

## 568 A Proofs

569 **Proof of Prop. 4.1.** Substituting  $v_{KL}$  into the definition of  $\Delta_v(S, j, \mathbf{x})$  gives:

$$\Delta_{KL}(S, j, \mathbf{x}) = -D_{KL}(p_{Y|\mathbf{x}} \parallel p_{Y|\mathbf{x}_S, x_j}) + D_{KL}(p_{Y|\mathbf{x}} \parallel p_{Y|\mathbf{x}_S}).$$

570 Rearranging and using the definition of KL-divergence, we have:

$$\Delta_{KL}(S, j, \mathbf{x}) = \mathbb{E}_{Y|\mathbf{x}} [\log p(y | \mathbf{x}) - \log p(y | \mathbf{x}_S)] - \mathbb{E}_{Y|\mathbf{x}} [\log p(y | \mathbf{x}) - \log p(y | \mathbf{x}_S, x_j)].$$

571 Cleaning up in steps:

$$\begin{aligned} \Delta_{KL}(S, j, \mathbf{x}) &= \mathbb{E}_{Y|\mathbf{x}} [\log p(y | \mathbf{x}) - \log p(y | \mathbf{x}_S) - \log p(y | \mathbf{x}) + \log p(y | \mathbf{x}_S, x_j)] \\ &= \mathbb{E}_{Y|\mathbf{x}} [\log p(y | \mathbf{x}_S, x_j) - \log p(y | \mathbf{x}_S)] \\ &= \int_{\mathcal{Y}} p(y | \mathbf{x}) \log \frac{p(y | \mathbf{x}_S, x_j)}{p(y | \mathbf{x}_S)} dy. \end{aligned}$$

572 Substituting  $v_{CE}$  into the definition of  $\Delta_v(S, j, \mathbf{x})$  gives:

$$\Delta_{CE}(S, j, \mathbf{x}) = -H(p_{Y|\mathbf{x}}, p_{Y|\mathbf{x}_S, x_j}) + H(p_{Y|\mathbf{x}}, p_{Y|\mathbf{x}_S}).$$

573 Rearranging and using the definition of cross entropy, we have:

$$\begin{aligned} \Delta_{CE}(S, j, \mathbf{x}) &= H(p_{Y|\mathbf{x}}, p_{Y|\mathbf{x}_S}) - H(p_{Y|\mathbf{x}}, p_{Y|\mathbf{x}_S \cup \{j\}}) \\ &= \mathbb{E}_{Y|\mathbf{x}} [-\log p(y | \mathbf{x}_S)] - \mathbb{E}_{Y|\mathbf{x}} [-\log p(y | \mathbf{x}_S, x_j)] \\ &= \mathbb{E}_{Y|\mathbf{x}} [\log p(y | \mathbf{x}_S, x_j) - \log p(y | \mathbf{x}_S)] \\ &= \int_{\mathcal{Y}} p(y | \mathbf{x}) \log \frac{p(y | \mathbf{x}_S, x_j)}{p(y | \mathbf{x}_S)} dy. \end{aligned}$$

574 **Proof of Prop. 4.2.** Since the Shapley value  $\phi_v(j, \mathbf{x})$  is just the expectation of  $\Delta_v(S, j, \mathbf{x})$  under a certain  
575 distribution on coalitions  $S \subseteq [d] \setminus \{j\}$  (see Eq. 1), it follows from Prop. 4.1 that feature attributions will  
576 be identical under  $v_{KL}$  and  $v_{CE}$ . To show that resulting Shapley values sum to the KL-divergence between  
577  $p(Y | \mathbf{x})$  and  $p(Y)$ , we exploit the efficiency property:

$$\begin{aligned} \sum_{j=1}^d \phi_{KL}(j, \mathbf{x}) &= v_{KL}([d], \mathbf{x}) - v_{KL}(\emptyset, \mathbf{x}) \\ &= -D_{KL}(p_{Y|\mathbf{x}} \parallel p_{Y|\mathbf{x}}) + D_{KL}(p_{Y|\mathbf{x}} \parallel p_Y) \\ &= D_{KL}(p_{Y|\mathbf{x}} \parallel p_Y). \end{aligned}$$

578 The last step exploits Gibbs's inequality, according to which  $D_{KL}(p \parallel q) \geq 0$ , with  $D_{KL}(p \parallel q) = 0$  iff  $p = q$ .

579 **Proof of Prop. 4.3.** Substituting  $v_{IG}$  into the definition of  $\Delta_v(S, j, \mathbf{x})$  gives:

$$\begin{aligned} \Delta_{IG}(S, j, \mathbf{x}) &= -H(Y | \mathbf{x}_S, x_j) + H(Y | \mathbf{x}_S) \\ &= H(Y | \mathbf{x}_S) - H(Y | \mathbf{x}_S, x_j) \\ &= I(Y; x_j | \mathbf{x}_S) \\ &= \int_{\mathcal{Y}} p(y, x_j | \mathbf{x}_S) \log \frac{p(y, x_j | \mathbf{x}_S)}{p(y | \mathbf{x}_S) p(x_j | \mathbf{x}_S)} dy. \end{aligned}$$

580 In the penultimate line, we exploit the equality  $I(Y; X) = H(Y) - H(Y | X)$ , by which we define mutual  
581 information (see Appx. B.1).

582 **Proof of Prop. 4.4.** We once again rely on efficiency and the definition of mutual information in terms of  
583 marginal and conditional entropy:

$$\begin{aligned} \sum_{j=1}^d \phi_{IG}(j, \mathbf{x}) &= v_{IG}([d], \mathbf{x}) - v_{IG}(\emptyset, \mathbf{x}) \\ &= -H(Y | \mathbf{x}) + H(Y) \\ &= H(Y) - H(Y | \mathbf{x}) \\ &= I(Y; \mathbf{x}). \end{aligned}$$

584 **Proof of Thm. 4.5.** Begin with item (a). Note that the conditional independence statement  $Y \perp\!\!\!\perp X_j \mid \mathbf{X}_S$   
585 holds iff, for all points  $(\mathbf{x}, y) \sim \mathcal{D}$ , we have:

$$p(y \mid \mathbf{x}_S, x_j) = p(y \mid \mathbf{x}_S) \quad \text{and} \quad p(y, x_j \mid \mathbf{x}_S) = p(y \mid \mathbf{x}_S) p(x_j \mid \mathbf{x}_S).$$

586 The former guarantees that marginal payouts evaluate to zero for  $v \in \{v_{KL}, v_{CE}\}$ ; the latter does the same  
587 for  $v \in \{v_{IG}, v_H\}$ . This follows because the log ratio in each formula evaluates to zero when numerator and  
588 denominator are equal.

589 Of course, conditional independence is also sufficient for zero marginal payout with more familiar value  
590 functions such as  $v_0$ . But item (a) makes an additional claim—that the *converse* holds as well, i.e. that  
591 conditional independence is *necessary* for zero marginal payout across all  $\mathbf{x}$ . This follows from the definitions  
592 of the value functions themselves. Observe:

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} [\Delta_{KL}(S, j, \mathbf{x})] &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[ \log \frac{p(y \mid \mathbf{x}_S, x_j)}{p(y \mid \mathbf{x}_S)} \right] \\ &= \mathbb{E}_{\mathcal{D}_X} \left[ \mathbb{E}_{Y \mid \mathbf{x}_S, x_j} \left[ \log \frac{p(y \mid \mathbf{x}_S, x_j)}{p(y \mid \mathbf{x}_S)} \right] \right] \\ &= \mathbb{E}_{\mathcal{D}_X} [D_{KL}(p_{Y \mid \mathbf{x}_S, x_j} \parallel p_{Y \mid \mathbf{x}_S})] \end{aligned}$$

593 By Gibbs's inequality, the KL-divergence between two distributions is zero iff they are equal, so setting this  
594 value to zero for all  $\mathbf{x}$  satisfies the first definition of conditional independence above. For the latter, we simply  
595 point out that:

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}_X} [\Delta_{IG}(S, j, \mathbf{x})] = I(Y; X_j \mid \mathbf{X}_S).$$

596 Since conditional mutual information equals zero iff the relevant variables are conditionally independent, this  
597 satisfies the second definition above.

598 Item (b) states that CSI, which is strictly weaker than standard conditional independence, is also sufficient for  
599 zero marginal payout at a given point  $\mathbf{x}$ . This follows directly from the sufficiency argument above.

600 The converse relationship is more complex, however. Call a distribution *conspiratorial* if there exists some  
601  $S, j, \mathbf{x}$  such that  $\Delta_v(S, j, \mathbf{x}) = 0 \wedge Y \not\perp\!\!\!\perp x_j \mid \mathbf{x}_S$  for some  $v \in \{v_{KL}, v_{CE}, v_{IG}, v_H\}$ . Such distributions are  
602 so named because the relevant probabilities must coordinate in a very specific way to guarantee summation to  
603 zero as we marginalize over  $\mathcal{Y}$ . As a concrete example, consider the following data generating process:

$$X \sim \text{Bern}(0.5), \quad Z \sim \text{Bern}(0.5), \quad Y \sim \text{Bern}(0.3 + 0.4X - 0.2Z).$$

604 What is the contribution of  $X$  to coalition  $S = \emptyset$  when  $X = 1$  and  $Z = 1$ ? In this case, we have neither global  
605 nor context-specific independence, i.e.  $Y \not\perp\!\!\!\perp x$ . Yet, evaluating the payoffs in a KL-divergence game, we have:

$$\begin{aligned} \Delta_{KL}(S, j, \mathbf{x}) &= \sum_y P(y \mid X = 1, Z = 1) \log \frac{P(y \mid X = 1)}{P(y)} \\ &= 0.5 \log \frac{0.4}{0.6} + 0.5 \log \frac{0.6}{0.4} \\ &= 0. \end{aligned}$$

606 In this case, we find that negative and positive values of the log ratio cancel out exactly as we marginalize over  
607  $\mathcal{Y}$ . (Similar examples can be constructed for  $v_{IG}$  and  $v_H$ .) This shows that CSI is sufficient but not necessary  
608 for  $\Delta_v(S, j, \mathbf{x}) = 0$ .

609 However, just because conspiratorial distributions are possible does not mean that they are common. Item (c)  
610 states that the set of all such distributions has Lebesgue measure zero. Our proof strategy here follows that of  
611 Meek [44], who demonstrates a similar result in the case of *unfaithful* distributions, i.e. those whose (conditional)  
612 independencies are not entailed by the data's underlying graphical structure. This is an important topic in the  
613 causal discovery literature (see, e.g., [80, 81]).

614 For simplicity, assume a discrete state space  $\mathcal{X} \times \mathcal{Y}$ . Fix some  $S, j$  such that  $Y \not\perp\!\!\!\perp x_j \mid \mathbf{x}_S$ . Let  $C$  be the number  
615 of possible outcomes,  $\mathcal{Y} = \{y_1, \dots, y_C\}$ . Define vectors  $\mathbf{p}, \mathbf{r}$  of length  $C$  such that, for each  $c \in [C]$ :

$$p_c = p(y_c \mid \mathbf{x}), \quad r_c = \log \frac{p(y_c \mid \mathbf{x}_S, x_j)}{p(y_c \mid \mathbf{x}_S)}.$$

616 (Technically, we only require  $C - 1$  entries to fully describe these conditional distributions, but there is no penalty  
617 for overparametrization here.) By the assumption of local conditional dependence, we know that  $\|\mathbf{r}\|_0 > 0$ . Yet  
618 for our conspiracy to obtain, the inner product of these vectors must satisfy  $\mathbf{p} \cdot \mathbf{r} = 0$ . A well-known algebraic  
619 lemma of Okamoto [48] states that if a polynomial constraint is non-trivial (i.e., if there exists some  $\mathbf{p}, \mathbf{r}$  for  
620 which it does not hold), then the subset of parameters for which it does hold has Lebesgue measure zero. Since  
621 the conspiracy requires nontrivial constraints that are linear in the parameters  $\mathbf{p}, \mathbf{r}$ , we conclude that the set of  
622 conspiratorial distributions has Lebesgue measure zero.

623 **Proof of Thm. 5.1.** Our proof is an application of the split conformal method (see [38] Thm. 2.2]). Whereas  
624 that method was designed to bound the distance between predicted and observed outcomes for a regression  
625 task, we effectively treat the mean Shapley value as a constant outcome to measure the concentration of feature  
626 attributions. To achieve this, we replace out-of-sample absolute residuals with out-of-sample Shapley values and  
627 labels with the mean Shapley value. With these substitutions in place, the result follows immediately from the  
628 symmetry of  $\phi(j, \mathbf{x}^{(i+1)})$  and  $\phi(j, \mathbf{x}^{(i)})$ ,  $i \in \mathcal{I}_2$ , which is itself a direct implication of the i.i.d. assumption<sup>1</sup>  
629 Since the margin is calculated so as to cover  $(1 - \alpha) \times 100\%$  of the distribution, it is unlikely that new samples  
630 will fall outside this region. Specifically, such exceptions occur with probability at most  $\alpha$ . This amounts to a  
631 sort of PAC guarantee, i.e. that Shapley values will be within radius  $\tau_j$  of their mean  $\mu_j$  with probability at least  
632  $1 - \alpha$ .

## 633 B Addenda

634 This section includes extra background material on information theory and Shapley values.

### 635 B.1 Information Theory

636 Let  $p, q$  be two probability distributions over the same  $\sigma$ -algebra of events. (In the continuous case, we  
637 additionally require that  $p, q$  be absolutely continuous with respect to Lebesgue measure.) The *entropy* of  
638  $p$  is defined as  $H(p) := \mathbb{E}_p[-\log p]$ , i.e. the expected number of bits required to encode the distribution<sup>2</sup>  
639 The *cross entropy* of  $p$  and  $q$  is defined as  $H(p, q) := \mathbb{E}_p[-\log q]$ , i.e. the expected number of bits required  
640 to encode samples from  $p$  using code optimized for  $q$ . The *KL-divergence* between  $p$  and  $q$  is defined as  
641  $D_{KL}(p \parallel q) := \mathbb{E}_p[\log p/q]$ , i.e. the cost in bits of modeling  $p$  with  $q$ . These three quantities are related by the  
642 formula  $D_{KL}(p \parallel q) = H(p, q) - H(p)$ . The reduction in  $Y$ 's uncertainty attributable to  $X$  is also called the  
643 *mutual information*,  $I(Y; X) := H(Y) - H(Y | X)$ . This quantity is nonnegative, with  $I(Y; X) = 0$  if and  
644 only if the variables are independent.

645 However, conditioning on a specific value of  $X$  may increase uncertainty in  $Y$ , in which case the local conditional  
646 entropy exceeds the marginal. Thus it is possible that  $H(Y | x) > H(Y)$  for some  $x \in \mathcal{X}$ . For example,  
647 consider the following data generating process:

$$X \sim \text{Bern}(0.8), \quad Y \sim \text{Bern}(0.5 + 0.25X).$$

648 In this case, we have  $P(Y = 1) = 0.7$ ,  $P(Y = 1 | X = 0) = 0.5$ , and  $P(Y = 1 | X = 1) = 0.75$ . It is  
649 easy to see that even though the marginal entropy  $H(Y)$  exceeds the global conditional entropy  $H(Y | X)$ , the  
650 local entropy at  $X = 0$  is larger than either quantity,  $H(Y | X = 0) > H(Y) > H(Y | X)$ . In other words,  
651 conditioning on the event  $X = 0$  increases our uncertainty about  $Y$ .

652 Similarly, there may be cases in which  $I(Y; X | Z) > 0$ , but  $I(Y; X | z) = 0$ . This is what Boutilier et al.  
653 [7] call *context-specific independence* (CSI). For instance, if  $X, Z \in \{0, 1\}^2$  and  $Y := X \vee Z$ , then we have  
654  $Y \perp\!\!\!\perp X | Z$ , but  $Y \not\perp\!\!\!\perp X | (Z = 1)$  since  $Y$ 's value is determined as soon as we know that either parent is 1.

### 655 B.2 The Shapley Axioms

656 For completeness, we here list the Shapley axioms.

657 **Efficiency.** Shapley values sum to the difference in payoff between complete and null coalitions:

$$\sum_{j=1}^d \phi(j, \mathbf{x}) = v([d], \mathbf{x}) - v(\emptyset, \mathbf{x}).$$

658 **Symmetry.** If two players make identical contributions to all coalitions, then their Shapley values are equal:

$$\forall S \subseteq [d] \setminus \{i, j\} : v(S \cup \{i\}, \mathbf{x}) = v(S \cup \{j\}, \mathbf{x}) \Rightarrow \phi(i, \mathbf{x}) = \phi(j, \mathbf{x}).$$

659 **Sensitivity.** If a player makes zero contribution to all coalitions, then its Shapley value is zero:

$$\forall S \subseteq [d] \setminus \{j\} : v(S \cup \{j\}, \mathbf{x}) = v(S, \mathbf{x}) \Rightarrow \phi(j, \mathbf{x}) = 0.$$

<sup>1</sup>Note that conformal inference relies on the weaker assumption of exchangeability. However, since we operate in the standard i.i.d. setting of statistical learning theory (see Sect. 3), exchangeability naturally follows.

<sup>2</sup>Though the term ‘‘bit’’ is technically reserved for units of information measured with logarithmic base 2, we use the word somewhat more loosely to refer to any unit of information.

660 **Linearity.** The Shapley value for a convex combination of games can be decomposed into a convex combina-  
 661 tion of Shapley values. For any  $a, b \in \mathbb{R}$  and value functions  $v_1, v_2$ , we have:

$$\phi_{a \cdot v_1 + b \cdot v_2}(j, \mathbf{x}) = a\phi_{v_1}(j, \mathbf{x}) + b\phi_{v_2}(j, \mathbf{x}).$$

## 662 C Experiments

### 663 C.1 Datasets.

664 The MNIST dataset is available online<sup>3</sup>. The IMDB dataset is available on Kaggle<sup>4</sup>. The BreastCancer,  
 665 Diabetes, Ionosphere, and Sonar datasets are all distributed in the mlbench package, which is available on  
 666 CRAN<sup>5</sup>.

### 667 C.2 Models.

668 All neural network training was conducted in PyTorch<sup>6</sup>. We use a standard convolutional neural network  
 669 for the MNIST experiment, including convolutions, max pooling, and batch norm. We use ReLU activations,  
 670 cross entropy loss, and optimize with Adam<sup>7</sup>. For the IMDB experiment, we use a pre-trained BERT model  
 671 from the Hugging Face transformers library<sup>8</sup>. All hyperparameters are set to their default values. All XGBoost  
 672 models are trained with the default hyperparameters, with the number of training rounds cited in the text.

### 673 C.3 Coverage

674 To empirically test our conformal coverage guarantee, we compute means and margins on out-of-sample Shapley  
 675 values for the modified Friedman benchmark. Results for conditional expectation and conditional variance are  
 676 reported in Table 1 with target level  $\alpha = 0.1$ . Note that what constitutes a “small” or “large” margin is context  
 677 dependent. The conditional variance model is fit to  $\epsilon_y^2$ , which has a tighter range than  $Z$ , leading to smaller  
 678 Shapley values on average. However, nominal coverage is very close to the target 90% throughout, illustrating  
 679 how the conformal method can be used for feature selection and outlier detection.

Table 1: Means, margins, and nominal coverage at  $\alpha = 0.1$  for Shapley values from the conditional mean and conditional variance models. Results are averaged over 50 replicates.

Feature	Mean			Variance		
	$\mu$	$\tau$	Coverage	$\mu$	$r$	Coverage
$X_1$	-0.002	0.066	0.899	-0.009	0.505	0.898
$X_2$	0.008	0.141	0.898	-0.001	0.435	0.900
$X_3$	0.002	0.084	0.899	0.001	0.278	0.898
$X_4$	-0.004	0.098	0.901	-0.006	0.727	0.900
$X_5$	-0.004	0.092	0.905	0.020	0.333	0.902
$X_6$	-0.162	3.637	0.903	-0.001	0.060	0.900
$X_7$	-0.032	3.555	0.901	0.003	0.049	0.899
$X_8$	-0.027	1.981	0.898	0.001	0.055	0.900
$X_9$	0.190	4.114	0.898	-0.002	0.053	0.899
$X_{10}$	-0.044	1.952	0.903	-0.001	0.053	0.900

<sup>3</sup><http://yann.lecun.com/exdb/mnist/>.

<sup>4</sup><https://www.kaggle.com/datasets/lakshmi25npathi/imdb-dataset-of-50k-movie-reviews>.

<sup>5</sup><https://cran.r-project.org/web/packages/mlbench/index.html>.

<sup>6</sup>[https://huggingface.co/docs/transformers/model\\_doc/bert](https://huggingface.co/docs/transformers/model_doc/bert).