IS IT WORTH IT TO COLLECT MISSING VALUES?: THE MISSING VALUE UNCERTAINTY PROBLEM

Anonymous authors

Paper under double-blind review

ABSTRACT

In high-stakes domains like healthcare, operators often face the critical decision of whether to act on incomplete information or incur costs to collect missing values. Existing methods typically focus on imputing missing data or quantifying model uncertainty, but they do not directly assess the stability of a prediction if missing features were to be revealed. To address this gap, we introduce a framework for Missing Value Uncertainty (MVU), which is the distribution of predictions induced by incomplete inputs. We formalize the problem by defining *hard confidence*: the probability that a prediction will not change after collecting the missing data. We propose a novel Direct Missing Value (DMV) to efficiently estimate the MVU distribution, bypassing the need for expensive Monte Carlo sampling. Second, we introduce the Missing Value Calibration Error (MVCE), a new metric specifically designed to evaluate the calibration of hard confidence values, and a post-hoc calibration procedure to improve MVU estimation. We showcase our method and metric on synthetic and real-world datasets.

1 Introduction

In high-stakes domains such as healthcare and security, decisions are often made with incomplete information. This raises a critical and practical question for a human operator: **Is it worth the cost and effort to collect missing input values for a specific instance at inference time?** For example, a doctor with a patient's initial lab results must decide whether this information is sufficient to make a diagnosis or if more costly and invasive tests are required. Similarly, a security analyst observing a potential threat with data from a partially failed sensor network must determine whether to act immediately or deploy resources to gather more information. If additional information is unlikely to alter the optimal course of action, an operator can proceed without collecting missing values. However, if the missing values could significantly change the decision, the most prudent action is to collect them first. This paper centers on developing a framework to help an operator make this crucial decision.

Prior work has addressed aspects of this problem but fails to directly answer the central question. Missing value literature, for instance, primarily focuses on developing methods like imputation to handle missing data and make the best possible prediction given the observed values (Little & Rubin, 2019; Azur et al., 2011). While useful, these techniques do not quantify the uncertainty introduced by the missing features, leaving the operator unsure of how the prediction might change if the missing values were revealed.

Separately, research in uncertainty quantification has focused on a different question: "How likely is the prediction to be correct?". This question is useful for deciding whether to accept a model's prediction, corresponding to the notion of **soft confidence** that we explain in Section 2). However, it does not inform an operator about the stability of the prediction. Our work instead focuses on a notion called **hard confidence**: the probability that a prediction will *not* change if missing values are collected (see Section 2 for formal definition). Furthermore, existing work on epistemic uncertainty—i.e., uncertainty that can be reduced by more information—typically addresses model uncertainty arising from limited training data (Liu et al., 2019). While this may inform a decision maker if more training data is needed, it is orthogonal to the uncertainty from missing features for a single test-time instance, where the solution is to collect more features, not more training data. Given all this, to our knowledge, *no prior work has formalized, estimated, or evaluated the uncertainty stemming*

specifically from missing values at inference time to aid in deciding whether to collect missing values or not. This represents a critical gap, as it leaves a practical decision-making problem without a principled solution.

To fill this gap, we propose a complete framework for analyzing Missing Value Uncertainty (MVU), the distribution of predictions induced by missing inputs. We formalize the decision-making problem through the lens of hard confidence, which directly quantifies prediction stability. We establish imputation and Monte-Carlo baselines for estimating MVU and propose a novel explicit method for estimating MVU: the Direct Missing Value (DMV) estimator, which is significantly more efficient by circumventing sampling of the missing values, particularly in high dimensions. To evaluate these methods, we develop the Missing Value Calibration Error (MVCE), a new metric designed specifically for assessing the calibration of hard confidence values. Finally, we introduce a post-hoc procedure to improve the calibration of any MVU estimation method. Our main contributions are:

- We formalize the problem of missing value uncertainty (MVU) to aid operator decisions about collecting more information.
- We develop DMV, a novel and efficient approach that directly estimates MVU without requiring expensive Monte Carlo sampling.
- We define a novel metric, the Missing Value Calibration Error (MVCE), for evaluating MVU estimation methods on any dataset and propose a MVCE-based post-hoc calibration approach that can improve an MVU method after training.
- We empirically validate our methods on both synthetic and real-world datasets.

2 SOFT AND HARD VOTING CLASSIFICATION RULES AND CONFIDENCES

Given a distribution of predicted probabilities $p(\Phi)$, which implicitly represents a weighted ensemble of model predictions, we consider two decision rules: soft and hard voting. **Soft voting** averages class probabilities directly and corresponds to the standard Bayes optimal classification rule. In contrast, **hard voting** averages the argmax predictions from each model, forming a robust ensemble classifier. While soft voting is well-studied, we focus on the under-explored hard voting approach, as its resulting confidence value is more actionable for deciding whether to collect missing data. Though we apply this approach to missing value uncertainty, it may be of independent interest for other types of epistemic uncertainty. We first discuss these rules generically before specializing to the missing value case in the next section.

Notation For discrete class distributions, we will often use p(Y) (or conditional variants like $p(Y|X_{\mathcal{O}})$) to denote either the distribution itself or the vector of probabilities [p(Y=1), p(Y=2), ..., p(Y=k)], which fully defines the distribution. In many cases these probabilities are defined by a vector or parameters ϕ such that $p(Y=j) \equiv \phi_j$.

2.1 Background: Bayes Optimal Classification via Soft Voting

Given a distribution of predictions $p(\Phi)$ where Φ represent class probabilities, the soft-voting rule uses Φ directly as a soft vote and averages over the $p(\Phi)$ distribution:

$$y_{\text{soft}} \triangleq \arg\max_{j} \mathbb{E}_{p(\Phi)}[\Phi]_{j} \equiv \arg\max_{j} p(Y=j), \ c_{\text{soft}} \triangleq \max_{j} \mathbb{E}_{p(\Phi)}[\Phi]_{j} \equiv \max_{j} p(Y=j), \ (1)$$

where $p(Y=j)=\mathbb{E}_{p(\Phi)}[p(Y=j;\Phi)]=\mathbb{E}_{p(\Phi)}[\Phi]_j$ is the marginal probability of Y when marginalizing over the uncertainty in $p(\Phi)$. This soft voting classification is equivalent to the *Bayes optimal classification rule*. The Bayes rule is the most common way to classify because it minimizes the misclassification error. The confidence value $c_{\rm soft}$ is simply the probability that the selected class is correct. When given to a human operator, they could use $c_{\rm soft}$ to determine whether they should accept or ignore the prediction depending on the confidence level required. If the confidence is high, the operator could confidently accept the prediction. However, if the confidence is low (e.g., close to 1/k), then the operator should simply ignore the prediction since it is uninformative. Thus, from a practical standpoint, it mostly provides useful and actionable information if the confidence is high.

2.2 CUMULATIVE PROBABILITY CLASSIFICATION VIA HARD VOTING

We propose to use hard voting (i.e., majority voting) as an alternative and complementary classification rule that yields distinct information compared to soft voting confidence. In particular, we will explain why it is useful for deciding whether to collect more information or not. The hard-voting rule uses the argmax of Φ (i.e., a single class hard vote) and averages these hard votes over the $p(\Phi)$ distribution:

$$y_{\text{hard}} \triangleq \arg\max_{j} \mathbb{E}_{p(\Phi)}[\text{OneHotArgmax}(\Phi)]_{j} \equiv \arg\max_{j} p(Y_{\text{vote}} = j)$$
 (2)

$$\equiv \arg\max_{j} \Pr(\bigcap_{j' \neq j} \Phi_j \ge \Phi_{j'}), \tag{3}$$

$$c_{\text{hard}} \triangleq \max_{i} \mathbb{E}_{p(\Phi)}[\text{OneHotArgmax}(\Phi)]_{i} \equiv \max_{i} p(Y_{\text{vote}} = j),$$
 (4)

Hard Voting Confidence Values For Deciding Whether to Collect More Information In the context of epistemic uncertainty where the uncertainty could be reduced by collecting more information, the hard voting confidence values provide the probability that the prediction would stay the same if all missing values were revealed. If the hard confidence is high, then gathering more information will not change the predicted class. If the hard confidence is low, then it means that more information could change the predicted class.

Comparing Soft Voting and Hard Voting Soft and hard voting confidences capture complementary information for decision-making: soft confidence helps determine whether to accept a prediction as being accurate, while hard confidence informs the decision to collect more information. A few distinctions are:

- Behavior with No Uncertainty: For a degenerate distribution $p(\Phi)$ where epistemic uncertainty is zero, the predictions are identical. However, soft confidence can vary between 1/k and 1, while hard confidence is always 1, correctly reflecting that no new information will change the outcome.
- Sensitivity to Variance: Soft confidence depends only on the mean of the distribution $p(\Phi)$ and is insensitive to its variance. In contrast, hard confidence decreases as the variance increases, thereby capturing the spread of the epistemic uncertainty, not just its central tendency.

3 Missing Value Uncertainty (MVU)

Having discussed uncertainty from a generic distribution of predictions $p(\Phi)$, we now formalize the Missing Value Uncertainty (MVU) distribution, which arises from incomplete inputs at inference time. To distinguish MVU from standard epistemic uncertainty, we briefly contrast it with the uncertainty stemming from finite training data.

Standard epistemic uncertainty typically refers to model uncertainty from a finite training set \mathcal{D}_n . In a Bayesian context, this induces a posterior over model parameters, $p(\Theta|\mathcal{D}_n)$, which in turn creates a distribution of predictions. This uncertainty is reducible, as it vanishes in the limit of infinite training samples $(n \to \infty)$. In contrast, MVU is an orthogonal form of epistemic uncertainty induced by missing features for a single test-time instance. It is reducible not by collecting more training data, but by observing the missing features of that specific instance. We now formalize this concept.

Definition 1 (Missing Value Uncertainty Distribution). Given a joint distribution p(X, Y), let us define the true class distribution $\pi(x)$ given complete inputs and the corresponding random variable Φ as:

$$\pi(x) \triangleq p(Y|X=x), \qquad \Phi \triangleq \pi(X)$$
 (5)

where $\pi: \mathcal{X} \to \Delta^{|\mathcal{Y}|-1}$ is a deterministic function mapping a complete input to a probability vector on the simplex and Φ is the random variable representing these class probabilities given the random input X. Given these, the Missing Value Uncertainty (MVU) distribution given an observed set of input values $X_{\mathcal{O}}$ is defined as:

$$p(\Phi|X_{\mathcal{O}} = x_{\mathcal{O}}) \equiv \pi_{\sharp}^{\mathcal{O}} p(X_{\mathcal{M}}|X_{\mathcal{O}} = x_{\mathcal{O}}), \tag{6}$$

where $\pi^{\mathcal{O}}(x_{\mathcal{M}}) \triangleq \pi(x_{\mathcal{M}}, x_{\mathcal{O}})$ is π conditioned on the observed inputs and $\pi^{\mathcal{O}}_{\sharp}$ denotes the measure pushforward operator.

It should be noted that the target MVU distribution $p(\Phi|X_{\mathcal{O}}=x_{\mathcal{O}})$ is well-defined for any joint distribution p(X,Y). Furthermore, note that π is the optimal probabilistic classifier given complete inputs since it directly gives the class distribution conditioned on a complete input x. Finally, compared to epistemic uncertainty induced by finite samples, this uncertainty is induced by incomplete or partial inputs. This uncertainty distribution becomes degenerate (i.e., perfectly certain) when all inputs are observed. Thus, it aligns naturally with the view of epistemic uncertainty that it becomes zero when all information is revealed.

Soft and Hard Voting for MVU Applying the voting rules from Section 2 to the MVU distribution provides distinct types of information. Soft voting corresponds to the standard Bayes optimal classification given the observed features: $\arg\max_j p(Y=j|X_{\mathcal{O}}=x_{\mathcal{O}})$. As previously discussed, the resulting soft confidence does not inform an operator whether collecting more inputs would be useful. Hard voting, however, directly addresses this problem. The hard confidence derived from $\arg\max_j p(Y_{\text{hard}}=j|X_{\mathcal{O}}=x_{\mathcal{O}})$ represents the probability that the prediction will *not* change if the missing values are revealed. A high hard confidence therefore suggests that the decision is stable and collecting more data is unnecessary, a critical insight for operators in fields like medical diagnosis or sensor networks where acquiring more information is costly.

4 METHODS FOR MVU ESTIMATION

Having established hard voting confidence in Section 2 and defining missing value uncertainty in Section 3, we now consider ways to estimate the MVU distribution $p(\Phi|X_{\mathcal{O}})$ for a distribution of $X_{\mathcal{O}}$. Because there are no prior methods that focus on estimating MVU to the authors' best knowledge, we first propose two natural baselines based on missing value imputation and Monte Carlo estimates using generative models. After establishing these baselines, we then introduce our novel Direct Missing Value (DMV) approach for directly estimating the MVU distribution without requiring imputation or sampling.

4.1 BASELINE MVU METHODS

Imputation with Simple Variance Leveraging prior work, we can handle missing values using a simple imputation approach such as zero imputation or mean imputation (Little & Rubin, 2019). However, imputation approaches are limited by the classifier's inability to estimate missing value uncertainty, so our best bet is a simple heuristic to estimate variance then use method of moments to estimate distribution parameters. A naive approach would be to choose some constant variance regardless of the sample, though with any constant other than 0 (which induces a degenerate distribution), we run the risk of that uncertainty being too large for the predicted mean (as distributions over probabilities tend to have restrictions on maximum variance). A better heuristic for estimating Dirichlet parameters is to simply scale the predicted probability ϕ by some constant to produce α ; this approach guarantees any positive scaling constant will produce a valid distribution, though it may be unintuitive to set the constant. We can make the scaling constant more intuitive by instead scaling the maximum variance - Dirichlet variance cannot be greater than or equal to $\phi \cdot (1-\phi)$ without producing non-positive α values. Leveraging this knowledge, we can use a variance of $\phi \cdot (1-\phi) \cdot s$ where s is between 0 and 1. One flaw with the simple variance approaches is our uncertainty does not change with respect to the specific missing features; we get the exact same uncertainty for a fully observed input as an input that imputed those same values. This will lead to the model being overconfident under large numbers of missing features, and underconfident with fully observed data.

Monte Carlo Approximation A simple MVU approach is to use Monte Carlo samples of the missing values given observed values, i.e., samples from $p(X_M|x_O)$, to empirically estimate the

MVU distribution. This is motivated by the fact that the MVU distribution is the pushforward of the missing value distribution $p(X_{\mathcal{M}}|x_{\mathcal{O}})$ in Equation (6). At a high enough number of samples, this will produce an approximation close to the true MVU distribution. The key challenge involves sampling from the missing values given an (arbitrary) set of observed values $x_{\mathcal{O}}$ and unknown missingness pattern. Lower-dimensional cases could use a multivariate normal distribution, which has a simple closed form conditional distribution; however, in very high-dimensional settings such as image data, this is a poor approximation. Thus, we propose leveraging prior work on image inpainting in higher dimension cases, such as diffusion models (Zhang et al., 2023). Using the diffusion model to directly sample from $p(\Phi|x_{\mathcal{O}})$ is infeasible as such a large number of samples would take too long, so instead we recommend sampling $p(X_{\mathcal{M}}|x_{\mathcal{O}})$ to estimate $p(\Phi|x_{\mathcal{O}})$, then sampling the estimated distribution. However, this Monte Carlo approach is still very expensive as it requires sampling from high-dimensional conditional models. This motivates our proposed method which is a much more efficient alternative that does not require sampling from high-dimensional distributions.

4.2 DIRECT MISSING VALUE UNCERTAINTY (DMV)

We propose a novel direct estimator of the MVU distributions based on minimizing the KL divergence (or equivalently the negative log likelihood (NLL)) between the true and estimated MVU distributions: $\arg\min_{\psi}\mathbb{E}_{X_{\mathcal{O}}}[\mathrm{KL}(p(\Phi|X_{\mathcal{O}}),\hat{p}_{\psi}(\Phi|X_{\mathcal{O}}))] \equiv \arg\min_{\psi}\mathbb{E}_{X_{\mathcal{O}}}[\mathbb{E}_{\Phi|X_{\mathcal{O}}}[-\log\hat{p}_{\psi}(\Phi|X_{\mathcal{O}}))]]$, where ψ represents the model parameters. While at first this might seem like a standard problem, the challenge is that Φ is latent rather than observed. Thus, we cannot directly estimate the NLL given only samples from $X_{\mathcal{O}}$. However, if we know $\pi(x) \triangleq p(Y|X=x)$ from Equation (5), then we can convert this NLL objective to an objective that only requires complete training samples, i.e., samples with all features (proof in appendix).

Proposition 1. Given the optimal probabilistic predictor π from Equation (5) and any set of observed feature indices \mathcal{O} , the following holds, where γ is a constant w.r.t. ψ but does depend on π :

$$\mathbb{E}_{X_{\mathcal{O}}}[\mathrm{KL}(p(\Phi|X_{\mathcal{O}}), \hat{p}_{\psi}(\Phi|X_{\mathcal{O}}))] = \mathbb{E}_{X_{\mathcal{O}}, X_{\mathcal{M}}}[-\log \hat{p}_{\psi}(\pi(X_{\mathcal{O}}, X_{\mathcal{M}}) \mid X_{\mathcal{O}})] + \gamma_{\pi}. \tag{7}$$

The right hand side of Proposition 1 gives a natural way to directly estimate the uncertainty distribution given only complete samples and an estimate of the optimal complete predictor π . Importantly, this approach can directly estimate the uncertainty distribution while elegantly bypassing the need to do conditional sampling of missing values given observed values.

Estimating both the optimal classifier and the MVU distributions from training data. Given the above result, we propose a natural two-stage approach to estimating MVU distributions. First, we estimate the optimal predictor π using standard supervised learning on the complete dataset. Second, we plug-in this estimated predictor into the objective above to learn the final MVU predictor. At first glance, it may seem that you could simply minimize the linear combination of the standard cross entropy loss for $\hat{\pi}_{\theta}$ and the objective above for \hat{p}_{ψ} :

$$\arg\min_{\theta,\psi} \left(\mathbb{E}_{X,Y}[\ell_{\text{CE}}(\hat{\pi}_{\theta}(X),Y)] + \mathbb{E}_{X_{\mathcal{O}},X_{\mathcal{M}}}[-\log \hat{p}_{\psi}(\hat{\pi}_{\theta}(X_{\mathcal{O}},X_{\mathcal{M}}) \mid X_{\mathcal{O}})] + \gamma_{\hat{\pi}_{\theta}} \right), \quad (8)$$

where $\ell_{\rm CE}$ is the standard cross entropy loss, θ are the parameters for the predictor on complete inputs (i.e., no missing values) and ψ are the parameters of the MVU predictor given observed inputs $X_{\mathcal{O}}$. However, Equation (8) would incorrectly ignore the fact that the constant in Proposition 1 depends on $\hat{\pi}_{\theta}$ and thus in this case would depend on θ . Moreover, the γ_{π} term is not possible to approximate since we do not know the density of $p(\Phi \mid X_{\mathcal{O}})$. Therefore, we propose a bi-level optimization problem that will be valid because the lower-level problem assumes that the corresponding upper level variables are fixed:

$$\min_{\theta,\psi} \ \mathbb{E}_{X,Y}[\ell_{\mathrm{CE}}(\hat{\pi}_{\theta}(X),Y)], \quad \text{s.t. } \psi \in \arg\min_{\tilde{\psi}} \ \mathbb{E}_{X_{\mathcal{O}},X_{\mathcal{M}}}[-\log \hat{p}_{\tilde{\psi}}(\hat{\pi}_{\theta}(X_{\mathcal{O}},X_{\mathcal{M}}) \mid X_{\mathcal{O}})]. \quad (9)$$

This bi-level problem avoids the previous issue and can be easily decomposed into a two stage optimization problem: (1) Optimize a classifier $\hat{\pi}_{\theta}$ on the *complete data* via standard supervised learning, and then (2) optimize an uncertainty model $\hat{p}_{\psi}(\Phi \mid X_{\mathcal{O}})$ assuming that $\hat{\pi}_{\theta}$ is fixed leveraging Proposition 1. The beauty of this approach is that the optimization problems are completely decoupled. Notably, it is possible to use a large *pretrained* model for $\hat{\pi}_{\theta}$, including a foundation model.

¹In the appendix, we discuss a regularized bi-level problem that cannot be decoupled into two stages.

Theoretic Guarantee for DMV The theoretic version of this bi-level optimization matches the true MVU distributions (proof in appendix), which establishes the theoretic guarantees for DMV.

Proposition 2. Let the non-parametric version of Equation (9) be defined as:

$$q^*, f^* \triangleq \mathop{\arg\min}_{q,f} \mathop{\mathbb{E}}_{X,Y}[\ell_{\mathrm{CE}}(f(X),Y)], \quad \textit{s.t. } q \in \mathop{\arg\min}_{\tilde{q}} \mathop{\mathbb{E}}_{X_{\mathcal{O}},X_{\mathcal{M}}}[-\log \tilde{q}(f(X_{\mathcal{O}},X_{\mathcal{M}}) \mid X_{\mathcal{O}})]\,,$$

where q and f are general non-parametric functions. The optimal solution q^* corresponds to the true MVU distributions, i.e., $q^*(\Phi|X_{\mathcal{O}}) = p(\Phi|X_{\mathcal{O}})$.

Approximation of \hat{p}_{ψ} Our DMV approach could use any approximation for \hat{p}_{ψ} including a Dirichlet distribution, a mixture of Dirichlet distributions, or even any real-valued distribution that is projected onto the probability simplex. For example, if $Z \sim \mathrm{Normal}(\mu, \Sigma)$ comes from a multivariate normal distribution, then $\Phi \triangleq \mathrm{softmax}(Z)$ is a distribution on the simplex. However, because the Dirichlet density is known in closed form, we will use a Dirichlet approximation in our experiments where the strength parameters α are predicted by a function $g_{\psi}(X_{\mathcal{O}})$ given observed values $X_{\mathcal{O}}$, i.e., $\hat{p}_{\psi}(\Phi|X_{\mathcal{O}}) = p_{\mathrm{Dir}}(\Phi|\alpha = g_{\psi}(X_{\mathcal{O}}))$. Importantly, g_{ψ} must be able to handle arbitrary observed features including no observed features or all features. Additionally to be proper strength parameters for a Dirichlet, g_{ψ} must be strictly positive, which we will enforce by simply applying an exponential final activation function to a standard NN.

5 EVALUATING MVU VIA MISSING VALUE CALIBRATION ERROR (MVCE)

As stated in the introduction, the key decision question is: "Is it worth it to collect missing values (for this specific test sample)?" The answer to this question depends heavily on the specific real-world problem context such as the cost of revealing missing values (e.g., doing a biospy is significantly more expensive than collecting blood pressure) and the costs of being wrong (e.g., performing surgery if there is no cancer has high cost). To formally evaluate the decision problem, we would have to specify a decision process with all associated actions, costs/rewards, world environment, etc. Thus, instead of focusing on a particular scenario, we aim for a problem-agnostic evaluation of MVU models by asking an uncertainty calibration question w.r.t. hard confidence: "When my model predicts a hard confidence of c% on a partially observed input $x_{\mathcal{O}}$, is the prediction on the fully observed x the same c% of the time?" This is similar but distinct from the standard uncertainty calibration question corresponding to soft confidence which is: "When my model predicts a soft confidence of c%, does it match the true class c% of the time?" Specifically, the hard confidence calibration gives the probability that the prediction will not change, while the soft confidence calibration gives the probability the class is correct. While soft confidence values can be naturally evaluated using the well-known Expected Calibration Error (ECE) (Naeini et al., 2015; Guo et al., 2017) (see review of ECE in the appendix), evaluating hard confidence values for MVU requires some adaptation.

5.1 MISSING VALUE CALIBRATION ERROR FOR HARD CONFIDENCE EVALUATION

Because hard confidence aims to quantify the probability that the prediction will stay the same, the key idea is that we will simulate the phenomena of revealing all missing values using complete training data. Intuitively, we simulate a partially observed input $x_{\mathcal{O}}$ by dropping features from the full input x and then compare with the prediction on the complete x. Compared to ECE, we are not estimating whether the prediction is accurate but whether the prediction *changed* on average after revealing all missing values. As a reminder, the aim of hard confidences is to help operators know whether collecting missing values would be useful or not. Given this, like ECE, MVCE first partitions the dataset into bins \mathcal{B} based on the *hard* confidence values $\hat{c}_{\text{hard},i}$.

Definition 2 (Missing Value Calibration Error (MVCE)). Given a dataset of labeled pairs, their computed hard predictions and hard confidences, and a partition of the dataset \mathcal{B} , MVCE is defined as: $\text{MVCE} = \sum_{B \in \mathcal{B}} \frac{|B|}{|\mathcal{D}|} |\cos(B) - \bar{c}_{\text{hard}}^{(\bar{\mathcal{O}})}(B)|$, where $\bar{c}_{\text{hard}}^{(\bar{\mathcal{O}})}(B) \triangleq \frac{1}{|B|} \sum_{i \in B} \hat{c}_{\text{hard},i}^{(\mathcal{O})}$ is the average hard confidence on the (simulated) partial input $x_{\mathcal{O},i}$ in bin B and the consistency of bin B is defined as $\cos(B) \triangleq \frac{1}{|B|} \sum_{i \in B} \mathbb{1}(\hat{y}_{\text{hard},i}^{(\mathcal{O})} = \hat{y}_{\text{hard},i})$, where $\hat{y}_{\text{hard},i}^{(\mathcal{O})}$ is the prediction on the partial input $x_{\mathcal{O},i}$ and $\hat{y}_{\text{hard},i}$ is the prediction on the complete input x_i .

The *consistency* term captures the idea of how often the prediction changed when all the missing features were revealed. This is the key difference for evaluating hard confidence values for MVU

while the other parts are similar to ECE. This MVCE definition can also naturally be generalized to the case of cost-sensitive classification defined in Appendix A.1 by using the corresponding cost-sensitive predictions and confidence values. We use a simple partitioning scheme of diving confidence values into a number of equal range bins, which is the common partitioning scheme for ECE.

Relation to Other Metrics Our MVCE metric isolates the evaluation of the hard confidence values for MVU. As such, it does not evaluate the accuracy or standard calibration of the classifier. For this, one can simply use standard metrics like accuracy and ECE to evaluate the performance and calibration of the classifier. This is similar to how a classifier may be accurate but not calibrated or calibrated but not accurate. In practice, we recommend using accuracy, ECE and our MVCE metric so that all aspects of the system can be properly evaluated, but we focus on the evaluation of hard confidence values in this paper.

5.2 POST-HOC CALIBRATION OF MVU VIA MVCE

Similar to prior post-hoc calibration methods, we propose a post-hoc calibration method to improve MVU distribution estimate after training using the MVCE metric. While traditional methods for first-order uncertainty calibration typically adjust the predicted class probabilities (or soft confidences in our context) (Bengs et al., 2022), our goal is to improve the estimate of hard confidences, which depend on the variance. Therefore, we must calibrate the predicted MVU distribution $p(\Phi|x_{\mathcal{O}})$ directly, which will *implicitly* calibrate our confidence values.

Our post-hoc MVU adjustment approach can be viewed as a type of post-processing of the estimated MVU distribution, i.e., $\hat{p}_{\psi,\lambda}(\Phi|x_{\mathcal{O}}) = \Omega(\hat{p}_{\psi}(\Phi|x_{\mathcal{O}}),\lambda)$, where Ω modifies the MVU distribution based on parameters λ , ideally by changing either the variance or entropy without changing the mean. As one example which we will use in experiments, when the MVU is a Dirichlet distribution, i.e., $\hat{p}_{\psi}(\Phi|x_{\mathcal{O}}) = p_{\mathrm{Dir}}(\Phi|\alpha = g_{\psi}(x_{\mathcal{O}}))$, we can simply scale the predicted Dirichlet strengths by a positive scalar λ , i.e., $\hat{p}_{\psi,\lambda} = p_{\mathrm{Dir}}(\Phi|\alpha = \lambda g_{\psi}(x_{\mathcal{O}}))$. Thus, our post-hoc calibration approach can be defined as: $\lambda^*(\psi) = \arg\min_{\lambda} \mathrm{MVCE}(\hat{p}_{\psi,\lambda}(\Phi|x_{\mathcal{O}}))$. While the objective is non-differentiable due to the binning of \mathcal{B} , we can use zero-th order optimization to choose λ such as a simple grid search for low-dimensional λ or Bayesian optimization approaches. Furthermore, if robustness to multiple cost functions is important (Appendix A.1), we could change the objective to an expectation over cost functions. For example, let $W \sim \mathrm{Dirichlet}(\alpha=1)$ be drawn from the uniform distribution over the probability simplex, then the objective to optimize could be $\mathbb{E}_W[\mathrm{MVCE}_W(\hat{p}_{\psi,\lambda}(\Phi|x_{\mathcal{O}})),$ where MVCE_W denotes the cost-sensitive version of MVCE.

6 RELATED WORKS

Missing Values Missing values are features within a sample that are unobserved where having an observed value would be useful for evaluation (Little & Rubin, 2019). We are particularly focused on quantifying the impact of missing values on the prediction. While some work considers incomplete training data, we currently only consider missing values during evaluation, which we assume are missing at random. One approach to handling missing values is imputation (Little & Rubin, 2019; Khosravi et al., 2019), though this lacks a mechanism for estimating uncertainty. Some approaches can produce multiple imputations conditioned on the observation which can be leveraged through Monte Carlo to estimate uncertainty; image inpainting is a notable example for high dimensional spaces (Ma et al., 2018; Zhang et al., 2023; Liu et al., 2023). Other approaches modify the classifier to allow missing inputs, which may even be able to estimate missing value uncertainty directly (Khosravi et al., 2019), though these often impose restrictions on the model architecture or input space. Our goal is to directly handle missing values with uncertainty without such restrictions.

Uncertainty Aleatoric uncertainty is any source of variability outside our model (notably unmeasurable), which we simply model as random behavior. Epistemic uncertainty is a reducible form of uncertainty in the model due to a lack of data, which can be further decomposed based on the type of data; commonly parameter uncertainty is considered (Liu et al., 2019). We focus on missing value uncertainty (MVU). While MVU is reduced as $x_{\mathcal{O}} \to x$, it is not always possible to collect additional features, which may classify it as either aleatoric or epistemic depending on modeling assumptions. Some uncertainty work considers robustness to missing values (Zaffran et al., 2023), though they do not directly report MVU. Some prior works such as Bayesian Inference (Gelman et al., 1995),

Evidential Deep Learning (Sensoy et al., 2018), and other second-order uncertainty approaches (Sale et al., 2023; Bengs et al., 2022) can report prediction uncertainty, though this is limited to aleatoric or other types of epistemic such as parametric. To our knowledge prior work has not investigated both estimating and reporting MVU.

Model Calibration Calibration in machine learning refers to the alignment between a model's predicted probabilities and the true likelihood of outcomes, ensuring that confidence scores accurately reflect real-world chances of correctness. Expected Calibration Error (ECE) is a key metric used to assess the calibration of probabilistic classifiers by quantifying the difference between the predicted probabilities and actual outcomes (Naeini et al., 2015; Guo et al., 2017). Nixon et al. (2019) discuss the shortcomings of the ECE metric and introduces Adaptive Calibration Error (ACE) and Thresholded ACE (TACE) to address its limitations, particularly in multi-class settings. Vaicenavicius et al. (2019) build on this by proposing a general theoretical framework for evaluating the calibration of probabilistic classifiers. Calibration properties of modern neural networks are analyzed by Minderer et al. (2021). However, these calibration works are directed towards directly calibrating the predicted probabilities, rather than other types of uncertainty over the prediction.

7 EXPERIMENTS

First, we aim to justify our missing value calibration error (MVCE) metric through synthetic data, showing that the metric correctly penalizes approaches to estimating missing value uncertainty (MVU) that differ from ground truth. Next, we consider CelebA as a real world dataset with 3 intuitive features, first demonstrating MVCE works in real world data using prior techniques for accurately estimating missing features, and then demonstrating DMV works as a simpler approach to solving the MVU problem. Finally, we use DMV on multiclass setups using MNIST, CIFAR10, and StarcraftCIFAR10, demonstrating it outperforms other simpler approaches to handling missing values including additional heuristics for estimating missing value imputation variance and a classifier that directly handles missing values. Due to the randomness in methods and mask generation, the results subtly varied each run, thus we ran multiple trials of each experiment and average the results; Appendix D contains additional results including standard deviations of the averages.

Validating MVCE Metric With Synthetic **Data** To validate our MVCE metric, we use synthetic data. We generated a set of X values from a bi-variate normal distribution then marked one of the two variables missing and evaluated different generators $p(X_{\mathcal{M}}|x_{\mathcal{O}})$ against a simple linear classifier. The expectation is since the ground truth generator produced the data, it should produce the best estimates of missing features. Table 1 shows results of this experiment, with the metric correctly penalizing generators that were modified to no longer match ground truth. Calibration manages to reduce MVCE for models more distant from ground truth, but notably does not trivially minimize leaving ground truth as the best.

Table 1: With synthetic data, the ground truth model leads to the lowest MVCE as expected. Mutating the generator or using a heuristic for estimating MVU both increase MVCE. Post-hoc calibration can greatly reduce MVCE, though it is not a substitute for improving the model.

Method \ Missing	Uncalibrated	Calibrated	Change
Ground Truth Generator	0.0066	0.0066	0%
Generator Correlation 0.7 → 0.6	0.0116	0.0116	0%
Generator Correlation 0.7 → 0.8	0.0204	0.0111	-46%
Generator Covariance × 0.25	0.0485	0.0082	-83%
Generator Covariance × 4.0	0.0773	0.0071	-91%
Imputation, 0.99 Max Variance	0.1001	0.0994	-1%

Real World Data Using Image Inpainting To test MVCE on real world data, we wish to compare to a well-established baseline for estimating missing values. To do this, we make use of a pretrained CoPaint model (Zhang et al., 2023) that can inpaint CelebA-HQ images (Karras et al., 2017). Since CelebA-HQ lacks attributes for classification, we obtain them through CelebAMask-HQ (Lee et al., 2019). We mask out either the top half or bottom half of the image, then use CoPaint to predict different possibilities for Monte Carlo sampling, (see subsection 4.1). The samples are passed through independently trained ResNet18 classifiers (He et al., 2016) trained to predict three features: Blond Hair (visible in both halves), Eyeglasses (primarily in the top), and smiling (primarily in the bottom). As seen in Table 2, exploring more possibilities by taking more diffusion samples leads to both improved consistency and confidence estimates. However, taking a large number of samples ends up being too slow to practically deploy. To address the performance issue, we make use of DMV

(subsection 4.2), along with some simple imputation baselines with 0, 0.5, and 0.99 scaling of the max variance (subsection 4.1). DMV produced notably faster results than the diffusion approach even at low sample counts, with MVCE comparable to that of the diffusion approach. While in some cases the 0 imputation heuristic could produce similar performance, this usually comes at a tradeoff between minimizing MVCE and maximizing consistency; DMV does not require such a tradeoff.

Table 2: On CelebA data, high number of diffusion samples produce the best result, but this takes too long to use in practice. DMV provides comparable MVCE while running much more efficiently than even single sample baselines, and performs well across all features unlike the 0 imputation approach. Reported times are for a single sample on the blond hair experiment (features ran in comparable times). The 0 variance methods cannot be calibrated as our calibration approach scales the variance.

Method \ Missing	Time	Fe	Feature: Blond Hair			Feature: Smiling				Feature: Eyeglasses			
Wiethou \ Wissing	(seconds)	Cons.	MVCE	Calibra	ted	Cons.	MVCE	Calibra	ated	Cons.	MVCE	Calibr	ated
Diffusion - 1 Sample, 0 Variance	182	0.9478	0.0522	-	-	0.8643	0.1357	-	-	0.9942	0.0058	-	-
Diffusion - 1 Sample, 0.5 Max Variance	182	0.9479	0.0327	0.0309	-5%	0.8643	0.1032	0.0988	-4%	0.9942	0.0045	0.0040	-11%
Diffusion - 1 Sample, 0.99 Max Variance	182	0.9446	0.0429	0.0282	-34%	0.8668	0.0976	0.0968	-1%	0.9936	0.0069	0.0041	-41%
Diffusion - 3 Samples - Empirical Variance	547	0.9494	0.0274	0.0234	-14%	0.8451	0.0858	0.0776	-10%	0.9945	0.0051	0.0046	-11%
Diffusion - 30 Samples -Empirical Variance	5465	0.9651	0.0068	0.0058	-15%	0.8689	0.0403	0.0383	-5%	0.9938	0.0046	0.0044	-4%
0 Imputation, 0 Variance	0.547	0.9173	0.0827	-	-	0.7410	0.2590	-	-	0.7853	0.2147	-	-
0 Imputation, 0.5 Max Variance	0.538	0.9171	0.0634	0.0606	-4%	0.7410	0.2172	0.2122	-2%	0.7853	0.0875	0.0793	-9%
0 Imputation, 0.99 Max Variance	0.619	0.9182	0.0687	0.0555	-19%	0.7576	0.0618	0.0618	0%	0.9078	0.0946	0.0668	-29%
Direct Missing Value (DMV)	0.901	0.9632	0.0287	0.0060	-79%	0.8043	0.0704	0.0683	-3%	0.9946	0.0028	0.0028	0%

Table 3: The Direct Missing Value approach performance well across all three datasets, notably having both high consistency while also having low MVCE. While the missing value robust classifier could produce good consistency, reducing MVCE through uncertainty heuristics produced a tradeoff due to the reduced confidence being applied to all samples. Mean imputation was generally not viable in all cases.

Method \ Missing		MNIST				CIFAR10				StarcraftCIFAR10				
Wiethod \ Wilssing	Cons.	MVCE	Calibrated		Cons.	MVCE	Calibr	Calibrated		MVCE	Calibr	ated		
Mean Imputation, 0 Variance	0.5877	0.4123	-	-	0.3519	0.6481	-	-	0.7811	0.2189	-	-		
Mean Imputation, 0.5 Max Variance	0.5886	0.1870	0.1663	-11%	0.3547	0.3121	0.2855	-9%	0.7789	0.0494	0.0214	-57%		
Mean Imputation, Scale Probability	0.5856	0.3121	0.1880	-40%	0.3536	0.4768	0.3146	-34%	0.7737	0.1023	0.0216	-79%		
Missing Robust, 0 Variance	0.9724	0.0276	-	-	0.8715	0.1285	-	-	0.9158	0.0842	-	-		
Missing Robust, 0.5 Max Variance	0.5555	0.2553	0.2560	0%	0.5786	0.0523	0.0542	4%	0.8708	0.1128	0.0225	-80%		
Missing Robust, Scale Probability	0.8010	0.0194	0.0095	-51%	0.6003	0.0547	0.0336	-39%	0.8744	0.0222	0.0222	0%		
Direct Missing Value	0.9320	0.0729	0.0353	-52%	0.8286	0.0457	0.0457	0%	0.9229	0.0200	0.0114	-43%		

Comparison of DMV to Simple Missing Value Approaches on Multiclass Datasets While DMV managed to produce comparable performance to the diffusion model on CelebA data, the 0-imputation approach was also able to achieve some good results, which is partially obscured by the binary classification task. Thus, to fully show the benefit of DMV, we test it on MNIST (Deng, 2012), CIFAR10 (Krizhevsky et al., 2009) and StarCraftCIFAR10 (Kulinski et al., 2023) giving us 10 classes. For masking, we divided the image into a 4x4 grid of "sensors", and had a random chance for each sensor to be dropped during evaluation. We compare DMV to two baselines for handling missing values: mean imputation, and a classifier robust to missing values. For both approaches, we consider three heuristics for estimating uncertainty as described in subsection 4.1. As shown in Table 3, DMV was able to consistently perform well across all three datasets, while the performance of the other approaches was far less consistent. Mean imputation notably only produced acceptable consistency on Starcraft, regardless of variance approach. Meanwhile, a classifier simply robust to missing values has a tradeoff between good consistency and good MVCE based on the variance approach.

Conclusion In this paper, we addressed the challenge of whether it is worth the cost of collecting missing inputs based on information at inference time. We leveraged the concept of missing value uncertainty (MVU) to produce hard confidence, which can be used by an operator to estimate whether collecting missing features will not change the prediction. We proposed a Direct Missing Value method to estimate MVU, and developed a metric to evaluate and calibrate MVU estimates. Finally, we demonstrated our metric and uncertainty estimation methods work through use of several real-world datasets. Our work is not meant to replace existing approaches to estimating uncertainty, rather we recommend using it alongside existing epistemic and aleatoric uncertainty methods.

REFERENCES

- Anastasios N Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *arXiv* preprint arXiv:2107.07511, 2021.
- 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*, 20(1):40–49, 2011.
- Viktor Bengs, Eyke Hüllermeier, and Willem Waegeman. Pitfalls of epistemic uncertainty quantification through loss minimisation. *Advances in Neural Information Processing Systems*, 35: 29205–29216, 2022.
 - Li Deng. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE signal processing magazine*, 29(6):141–142, 2012.
 - Andrew Gelman, John B Carlin, Hal S Stern, and Donald B Rubin. *Bayesian data analysis*. Chapman and Hall/CRC, 1995.
 - Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
 - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
 - Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.
 - Pasha Khosravi, YooJung Choi, Yitao Liang, Antonio Vergari, and Guy Van den Broeck. On tractable computation of expected predictions. *Advances in Neural Information Processing Systems*, 32, 2019.
 - Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
 - Sean Kulinski, Nicholas R Waytowich, James Z Hare, and David I Inouye. Starcraftimage: A dataset for prototyping spatial reasoning methods for multi-agent environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22004–22013, 2023.
 - Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. corr abs/1907.11922 (2019). *arXiv preprint arXiv:1907.11922*, 2019.
 - Yoad Lewenberg, Yoram Bachrach, Ulrich Paquet, and Jeffrey Rosenschein. Knowing what to ask: A bayesian active learning approach to the surveying problem. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
 - Roderick JA Little and Donald B Rubin. *Statistical analysis with missing data*, volume 793. John Wiley & Sons, 2019.
 - Anji Liu, Mathias Niepert, and Guy Van den Broeck. Image inpainting via tractable steering of diffusion models. *arXiv preprint arXiv:2401.03349*, 2023.
 - Jeremiah Liu, John Paisley, Marianthi-Anna Kioumourtzoglou, and Brent Coull. Accurate uncertainty estimation and decomposition in ensemble learning. *Advances in neural information processing systems*, 32, 2019.
 - Chao Ma, Sebastian Tschiatschek, Konstantina Palla, José Miguel Hernández-Lobato, Sebastian Nowozin, and Cheng Zhang. Eddi: Efficient dynamic discovery of high-value information with partial vae. *arXiv preprint arXiv:1809.11142*, 2018.
- Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Hubis, Xiaohua Zhai, Neil Houlsby, Dustin Tran, and Mario Lucic. Revisiting the calibration of modern neural networks. *Advances in Neural Information Processing Systems*, 34:15682–15694, 2021.

Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated proba-bilities using bayesian binning. In Proceedings of the AAAI conference on artificial intelligence, volume 29, 2015. Jeremy Nixon, Michael W Dusenberry, Linchuan Zhang, Ghassen Jerfel, and Dustin Tran. Measuring calibration in deep learning. In CVPR workshops, volume 2, 2019. Yaniv Romano, Evan Patterson, and Emmanuel Candes. Conformalized quantile regression. Advances in neural information processing systems, 32, 2019. Yusuf Sale, Viktor Bengs, Michele Caprio, and Eyke Hüllermeier. Second-order uncertainty quantifi-cation: A distance-based approach. In Forty-first International Conference on Machine Learning, 2023. Murat Sensoy, Lance Kaplan, and Melih Kandemir. Evidential deep learning to quantify classification uncertainty. Advances in neural information processing systems, 31, 2018. Burr Settles. Active learning literature survey. 2009. Juozas Vaicenavicius, David Widmann, Carl Andersson, Fredrik Lindsten, Jacob Roll, and Thomas Schön. Evaluating model calibration in classification. In The 22nd international conference on artificial intelligence and statistics, pp. 3459–3467. PMLR, 2019. Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. Algorithmic learning in a random world. Springer, 2005. Margaux Zaffran, Aymeric Dieuleveut, Julie Josse, and Yaniv Romano. Conformal prediction with missing values. In International Conference on Machine Learning, pp. 40578–40604. PMLR, 2023. Guanhua Zhang, Jiabao Ji, Yang Zhang, Mo Yu, Tommi S Jaakkola, and Shiyu Chang. Towards coherent image inpainting using denoising diffusion implicit models. The Fortieth International Conference on Machine Learning, 2023.

The supplementary material is organized as follows:

- Appendix A Additional Content: We provide some additional content that fit did not fit
 in the main paper including the cost-sensitive classification generalization of hard voting
 and a review of ECE.
- **Appendix B Proofs:** Additional proofs that were not essential in the main paper.
- Appendix C Extension of Bi-Level DMV Problem with Regularization: An extension of the bi-level optimization approach from Section 4.2.
- Appendix D Additional Experimental Results: Additional experimental results.
 - Appendix D.1 CelebA: Extended results for the CelebA dataset, including reports
 of consistency, MVCE after calibration, standard deviations for experiments, and
 information on the diffusion model.
 - Appendix D.2 Multiclass Datasets: Additional information for the MNIST, CIFAR10 and StarCraftCIFAR10 datasets, notably including standard deviations for experiments.
 - Appendix D.3 Synthetic: Full details on the synthetic dataset and how it was generated.
- Appendix E LLM Usage: Details how LLMs were used in this paper.

A ADDITIONAL CONTENT

A.1 COST-SENSITIVE HARD VOTING CLASSIFICATION RULE

The ensembling classification framework above can generalized to the case where different misclassifications have different costs defined by a cost function: $\ell(y,\hat{y})$ where y is the true label and \hat{y} is the predicted label. For example, in medical tests, false positives may have a much lower cost than false negatives. Thus, we generalize the local Bayes optimal classification rule used to generate each hard vote using a minimum cost classification rule: $y = \arg\min_j \mathbb{E}_{p(Y;\phi)}[\ell(Y,j)]$, where $p(Y;\phi)$ denotes the categorical distribution with parameter ϕ . The Bayes optimal classification rule used in the previous sections can be seen as minimizing the 0-1 loss function: $y = \arg\min_j \mathbb{E}_{p(Y;\phi)}[\ell_{0-1}(Y,j)] = \arg\max_j \phi_j$. For a simple generalization, we can use the cost-sensitive loss that has different costs for misclassification based on the true class: $\ell_w(y,\hat{y}) \triangleq w_y \mathbb{1}(y \neq y')$, where w_y is the cost associated with misclassifying an sample from class y. This yields the following cost-sensitive classification rule: $y = \arg\min_j \mathbb{E}_{p(Y;\phi)}[\ell_w(Y,j)] = \arg\max_j w_j \phi_j$. Intuitively, this changes the classification boundary on the simplex from being in the center (corresponding to a simple argmax) to being off the center depending on w. When used in the hard voting method to determine the votes, this cost-sensitive classification rule will produce a different predictions and confidences:

$$y_{\text{hard}}^{(w)} \triangleq \arg\max_{j} \mathbb{E}_{p(\Phi)}[\text{OneHotArgmax}(w \odot \Phi)]_{j} \equiv \arg\max_{j} p(Y_{\text{vote}}^{(w)} = j), \tag{10}$$

$$\equiv \arg\max_{j} \Pr(\bigcap_{j' \neq j} w_j \Phi_j \ge w_{j'} \Phi_{j'}), \tag{11}$$

$$c_{\text{hard}}^{(w)} \triangleq \max_{j} \mathbb{E}_{p(\Phi)}[\text{OneHotArgmax}(w \odot \Phi)]_{j} \equiv \max_{j} p(Y_{\text{vote}}^{(w)} = j). \tag{12}$$

A.2 BACKGROUND: EXPECTED CALIBRATION ERROR FOR SOFT CONFIDENCE EVALUATION

The basic intuition of soft confidence values is that if the confidence is $c_{\rm soft}$ percent, then the classifier should be correct $c_{\rm soft}$ percent of the time. Ideally, a calibration metric would compare the confidence value to the empirical accuracy conditioned on a specific input. However, given that almost every input is unique, the accuracy is impossible to estimate accurately. Thus, ECE bins samples based on their predicted confidence values and then estimates the empirical accuracy within each bin. The final ECE score is the average over the difference of accuracy and average confidence in each bin. Formally, given a test dataset $\mathcal{D} \equiv \{(x_i,y_i)\}_{i=1}^n$ and the computed the soft predictions $\hat{y}_{\rm soft,i}$ and soft confidence values $\hat{c}_{\rm soft,i}$ for each sample, ECE first partitions the dataset into bins \mathcal{B} based on the soft confidence values $\hat{c}_{\rm soft,i}$. Then, the ECE is defined as: $\mathrm{ECE} = \sum_{B \in \mathcal{B}} \frac{|B|}{|\mathcal{D}|} |\mathrm{acc}(B) - \bar{c}_{\rm soft}(B)|$, where $\mathrm{acc}(B) \triangleq \frac{1}{|B|} \sum_{i \in B} \mathbb{1}(\hat{y}_{\rm soft,i} = y_i)$ is the empirical classification accuracy and $\bar{c}_{\rm soft}(B) \triangleq \frac{1}{|B|} \sum_{i \in B} \hat{c}_{\rm soft,i}$ is the average soft confidence. The fact that ECE compares the average confidence

to accuracy is related to the fact that the soft decision rule and corresponding confidences are based on the Bayes optimal decision rule, which provides the best accuracy among all possible classifiers.

ADDITIONAL RELATED WORKS A.3

This section covers additional works adjacent to ours that are less directly relevant to our methods.

Conformal Prediction Conformal prediction is an uncertainty quantification method that can provide distribution-free, statistical guarantees on predictions for any underlying model (Vovk et al., 2005). Instead of a single prediction, it produces a prediction set (for classification) or interval (for regression) that is guaranteed to contain the true outcome with a user-specified probability (e.g., 90%), without depending on assumptions about underlying data distribution or model correctness (Angelopoulos & Bates, 2021). Recent work has adapted the area to missing values (Zaffran et al., 2023) and high-dimensional settings (Romano et al., 2019; Zaffran et al., 2023). However, these methods aim to keep the uncertainty prediction valid despite missing values instead of analyzing the uncertainty added specifically due to missing values to determine if there is insufficient information.

Active Learning Active learning is a field in machine learning where data is left unlabeled, and the model attempts to determine the most useful additional samples to label (Settles, 2009). A problem in that field of particular relevance to our work is the active surveying problem; survey responses are modeled as a set of questions with incomplete answers, and the model must decide which question is most useful to ask next (Lewenberg et al., 2017; Ma et al., 2018). This setup naturally fits as a missing values problem, though notably work in this area is focused on selecting features based on information in the feature distribution $p(X_{\mathcal{M}}|X_{\mathcal{O}})$ rather than measuring uncertainty in the prediction due to missing features $p(Y|X_{\mathcal{O}})$.

PROOFS В

MISCELLANEOUS PROOF(S)

Proof of the equivalence between 0-1 loss minimization and Bayes optimal classification rule:

$$y = \arg\min_{j} \mathbb{E}_{p(Y;\phi)}[\ell_{0-1}(Y,j)]$$
 (13)

$$y = \underset{j}{\operatorname{arg \, min}} \, \mathbb{E}_{p(Y;\phi)}[\ell_{0-1}(Y,j)]$$

$$= \underset{j}{\operatorname{arg \, min}} \, \mathbb{E}_{p(Y;\phi)}[\mathbb{1}(Y \neq j)]$$

$$(13)$$

$$= \arg\min_{j} \sum_{j' \neq j} \phi_{j'} \tag{15}$$

$$= \underset{j}{\operatorname{arg \, min}} \ 1 - \phi_j$$

$$= \underset{j}{\operatorname{arg \, max}} \ \phi_j \,.$$

$$(16)$$

$$= \underset{j}{\operatorname{arg}} \max_{j} \phi_{j} \,. \tag{17}$$

Proof that cost-sensitive classification is a weighted version of the Bayes optimal classification rule:

$$y = \arg\min_{j} \mathbb{E}_{p(Y;\phi)}[\ell_w(Y,j)]$$
 (18)

$$y = \underset{j}{\operatorname{arg \, min}} \, \mathbb{E}_{p(Y;\phi)}[\ell_w(Y,j)]$$

$$= \underset{j}{\operatorname{arg \, min}} \, \mathbb{E}_{p(Y;\phi)}[w_Y \mathbb{1}(Y \neq j)]$$
(18)

$$= \arg\min_{j} \sum_{j' \neq j} w_{j'} \phi_{j'} \tag{20}$$

$$= \arg\min_{j} \sum_{j'} w_{j'} \phi_{j'} - w_{j} \phi_{j} \tag{21}$$

$$= \underset{j}{\operatorname{arg\,min}} - w_j \phi_j \tag{22}$$

$$= \underset{j}{\operatorname{arg\,max}} \ w_j \phi_j \,. \tag{23}$$

B.2 MINIMIZING DMV OBJECTIVE

We prove that minimizing the DMV objective from Section 4.2 is equivalent to an objective that only requires complete samples, i.e., samples in the training data have all the features

Proof. We can simply use the definition of KL divergence along with LOTUS to derive the result:

$$\underset{p(X_{\mathcal{O}})}{\mathbb{E}} \left[\text{KL}(p_f(\Phi \mid X_{\mathcal{O}}), \hat{p}_{\psi}(\Phi \mid X_{\mathcal{O}})) \right] \tag{24}$$

$$= \underset{p(X_{\mathcal{O}})}{\mathbb{E}} \left[\underset{p_f(\Phi \mid X_{\mathcal{O}})}{\mathbb{E}} \left[-\log \hat{p}_{\psi}(\Phi \mid X_{\mathcal{O}}) \right) \right] \right] + \gamma_f \tag{25}$$

$$= \underset{p(X_{\mathcal{O}})}{\mathbb{E}} \left[\underset{p(X_{\mathcal{M}}|X_{\mathcal{O}})}{\mathbb{E}} \left[-\log \hat{p}_{\psi}(f(X_{\mathcal{O}}, X_{\mathcal{M}}) \mid X_{\mathcal{O}}) \right] \right] + \gamma_f$$
 (26)

$$= \underset{p(X_{\mathcal{O}}, X_{\mathcal{M}})}{\mathbb{E}} \left[-\log \hat{p}_{\psi}(f(X_{\mathcal{O}}, X_{\mathcal{M}})|X_{\mathcal{O}}) \right] + \gamma_f, \tag{27}$$

where (25) is by the definition of KL where $\gamma_f = \mathbb{E}_{p(X_{\mathcal{O}})}[\mathbb{E}_{p_f(\Phi|X_{\mathcal{O}})}[\log p_f(\Phi|X_{\mathcal{O}})]]$, (26) is by the law of the unconscious statistician (LOTUS), and the last is simply by combining the distributions. Importantly, note that γ_f does depend on f, and thus f must be fixed in the optimization problem for γ_f to be a constant.

B.3 OPTIMAL SOLUTION TO NON-PARAMETERIC PROBLEM

We prove the theoretic version of our bi-level optimization from Section 4.2 would result in the true missing value uncertainty

Proof. Since the upper problem is decoupled, this is simply standard cross entropy minimization, which is equivalent to KL divergence minimization between f(X) and p(Y|X). Thus, the non-parametric f(X) if solved perfectly will be equal to p(Y|X), i.e., $p_{f^*}(Y|X) = p(Y|X)$.

Given that $f^*(X) = p(Y|X)$, we then invoke Proposition 1 on the lower level problem:

$$\underset{\tilde{q}}{\operatorname{arg\,min}} \quad \underset{X_{\mathcal{O}}, X_{\mathcal{M}}}{\mathbb{E}} \left[-\log \tilde{q}(f^*(X_{\mathcal{O}}, X_{\mathcal{M}}) \mid X_{\mathcal{O}}) \right] \tag{28}$$

$$= \underset{\tilde{q}}{\operatorname{arg\,min}} \quad \underset{X_{\mathcal{O}}}{\mathbb{E}} \left[\operatorname{KL}(p_{f^*}(\Phi|X_{\mathcal{O}}), \tilde{q}(\Phi|X_{\mathcal{O}})) \right] \tag{29}$$

$$= \arg\min_{\tilde{q}} \ \mathbb{E}_{X_{\mathcal{O}}} [\mathrm{KL}(p(\Phi|X_{\mathcal{O}}), \tilde{q}(\Phi|X_{\mathcal{O}}))] \tag{30}$$

$$= p(\Phi|X_{\mathcal{O}}), \tag{31}$$

where the first is by Proposition 1, the second is by the fact that f^* is optimal, and the last is by the property of KL divergence that it is minimized if and only if the distributions are equal.

C EXTENSION OF BI-LEVEL DMV PROBLEM WITH REGULARIZATION

While in general the two stage approach is simple and elegant, the learned model might not satisfy a natural constraint that the uncertainty model should converge to a Dirac delta if given fully observed features, i.e., if $\mathcal{O} = \mathcal{F}$, then $p_{\theta}(\Phi|X_{\mathcal{O}} = X)$ should converge to a Dirac delta at $\hat{\pi}_{\theta}(X)$. To enforce this natural constraint, we can ensure that the mean of the distribution is equal to $\hat{\pi}_{\theta}(X)$ and that the entropy of the distribution is minimized. Specifically, let us define the following loss:

$$\ell_{\text{reg}}(x, \psi, \theta) := \left(\ell_{\text{KL}}(\hat{\pi}_{\theta}(x), \mathbb{E}_{p_{\psi}}[\Phi | X_{\mathcal{O}} = x]) + H(p_{\psi}(\Phi | X_{\mathcal{O}} = x))\right)$$
(32)

where x is a complete feature instance, $\ell_{\mathrm{KL}}(p,q)$ is the KL divergence between two probability vectors p and q, and H is the standard entropy formulation. Unlike the uncertainty estimation term, this objective function does not depend on a fixed $\hat{\pi}_{\theta}$ and thus can be added to both the upper and lower optimization problems:

$$\min_{\psi,\theta} \underset{p(X,Y)}{\mathbb{E}} \left[\ell_{\text{CE}}(\hat{\pi}_{\theta}(X), Y) + \lambda \ell_{\text{reg}}(X, \psi, \theta) \right]$$
(33)

$$\text{s.t. } \psi \in \arg\min_{\tilde{\psi}} \ \Big(\underset{p(X_{\mathcal{O}}, X_{\mathcal{M}})}{\mathbb{E}} [-\log \hat{p}_{\tilde{\psi}}(\hat{\pi}_{\theta}(X_{\mathcal{O}}, X_{\mathcal{M}}) \mid X_{\mathcal{O}})] + \lambda \ell_{\text{reg}}(X, \psi, \theta)] \Big).$$

Unlike the previous bi-level problem, this problem does not decompose and thus must use more advanced bi-level optimization strategies such as alternating optimization.

D ADDITIONAL EXPERIMENTAL RESULTS

For real world each experiment, we reported two values: MVCE and consistency, as described in Section 3. We used a simple zero-one cost function for estimating confidence, with different mutators based on the dataset. All experiments were run in Ubuntu servers using NVIDIA RTX A5000 GPUs.

MVCE MVCE is our primary metric for comparing methods, reporting how closely the estimate of confidence matches to the likelihood the prediction will change as more information is revealed. Lower MVCE indicates the method is well calibrated and thus gives good confidence estimates. In our experiments, we computed MVCE using 10 bins. When computing MVCE we ran 4 trials of each experiment (or 10 for synthetic) and reported the average of the metric and its standard deviation in order to minimize sources of randomness in the experiment.

Consistency Consistency is the proportion of samples where the prediction under missingness matches the prediction for fully observed data. While ideally consistency would be close to 1, when essential features are missing we expect lower consistency. Consistency in the results below is computed with 1 bin (i.e. consistency for the entire dataset). Like MVCE, we also ran 4 trials of each experiment and reported the average result in the tables below.

Table 4: Reported first-order accuracy and expected calibration loss for all models on clean data. For most models, we had high accuracy on the relevant feature, and typically the DMV model would reduce expected calibration error with minimal change to accuracy. The robust classifiers had more limited training, which sometimes led to a small accuracy drop, though since our work is focused on evaluating consistency and hard confidences this drop is less relevant.

Dataset	Model	Classes	Accuracy	ECE
CelebA - Blond Hair	Classifier	2	93.80%	0.0397
CelebA - Blond Hair	DMV	2	94.45%	0.0046
CelebA - Eyeglasses	Classifier	2	99.25%	0.0053
CelebA - Eyeglasses	DMV	2	99.50%	0.0035
CelebA - Smiling	Classifier	2	91.85%	0.0462
CelebA - Smiling	DMV	2	92.85%	0.0118
MNIST	Classifier	10	99.48%	0.0009
MNIST	Robust	10	79.52%	0.0221
MNIST	DMV	10	98.54%	0.0233
CIFAR10	Classifier	10	71.77%	0.0753
CIFAR10	Robust	10	54.70%	0.0623
CIFAR10	DMV	10	79.99%	0.0561
StarCraftCIFAR10	Classifier	10	79.90%	0.0312
StarCraftCIFAR10	Robust	10	72.57%	0.0401
StarCraftCIFAR10	DMV	10	79.78%	0.0297

Classifiers For all non-synthetic datasets, we used a modified ResNet18 for predictions, replacing the final layer with an appropriately sized layer for the number of classes. For the Direct Missing Value model and the robust to missing values classifier, we additionally replaced the first layer to take 4 channels as an input instead of 3. DMV additionally changed the final activation function from sigmoid (2 classes) or softmax (3 or more classes) to exponential. Both the standard and robust to missing values classifier were trained using cross entropy loss. We used the SVG optimizer for the standard classifier and Adam for the robust to missing values classifier. Training bash files can be found in the in the codebase for more details.

Hyperparameters

For full details on experiment hyperparameters, we have included .json files containing the arguments used to our training and evaluation scripts (including random seeds). See the README for details on locating them, and the scripts folder for example bash files to run the scripts.

Experiment Setup We used a random crop from the larger image size to 224x244 for our model inputs at training time, and a center crop to 224x244 at testing time. For the DMV and the robust to missing values classifier, we replaced the first layer of PyTorch's pretrained ResNet18. Additionally, we replaced the final layer to adjust the class size. All other layers had weights copied over from the pretrained weights found in PyTorch. For our mutator used in both training the DMV and robust classifier and evaluating MVCE, we divided the image into a 4x4 grid of 56x56 pixel regions and gave each region a 50% chance to be removed each time a sample was fetched. For all three datasets, we cached the average image to use for mean imputation. Both mean imputation and the robust to missing values classifier used three different approaches to estimating variance: 0 variance (a single

point prediction), taking half of the maximum possible variance for a Dirichlet distribution, and scaling the predicted probabilities by a factor of 10 (which is multiplied by the calibration constant).

Table 5: Overall, while the diffusion model approach tends to perform the best given enough samples, the direct missing value approach performs nearly as well but at a fraction of the time. This makes it the most practical model to deploy practically. Calibrating applies a notable improvement to all methods but does not altar the model that performs best on average. Note that 0 variance is not calibrated as our calibration approach scales the variance and anything times 0 remains 0.

CelebA - Method \ Missing	Time	Consistency			MVCE			Calibrated MVCE			E
Diffusion - 1 Sample, 0 Variance	182	0.9354	(SD	0.0000)	0.0646	(SD	0.0000)	-		-	-
Diffusion - 1 Sample, 0.5 Max Variance	182	0.9355	(SD	0.0005)	0.0468	(SD	0.0007)	0.0446	(SD	0.0009)	-5%
Diffusion - 1 Sample, 0.99 Max Variance	182	0.9350	(SD	0.0005)	0.0491	(SD	0.0012)	0.0431	(SD	0.0005)	-12%
Diffusion - 3 Samples - Empirical Variance	547	0.9297	(SD	0.0005)	0.0394	(SD	0.0006)	0.0352	(SD	0.0006)	-11%
Diffusion - 30 Samples -Empirical Variance	5465	0.9426	(SD	0.0009)	0.0172	(SD	0.0008)	0.0162	(SD	0.0007)	-6%
0 Imputation, 0 Variance	0.547	0.8145	(SD	0.0000)	0.1855	(SD	0.0000)	-		-	
0 Imputation, 0.5 Max Variance	0.538	0.8145	(SD	0.0008)	0.1227	(SD	0.0007)	0.1174	(SD	0.0011)	-4%
0 Imputation, 0.99 Max Variance	0.619	0.8612	(SD	0.0012)	0.0750	(SD	0.0025)	0.0614	(SD	0.0022)	-18%
Direct Missing Value (DMV)	0.901	0.9207	(SD	0.0013)	0.0340	(SD	0.0015)	0.0257	(SD	0.0016)	-24%

D.1 CELEBA

We made use of the CelebA-HQ dataset (Karras et al., 2017), which contains 10,000 64x64 color images of celebrity faces. We choose this dataset as it was easy to form intuitions about the relationship between missing masks and the CelebA features. Additionally, it was easy to locate pre-trained diffusion models for the dataset. Since CelebA-HQ was designed for training generative models, it lacks features, so we obtained the feature label information from (Lee et al., 2019). Experiments on the CelebA dataset were run on three different target features: Blond Hair, Eyeglasses, and Smiling. A summary of results on the CelebA dataset are shown in Table 5. Tables 6, 9, and 12 demonstrate the consistency achieved by each classifier on the relevant feature. Tables 7, 10, and 13 show a detailed breakdown of MVCE computed on each of the three features. Tables 8, 11, and 14 show a detailed breakdown of the improvements made by calibration on MVCE.

Experiment Setup An independant ResNet18 classifier with pre-trained initial weights was fine-tuned to predict each feature, making the experiments on CelebA all two class: either the feature is present or absent. We did not train the classifier with robustness to missing values. We used the 64x64 pixel versions of the CelebA images as inputs with no rescaling during preprocessing to better match the output resolution of the pretrained diffusion model; taking advantage of the fact ResNet18 allows variable sized inputs. The DMV model was a similarly constructed ResNet18 fine-tuned with a mutator that randomly selected a missing mask between fully observed, fully missing, top half missing, and bottom half missing. At test time, we used a single mask for the entire experiment, either top half missing or bottom half missing. In the tables below, the words "top" or "bottom" always refer to the half missing.

Diffusion For the CelebA dataset, we could take advantage of a pretrained diffusion model to perform experiments. This pretrained diffusion model allowed us to use the Monte Carlo approximation for Missing Value Uncertainty estimation. We could not do the same on other datasets due to the lack of a diffusion model, choosing to forego training more models after realizing the diffusion model approaches are too slow in practice to use. Instead, they serve as a baseline to demonstrate whether DMV is a viable approach.

Sample Cache To reduce the time it takes to run multiple experiments, we cached 30 samples from the diffusion model for each mask on each of the 2000 test samples, along with each of the 2000 validation samples for the sake of calibration. This allowed the time to run each experiment with the diffusion model to be relatively close to that of the non-diffusion approaches at the cost of requiring several months to pre-generate all the samples. The single sample times reported in Table 5 for any diffusion approaches takes the time to compute that number of samples from the cache generation and adds it to the time to run that particular method. This caching approach is likely not representative of how the model would be used when deployed, but we do not believe it had any significant impact on the results beyond a small reduction of randomness when running multiple trials. Leveraging this

cache limited us to just the two masks in our experiments, though this is not a practical limitation to either method as both DMV and the diffusion model can handle any arbitrary missing feature with the right training.

Calibration

We calibrated models by making use of the validation dataset split, ensuring that test data remains unseen. Evaluation of calibration is done on the same test data as the original evaluation of MVCE. For the expectation over cost functions, we made a set of cost functions with 0 cost when y=a,t loss when y=0, a=1, and 1-t cost when y=1, a=0. To perform the expectation, we created a set of t values from 0.1 to 0.9 in 0.1 increments, and then randomly choose a t value in every batch while computing MVCE.

Table 6: The blond hair feature is typically visible in both halves of the image, so regardless of the half that is missing it is not difficult to make consistent predictions with any method.

Blond Hair - Method \ Missing	Top Consistency	Bottom Consistency	Average Consistency
Diffusion - 1 Sample, 0 Variance	0.9335 (SD 0.0000)	0.9620 (SD 0.0000)	0.9478 (SD 0.0000)
Diffusion - 1 Sample, 0.5 Max Variance	0.9341 (SD 0.0005)	0.9618 (SD 0.0012)	0.9479 (SD 0.0008)
Diffusion - 1 Sample, 0.99 Max Variance	0.9352 (SD 0.0005)	0.9540 (SD 0.0007)	0.9446 (SD 0.0006)
Diffusion - 3 Samples - Empirical Variance	0.9277 (SD 0.0012)	0.9711 (SD 0.0005)	0.9494 (SD 0.0008)
Diffusion - 30 Samples -Empirical Variance	0.9572 (SD 0.0013)	0.9729 (SD 0.0005)	0.9651 (SD 0.0009)
0 Imputation, 0 Variance	0.8930 (SD 0.0000)	0.9415 (SD 0.0000)	0.9173 (SD 0.0000)
0 Imputation, 0.5 Max Variance	0.8925 (SD 0.0004)	0.9417 (SD 0.0003)	0.9171 (SD 0.0003)
0 Imputation, 0.99 Max Variance	0.8840 (SD 0.0004)	0.9525 (SD 0.0006)	0.9182 (SD 0.0005)
Direct Missing Value (DMV)	0.9480 (SD 0.0013)	0.9784 (SD 0.0005)	0.9632 (SD 0.0009)

Table 7: Since consistency is high across the board, the high MVCE from many of the baseline methods suggests in most cases they underestimate confidence. DMV notably only performs as good as the three sample diffusion approach for this feature, which is still highly competitive for how quickly it can run.

Blond Hair - Method \ Missing	Top MVCE	Bottom Missing	Average MVCE
Diffusion - 1 Sample, 0 Variance	0.0665 (SD 0.0000)	0.0380 (SD 0.0000)	0.0522 (SD 0.0000)
Diffusion - 1 Sample, 0.5 Max Variance	0.0470 (SD 0.0003)	0.0183 (SD 0.0006)	0.0327 (SD 0.0004)
Diffusion - 1 Sample, 0.99 Max Variance	0.0494 (SD 0.0011)	0.0364 (SD 0.0018)	0.0429 (SD 0.0014)
Diffusion - 3 Samples - Empirical Variance	0.0397 (SD 0.0006)	0.0151 (SD 0.0003)	0.0274 (SD 0.0004)
Diffusion - 30 Samples -Empirical Variance	0.0098 (SD 0.0007)	0.0038 (SD 0.0005)	0.0068 (SD 0.0006)
0 Imputation, 0 Variance	0.1070 (SD 0.0000)	0.0585 (SD 0.0000)	0.0827 (SD 0.0000)
0 Imputation, 0.5 Max Variance	0.0984 (SD 0.0003)	0.0284 (SD 0.0005)	0.0634 (SD 0.0004)
0 Imputation, 0.99 Max Variance	0.1068 (SD 0.0009)	0.0307 (SD 0.0011)	0.0687 (SD 0.0009)
Direct Missing Value (DMV)	0.0318 (SD 0.0014)	0.0257 (SD 0.0004)	0.0287 (SD 0.0010)

Table 8: Calibration overall reduces the MVCE for all methods with the most notable benefit on the DMV model by increasing its overall confidence. This notably brings its performance inline with the diffusion model approach while still performing much quicker.

Blond Hair - Calibrated Method \ Missing	Scale		Top	MVCE		В	otto	m Missir	ng	Δ.	vera	ge MVCE	
Diffusion - 1 Sample, 0.5 Max Variance	0.5	0.0453	(SD	0.0002)	-4%	0.0165	(SD	0.0009) -10%	0.0309	(SD	0.0006)	-5%
Diffusion - 1 Sample, 0.99 Max Variance	5	0.0427	(SD	0.0004)	-14%	0.0138	(SD	0.0008) -62%	0.0282	(SD	0.0006)	-34%
Diffusion - 3 Samples - Empirical Variance	0.1	0.0347	(SD	0.0006)	-13%	0.0121	(SD	0.0006) -20%	0.0234	(SD	0.0005)	-14%
Diffusion - 30 Samples -Empirical Variance	0.25	0.0066	(SD	0.0012)	-33%	0.0050	(SD	0.0007) 30%	0.0058	(SD	0.0009)	-15%
0 Imputation, 0.5 Max Variance	0.25	0.0972	(SD	0.0006)	-1%	0.0240	(SD	0.0006) -15%	0.0606	(SD	0.0006)	-4%
0 Imputation, 0.99 Max Variance	2.5	0.0981	(SD	0.0006)	-8%	0.0129	(SD	0.0003) -58%	0.0555	(SD	0.0004)	-19%
Direct Missing Value (DMV)	7.5	0.0061	(SD	0.0015)	-81%	0.0060	(SD	0.0008) -77%	0.0060	(SD	0.0011)	-79%

Table 9: While the eyeglasses feature is primarily in the top of the image, for many samples in the dataset the feature is partially included in the bottom making it easier to predict the feature without the top half. There is still some notable loss of consistency when the top half is missing for zero imputation as the classifier was likely relying on those pixels for the feature.

Eyeglasses - Method \ Missing	Top Consistency	Bottom Consistency	Average Consistency
Diffusion - 1 Sample, 0 Variance	0.9945 (SD 0.0000)	0.9940 (SD 0.0000)	0.9942 (SD 0.0000)
Diffusion - 1 Sample, 0.5 Max Variance	0.9945 (SD 0.0000)	0.9940 (SD 0.0000)	0.9942 (SD 0.0000)
Diffusion - 1 Sample, 0.99 Max Variance	0.9945 (SD 0.0000)	0.9927 (SD 0.0005)	0.9936 (SD 0.0003)
Diffusion - 3 Samples - Empirical Variance	0.9955 (SD 0.0000)	0.9935 (SD 0.0000)	0.9945 (SD 0.0000)
Diffusion - 30 Samples -Empirical Variance	0.9925 (SD 0.0000)	0.9950 (SD 0.0000)	0.9938 (SD 0.0000)
0 Imputation, 0 Variance	0.7095 (SD 0.0000)	0.8610 (SD 0.0000)	0.7853 (SD 0.0000)
0 Imputation, 0.5 Max Variance	0.7091 (SD 0.0019)	0.8615 (SD 0.0007)	0.7853 (SD 0.0013)
0 Imputation, 0.99 Max Variance	0.8729 (SD 0.0013)	0.9427 (SD 0.0019)	0.9078 (SD 0.0015)
Direct Missing Value (DMV)	0.9952 (SD 0.0003)	0.9940 (SD 0.0000)	0.9946 (SD 0.0002)

Table 10: Since this model overall had high consistency, MVCE was low across all diffusion approaches, with the empirical occasionally getting the best estimate of confidence. The DMV model trained over this feature was also notably high quality, even outperforming the diffusion methods, likely due to more successful hyperparameter tuning for this feature.

Eyeglasses - Method \ Missing	Top MVCE	Bottom Missing	Average MVCE
Diffusion - 1 Sample, 0 Variance	0.0055 (SD 0.0000)	0.0060 (SD 0.0000)	0.0058 (SD 0.0000)
Diffusion - 1 Sample, 0.5 Max Variance	0.0041 (SD 0.0001)	0.0049 (SD 0.0002)	0.0045 (SD 0.0002)
Diffusion - 1 Sample, 0.99 Max Variance	0.0077 (SD 0.0003)	0.0061 (SD 0.0009)	0.0069 (SD 0.0006)
Diffusion - 3 Samples - Empirical Variance	0.0048 (SD 0.0000)	0.0054 (SD 0.0004)	0.0051 (SD 0.0003)
Diffusion - 30 Samples -Empirical Variance	0.0052 (SD 0.0004)	0.0039 (SD 0.0002)	0.0046 (SD 0.0003)
0 Imputation, 0 Variance	0.2905 (SD 0.0000)	0.1390 (SD 0.0000)	0.2147 (SD 0.0000)
0 Imputation, 0.5 Max Variance	0.1346 (SD 0.0016)	0.0404 (SD 0.0002)	0.0875 (SD 0.0010)
0 Imputation, 0.99 Max Variance	0.1178 (SD 0.0036)	0.0714 (SD 0.0021)	0.0946 (SD 0.0027)
Direct Missing Value (DMV)	0.0022 (SD 0.0003)	0.0033 (SD 0.0001)	0.0028 (SD 0.0002)

Table 11: Calibration was able to improve most models, with the largest improvement on several of the high variance models with high scaling values, increasing the confidence (and effectively reducing the variance). The DMV model in this case calibrated to a scale of 1, meaning no change.

Eyeglasses - Calibrated Method \ Missing	Scale		Top	MVCE		В	otto	m Missir	ng		Ave	rage MVCI	:
Diffusion - 1 Sample, 0.5 Max Variance	0.1	0.0033	(SD	0.0000)	-21%	0.0048	(SD	0.0000) -2	% 0.00	40 (S	D 0.0000)	-11%
Diffusion - 1 Sample, 0.99 Max Variance	7.5	0.0035	(SD	0.0003)	-55%	0.0047	(SD	0.0001) -23	% 0.00	41 (S	D 0.0002	-41%
Diffusion - 3 Samples - Empirical Variance	0.1	0.0041	(SD	0.0000)	-14%	0.0050	(SD	0.0000) -8	% 0.00	46 (S	D 0.0000)	-11%
Diffusion - 30 Samples -Empirical Variance	0.25	0.0048	(SD	0.0001)	-9%	0.0041	(SD	0.0003) 3	% 0.00	44 (S	D 0.0002	-4%
0 Imputation, 0.5 Max Variance	0.25	0.1209	(SD	0.0023)	-10%	0.0376	(SD	0.0017) -7	% 0.07	93 (S	D 0.0018)	-9%
0 Imputation, 0.99 Max Variance	2.5	0.0926	(SD	0.0025)	-21%	0.0410	(SD	0.0016) -43	% 0.06	68 (S	D 0.0019)	-29%
Direct Missing Value (DMV)	1.0	0.0022	(SD	0.0003)	0%	0.0033	(SD	0.0001) (% 0.00	28 (S	D 0.0002	0%

Table 12: Smiling is a difficult feature to predict when the bottom half is missing, leading to lower consistency when using that mask; especially in zero-imputation.

Smiling - Method \ Missing	Top Consistency	Bottom Consistency	Average Consistency
Diffusion - 1 Sample, 0 Variance	0.9610 (SD 0.0000)	0.7675 (SD 0.0000)	0.8643 (SD 0.0000)
Diffusion - 1 Sample, 0.5 Max Variance	0.9608 (SD 0.0003)	0.7679 (SD 0.0003)	0.8643 (SD 0.0003)
Diffusion - 1 Sample, 0.99 Max Variance	0.9609 (SD 0.0009)	0.7726 (SD 0.0005)	0.8668 (SD 0.0006)
Diffusion - 3 Samples - Empirical Variance	0.9559 (SD 0.0005)	0.7343 (SD 0.0003)	0.8451 (SD 0.0004)
Diffusion - 30 Samples -Empirical Variance	0.9481 (SD 0.0009)	0.7898 (SD 0.0018)	0.8689 (SD 0.0013)
0 Imputation, 0 Variance	0.9480 (SD 0.0000)	0.5340 (SD 0.0000)	0.7410 (SD 0.0000)
0 Imputation, 0.5 Max Variance	0.9475 (SD 0.0004)	0.5345 (SD 0.0010)	0.7410 (SD 0.0007)
0 Imputation, 0.99 Max Variance	0.9534 (SD 0.0014)	0.5617 (SD 0.0016)	0.7576 (SD 0.0014)
Direct Missing Value (DMV)	0.9836 (SD 0.0009)	0.6250 (SD 0.0033)	0.8043 (SD 0.0023)

Table 13: Due to the amount of information lost when the bottom is missing, maximizing the variance ends up being one of the best non-diffusion approaches. DMV has comparable average performance to this estimate, as it does not overestimate the variance as much when the top half is missing.

Smiling - Method \ Missing	Top MVCE		Botto	om Missing	Average MVCE			
Diffusion - 1 Sample, 0 Variance	0.0390 (SD 0.00	00)	0.2325	(SD 0.0000)	0.1357	(SD 0.0000)		
Diffusion - 1 Sample, 0.5 Max Variance	0.0105 (SD 0.00	17)	0.1960	(SD 0.0007)	0.1032	(SD 0.0012)		
Diffusion - 1 Sample, 0.99 Max Variance	0.0414 (SD 0.00	16)	0.1537	(SD 0.0019)	0.0976	(SD 0.0016)		
Diffusion - 3 Samples - Empirical Variance	0.0227 (SD 0.00	08)	0.1490	(SD 0.0014)	0.0858	(SD 0.0011)		
Diffusion - 30 Samples -Empirical Variance	0.0352 (SD 0.00	07)	0.0455	(SD 0.0018)	0.0403	(SD 0.0013)		
0 Imputation, 0 Variance	0.0520 (SD 0.00	00)	0.4660	(SD 0.0000)	0.2590	(SD 0.0000)		
0 Imputation, 0.5 Max Variance	0.0253 (SD 0.00	02)	0.4091	(SD 0.0010)	0.2172	(SD 0.0006)		
0 Imputation, 0.99 Max Variance	0.0677 (SD 0.00	11)	0.0559	(SD 0.0052)	0.0618	(SD 0.0035)		
Direct Missing Value (DMV)	0.0364 (SD 0.00	10)	0.1043	(SD 0.0038)	0.0704	(SD 0.0026)		

Table 14: Calibration on the smiling feature provided minimal benefits overall, notably not changing the best method for any of the categories.

Smiling - Calibrated Method \ Missing	Scale	Top MVCE				Bottom Missing					Average MVCE			
Diffusion - 1 Sample, 0.5 Max Variance	0.1	0.0074	(SD	0.0023)	-29%	0.1902	(SD	0.0004)	-3%	0.0988	(SD	0.0015)	-4%
Diffusion - 1 Sample, 0.99 Max Variance	2.5	0.0146	(SD	0.0008)	-65%	0.1791	(SD	0.0006)	17%	0.0968	(SD	0.0006)	-1%
Diffusion - 3 Samples - Empirical Variance	0.1	0.0162	(SD	0.0014)	-29%	0.1390	(SD	0.0001)	-7%	0.0776	(SD	0.0009)	-10%
Diffusion - 30 Samples -Empirical Variance	0.75	0.0319	(SD	0.0006)	-9%	0.0448	(SD	0.0012)	-2%	0.0383	(SD	0.0009)	-5%
0 Imputation, 0.5 Max Variance	0.25	0.0213	(SD	0.0004)	-16%	0.4032	(SD	0.0007)	-1%	0.2122	(SD	0.0005)	-2%
0 Imputation, 0.99 Max Variance	1.0	0.0677	(SD	0.0011)	0%	0.0559	(SD	0.0052)	0%	0.0618	(SD	0.0035)	0%
Direct Missing Value (DMV)	2.5	0.0281	(SD	0.0010)	-23%	0.1086	(SD	0.0038)	4%	0.0683	(SD	0.0026)	-3%

D.2 MULTICLASS DATASETS

MNIST

For a simple baseline, we used MNIST (Deng, 2012), which consists of hand-drawn black and white 28x28 digits between 0 and 9. Class labels are easy to interpret, and the single channel makes it easier to learn the models, though it is also a less challenging problem under missing values. Results for the experiments on MNIST are shown in Table 15.

Table 15: MNIST was an easier dataset to predict with missing values, leading to few cases where more information was needed to make good predictions (and this most models were under-confident). Despite this, DMV was still comparable, and could be calibrated to perform nearly the same as the robust classifier. While it was possible to similarly calibrate the robust classifier with uncertainty heuristics, the under confident predictions reduced the uncertainty significantly. Overall,

MNIST Method \ Missing	Consiste	ency		MVC	CE	Scale	le Calibrated MVCE		
Mean Imputation, 0 Variance	0.5877 (SD	0.0041)	0.4123	(SD	0.0041)	-		-	-
Mean Imputation, 0.5 Max Variance	0.5886 (SD	0.0076)	0.1870	(SD	0.0074)	0.1	0.1663	(SD 0.0025)	-11%
Mean Imputation, Scale Probability	0.5856 (SD	0.0077)	0.3121	(SD	0.0066)	0.1	0.1880	(SD 0.0034)	-40%
Missing Robust, 0 Variance	0.9724 (SD	0.0013)	0.0276	(SD	0.0013)	-		-	-
Missing Robust, 0.5 Max Variance	0.5555 (SD	0.0026)	0.2553	(SD	0.0024)	0.3	0.2560	(SD 0.0057)	0%
Missing Robust, Scale Probability	0.8010 (SD	0.0023)	0.0194	(SD	0.0021)	0.1	0.0095	(SD 0.0015)	-51%
Direct Missing Value (DMV)	0.9320 (SD	0.0013)	0.0729	(SD	0.0027)	10.0	0.0353	(SD 0.0016)	-52%

CIFAR10 (Krizhevsky et al., 2009) is a well known baseline in Machine Learning, containing 60,000 32x32 color images of airplanes, automobiles, birds, cats, deer, dogs, frogs, horses, ships, and trucks. Like CelebA, the class labels are easy to interpret, though there is typically a much less obvious relationship between particular regions in the image and the prediction making CIFAR10 a good intermediate difficulty experiment. Additionally, the multi-class setup provides some difficulties that were not seen in single class. Results for the experiments on CIFAR10 are shown in Table 16.

Table 16: The DMV method and the robust classifier both have comparable consistency, though DMV notably produces better confidence estimates as its able to change confidence with respect to each sample. Reducing MVCE on the robust classifier through uncertainty heuristics leads to a notable drop in consitency making it not viable in practice. Overall, this dataset shows the value of DMV for estimating uncertainty that indicates the prediction likely shifted.

CIFAR10 Method \ Missing	Cor	ısist	ency		MV	CE	Scale Calibrated MVCE			E	
Mean Imputation, 0 Variance	0.3519	(SD	0.0022)	0.6481	(SD	0.0022)	-		-		-
Mean Imputation, 0.5 Max Variance	0.3547	(SD	0.0024)	0.3121	(SD	0.0013)	0.1	0.2855	(SD	0.0032)	-9%
Mean Imputation, Scale Probability	0.3536	(SD	0.0044)	0.4768	(SD	0.0048)	0.1	0.3146	(SD	0.0042)	-34%
Missing Robust, 0 Variance	0.8715	(SD	0.0045)	0.1285	(SD	0.0045)	-		-		-
Missing Robust, 0.5 Max Variance	0.5786	(SD	0.0062)	0.0523	(SD	0.0017)	0.5	0.0542	(SD	0.0017)	4%
Missing Robust, Scale Probability	0.6003	(SD	0.0027)	0.0547	(SD	0.0033)	0.1	0.0336	(SD	0.0016)	-39%
Direct Missing Value (DMV)	0.8286	(SD	0.0014)	0.0457	(SD	0.0051)	1.0	0.0457	(SD	0.0051)	0%

StarCraftCIFAR10

StarCraftCIFAR10 (Kulinski et al., 2023) is a dataset meant to simulate a battlefield scenario with data created from replays of the game StarCraft II. The dataset follows the same format as CIFAR10, though the classes are replaced with 5 maps and a time of game (either beginning or end). The map part of the class is easily interpretable for those familiar with the game, while the time of game tends to be more difficult for a human to identify making this a more difficult dataset. It was chosen partly for the similarity in format to CIFAR10, and partly to expand upon our sensor network motivation. Results for the experiments on StarCraftCIFAR10 are shown in Table 17.

Experiment Setup

We fine-tuned a modified pretrained ResNet18 model for both classifiers and DMV for both CIFAR10 and StarCraftCIFAR10. Since MNIST are only a single channel, we modified the ResNet18 model to

Table 17: The map prediction for StarCraftCIFAR10 makes it easier to achieve high levels of consistency even using simple heuristics. The time of game prediction however is more difficult than either of the prior prediction tasks, which requires a wholistic approach to missing values to fully handle well.

StarCraft Method \ Missing	Co	nsist	ency		ΜV	CE	Scale	Calibrated MVCE			
Mean Imputation, 0 Variance	0.7811	(SD	0.0040)	0.2189	(SD	0.0040)	-		-		-
Mean Imputation, 0.5 Max Variance	0.7789	(SD	0.0017)	0.0494	(SD	0.0023)	2.5	0.0214	(SD	0.0038)	-57%
Mean Imputation, Scale Probability	0.7737	(SD	0.0051)	0.1023	(SD	0.0027)	0.25	0.0216	(SD	0.0014)	-79%
Missing Robust, 0 Variance	0.9158	(SD	0.0009)	0.0842	(SD	0.0009)	-		-		-
Missing Robust, 0.5 Max Variance	0.8708	(SD	0.0026)	0.1128	(SD	0.0026)	10.0	0.0225	(SD	0.0030)	-80%
Missing Robust, Scale Probability	0.8744	(SD	0.0030)	0.0222	(SD	0.0031)	1.0	0.0222	(SD	0.0031)	0%
Direct Missing Value (DMV)	0.9229	(SD	0.0014)	0.0200	(SD	0.0019)	0.5	0.0114	(SD	0.0012)	-43%

use a single channel input for the standard classifier, and a 2 channel input for the DMV and robust classifiers. Other datasets used 3 channels for standard classifiers and 4 for the DMV and robust classifier. We rescaled the images to 224x244 during preprocessing to better match the expected input size for ResNet18. For our mutator used in both training the DMV and evaluating MVCE, we divided the image into a 4x4 grid of 56x56 pixel regions and gave each region a 50% chance to be removed each time a sample was fetched. As we lacked a pretrained diffusion model this datasets and it was too slow to use in practice, we skipped all diffusion related methods for experiments on both datasets.

Calibration

We calibrated models by making use of the validation dataset split, ensuring that test data remains unseen. Evaluation of calibration is done on the same test data as the original evaluation of MVCE. Since the cost function becomes immensely more complex with more than 2 variables, we simply calibrated the model on a single zero-one cost function. This means CIFAR10 and StarCraftCIFAR10 both likely get slightly better results from calibration as we know the testing environment.

D.3 SYNTHETIC

For the sake of validating missing value calibration error along with our post-hoc calibration, we created a simple Gaussian Distribution with two variables, X_1 and X_2 , using $\mathbb{E}[X]=(0,0)$, $\mathrm{Var}(X_1)=0.3$, $\mathrm{Var}(X_2)=1$, and $\mathrm{corr}(X_1,X_2)=0.7$. This distribution is visualized in Figure 1. From this setup, we generated 10000 samples and treated those as our ground truth testing dataset. In addition, we constructed a simple classifier as $\mathrm{sigmoid}(x_1+x_2-1)$ which was used for all experiments. Since labels are not required to compute missing value calibration error, we did not need to generate ground truth labels. The results of these experiments are shown in Table 18 and Table 19.

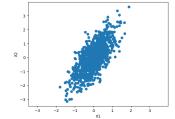


Figure 1: Example distribution of random samples from the synthetic dataset generated to test Missing Value Prediction Uncertainty.

Experiment Setup

We started with the ground truth generator, which we expect minimize MVCE as it can produce confidence estimates consistent with the ground truth data. For comparison, we mutated the generator by increasing or decreasing the correlation between the two variables, and scaling up or scaling down the entire covariance matrix. In

addition, we compared against a single sample imputation, which is comparable to taking a single sample from the diffusion model on the CelebA dataset.

Calibration In order to verify post-hoc calibration worked, we first calibrated the synthetic dataset. Like the CelebA method, we calibrated using a range of values. To simulate the testing environment, we calibrated using a second set of 10000 samples as our calibration dataset. For the expectation over cost functions, we made a set of cost functions with 0 cost when y = a, t loss when y = 0, a = 1, and 1 - t cost when y = 1, a = 0. To perform the expectation, we created a set of t values from 0.1 to 0.9 in 0.1 increments, and then randomly choose a t value in every batch while computing MVCE.

Table 18: As expected, the ground truth generator minimizes MVCE. Generators that are closer to ground truth such as a small change to correlation produce comparabe but still higher MVCE. Models further away such as a full covariance scale further increase MVCE. Heuristics such as imputation lead to the worst results.

Method \ Missing	Uncalibrated										
Wethou \ Whissing	X1	X2	Average								
Ground Truth Generator	0.0052 (SD 0.0005)	0.0081 (SD 0.0015)	0.0066 (SD 0.0011)								
Generator Correlation 0.7 → 0.6	0.0066 (SD 0.0009)	0.0166 (SD 0.0010)	0.0116 (SD 0.0009)								
Generator Correlation 0.7 → 0.8	0.0167 (SD 0.0007)	0.0241 (SD 0.0008)	0.0204 (SD 0.0008)								
Generator Covariance × 0.25	0.0369 (SD 0.0007)	0.0601 (SD 0.0014)	0.0485 (SD 0.0011)								
Generator Covariance × 4.0	0.0566 (SD 0.0009)	0.0981 (SD 0.0021)	0.0773 (SD 0.0015)								
1 Sample Imputation, 0.99 Max Variance	0.1393 (SD 0.0019)	0.0609 (SD 0.0034)	0.1001 (SD 0.0027)								

Table 19: Calibration brings the scaled covariance approaches much more inline with the ground truth, effectively reversing the scale. This ends up being sufficient to reduce the MVCE of one method to lower than the ground truth. It is interesting to note that the model with 0.25x covariance (effectively increasing confidence) has a scale that decreases confidence, and vice versa for the model with 4x covariance. It is also notable that the ground truth generator had a scaling constant of 1, suggesting it is already calibrated.

Method \ Missing		Ca	Calib	ration	Post-Hoc							
Wethou / Wissing	X1	X1			X2			age	X1	X2	Average	Calibration
Ground Truth Generator	0.0052 (SD	0.0006)	0.0081	(SD	0.0015)	0.0066	(SD	0.0011)	0%	0%	0%	1
Generator Correlation 0.7 → 0.6	0.0066 (SD	0.0008)	0.0166	(SD	0.0007)	0.0116	(SD	0.0008)	0%	0%	0%	1
Generator Correlation 0.7 → 0.8	0.0089 (SD	0.0005)	0.0133	(SD	0.0011)	0.0111	(SD	0.0008)	-47%	-45%	-46%	0.75
Generator Covariance × 0.25	0.0056 (SD	0.0007)	0.0107	(SD	0.0009)	0.0082	(SD	0.0008)	-85%	-82%	-83%	0.25
Generator Covariance × 4.0	0.0090 (SD	0.0005)	0.0053	(SD	0.0018)	0.0071	(SD	0.0013)	-84%	-95%	-91%	5
1 Sample Imputation, 0.99 Max Variance	0.1383 (SD	0.0023)	0.0604	(SD	0.0019)	0.0994	(SD	0.0021)	-1%	-1%	-1%	2.5

E LLM USAGE

We made use of LLMs to assist in revising paper contents such as editing, suggestions for notation, and getting initial feedback on ideas. Additionally, LLMs were used to help locate relevant related works. No section of the paper was written fully by LLMs, nor were they involved in any notable capacity in producing the code used for experiments.