# Fast Calibrated Explanations:
# Efficient and Uncertainty-Aware Explanations
# for Machine Learning Models

Tuwe Löfström[a], Fatima Rabia Yapicioglu[b], Alessandra Stramiglio[b], Helena Löfström[a], Fabio Vitali[b]

*[a]Jönköping AI Lab, Department of Computing, Jönköping University, Jönköping, Sweden*
*[b]DISI, Departement of Computer Science and Engineering, University of Bologna, Bologna, Italy*

## Abstract

This paper introduces Fast Calibrated Explanations, a method designed for generating rapid, uncertainty-aware explanations for machine learning models. By incorporating perturbation techniques from ConformaSight—a global explanation framework—into the core elements of Calibrated Explanations (CE), we achieve significant speedups. These core elements include local feature importance with calibrated predictions, both of which retain uncertainty quantification. While the new method sacrifices a small degree of detail, it excels in computational efficiency, making it ideal for high-stakes, real-time applications. Fast Calibrated Explanations are applicable to probabilistic explanations in classification and thresholded regression tasks, where they provide the likelihood of a target being above or below a user-defined threshold. This approach maintains the versatility of CE for both classification and probabilistic regression, making it suitable for a range of predictive tasks where uncertainty quantification is crucial.

*Keywords:* Uncertainty Quantification, Calibrated Explanations, ConformaSight, Explainable AI

*Email address:* `tuwe.lofstrom@ju.se` (Tuwe Löfström)

## 1. Introduction

Artificial Intelligence (AI) is becoming an integral part of modern society, influencing everything from retail recommendations to medical diagnosis predictions [1, 2] and even defence strategies [3]. In many predictive tasks, AI systems are typically developed as models trained via Machine Learning (ML) algorithms. While these models often achieve remarkable accuracy, such as outperforming medical professionals in cancer detection [4], they are neither flawless nor purely objective. Their performance is inherently tied to the data and algorithms used for training, making the outcomes sensitive to these factors.

As AI becomes more ubiquitous, the need for explainability in AI models has grown. In high-stakes applications, such as healthcare or autonomous driving, it is crucial that AI systems provide not only accurate predictions but also transparent reasoning behind those predictions. Explainable AI (XAI) aims to make AI's decision-making processes more understandable to humans, which enhances trust and facilitates informed decision-making by users. Furthermore, explainable models allow developers and stakeholders to identify potential weaknesses, limitations, or unintended consequences of the system. By providing insights into how the AI arrives at its conclusions, XAI helps bridge the gap between complex model operations and human interpretability, ultimately fostering more reliable and accountable AI systems across various domains [5].

When making decisions based on Machine Learning (ML) models, it is essential to account for the inherent uncertainty in their predictions [6, 7]. While many models provide point estimates, these alone fail to capture the degree of confidence in a given prediction, which can be crucial, particularly in safety-critical applications. Understanding uncertainty helps quantify the reliability of the model's output and aids in making more informed decisions.

Uncertainty in ML models can generally be categorized into two types: aleatoric and epistemic uncertainty [8]. Aleatoric uncertainty, also known as statistical or irreducible uncertainty, arises from inherent noise in the data. It represents variability in outcomes due to factors that cannot be explained by the model, such as measurement

2

errors or inherent randomness in the data-generating process. This type of uncertainty cannot be reduced by gathering more data because it is intrinsic to the task.

In contrast, epistemic uncertainty reflects the model's lack of knowledge, typically caused by limited or insufficient data. It is also known as reducible uncertainty because it can be decreased by acquiring more data or improving the model's capacity. Epistemic uncertainty is particularly important in situations where the model is making predictions on out-of-distribution or novel examples—scenarios where the model may be more prone to errors.

In light of these uncertainties, methods like Conformal Prediction (CP) [9] offer a principled framework for uncertainty quantification. CP is a distribution-free, model-agnostic approach that provides reliable confidence intervals for predictions. By not assuming any particular distribution for the data, CP can generate prediction intervals that account for both aleatoric and epistemic uncertainties, making it highly adaptable across different ML tasks and models.

Moreover, for classification tasks or probabilistic estimates, Venn Predictors, an extension of conformal prediction, are particularly useful. Venn Predictors provide a way to output multiple probability estimates. Unlike typical methods that output a single probability estimate (e.g., from softmax layers in neural networks), Venn Predictors yield a set of probabilities that, based on the current available data, ensures one of the probabilities will be correct. This added layer of uncertainty quantification offers more robust probabilistic estimates, which can be crucial in risk-sensitive applications like medical diagnostics or autonomous systems.

The strength of conformal methods, including Venn predictors, lies in their ability to complement standard ML outputs by offering uncertainty estimates that are both mathematically rigorous and practically useful. This enhances the interpretability of AI systems and provides decision-makers with a clearer understanding of the model's confidence, making these methods preferable for applications where understanding uncertainty is as important as the prediction itself.

A natural development has been to integrate conformal methods into explanation techniques, aiming to provide not only accurate predictions but also transparent and trustworthy explanations. Calibrated Explanations [10] exemplify this approach by

embedding uncertainty directly into the explanation process through calibration. By aligning the confidence of a model's prediction with the reliability of its explanation, this method are offering users a calibrated view of when the model's reasoning can be trusted. Calibrated Explanations provide local explanations with uncertainty quantification of both prediction and feature importances for individual instances and is applicable to both classification and regression.

Similarly, ConformaSight [11] builds on conformal prediction to provide a global explanation framework. It ensures robust explanation sets with guaranteed coverage, making it applicable across different models and robust even in the presence of noisy or perturbed data. This allows for consistent, high-confidence explanations, even in unpredictable environments.

This paper aims to combine the strengths of local explanations from Calibrated Explanations — which maintain applicability across classification and thresholded regression scenarios — with the perturbation approach employed in ConformaSight. Specifically, we preserve the flexibility of explaining predictions for both classification and thresholded regression outputs (e.g., explaining the probability of exceeding a threshold). This is enhanced by incorporating ConformaSight's perturbation approach, enabling the generation of computationally efficient factual explanations providing feature importance.

The resulting approach offers extremely fast explanations that can be highly advantageous in real-time decision-making processes, particularly where the outputs of these explanations are fed into subsequent steps or systems. Performance is essential for scenarios like machine teaching [12, 13], where explanation algorithms are required to function in real-time, ideally on resource-constrained platforms such as mobile devices. In scenarios requiring critical decisions within a short time frame, such as emergency response or automated monitoring, the ability to generate explanations quickly ensures that the system can continue operating seamlessly. Essentially, these fast explanations become an integral part of a continuous decision-making loop, where real-time insights flow directly into other processes for immediate action. Additionally, this could especially be crucial in time-series or sequential data tasks, where the significance of features may shift dynamically across different time points. In such cases, generating

4

a single static explanation for an entire series would be misleading, as the importance of each feature evolves over time. Hence, fast calibrated explanations can adapt to the changing context at each timestamp to ensure more accurate insights, enhancing the reliability of real-time predictive models.

In addition to delivering rapid explanations with feature importance, it also quantifies uncertainty for both the overall prediction and the individual contribution of each feature. This allows decision-makers to interpret not only the model's output but also the reliability of each contributing factor, providing deeper insights for risk-sensitive applications where understanding both prediction and explanation uncertainty is critical.

In the next section, thorough descriptions of the building blocks for this paper are provided. In Section 3, our contribution is described. The experimental setup is described in Section 4 whereas the results provide both evaluation results and a demonstration of its applicability. Section 6 wraps up the paper with a concluding discussion and pointers for future work.

## 2. Background

### 2.1. Post-Hoc Explanation Methods

In machine learning (ML), there are two main approaches to generating explanations. One method involves using inherently interpretable and transparent models, also known as Interpretable AI, which describes the internals of a system in a way that is understandable to humans [14]. Alternatively, post-hoc techniques can be used to explain complex models, making them suitable for explaining black-box models as well [15].

Post-hoc explanations involve creating simpler models that clarify how a complex model's predictions are linked to input features. These explanations can be either local, focusing on a single instance, or global, providing insights into the overall model. They often incorporate visual aids like feature importance plots, pixel representations, or word clouds to emphasize the key features, pixels, or words influencing the model's predictions. The importance of making black-box models explainable is highlighted

by [16], which discusses various interpretability methods for machine learning models. Moradi et al. [17] propose the Confident Itemsets Explanation (CIE) method, which uses highly correlated feature values to explain black-box classifiers by discretizing the decision space into smaller subspaces. Another interesting approach is introduced by [18], which focuses on the quality of explanations in AI systems, particularly in the medical field, by introducing the concept of causability in addition to explainability.

Research on proposed post-hoc methods is diverse and frequently tailored to specific tasks. For instance, [19] provides a visual explanation of black-box models by localizing the region responsible for a prediction. This approach is tested on a classification task by localizing entire object classes within an image. Another popular approach is the use of counterfactual explanations, which, similar to perturbation methods, vary the input space to understand how it affects the output. This approach, employed in a statistical fashion, is used by [20] to produce human-friendly interpretations on classification tasks.

Post-hoc explanations for classification and regression have some distinguishing characteristics due to the nature of the insights they offer. In classification, explanations involve predicting the class an instance belongs to from predefined classes, with probability estimates reflecting the model's confidence for each class. Techniques like SHAP [21], LIME [22], and Anchor [23] explore factors contributing to class assignment, often using feature importance, such as words in textual data or pixels in images. However, both SHAP and LIME can also be used to generate explanations for regression models. In regression, the focus is on predicting numerical values associated with instances without predefined classes. Explanations for regression models normally adapt techniques designed for classifiers by attributing features to predicted outputs.

Local explanations in classification often rely on probability estimates, which most machine learning models can generate to indicate the likelihood of each class. These probability estimates are commonly interpreted as a measure of prediction confidence. For example, in a binary classification scenario, a model predicting an instance belonging to the positive class with a probability estimate of 0.8 is considered more confident than one predicting the same instance with a probability estimate of 0.65. Probability estimates form the basis for local explanations in classification tasks, where the con-

fidence level indicates the likelihood of the positive class being the true class for that instance.

## 2.2. Calibration and Uncertainty Quantification

Basing decisions on accurate information is crucial in decision-making, placing an additional layer of requirements on predictive models to provide well-calibrated predictions and guarantees.

Conformal Prediction (CP) [9] is a distribution-free framework that offers prediction regions with guaranteed coverage, whose value and effectiveness have been demonstrated in numerous studies [24] [25]. Errors occur when the true target falls outside the predicted region. However, conformal predictors maintain automatic validity under exchangeability, resulting in an error rate of $\epsilon$ over time. Conformal regression (CR) provides prediction intervals with user-decided guaranteed coverage, and conformal predictive systems (CPS) [26] provide a conformal predictive distribution (CPD). The CPDs can be queried for intervals with guaranteed coverage, similar to but more dynamic than CR. This is done by defining intervals based on percentiles in the distribution so that a symmetric interval with 90% coverage can be achieved using the percentiles [5, 95]. The CPDs can also be queried for the probability of the actual instance value being below a user-given threshold, corresponding to the percentile of the threshold value in the distribution.

For classification, the focus is generally on the calibration of the probability estimates produced by the classifier, which can be defined as follows:

$$p(c \mid p^c) \approx p^c, \tag{1}$$

where $p^c$ represents the probability estimate for a particular class label $c$. This means that a well-calibrated model produces predicted probabilities that match observed accuracy. Consequently, whenever a model assigns a probability estimate of 0.9 to a label, the accuracy for that label should be approximately 90%.

It is well-known that many predictive models produce poorly calibrated probability estimates [27]. An external calibration method can be applied to calibrate a poorly

calibrated model using a separate portion of the labelled data, called the calibration set, to adjust the predicted probabilities.

The CP framework defines Venn [28] and Venn-Abers (VA) [29] predictors that produce multi-probabilistic predictions in the form of confidence-based probability intervals. Venn prediction involves a Venn taxonomy, categorising calibration data for probability estimation. The estimated probability for test instances falling into a category is the relative frequency of each class label among all calibration instances (including the test instance) in that category. Defining a proper Venn taxonomy can be challenging, which is the strength of VA.

### 2.2.1. Venn-Abers Calibration

VA Calibration offer automated taxonomy optimisation using isotonic regression, resulting in dynamic probability intervals for binary classification. Since the probability interval includes the well-calibrated probability estimates for the true class label being both negative (lower bound) and positive (upper bound), and the instance must be one or the other, it follows that the interval must contain the true probability. The problem from a predictive perspective is that the true class label is not known. However, the width and location of the interval can provide a lot of information. A smaller interval indicates higher certainty about the prediction, while a larger interval indicates more uncertainty. Since it is often impractical to have only an interval to indicate the probability estimate of the positive class, it is common to use a regularisation of the interval as an estimate for the positive class.

To define a VA predictor predicting a test object $x_{n+1}$, let $Z = \{z_1, \ldots, z_n\}$, where $n = l+q$, be a training set. Each instance $z_i = (x_i, y_i)$ consists of two parts, an object $x_i$ and a target $y_i$. Normally, calibration requires a separate calibration set, motivating a split of the training set into a proper training set $Z_l$ with $l$ instances and a calibration set $Z_q = \{z_1, \ldots, z_q\}$[1]. A scoring classifier is trained on $Z_l$ to compute $s$ for $\{x_1, \ldots, x_q, x_{n+1}\}$. The score $s$ is defined as the probability estimate for the positive class from a classifier

---

[1]As we assume random ordering, the calibration set is indexed $1, \ldots, q$ rather than $l+1, \ldots, n$, for indexing convenience.

*h*. Inductive VA prediction follows these steps:

1. Use $\{(s_1, y_1), \ldots, (s_q, y_q), (s_{n+1}, y_{n+1} = 0)\}$ to derive the isotonic calibrator $g_0$ and use $\{(s_1, y_1), \ldots, (s_q, y_q), (s_{n+1}, y_{n+1} = 1)\}$ to derive the isotonic calibrator $g_1$.

2. The probability interval for $y_{n+1} = 1$ is defined as $[g_0(s_{n+1}), g_1(s_{n+1})]$ (hereafter referred to as $[p_{low}, p_{high}]$, representing the lower and upper bounds of the interval).

3. The regularised probability estimate for $y_{n+1} = 1$, minimising the log loss [29], can be defined as:

$$p = \frac{p_{high}}{1 - p_{low} + p_{high}} \tag{2}$$

In summary, VA produces a calibrated (regularised) probability estimate $p$ together with a probability interval with a lower and upper bound $[p_{low}, p_{high}]$.

### 2.2.2. Conformal Predictive Systems

Conformal Predictive Systems produce CPDs for each test object $x_{n+1}$ when the target domain is numeric (i.e. regression). To define a CPS, assume the existence of an underlying regression model $h$ trained using $Z_l$. Like all conformal predictors, CPS relies on nonconformity scores $\alpha$, defining the strangeness of an instance. Unlike CR, where the nonconformity is usually defined as the absolute error $\alpha_i = |y_i - h(x_i)|$, CPS defines nonconformity using the signed errors $\alpha_i = y_i - h(x_i)$. The prediction for a test instance $x_{n+1}$ then becomes the following CPD:

$$CPD(y) = \begin{cases} \frac{i+\tau}{q+1}, & \text{if } y \in (C_{(i)}, C_{(i+1)}), \quad \text{for } i \in \{0, ..., q\} \\ \frac{i'-1+(i''-i'+2)\tau}{q+1}, & \text{if } y = C_{(i)}, \quad \text{for } i \in \{1, ..., q\} \end{cases} \tag{3}$$

where $C_{(1)}, \ldots, C_{(q)}$ are obtained from the calibration scores $\alpha_1, \ldots, \alpha_q$, sorted in increasing order:

$$C_{(i)} = h(x_{n+1}) + \alpha_i \tag{4}$$

with $C_{(0)} = -\infty$ and $C_{(q+1)} = \infty$. In case of a tie, $\tau$ is sampled from the uniform distribution $U(0, 1)$, and its role is to allow the $p$-values of target values to be uniformly distributed, $i''$ is the highest index such that $y = C_{(i'')}$, while $i'$ is the lowest index such that $y = C_{(i')}$.

The following cases provide some further intuition on how a CPD can be used:

- Obtaining a two-sided symmetric prediction interval for a chosen significance level $\epsilon$ can be done by $[C_{\lfloor(\epsilon/2)(q+1)\rfloor}, C_{\lceil(1-\epsilon/2)(q+1)\rceil}]$. Since the CPS has guaranteed coverage, the expected error of the obtained interval will be $\epsilon$ in the long run. Asymmetric prediction intervals are possible by selecting percentiles for the lower ($p^{low}$) and higher ($p^{high}$) bounds of the interval. The guaranteed coverage of the interval will be $\epsilon = p^{high} - p^{low}$.

- Still using the significance level $\epsilon$, a lower-bounded one-sided prediction interval can be obtained by $[C_{\lfloor\epsilon(q+1)\rfloor}, \infty]$, and an upper-bounded one-sided prediction interval can be obtained by $[-\infty, C_{\lceil(1-\epsilon)(q+1)\rceil}]$. The coverage guarantees still apply.

- Similarly, a point prediction corresponding to the median of the distribution can be obtained by $(C_{\lceil0.5(q+1)\rceil} + C_{\lfloor0.5(q+1)\rfloor})/2$. The median prediction can be seen as a calibration of the underlying model's prediction. Unless the model is biased, the median will tend to be very close to the prediction of the underlying model.

- For a specific threshold $t$, the distribution can return the estimated probability $p(C \leq t)$. Thus, it is possible to get the probability of the true target being below the threshold $t$.

A CPS offers richer opportunities to define intervals and probabilities through querying the CPD compared to CR. A particular strength is the ability to calibrate the underlying model. For example, if the underlying model is consistently overly optimistic, the median from the CPS will adjust for that and provide a calibrated prediction that is better adjusted to reality.

### 2.3. Calibrated Explanations

Calibrated Explanations is a recently released[2] local explanation method supporting both classification and regression, providing feature importance with uncertainty

---

[2]Calibrated Explanations can be installed using, e.g., `pip install calibrated-explanations` or accessed at github.com/Moffran/calibrated_explanations.

quantification [10, 30]. Calibrated Explanations produce instance-based explanations, and a *factual explanation* is composed of a *calibrated prediction* from the underlying model accompanied by an *uncertainty interval* and a collection of *factual feature rules*, each composed of a *feature weight with an uncertainty interval* and a *factual condition*, covering that feature's instance value. Calibrated Explanations support both binary and multi-class classification. In binary classification, the explanation explains the calibrated probability estimate (and its level of uncertainty) for the positive class, whereas in multi-class classification, the most probable class (after calibration) is considered the positive class and all other classes are treated as the negative class, i.e., not the predicted class. For regression, there are two alternative use cases:

1. The regression explanation explains a calibrated estimate of the prediction from the regressor, with a confidence interval covering the true target with a user-assigned level of confidence.

2. The thresholded explanation explains the calibrated probability estimate (and its level of uncertainty) for the calibrated estimate of the prediction being below a user-given threshold.

The algorithm's core is agnostic to whether it is a classification or regression problem since it is defined based on a numeric estimate and a lower and an upper bound defining an uncertainty interval for the numeric estimate. For classification, the probability estimate for the positive class is calibrated using a VA calibrator [29], producing a lower and an upper bound for the calibrated probability estimate (using a regularised mean of these bounds as the numeric estimate). For regression, a Conformal Predictive System (CPS) [31], producing a Conformal Predictive Distribution (CPD), is used as a calibrator of the underlying model. For the first use case, explaining the prediction value, the numeric estimate is the median from the CPD, and the lower and upper bounds are represented by user-selected percentiles in the CPD, defining the interval with guaranteed coverage. For the second use case, explaining the probability of being below a user-given threshold, the percentile in the CPD representing the threshold position is used as a probability estimate (similar to classification) upon which a VA calibrator is applied. For details on how thresholded regression works, see the original

11

regression paper by Löfström et al. [30].

Calibrated Explanations assume the existence of a predictive model $h$, trained using the proper training set $Z_l$, outputting a numeric value when predicting an object $h(x_i)$. For classification, the model is a scoring classifier, producing probability estimates for the positive class. For regression, it is an ordinary regressor predicting the expected value. Algorithm 1 describes how Calibrated Explanations creates a factual explanation of $x^3$.

For further details on the algorithm and how it is applied to classification and the two use cases for regression, see [10] and [30].

### 2.4. Perturbation Based Explanations

Yapicioglu et al. [11] presents ConformaSight, a unique explanation approach based on conformal prediction methodology that provides insightful and resilient explanations regardless of the underlying data distribution. Its purpose is to give explanations for set-type predictions that conformal predictors create. ConformaSight highlights the influence of the calibration process on prediction outputs, in contrast to conventional explanation approaches that mainly concentrate on feature importance.

In classification tasks, it is essential to understand the factors that influence the formation of prediction sets within conformal prediction frameworks to enhance model interpretability and trust[32]. Analyzing metrics such as weighted coverage and weighted set size provides valuable insights into the model's uncertainty representation [33]. Weighted coverage measures the proportion of instances that are correctly classified and included within prediction sets, thereby indicating the model's reliability in identifying uncertain regions relative to class distribution [34]. On the other hand, weighted set size offers critical information about the granularity of uncertainty representation, reflecting the average number of instances within prediction sets while accounting for class imbalance. By examining these metrics, researchers can gain a deeper understanding of the relationship between model predictions and input features.

---

[3]The index $n + 1$ is omitted to reduce clutter.

**Algorithm 1** Factual Calibrated Explanations

1: **Input:** Fitted model $h$, calibrator, test object $x$

2: **Output:** Factual explanation of $x$

3: **if** Classification **then**

4:     Use VA as calibrator to produce calibrated probability estimate $\varphi = p$ and uncertainty interval $\left[\varphi_{low} = p_{low}, \varphi_{high} = p_{high}\right]$.

5: **end if**

6: **if** Regression **then**

7:     Use CPS as calibrator and let $\varphi$ be the median and $\left[\varphi_{low}, \varphi_{high}\right]$ is either a one- or two-sided interval as described above.

8: **end if**

9: **for** each feature $f \in F$ **do**

10:     Changing the value of feature $f$, one at a time in a systematic way, producing slightly perturbed versions of object $x$, the calibrator can be used to estimate the (averaged) prediction $\varphi_f$ and uncertainty intervals $\left[\varphi_{low\_f}, \varphi_{high\_f}\right]$.

11:     The feature importance for feature $f$ is defined as the difference between the calibrated prediction $\varphi$, achieved on the original object $x$, and the estimated (averaged) calibrated prediction $\varphi_f$, achieved on the perturbed versions of $x$.

12:     The uncertainty intervals for the feature importance are defined analogously by calculating the difference between $\varphi$ and the uncertainty intervals $\left[\varphi_{low\_f}, \varphi_{high\_f}\right]$ for the perturbed versions of $x$.

13:     A factual feature rule is formed, with a factual condition defined as `feature = categorical instance value`, for categorical features, or `feature ≤ threshold` or `feature > threshold`, for numerical features. The `threshold` is defined so that the factual condition incorporates the numerical instance value for that feature. Since the factual condition must always include the feature value, only one factual condition is formed for each feature.

14: **end for**

15: **return** A factual explanation, composed of a *calibrated prediction* $\varphi$ from the underlying model accompanied by an *uncertainty interval* $\left[\varphi_{low}, \varphi_{high}\right]$ and the collection of *factual feature rules*.

Furthermore, ConformaSight focuses on how the calibration process affects prediction intervals, which goes beyond conventional feature importance explanations. A feature in the calibration set is important in influencing the model's confidence in its predictions when it has a significant impact on the coverage of these intervals. This method takes advantage of the finding that, when the conformal prediction algorithm is run multiple times with distinct calibration datasets, the coverage will differ over a limitless number of validation points and yet satisfy the $1 - \alpha$ (error rate) minimum coverage requirement [35]. This variability highlights the significance of some properties in the calibration set, especially those that have a large impact on the coverage of prediction intervals, suggesting that these factors are critical in determining the model's level of confidence in its predictions.

As a result, ConformaSight provides explanations that highlight the significance of features as well as the calibration process's impact on forecast accuracy. Metrics like *weighted coverage* and *weighted set size* are used to achieve this; they are particularly useful in conformal classification situations [9], offering a more thorough comprehension of the dynamics of the model.

Central to ConformaSight is the idea of leveraging distributional changes within the calibration set to systematically detect variations that highlight feature importance. These changes are introduced through perturbations, mathematically defined in Sections 2.4.1, 2.4.2, and 2.4.3.

### 2.4.1. Definition 1: Permutation-based Perturbations

Let the set of categorical features be represented by $C_F$. Let $x_{C_f}$ represent a categorical feature with $c$ different categories for each $C_f \in C_F$. By arbitrarily permuting the values of $x'_{C_f}$ $k$ times, the permutation-based perturbation function `permute`$(x_{C_f}, \ k)$ creates a perturbed variant of $x_{C_f}$, called $x'_{C_f}$. Formally, this procedure can be explained as:

$$x'_{C_f} = \texttt{permute}(x_{C_f}, k) \tag{5}$$

### 2.4.2. Definition 2: Gaussian Noise Perturbations

Define the set of numerical characteristics by $N_F$. Consider a numerical feature $x_{N_f}$ with a standard deviation $\sigma$ for each $N_f \in N_F$. The Gaussian noise perturbation algorithm takes a severity parameter $s$ and adds noise to $x_{N_f}$, sampled from a normal distribution; this produces the perturbed feature $x'_{N_f}$. Theoretically this procedure is defined as:

$$x'_{N_f} = x_{N_f} + \eta(s), \quad \text{where } \eta(s) \sim \texttt{Normal}(0, s \times \sigma) \tag{6}$$

### 2.4.3. Definition 3: Uniform Noise Perturbations

The uniform noise perturbation function adds noise from a uniform distribution to a numerical feature $x_{N_f}$. The definition of the perturbed feature $x'_{N_f}$ given a severity parameter $s$ is as outlined below:

$$x'_{N_f} = x_{N_f} + \eta(s), \quad \text{where } \eta(s) \sim \texttt{Uniform}(-s \times R_{N_f}, s \times R_{N_f}) \tag{7}$$

where $R_{N_f}$ represents the range of values in $x_{N_f}$. The uniform noise perturbation introduces variability across the dataset, facilitating the exploration of various distributional shifts.

By implementing these perturbations on a feature-by-feature basis, ConformaSight assesses the changes in the model's behavior relative to the original calibration set. This alteration acts as a measure of each feature's significance, offering enhanced insights into how individual features impact the model's predictions.

## 3. Proposed Solution

This paper aims to incorporate the perturbation approach, constituting a core element in ConformaSight, into the explanation method Calibrated Explanations, making it possible to extract a new form of explanations, providing fast feature importance generation without rule conditions.

As described in 2.3, perturbations are done for each feature of the test instance at explanation time in Calibrated Explanations. The solution that we propose in this paper

performs all perturbations on the calibration set at initialisation of the Fast Calibrated Explanations, resulting in some additional overhead once when initialising Fast Calibrated Explanations, while avoiding any perturbations at explanation time. Compared to the explanations in Calibrated Explanations, a main difference with the solution proposed here is that there is no rule conditions. Instead, each feature is assigned a feature weight determining the relative importance of that feature compared to other features. As such, the provided explanations are factual, conveying the feature importance per instance. The resulting solution provides very fast explanations with feature weights that can be analysed per instance.

More formally, the solution that we propose can be divided into two stages: 1) initialisation, and 2) explanation. The initialisation stage 1) is described in Algorithm 2 while the explanation stage 2) is described in Algorithm 3.

---

**Algorithm 2** Initialisation of Fast Calibrated Explainer

1: **Input:** Calibration set $Z_q$, factor $k$, severity $s$, noise $\eta$

2: **Output:** An initialised Fast Calibrated Explainer

3: Multiply $Z_q$ with factor $k$ resulting in $\mathbf{Z}_q = [Z_q]_{i=1}^k$

4: **for** each feature $f \in F$ **do**

5:     Permute the $k$ copies of $X_f$ to create a permuted $\mathbf{X}'_f$ using Equation (5), (6), or (7).

6:     Initiate a calibrator $C_f$ using the multiplied calibration set $\mathbf{Z}_q$, substituting the $k$ original $X_f$ with $\mathbf{X}'_f$. The kind of calibrator is either a VA for classification or a combination of CPS and VA for thresholded regression explanations.

7: **end for**

8: Initiate a base calibrator $C$ using the original calibration set $Z_q$

9: **return** A new Fast Calibrated Explainer with all permutations and calibrators stored

---

Since the perturbation is performed at initialisation and is using a CPD when used for regression, all explanations from a single model that use the same uncertainty interval, i.e., the same percentiles, will result feature weights that have exactly the same size in relation to all other feature weights. The only thing that will differ from instance

**Algorithm 3** Explanation using Fast Calibrated Explainer

1: **Input:** An initialised Fast Calibrated Explainer, a test object $x$

2: **Output:** An explanation of $x$

3: Use the base calibrator $C$ to produce calibrated estimate $\varphi$ and uncertainty interval $\left[\varphi_{low}, \varphi_{high}\right]$ for $x$.

4: **for** each feature $f \in F$ **do**

5:     Estimate the prediction $\varphi_f$ and uncertainty intervals $\left[\varphi_{low\_f}, \varphi_{high\_f}\right]$ using $C_f$.

6:     The feature weights for feature $f$ is defined as the difference between the calibrated prediction $\varphi$, achieved on the original object $x$, and the estimated (averaged) calibrated prediction $\varphi_f$.

7:     The uncertainty intervals for the feature weights are defined analogously by calculating the difference between $\varphi$ and the uncertainty intervals $\left[\varphi_{low\_f}, \varphi_{high\_f}\right]$.

8: **end for**

9: **return** An explanation, composed of a *calibrated prediction* $\varphi$ from the underlying model accompanied by an *uncertainty interval* $\left[\varphi_{low}, \varphi_{high}\right]$ and the collection of *feature weights* with *feature weight uncertainties*.

to instance is the scale of the feature weights, making Fast Calibrated Explanations for standard regression clearly less useful. The same issue does not exist for probabilistic explanations (applicable to both classification and thresholded regression).

The similarities and differences between Calibrated Explanations and Fast Calibrated Explanations are summarised in Table 1.

## 4. Experimental Setup

The evaluation is divided into two parts. The first part contains a comparative evaluation between the Fast Calibrated Explanations, our proposed solution, and Calibrated Explanations on classification and regression problems. The evaluation is performed separately for classification and regression, focusing on complementary aspects. For code, see the calibrated_explanations/evaluation/FastCE folder in the repository.

For classification, an ablation study of the impact from possible permutation parameters is performed. The evaluation covers both computational cost and how the mean variance of feature weights vary across different parameter settings. The ablation study includes various parameter values for the Fast Calibrated Explanations including both forms of noise type (*uniform* and *gaussian*), four different scaling factors $(1, 3, 5, 10)$, and five different severity values $(0, 0.25, 0.5, 0.75, 1)$. Computational time and mean variance per feature importance is reported for 25 binary classification data sets. Each data set was split into 50% training data, 25% calibration data and 25% test data and the underlying model was a `RandomForestClassifier` with 100 trees[4].

For regression, focus is on computational speed, stability and robustness in comparison with SHAP and LIME applied on calibrated models (using the CPS as calibrator of the underlying model). All target values were min-max normalised to the range $[0, 1]$. Each data set was split using 200 calibration instances, 100 test instances, and the remaining instances as training set. For each one of the 31 data sets, six different setups where evaluated, three for standard regression and two for thresholded regression (using 0.5 as threshold). Two setups applied SHAP and LIME on a calibrated model, two

---

[4]To run the classification experiment, run first *Perturbation_Experiment_Ablation.py* to generate results and then run the notebook *Perturbation_Analysis_Times.ipynb*.

Table 1: Comparison between Calibrated Explanations and Fast Calibrated Explanations

|  | Calibrated Explanations | Fast Calibrated Explanations |
|---|---|---|
| Perturbation | Perturbation is done at explanation time on the test instance | Perturbation is done at initialisation time on the calibration set, requiring additional memory space for the perturbed calibration set |
| Expressiveness | Each rule contains a condition for which the feature weight applies | No condition is used, with the implication that the feature weight is less expressive |
| Interpretation | Each rule condition clearly convey when the feature weight apply, providing clear cues for interpretation | The feature weights provide insights on how much and in which direction a feature affects the prediction, with positive weights favouring the positive class and vice versa |
| Regression | Supports both standard regression and thresholded regression | Is only clearly useful for thresholded regression |
| Uncertainty | Uncertainty quantification is provided for both prediction and feature weights | Uncertainty quantification is provided for both prediction and feature weights |
| Alternatives | Alternative explanations, indicating what prediction the model would output if the feature is altered in accordance with the condition, can be extracted | No alternative explanations can be extracted |
| Conjunctions | Conjunctive explanations, indicating the joint impact from several conditional rules, are possible | No conjunctions are available |
| Speed | Since perturbations are done on the test instance and since a suitable rule condition must be identified, the explanation time is penalised | Since all perturbations have already been done on the calibration set and no condition is used, explanations can be generated much faster |

setups used Calibrated Explanations for standard and thresholded regression and two setups used Fast Calibrated Explanations for standard and thresholded regression. Fast Calibrated Explanations is evaluated for standard regression, to indicate the potential for that kind of context, even though these explanations are deemed less useful. The underlying model was a `RandomForestRegressor` which was calibrated using a CPS, either explicitly (for LIME and SHAP) using the `ConformalPredictiveSystem` class from *crepes* or implicitly as part of using the `CalibratedExplainer` class[5]. The following metrics are evaluated:

- *Stability* means that multiple runs on the same instance and model should produce consistent results. Stability is evaluated by generating explanations for the same predicted instances 5 times with different random seeds (using the iteration counter as random seed). The random seed is used to initialise the `numpy.random.seed()` and by the discretizers. The largest variance in feature weight (or feature prediction estimate) can be expected among the most important features (by definition of having higher absolute weights). The top feature for each test instance is identified as the feature being most important most often in the 100 runs (i.e., the mode of the feature ranks defined by the absolute feature weight). The variance for the top feature is measured over the 100 runs and the mean variance among the test instances is reported.

- *Robustness* means that small variations in the input should not result in large variations in the explanations. Robustness is measured in a similar way as stability, but with the training and calibration set being randomly drawn and a new model being fitted for each run, creating a natural variation in the predictions of the same instances without having to construct artificial instances. Again, the variance of the top feature is used to measure robustness. The same setups as for stability are used except that each run use a new model and calibration set and that the random seed was set to 42 in all experiments.

---

[5]To run the regression experiment, run first *Regression_Experiment.py* to generate results and then run the notebook *Regression_Analysis.ipynb*.

- *Explanation time* is compared between the setups regarding explanation generation times (in seconds per instance). It is only the method call resulting in an explanation that is measured. Any overhead in initiating the explainer class is not considered).

Table 2: Explanation time in seconds per instance for binary classification datasets

| Dataset | #Features | CE | FCE |
|---|---|---|---|
| colic | 60 | .369 | .004 |
| creditA | 43 | .594 | .002 |
| diabetes | 9 | .034 | .000 |
| german | 28 | .210 | .001 |
| haberman | 4 | .004 | .001 |
| heartC | 23 | .058 | .002 |
| heartH | 21 | .045 | .002 |
| heartS | 14 | .032 | .001 |
| hepati | 20 | .034 | .003 |
| iono | 34 | .236 | .003 |
| je4042 | 9 | .021 | .001 |
| je4243 | 9 | .028 | .001 |
| kc1 | 22 | .340 | .001 |
| kc2 | 22 | .121 | .002 |
| kc3 | 40 | .475 | .003 |
| liver | 7 | .011 | .001 |
| pc1req | 9 | .009 | .002 |
| pc4 | 38 | 1.336 | .002 |
| sonar | 61 | .833 | .007 |
| spect | 23 | .032 | .003 |
| spectf | 45 | .421 | .004 |
| transfusion | 5 | .008 | .000 |
| ttt | 28 | .168 | .001 |
| vote | 17 | .045 | .001 |
| wbc | 10 | .038 | .001 |
| **Mean** | **24.0** | **.220** | **.002** |

## 5. Experimental Results

### 5.1. Performance Evaluation for Classification and Regression

#### 5.1.1. Comparison of Computation Time

The code used for the evaluation of Calibrated Explanations (CE) and Fast Calibrated Explanations (FCE) can be found in the FastCE folder in the repository. As the difference in initialisation time is almost negligible, the results are not presented here but can be found in the folder above. The explanation time clearly differs between Calibrated Explanations (CE) and Fast Calibrated Explanations (FCE). However, as the explanation time between different parameter settings for FCE does not differ much, only the default parameters (noise type=*uniform*, scale factor=5, and severity=0.5) are compared with CE in Table 2.

Table 3: Average explanation speedup factor for Fast Calibrated Explanations compared to Calibrated Explanations for different parameter settings

| Noise Type | gaussian | | | | uniform | | | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| Scale Factor | 1 | 3 | 5 | 10 | 1 | 3 | 5 | 10 | |
| Severity | | | | | | | | | Mean |
| 0 | 114 | 105 | 110 | 113 | 120 | 110 | 107 | 115 | **112** |
| 0.25 | 96 | 110 | 112 | 118 | 103 | 109 | 110 | 114 | **109** |
| 0.5 | 94 | 115 | 112 | 114 | 110 | 116 | <u>112</u> | 117 | **111** |
| 0.75 | 114 | 108 | 116 | 112 | 109 | 114 | 109 | 117 | **113** |
| 1 | 116 | 109 | 114 | 113 | 115 | 116 | 115 | 110 | **114** |
| **Mean** | **107** | **109** | **113** | **114** | **111** | **113** | **111** | **115** | **112** |

As can be seen, the explanation time in seconds per instance is on average more than 100 times higher for CE compared to FCE. It is worth noting that the average explanation time for CE is heavily influenced by a few really costly data sets. The median explanation time for CE is 0.045, being only about 30 times higher for CE compared to FCE. The average speedup across the various evaluated setups is shown in Table 3 (with the underlined result representing the result in Table 2). The average speedup is similar across the various evaluated settings, with some minor deviations.

Table 4: Average initialisation computational cost factor for Fast Calibrated Explanations compared to Calibrated Explanations for different parameter settings

| Noise Type | gaussian | | | | uniform | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Scale Factor | 1 | 3 | 5 | 10 | 1 | 3 | 5 | 10 | |
| Severity | | | | | | | | | Mean |
| 0 | 20.9 | 21.9 | 23.2 | 26.7 | 20.3 | 22.4 | 23.2 | 26.5 | 23.1 |
| 0.25 | 24.1 | 21.7 | 24.1 | 26.7 | 21.1 | 23.1 | 23.9 | 27.4 | 24.0 |
| 0.5 | 21.0 | 21.9 | 23.8 | 27.7 | 21.7 | 22.4 | 24.8 | 27.0 | 23.8 |
| 0.75 | 20.2 | 22.2 | 24.4 | 27.4 | 20.7 | 22.7 | 24.8 | 26.9 | 23.7 |
| 1 | 19.7 | 22.1 | 23.3 | 27.3 | 19.8 | 22.5 | 24.2 | 27.7 | 23.3 |
| Mean | 21.2 | 22.0 | 23.8 | 27.2 | 20.7 | 22.6 | 24.2 | 27.1 | 23.6 |

So, what is the trade-off that we need to make to get this speedup? Table 4 show the computational cost in terms of how many times faster the initialisation of Calibrated Explanations is compared to Fast Calibrated Explanations. The average initialisation time for Calibrated Explanations is $8.3e - 5$, as a comparison.

On average, the one-time cost for initialisation is that it takes about 25 times longer for Fast Calibrated Explanations compared to Calibrated Explanations. It means that for the default setup, the average initialisation time for Fast Calibrated Explanations across all data sets was 0.0013 seconds, whereas the average explanation time per instance across all data sets is 0.0020 seconds. So even if the initialisation time for Fast Calibrated Explanations is substantially larger than for Calibrated Explanations, it is still negligible considering it is only done once.

In the evaluation for regression, a comparison is made with calibrated SHAP and LIME. The explanation time in seconds per instance is tabulated in Table 5.

The speedup is not as impressive for regression as it was for classification, but it is still substantial. FCE is almost 19 times faster for standard regression (CE) and Probabilistic Fast Calibrated Explanations (PFCE) is more than 8 times faster than Probabilistic Calibrated Explanations (PCE). FCE is more than 75 times faster than calibrated LIME and more than 200 times faster than calibrated SHAP. The relatively smaller speedup observed in regression can be attributed to the shorter baseline explanation times for CE in regression tasks when compared to classification scenarios.

Table 5: Explanation time in seconds per instance for regression datasets

| Dataset | LIME | SHAP | CE | FCE | PCE | PFCE |
|---|---|---|---|---|---|---|
| abalone | .071 | .239 | .021 | .001 | .022 | .003 |
| anacalt | .044 | .044 | .014 | .001 | .015 | .002 |
| bank8fh | .109 | .267 | .022 | .002 | .024 | .004 |
| bank8fm | .102 | .252 | .021 | .001 | .023 | .004 |
| bank8nh | .108 | .269 | .022 | .001 | .023 | .003 |
| bank8nm | .107 | .265 | .022 | .001 | .023 | .003 |
| comp | .093 | .632 | .041 | .002 | .042 | .004 |
| concrete | .069 | .180 | .020 | .001 | .021 | .003 |
| cooling | .071 | .155 | .019 | .001 | .019 | .002 |
| deltaA | .114 | .052 | .011 | .001 | .013 | .002 |
| deltaE | .141 | .070 | .014 | .001 | .015 | .003 |
| friedm | .098 | .042 | .010 | .001 | .012 | .002 |
| heating | .071 | .154 | .019 | .001 | .020 | .002 |
| kin8fh | .113 | .313 | .024 | .002 | .027 | .004 |
| kin8fm | .117 | .326 | .025 | .002 | .029 | .004 |
| kin8nh | .116 | .320 | .024 | .002 | .026 | .004 |
| kin8nm | .116 | .334 | .025 | .002 | .025 | .005 |
| laser | .083 | .027 | .008 | .001 | .008 | .002 |
| mg | .108 | .071 | .014 | .001 | .015 | .002 |
| mortage | .076 | .459 | .064 | .002 | .066 | .003 |
| plastic | .063 | .016 | .003 | .000 | .004 | .001 |
| puma8fh | .115 | .291 | .023 | .001 | .024 | .004 |
| puma8fm | .108 | .265 | .022 | .002 | .025 | .004 |
| puma8nh | .108 | .276 | .022 | .002 | .025 | .004 |
| puma8nm | .108 | .271 | .022 | .002 | .024 | .004 |
| quakes | .077 | .023 | .005 | .001 | .006 | .001 |
| stock | .067 | .361 | .026 | .001 | .027 | .003 |
| treasury | .075 | .448 | .063 | .001 | .065 | .003 |
| wineRed | .075 | .485 | .036 | .001 | .037 | .003 |
| wineWhite | .087 | .579 | .035 | .001 | .037 | .004 |
| wizmir | .073 | .231 | .026 | .001 | .028 | .003 |
| **Mean** | **.093** | **.249** | **.023** | **.001** | **.025** | **.003** |

### 5.1.2. Stability and Robustness

Both stability and robustness have been evaluated for the regression data. The mean and median stability and robustness aggregated over all regression data sets are shown in Table 6. The motivation for using both aggregation methods is that some data sets deviated drastically from the general picture for some methods, having a huge impact on the mean but not on the median. For extremely low variation (less than $1e-30$), the value was set to 0. Detailed result per data set can be found in the evaluation folder. The robustness must be compared to the amount of variability in the predictions from the underlying model, which was $1.7e-4$ on average.

Table 6: Mean and median stability and robustness aggregated over all regression data sets.

|  |  | LIME | SHAP | CE | FCE | PCE | PFCE |
|---|---|---|---|---|---|---|---|
| **Stability** | Mean | 1.3e-5 | 1.6e-7 | 0 | 1.6e-5 | 4.8e-3 | 4.2e-3 |
|  | Median | 8.8e-6 | 0 | 0 | 6.0e-6 | 3.9e-3 | 3.8e-3 |
| **Robustness** | Mean | 1.1e-3 | 2.5e-4 | 7.1e-3 | 4.8e-5 | 2.3e-2 | 6.7e-3 |
|  | Median | 5.6e-4 | 1.1e-4 | 2.9e-3 | 3.2e-5 | 1.9e-2 | 6.4e-3 |

Lets start by considering the standard regression results from LIME, SHAP, CE and FCE, which can all be compared. Considering the results reported in [30], the results for LIME, SHAP and CE are as expected. It is clear that stability is worse for FCE than for CE and SHAP but comparable to LIME. Furthermore, the robustness of FCE is even better than all of the other three methods. Looking at the probabilistic results, it does not make sense to compare these results with the previously mentioned results, as their predictions are probabilities. Both PCE and PFCE are less stable and robust and the reason is related to the sensitivity of the probabilities derived from the CPD. The reason for the sensitivity is that a relatively small change in prediction can easily result in a comparably much larger change in probability for exceeding the threshold, especially if the target is close to the threshold (which is set to 0.5, i.e., the mid-point in the interval of possible target values). Results are comparable between PCE and PFCE.

### 5.2. Demonstration

In the demonstration below, a few different data sets are included, to show examples from different use cases. The examples will illustrate how the explanations may

```python
from calibrated_explanations import WrapCalibratedExplainer
# Load and pre-process your data
# Divide it into proper training, calibration, and test sets

# Initialize the WrapCalibratedExplainer with your model
classifier = WrapCalibratedExplainer(ClassifierOfYourChoice())
regressor = WrapCalibratedExplainer(RegressorOfYourChoice())

# Train your model using the proper training set
classifier.fit(X_proper_training_cl, y_proper_training_cl)
regressor.fit(X_proper_training_reg, y_proper_training_reg)

# Calibrate your model using the calibration set
# Ensure fast initialisation by assigning fast=True
classifier.calibrate(X_calibration_cl, y_calibration_cl, fast=True)
regressor.calibrate(X_calibration_reg, y_calibration_reg, fast=True)

# Create and plot fast explanations for classification
explanations = classifier.explain_fast(X_test_cl)
explanations.plot(uncertainty=True)

# Create and plot fast explanations for thresholded regression
my_threshold = 500
explanations = regressor.explain_fast(X_test_reg,
                                      threshold=my_threshold)
explanations.plot(uncertainty=True)
```

Figure 1: Code example on using *calibrated-explanations* for fast explanations

look and how they can be understood. The demonstration will include both binary and multi-class data sets, as well as a regression data set for which a thresholded regression explanation is given. Since the calibrated probability of the test instances can be retrieved by the explainer and very certain predictions without much uncertainty will generally not prompt further investigation, all the examples are for less certain predictions. In order to produce an explanation plot conveying uncertainty, code similar to the example in Figure 1 can be used.



Figure 2: A fast explanation for an instance from the Wine data set

To begin the demonstration, the binary classification *wine* data set is used. Figure 2 shows an instance which is predicted as likely to be positive (indicated by the red bar at the top). The plot box at the lower part of the figure provides the contribution by each of the features to the prediction. To the right of the plot box, the feature names are written and to the right of the plot box, the actual feature values of the instance are shown. The interpretation is that the instance value for *Volatile Acidity* make the probability for the positive class be higher, and the instance value for *Citric Acid* make the probability for the positive class be slightly lower, and so on. Each of the bars represent the importance of a particular feature and is indicative of that features importance for this instance. The grey area in the background correspond to the light red area at the top, indicating the uncertainty interval of the calibrated probability for the positive class. The lighter red
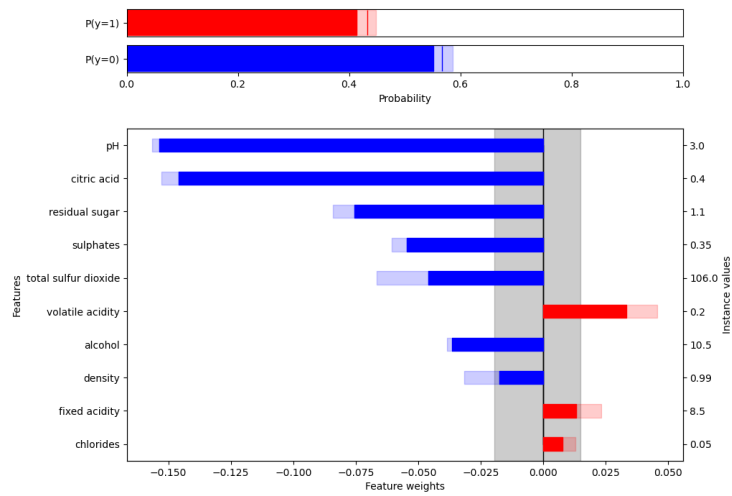
Figure 3: A fast explanation for another instance from the Wine data set

or blue areas on each feature weight bar indicate the uncertainty interval for the feature weight. This particular instance is in fact class 1, which several of the features indicate.

Another example for the same data set can be seen in Figure 3. In this example, the presence of several features each contribute (individually) to lower probability for the positive class (indicated by negative weights, which is shown as blue bars). This particular instance is in fact class 0.

Looking instead at an example from the multi-class *Glass* data set, the explanation provides the probability for the class being most likely after calibration. As is exemplified in the paper introducing multi-class Calibrated Explanations [36], it is advisable to compare the explanations against a confusion matrix, to get an understanding of the typical errors made by the model. Here, it suffices to say that both precision and recall for the class *build wind non-float* is just above 0.7. In this case, as seen in Figure 4, several of the features are blue, indicating that they contribute to lower probability for the predicted class, strengthening the indication provided by the probability against the predicted class. One of the feature weights, for *Na*, has a large degree of uncertainty, indicating that the weight can be either rather low or rather high.

The final two examples are taken from the *California Housing* data set with the threshold set to 490 (which is the median house price in the data set). The example in
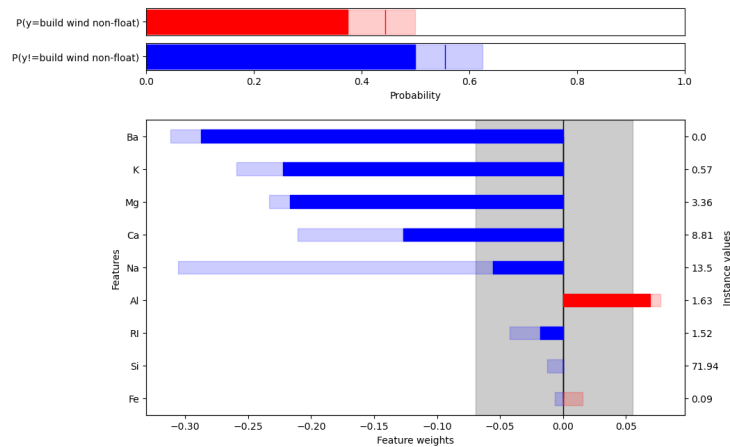
Figure 4: A fast explanation for an instance from the Glass data set

Figure 5 is likely to be above the threshold with a fairly low uncertainty. One of the features, the *grade*, is favouring the prediction in the sense that if this feature would have been randomly assigned (as it is in the calibration set), it would have resulted in reduced probability $\mathcal{P}(y > 490)$. Similarly, there are several features that indicates the opposite impact, even if several of them indicate a lot of uncertainty. It is, however, important to remember that the feature weights are not cumulative, their impact cannot be stacked upon each other. Instead, each feature must be considered by itself. All in all, this is a difficult explanation to interpret, as there are several features that indicate that they favour a lower price. At the same time, the initial probability $\mathcal{P}(y > 490)$ is still rather high, reducing $\mathcal{P}(y > 490)$ with up to 20 percentage points (the highest positive weights), would still favour a higher price. In this case, the house was sold for $600K$.

The final example, seen in Figure 6, is also from the same data set, using the same threshold. Here the probability is high for a lower price, even though the uncertainty is also high. One of the feature weights is considerably larger than the other, indicating that the location along the *latitude* is an important feature for the prediction of this instance, strongly favouring the likelihood for a low price. None of the other feature weights have nearly the same impact. In this case, the house was sold for $350K$, which seem reasonable given the indication of the model and the explanation.
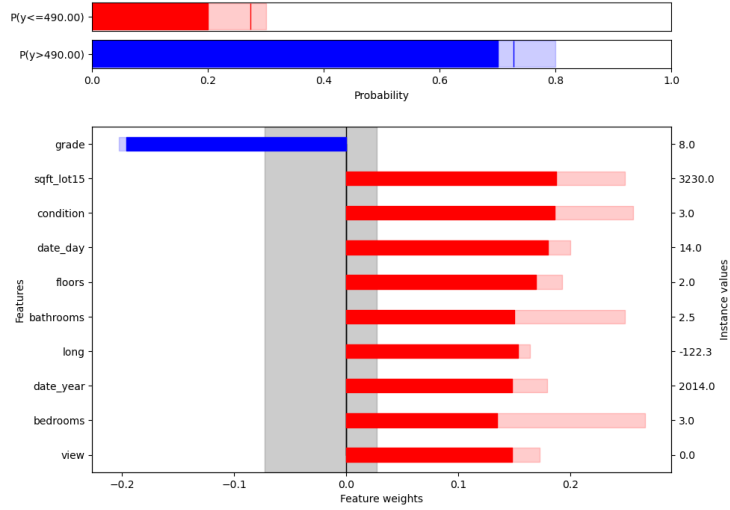
29

Figure 5: A fast explanation for an instance from the Housing data set, with the threshold 490
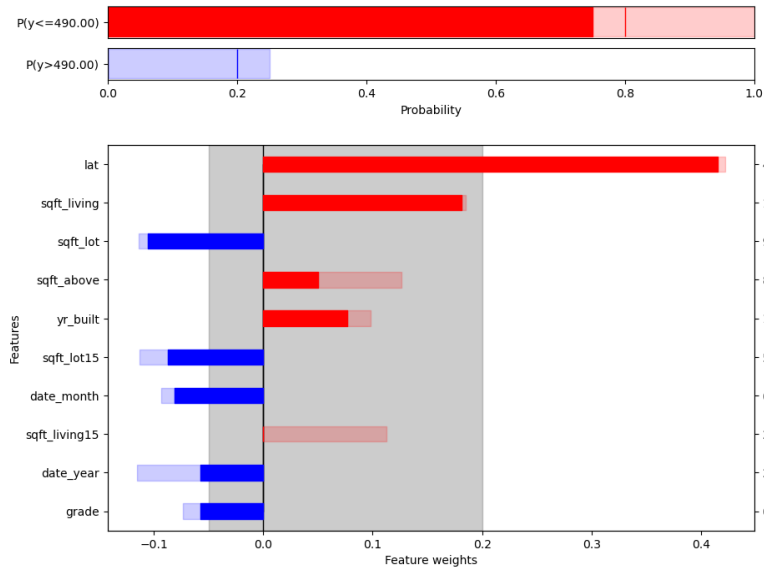


Figure 6: A fast explanation for another instance from the Housing data set, with the threshold 490

These examples have served to indicate how fast explanations can be used to indicate which features that are important. However, it is worth noting that this is all they can do, they cannot provide the user with any additional insights on why or in what way a feature is important. Compared with Calibrated Explanations, having the possibility to answer the why and how questions with their slower explanations, Fast Calibrated Explanations have the benefit of being very fast at the cost of being less informative.

## 6. Concluding Discussion

This paper address a common drawback among explanation methods, namely the computational overhead of explaining an instance. The proposed solution combines fundamental building blocks from two recently proposed explanation methods: Calibrated Explanations and ConformaSight. The proposed solution have used the perturbation strategy used in ConformaSight in combination with the explanation engine in Calibrated Explanations, allowing really fast local explanations with uncertainty quantification for both classification and thresholded regression. The uncertainty quantification extend both the calibrated predictions and the provided feature weights. Fast Calibrated Explanations is able to take advantage of Calibrated Explanations support for both binary and multi-class classification, as well as thresholded regression, providing the probability of the true target being above a user-given threshold.

Possible directions for future work include considering the issue of perturbations outside the natural scope of the data set as well as ways of speeding up Calibrated Explanations using insights from Fast Calibrated Explanations. Another important area for future work is to consider the decision-making aspect, focusing on situations where fast explanations are critical, exploring how our solution can help create trustworthy explanations in such use cases. Currently, Fast Calibrated Explanations does not convey much insights for standard regression, which should be addressed in future development.

31

## Acknowledgement

## Declaration of generative AI and AI-assisted technologies in the writing process

During the preparation of this work the author(s) used ChatGPT and Grammarly in order to improve language and style. After using these tools/services, the author(s) reviewed and edited the content as needed and take(s) full responsibility for the content of the published article.

## References

[1] S. Das, G. K. Nayak, L. Saba, M. Kalra, J. S. Suri, S. Saxena, An artificial intelligence framework and its bias for brain tumor segmentation: A narrative review, Computers in biology and medicine 143 (2022) 105273.

[2] A. Albahri, A. M. Duhaim, M. A. Fadhel, A. Alnoor, N. S. Baqer, L. Alzubaidi, O. Albahri, A. Alamoodi, J. Bai, A. Salhi, et al., A systematic review of trustworthy and explainable artificial intelligence in healthcare: Assessment of quality, bias risk, and data fusion, Information Fusion (2023).

[3] K. Devitt, M. Gan, J. Scholz, R. Bolia, A method for ethical ai in defence (2021).

[4] J. Zou, L. Schiebinger, Ai can be sexist and racist–it's time to make it fair, Nature 559 (2018) 324–327.

[5] D. Gunning, D. W. Aha, Darpa's explainable artificial intelligence program, AI Magazine 40 (2019) 44–58.

[6] Y. Romano, R. F. Barber, C. Sabatti, E. Candès, With malice toward none: Assessing uncertainty via equalized coverage, Harvard Data Science Review 2 (2020) 4.

[7] F. Wang, L. Cheng, R. Guo, K. Liu, P. S. Yu, Equal opportunity of coverage in fair regression, Advances in Neural Information Processing Systems 36 (2024).

[8] E. Hüllermeier, W. Waegeman, Aleatoric and epistemic uncertainty in machine learning: An introduction to concepts and methods, Machine learning 110 (2021) 457–506.

[9] V. Vovk, A. Gammerman, G. Shafer, Algorithmic Learning in a Random World, Springer-Verlag, Berlin, Heidelberg, 2005.

[10] H. Löfström, T. Löfström, U. Johansson, C. Sönströd, Calibrated explanations: With uncertainty information and counterfactuals, Expert Systems with Applications (2024) 123154.

[11] F. R. Yapicioglu, A. Stramiglio, F. Vitali, Conformasight: Conformal prediction-based global and model-agnostic explainability framework, in: World Conference on Explainable Artificial Intelligence, Springer, 2024, pp. 270–293.

[12] Y. Goyal, Z. Wu, J. Ernst, D. Batra, D. Parikh, S. Lee, Counterfactual visual explanations, in: International Conference on Machine Learning, PMLR, 2019, pp. 2376–2384.

[13] X. Zhu, A. Singla, S. Zilles, A. N. Rafferty, An overview of machine teaching, arXiv preprint arXiv:1801.05927 (2018).

[14] L. H. Gilpin, D. Bau, B. Z. Yuan, A. Bajwa, M. Specter, L. Kagal, Explaining explanations: An overview of interpretability of machine learning, in: 2018 IEEE 5th International Conference on data science and advanced analytics (DSAA), IEEE, 2018, pp. 80–89.

[15] E. M. Kenny, C. Ford, M. S. Quinn, M. Keane, Explaining black-box classifiers using post-hoc explanations-by-example: The effect of explanations and error-rates in xai user studies, Artif. Intell. 294 (2021) 103459.

[16] C. Molnar, Interpretable Machine Learning, 2 ed., 2022. URL: `https://christophm.github.io/interpretable-ml-book`.

[17] M. Moradi, M. Samwald, Post-hoc explanation of black-box classifiers using confident itemsets, Expert Systems with Applications 165 (2021) 113941.

[18] A. Holzinger, G. Langs, H. Denk, K. Zatloukal, H. Müller, Causability and explainability of artificial intelligence in medicine, Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery 9 (2019) e1312.

[19] S. M. Muddamsetty, M. N. Jahromi, A. E. Ciontos, L. M. Fenoy, T. B. Moeslund, Visual explanation of black-box model: Similarity difference and uniqueness (sidu) method, Pattern Recognition 127 (2022).

[20] H.-G. Jung, S.-H. Kang, H.-D. Kim, D.-O. Won, S.-W. Lee, Counterfactual explanation based on gradual construction for deep networks, Pattern Recognition 132 (2022).

[21] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, in: Proceedings of the 31st international conference on neural information processing systems, 2017, pp. 4768–4777.

[22] M. T. Ribeiro, S. Singh, C. Guestrin, Why should i trust you?: Explaining the predictions of any classifier, in: Proceedings of the 22nd ACM SIGKDD, KDD '16, ACM, 2016, pp. 1135–1144.

[23] M. T. Ribeiro, S. Singh, C. Guestrin, Anchors: High-precision model-agnostic explanations, in: Proceedings of the AAAI conference on artificial intelligence, volume 32, 2018.

[24] P. Toccaceli, Introduction to conformal predictors, Pattern Recognition 124 (2022).

[25] R. F. Barber, E. J. Candès, A. Ramdas, R. J. Tibshirani, Conformal prediction beyond exchangeability, The Annals of Statistics (2022).

[26] V. Vovk, J. Shen, V. Manokhin, M.-g. Xie, Nonparametric predictive distributions based on conformal prediction, in: Conformal and probabilistic prediction and applications, PMLR, 2017, pp. 82–102.

[27] B. Van Calster, D. J. McLernon, M. Van Smeden, L. Wynants, E. W. Steyerberg, Calibration: the achilles heel of predictive analytics, BMC medicine 17 (2019) 1–7.

[28] V. Vovk, G. Shafer, I. Nouretdinov, Self-calibrating probability forecasting, in: Advances in Neural Information Processing Systems, 2004, pp. 1133–1140.

[29] V. Vovk, I. Petej, Venn-abers predictors, in: Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence, 2014, pp. 829–838.

[30] T. Löfström, H. Löfström, U. Johansson, C. Sönströd, R. Matela, Calibrated explanations for regression, Machine Learning (2024).

[31] V. Vovk, I. Petej, I. Nouretdinov, V. Manokhin, A. Gammerman, Computationally efficient versions of conformal predictive distributions, Neurocomputing 397 (2020) 292–308.

[32] V. Vovk, A. Gammerman, G. Shafer, Algorithmic Learning in a Random World, Springer-Verlag New York, Inc., 2005.

[33] A. N. Angelopoulos, S. Bates, Conformal prediction: A user-friendly introduction, Foundations and Trends® in Machine Learning 16 (2023) 494–591.

[34] G. Shafer, V. Vovk, A tutorial on conformal prediction., Journal of Machine Learning Research 9 (2008).

[35] A. N. Angelopoulos, S. Bates, A gentle introduction to conformal prediction and distribution-free uncertainty quantification, arXiv preprint arXiv:2107.07511 (2021).

[36] T. Löfström, H. Löfström, U. Johansson, Calibrated explanations for multi-class, volume 230 of *Proceedings of Machine Learning Research*, PMLR, 2024, pp. 175–194.