

MISSING VALUE UNCERTAINTY: WILL COLLECTING MISSING VALUES CHANGE THE PREDICTION?

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 if [collecting missing values is likely to change the prediction](#). Existing methods typically focus on imputing missing values or quantifying model uncertainty, but they do not directly assess the stability of a prediction if missing values 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. [We start to answer this question by asking whether collecting missing values is likely to change the prediction](#). 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 of [quantifying when more information may change the prediction](#). Missing value literature, for instance, primarily focuses on developing methods like imputation to handle missing values 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 values, 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 values 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 collecting missing values may be beneficial*. 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 values. 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. *Note that while typically the term "voting" references an ensemble setup, for our work we will apply it instead to a distribution of predictions, which can be either sampled or used directly to compute the desired queries.*

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 $[p(Y = 1), p(Y = 2), \dots, p(Y = k)]$ for all k classes, 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:

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

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 $p(\Phi)$. *When properly calibrated*, 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_{soft} is simply the probability that the selected class is correct. When given to a human operator, they could use c_{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:

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

$$\equiv \arg \max_j \Pr(\bigcap_{j' \neq j} \Phi_j \geq \Phi_{j'}), \quad (3)$$

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

where $\text{OneHotArgmax}(\phi) \triangleq \text{OneHot}(\arg \max_j \phi_j)$ is the one-hot encoding of the argmax function and $Y_{\text{vote}} \triangleq \arg \max_j \Phi_j$ is a random variable corresponding to the hard vote of a single Φ (where each Φ intuitively represents a model in the ensemble). Note that the objective is equivalent to the cumulative probability of the distribution in a region defined by the inequalities $\bigcap_{j' \neq j} \Phi_j \geq \Phi_{j'}$. Another interpretation is that each model performs the Bayes optimal classification *locally* using its own belief and then we take an average over these local classifications. Because hard voting only considers the index of the largest probability, it ignores the weightings and thus is naturally more robust to miscalibrated model predicted probabilities. We also generalize this to the case of *cost-sensitive classification* in the appendix.

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 even 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 class prediction as likely correct**, while hard confidence informs the decision to **collect more information** to reduce epistemic uncertainty. 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 ($n \rightarrow \infty$). In contrast, MVU is an orthogonal form of epistemic uncertainty induced by missing values for a single test-time instance. It is reducible not by collecting more training data, but by observing the missing values 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*

162 Φ as:

$$163 \pi(x) \triangleq p(Y|X = x), \quad \Phi \triangleq \pi(X) \quad (5)$$

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

$$168 p(\Phi|X_{\mathcal{O}} = x_{\mathcal{O}}) \equiv \pi_{\sharp}^{\mathcal{O}} p(X_{\mathcal{M}}|X_{\mathcal{O}} = x_{\mathcal{O}}), \quad (6)$$

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

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

178 **Soft and Hard Voting for MVU** Applying the voting rules from Section 2 to the MVU distribution
 179 provides distinct types of information. Soft voting corresponds to the standard Bayes optimal
 180 classification given the observed features: $\arg \max_j p(Y = j|X_{\mathcal{O}} = x_{\mathcal{O}})$. As previously discussed,
 181 the resulting soft confidence does not inform an operator whether collecting more inputs would be
 182 useful. Hard voting, however, directly addresses this problem. The hard confidence derived from
 183 $\arg \max_j p(Y_{\text{hard}} = j|X_{\mathcal{O}} = x_{\mathcal{O}})$ represents the probability that the prediction will *not* change if
 184 the missing values are revealed. A high hard confidence therefore suggests that the decision is stable
 185 and collecting more data is unnecessary, a critical insight for operators in fields like medical diagnosis
 186 or sensor networks where acquiring more information is costly.

191 4 METHODS FOR MVU ESTIMATION

192 Having established hard voting confidence in Section 2 and defining missing value uncertainty in
 193 Section 3, we now consider ways to estimate the MVU distribution $p(\Phi|X_{\mathcal{O}})$ for a distribution of
 194 $X_{\mathcal{O}}$. [To the author’s best knowledge](#), there are no prior methods that focus on estimating MVU, thus
 195 we first propose two natural baselines based on missing value imputation and Monte Carlo estimates
 196 using generative models. After establishing these baselines, we then introduce our novel Direct
 197 Missing Value (DMV) approach for directly estimating the MVU distribution without requiring
 198 imputation or sampling.

201 4.1 BASELINE MVU METHODS

202 **Imputation with Simple Variance** Leveraging prior work, we can handle missing values using
 203 a simple imputation approach such as zero imputation or mean imputation (Little & Rubin, 2019).
 204 However, imputation approaches are limited by the classifier’s inability to estimate missing value
 205 uncertainty, so our best bet is a simple heuristic to estimate variance then use method of moments
 206 to estimate distribution parameters. A naive approach would be to choose some constant variance
 207 regardless of the sample, though with any constant other than 0 (which induces a degenerate distribu-
 208 tion), we run the risk of that uncertainty being too large for the predicted mean (as distributions over
 209 probabilities tend to have restrictions on maximum variance). A better heuristic [would be to create a](#)
 210 [Dirichlet distribution by simply scaling](#) the predicted probability ϕ by some constant to produce α ;
 211 this approach guarantees any positive scaling constant will produce a valid distribution, though it may
 212 be unintuitive to set the constant. We can make the scaling constant more intuitive by instead scaling
 213 the maximum variance - Dirichlet variance cannot be greater than or equal to $\phi \cdot (1 - \phi)$ without
 214 producing non-positive α values. Leveraging this knowledge, we can use a variance of $\phi \cdot (1 - \phi) \cdot s$
 215 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 features that are missing; 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 values, 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 $p(X_{\mathcal{M}}|x_{\mathcal{O}})$, 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, [making it impractical for real-time predictions](#). 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: $\arg \min_{\psi} \mathbb{E}_{X_{\mathcal{O}}}[\text{KL}(p(\Phi|X_{\mathcal{O}}), \hat{p}_{\psi}(\Phi|X_{\mathcal{O}}))] \equiv \arg \min_{\psi} \mathbb{E}_{X_{\mathcal{O}}}[\mathbb{E}_{\Phi|X_{\mathcal{O}}}[-\log \hat{p}_{\psi}(\Phi|X_{\mathcal{O}})]]$, 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 :*

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

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} :

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

where ℓ_{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 $p(\Phi | X_{\mathcal{O}})$. 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:

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

This bi-level problem avoids the previous issue and can be easily decomposed into a two stage optimization problem:¹ (1) Optimize a classifier $\hat{\pi}_{\theta}$ on the *complete data* via standard supervised

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

learning, and then (2) optimize an uncertainty model $\hat{p}_\psi(\Phi | X_{\mathcal{O}})$ 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.

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:*

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

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 , provided it can both be sampled (to compute MVU and MVCE) and evaluate the log-likelihood of ϕ (for minimizing our loss function during training). In our experiments, we leverage the Dirichlet distribution as it has a known closed form density parameterized on strength parameters α . We directly estimate α by a learnable function $g_\psi(X_{\mathcal{O}})$ given observed values $X_{\mathcal{O}}$, i.e., $\hat{p}_\psi(\Phi | X_{\mathcal{O}}) = p_{\text{Dir}}(\Phi | \alpha = g_\psi(X_{\mathcal{O}}))$. g_ψ can use any standard classification architecture, simply requiring switching the final activation function from producing probability values ϕ to strictly positive strengths α (e.g. if the final activation function was softmax, we can instead use exponential). We additionally need our architecture to handle both partially and fully observed inputs; for our experiments we augment the input with the mask information (making the corresponding change to the architecture input dimensions), then zero missing values. For example, if the inputs are 3 channel RGB images, then we can add a 4th mask channel, which contains 1s for missing values and 0s otherwise. Finally, we train this model using the DMV bi-level optimization problem described above. Note these specific approaches (Dirichlet approximation and mask channel) are not required to use the DMV bi-level optimization, other approaches may be used to handle missing inputs and to produce a distribution as an output without losing compatibility with the rest of our work.

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 biopsy 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: $\text{MVCE} = \sum_{B \in \mathcal{B}} \frac{|B|}{|\mathcal{D}|} |\text{cons}(B) - \bar{c}_{\text{hard}}^{(\odot)}(B)|$, where $\bar{c}_{\text{hard}}^{(\odot)}(B) \triangleq \frac{1}{|B|} \sum_{i \in B} \hat{c}_{\text{hard},i}^{(\odot)}$ 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 $\text{cons}(B) \triangleq \frac{1}{|B|} \sum_{i \in B} \mathbb{1}(\hat{y}_i^{(\odot)} = \hat{y}_i)$, where $\hat{y}_i^{(\odot)}$ is the prediction on the partial input $x_{\mathcal{O},i}$ and \hat{y}_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 values were revealed. This is the key difference for evaluating hard confidence values for MVU while the other parts are similar to ECE. Note that since consistency does not consider the true label, it will not evaluate whether the model is accurate; for best results MVCE should be employed alongside traditional accuracy and ECE to fully evaluate a model. Unlike confidence $\bar{c}_{\text{hard}}^{(\odot)}(B)$, consistency $\text{cons}(B)$ can be computed using either y_{soft} and y_{hard} ; for the fully observed prediction they are interchangeable. For our experiments, we use y_{soft} as it is notably less impacted by variance than y_{hard} leading to less noise in the results. Additionally, it makes more sense to quantify whether y_{soft} may shift as that is prediction approach an operator is more likely to use. 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., $\hat{p}_{\psi}(\Phi|x_{\mathcal{O}}) = p_{\text{Dir}}(\Phi|\alpha = g_{\psi}(x_{\mathcal{O}}))$, we can simply scale the predicted Dirichlet strengths by a positive scalar λ , i.e., $\hat{p}_{\psi,\lambda} = p_{\text{Dir}}(\Phi|\alpha = \lambda g_{\psi}(x_{\mathcal{O}}))$. **This scaling constant helps in the case where the model on average predicts too much or too little uncertainty; other cases would require a more advanced approach which is left to future work.** Thus, our post-hoc calibration approach can be defined as: $\lambda^*(\psi) = \arg \min_{\lambda} \text{MVCE}(\hat{p}_{\psi,\lambda}(\Phi|x_{\mathcal{O}}))$. 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 $W \sim \text{Dirichlet}(\alpha = 1)$ be drawn from the uniform distribution over the probability simplex, then the objective to optimize could be $\mathbb{E}_W[\text{MVCE}_W(\hat{p}_{\psi,\lambda}(\Phi|x_{\mathcal{O}}))]$, 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

378 focused on quantifying the impact of missing values on the prediction. While some work considers
 379 incomplete training data, we currently only consider missing values during evaluation, which we
 380 assume are [missing completely at random though this assumption is not required to use for our](#)
 381 [work](#). [Our evaluation considers both MCAR and missing not at random setups](#). One approach to
 382 handling missing values is imputation (Little & Rubin, 2019; Khosravi et al., 2019), though this
 383 lacks a mechanism for estimating uncertainty. Some approaches can produce multiple imputations
 384 conditioned on the observation which can be leveraged through Monte Carlo to estimate uncertainty;
 385 image inpainting is a notable example for high dimensional spaces (Ma et al., 2018; Zhang et al.,
 386 2023; Liu et al., 2023). Other approaches modify the classifier to allow missing inputs, which may
 387 even be able to estimate missing value uncertainty directly (Khosravi et al., 2019), though these often
 388 impose restrictions on the model architecture or input space. Our goal is to directly handle missing
 389 values with uncertainty without such restrictions.

390 **Uncertainty** Aleatoric uncertainty is any source of variability outside our model (notably unmea-
 391 surable), which we simply model as random behavior. Epistemic uncertainty is a reducible form of
 392 uncertainty in the model due to a lack of data, which can be further decomposed based on the type of
 393 data; commonly parameter uncertainty is considered (Liu et al., 2019). We focus on missing value
 394 uncertainty (MVU). While MVU is reduced as $x_{\mathcal{O}} \rightarrow x$, it is not always possible to collect additional
 395 features, which may classify it as either aleatoric or epistemic depending on modeling assumptions.
 396 Some uncertainty work considers robustness to missing values (Zaffran et al., 2023), though they
 397 do not directly report MVU. Some prior works such as Bayesian Inference (Gelman et al., 1995),
 398 Evidential Deep Learning (Sensoy et al., 2018), and other second-order uncertainty approaches (Sale
 399 et al., 2023; Bengs et al., 2022) can report prediction uncertainty, though this is limited to aleatoric or
 400 other types of epistemic such as parametric. To our knowledge prior work has not investigated both
 401 estimating and reporting MVU.

402 **Model Calibration** Calibration in machine learning refers to the alignment between a model’s
 403 predicted probabilities and the true likelihood of outcomes, ensuring that confidence scores accurately
 404 reflect real-world chances of correctness. Expected Calibration Error (ECE) is a key metric used to
 405 assess the calibration of probabilistic classifiers by quantifying the difference between the predicted
 406 probabilities and actual outcomes (Naeini et al., 2015; Guo et al., 2017). Nixon et al. (2019) discuss the
 407 shortcomings of the ECE metric and introduces Adaptive Calibration Error (ACE) and Thresholded
 408 ACE (TACE) to address its limitations, particularly in multi-class settings. Vaicenavicius et al.
 409 (2019) build on this by proposing a general theoretical framework for evaluating the calibration of
 410 probabilistic classifiers. Calibration properties of modern neural networks are analyzed by Minderer
 411 et al. (2021). However, these calibration works are directed towards directly calibrating the predicted
 412 probabilities, rather than other types of uncertainty over the prediction.

413 7 EXPERIMENTS

414
 415
 416 First, we aim to validate our missing value calibration error (MVCE) metric through synthetic data,
 417 showing that the metric correctly penalizes approaches to estimating MVU that differ from known
 418 ground truth uncertainty (which is unknown for real datasets). Next, we compare MVU baselines,
 419 DMV and post-hoc calibration on the real-world CelebA dataset for three binary classification
 420 problems (i.e., blonde hair, eyeglasses, and smiling). We show that DMV has comparable MVCE to
 421 expensive diffusion sampling baselines while being 100x or more faster. [Additionally, we show that](#)
 422 [the best methods do not trivially minimize MVCE by comparing their accuracy using the prediction](#)
 423 [on mutated data \(M. Acc.\) to their accuracy using the prediction on clean data \(C. Acc.\)](#). Finally, we
 424 evaluate DMV compared to efficient baselines on multiclass classification using MNIST, CIFAR10,
 425 and StarcraftCIFAR10, demonstrating it outperforms other approaches in estimating MVU. [We](#)
 426 [implement DMV using a ResNet18 classifier, with the input layer augmented to include a mask](#)
 427 [channel and the final activation function swapped to output Dirichlet weights; a similar ResNet18](#)
 428 [classifier is used for baselines and the robust to missingness classifier \(latter uses the input layer](#)
 429 [augmentation only\)](#). For most methods, our goal was to estimate the parameters of a Dirichlet
 430 [distribution, which we then used with hard voting \(subsection 2.2\) via 1000 sample Monte Carlo to](#)
 431 [compute confidences for MVCE](#). Due to the randomness in [sampling](#) and mask generation, we ran
 multiple trials of each experiment and average the results. Appendix D contains additional details on
 experiments along with additional results including standard deviations.

Validating MVCE Metric With Synthetic Data To briefly validate our MVCE metric, we create a synthetic dataset where X is sampled from a bivariate normal distribution and $p(Y = 1|x) = \sigma(\theta^T x + \theta_0)$ is a simple linear classifier defined by θ and θ_0 . We randomly set one of the inputs to be missing. *Method Details:* The ground-truth MVU can be estimated via Monte Carlo by using the fact that the MVU distribution is simply the pushforward of the missing conditional distribution $p(X_{\mathcal{M}}|x_{\mathcal{O}})$ (see Equation (6)). We compare the ground-truth missing conditional distribution with perturbations of this distribution to check that the MVCE increases if the MVU distribution shifts. *Results:* Table 1 shows results of this experiment, with the metric correctly penalizing generators with higher or lower variances as expected. Calibration manages to reduce MVCE for models more distant from ground truth, but is not a substitute for improving the generator.

Comparing Diffusion-Based, Scaled Max Variance and DMV Methods on CelebA Dataset We now compare MVU methods on the real-world CelebA-HQ dataset (Karras et al., 2017). Since CelebA-HQ lacks attributes for classification, we obtain them through CelebAMask-HQ (Lee et al., 2019).

The task is to predict three semantic attributes which are mostly on the top or bottom: **Blond Hair** (visible primarily in the top), **Eyeglasses** (partially visible in both halves), and **smiling** (primarily in the bottom). We mask out either the top half or bottom half of the image, making this experiment MNAR during evaluation. This allows us to intuitively understand when the MVU should be high vs low, e.g., if the bottom is masked, then smiling should be relatively uncertain compared to the top being masked.

Method Details: For the Monte Carlo baseline, 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 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 the three semantic attributes. We then use the prediction samples to estimate the parameters of a Dirichlet distribution. For computationally efficient methods, we use DMV (subsection 4.2) with ResNet18 architecture slightly modified to allow for a channel corresponding to the missing mask to directly estimate the Dirichlet parameters. We train DMV using a randomly selected masks from top missing, bottom missing, both missing, or neither missing, making it MCAR during training time. We also use a simple mean imputation baselines with three methods for estimating variance: 0, 0.5 max variance, and directly scaling the probability values by 10 (subsection 4.1).

Results: As seen in Table 2, exploring more possibilities by taking more diffusion samples leads to both improved mutated accuracy and confidence estimates. However, taking a large number of samples ends up being too slow to practically deploy. 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 mean imputation heuristic could produce similar performance, this usually comes at a tradeoff between minimizing MVCE and maximizing accuracy under missing values; DMV does not require such a tradeoff.

Comparison of DMV to Computationally Efficient Missing Value Approaches on Multiclass Datasets While DMV managed to produce comparable performance to the diffusion model on CelebA data, the binary classification task lead to some misleading high performance of the naive baselines such as mean-imputation. Thus, to fully show the benefit of DMV, we test on three different 10-class datasets. MNIST (Deng, 2012) contains hand-drawn digits, with class labels for the 10 digits. CIFAR10 (Krizhevsky et al., 2009) contains low resolution images belonging to various classes, such as cat, airplane, or truck. StarCraftCIFAR10 (Kulinski et al., 2023) contains spacial data representing the location of units and environment features, and asks us to predict the map out of 5 options and whether the game is in the start or end. For all three datasets, we masked them by dividing the image

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

Method	MVCE	Calibrated
1 Sample Imputation, 0 Variance	0.1376	- -
1 Sample Imputation, 0.99 Max Variance	0.1004	0.0841 -16%
Mean Imputation, 0 Variance	0.0976	- -
Mean Imputation, 0.99 Max Variance	0.1478	0.0991 -33%
Generator Correlation 0.7 \rightarrow 0.6	0.0107	0.0107 0%
Generator Correlation 0.7 \rightarrow 0.8	0.0203	0.0076 -63%
Generator Covariance \times 0.25	0.0488	0.0086 -82%
Generator Covariance \times 4.0	0.0771	0.0077 -90%
Ground Truth Generator	0.0068	0.0068 0%

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

Method \ Missing	Time (seconds)	Feature: Blond Hair				Feature: Smiling				Feature: Eyeglasses						
		C. Acc.	M. Acc.	MVCE	Calibrated	C. Acc.	M. Acc.	MVCE	Calibrated	C. Acc.	M. Acc.	MVCE	Calibrated			
Diffusion - 1 Sample, 0 Variance	182	93.8%	92.1%	0.0522	-	91.9%	84.1%	0.1357	-	99.3%	99.0%	0.0058	-			
Diffusion - 1 Sample, 0.5 Max Variance	182	93.8%	92.1%	0.0328	0.0296	-10%	91.9%	84.1%	0.1028	0.0990	-4%	99.3%	99.0%	0.0045	0.0040	-11%
Diffusion - 1 Sample, Scale Probability	182	93.8%	92.1%	0.0448	0.0326	-27%	91.9%	84.1%	0.1226	0.1028	-16%	99.3%	99.0%	0.0058	0.0045	-22%
Diffusion - 3 Samples - Empirical Variance	547	93.8%	92.4%	0.0205	0.0176	-14%	91.9%	85.6%	0.0555	0.0530	-5%	99.3%	99.0%	0.0037	0.0037	-1%
Diffusion - 30 Samples - Empirical Variance	5465	93.8%	93.2%	0.0062	0.0063	2%	91.9%	86.4%	0.0163	0.0144	-12%	99.3%	99.0%	0.0021	0.0033	54%
Mean Imputation, 0 Variance	0.547	93.8%	89.4%	0.0775	-	91.9%	70.5%	0.2633	-	99.3%	97.5%	0.0220	-			
Mean Imputation, 0.5 Max Variance	0.538	93.8%	89.4%	0.0581	0.0548	-6%	91.9%	70.5%	0.2335	0.2281	-2%	99.3%	97.5%	0.0111	0.0109	-2%
Mean Imputation, Scale Probability	0.619	93.8%	89.4%	0.0704	0.0583	-17%	91.9%	70.5%	0.2548	0.2335	-8%	99.3%	97.5%	0.0192	0.0113	-41%
Direct Missing Value (DMV)	0.901	94.5%	93.0%	0.0286	0.0059	-79%	92.8%	77.1%	0.0711	0.0695	-2%	99.5%	99.1%	0.0029	0.0036	23%

Table 3: The Direct Missing Value approach performance well across all three datasets, having low MVCE without a significant drop in accuracy. While the missing value robust classifier could produce good accuracy, uncertainty heuristics only work when behavior is mostly constant across the whole dataset, limiting high performance to MNIST. Mean imputation produced too significant of an accuracy drop in all cases.

Method \ Missing	MNIST				CIFAR10				StarCraftCIFAR10						
	C Acc.	M Acc.	MVCE	Calibrated	C Acc.	M Acc.	MVCE	Calibrated	C Acc.	M Acc.	MVCE	Calibrated			
Mean Imputation, 0 Variance	99.5%	58.7%	0.4116	-	71.8%	32.1%	0.6471	-	79.9%	70.5%	0.2231	-			
Mean Imputation, 0.5 Max Variance	99.5%	58.6%	0.1900	0.1644	-14%	71.8%	32.4%	0.3092	0.2855	-8%	79.9%	70.6%	0.0493	0.0214	-57%
Mean Imputation, Scale Probability	99.5%	59.0%	0.3053	0.1903	-38%	71.8%	32.0%	0.4776	0.3115	-35%	79.9%	70.6%	0.0989	0.0233	-76%
Missing Robust, 0 Variance	99.6%	96.9%	0.0296	-	82.3%	74.7%	0.1604	-	71.8%	71.3%	0.1278	-			
Missing Robust, 0.5 Max Variance	99.6%	96.9%	0.0031	0.0031	0%	82.3%	74.7%	0.0804	0.0740	-8%	71.8%	71.2%	0.0809	0.0582	-28%
Missing Robust, Scale Probability	99.6%	96.9%	0.0168	0.0048	-71%	82.3%	74.7%	0.1349	0.0807	-40%	71.8%	71.2%	0.0937	0.0571	-39%
Direct Missing Value	98.5%	92.7%	0.0736	0.0354	-52%	80.0%	72.3%	0.0455	0.0455	0%	79.8%	78.5%	0.0201	0.0113	-44%

into a 4x4 grid of "sensors", and had an independent random chance for each sensor to be dropped during evaluation, making the experiment MCAR during evaluation.

Method Details: Since the Monte Carlo baseline was too expensive to practically deploy, we skipped it for the multi-class datasets. We instead compare DMV to two baselines for handling missing values: mean imputation, and a classifier robust to missing values (effectively, the input augmentation of DMV without the output or the training objective). Both DMV and the robust to missing classifier use the evaluation mask setup for training, making it MCAR. Both baselines used three methods for estimating variance: 0, 0.5 max variance, and directly scaling the probability values by 10 (subsection 4.1). For post-hoc calibration (subsection 5.2), we used the same mask setup on validation data to select between 9 different calibration constants.

Results: 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 accuracy 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 values will not change the prediction. We proposed a Direct Missing Value method to estimate MVU, and developed the MVCE 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.

540 REFERENCES

- 541 Anastasios N Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and
542 distribution-free uncertainty quantification. *arXiv preprint arXiv:2107.07511*, 2021.
- 543
544 Melissa J Azur, Elizabeth A Stuart, Constantine Frangakis, and Philip J Leaf. Multiple imputation by
545 chained equations: what is it and how does it work? *International journal of methods in psychiatric*
546 *research*, 20(1):40–49, 2011.
- 547
548 Viktor Bengs, Eyke Hüllermeier, and Willem Waegeman. Pitfalls of epistemic uncertainty quan-
549 tification through loss minimisation. *Advances in Neural Information Processing Systems*, 35:
550 29205–29216, 2022.
- 551
552 Li Deng. The mnist database of handwritten digit images for machine learning research [best of the
553 web]. *IEEE signal processing magazine*, 29(6):141–142, 2012.
- 554
555 Andrew Gelman, John B Carlin, Hal S Stern, and Donald B Rubin. *Bayesian data analysis*. Chapman
556 and Hall/CRC, 1995.
- 557
558 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
559 networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
- 560
561 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
562 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
563 pp. 770–778, 2016.
- 564
565 Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for
566 improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.
- 567
568 Pasha Khosravi, YooJung Choi, Yitao Liang, Antonio Vergari, and Guy Van den Broeck. On tractable
569 computation of expected predictions. *Advances in Neural Information Processing Systems*, 32,
570 2019.
- 571
572 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 573
574 Sean Kulinski, Nicholas R Waytowich, James Z Hare, and David I Inouye. Starcraftimage: A dataset
575 for prototyping spatial reasoning methods for multi-agent environments. In *Proceedings of the*
576 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22004–22013, 2023.
- 577
578 Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive
579 facial image manipulation. corr abs/1907.11922 (2019). *arXiv preprint arXiv:1907.11922*, 2019.
- 580
581 Yoad Lewenberg, Yoram Bachrach, Ulrich Paquet, and Jeffrey Rosenschein. Knowing what to
582 ask: A bayesian active learning approach to the surveying problem. In *Proceedings of the AAAI*
583 *Conference on Artificial Intelligence*, volume 31, 2017.
- 584
585 Roderick JA Little and Donald B Rubin. *Statistical analysis with missing data*, volume 793. John
586 Wiley & Sons, 2019.
- 587
588 Anji Liu, Mathias Niepert, and Guy Van den Broeck. Image inpainting via tractable steering of
589 diffusion models. *arXiv preprint arXiv:2401.03349*, 2023.
- 590
591 Jeremiah Liu, John Paisley, Marianthi-Anna Kioumourtoglou, and Brent Coull. Accurate uncertainty
592 estimation and decomposition in ensemble learning. *Advances in neural information processing*
593 *systems*, 32, 2019.
- 588
589 Chao Ma, Sebastian Tschitschek, Konstantina Palla, José Miguel Hernández-Lobato, Sebastian
590 Nowozin, and Cheng Zhang. Eddi: Efficient dynamic discovery of high-value information with
591 partial vae. *arXiv preprint arXiv:1809.11142*, 2018.
- 592
593 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.

- 594 Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated proba-
595 bilities using bayesian binning. In *Proceedings of the AAAI conference on artificial intelligence*,
596 volume 29, 2015.
- 597 Jeremy Nixon, Michael W Dusenberry, Linchuan Zhang, Ghassen Jerfel, and Dustin Tran. Measuring
598 calibration in deep learning. In *CVPR workshops*, volume 2, 2019.
- 600 Yaniv Romano, Evan Patterson, and Emmanuel Candes. Conformalized quantile regression. *Advances*
601 *in neural information processing systems*, 32, 2019.
- 602 Yusuf Sale, Viktor Bengs, Michele Caprio, and Eyke Hüllermeier. Second-order uncertainty quantifi-
603 cation: A distance-based approach. In *Forty-first International Conference on Machine Learning*,
604 2023.
- 606 Murat Sensoy, Lance Kaplan, and Melih Kandemir. Evidential deep learning to quantify classification
607 uncertainty. *Advances in neural information processing systems*, 31, 2018.
- 608 Burr Settles. Active learning literature survey. 2009.
- 609 Juozas Vaicenavicius, David Widmann, Carl Andersson, Fredrik Lindsten, Jacob Roll, and Thomas
610 Schön. Evaluating model calibration in classification. In *The 22nd international conference on*
611 *artificial intelligence and statistics*, pp. 3459–3467. PMLR, 2019.
- 613 Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic learning in a random world*.
614 Springer, 2005.
- 616 Margaux Zaffran, Aymeric Dieuleveut, Julie Josse, and Yaniv Romano. Conformal prediction with
617 missing values. In *International Conference on Machine Learning*, pp. 40578–40604. PMLR,
618 2023.
- 619 Guanhua Zhang, Jiabao Ji, Yang Zhang, Mo Yu, Tommi S Jaakkola, and Shiyu Chang. Towards
620 coherent image inpainting using denoising diffusion implicit models. *The Fortieth International*
621 *Conference on Machine Learning*, 2023.
- 622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

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 be 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: $\ell_w(y, \hat{y}) \triangleq w_y \mathbb{1}(y \neq \hat{y})$, where w_y is the cost associated with misclassifying a 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:

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

$$\equiv \arg \max_j \Pr(\bigcap_{j' \neq j} w_j \Phi_j \geq w_{j'} \Phi_{j'}), \quad (11)$$

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

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

The basic intuition of soft confidence values is that if the confidence is c_{soft} percent, then the classifier should be correct c_{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 soft predictions $\hat{y}_{\text{soft}, i}$ and soft confidence values $\hat{c}_{\text{soft}, i}$ for each sample, ECE first partitions the dataset into bins \mathcal{B} based on the soft confidence values $\hat{c}_{\text{soft}, i}$. Then, the ECE is defined as: $\text{ECE} = \sum_{B \in \mathcal{B}} \frac{|B|}{|\mathcal{D}|} |\text{acc}(B) - \bar{c}_{\text{soft}}(B)|$, where $\text{acc}(B) \triangleq \frac{1}{|B|} \sum_{i \in B} \mathbb{1}(\hat{y}_{\text{soft}, i} = y_i)$ is the empirical classification accuracy and $\bar{c}_{\text{soft}}(B) \triangleq \frac{1}{|B|} \sum_{i \in B} \hat{c}_{\text{soft}, i}$ is the average soft confidence. The fact that ECE compares the average confidence

702 to accuracy is related to the fact that the soft decision rule and corresponding confidences are based
 703 on the Bayes optimal decision rule, which provides the best accuracy among all possible classifiers.
 704

705 A.3 ADDITIONAL RELATED WORKS

706 This section covers additional works adjacent to ours that are less directly relevant to our methods.
 707

708 **Conformal Prediction** Conformal prediction is an uncertainty quantification method that can
 709 provide distribution-free, statistical guarantees on predictions for any underlying model (Vovk et al.,
 710 2005). Instead of a single prediction, it produces a prediction set (for classification) or interval (for
 711 regression) that is guaranteed to contain the true outcome with a user-specified probability (e.g.,
 712 90%), without depending on assumptions about underlying data distribution or model correctness
 713 (Angelopoulos & Bates, 2021). Recent work has adapted the area to missing values (Zaffran et al.,
 714 2023) and high-dimensional settings (Romano et al., 2019; Zaffran et al., 2023). However, these
 715 methods aim to keep the uncertainty prediction valid despite missing values instead of analyzing the
 716 uncertainty added specifically due to missing values to determine if there is insufficient information.
 717

718 **Active Learning** Active learning is a field in machine learning where data is left unlabeled, and the
 719 model attempts to determine the most useful additional samples to label (Settles, 2009). A problem
 720 in that field of particular relevance to our work is the active surveying problem; survey responses are
 721 modeled as a set of questions with incomplete answers, and the model must decide which question
 722 is most useful to ask next (Lewenberg et al., 2017; Ma et al., 2018). This setup naturally fits as a
 723 missing values problem, though notably work in this area is focused on selecting features based on
 724 information in the feature distribution $p(X_{\mathcal{M}}|X_{\mathcal{O}})$ rather than measuring uncertainty in the prediction
 725 due to missing features $p(Y|X_{\mathcal{O}})$.
 726

727 B PROOFS

728 B.1 MISCELLANEOUS PROOF(S)

729 Proof of the equivalence between 0-1 loss minimization and Bayes optimal classification rule:
 730

$$731 y = \arg \min_j \mathbb{E}_{p(Y;\phi)}[\ell_{0-1}(Y, j)] \quad (13)$$

$$732 = \arg \min_j \mathbb{E}_{p(Y;\phi)}[\mathbb{1}(Y \neq j)] \quad (14)$$

$$733 = \arg \min_j \sum_{j' \neq j} \phi_{j'} \quad (15)$$

$$734 = \arg \min_j 1 - \phi_j \quad (16)$$

$$735 = \arg \max_j \phi_j. \quad (17)$$

736 Proof that cost-sensitive classification is a weighted version of the Bayes optimal classification rule:
 737

$$738 y = \arg \min_j \mathbb{E}_{p(Y;\phi)}[\ell_w(Y, j)] \quad (18)$$

$$739 = \arg \min_j \mathbb{E}_{p(Y;\phi)}[w_Y \mathbb{1}(Y \neq j)] \quad (19)$$

$$740 = \arg \min_j \sum_{j' \neq j} w_{j'} \phi_{j'} \quad (20)$$

$$741 = \arg \min_j \sum_{j'} w_{j'} \phi_{j'} - w_j \phi_j \quad (21)$$

$$742 = \arg \min_j -w_j \phi_j \quad (22)$$

$$743 = \arg \max_j w_j \phi_j. \quad (23)$$

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:

$$\mathbb{E}_{p(X_{\mathcal{O}})} [\text{KL}(p_f(\Phi | X_{\mathcal{O}}), \hat{p}_\psi(\Phi | X_{\mathcal{O}}))] \quad (24)$$

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

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

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

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. \square

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:

$$\arg \min_{\tilde{q}} \mathbb{E}_{X_{\mathcal{O}}, X_{\mathcal{M}}} [-\log \tilde{q}(f^*(X_{\mathcal{O}}, X_{\mathcal{M}}) | X_{\mathcal{O}})] \quad (28)$$

$$= \arg \min_{\tilde{q}} \mathbb{E}_{X_{\mathcal{O}}} [\text{KL}(p_{f^*}(\Phi|X_{\mathcal{O}}), \tilde{q}(\Phi|X_{\mathcal{O}}))] \quad (29)$$

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

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

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. \square

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:

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

where x is a complete feature instance, $\ell_{\text{KL}}(p, q)$ 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:

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

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

810 Unlike the previous bi-level problem, this problem does not decompose and thus must use more
 811 advanced bi-level optimization strategies such as alternating optimization.
 812

813 814 815 816 D ADDITIONAL EXPERIMENTAL RESULTS

817
 818
 819
 820 For real world each experiment, we reported two
 821 values: MVCE and consistency, as described in
 822 Section 3. We used a simple zero-one cost func-
 823 tion for estimating confidence, with different
 824 mutators based on the dataset. All experiments
 825 were run in Ubuntu servers using NVIDIA RTX
 826 A5000 GPUs.

827 **MVCE** MVCE is our primary metric for compar-
 828 ing methods, reporting how closely the esti-
 829 mate of confidence matches to the likelihood the
 830 prediction will change as more information is
 831 revealed. Lower MVCE indicates the method is
 832 well calibrated and thus gives good confidence
 833 estimates. In our experiments, we computed
 834 MVCE using 10 bins. When computing MVCE
 835 we ran 4 trials of each experiment (or 10 for
 836 synthetic) and reported the average of the metric
 837 and its standard deviation in order to minimize
 sources of randomness in the experiment.

838 **Classifiers** For all non-synthetic datasets, we
 839 used a modified ResNet18 for predictions, re-
 840 placing the final layer with an appropriately
 841 sized layer for the number of classes. For the Di-
 842 rect Missing Value model and the robust to miss-
 843 ing values classifier, we additionally replaced
 844 the first layer to take 4 channels as an input in-
 845 stead of 3. DMV additionally changed the final activation function from sigmoid (2 classes) or
 846 softmax (3 or more classes) to exponential. Both the standard and robust to missing values classifier
 847 were trained using cross entropy loss. We used the SVG optimizer for the standard classifier and
 848 Adam for the robust to missing values classifier. Training bash files can be found in the in the
 849 codebase for more details.

850 Hyperparameters

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

854 **Experiment Setup** We used a random crop from the larger image size to 224x224 for our model
 855 inputs at training time, and a center crop to 224x224 at testing time. For the DMV and the robust to
 856 missing values classifier, we replaced the first layer of PyTorch’s pretrained ResNet18. Additionally,
 857 we replaced the final layer to adjust the class size. All other layers had weights copied over from
 858 the pretrained weights found in PyTorch. For our mutator used in both training the DMV and robust
 859 classifier and evaluating MVCE, we divided the image into a 4x4 grid of 56x56 pixel regions and
 860 gave each region a 50% chance to be removed each time a sample was fetched. For all three datasets,
 861 we cached the average image to use for mean imputation. Both mean imputation and the robust to
 862 missing values classifier used three different approaches to estimating variance: 0 variance (a single
 863 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 4: Reported accuracy and expected calibration loss for all models on clean data. Despite the difference in training, neither DMV nor robust to missing classifiers have significantly different accuracy. ECE changed between models, though this difference could likely be mitigated using traditional calibration techniques.

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	99.59%	0.0026
MNIST	DMV	10	98.54%	0.0233
CIFAR10	Classifier	10	71.77%	0.0753
CIFAR10	Robust	10	82.35%	0.0478
CIFAR10	DMV	10	79.99%	0.0561
StarCraftCIFAR10	Classifier	10	79.90%	0.0312
StarCraftCIFAR10	Robust	10	71.84%	0.0830
StarCraftCIFAR10	DMV	10	79.78%	0.0297

864 D.1 CELEBA

865
866 We made use of the CelebA-HQ dataset (Karras et al., 2017), which contains 10,000 64x64 color
867 images of celebrity faces. We choose this dataset as it was easy to form intuitions about the relationship
868 between missing masks and the CelebA features. Additionally, it was easy to locate pre-trained
869 diffusion models for the dataset. Since CelebA-HQ was designed for training generative models, it
870 lacks features, so we obtained the feature label information from (Lee et al., 2019). Experiments on
871 the CelebA dataset were run on three different target features: Blond Hair, Eyeglasses, and Smiling.
872 Tables 5, 6, and 7 show CelebA results on the blond hair feature. Tables 8, 9, and 10 show CelebA
873 results on the eyeglasses feature. Tables 11, 12, and 13 show CelebA results on the smiling feature.

874 **Experiment Setup** An independant ResNet18 classifier with pre-trained initial weights was fine-
875 tuned to predict each feature, making the experiments on CelebA all two class: either the feature
876 is present or absent. We did not train the classifier with robustness to missing values. We used the
877 64x64 pixel versions of the CelebA images as inputs with no rescaling during preprocessing to better
878 match the output resolution of the pretrained diffusion model; taking advantage of the fact ResNet18
879 allows variable sized inputs. The DMV model was a similarly constructed ResNet18 fine-tuned with
880 a mutator that randomly selected a missing mask between fully observed, fully missing, top half
881 missing, and bottom half missing. At test time, we used a single mask for the entire experiment,
882 either top half missing or bottom half missing. In the tables below, the words "top" or "bottom"
883 always refer to the half missing.

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

890 **Sample Cache** To reduce the time it takes to run multiple experiments, we cached 30 samples from
891 the diffusion model for each mask on each of the 2000 test samples, along with each of the 2000
892 validation samples for the sake of calibration. This allowed the time to run each experiment with the
893 diffusion model to be relatively close to that of the non-diffusion approaches at the cost of requiring
894 several months to pre-generate all the samples. The single sample times reported in Table 2 for any
895 diffusion approaches takes the time to compute that number of samples from the cache generation and
896 adds it to the time to run that particular method. This caching approach is likely not representative of
897 how the model would be used when deployed, but we do not believe it had any significant impact on
898 the results beyond a small reduction of randomness when running multiple trials. Leveraging this
899 cache limited us to just the two masks in our experiments, though this is not a practical limitation to
900 either method as both DMV and the diffusion model can handle any arbitrary missing feature with
901 the right training.

902 Calibration

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

909
910
911
912
913
914
915
916
917

Table 5: Blond Hair is a feature that is often visible in both halves of the image, though may sometimes only be visible in the top half. This causes a smaller average accuracy drop across the baselines. Diffusion and DMV both give the best estimate of the confidence for that drop.

CelebA Blond Hair - Mask Average	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Diffusion - 1 Sample, 0 Variance	93.8%	92.1% (SD 0.0000)	0.0522 (SD 0.0000)	-	-
Diffusion - 1 Sample, 0.5 Max Variance	93.8%	92.1% (SD 0.0000)	0.0328 (SD 0.0005)	0.1	0.0296 (SD 0.0008) -9.7%
Diffusion - 1 Sample, Scale Probability	93.8%	92.1% (SD 0.0000)	0.0448 (SD 0.0001)	0.1	0.0326 (SD 0.0005) -27.2%
Diffusion - 3 Samples - Empirical Variance	93.8%	92.4% (SD 0.0000)	0.0205 (SD 0.0003)	0.25	0.0176 (SD 0.0008) -14.1%
Diffusion - 30 Samples - Empirical Variance	93.8%	93.2% (SD 0.0000)	0.0062 (SD 0.0004)	0.5	0.0063 (SD 0.0007) 2.0%
Mean Imputation, 0 Variance	93.8%	89.4% (SD 0.0000)	0.0775 (SD 0.0000)	-	-
Mean Imputation, 0.5 Max Variance	93.8%	89.4% (SD 0.0000)	0.0581 (SD 0.0009)	0.1	0.0548 (SD 0.0006) -5.7%
Mean Imputation, Scale Probability	93.8%	89.4% (SD 0.0000)	0.0704 (SD 0.0003)	0.1	0.0583 (SD 0.0004) -17.2%
Direct Missing Value (DMV)	94.5%	93.0% (SD 0.0000)	0.0286 (SD 0.0002)	7.5	0.0059 (SD 0.0011) -79.2%

Table 6: While blond hair is typically visible in both halves, it is more likely to be in the top half, causing a small accuracy drop for the imputation methods. This is less notable than the smiling method, and does not impact the ranking of methods.

CelebA Blond Hair - Top Missing	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Diffusion - 1 Sample, 0 Variance	93.8%	90.8% (SD 0.0000)	0.0665 (SD 0.0000)	-	-
Diffusion - 1 Sample, 0.5 Max Variance	93.8%	90.8% (SD 0.0000)	0.0472 (SD 0.0001)	0.1	0.0443 (SD 0.0001) -6.1%
Diffusion - 1 Sample, Scale Probability	93.8%	90.8% (SD 0.0000)	0.0599 (SD 0.0000)	0.1	0.0473 (SD 0.0001) -21.1%
Diffusion - 3 Samples - Empirical Variance	93.8%	91.0% (SD 0.0000)	0.0253 (SD 0.0001)	0.25	0.0224 (SD 0.0009) -11.6%
Diffusion - 30 Samples - Empirical Variance	93.8%	92.5% (SD 0.0000)	0.0083 (SD 0.0003)	0.5	0.0073 (SD 0.0008) -12.0%
Mean Imputation, 0 Variance	93.8%	85.7% (SD 0.0000)	0.1005 (SD 0.0000)	-	-
Mean Imputation, 0.5 Max Variance	93.8%	85.7% (SD 0.0000)	0.0901 (SD 0.0010)	0.1	0.0880 (SD 0.0001) -2.4%
Mean Imputation, Scale Probability	93.8%	85.7% (SD 0.0000)	0.0978 (SD 0.0001)	0.1	0.0901 (SD 0.0003) -7.9%
Direct Missing Value (DMV)	94.5%	92.0% (SD 0.0000)	0.0323 (SD 0.0002)	7.5	0.0058 (SD 0.0017) -82.0%

Table 7: Blond hair is almost always visible in the top half, making the accuracy remain high even among the baselines. Despite this, they often still get low MVCE, in this case due to under-estimating the confidence.

CelebA Blond Hair - Bottom Missing	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Diffusion - 1 Sample, 0 Variance	93.8%	93.5% (SD 0.0000)	0.0380 (SD 0.0000)	-	-
Diffusion - 1 Sample, 0.5 Max Variance	93.8%	93.5% (SD 0.0000)	0.0185 (SD 0.0007)	0.1	0.0150 (SD 0.0012) -18.8%
Diffusion - 1 Sample, Scale Probability	93.8%	93.5% (SD 0.0000)	0.0297 (SD 0.0002)	0.1	0.0179 (SD 0.0007) -39.7%
Diffusion - 3 Samples - Empirical Variance	93.8%	93.8% (SD 0.0000)	0.0156 (SD 0.0004)	0.25	0.0127 (SD 0.0009) -18.2%
Diffusion - 30 Samples - Empirical Variance	93.8%	93.8% (SD 0.0000)	0.0040 (SD 0.0005)	0.5	0.0053 (SD 0.0006) 30.9%
Mean Imputation, 0 Variance	93.8%	93.0% (SD 0.0000)	0.0545 (SD 0.0000)	-	-
Mean Imputation, 0.5 Max Variance	93.8%	93.0% (SD 0.0000)	0.0261 (SD 0.0009)	0.1	0.0216 (SD 0.0010) -17.2%
Mean Imputation, Scale Probability	93.8%	93.0% (SD 0.0000)	0.0430 (SD 0.0004)	0.1	0.0265 (SD 0.0005) -38.4%
Direct Missing Value (DMV)	94.5%	94.0% (SD 0.0000)	0.0249 (SD 0.0002)	7.5	0.0060 (SD 0.0005) -75.7%

Table 8: Eyeglasses are often visible in both halves of the image, causing a minimal drop in accuracy regardless of mask. This leads to MVCE primarily penalizing underestimated confidence.

CelebA Eyeglasses - Mask Average	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Diffusion - 1 Sample, 0 Variance	99.3%	99.0% (SD 0.0000)	0.0058 (SD 0.0000)	-	-
Diffusion - 1 Sample, 0.5 Max Variance	99.3%	99.0% (SD 0.0000)	0.0045 (SD 0.0002)	0.1	0.0040 (SD 0.0001) -10.9%
Diffusion - 1 Sample, Scale Probability	99.3%	99.0% (SD 0.0000)	0.0058 (SD 0.0001)	0.1	0.0045 (SD 0.0002) -22.4%
Diffusion - 3 Samples - Empirical Variance	99.3%	99.0% (SD 0.0000)	0.0037 (SD 0.0002)	0.5	0.0037 (SD 0.0001) -1.1%
Diffusion - 30 Samples - Empirical Variance	99.3%	99.0% (SD 0.0000)	0.0021 (SD 0.0001)	7.5	0.0033 (SD 0.0003) 54.4%
Mean Imputation, 0 Variance	99.3%	97.5% (SD 0.0000)	0.0220 (SD 0.0000)	-	-
Mean Imputation, 0.5 Max Variance	99.3%	97.5% (SD 0.0000)	0.0111 (SD 0.0003)	0.75	0.0109 (SD 0.0004) -1.7%
Mean Imputation, Scale Probability	99.3%	97.5% (SD 0.0000)	0.0192 (SD 0.0003)	0.1	0.0113 (SD 0.0004) -41.0%
Direct Missing Value (DMV)	99.5%	99.1% (SD 0.0000)	0.0029 (SD 0.0001)	5	0.0036 (SD 0.0001) 23.4%

Table 9: While intuitively eyeglasses are a feature in the top half, the edges are often on the boundary allowing high accuracy from just the bottom half as well.

CelebA Eyeglasses - Top Missing	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Diffusion - 1 Sample, 0 Variance	99.3%	99.1% (SD 0.0000)	0.0055 (SD 0.0000)	-	-
Diffusion - 1 Sample, 0.5 Max Variance	99.3%	99.1% (SD 0.0000)	0.0041 (SD 0.0002)	0.1	0.0033 (SD 0.0000) -19.5%
Diffusion - 1 Sample, Scale Probability	99.3%	99.1% (SD 0.0000)	0.0053 (SD 0.0001)	0.1	0.0040 (SD 0.0003) -25.8%
Diffusion - 3 Samples - Empirical Variance	99.3%	99.0% (SD 0.0000)	0.0034 (SD 0.0002)	0.5	0.0031 (SD 0.0001) -9.4%
Diffusion - 30 Samples - Empirical Variance	99.3%	99.0% (SD 0.0000)	0.0024 (SD 0.0002)	7.5	0.0042 (SD 0.0004) 75.1%
Mean Imputation, 0 Variance	99.3%	97.6% (SD 0.0000)	0.0230 (SD 0.0000)	-	-
Mean Imputation, 0.5 Max Variance	99.3%	97.6% (SD 0.0000)	0.0062 (SD 0.0003)	0.75	0.0061 (SD 0.0006) -1.5%
Mean Imputation, Scale Probability	99.3%	97.6% (SD 0.0000)	0.0182 (SD 0.0000)	0.1	0.0067 (SD 0.0005) -63.4%
Direct Missing Value (DMV)	99.5%	99.0% (SD 0.0000)	0.0025 (SD 0.0001)	5	0.0026 (SD 0.0001) 4.8%

Table 10: Eyeglasses are clearly visible in the top half of the image, allowing high accuracy regardless of method of handling missing values. The direct methods manage to minimize MVCE by not underestimating confidence.

CelebA Eyeglasses - Bottom Missing	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Diffusion - 1 Sample, 0 Variance	99.3%	98.9% (SD 0.0000)	0.0060 (SD 0.0000)	-	-
Diffusion - 1 Sample, 0.5 Max Variance	99.3%	98.9% (SD 0.0000)	0.0050 (SD 0.0002)	0.1	0.0048 (SD 0.0001) -3.9%
Diffusion - 1 Sample, Scale Probability	99.3%	98.9% (SD 0.0000)	0.0062 (SD 0.0000)	0.1	0.0050 (SD 0.0000) -19.5%
Diffusion - 3 Samples - Empirical Variance	99.3%	99.0% (SD 0.0000)	0.0040 (SD 0.0002)	0.5	0.0043 (SD 0.0002) 5.9%
Diffusion - 30 Samples - Empirical Variance	99.3%	99.0% (SD 0.0000)	0.0018 (SD 0.0001)	7.5	0.0023 (SD 0.0001) 26.9%
Mean Imputation, 0 Variance	99.3%	97.5% (SD 0.0000)	0.0210 (SD 0.0000)	-	-
Mean Imputation, 0.5 Max Variance	99.3%	97.5% (SD 0.0000)	0.0160 (SD 0.0003)	0.75	0.0157 (SD 0.0002) -1.8%
Mean Imputation, Scale Probability	99.3%	97.5% (SD 0.0000)	0.0202 (SD 0.0004)	0.1	0.0160 (SD 0.0003) -20.8%
Direct Missing Value (DMV)	99.5%	99.1% (SD 0.0000)	0.0033 (SD 0.0001)	5	0.0045 (SD 0.0001) 37.3%

Table 11: The smiling feature has the worst performance of DMV out of all three features, though it is still on average better than the baselines. This suggests a need for further optimization of the model.

CelebA Smiling - Mask Average	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Diffusion - 1 Sample, 0 Variance	91.9%	84.1% (SD 0.0000)	0.1357 (SD 0.0000)	-	-
Diffusion - 1 Sample, 0.5 Max Variance	91.9%	84.1% (SD 0.0000)	0.1028 (SD 0.0015)	0.1	0.0990 (SD 0.0007) -3.7%
Diffusion - 1 Sample, Scale Probability	91.9%	84.1% (SD 0.0000)	0.1226 (SD 0.0002)	0.1	0.1028 (SD 0.0005) -16.2%
Diffusion - 3 Samples - Empirical Variance	91.9%	85.6% (SD 0.0000)	0.0555 (SD 0.0008)	0.25	0.0530 (SD 0.0006) -4.6%
Diffusion - 30 Samples - Empirical Variance	91.9%	86.4% (SD 0.0000)	0.0163 (SD 0.0018)	0.5	0.0144 (SD 0.0028) -11.5%
Mean Imputation, 0 Variance	91.9%	70.5% (SD 0.0000)	0.2633 (SD 0.0000)	-	-
Mean Imputation, 0.5 Max Variance	91.9%	70.5% (SD 0.0000)	0.2335 (SD 0.0007)	0.1	0.2281 (SD 0.0002) -2.3%
Mean Imputation, Scale Probability	91.9%	70.5% (SD 0.0000)	0.2548 (SD 0.0003)	0.1	0.2335 (SD 0.0006) -8.4%
Direct Missing Value (DMV)	92.8%	77.1% (SD 0.0000)	0.0711 (SD 0.0004)	2.5	0.0695 (SD 0.0010) -2.2%

Table 12: Since smiling is primarily in the bottom half, the top half missing does not significantly impact the prediction leading to a high accuracy prediction. Interestingly, the simple heuristic of half the maximum variance manages to minimize MVCE. For Monte Carlo, this suggests too much of the top half is used in the prediction, while DMV may simply estimate too little confidence overall.

CelebA Smiling - Top Missing	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Diffusion - 1 Sample, 0 Variance	91.9%	91.0% (SD 0.0000)	0.0390 (SD 0.0000)	-	-
Diffusion - 1 Sample, 0.5 Max Variance	91.9%	91.0% (SD 0.0000)	0.0096 (SD 0.0023)	0.1	0.0072 (SD 0.0011) -24.8%
Diffusion - 1 Sample, Scale Probability	91.9%	91.0% (SD 0.0000)	0.0265 (SD 0.0003)	0.1	0.0095 (SD 0.0008) -64.2%
Diffusion - 3 Samples - Empirical Variance	91.9%	91.8% (SD 0.0000)	0.0152 (SD 0.0005)	0.25	0.0090 (SD 0.0005) -40.6%
Diffusion - 30 Samples - Empirical Variance	91.9%	91.8% (SD 0.0000)	0.0126 (SD 0.0004)	0.5	0.0088 (SD 0.0011) -30.0%
Mean Imputation, 0 Variance	91.9%	91.1% (SD 0.0000)	0.0450 (SD 0.0000)	-	-
Mean Imputation, 0.5 Max Variance	91.9%	91.1% (SD 0.0000)	0.0218 (SD 0.0002)	0.1	0.0173 (SD 0.0001) -20.6%
Mean Imputation, Scale Probability	91.9%	91.1% (SD 0.0000)	0.0382 (SD 0.0002)	0.1	0.0218 (SD 0.0004) -43.0%
Direct Missing Value (DMV)	92.8%	92.7% (SD 0.0000)	0.0363 (SD 0.0002)	2.5	0.0277 (SD 0.0002) -23.6%

Table 13: Without the bottom half, the main indicator of smiling is missing, leading to a significant accuracy drop among baselines. DMV in this case is not ideally optimized, making it barely better than imputation.

CelebA Smiling - Bottom Missing	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Diffusion - 1 Sample, 0 Variance	91.9%	77.2% (SD 0.0000)	0.2325 (SD 0.0000)	-	-
Diffusion - 1 Sample, 0.5 Max Variance	91.9%	77.2% (SD 0.0000)	0.1960 (SD 0.0002)	0.1	0.1908 (SD 0.0001) -2.7%
Diffusion - 1 Sample, Scale Probability	91.9%	77.2% (SD 0.0000)	0.2188 (SD 0.0001)	0.1	0.1960 (SD 0.0002) -10.4%
Diffusion - 3 Samples - Empirical Variance	91.9%	79.5% (SD 0.0000)	0.0959 (SD 0.0011)	0.25	0.0969 (SD 0.0007) 1.1%
Diffusion - 30 Samples - Empirical Variance	91.9%	81.0% (SD 0.0000)	0.0199 (SD 0.0028)	0.5	0.0199 (SD 0.0042) 0.2%
Mean Imputation, 0 Variance	91.9%	49.9% (SD 0.0000)	0.4815 (SD 0.0000)	-	-
Mean Imputation, 0.5 Max Variance	91.9%	49.9% (SD 0.0000)	0.4453 (SD 0.0010)	0.1	0.4390 (SD 0.0004) -1.4%
Mean Imputation, Scale Probability	91.9%	49.9% (SD 0.0000)	0.4714 (SD 0.0004)	0.1	0.4453 (SD 0.0009) -5.5%
Direct Missing Value (DMV)	92.8%	61.4% (SD 0.0000)	0.1058 (SD 0.0006)	2.5	0.1113 (SD 0.0016) 5.2%

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 14.

Table 14: MNIST was an easier dataset to predict with missing values, leading to few cases where more information was needed to make good predictions (and thus 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, this dataset ended up as too easy of a task for traditional neural network to handle so its harder to see the benefit from MVU.

MNIST Method \ Missing	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Mean Imputation, 0 Variance	99.48%	58.7% (SD 0.0024)	0.4116 (SD 0.0027)	-	-
Mean Imputation, 0.5 Max Variance	99.48%	58.6% (SD 0.0022)	0.1900 (SD 0.0025)	0.1	0.1644 (SD 0.0022) -14%
Mean Imputation, Scale Probability	99.48%	59.0% (SD 0.0021)	0.3053 (SD 0.0013)	0.1	0.1903 (SD 0.0044) -38%
Missing Robust, 0 Variance	99.59%	96.9% (SD 0.0017)	0.0296 (SD 0.0017)	-	-
Missing Robust, 0.5 Max Variance	99.59%	96.9% (SD 0.0005)	0.0031 (SD 0.0006)	1.0	0.0031 (SD 0.0006) 0%
Missing Robust, Scale Probability	99.59%	96.9% (SD 0.0004)	0.0168 (SD 0.0005)	0.1	0.0048 (SD 0.0004) -71%
Direct Missing Value (DMV)	98.54%	92.7% (SD 0.0016)	0.0736 (SD 0.0027)	10.0	0.0354 (SD 0.0010) -52%

CIFAR10 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 15.

Table 15: The DMV method and the robust classifier both have comparable mutated accuracy, 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 notably better MVCE, but it still did not compare. Overall, this dataset shows the value of DMV for estimating uncertainty that indicates the prediction likely shifted.

CIFAR10 Method \ Missing	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Mean Imputation, 0 Variance	71.77%	32.1% (SD 0.0055)	0.6471 (SD 0.0030)	-	-
Mean Imputation, 0.5 Max Variance	71.77%	32.4% (SD 0.0050)	0.3092 (SD 0.0056)	0.1	0.2855 (SD 0.0030) -8%
Mean Imputation, Scale Probability	71.77%	32.0% (SD 0.0027)	0.4776 (SD 0.0026)	0.1	0.3115 (SD 0.0033) -35%
Missing Robust, 0 Variance	82.35%	74.7% (SD 0.0050)	0.1604 (SD 0.0033)	-	-
Missing Robust, 0.5 Max Variance	82.35%	74.7% (SD 0.0022)	0.0804 (SD 0.0041)	0.1	0.0740 (SD 0.0009) -8%
Missing Robust, Scale Probability	82.35%	74.7% (SD 0.0014)	0.1349 (SD 0.0010)	0.1	0.0807 (SD 0.0017) -40%
Direct Missing Value (DMV)	79.99%	72.3% (SD 0.0011)	0.0455 (SD 0.0050)	1.0	0.0455 (SD 0.0050) 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 16.

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 16: The map prediction for StarCraftCIFAR10 makes it easier to achieve comparable accuracy on mutated data compared to clean data 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, hence accuracy being lower across the board compared to other datasets.

StarCraft Method \ Missing	Clean Acc.	Mutated Accuracy	MVCE	Scale	Calibrated MVCE
Mean Imputation, 0 Variance	79.90%	70.5% (SD 0.0039)	0.2231 (SD 0.0035)	-	-
Mean Imputation, 0.5 Max Variance	79.90%	70.6% (SD 0.0029)	0.0493 (SD 0.0046)	2.5	0.0214 (SD 0.0026) -57%
Mean Imputation, Scale Probability	79.90%	70.6% (SD 0.0078)	0.0989 (SD 0.0060)	0.25	0.0233 (SD 0.0044) -76%
Missing Robust, 0 Variance	71.84%	71.3% (SD 0.0031)	0.1278 (SD 0.0032)	-	-
Missing Robust, 0.5 Max Variance	71.84%	71.2% (SD 0.0032)	0.0809 (SD 0.0009)	2.5	0.0582 (SD 0.0032) -28%
Missing Robust, Scale Probability	71.84%	71.2% (SD 0.0021)	0.0937 (SD 0.0035)	0.25	0.0571 (SD 0.0008) -39%
Direct Missing Value (DMV)	79.78%	78.5% (SD 0.0026)	0.0201 (SD 0.0016)	0.5	0.0113 (SD 0.0016) -44%

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 224x224 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 $\mathbb{E}[X] = (0, 0)$, $\text{Var}(X_1) = 0.3$, $\text{Var}(X_2) = 1$, and $\text{corr}(X_1, X_2) = 0.7$. 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 $\text{sigmoid}(x_1 + x_2 - 1)$ 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 17 and Table 18.

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.

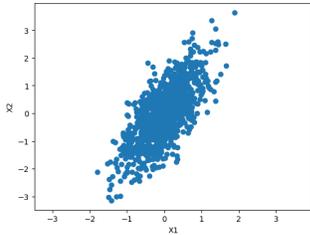


Figure 1: Example distribution of random samples from the synthetic dataset generated to test Missing Value Prediction Uncertainty.

1188 Table 17: Regardless of which feature is missing, we find ground truth overall outperforms the
 1189 alternatives. The imputation approaches with variance heuristics perform notably poorly. Generator
 1190 approaches approach the ground truth, but still see a notable increase in MVCE when considering
 1191 both features, especially before calibration.

Method \ Missing	Scale	X1			X2		
		MVCE	Calibrated MVCE		MVCE	Calibrated MVCE	
1 Sample Imputation, 0 Variance	-	0.0983 (SD 0.0018)	-	-	0.1768 (SD 0.0025)	-	-
1 Sample Imputation, 0.99 Max Variance	0.1	0.1392 (SD 0.0017)	0.1054 (SD 0.0026)	-24%	0.0616 (SD 0.0033)	0.0628 (SD 0.0018)	2%
Mean Imputation, 0 Variance	-	0.0718 (SD 0.0000)	-	-	0.1235 (SD 0.0000)	-	-
Mean Imputation, 0.99 Max Variance	0.1	0.1692 (SD 0.0001)	0.1164 (SD 0.0015)	-31%	0.1264 (SD 0.0002)	0.0817 (SD 0.0024)	-35%
Generator Correlation 0.7 \rightarrow 0.6	1	0.0054 (SD 0.0006)	0.0054 (SD 0.0006)	0%	0.0160 (SD 0.0008)	0.0160 (SD 0.0008)	0%
Generator Correlation 0.7 \rightarrow 0.8	0.5	0.0167 (SD 0.0006)	0.0058 (SD 0.0012)	-65%	0.0239 (SD 0.0004)	0.0094 (SD 0.0009)	-61%
Generator Covariance \times 0.25	0.25	0.0372 (SD 0.0007)	0.0057 (SD 0.0004)	-85%	0.0604 (SD 0.0007)	0.0115 (SD 0.0011)	-81%
Generator Covariance \times 4.0	5	0.0569 (SD 0.0013)	0.0092 (SD 0.0010)	-84%	0.0973 (SD 0.0016)	0.0062 (SD 0.0017)	-94%
Ground Truth Generator	1	0.0051 (SD 0.0007)	0.0051 (SD 0.0007)	0%	0.0085 (SD 0.0009)	0.0085 (SD 0.0009)	0%

1202 E LLM USAGE

1205 We made use of LLMs to assist in re-
 1206 vising paper contents such as editing,
 1207 suggestions for notation, and getting
 1208 initial feedback on ideas. Addition-
 1209 ally, LLMs were used to help locate
 1210 relevant related works. No section of
 1211 the paper was written fully by LLMs,
 1212 nor were they involved in any notable
 1213 capacity in producing the code used
 1214 for experiments.

1205 Table 18: On average, the ground truth generator is shown
 1206 to achieve the lowest MVCE, demonstrating that MVCE is
 1207 correctly evaluating our objective on this synthetic dataset.

Method \ Missing	Scale	Feature Average	
		MVCE	Calibrated MVCE
1 Sample Imputation, 0 Variance	-	0.1376 (SD 0.0021)	-
1 Sample Imputation, 0.99 Max Variance	0.1	0.1004 (SD 0.0025)	0.0841 (SD 0.0022) -16%
Mean Imputation, 0 Variance	-	0.0976 (SD 0.0000)	-
Mean Imputation, 0.99 Max Variance	0.1	0.1478 (SD 0.0002)	0.0991 (SD 0.0019) -33%
Generator Correlation 0.7 \rightarrow 0.6	1	0.0107 (SD 0.0007)	0.0107 (SD 0.0007) 0%
Generator Correlation 0.7 \rightarrow 0.8	0.5	0.0203 (SD 0.0005)	0.0076 (SD 0.0010) -63%
Generator Covariance \times 0.25	0.25	0.0488 (SD 0.0007)	0.0086 (SD 0.0008) -82%
Generator Covariance \times 4.0	5	0.0771 (SD 0.0014)	0.0077 (SD 0.0014) -90%
Ground Truth Generator	1	0.0068 (SD 0.0008)	0.0068 (SD 0.0008) 0%