

What to Expect of Classifiers? Reasoning about Logistic Regression with Missing Features

Pasha Khosravi, Yitao Liang, YooJung Choi and Guy Van den Broeck

University of California, Los Angeles

{pashak, yliang, yjchoi, guyvdb}@cs.ucla.edu

Abstract

While discriminative classifiers often yield strong predictive performance, missing feature values at prediction time can still be a challenge. Classifiers may not behave as expected under certain ways of substituting the missing values, since they inherently make assumptions about the data distribution they were trained on. In this paper, we propose a novel framework that classifies examples with missing features by computing the expected prediction with respect to a feature distribution. Moreover, we use geometric programming to learn a naive Bayes distribution that embeds a given logistic regression classifier and can efficiently take its expected predictions. Empirical evaluations show that our model achieves the same performance as the logistic regression with all features observed, and outperforms standard imputation techniques when features go missing during prediction time. Furthermore, we demonstrate that our method can be used to generate “sufficient explanations” of logistic regression classifications, by removing features that do not affect the classification.

1 Introduction

Missing values are pervasive in real-world machine learning applications. Learned classifiers usually assume all input features are known, but when a classifier is deployed, some features may not be available, or too difficult to acquire. This can be due to the noisy nature of our environment, unreliable or costly sensors, and numerous other challenges in gathering and managing data [Dekel and Shamir, 2008; Graham, 2012]. Consider autonomous driving for example, where blocked sensors may leave observations incomplete.

Moreover, it can be difficult to anticipate the many ways in which a learned classifier will be deployed. For example, the same classifier could be used by different doctors who choose to run different medical tests on their patients, and this information is not known at learning time. Nevertheless, even a small portion of missing data can severely affect the performance of well-trained models, leading to predictions that are very different from the ones made when all data is observed.

Certain machine learning models, such as probabilistic graphical models, provide a natural solution to missing features at prediction time, by formulating the problem as a probabilistic inference task [Koller and Friedman, 2009; Darwiche, 2009]. Unfortunately, with the increasing emphasis on predictive performance, the general consensus is that such generative models are not competitive as classifiers given fully-observed feature vectors, and that discriminatively-trained models are to be preferred [Ng and Jordan, 2002].

To alleviate the impact of missing data on discriminative classifiers, it is common practice to substitute the missing values with plausible ones [Schafer, 1999; Little and Rubin, 2014]. As we will argue later, a drawback for such imputation techniques is that they can make overly strong assumptions about the feature distribution. Furthermore, to be compatible with many types of classifiers, they tend to overlook how the multitude of possible imputed values would interact with the classifier at hand, and risk yielding biased predictions.

To better address this issue, we propose a principled framework of handling missing features by reasoning about the classifier’s *expected* output given the feature distribution. One obvious advantage is that it can be tailored to the given family of classifiers and feature distributions. In contrast, the popular mean imputation approach only coincides with the expected prediction for simple models (e.g. linear functions) under very strong independence assumptions about the feature distribution. We later show that calculating the expected predictions with respect to arbitrary feature distributions is computationally highly intractable. In order to make our framework more feasible in practice, we leverage generative-discriminative counterpart relationships to learn a joint distribution that can take expectations of its corresponding discriminative classifier. We call this problem *conformant learning*. Then, we develop an algorithm, based on geometric programming, for a well-known example of such relationship: naive Bayes and logistic regression. We call this specific algorithm naive conformant learning (NaCL).

Through an extensive empirical evaluation over five characteristically distinct datasets, we show that NaCL consistently achieves better estimates of the conditional probability, as measured by average cross entropy and classification accuracy, compared to commonly used imputation methods. Lastly, we conduct a short case study on how our framework can be applied to *explain* classifications.

2 The Expected Prediction Task

In this section, we describe our intuitive approach to making a prediction when features are missing and discuss how it relates to existing imputation methods. Then we study the computational hardness of our expected prediction task.

We use uppercase letters to denote features/variables and lowercase letters for their assignment. Sets of variables \mathbf{X} and their joint assignments \mathbf{x} are written in bold. For an assignment x to a binary variable X , we let \bar{x} denote its negation. Concatenation \mathbf{XY} denotes the union of disjoint sets. The set of all possible assignments to \mathbf{X} is denoted \mathcal{X} .

Suppose we have a model trained with features \mathbf{X} but are now faced with the challenge of making a prediction without knowing all values \mathbf{x} . In this situation, a common solution is to impute certain substitute values for the missing data (for example their mean) [Little and Rubin, 2014]. However, the features that were observed provide information not only about the class but also about the missing features, yet this information is typically not taken into account by popular methods such as mean imputation.

We propose a very natural alternative: to utilize the feature distribution to probabilistically reason about what a predictor is expected to return if it could observe the missing features.

Definition 1. Let $\mathcal{F} : \mathcal{X} \rightarrow \mathbb{R}$ be a predictor and P be a distribution over features \mathbf{X} . Given a partitioning of features $\mathbf{X} = \mathbf{Y}\mathbf{M}$ and an assignment \mathbf{y} to some of the features \mathbf{Y} , the *expected prediction task* is to compute

$$E_{\mathcal{F},P}(\mathbf{y}) = \mathbb{E}_{\mathbf{m} \sim P(\mathbf{M}|\mathbf{y})} [\mathcal{F}(\mathbf{y}\mathbf{m})].$$

The expected prediction task and (mean) imputation are related, but only under very restrictive assumptions.

Example 1. Let $\mathcal{F} : \mathcal{X} \rightarrow \mathbb{R}$ be a linear function. That is, $\mathcal{F}(\mathbf{x}) = \sum_{x \in \mathbf{x}} w_X x$ for some weights \mathbf{w} . Suppose P is a distribution over \mathbf{X} that assumes independence between features: $P(\mathbf{X}) = \prod_{X \in \mathbf{X}} P(X)$. Then, using linearity of expectation, the following holds for any partial observation \mathbf{y} :

$$\begin{aligned} E_{\mathcal{F},P}(\mathbf{y}) &= \mathbb{E}_{\mathbf{m} \sim P(\mathbf{M}|\mathbf{y})} \left[\sum_{y \in \mathbf{y}} w_Y y + \sum_{m \in \mathbf{m}} w_M m \right] \\ &= \sum_{y \in \mathbf{y}} w_Y y + \sum_{M \in \mathbf{M}} w_M \mathbb{E}_{m \sim P(M)} [m]. \end{aligned}$$

Hence, substituting the missing features with their means effectively computes the expected predictions of linear models if the independence assumption holds. Furthermore, if \mathcal{F} is the true conditional probability of the labels and features are generated by a fully-factorized $P(\mathbf{X})$, then classifying by comparing the expected prediction $E_{\mathcal{F},P}$ to 0.5 is Bayes optimal on the observed features. That is, an expected prediction higher than 0.5 means that the positive class is the most likely one given the observation, and thus minimizes expected loss, according to the distribution defined by \mathcal{F} and P .

Example 2. Consider a logistic regression model $\mathcal{G}(\mathbf{x}) = \text{sigmoid}(\mathcal{F}(\mathbf{x}))$ where \mathcal{F} is a linear function. Now, mean imputation no longer computes the expected prediction, even when the independence assumption in the previous example

holds. In particular, if \mathbf{y} is a partial observation such that $\mathcal{G}(\mathbf{y}\mathbf{m})$ is positive for all \mathbf{m} , then the mean-imputed prediction is an over-approximation of the expected prediction:

$$\begin{aligned} \mathcal{G}(\mathbf{y} \mathbb{E}[\mathbf{m}]) &= \text{sigmoid}(\mathbb{E}[\mathcal{F}(\mathbf{y}\mathbf{m})]) \\ &> \mathbb{E}[\text{sigmoid}(\mathcal{F}(\mathbf{y}\mathbf{m}))] = E_{\mathcal{G},P}(\mathbf{y}). \end{aligned}$$

This is due to Jensen’s inequality and concavity of the sigmoid function in the positive portion; conversely, it is an under-approximation in the negative cases.

Example 1 showed how to efficiently take the expectation of a linear function w.r.t. a fully factorized distribution. Unfortunately, the expected prediction task is in general computationally hard, even on simple classifiers and distributions.

Proposition 1. *Taking expectation of a nontrivial classifier w.r.t. a uniform distribution is #P-hard.*

Proof. Suppose our classifier tests whether a logical constraint holds between the input features. Then asking whether there exists a positive example is equivalent to SAT which is NP-hard. The expected classification on a uniform distribution is solving an even harder task, of counting solutions to the constraint, which is #P-hard [Roth, 1996]. \square

Next, consider the converse in which the classifier is trivial but the distribution is more general.

Proposition 2. *The expectation of a classifier that returns the value of a single feature w.r.t. a distribution represented by a probabilistic graphical model is #P-hard.*

Proof. Computing expectations of such classifier is as hard as computing marginals in the feature distribution, which is #P-hard for graphical models [Roth, 1996]. \square

Previous propositions showed that the expected prediction task stays intractable, even when we allow either the distribution or the classifier to be trivial.

Our next theorem states that the task is hard even for a relatively simple classifier and a tractable distribution.¹

Theorem 1. *Computing the expectation of a logistic regression classifier over a naive Bayes distribution is NP-hard.*

That is, the expected prediction task is hard even though logistic regression classification and probabilistic reasoning on naive Bayes models can both be done in linear time.

In summary, while the expected prediction task appears natural for dealing with missing data, its vast intractability poses a serious challenge, especially compared to efficient alternatives such as imputation. Next, we investigate specific ways of practically overcoming this challenge.

3 Joint Distributions as Classifiers

As shown previously, taking expectations is intractable for arbitrary pairs of classifiers and distributions. In this section, we propose *conformant learning* which aims to learn a joint distribution that encodes a given classifier as well as a feature distribution. On such distribution, the expected prediction task is well-defined as probabilistic inference, and is

¹All proofs can be found in Appendix A.

tractable for a large class of probabilistic models, including naive Bayes [Darwiche, 2009; Dechter, 2013].

We first describe how a joint distribution can be used as a classifier that inherently support missing features during prediction time. Given a distribution $P(\mathbf{X}, C)$ over the features \mathbf{X} and class variable C , we can classify a partial observation \mathbf{y} simply by computing the conditional probability $P(c|\mathbf{y})$ where c denotes the positive class.² In some sense, a joint distribution embeds a classifier $P(C|\mathbf{Y})$ for each subset of observed features \mathbf{Y} . In fact, computing $P(c|\mathbf{y})$ is *equivalent to computing the expected prediction* of classifier \mathcal{F} that outputs $P(c|\mathbf{x})$ for every \mathbf{x} , with respect to distribution $P(\mathbf{X})$:

$$\begin{aligned} P(c|\mathbf{y}) &= \sum_{\mathbf{m}} P(c, \mathbf{m} | \mathbf{y}) = \sum_{\mathbf{m}} P(c|\mathbf{m}\mathbf{y})P(\mathbf{m} | \mathbf{y}) \\ &= \mathbb{E}_{\mathbf{m} \sim P(\mathbf{M}|\mathbf{y})} [P(c|\mathbf{m}\mathbf{y})] = E_{\mathcal{F}, P}(\mathbf{y}). \end{aligned} \quad (1)$$

Nevertheless, the prevailing consensus is that in practice discriminatively training a classifier $P(C|\mathbf{X})$ should be preferred to generatively learning $P(\mathbf{X}, C)$, because it tends to achieve higher classification accuracy [Bouchard and Triggs, 2004; Ulusoy and Bishop, 2005].

There are many generative-discriminative pairs obtained from fitting the same family of probabilistic models to optimize either the joint or conditional likelihood [Jaakkola and Haussler, 1999; Sutton and McCallum, 2012; Liang and Van den Broeck, 2019], including naive Bayes and logistic regression [Ng and Jordan, 2002]. We formally describe such relationship as follows:

Definition 2. We say $P(\mathbf{X}, C)$ *conforms* with $\mathcal{F} : \mathcal{X} \rightarrow [0, 1]$ if their classifications agree: $P(c|\mathbf{x}) = \mathcal{F}(\mathbf{x})$ for all \mathbf{x} .

Next, let us study naive Bayes models in detail as they support efficient inference and thus are a good target distribution to leverage for the expected prediction task. Naive Bayes models assume that features are mutually independent given the class; that is, its joint distribution is $P(\mathbf{X}, C) = P(C) \prod_{X \in \mathbf{X}} P(X|C)$. Under such assumption, marginal inference takes linear time, and so does computing expectations under missing features as in Equation 1.

Logistic regression is the discriminative counterpart to naive Bayes. It has parameters \mathbf{w} and posits that³

$$\mathcal{F}(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}.$$

Any naive Bayes classifier can be translated to an equivalent logistic regression classifier on fully observed features.

Lemma 1. *Given a naive Bayes distribution P , there is a unique logistic regression model \mathcal{F} that it conforms with. Such logistic regression model has the following weights:*

$$\begin{aligned} w_0 &= \log \frac{P(c)}{P(\bar{c})} + \sum_{i=1}^n \log \frac{P(\bar{x}_i | c)}{P(\bar{x}_i | \bar{c})} \\ w_i &= \log \frac{P(x_i | c)}{P(x_i | \bar{c})} \cdot \frac{P(\bar{x}_i | \bar{c})}{P(\bar{x}_i | c)}, \quad i = 1, \dots, n \end{aligned}$$

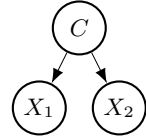
²We assume binary class for conciseness, but our approach easily generalizes to multiclass. Details can be found in Appendix B.

³Here, \mathbf{x} also includes a dummy feature that is always 1 to correspond with the bias parameter w_0 .

$$\mathbf{w} = \begin{bmatrix} -1.16 \\ 2.23 \\ -0.20 \end{bmatrix}$$

X_1	X_2	$\mathcal{F}(x_1, x_2)$
1	1	0.70
1	0	0.74
0	1	0.20
0	0	0.24

(a) Logistic regression weights and resulting predictions.



$P_1(c)$	C	$P_1(x_1 C)$	C	$P_1(x_2 C)$
0.5	1	0.8	1	0.45
	0	0.3	0	0.5
$P_2(c)$	C	$P_2(x_1 C)$	C	$P_2(x_2 C)$
0.36	1	0.6	1	0.9
	0	0.14	0	0.92

(b) Two naive Bayes distributions with same structure.

Figure 1: Logistic regression $\mathcal{F}(\mathbf{x}) = \text{sigmoid}(\mathbf{w}^T \mathbf{x})$ and two conformant naive Bayes models.

Here, x, \bar{x} denote $X=1, X=0$ respectively.

The lemma above can be drawn from Roos *et al.* [2005] and Ng and Jordan [2002]. Complete proofs of all lemmas can be found in Appendix A.

Consider for example the naive Bayes (NB) distribution P_1 in Figure 1b. For all possible feature observations, the NB classification $P_1(c|\mathbf{x})$ is equal to that of logistic regression (LR) \mathcal{F} in Figure 1a, whose weights are as given by above lemma (i.e., P_1 conforms with \mathcal{F}). Furthermore, distribution P_2 also translates into the same logistic regression. In fact, there can be infinitely many such naive Bayes distributions.

Lemma 2. *Given a logistic regression \mathcal{F} and $\theta \in (0, 1)^n$, there exists a unique naive Bayes model P such that*

$$\begin{aligned} P(c|\mathbf{x}) &= \mathcal{F}(\mathbf{x}), \quad \forall \mathbf{x} \\ P(x_i | c) &= \theta_i, \quad i = 1, \dots, n. \end{aligned}$$

That is, given a logistic regression there are infinitely many naive Bayes models that conform with it. Moreover, after fixing values for n parameters of the NB model, there is a uniquely corresponding naive Bayes model.

We can expect this phenomenon to generalize to other generative-discriminative pairs; given a conditional distribution $P(C|\mathbf{X})$ there are many possible feature distributions $P(\mathbf{X})$ to define a joint distribution $P(\mathbf{X}, C)$. For instance, distributions P_1 and P_2 in Figure 1b assign different probabilities to feature observations; $P_1(\bar{x}_1, \bar{x}_2) = 0.23$ whereas $P_2(\bar{x}_1, \bar{x}_2) = 0.06$. Hence, we wish to define which one of these models is the “best”. Naturally, we choose the one that best explains a given dataset of feature observations.⁴

Definition 3. Let $\mathcal{F} : \mathcal{X} \rightarrow [0, 1]$ be a discriminative classifier and D be a dataset where each example d is a joint assignment to \mathbf{X} . Given a family of distributions over C and \mathbf{X} , let \mathcal{P} denote the subset of them that conforms with \mathcal{F} . Then *conformant learning* on D is to solve the following optimization:

⁴Here we assume i.i.d. sampled data. If a true distribution is known, we can equivalently minimize the KL-divergence to it.

$$\arg \max_{P \in \mathcal{P}} \prod_{d \in D} P(d) = \arg \max_{P \in \mathcal{P}} \prod_{d=(\mathbf{x}) \in D} \sum_c P(\mathbf{x}, c). \quad (2)$$

The learned model thus conforms with \mathcal{F} and defines a feature distribution; therefore, we can take the expectation of \mathcal{F} via probabilistic inference. In other words, it attains the desired classification performance of the given discriminative model while also returning sophisticated predictions under missing features. Specifically, conformant naive Bayes models can be used to efficiently take expectations of logistic regression classifiers. Note that this does not contradict Theorem 1 which considers arbitrary pairs of LR and NB models.

4 Naive Conformant Learning

In this section, we study a special case of conformant learning – naive conformant learning (NaCL), and show how it can be solved as a geometric program.

A naive Bayes distribution is defined by a parameter set θ that consists of $\theta_c, \theta_{\bar{c}}$, and $\theta_{x|c}, \theta_{x|\bar{c}}$ for all x . *Naive conformant learning* outputs the naive Bayes distribution P_θ that maximizes the (marginal) likelihood given the dataset and conforms with a given logistic regression model \mathcal{F} .

We will next show that above problem can be formulated as a *geometric program*, an optimization problem of the form:

$$\begin{aligned} \min f_0(x) \\ \text{s.t. } f_i(x) \leq 1, \quad i = 1 \dots m \\ g_i(x) = 1, \quad i = 1 \dots p \end{aligned}$$

where each f_i is a posynomial and g_i monomial. A *monomial* is a function of the form $bx_1^{a_1} \dots x_n^{a_n}$ defined over positive real variables x_1, \dots, x_n where $b > 0$ and $a_i \in \mathbb{R}$. A *posynomial* is a sum of monomials. Every geometric program can be transformed into an equivalent convex program through change of variables, and thus its global optimum can be found efficiently [Boyd *et al.*, 2007].

To maximize the likelihood, we instead minimize its inverse. Let $n(\mathbf{x})$ denote the number of times an assignment \mathbf{x} appears in dataset D . Then the objective function is:

$$\prod_{d \in D} P_\theta(d)^{-1} = \prod_{\mathbf{x}} P_\theta(\mathbf{x})^{-n(\mathbf{x})} = \prod_{\mathbf{x}} \left(\sum_c \prod_{x \in \mathbf{x}} \theta_{x|c} \theta_c \right)^{-n(\mathbf{x})}.$$

Above formula, directly expanded, is not a posynomial. In order to express it as a posynomial we consider an auxiliary dataset D' constructed from D as follows: for each data point $d_j = (\mathbf{x}) \in D$, there are $d'_{j,c} = (\mathbf{x}, c) \in D'$ with weight $\alpha_{j,c}$ for all values of c . If the weights are such that $\alpha_{j,c} \geq 0$ and $\sum_c \alpha_{j,c} = 1$ for all j , then the inverse of the expected joint likelihood given the new dataset D' is

$$\begin{aligned} & \prod_{d'_{j,c}=(\mathbf{x},c) \in D'} P_\theta(\mathbf{x}, c)^{-\alpha_{j,c}} \\ &= \prod_{d_j=(\mathbf{x}) \in D} \prod_c P_\theta(\mathbf{x})^{-\alpha_{j,c}} P_\theta(c | \mathbf{x})^{-\alpha_{j,c}} \\ &= \prod_{d \in D} P_\theta(d)^{-1} \cdot \prod_{d_j=(\mathbf{x}) \in D, c} P_\theta(c | \mathbf{x})^{-\alpha_{j,c}}. \quad (3) \end{aligned}$$

DATASETS	SIZE	# CLASSES: DIST.	# FEATURES	FEATURE TYPES
MNIST	60K	10: BALANCED	784	INT. PIXEL VALUE
FASHION	60K	10: BALANCED	784	INT. PIXEL VALUE
COVTYPE	581K	7: UNBALANCED	54	CONT. & CATEGORICAL
ADULT	49K	2: UNBALANCED	14	INT. & CATEGORICAL
SPLICE	3K	3: UNBALANCED	61	CATEGORICAL

Table 2: Summary of our testbed.

For any $P_\theta \in \mathcal{P}$, the conditional distribution $P_\theta(C | \mathbf{X})$ is fixed by the logistic regression model; in other words, the last product term in Equation 3 is a constant. Therefore, maximizing the expected (joint) likelihood on a completed dataset must also maximize the marginal likelihood, which is our original objective. Intuitively, maximizing the joint likelihood on any dataset is equivalent to maximizing the marginal likelihood $P(\mathbf{X})$ if the conditional distribution $P(C | \mathbf{X})$ is fixed. Now our objective function can be written as a monomial in terms of the parameters:

$$\prod_{d'_{j,c} \in D'} P_\theta(d'_{j,c})^{-\alpha_{j,c}} = \prod_{d'_{j,c}=(\mathbf{x},c) \in D'} \left(\theta_c \prod_{x \in \mathbf{x}} \theta_{x|c} \right)^{-\alpha_{j,c}}. \quad (4)$$

Next, we express the set of conformant naive Bayes distributions \mathcal{P} as geometric program constraints in terms of θ . An NB model P_θ conforms with an LR \mathcal{F} if and only if its corresponding logistic regression weights, according to Lemma 1, match those of \mathcal{F} . Hence, $P_\theta \in \mathcal{P}$ precisely when

$$e^{w_0} \theta_c^{-1} \theta_{\bar{c}} \prod_{i=1}^n \theta_{x_i|c}^{-1} \theta_{x_i|\bar{c}} = 1 \quad (5)$$

$$e^{w_i} \theta_{x_i|c}^{-1} \theta_{x_i|\bar{c}} \theta_{x_i|c} \theta_{x_i|\bar{c}}^{-1} = 1, \quad \forall i \quad (6)$$

We also need to ensure that the parameters define a valid probability distribution (e.g., $\theta_c + \theta_{\bar{c}} = 1$). Because such posynomial equalities are not valid geometric program constraints, we instead relax them to posynomial inequalities:⁵

$$\theta_c + \theta_{\bar{c}} \leq 1, \theta_{x_i|c} + \theta_{x_i|\bar{c}} \leq 1, \theta_{x_i|\bar{c}} + \theta_{x_i|c} \leq 1, \quad \forall i \quad (7)$$

Putting everything together, naive conformant learning can be solved as a geometric program whose objective function is given by Equation 4 and constraints by Equations 5 – 7. We used the GPkit⁶ library to solve our geometric programs.

5 Empirical Evaluation

In this section, we empirically evaluate the performance of naive conformant learning (NaCL) and provide a detailed discussion of our method’s advantages over existing imputation approaches in practice.⁷ More specifically, we want to answer the following questions:

⁵The learned parameters may not sum to 1. They can still be interpreted as a multi-valued NB with same likelihood that conforms with \mathcal{F} . These constraints were always active in our experiments.

⁶<https://gpkit.readthedocs.io>

⁷Our implementation of the algorithm and experiments are available at <https://github.com/UCLA-StarAI/NaCL>.

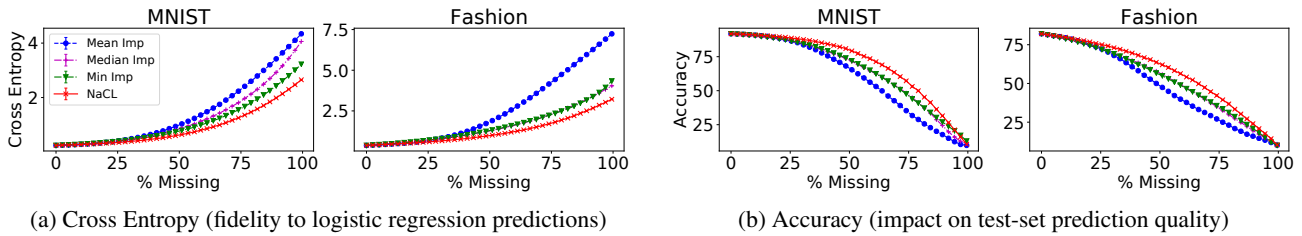


Figure 3: Standard image classification datasets: comparison of average cross entropy to the original predictions and classification accuracies between naive conformant learning (NaCL) and commonly used imputation methods. The conditional probabilities from NaCL are consistently closest to the full-data predictions, and NaCL consistently outperforms other methods with different percentages missing features.

CROSS ENTROPY	COVTYPE				ADULT				SPlice			
	20%	40%	60%	80%	20%	40%	60%	80%	20%	40%	60%	80%
MIN IMPUTATION	12.8	15.5	20.7	29.4	41.0	49.2	55.4	59.6	71.3	81.8	97.2	117.2
MAX IMPUTATION	52.6	89.1	133.5	187.7	84.2	114.6	125.3	114.5	70.5	78.3	89.3	103.2
MEAN IMPUTATION	12.8	15.6	21.2	30.5	34.1	38.7	44.8	52.6	69.2	74.7	82.4	92.3
MEDIAN IMPUTATION	12.8	15.7	21.4	30.8	35.3	41.2	48.6	57.8	70.0	75.7	83.0	92.0
NAIVE CONFORMANT LEARNING	12.6	14.8	18.9	25.8	33.6	37.0	41.2	46.6	69.1	74.7	82.8	94.0

Table 4: Three unbalanced UCI datasets with categorical features: comparison of average cross entropy to the original predictions between naive conformant learning (NaCL) and commonly used imputation methods. The closest are denoted in bold.

WEIGHTED F1	COVTYPE				ADULT				SPlice			
	20%	40%	60%	80%	20%	40%	60%	80%	20%	40%	60%	80%
MIN IMPUTATION	64.0	58.1	52.2	46.1	81.7	79.3	77.5	76.0	86.9	69.8	49.2	38.8
MAX IMPUTATION	49.8	44.4	41.6	37.3	81.7	79.3	77.4	76.0	86.9	69.8	49.1	38.8
MEAN IMPUTATION	64.0	58.0	52.2	46.3	82.9	79.8	75.3	70.7	91.8	82.3	66.2	45.7
MEDIAN IMPUTATION	64.0	58.1	52.2	46.1	82.7	79.2	74.8	70.5	89.4	77.6	59.5	42.5
NAIVE CONFORMANT LEARNING	66.1	61.7	56.9	51.7	83.4	81.2	77.9	73.5	93.3	87.2	76.6	59.1

Table 5: Three unbalanced UCI datasets with categorical features: comparison of weighted F1 scores between naive conformant learning (NaCL) and commonly used imputation. The highest are denoted in bold.

- Q1** Does NaCL reliably estimate the probabilities of the original logistic regression with full data? How do these estimates compare to those from imputation techniques, including ones that also model a feature distribution?
- Q2** Do higher-quality expectations of a logistic regression classifier result in higher accuracy on test data?
- Q3** Does NaCL retain logistic regression’s higher predictive accuracy over unconstrained naive Bayes?

5.1 Experimental Setup

To demonstrate the generality of our method, we construct a 5-dataset testbed suite that covers assorted configurations [Yann *et al.*, 2009; Xiao *et al.*, 2017; Blackard and Dean, 1999; Dua and Karra Taniskidou, 2017; Noordewier *et al.*, 1991]; see Table 2. The suite ranges from image classification to DNA sequence recognition; from fully balanced labels to > 75% of samples belonging to a single class; from continuous to categorical features with up to 40 different values. For datasets with no predefined test set, we construct one by a 80 : 20 split. As our method assumes binary inputs, we transform categorical features through one-hot encodings

and binarize continuous ones based on whether they are 0.05 standard deviation above their respective mean.

Our algorithm takes as input a logistic regression model which we trained using fully observed training data. During prediction time, we make the features go missing uniformly at random based on a set missingness percentage, which corresponds to a missing completely at random (MCAR) mechanism [Little and Rubin, 2014]. We repeat all experiments for 10 (resp. 100) runs on MNIST, Fashion, and CovType (resp. Adult and Splice) and report the average.

5.2 Fidelity to the Original Predictions

The optimal method to deal with missing values would be one that enables the original classifier to act as if no features were missing. In other words, we want the predictions to be affected as little as possible by the missingness. As such, we evaluate the similarity between predictions made with and without missingness, measured by the average cross entropy. The results are reported in Figure 3a⁸ and Table 4; the error bars were too small to be visibly noticeable and were omitted in

⁸Max imputation results are dismissed as they are orders of magnitude worse than the rest.

the plots. In general, our method outperforms all the baselines by a significant margin, demonstrating the superiority of the expected predictions produced by our method.

We also compare NaCL with two imputation methods that consider the feature distribution, namely EM [Dempster *et al.*, 1977] and MICE [Buuren and Groothuis-Oudshoorn, 2010]. EM imputation reports the second-to-worst average cross entropies and MICE’s results are very similar to those of mean imputation when 1% of features are missing. Due to the fact that both EM and MICE are excessively time-consuming to run and their imputed values are no better quality than more lightweight alternatives, we do not compare with them in the rest of the experiments. We would like to especially emphasize this comparison; it demonstrates that directly leveraging feature distributions without also considering how the imputed values impact the classifier may lead to unsatisfactory predictions, further justifying the need for solving the expected prediction task and conformant learning. This also concludes our answer to Q1.

5.3 Classification Accuracy

Encouraged by the fact that NaCL produces more reliable estimates of the conditional probability of the original logistic regression, we further investigate how much it helps achieve better classification accuracy under different percentages of missing features (i.e., Q2). As suggested by Figure 3b and Table 5,⁹ NaCL consistently outperforms all other methods except on the Adult dataset with 80% of the features missing.

Lastly, to answer Q3 we compare NaCL to a maximum-likelihood naive Bayes model.¹⁰ In all datasets except Splice, logistic regression achieves higher classification accuracy than naive Bayes with fully observed features. NaCL maintains this advantage until 40% of the features go missing, further demonstrating the effectiveness of our method. Note that these four datasets have a large number of samples, which is consistent with the prevailing consensus that discriminative learners are better classifiers given a sufficiently large number of samples [Ng and Jordan, 2002].

6 Case Study: Sufficient Explanations

In this section we briefly discuss utilizing conformant learning to explain classifications and show some empirical examples as a proof of concept.

On a high level, the task of explaining a particular classification can be thought of as quantifying the “importance” of each feature and choosing a small subset of the most important features as the explanation. Linear models are widely considered easy to interpret, and thus many explanation methods learn a linear model that is closely faithful to the original one, and then use the learned model to assign importance to features [Ribeiro *et al.*, 2016; Lundberg and Lee, 2017; Shrikumar *et al.*, 2017]. These methods often assume a black-box setting, and to generate explanations they internally eval-

⁹We report weighted F1 scores as the datasets are unbalanced.

¹⁰We do not report the full set of results in the table because maximum-likelihood learning of naive Bayes optimizes for a different loss and effectively solves a different task than NaCL and the imputation methods.

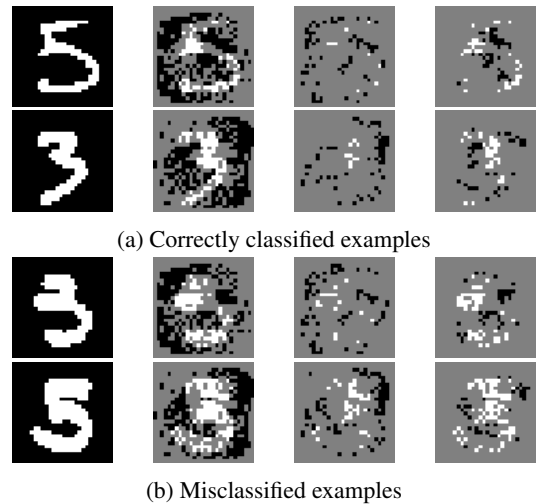


Figure 6: Explanations for MNIST classifications. Grey features were not chosen as explanations; white/black are the true color of chosen features. From left to right: 1) all features; 2) support features; 3) top- k support features; 4) sufficient explanation of size k .

uate the predictor on multiple perturbations of the given instance. A caveat is that the perturbed values may have a very low probability on the distribution the classifier was trained on. This can lead to unexpected results as machine learning models typically only guarantee generalization if both train and test data are drawn from the same distribution.

Instead we propose to leverage the feature distribution in producing explanations. To explain a given binary classifier, we consider a small subset of feature observations that is sufficient to get the same classification, in expectation w.r.t. a feature distribution. Next, we formally define our method:

Definition 4. (Support and Opposing Features) Given \mathcal{F} , P , and \mathbf{x} , we partition the given feature observations into two sets. The first set consists of the *support* features that contribute towards the classification of $\mathcal{F}(\mathbf{x})$:

$$\mathbf{x}_+ = \begin{cases} \{x \in \mathbf{x} : E_{\mathcal{F}, P}(\mathbf{x} \setminus x) \leq \mathcal{F}(\mathbf{x})\} & \text{if } \mathcal{F}(\mathbf{x}) \geq 0.5, \\ \{x \in \mathbf{x} : E_{\mathcal{F}, P}(\mathbf{x} \setminus x) > \mathcal{F}(\mathbf{x})\} & \text{otherwise.} \end{cases}$$

The rest are the *opposing* features that provide evidence against the current classification: $\mathbf{x}_- = \mathbf{x} \setminus \mathbf{x}_+$.

Definition 5. *Sufficient explanation* of $\mathcal{F}(\mathbf{x})$ with respect to P is defined as the following:

$$\begin{aligned} & \arg \min_{\mathbf{e} \subseteq \mathbf{x}_+} |\mathbf{e}| \\ & \text{s.t. } \text{sgn}(E_{\mathcal{F}, P}(\mathbf{e}\mathbf{x}_-) - 0.5) = \text{sgn}(\mathcal{F}(\mathbf{x}) - 0.5) \end{aligned}$$

Intuitively, this is the smallest set of support features that, in expectation, result in the same classification despite all the evidence to the contrary. In other words, we explain a classification using the strongest evidence towards it.

For a qualitative evaluation, we generate sufficient explanations on instances of a binary logistic regression task for MNIST digits 5 and 3; see the last column of Figure 6. Take the first example in Figure 6a: the white pixels selected as sufficient explanation show that the digit should be a 5. Also notice the black pixels in the explanation: they express how the

absence of white pixels significantly contributes to the classification, especially in parts of the image where they would be expected for the opposing class. Similarly, the black pixels in the first example in Figure 6b look like a 3, and the white pixels in the explanation look like a 5, explaining why this 3 was misclassified as a 5. We further compare our approach to an alternative one that selects a subset of support features based on their logistic regression weights; see the third column of Figure 6. It selects features that will cause a large difference in prediction if the value was flipped, as opposed to missing, which is what sufficient explanation considers.

7 Related Work

There have been many approaches developed to classify with missing values, which can broadly be grouped into two different types. The first one focuses on increasing classifiers’ inherent robustness to feature corruption, which includes missingness. A common way to achieve such robustness is to spread the importance weights more evenly among features [Globerson and Roweis, 2006; Dekel and Shamir, 2008; Xia *et al.*, 2017]. One downside of this approach is that the trained classifier may not achieve its best possible performance if no features go missing.

The second one investigates how to impute the missing values. In essence, imputation is a form of reasoning about missing values from observed ones [Sharpe and Solly, 1995; Batista *et al.*, 2002; McKnight *et al.*, 2007]. An iterative process is commonly used during this reasoning process [Buuren and Groothuis-Oudshoorn, 2010]. Some recent works also adapt auto-encoders and GANs for the task [Costa and *et al.*, 2018; Mattei and Frellsen, 2019]. Some of these works can be incorporated into a framework called multiple imputations to reflect and better bound one’s uncertainty [Schafer, 1999; Azur *et al.*, 2011; Gondara and Wang, 2018]. These existing methods focus on substituting missing values with those closer to the ground truth, but do not model how the imputed values interact with the trained classifier. On the other hand, our proposed method explicitly reasons about what the classifier is expected to return.

We are among the first to incorporate feature distributions to generate explanations. Notable recent work along this line proposes to maximize the mutual information between selected features and the class [Chen *et al.*, 2018]. To more explicitly leverage a feature distribution one can explain a classification by a subset of features that maximally affect the classifier output, when its values are substituted by in-fills sampled from the feature distribution conditioned on the rest of the features [Chang *et al.*, 2019]. This contrasts with our method which studies the affect of certain features on a classifier by marginalizing, rather than sampling.

8 Conclusion & Future Work

In this paper we introduced the expected prediction task, a principled approach to predicting with missing features. It leverages a feature distribution to reason about what a classifier is expected to return if it could observe all features. We then proposed conformant learning to learn joint distributions that conform with and can take expectations of discrim-

inative classifiers. A special instance of it—naive conformant learning—was shown empirically to outperform many existing imputation methods. For future work, we would like to explore conformant learning for other generative-discriminative pairs of models, and extend NaCL to real-valued features.

Acknowledgements

This work is partially supported by NSF grants #IIS-1657613, #IIS-1633857, #CCF-1837129, DARPA XAI grant #N66001-17-2-4032, NEC Research, and gifts from Intel and Facebook Research.

A Proofs

A.1 Proof of Theorem 1

The proof is by reduction from computing the same-decision probability, whose decision problem **D-SDP** was shown to be NP-hard. [Chen *et al.*, 2013]

Given a naive Bayes distribution $P(\cdot)$ over variables C and \mathbf{X} , a threshold T , and a probability p , **D-SDP** asks: is the same-decision probability $\sum_{\mathbf{x}} \mathbb{I}(P(c|\mathbf{x}) > T)P(\mathbf{x})$ greater than p ? Here, $\mathbb{I}(\cdot)$ denotes an indicator function which returns 1 if the enclosed expression is true, and 0 otherwise.

Using Lemma 1 we can efficiently translate a naive Bayes model P to a logistic regression with a weight function $w(\cdot)$ such that

$$P(c|\mathbf{x}) = \frac{1}{1 + e^{-w(\mathbf{x})}}.$$

Note that $P(c|\mathbf{x}) > T$ iff $w(\mathbf{x}) > -\log(\frac{1}{T} - 1)$. Then we construct another logistic regression with weight function

$$w'(\mathbf{x}) = n \cdot \left(w(\mathbf{x}) + \log\left(\frac{1}{T} - 1\right) \right),$$

for some positive constant n . As w is a linear model, w' is also linear, and $w'(\mathbf{x}) > 0$ iff $P(c|\mathbf{x}) > T$. As n grows, $w'(\cdot)$ approaches ∞ and $-\infty$ for positive and negative examples, respectively. Hence, this logistic regression model outputs 1 if $P(c|\mathbf{x}) > T$ and 0 otherwise, effectively being an indicator function. Therefore, the expectation of such classifier over $P(\mathbf{X})$ is equal to the same-decision probability of \mathbf{X} . \square

A.2 Proof of Lemma 1

We want to prove there is a unique \mathcal{F} such that $\mathcal{F}(\mathbf{x}) = P(c|\mathbf{x})$ for all \mathbf{x} given naive Bayes distribution P . Using Bayes’ rule and algebraic manipulation, we get:

$$\begin{aligned} P(c|\mathbf{x}) &= \frac{P(\mathbf{x}|c)P(c)}{P(\mathbf{x}|c)P(c) + P(\mathbf{x}|\bar{c})P(\bar{c})} \\ &= \frac{1}{1 + \frac{P(\mathbf{x}|\bar{c})P(\bar{c})}{P(\mathbf{x}|c)P(c)}} = \frac{1}{1 + \exp\left[-\log\frac{P(\mathbf{x}|c)P(c)}{P(\mathbf{x}|\bar{c})P(\bar{c})}\right]} \end{aligned}$$

For any input \mathbf{x} , we want above quantity to be equal to $\mathcal{F}(\mathbf{x}) = 1/(1 + \exp[-\sum_i w_i x_i])$. In other words, we need:

$$\log\frac{P(\mathbf{x}|c)P(c)}{P(\mathbf{x}|\bar{c})P(\bar{c})} = \sum_i w_i x_i$$

Using naive Bayes assumption, we arrive at:

$$\sum_{i=0}^n w_i x_i = \log \frac{P(c)}{P(\bar{c})} + \sum_{i=1}^n \log \frac{P(x_i | c)}{P(x_i | \bar{c})} \quad (1)$$

Now we want the RHS of Equation 1 to be a linear function of x_i 's, so we do the following substitution for $i > 0$ assuming binary features:

$$\begin{aligned} \log \frac{P(x_i | c)}{P(x_i | \bar{c})} &= (x_i) \cdot \log \frac{P(x_i = 1 | c)}{P(x_i = 1 | \bar{c})} \\ &+ (1 - x_i) \cdot \log \frac{P(x_i = 0 | c)}{P(x_i = 0 | \bar{c})} \end{aligned} \quad (2)$$

By combining Equations 1 and 2 we get the weights in Lemma 1 by simple algebraic manipulations. To solve for the bias term w_0 we plug in $x_i = 0$ for all $i > 0$. To compute w_i for a non-zero i we take the coefficient of x_i in Equation 2. \square

A.3 Proof of Lemma 2

Through the same algebraic manipulation as before, we get the same equations as in Lemma 1 with the only difference being that we are now solving the parameters of a naive Bayes model rather than weights of the logistic regression model. Intuitively, because the NB model has $2n + 1$ free parameters but the LR model only has $n + 1$ parameters, we expect some degree of freedom. To get rid of the freedom and get a unique solution, we fix the values for n parameters as follows:

$$P(x_i = 1 | c) = \theta_i \quad (3)$$

Without loss of generality we have fixed the parameter values for positive features. One can equally set the values in other ways as long as one parameter value per feature is fixed.

Now there is a unique naive Bayes model that matches the logistic regression classifier and also agrees with Equation 3. That is, the remaining $n + 1$ parameter values are given by the LR parameters, resulting in the following parameters for such naive Bayes model:

$$\begin{aligned} P(x_i | c) &= \theta_i, & P(\bar{x}_i | c) &= 1 - P(x_i | c), \\ P(x_i | \bar{c}) &= \frac{1}{1 + e^{w_i \frac{1 - \theta_i}{\theta_i}}}, & P(\bar{x}_i | \bar{c}) &= 1 - P(x_i | \bar{c}), \\ P(c) &= \text{sigmoid} \left(w_0 - \sum_{i=1}^n \log \frac{P(\bar{x}_i | c)}{P(\bar{x}_i | \bar{c})} \right). \end{aligned}$$

\square

B Beyond Binary Classification: Multiclass

In the paper, we studied logistic regression and conformant naive Bayes models assuming binary classification. We now show that our method can easily be extended to multiclass. We first modify our notation of logistic regression and naive Bayes models to allow for an arbitrary number of classes.

Definition 6. (Multiclass Classifiers) Suppose we have a classifier with K classes, each denoted by c_k ($k \in [0, K - 1]$).

Then \mathcal{F}_k denotes the conditional probability for class c_k in logistic regression, and P a naive Bayes distribution defined as:

$$\begin{aligned} \mathcal{F}_k(\mathbf{x}) &= \frac{e^{W_k \cdot \mathbf{x}}}{\sum_j e^{W_j \cdot \mathbf{x}}} \\ P(c_k | \mathbf{x}) &= \frac{P(c_k) \prod_{i=1}^n P(x_i | c_k)}{\sum_j P(c_j) \prod_{i=1}^n P(x_i | c_j)} \end{aligned}$$

We say P conforms with \mathcal{F} if their predictions agree for all classes: $P(c_k | \mathbf{x}) = \mathcal{F}_k(\mathbf{x})$ for all k and \mathbf{x} .

Next, we describe naive conformant learning for multiclass. Instead of directly matching the predictions of P and \mathcal{F} for all classes, we match their ratios in order to simplify equations going forward. Moreover, to get the same classifiers it suffices to divide by the probability of only one class, so without loss of generality we set the following be true.

$$\frac{\mathcal{F}_k(\mathbf{x})}{\mathcal{F}_0(\mathbf{x})} = \frac{P(c_k | \mathbf{x})}{P(c_0 | \mathbf{x})}, \quad \forall k \in [1, K - 1]$$

Using Definition 6, this leads to

$$\frac{e^{-W_k \cdot \mathbf{x}}}{e^{-W_0 \cdot \mathbf{x}}} = \frac{P(c_k) \prod_{i=1}^n P(x_i | c_k)}{P(c_0) \prod_{i=1}^n P(x_i | c_0)}, \quad \forall k \in [1, K - 1].$$

The parameters of a multiclass naive Bayes are: $\theta_{c_k} = P(c_k)$, $\theta_{x_i | c_k} = P(x_i = 1 | c_k)$, and $\theta_{\bar{x}_i | c_k} = P(x_i = 0 | c_k)$. Then, we get the following constraints for NaCL through similar algebraic manipulations as in the binary case:

$$\begin{aligned} \sum_k \theta_{c_k} &= 1 & (4) \\ \theta_{x_i | c_k} + \theta_{\bar{x}_i | c_k} &= 1, \quad \forall i, k > 0 & (5) \\ e^{w_{k,i} - w_{0,i}} \theta_{x_i | c_k}^{-1} \theta_{\bar{x}_i | c_k} \theta_{x_i | c_0} \theta_{\bar{x}_i | c_0}^{-1} &= 1, \quad \forall i, k > 0 \\ e^{w_{k,0} - w_{0,0}} \theta_{c_k}^{-1} \theta_{c_0} \prod_{i=1}^n \theta_{\bar{x}_i | c_k}^{-1} \theta_{\bar{x}_i | c_0} &= 1, \quad \forall k > 0 \end{aligned}$$

Again, we relax Equations 4 and 5 to inequalities to obtain valid geometric program constraints. The rest of the method stays the same: we maximize the marginal likelihood with above constraints by minimizing the inverse of the joint likelihood on a completed dataset, as described in Section 4.

References

- [Azur *et al.*, 2011] Melissa J Azur, Elizabeth A Stuart, Constantine Frangakis, and Philip J Leaf. Multiple imputation by chained equations: what is it and how does it work? *International journal of methods in psychiatric research*, 2011.
- [Batista *et al.*, 2002] Gustavo EAPA Batista, Maria Carolina Monard, et al. A study of k-nearest neighbour as an imputation method. *HIS*, 2002.
- [Blackard and Dean, 1999] Jock A Blackard and Denis J Dean. Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables. *Computers and electronics in agriculture*, 1999.

- [Bouchard and Triggs, 2004] Guillaume Bouchard and Bill Triggs. The tradeoff between generative and discriminative classifiers. In *COMPSTAT*, 2004.
- [Boyd *et al.*, 2007] Stephen Boyd, Seung-Jean Kim, Lieven Vandenbergh, and Arash Hassibi. A tutorial on geometric programming. *Optimization and engineering*, 2007.
- [Buuren and Groothuis-Oudshoorn, 2010] Stef van Buuren and Karin Groothuis-Oudshoorn. MICE: Multivariate imputation by chained equations in R. *Journal of statistical software*, 2010.
- [Chang *et al.*, 2019] Chun-Hao Chang, Elliot Creager, Anna Goldenberg, and David Duvenaud. Explaining image classifiers by counterfactual generation. In *ICLR*, 2019.
- [Chen *et al.*, 2013] Suming Jeremiah Chen, Arthur Choi, and Adnan Darwiche. An exact algorithm for computing the same-decision probability. In *IJCAI*, 2013.
- [Chen *et al.*, 2018] Jianbo Chen, Le Song, Martin Wainwright, and Michael Jordan. Learning to explain: An information-theoretic perspective on model interpretation. In *ICML*, 2018.
- [Costa and et al., 2018] Adriana Fonseca Costa and et al. Missing data imputation via denoising autoencoders: the untold story. In *Int'l Symposium on Intelligent Data Analysis*, 2018.
- [Darwiche, 2009] Adnan Darwiche. *Modeling and reasoning with Bayesian networks*. Cambridge University Press, 2009.
- [Dechter, 2013] Rina Dechter. Reasoning with probabilistic and deterministic graphical models: Exact algorithms. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 2013.
- [Dekel and Shamir, 2008] Ofer Dekel and Ohad Shamir. Learning to classify with missing and corrupted features. In *ICML*, 2008.
- [Dempster *et al.*, 1977] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 1977.
- [Dua and Karra Taniskidou, 2017] Dheeru Dua and Efi Karra Taniskidou. UCI machine learning repository, 2017.
- [Globerson and Roweis, 2006] Amir Globerson and Sam Roweis. Nightmare at test time: Robust learning by feature deletion. In *ICML*, 2006.
- [Gondara and Wang, 2018] Lovedeep Gondara and Ke Wang. Mida: Multiple imputation using denoising autoencoders. In *SIGKDD*, 2018.
- [Graham, 2012] John W Graham. *Missing data: Analysis and design*. Springer Science & Business Media, 2012.
- [Jaakkola and Haussler, 1999] Tommi Jaakkola and David Haussler. Exploiting generative models in discriminative classifiers. In *NeurIPS*, 1999.
- [Koller and Friedman, 2009] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [Liang and Van den Broeck, 2019] Yitao Liang and Guy Van den Broeck. Learning logistic circuits. In *AAAI*, 2019.
- [Little and Rubin, 2014] Roderick JA Little and Donald B Rubin. *Statistical analysis with missing data*. 2014.
- [Lundberg and Lee, 2017] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *NeurIPS*, 2017.
- [Mattei and Frellsen, 2019] Pierre-Alexandre Mattei and Jes Frellsen. MIWAE: Deep generative modelling and imputation of incomplete data sets. In *ICML*, 2019.
- [McKnight *et al.*, 2007] Patrick E McKnight, Katherine M McKnight, Souraya Sidani, and Aurelio Jose Figueredo. *Missing data: A gentle introduction*. Guilford Press, 2007.
- [Ng and Jordan, 2002] Andrew Y Ng and Michael I Jordan. On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes. In *NeurIPS*, 2002.
- [Noordewier *et al.*, 1991] Michiel O Noordewier, Geoffrey G Towell, and Jude W Shavlik. Training knowledge-based neural networks to recognize genes in dna sequences. In *NeurIPS*, 1991.
- [Ribeiro *et al.*, 2016] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should I trust you?: Explaining the predictions of any classifier. In *SIGKDD*, 2016.
- [Roos *et al.*, 2005] Teemu Roos, Hannes Wettig, Peter Grünwald, Petri Myllymäki, and Henry Tirri. On discriminative bayesian network classifiers and logistic regression. *Machine Learning*, 59(3), Jun 2005.
- [Roth, 1996] Dan Roth. On the hardness of approximate reasoning. *Artificial Intelligence*, 1996.
- [Schafer, 1999] Joseph L Schafer. Multiple imputation: a primer. *Statistical methods in medical research*, 1999.
- [Sharpe and Solly, 1995] Peter K. Sharpe and RJ Solly. Dealing with missing values in neural network-based diagnostic systems. *Neural Computing & Applications*, 1995.
- [Shrikumar *et al.*, 2017] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *ICML*, 2017.
- [Sutton and McCallum, 2012] Charles Sutton and Andrew McCallum. An introduction to conditional random fields. *Foundations and Trends in Machine Learning*, (4), 2012.
- [Ulusoy and Bishop, 2005] Ilkay Ulusoy and Christopher M Bishop. Generative versus discriminative methods for object recognition. In *Proceedings of CVPR*, 2005.
- [Xia *et al.*, 2017] Jing Xia, Shengyu Zhang, Guolong Cai, Li Li, Qing Pan, Jing Yan, and Gangmin Ning. Adjusted weight voting algorithm for random forests in handling missing values. *Pattern Recognition*, 2017.
- [Xiao *et al.*, 2017] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. 2017.
- [Yann *et al.*, 2009] LeCun Yann, Cortes Corinna, and Christopher JC Burges. The MNIST database of handwritten digits, 2009.