysis with missing data*, volume 793. John Wiley &  
895 Sons, 2019.
- 896 [5] Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning with noisy  
897 labels. *Advances in neural information processing systems*, 26, 2013.
- 898 [34] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making deep  
899 neural networks robust to label noise: A loss correction approach. In *Proceedings of the IEEE conference*  
900 *on computer vision and pattern recognition*, pages 1944–1952, 2017.
- 901 [7] Brendan Van Rooyen and Robert C Williamson. A theory of learning with corrupted labels. *J. Mach.*  
902 *Learn. Res.*, 18(1):8501–8550, 2017.
- 903 [32] Dana Angluin and Philip Laird. Learning from noisy examples. *Machine Learning*, 2:343–370, 1988.
- 904 [33] Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of*  
905 *the eleventh annual conference on Computational learning theory*, pages 92–100, 1998.
- 906 [12] Brendan Van Rooyen, Aditya Menon, and Robert C Williamson. Learning with symmetric label noise:  
907 The importance of being unhinged. *Advances in neural information processing systems*, 28, 2015.
- 908 [13] Takashi Ishida, Gang Niu, Weihua Hu, and Masashi Sugiyama. Learning from complementary labels.  
909 *Advances in neural information processing systems*, 30, 2017.
- 910 [14] Takashi Ishida, Gang Niu, Aditya Menon, and Masashi Sugiyama. Complementary-label learning for  
911 arbitrary losses and models. In *International Conference on Machine Learning*, pages 2971–2980. PMLR,  
912 2019.
- 913 [35] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. A neural algorithm of artistic style. *arXiv preprint*  
914 *arXiv:1508.06576*, 2015.
- 915 [36] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-  
916 resolution. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands,*  
917 *October 11–14, 2016, Proceedings, Part II 14*, pages 694–711. Springer, 2016.
- 918 [37] Eric Grinstead, Ngoc QK Duong, Alexey Ozerov, and Patrick Pérez. Audio style transfer. In *2018 IEEE*  
919 *international conference on acoustics, speech and signal processing (ICASSP)*, pages 586–590. IEEE,  
920 2018.
- 921 [18] Aritra Ghosh, Naresh Manwani, and PS Sastry. Making risk minimization tolerant to label noise. *Neuro-*  
922 *computing*, 160:93–107, 2015.
- 923 [8] Aditya Krishna Menon, Brendan Van Rooyen, and Nagarajan Natarajan. Learning from binary labels with  
924 instance-dependent noise. *Machine Learning*, 107(8):1561–1595, 2018.
- 925 [43] Jiacheng Cheng, Tongliang Liu, Kotagiri Ramamohanarao, and Dacheng Tao. Learning with bounded  
926 instance and label-dependent label noise. In *International Conference on Machine Learning*, pages  
927 1789–1799. PMLR, 2020.
- 928 [44] Yu Yao, Tongliang Liu, Mingming Gong, Bo Han, Gang Niu, and Kun Zhang. Instance-dependent label-  
929 noise learning under a structural causal model. *Advances in Neural Information Processing Systems*, 34:  
930 4409–4420, 2021.
- 931 [45] Qizhou Wang, Bo Han, Tongliang Liu, Gang Niu, Jian Yang, and Chen Gong. Tackling instance-dependent  
932 label noise via a universal probabilistic model. In *Proceedings of the AAAI Conference on Artificial*  
933 *Intelligence*, volume 35, pages 10183–10191, 2021.
- 934 [23] Jun Du and Zhihua Cai. Modelling class noise with symmetric and asymmetric distributions. In *Proceedings*  
935 *of the AAAI Conference on Artificial Intelligence*, volume 29, 2015.
- 936 [46] Nathalie Japkowicz and Shaju Stephen. The class imbalance problem: A systematic study. *Intelligent data*  
937 *analysis*, 6(5):429–449, 2002.

- 938 [47] Haibo He and Edwardo A Garcia. Learning from imbalanced data. *IEEE Transactions on knowledge and*  
939 *data engineering*, 21(9):1263–1284, 2009.
- 940 [48] Mateusz Buda, Atsuto Maki, and Maciej A Mazurowski. A systematic study of the class imbalance  
941 problem in convolutional neural networks. *Neural networks*, 106:249–259, 2018.
- 942 [49] Zachary Lipton, Yu-Xiang Wang, and Alexander Smola. Detecting and correcting for label shift with black  
943 box predictors. In *International conference on machine learning*, pages 3122–3130. PMLR, 2018.
- 944 [6] Aditya Menon, Brendan Van Rooyen, Cheng Soon Ong, and Bob Williamson. Learning from corrupted  
945 binary labels via class-probability estimation. In *International conference on machine learning*, pages  
946 125–134. PMLR, 2015.
- 947 [50] Gilles Blanchard, Marek Flaska, Gregory Handy, Sara Pozzi, and Clayton Scott. Classification with  
948 asymmetric label noise: Consistency and maximal denoising. *Electronic Journal of Statistics*, 10(2):  
949 2780–2824, 2016.
- 950 [51] Julian Katz-Samuels, Gilles Blanchard, and Clayton Scott. Decontamination of mutual contamination  
951 models. *Journal of machine learning research*, 20(41), 2019.
- 952 [3] Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood  
953 function. *Journal of statistical planning and inference*, 90(2):227–244, 2000.
- 954 [52] Arthur Gretton, Alex Smola, Jiayuan Huang, Marcel Schmittfull, Karsten Borgwardt, and Bernhard  
955 Schölkopf. Covariate shift by kernel mean matching. *Dataset shift in machine learning*, 3(4):5, 2009.
- 956 [53] Masashi Sugiyama and Motoaki Kawanabe. *Machine learning in non-stationary environments: Introduc-*  
957 *tion to covariate shift adaptation*. MIT press, 2012.
- 958 [54] Tianyi Zhang, Ikko Yamane, Nan Lu, and Masashi Sugiyama. A one-step approach to covariate shift  
959 adaptation. In *Asian Conference on Machine Learning*, pages 65–80. PMLR, 2020.
- 960 [55] Kun Zhang, Bernhard Schölkopf, Krikamol Muandet, and Zhikun Wang. Domain adaptation under target  
961 and conditional shift. In *International conference on machine learning*, pages 819–827. PMLR, 2013.
- 962 [56] Mingming Gong, Kun Zhang, Tongliang Liu, Dacheng Tao, Clark Glymour, and Bernhard Schölkopf.  
963 Domain adaptation with conditional transferable components. In *International conference on machine*  
964 *learning*, pages 2839–2848. PMLR, 2016.
- 965 [57] Xiyu Yu, Tongliang Liu, Mingming Gong, Kun Zhang, Kayhan Batmanghelich, and Dacheng Tao. Label-  
966 noise robust domain adaptation. In *International conference on machine learning*, pages 10913–10924.  
967 PMLR, 2020.
- 968 [17] Jose G Moreno-Torres, Troy Raeder, Rocío Alaiz-Rodríguez, Nitesh V Chawla, and Francisco Herrera. A  
969 unifying view on dataset shift in classification. *Pattern recognition*, 45(1):521–530, 2012.
- 970 [27] Kenta Cho and Bart Jacobs. Disintegration and bayesian inversion via string diagrams. *Mathematical*  
971 *Structures in Computer Science*, 29(7):938–971, 2019.
- 972 [26] Fredrik Dahlqvist, Vincent Danos, Ilias Garnier, and Ohad Kammar. Bayesian inversion by  $\omega$ -complete  
973 cone duality. In *27th International Conference on Concurrency Theory*, 2016.
- 974 [41] Saunders Mac Lane. *Categories for the working mathematician*, volume 5. Springer Science & Business  
975 Media, 2013.
- 976 [42] Erhan Çinlar. *Probability and Stochastics*. Springer, 2011.
- 977 [20] Erik Torgersen. *Comparison of statistical experiments*, volume 36. Cambridge University Press, 1991.