

PREDICTIVE INFERENCE WITH FEATURE CONFORMAL PREDICTION

Jiaye Teng^{1,3,4,*}, Chuan Wen^{1,3,4,*}, Dinghuai Zhang^{2,*},
Yoshua Bengio², Yang Gao^{1,3,4}, Yang Yuan^{1,3,4,†}

¹Institute for Interdisciplinary Information Sciences, Tsinghua University

²Mila - Quebec AI Institute

³Shanghai Artificial Intelligence Laboratory

⁴Shanghai Qi Zhi Institute

{tjy20, cwen20}@mails.tsinghua.edu.cn, dinghuai.zhang@mila.quebec

ABSTRACT

Conformal prediction is a distribution-free technique for establishing valid prediction intervals. Although conventionally people conduct conformal prediction in the output space, this is not the only possibility. In this paper, we propose feature conformal prediction, which extends the scope of conformal prediction to semantic feature spaces by leveraging the inductive bias of deep representation learning. From a theoretical perspective, we demonstrate that feature conformal prediction provably outperforms regular conformal prediction under mild assumptions. Our approach could be combined with not only vanilla conformal prediction, but also other adaptive conformal prediction methods. Apart from experiments on existing predictive inference benchmarks, we also demonstrate the state-of-the-art performance of the proposed methods on *large-scale* tasks such as ImageNet classification and Cityscapes image segmentation. The code is available in <https://github.com/AlvinWen428/FeatureCP>.

1 INTRODUCTION

Although machine learning models work well in numerous fields (Silver et al., 2017; Devlin et al., 2019; Brown et al., 2020), they usually suffer from over-confidence issues, yielding unsatisfactory uncertainty estimates (Guo et al., 2017a; Chen et al., 2021; Gawlikowski et al., 2021). To tackle the uncertainty issues, people have developed a multitude of uncertainty quantification techniques, including calibration (Guo et al., 2017b; Minderer et al., 2021), Bayesian neural networks (Smith, 2014; Blundell et al., 2015), and many others (Sullivan, 2015).

Among different uncertainty quantification techniques, *conformal prediction* (CP) stands out due to its simplicity and low computational cost properties (Vovk et al., 2005; Shafer & Vovk, 2008; Angelopoulos & Bates, 2021). Intuitively, conformal prediction first splits the dataset into a training fold and a calibration fold, then trains a machine learning model on the training fold, and finally constructs the confidence band via a non-conformity score on the calibration fold. Notably, the confidence band obtained by conformal prediction is *guaranteed* due to the exchangeability assumption in the data. With such a guarantee, conformal prediction has been shown to perform promisingly on numerous realistic applications (Lei & Candès, 2021b; Angelopoulos et al., 2022).

Despite its remarkable effectiveness, vanilla conformal prediction (vanilla CP) is only deployed in the output space, which is not the only possibility. As an alternative, feature space in deep learning stands out due to its powerful inductive bias of deep representation. Take the image segmentation problem as an example. In such problems, we anticipate a predictive model to be certain in the informative regions (*e.g.*, have clear objects), while uncertain elsewhere. Since different images would possess different object boundary regions, it is inappropriate to return the same uncertainty for different positions, as standard conformal prediction does. Nonetheless, if we instead employ conformal

*Equal Contribution.

†Correspond to yuanyang@mail.tsinghua.edu.

Figure 1: Illustration of vanilla CP (left) vs Feature CP (right). Feature CP operates in the semantic feature space, as opposed to the commonly adopted output space. These methods are described in further detail in Sections 3 and 4.

prediction on the more meaningful feature space, albeit all images have the same uncertainty on this intermediate space, the pixels would exhibit effectively different uncertainty in the output space after a non-trivial non-linear transformation (see Figure 3).

In this work, we thus propose the Feature Conformal Prediction (Feature CP) framework, which deploys conformal prediction in the feature space rather than the output space (see Figure 1). However, there are still two issues unsolved for performing Feature CP: (a) commonly used non-conformity scores require a ground truth term, but here the ground truth in feature space is not given; and (b) transferring the confidence band in the feature space to the output space is non-trivial. To solve problem (a), we propose a new non-conformity score based on the notation surrogate feature, which replaces the ground truth term in previous non-conformity scores. As for (b), we propose two methods: Band Estimation which calculates the upper bound of the confidence band, together with Band Detection to determine whether a response locates in the confidence band. More interestingly, feature-level techniques are pretty general and can be deployed into other distribution-free inference algorithms, e.g., conformalized quantile regression (CQR). This shows the great potential application impact of the proposed Feature CP methodology (see the discussion in Appendix B.4).

From a theoretical perspective, we demonstrate Feature CP is provably more efficient, in the sense that it yields shorter confidence bands than vanilla CP, given that the feature space meets cubic conditions. Here the cubic conditions sketch the properties of feature space from three perspectives, including length preserving, expansion, and quantile stability (see Theorem 6). At a colloquial level, the cubic conditions assume the feature space has a smaller distance between individual non-conformity scores and their quantiles, which reduces the cost of the quantile operation. We empirically validate that the feature space in deep learning satisfies the cubic conditions, thus resulting in a better confidence band with a shorter length (See Figure 2) according to our theoretical analysis.

Our contributions can be summarized as follows:

- We propose Feature CP, together with a corresponding non-conformity score and an uncertainty band estimation method. The proposed method no longer treats the trained model as a black box but exploits the semantic feature space information. What's more, our approach could be directly deployed with any pretrained model as a plug-in component, without the need of re-training under specially designed learning criteria.
- Theoretical evidence guarantees that Feature CP is both (a) efficient, where it yields shorter confidence bands, and (b) effective, where the empirical coverage provably exceeds the given confidence level, under reasonable assumptions.
- We conduct extensive experiments under both synthetic and realistic settings (e.g., pixel-level image segmentation) to corroborate the effectiveness of the proposed algorithm. Besides, we demonstrate the universal applicability of our method by deploying feature-level operations to improve other adaptive conformal prediction methods such as CQR.

2 RELATED WORK

Conformal prediction is a statistical framework dealing with uncertainty quantification problems (Vovk et al., 2005; Shafer & Vovk, 2008; Nourtdinov et al., 2011; Barber et al., 2020; Angelopoulos & Bates, 2021). The research on conformal prediction can be roughly split into the following branches. The first line of work focuses on relaxing the assumptions about data distribution in conformal prediction, e.g., exchangeability (Tibshirani et al., 2019; Hu & Lei, 2020; Podkopaev & Ramdas, 2021; Barber et al., 2022). The second line aims at improving the efficiency of conformal prediction (Romano et al., 2020b; Sesia & Candès, 2020; Izbicki et al., 2020a; Yang & Kuchibhotla, 2021; Stutz et al., 2021). The third line tries to generalize conformal prediction to different settings, e.g., quantile regression (Romano et al., 2019), Nearest Neighbors (Papadopoulos et al., 2011), density estimator (Izbicki et al., 2020b), survival analysis (Teng et al., 2021), eCand et al., 2021), or conditional histogram regression (Sesia & Romano, 2021). There are also works combining conformal prediction with other machine learning topics, such as functional data (Lei et al., 2013), treatment effects (Lei & Candès, 2021a), time series analysis (Xu & Xie, 2021), online learning (Gibbs & Candès, 2021), adversarial robustness (Gendler et al., 2022), and many others.

Besides conformal prediction, there are many other uncertainty quantification techniques, including calibration (Guo et al., 2017a; Kuleshov et al., 2018; Nixon et al., 2019) and Bayesian-based techniques (Blundell et al., 2015; Hernández-Lobato & Adams, 2015; Li & Gal, 2017). Different from the above techniques, conformal prediction is appealing due to its simplicity, computationally free, and model-free properties.

Image segmentation is a traditional task in computer vision, which focuses on partitioning images into different semantic segments (Haralick & Shapiro, 1985; Senthilkumaran & Rajesh, 2009; Minaee et al., 2020). A line of researches applies conformal prediction with some threshold output for all pixels (Angelopoulos & Bates, 2021; Bates et al., 2021), or focus on the risk control tasks (Angelopoulos et al., 2021a). Different from previous approaches, our method first achieves meaningful pixel-level conformal prediction results to the best of our knowledge.

3 PRELIMINARIES

Predictive inference. Let $(X; Y) \sim P$ denotes a random data pair, an image and its segmentation map. Given a significance level, we aim to construct a confidence band $\hat{C}_1(X)$, such that

$$P(Y \in \hat{C}_1(X)) \geq 1 - \alpha. \quad (1)$$

There is a tradeoff between efficiency and effectiveness, since one can always set $\hat{C}_1(X)$ to be infinitely large to satisfy Equation (1). In practice, we wish the measure of the confidence band, (its length) can be as small as possible, given that the coverage in Equation (1) holds.

Dataset. Let $D = \{(X_i; Y_i)\}_{i \in \mathcal{I}}$ denotes the dataset, where \mathcal{I} denotes the set of data index and $(X_i; Y_i)$ denotes a sample pair following the distribution. Typically, conformal prediction requires that data in D satisfies exchangeability (see below) rather than the stronger i.i.d. (independent and identically distributed) condition. We use n to represent the cardinality of a set. Conformal prediction needs to first randomly split the dataset into a training fold $D_{tr} = \{(X_i; Y_i)\}_{i \in \mathcal{I}_{tr}}$ and a calibration fold $D_{ca} = \{(X_i; Y_i)\}_{i \in \mathcal{I}_{ca}}$, where $|\mathcal{I}_{tr}| = n_{tr}$ and $|\mathcal{I}_{ca}| = n_{ca}$. We denote the test point as $(X^0; Y^0)$, which is also sampled from the distribution.

Training process. During the training process, we train a machine learning model denoted by \hat{g} (e.g., neural network) with the training fold D_{tr} . For the ease of the following discussion, we rewrite the model as $\hat{g} = \hat{f} \circ \hat{h}$, where \hat{f} denotes the feature function (e.g., first several layers in neural networks) and \hat{h} denotes the prediction head (e.g., last several layers in neural networks).

Calibration process. Different from usual machine learning methods, conformal prediction has an additional calibration process. Specifically, we calculate non-conformity scores $v_i = s(X_i; Y_i; \hat{g})$ based on the calibration fold D_{ca} , where $s(\cdot; \cdot; \cdot)$ is a function informally measuring how the model fits the ground truth. The simplest form of non-conformity scores is $v_i = |Y_i - \hat{g}(X_i)|$. One could adjust the form of the non-conformity score according to different contexts (Romano et al. (2019); Teng et al. (2021)). Based on the selected non-conformity score, a matching confidence band could be subsequently created.

Algorithm 1 Conformal Prediction

Require: Desired confidence level α , dataset $\mathcal{D} = \{(X_i; Y_i)\}_{i=1}^n$, test point (X^0, Y^0) , non-conformity score function $s(\cdot)$

- 1: Randomly split the dataset \mathcal{D} into a training fold $\mathcal{D}_{tr} = \{(X_i; Y_i)\}_{i=1}^{n-1}$ and a calibration fold $\mathcal{D}_{ca} = \{(X_i; Y_i)\}_{i=1}^1$;
- 2: Train a base machine learning model \hat{f} with \mathcal{D}_{tr} to estimate the response \hat{Y} ;
- 3: For each $(X_i; Y_i) \in \mathcal{D}_{ca}$, calculate its non-conformity score $s_i = s(X_i; Y_i; \hat{f})$;
- 4: Calculate the $(1 - \alpha)$ -th quantile $Q_{1-\alpha}$ of the distribution $\frac{1}{|\mathcal{D}_{ca}|} \sum_{i=1}^{|\mathcal{D}_{ca}|} s_i$.

Ensure: $C_1(X^0) = \{Y : s(X^0; Y; \hat{f}) \leq Q_{1-\alpha}\}$.

We present vanilla CP in Algorithm 1. Moreover, we demonstrate its theoretical guarantee in Proposition 2, based on the following notation of exchangeability in Assumption 1.

Assumption 1 (exchangeability) Assume that the calibration data $\{(X_i; Y_i)\}_{i=1}^{|\mathcal{D}_{ca}|}$ and the test point (X^0, Y^0) are exchangeable. Formally, denote $Z_i = (X_i; Y_i)$, $i = 1, \dots, |\mathcal{D}_{ca}| + 1$, as the above data pair, then Z_i are exchangeable if arbitrary permutation leads to the same distribution, i.e.,

$$(Z_1; \dots; Z_{|\mathcal{D}_{ca}|+1}) \stackrel{d}{=} (Z_{(1)}; \dots; Z_{(|\mathcal{D}_{ca}|+1)}) \quad (2)$$

with arbitrary permutation σ over $\{1, \dots, |\mathcal{D}_{ca}| + 1\}$.

Note that Assumption 1 is weaker than the i.i.d. assumption. Therefore, it is reasonable to assume the exchangeability condition to hold in practice. Based on the exchangeability assumption, one can show the following theorem, indicating that conformal prediction indeed returns a valid confidence band, which satisfies Equation (1).

Theorem 2 (theoretical guarantee for conformal prediction) (Law, 2006; Lei et al., 2018; Tibshirani et al., 2019) Under Assumption 1, the confidence band $C_1(X^0)$ returned by Algorithm 1 satisfies

$$P(Y^0 \in C_1(X^0)) \geq 1 - \alpha.$$

4 METHODOLOGY

In this section, we broaden the concept of conformal prediction using feature-level operations. This extends the scope of conformal prediction and makes it more flexible. We analyze the algorithm components and details in Section 4.1 and Section 4.2. The algorithm is naturally summarized in Section 4.3. We remark that although in this work we discuss Feature CP under regression regimes for simplicity's sake, one can easily extend the idea to classification problems.

4.1 NON-CONFORMITY SCORE

Conformal prediction necessitates a non-conformity score to measure the conformity between prediction and ground truth. Traditional conformal prediction usually uses norm-based non-conformity score due to its simplicity, e.g., $s(X; Y; \hat{f}) = \|Y - \hat{f}(X)\|$, where Y is the provided ground truth target label. Nonetheless, we have no access to the given target features if we want to conduct conformal prediction at the feature level. To this end, we introduce surrogate features (see Definition 3), which could serve as the role of ground truth in Feature CP.

Definition 3 (Surrogate feature) Consider a trained neural network $\hat{f} = \mathcal{G} \circ \hat{f}$ where \mathcal{G} denotes the composition operator. For a sample $(X; Y)$, we define $\hat{v} = \hat{f}(X)$ to be the trained feature. Besides, we define the surrogate feature to be any feature v such that $\mathcal{G}(v) = Y$.

In contrast to commonly adopted regression or classification scenarios where the label is unidimensional, the dimensionality of features could be much larger. We thus define a corresponding non-conformity score based on the surrogate feature as follows:

$$s(X; Y; \hat{f}) = \inf_{v: \mathcal{G}(v) = Y} \|v - \hat{f}(X)\| \quad (3)$$

¹We use δ_u to represent a Dirac Delta function (distribution) at point

It is usually complicated to calculate the score in Equation 3 due to the in num operator. Therefore, we design Algorithm 2 to calculate an upper bound of the non-conformity score. Although the exact in num is hard to achieve in practice, we can apply gradient descent starting from the trained feature \hat{v} to find a surrogate feature around it. In order to demonstrate the reasonability of this algorithm, we analyze the non-conformity score distribution with realistic data in Appendix B.9.

4.2 BAND ESTIMATION AND BAND DETECTION

Utilizing the non-conformity score derived in Section 4.1, one could derive a confidence band in the feature space. In this section, we mainly focus on how to transfer the confidence band in feature space to the output space, calculating the set

$$f(\hat{g}(v) : kv - \hat{v}k \leq Q_1 - g; \quad (4)$$

where \hat{v} is the trained feature, \hat{g} is the prediction head, and Q_1 is derived based on the calibration set (even though slightly different, we refer to step 4 in Algorithm 1 for the notion of Q_1 ; a formal discussion of it is deferred to Algorithm 3).

Since the prediction head is usually highly non-linear, the exact confidence band is hard to represent explicitly. Consequently, we provide two approaches: Band Estimation which aims at estimating the upper bound of the confidence band, and Band Detection which aims at identifying whether a response falls inside the confidence interval. We next crystallize the two methods.

Band Estimation. We model the Band Estimation problem as a perturbation analysis one, where we regard Equation (4) as a perturbation of the trained feature \hat{v} , and analyze the output bounds of prediction head. In this work, we apply linear relaxation based perturbation analysis (LiPRA) (Xu et al., 2020) to tackle this problem under deep neural network regimes. LiPRA transforms the certification problem as a linear programming problem, and solves it accordingly. The relaxation would result in a relatively looser interval than the actual band, so this method would give an upper bound estimation of the exact band length.

Band Detection. Band Estimation could potentially end up with loose inference results. Typically, we are only interested in determining whether a point $\hat{g}(X^0)$ is in the confidence band $\mathcal{C}(X^0)$ for a test sample X^0 . To achieve this goal, we first apply Algorithm 2 using data point (X^0, Y) , which returns a non-conformity score s . We then test whether the score is smaller than quantile Q_1 on the calibration set (see Equation (4)). If so, we deduce that $\hat{g}(X^0)$ (or vice versa if not).

4.3 FEATURE CONFORMAL PREDICTION

Based on the above discussion, we summarize Feature CP in Algorithm 3. Different from vanilla CP (see Algorithm 1), Feature CP uses a different non-conformity score based on surrogate features, and we need an additional Band Estimation or Band Detection (step 5) to transfer the band from feature space to output space.

We then discuss two intriguing strengths of Feature CP. First, the proposed technique is universal and could improve other advanced adaptive conformal inference techniques utilizing the inductive bias of learned feature space. Specifically, we propose Feature CQR with insights from CQR (Romano et al., 2019), a prominent adaptive conformal prediction method with remarkable performance, to demonstrate the universality of our technique. We relegate related algorithmic details to Section B.4. Second, although methods such as CQR require specialized training criteria (quantile regression) for the predictive models, Feature CP could be directly applied to any given pretrained model and

²We here show the Feature CP algorithm with Band Estimation. We defer the practice details of using Band Detection in step 5 of Algorithm 3 to Appendix.

Algorithm 3 Feature Conformal Prediction

Require: Level α , dataset $\mathcal{D} = \{(X_i; Y_i)\}_{i \in \mathcal{I}}$, test point X^0 ;

- 1: Randomly split the dataset into a training fold \mathcal{D}_{tr} , $\{(X_i; Y_i)\}_{i \in \mathcal{I}_{\text{tr}}}$ together with a calibration fold \mathcal{D}_{ca} , $\{(X_i; Y_i)\}_{i \in \mathcal{I}_{\text{ca}}}$;
- 2: Train a base machine learning model $f^{\wedge}(\cdot)$ using \mathcal{D}_{tr} to estimate the response;
- 3: For each $i \in \mathcal{I}_{\text{ca}}$, calculate the non-conformity score based on Algorithm 2;
- 4: Calculate the $(1 - \alpha)$ -th quantile $Q_{1-\alpha}$ of the distribution $\frac{1}{|\mathcal{I}_{\text{ca}}|} \sum_{i \in \mathcal{I}_{\text{ca}}} V_i + 1$;
- 5: Apply Band Estimation on test data feature (X^0) with perturbation $Q_{1-\alpha}$ and prediction head g , which returns $C_1^{\text{fcp}}(X)$;

Ensure: $C_1^{\text{fcp}}(X)$.

could still give meaningful adaptive interval estimates. This trait facilitates the usage of our method with large pretrained models, which is common in modern language and vision tasks.

4.4 THEORETICAL GUARANTEE

This section presents theoretical guarantees for Feature CP regarding coverage (effectiveness) and band length (efficiency). We provide an informal statement of the theorem below and defer the complete details to Appendix A.1.

Theorem 4 (Informal Theorem on the Efficiency of Feature CP) Under mild assumptions, if the following cubic conditions hold:

1. Length Preservation. Feature CP does not cost much loss in feature space.
2. Expansion. The Band Estimation operator expands the differences between individual length and their quantiles.
3. Quantile Stability. The band length is stable in both the feature space and the output space for a given calibration set.

then Feature CP outperforms vanilla CP in terms of average band length.

The intuition of Theorem 4 is as follows: Firstly, Feature CP and Vanilla CP take quantile operations in different spaces, and the Expansion condition guarantees that the quantile step costs less in Feature CP. However, there may be an efficiency loss when transferring the band from feature space to output space. Fortunately, it is controllable under the Length Preservation condition. The Quantile Stability condition ensures that the band is generalizable from the calibration fold to the test samples. We provide the detailed theorem in Appendix A.1 and empirically validate the cubic conditions in Appendix B.2.

5 EXPERIMENTS

We conduct experiments on synthetic and real-world datasets, mainly to show that Feature CP is (a) effective, i.e., it could return valid confidence bands with empirical coverage larger than α ; (b) efficient, i.e., it could return shorter confidence bands than vanilla CP.

5.1 SETUP

Datasets. We consider both synthetic datasets and real-world datasets, including (a) realistic uni-dimensional target datasets: five datasets from UCI machine learning repository (Asuncion, 2007): physicochemical properties of protein tertiary structure (bio), bike sharing (bike), community and crimes (community) and Facebook comment volume variants one and two (facebook 1/2), five datasets from other sources: blog feedback (blog) (Buza, 2014), Tennessee’s student teacher achievement ratio (star) (Achilles et al., 2008), and medical expenditure panel survey (meps19–21) (Cohen et al., 2009); (b) synthetic multi-dimensional target datasets $Y = WX + b$, where $X \in [0; 1]^{100}$; $Y \in \mathbb{R}^{10}$,

Figure 2: Performance on tasks with one-dimensional target $(0:1)$. Left: Empirical coverage. Right: Confidence interval length, where smaller value is better. The proposed Feature CP and Feature CQR could consistently achieve shorter bands while maintaining a good coverage performance.

Table 1: Performance of both methods on multi-dimensional regression benchmarks (1), where “w-length” denotes the weighted confidence band length.

DATASET	SYNTHETIC				CITYSCAPES					
	COVERAGE		LENGTH		COVERAGE		LENGTH		W-LENGTH	
BASELINE	89:91	1:03	0:401	0:01	91:41	0:51	40:15	0:02	40:15	0:02
FEATURE CP	90:13	0:59	0.373	0:05	90:77	0:91	1.032	0:01	0.906	0:01

follows the standard Gaussian distribution, adds a fixed randomly generated matrix; and (c) real-world semantic segmentation dataset: Cityscapes (Cordts et al., 2016), where we transform the original pixel-wise classification problem into a high-dimensional pixel-wise regression problem. We also extend Feature CP to classification problems and test on the ImageNet (Deng et al., 2009) dataset. We defer more related details to Appendix B.1.

Algorithms. We compare the proposed Feature CP against the vanilla conformal baselines without further specified, which directly deploy conformal inference on the output space. For both methods, we use $\hat{\epsilon}_1$ -type non-conformity score, namely $\hat{\epsilon}_1(X; Y) = kY - (X)k_1$.

Evaluation. We adopt the following metrics to evaluate algorithmic empirical performance.

Empirical coverage (effectiveness) the empirical probability that a test point falls into the predicted confidence band. A good predictive inference method should achieve empirical coverage slightly larger than $1 - \alpha$ for a given significance level α . To calculate the coverage for Feature CP, we first apply Band Detection on the test point (X^0, Y^0) to detect whether Y^0 is in $C_1^{cp}(X^0)$, and then calculate its average value to obtain the empirical coverage.

Band length (efficiency) Given the empirical coverage being larger than $1 - \alpha$, we hope the confidence band to be as short as possible. The band length should be compared under the regime of empirical coverage being larger than $1 - \alpha$, otherwise one can always set the confidence band to empty to get a zero band length. Since the explicit expression for confidence bands is intractable for the proposed algorithm, we could only derive an estimated band length via Band Estimation. Concretely, we

Figure 3: Visualization of Feature CP in image segmentation. The brightness of the pixels in the third column measures the uncertainty of Feature CP, namely the length of confidence bands. The algorithm is more uncertain in the brighter regions. For better visualization, we rescale the interval length to the range $[0, 1]$. Feature CP is more uncertain in non-informative regions, which are object boundaries.

(a) Facebook1 (b) Synthetic (c) Cityscapes

Figure 4: Empirical coverage under different confidence levels. For a good conformal prediction method, the y -axis (i.e., empirical coverage minus α) should keep being above zero for different α . These three figures above show that Feature CP generally performs better than the baseline, in the sense that this difference is above zero most of the time.

First use Band Estimation to estimate the confidence interval, which returns a band with explicit formulation, and then calculate the average length across each dimension.

We formulate the metrics as follows. Let $(Y^{(1)}; \dots; Y^{(d)}) \in \mathbb{R}^d$ denotes the high dimensional response and $\mathcal{C}(X) \in \mathbb{R}^d$ denotes the obtained confidence interval, with length in each dimension forming a vector $|\mathcal{C}(X)| \in \mathbb{R}^d$. With the test set index being $i \in [n]$ and $[d] = \{1, \dots, d\}$, we calculate the empirical coverage and band length respectively as

$$\frac{1}{n} \sum_{i \in [n]} \mathbb{1}(Y_i \in \mathcal{C}(X_i)); \quad \frac{1}{n} \sum_{i \in [n]} \frac{1}{d} \sum_{j \in [d]} |\mathcal{C}(X_i)|^{(j)}.$$

5.2 RESULTS AND DISCUSSION

Effectiveness. We summarize the empirical coverage in Figure 2 (one-dimension response) and Table 1 (multi-dimension response). As Theorem 5 illustrates, the empirical coverage of Feature CP all exceeds the confidence level α , indicating that Feature CP is effective. Besides, Figure 4 demonstrates that that the effectiveness holds with different significance levels. For simple benchmarks such as facebook1 and synthetic data, both methods achieve similar coverage due to the simplicity; while for the harder Cityscapes segmentation task, the proposed method outperforms the baseline under many confidence levels.

Efficiency. We summarize the confidence band in Figure 2 (one-dimension response) and Table 1 (multi-dimension response). The band length presented is an estimated version via Band Estimation.

Table 2: Feature CP outperform previous methods on large-scale classification tasks (ImageNet, $\alpha = 0.1$). We compare our results with APS (Romano et al., 2020a) and RAPS (Angelopoulos et al., 2021b). The results of baselines are taken from Angelopoulos et al. (2021b).

METHOD MODEL	ACCURACY		COVERAGE				LENGTH				
	TOP-1	TOP-5	APS	RAPS	FEATURE CP	APS	RAPS	FEATURE CP	APS	RAPS	FEATURE CP
RESNET 18	0.698	0.891	0.900	0.900	0.902	0.0023	16:2	4:43	3.82	0.18	
RESNET 50	0.761	0.929	0.900	0.900	0.900	0.0047	12:3	2:57	2.14	0.03	
RESNET 101	0.774	0.936	0.900	0.900	0.900	0.0025	10:7	2:25	2.17	0.01	
RESNEXT 101	0.783	0.945	0.900	0.900	0.900	0.0030	19:7	2:00	1.80	0.08	
SHUFFLENET	0.694	0.883	0.900	0.900	0.901	0.0019	31:9	5:05	5.06	0.04	
VGG 16	0.716	0.904	0.901	0.900	0.901	0.0047	14:1	3:54	3.26	0.03	

Note that Feature CP outperforms the baseline in the sense that it achieves a shorter band length and thus a more efficient algorithm.

Comparison to CQR. The techniques proposed in this paper can be generalized to other conformal prediction techniques. As an example, we propose Feature CQR which is a feature-level generalized version of CQR, whose details are deferred to Appendix B.4. We display the comparison in Figure 2, where our method consistently outperforms CQR baseline by leveraging good representation. Besides, we evaluate the group coverage performance of CQR and Feature CQR in Appendix B.5, demonstrating that Feature CQR generally outperforms CQR in terms of condition coverage. One can also generalize the other existing techniques to feature versions, localized Conformal Prediction (Guan, 2019; Han et al., 2022).

Extension to classification tasks. The techniques in Feature CP can be generalized to classification tasks. We take the ImageNet classification dataset as an example. Table 2 shows the effectiveness of Feature CP in classification tasks. We experiment on various different architectures and use the same pretrained weights as the baselines. We provide additional details in Appendix B.7.

Ablation on the splitting point. We demonstrate that the coverage is robust to the splitting point in both small neural networks (Table 9 in Appendix) and large neural networks (Table 10 in Appendix). One can also use standard cross-validation to choose the splitting point. Since the output space is one of the special feature layers, such techniques always generalize the scope of vanilla CP.

Truthfulness. We visualize the segmentation results in Figure 3, which illustrates that Feature CP returns large bands (light region) on the non-informative regions (object boundaries) and small bands (dark region) on the informative regions. We do not show baseline visualization results since they return the same band in each dimension for each sample, and therefore does not contain much information. We also evaluate the performance with weighted band length, defined in Appendix B.1.

How does Feature CP benefit from a good deep representation? Here we provide some intuition on the success of Feature CP algorithm: we claim it is the usage of good (deep) representation that enables Feature CP to achieve better predictive inference. To validate this hypothesis, we contrast Feature CP against the baseline with an unlearned neural network (whose feature is not semantic as desired). This “random” variant of Feature CP does not outperform its vanilla counterpart with the same neural network, which confirms our hypothesis. We defer the results to Table 6 and related discussion to Appendix B.3.

More discussion. We analyze the failure (i.e., inefficient) reasons of vanilla CP in image segmentation task from the following two perspectives. Firstly, this paper aims to provide a provable coverage, namely, the confidence band should cover the ground truth for each pixel. Since vanilla CP returns the same band for different samples, the loss is pretty large such that the returned interval is large enough to cover the ground truth. Secondly, an intuitive explanation relates to our usage of conformal prediction. We choose the infinity norm because reporting the total band length requires the band length in each dimension. As a result, the non-conformity score is large as long as there exists one pixel that does not fit well, contributing to an unsatisfying band for vanilla CP.

ACKNOWLEDGMENTS

This work has been partly supported by the Ministry of Science and Technology of the People's Republic of China, the 2030 Innovation Megaprojects "Program on New Generation Artificial Intelligence" (Grant No. 2021AAA0150000). This work is also supported by a grant from the Guoqiang Institute, Tsinghua University.

REFERENCES

- CM Achilles, Helen Pate Bain, Fred Bellott, Jayne Boyd-Zaharias, Jeremy Finn, John Folger, John Johnston, and Elizabeth Word. Tennessee's student teacher achievement ratio (star) project. Harvard Dataverse, 2008.
- Anastasios N. Angelopoulos and Stephen Bates. A gentle introduction to conformal prediction and distribution-free uncertainty quantification. *CoRR*, abs/2107.07511, 2021. URL <https://arxiv.org/abs/2107.07511>
- Anastasios N Angelopoulos, Stephen Bates, Emmanuel J. Candès, Michael I Jordan, and Lihua Lei. Learn then test: Calibrating predictive algorithms to achieve risk control. *arXiv preprint arXiv:2110.01052*, 2021a.
- Anastasios N. Angelopoulos, Amit P. S. Kohli, Stephen Bates, Michael I. Jordan, Jitendra Malik, Thayer Alshaabi, Srigokul Upadhyayula, and Yaniv Romano. Image-to-image regression with distribution-free uncertainty quantification and applications in image-to-image regression. *CoRR*, abs/2202.05265, 2022. URL <https://arxiv.org/abs/2202.05265>
- Anastasios Nikolas Angelopoulos, Stephen Bates, Michael I. Jordan, and Jitendra Malik. Uncertainty sets for image classifiers using conformal prediction. *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021b. URL https://openreview.net/forum?id=eNdiU_DbM9
- Arthur U. Asuncion. Uci machine learning repository, university of california, irvine, school of information and computer sciences. 2007.
- Rina Foygel Barber, Emmanuel J. Candès, Aaditya Ramdas, and Ryan J. Tibshirani. The limits of distribution-free conditional predictive inference. *Information and Inference: A Journal of the IMA*, 2020.
- Rina Foygel Barber, Emmanuel J Candès, Aaditya Ramdas, and Ryan J Tibshirani. Conformal prediction beyond exchangeability. *arXiv preprint arXiv:2202.13415*, 2022.
- Stephen Bates, Anastasios Nikolas Angelopoulos, Lihua Lei, Jitendra Malik, and Michael I. Jordan. Distribution free, risk controlling prediction sets, 2021. URL <https://arxiv.org/abs/2101.02703>
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In Francis R. Bach and David M. Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 1613–1622. JMLR.org, 2015. URL <http://proceedings.mlr.press/v37/blundell15.html>
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>

Krisztian Buza. Feedback prediction for blogs. *Data analysis, machine learning and knowledge discovery* pp. 145–152. Springer, 2014.

Emmanuel J. Carron, Lihua Lei, and Zhimei Ren. Conformalized survival analysis. *arXiv preprint arXiv:2103.09763*, 2021.

Yanzhi Chen, Dinghuai Zhang, Michael U. Gutmann, Aaron Courville, and Zhanxing Zhu. Neural approximate sufficient statistics for implicit models. *International Conference on Learning Representation*, 2021. URL <https://openreview.net/forum?id=SRDuJssQud>.

Joel W. Cohen, Steven B. Cohen, and Jessica S. Banthin. The medical expenditure panel survey: A national information resource to support healthcare cost research and inform policy and practice. *Medical Care* 47:S44–S50, 2009.

Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, K. Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. *2009 IEEE Conference on Computer Vision and Pattern Recognition* pp. 248–255, 2009.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, 171–4186. Association for Computational Linguistics, 2019. doi: 10.18653/v1/n19-1423. URL <https://doi.org/10.18653/v1/n19-1423>.

Jakob Gawlikowski, Cedrique Rovile Njiteucheu Tassi, Mohsin Ali, Jongseok Lee, Matthias Humt, Jianxiang Feng, Anna M. Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, Muhammad Shahzad, Wen Yang, Richard Bamler, and Xiao Xiang Zhu. A survey of uncertainty in deep neural networks. *CoRR* abs/2107.03342, 2021. URL <https://arxiv.org/abs/2107.03342>.

Asaf Gendler, Tsui-Wei Weng, Luca Daniel, and Yaniv Romano. Adversarially robust conformal prediction. In *International Conference on Learning Representation*, 2022. URL <https://openreview.net/forum?id=9L1Bsl4wP1H>.

Isaac Gibbs and Emmanuel J. Carron. Adaptive conformal inference under distribution shift. In *NeurIPS*, 2021.

Leying Guan. Conformal prediction with localization. *arXiv preprint arXiv:1908.08558*, 2019.

Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1321–1330. PMLR, 2017a. URL <http://proceedings.mlr.press/v70/guo17a.html>.

Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. *ArXiv*, abs/1706.04599, 2017b.

Xing Han, Ziyang Tang, Joydeep Ghosh, and Qiang Liu. Split localized conformal prediction. *CoRR* abs/2206.13092, 2022. doi: 10.48550/arXiv.2206.13092. URL <https://doi.org/10.48550/arXiv.2206.13092>.

Robert M. Haralick and Linda G. Shapiro. Image segmentation techniques. *Comput. Vis. Graph. Image Process* 29(1):100–132, 1985. doi: 10.1016/S0734-189X(85)90153-7. URL [https://doi.org/10.1016/S0734-189X\(85\)90153-7](https://doi.org/10.1016/S0734-189X(85)90153-7).

- José Miguel Hernández-Lobato and Ryan P. Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. In Francis R. Bach and David M. Blei (eds), *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015 volume 37 of JMLR Workshop and Conference Proceedings*, pp. 1861–1869. JMLR.org, 2015. URL <http://proceedings.mlr.press/v37/hernandez-lobatoc15.html>
- Xiaoyu Hu and Jing Lei. A distribution-free test of covariate shift using conformal prediction. *Methodology* 2020.
- Rafael Izbicki, Gilson Shimizu, and Rafael B Stern. Cd-split and hpd-split: efficient conformal regions in high dimensions. *arXiv preprint arXiv:2007.12778*, 2020a.
- Rafael Izbicki, Gilson T. Shimizu, and Rafael Bassi Stern. Distribution-free conditional predictive bands using density estimation. *ArXiv*, abs/1910.05575, 2020b.
- Volodymyr Kuleshov, Nathan Fenner, and Stefano Ermon. Accurate uncertainties for deep learning using calibrated regression. In Jennifer G. Dy and Andreas Krause (eds), *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2801–2809. PMLR, 2018. URL <http://proceedings.mlr.press/v80/kuleshov18a.html>
- James Law. Review of "algorithmic learning in a random world by vovk, gammerman and shafer", springer, 2005, ISBN: 0-387-00152-2. *SIGACT News* 37(4):38–40, 2006. doi: 10.1145/1189056.1189065. URL <https://doi.org/10.1145/1189056.1189065>
- Jing Lei, Alessandro Rinaldo, and Larry A. Wasserman. A conformal prediction approach to explore functional data. *Annals of Mathematics and Artificial Intelligence* 74:29–43, 2013.
- Jing Lei, Max G'Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association* 113(523):1094–1111, 2018.
- Lihua Lei and Emmanuel J. Candès. Conformal inference of counterfactuals and individual treatment effects. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 2021a.
- Lihua Lei and Emmanuel J Candès. Conformal inference of counterfactuals and individual treatment effects. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 2021b.
- Yingzhen Li and Yarin Gal. Dropout inference in bayesian neural networks with alpha-divergences. In Doina Precup and Yee Whye Teh (eds), *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2052–2061. PMLR, 2017. URL <http://proceedings.mlr.press/v70/li17a.html>
- Shervin Minaee, Yuri Boykov, Fatih Porikli, Antonio Plaza, Nasser Kehtarnavaz, and Demetri Terzopoulos. Image segmentation using deep learning: A survey. *CoRR*, abs/2001.05566, 2020. URL <https://arxiv.org/abs/2001.05566>
- Matthias Minderer, Josip Djolonga, Rob Romijnders, Frances Ann Hubis, Xiaohua Zhai, Neil Houlsby, Dustin Tran, and Mario Lucic. Revisiting the calibration of modern neural networks. In *NeurIPS* 2021.
- Jeremy Nixon, Michael W. Dusenberry, Linchuan Zhang, Ghassen Jerfel, and Dustin Tran. Measuring calibration in deep learning. *IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 38–41. Computer Vision Foundation / IEEE, 2019. URL http://openaccess.thecvf.com/content_CVPRW_2019/html/Uncertainty_and_Robustness_in_Deep_Visual_Learning/Nixon_Measuring_Calibration_in_Deep_Learning_CVPRW_2019_paper.html
- Ilija Noutredinov, Sergi G. Costafreda, Alexander Gamberman, Alexey Ya. Chervonenkis, Vladimir Vovk, Vladimir Vapnik, and Cynthia H. Y. Fu. Machine learning classification with confidence: Application of transductive conformal predictors to mri-based diagnostic and prognostic markers in depression. *NeuroImage* 56(2):809–813, 2011. doi: 10.1016/j.neuroimage.2010.05.023. URL <https://doi.org/10.1016/j.neuroimage.2010.05.023>

- Harris Papadopoulos, Vladimir Vovk, and Alexander Gammerman. Regression conformal prediction with nearest neighbours. *Artif. Intell. Res.*, 40:815–840, 2011. URL <http://jair.org/papers/paper3198.html>
- Adam Paszke, Abhishek Chaurasia, Sangpil Kim, and Eugenio Culurciello. Enet: A deep neural network architecture for real-time semantic segmentation. ArXiv preprint arXiv:1606.02147, 2016.
- Aleksandr Podkopaev and Aaditya Ramdas. Distribution-free uncertainty quantification for classification under label shift. *ICAI*, 2021.
- Yaniv Romano, Evan Patterson, and Emmanuel J. Candès. Conformalized quantile regression. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Aubert, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 3538–3548, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/5103c3584b063c431bd1268e9b5e76fb-Abstract.html>
- Yaniv Romano, Matteo Sesia, and Emmanuel J. Candès. Classification with valid and adaptive coverage. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, pp. 244–255, 2020a. URL <https://proceedings.neurips.cc/paper/2020/hash/244edd7e85dc81602b7615cd705545f5-Abstract.html>
- Yaniv Romano, Matteo Sesia, and Emmanuel J. Candès. Classification with valid and adaptive coverage. *Methodology*, 2020b.
- N. Senthilkumaran and Reghunadhan Rajesh. Image segmentation - A survey of soft computing approaches. *ARTCom 2009, International Conference on Advances in Recent Technologies in Communication and Computing*, Kottayam, Kerala, India, 27-28 October, 2009, pp. 844–846. IEEE Computer Society, 2009. doi: 10.1109/ARTCom.2009.219. URL <https://doi.org/10.1109/ARTCom.2009.219>
- Matteo Sesia and Emmanuel J Candès. A comparison of some conformal quantile regression methods. *Stat*, 9(1):e261, 2020.
- Matteo Sesia and Yaniv Romano. Conformal prediction using conditional histograms. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 6304–6315, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/31b3b31a1c2f8a370206f11127c0dbd-Abstract.html>
- Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Mach. Learn. Res*, 9:371–421, 2008. URL <https://dl.acm.org/citation.cfm?id=1390693>
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy P. Lillicrap, Fan Hui, Laurent Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis. Mastering the game of go without human knowledge. *Nature*, 550(7676):354–359, 2017. doi: 10.1038/nature24270. URL <https://doi.org/10.1038/nature24270>
- Ralph C. Smith. *Uncertainty Quantification - Theory, Implementation, and Applications in Computational Science and Engineering*. SIAM, 2014. ISBN 978-1-611973-21-1. URL <http://bookstore.siam.org/cs12/>
- David Stutz, Krishnamurthy Dvijotham, Ali Taylan Cemgil, and A. Doucet. Learning optimal conformal classifiers. ArXiv, abs/2110.09192, 2021.
- Timothy John Sullivan. *Introduction to uncertainty quantification*, volume 63. Springer, 2015.

- Jiaye Teng, Zeren Tan, and Yang Yuan. T-SCI: A two-stage conformal inference algorithm with guaranteed coverage for cox-mlp. In Marina Meila and Tong Zhang (eds), Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event volume 139 of Proceedings of Machine Learning Research, pp. 10203–10213. PMLR, 2021. URL <http://proceedings.mlr.press/v139/teng21a.html>
- Ryan J. Tibshirani, Rina Foygel Barber, Emmanuel J. Candès, and Aaditya Ramdas. Conformal prediction under covariate shift. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Aubert, Emily B. Fox, and Roman Garnett (eds), Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp. 2526–2536, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/8fb21ee7a2207526da55a679f0332de2-Abstract.html>
- Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. Algorithmic learning in a random world. Springer Science & Business Media, 2005.
- Chen Xu and Yao Xie. Conformal prediction interval for dynamic time-series. ICML, 2021.
- Kaidi Xu, Zhouxing Shi, Huan Zhang, Yihan Wang, Kai-Wei Chang, Minlie Huang, Bhavya Kailkhura, Xue Lin, and Cho-Jui Hsieh. Automatic perturbation analysis for scalable certified robustness and beyond. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/0cbc5671ae26f67871cb914d81ef8fc1-Abstract.html>
- Yachong Yang and Arun Kumar Kuchibhotla. Finite-sample efficient conformal prediction. preprint arXiv:2104.13871, 2021.

Appendix

Section A provides the complete proofs, and Section B.1 provides experiment details.

A THEORETICAL PROOFS

We first show the formal version of Theorem 4 in Section A.1, and show its proof in Section A.2. Theorem 5 and Theorem 6 shows the effectiveness (empirical coverage) and the efficiency (band length) in Feature CP.

We additionally provide Theorem 9 (see Section A.3) and Theorem 15 (see Section A.4) to better validate our theorem, in terms of the length variance and convergence rate.

A.1 THEORETICAL GUARANTEE

This section provides theoretical guarantees for Feature CP regarding coverage (effectiveness) and band length (efficiency), starting from additional notations.

Notations. Let P denote the population distribution. Let $D_{ca} \subset P^n$ denote the calibration set with sample size n and sample index i_{ca} , where we overload the notation P^n to denote the distribution of a set with samples drawn from distribution P . Given the model f^* with feature extractor ϕ and prediction head g , we assume g is continuous. We also overload the notation $Q_1(V)$ to denote the $(1 - \alpha)$ -quantile of the set $V = \{f(\phi(x))\}$. Besides, $\mathbb{E}[\cdot]$ denote the mean of a set, and a set minus a real number denote the broadcast operation.

Vanilla CP. Let $V_{D_{ca}}^o = \{v_i^o\}_{i \in I_{ca}}$ denote the individual length in the output space for vanilla CP, given the calibration set D_{ca} . Concretely, $v_i^o = 2|y_i - \hat{y}_i|$ where y_i denotes the true response of sample i and \hat{y}_i denotes the corresponding prediction. Since vanilla CP returns band length with $(1 - \alpha)$ quantile of non-conformity score, the resulting average band length is derived by $\mathbb{E}[V_{D_{ca}}^o]$.

Feature CP. Let $V_{D_{ca}}^f = \{v_i^f\}_{i \in I_{ca}}$ be the individual length (or diameter in high dimensional cases) in the feature space for Feature CP, given the calibration set D_{ca} . To characterize the band length in the output space, we define $h(v; X)$ as the individual length on sample x in the output space, given the length v in the feature space, i.e., $h(v; X)$ represents the length of the set $\{u \in \mathbb{R}^k : \|u - f^*(X)\|_2 = v\}$. Due to the continuity assumption on function g on the above set is always simply-connected. We here omit the dependency of prediction head g for simplicity. The resulting band length in Feature CP is denoted by $\mathbb{E}_{(X^o, Y^o) \sim P} [H(Q_1(V_{D_{ca}}^f); X^o)]$. Without abuse of notations, operating on a dataset $(g, H(V_{D_{ca}}^f; D_{ca}))$ means operating on each data point $(v_i^f; X_i)$ in the set.

Coverage guarantee. We next provide theoretical guarantees for Feature CP in Theorem 5, which informally shows that under Assumption 1, the confidence band returned by Algorithm 3 is valid, meaning that the coverage is provably larger than $1 - \alpha$. We defer the whole proof to Appendix A.2.

Theorem 5 (theoretical guarantee for Feature CP). Under Assumption 1, for any $\alpha > 0$, the confidence band returned by Algorithm 3 satisfies:

$$P(Y^o \in C_1^{fcp}(X^o)) \geq 1 - \alpha;$$

where the probability is taken over the calibration fold and the testing point (X^o, Y^o) .

Length (efficiency) guarantee. We next show in Theorem 6 that Feature CP is provably more efficient than the vanilla CP.

Theorem 6 (Feature CP is provably more efficient). Assume that the non-conformity score is in norm-type. For the operator h , we assume a Holder assumption that there exist $0 < L > 0$ such that $|h(v; X) - h(u; X)| \leq L|v - u|$ for all X . Besides, we assume that there exist $0 < c > 0$, such that the feature space satisfies the following cubic conditions:

1. Length Preservation. Feature CP does not cost much loss in feature space in a quantile manner, namely $\mathbb{E}_{D_{ca}} [Q_1(H(V_{D_{ca}}^f; D_{ca}))] < \mathbb{E}_{D_{ca}} [Q_1(V_{D_{ca}}^o)] + \epsilon$.

2. Expansion. The operator $H(v; X)$ expands the differences between individual length and their quantiles, namely, $E_{D \sim P} M_j(Q_1(V_D^f); V_D^f) < E_{D \sim P} M_j(Q_1(H(V_D^f; D)); H(V_D^f; D))$ and $E_{D \sim P} M_j(Q_1(V_D^f); V_D^f) < 2 \max\{L; 1\} g(c = \frac{p}{n})^{\min\{f; 1\}}; 1g$.
3. Quantile Stability. Given a calibration set D_{ca} , the quantile of the band length is stable in both feature space and output space, namely, $E_{D \sim P} M_j(Q_1(V_D^f); Q_1(V_{D_{ca}}^f)) \leq \frac{c}{n}$ and $E_{D \sim P} M_j(Q_1(V_D^o); Q_1(V_{D_{ca}}^o)) \leq \frac{c}{n}$.

Then Feature CP provably outperforms vanilla CP in terms of average band length, namely,

$$E H(Q_1(V_{D_{ca}}^f); X^o) < Q_1(V_{D_{ca}}^o);$$

where the expectation is taken over the calibration fold and the testing point (X^o, Y^o) .

The cubic conditions used in Theorem 6 sketch the properties of feature space from different aspects. The first condition implies that the feature space is efficient for each individual, which holds when the band is generally not too large. The second condition is the core of the proof, which informally assumes that the difference between quantile and each individual is smaller in feature space. Therefore, conducting quantile operation would not harm the effectiveness (namely, step 4 in Algorithm 1 and step 4 in Algorithm 3), leading to the efficiency of Feature CP. The last condition helps generalize the results from the calibration set to the test set.

Proof of Theorem 6. We start the proof with Assumption 2, which claims that

$$E_{D \sim P} M_j(Q_1(V_D^f); V_D^f) < E_{D \sim P} M_j(Q_1(H(V_D^f; D)); H(V_D^f; D))$$

$$2 \max\{L; 1\} g(c = \frac{p}{n})^{\min\{f; 1\}}; 1g.$$

We rewrite it as

$$E_{D \sim P} M_j(H(V_D^f; D); H(V_D^f; D)) < E_{D \sim P} M_j(Q_1(H(V_D^f; D)); H(V_D^f; D))$$

$$2 \max\{L; 1\} g(c = \frac{p}{n})^{\min\{f; 1\}}; 1g - E_{D \sim P} M_j(Q_1(V_D^f); V_D^f);$$

Due to Holder condition, we have that $E_{D \sim P} M_j(Q_1(V_D^f); D) < M(H(V_D^f; D)) + L M_j(Q_1(V_D^f); V_D^f)$, therefore

$$E_{D \sim P} M_j(H(Q_1(V_D^f); D); H(Q_1(V_D^f); D)) < E_{D \sim P} M_j(Q_1(H(V_D^f; D)); H(V_D^f; D))$$

$$2 \max\{1; L\} g[c = \frac{p}{n}]^{\min\{f; 1\}}; 1g;$$

Therefore, due to assumption 1, we have that

$$E_{D \sim P} M_j(H(Q_1(V_D^f); D); H(Q_1(V_D^f); D)) < E_{D \sim P} M_j(Q_1(V_D^o); V_D^o)$$

$$2 \max\{1; L\} g[c = \frac{p}{n}]^{\min\{1; 1\}}; 1g;$$

Besides, according to the quantile stability assumption, we have $E_{D \sim P} M_j(H(Q_1(V_D^f); D); H(Q_1(V_D^f); D)) < L[c = \frac{p}{n}]$, and $E_{D \sim P} M_j(Q_1(V_D^o); Q_1(V_D^o)) \leq \frac{c}{n}$. Therefore,

$$E H(Q_1(V_{D_{ca}}^f); X^o)$$

$$= E_{D \sim P} M_j(H(Q_1(V_{D_{ca}}^f); D); H(Q_1(V_{D_{ca}}^f); D))$$

$$< Q_1(V_{D_{ca}}^o) + 2 \max\{1; L\} g[c = \frac{p}{n}]^{\min\{1; 1\}}; 1g + L[c = \frac{p}{n}] + \frac{c}{n}$$

$$< Q_1(V_{D_{ca}}^o);$$

□

A.1.1 EXAMPLE FOR THEOREM 6

This section provides an example for Theorem 6. The key information is that Feature CP loses less efficiency when conducting the quantile step.

Table 3: A concrete example for the comparison between Feature CP and Vanilla CP. IL_o denote the individual length in the feature and output space. $Q(\cdot)$ denote the quantile operator, and $H(\cdot)$ denote the operator that calculates output space length given the feature space length. We remark that the average band length returned by Feature CP (3.1) outperforms that of vanilla CP (4.0).

METHOD	VANILLA CP		FEATURE CP			
	IL_o	$Q(IL_o)$	IL_f	$H(IL_f)$	$Q(IL_f)$	$H(Q(IL_f))$
A	1.0	4.0	1.1	1.2	1.3	1.4
B	2.0	4.0	1.2	2.1	1.3	2.3
C	3.0	4.0	1.1	2.8	1.3	3.1
D	4.0	4.0	1.3	3.8	1.3	3.8
E	5.0	4.0	1.6	5.2	1.3	4.9
QUANTILE	4.0	/	1.3	/	/	/
AVERAGE	/	4.0	/	/	/	3.1

Assume the dataset has five samples labeled A, B, C, D, and E. When directly applying vanilla CP leads to individual length in the output space as 1; 2; 3; 4; 5, respectively. By taking 80% quantile (namely, $\alpha = 0.2$), the final confidence band returned by vanilla CP ($Q(\cdot)$) would be $Q_{0.8}(1; 2; 3; 4; 5) = 4$. Note that for any sample, the returned band length would be 4, and the final average band length is 4.

We next consider Feature CP. We assume that the individual length in the feature space is (1; 1; 1; 2; 1; 1; 1; 3; 1; 6, respectively). Due to the expansion condition (cubic condition #2), the difference between IL_f and $Q(IL_f)$ is smaller than that between IL_o and $Q(IL_o)$. Therefore, the quantile step costs less in Feature CP. Since IL_f is close to $Q(IL_f)$, their corresponding output length $H(IL_f)$, $H(Q(IL_f))$ are also close. Besides, to link conformal prediction and vanilla CP, the length preservation condition (cubic condition #1) ensures that $H(Q(IL_f))$ is close to $H(IL_f)$. Therefore, the final average length $H(Q(L_f))$ is close to the average length $H(IL_o)$, which is better than $Q(IL_o)$. Finally, the quantile stability condition (cubic condition #3) generalizes the results from the calibration set to the test set.

A.2 PROOF OF THEOREM 5

Theorem 5 (theoretical guarantee for Feature CP) Under Assumption 1, for any $\alpha > 0$, the confidence band returned by Algorithm 3 satisfies:

$$P(Y^0 \in C_1^{fcp}(X^0)) \geq 1 - \alpha;$$

where the probability is taken over the calibration fold and the testing point (X^0, Y^0) .

Proof of Theorem 5. The key to the proof is to derive the exchangeability of the non-conformity score, given that the data in the calibration fold and test fold are exchangeable (see Assumption 1).

For ease of notations, we denote the data points in the calibration fold and the test fold as $\mathcal{D} = \{(X_i; Y_i)\}_{i \in [m]}$, where m denotes the number of data points in both calibration fold and test fold. By Assumption 1, the data points \mathcal{D} are exchangeable.

The proof can be split into three parts. The first step is to show that for any function independent of \mathcal{D}^0 , $h(X_i; Y_i)$ are exchangeable. The second step is to show that the proposed score function satisfies the above requirements. And the third step is to show the theoretical guarantee based on the exchangeability of the non-conformity score.

We next prove the first step: for any given function $h: X \times Y \rightarrow \mathbb{R}$ that is independent of data points in \mathcal{D}^0 , we have that $h(X_i; Y_i)$ are exchangeable. Specifically, its CDF and its perturbation CDF

F_v is the same, given the training fold D_{tr} .

$$\begin{aligned} & F_v(u_1; \dots; u_n | D_{tr}) \\ &= P(h(X_1; Y_1) \approx u_1; \dots; h(X_n; Y_n) \approx u_n | D_{tr}) \\ &= P((X_1; Y_1) \in C_{h^{-1}}(u_1); \dots; (X_n; Y_n) \in C_{h^{-1}}(u_n) | D_{tr}) \\ &= P((X_{(1)}; Y_{(1)}) \in C_{h^{-1}}(u_1); \dots; (X_{(n)}; Y_{(n)}) \in C_{h^{-1}}(u_n) | D_{tr}) \\ &= P(h(X_{(1)}; Y_{(1)}) \approx u_1; \dots; h(X_{(n)}; Y_{(n)}) \approx u_n | D_{tr}) \\ &= F_v(u_1; \dots; u_n | D_{tr}); \end{aligned}$$

where \approx denotes a random perturbation, and $C_{h^{-1}}(u) = \{ (X; Y) : h(X; Y) \approx u \}$.

The second step is to show that the proposed non-conformity score function (See Equation 3) is independent of the data set. To show that, we note that the proposed score function in Equation (3) (we rewrite it in Equation (5)) is totally independent of data set. In that we only use the information of \hat{f} and \hat{g} which is dependent on the training fold D_{tr} , instead of D^0 .

$$s(X; Y; \hat{g}, \hat{f}) = \inf_{v: \hat{g}(v) = Y} \text{kv} - \hat{f}(X)k \quad (5)$$

Besides, note that when calculating the non-conformity score in Algorithm 3 for each testing data/calibration data, we do not access any information on the calibration folds for any other points. Therefore, the score function does not depend on the calibration fold or test fold. We finally remark that here we always state that the score function does not depend on the calibration fold or test fold, but its realizations $s(X; Y; \hat{g}, \hat{f})$ can depend on the two folds, $(X; Y) \in D^0$. This does not contrast with the requirement in the first step.

Therefore, combining the two steps leads to a conclusion that the non-conformity score is exchangeable. Finally, following Lemma 1 in Tibshirani et al. (2019), the theoretical guarantee holds under the exchangeability of non-conformity scores. \square

A.3 LENGTH VARIANCE GUARANTEE

The next Theorem 9 demonstrates that the length returned by Feature CP would be individually different. Specifically, the variance for the length is lower bounded by a constant. The essential intuition is that, for a non-linear function, the feature bands with the same length return different bands in output space. Before expressing the theorem, we first introduce a formal notation of length and other necessary assumptions. For ease of discussion, we define in Definition 7 a type of band length slightly different from the previous analysis. We assume R below, albeit our analysis can be directly extended to high-dimensional cases.

Definition 7 (band length) For a given feature v and any perturbation $\psi \in C_f(v) = \{ \psi : \text{kv} - \psi k \leq Q \}$ in the feature band, we define the band length in the output space as the maximum distance between $g(v)$ and $g(\psi)$, namely

$$L_o(v), \max_{\psi \in C_f(v)} |g(\psi) - g(v)|;$$

Besides, we require Assumption 8, which is about the smoothness of the prediction head

Assumption 8. Assume that the prediction head g is second order derivative and α -smooth, namely, $\|k^2 g(u)\| \leq M$ for all feasible u .

The following Theorem 9 indicates that the variance of the band length is lower bounded, meaning that the bands given by Feature CP are individually different.

Theorem 9. Under Assumption 8, if the band on the feature space is with radius Q , then the variance of band length on the output space satisfies:

$$E[L_o - EL_o]^2 \geq Q^2 - E[\|k^2 g(v)\| - E\|k^2 g(v)\|]^2 \geq MQ E\|k^2 g(v)\|;$$

From Theorem 9, the variance of the band length has a non-vacuous lower bound if

$$E[\|k^2 g(v)\| - E\|k^2 g(v)\|]^2 > MQ E\|k^2 g(v)\| \quad (6)$$

We next discuss the condition for Equation (6). For a linear function g , note that $E[\|g(v)\|_k^2] = E[\|g(v)\|_k^2] = 0$ and $M = 0$, thus does not meet Equation (6). But for any other non-linear function g , we at least have $E[\|g(v)\|_k^2] > 0$ and $M > 0$, and therefore there exists a τ such that Equation (6) holds. Hence, the band length in feature space must be individually different for a non-linear function and a small band length.

Proof of Theorem 9. We revisit the notation in the main text, where $f(X)$ denotes the feature, and $C_f(v) = \{v : \|v - v^0\|_k \leq Q\}$ denotes the confidence band returned in feature space. By Taylor Expansion, for any given $v \in C_f(v)$, there exists α^0 such that

$$g(v) - g(v^0) = r g'(v^0)(v - v^0) + \frac{1}{2} r^2 g''(v^0)(v - v^0)^2.$$

Due to Assumption 8, $\|r^2 g''(v^0)\|_k \leq M$. Therefore, for any $v \in C_f(v)$

$$\|g(v) - g(v^0)\|_k \leq \frac{1}{2} M Q^2.$$

On the one hand, by Cauchy Schwarz inequality, we have

$$L_o = \max_v \|g(v) - g(v^0)\|_k \leq \frac{1}{2} M Q^2.$$

On the other hand, by setting $v = v^0 + Qr g'(v^0)j$, we have that

$$L_o = \max_v \|g(v) - g(v^0)\|_k \geq \|g(v^0 + Qr g'(v^0)j) - g(v^0)\|_k = Q \|r g'(v^0)j\|_k \geq \frac{1}{2} M Q^2.$$

Therefore, we have that

$$L_o = Q \|r g'(v^0)j\|_k \geq \frac{1}{2} M Q^2.$$

We finally show the variance of the length, where the randomness is taken over the data

$$\begin{aligned} E[L_o - E L_o]^2 &= E[\|Qr g'(v^0)j - E[Qr g'(v^0)j]\|_k^2] + E[\|L_o - Qr g'(v^0)j\|_k^2] \\ &= E[\|Qr g'(v^0)j - E[Qr g'(v^0)j]\|_k^2] + E[\|L_o - Qr g'(v^0)j\|_k^2] \\ &\quad + 2E[\|Qr g'(v^0)j - E[Qr g'(v^0)j]\|_k \|L_o - Qr g'(v^0)j\|_k] \\ &= Q^2 E[\|r g'(v^0)j - E[r g'(v^0)j]\|_k^2] \\ &\quad + 2QE[\|r g'(v^0)j - E[r g'(v^0)j]\|_k \|L_o - Qr g'(v^0)j\|_k] \\ &\quad + Q^2 E[\|r g'(v^0)j - E[r g'(v^0)j]\|_k^2] + MQ^3 E[\|r g'(v^0)j - E[r g'(v^0)j]\|_k]. \end{aligned}$$

Besides, note that $E[\|r g'(v^0)j - E[r g'(v^0)j]\|_k] = E[\|r g'(v^0)j\|_k] - E[r g'(v^0)j]$. Therefore, we have that

$$E[L_o - E L_o]^2 = Q^2 E[\|r g'(v^0)j - E[r g'(v^0)j]\|_k^2] + MQ^3 E[\|r g'(v^0)j - E[r g'(v^0)j]\|_k].$$

□

A.4 THEORETICAL CONVERGENCE RATE

In this section, we prove the theoretical convergence rate for the width. Specifically, we derive that when the number of samples in the calibration fold goes to infinity, the width for the testing point converges to a fixed value. Before we introduce the main theorem, we introduce some necessary definitions. Without further clarification, we follow the notations in the main text.

Definition 10 (Precise Band) We define the precise band as

$$C_1^{\text{pre}} = \{v : \|v - v^0\|_k \leq Q_1\} \quad (7)$$

Definition 11 (Precise Exact Band) We define the exact precise band as

$$C_1^{\text{pre}} = \{v : \|v - v^0\|_k \leq Q_1\} \quad (8)$$

where Q_1 denotes the exact value such that

$$P(\|v - v^0\|_k \leq Q_1; g(v) = y) = 1 \quad (9)$$

Our goal is to prove that the band length (volume) $\mathcal{C}_1^{\text{pre}}$ (denoted by $V(\mathcal{C}_1^{\text{pre}})$) converges to $V(\mathcal{C}_1^{\text{pre}})$. We assume that the prediction head and the quantile function are both Lipschitz in Assumption 12 and Assumption 13.

Assumption 12 (Lipschitz for Prediction Head) Assume that for any v, v^0 , we have

$$\|g(v) - g(v^0)\| \leq L_1 \|v - v^0\|.$$

Assumption 13 (Lipschitz for Inverse Quantile Function) Denote the quantile function as

$$\text{Quantile}(Q_u) = P(g(v) = y) = u:$$

We assume that its inverse function is Lipschitz, that is to say,

$$\|\text{Quantile}^{-1}(u) - \text{Quantile}^{-1}(u^0)\| \leq L_2 \|u - u^0\|.$$

Besides, we assume that the region $\mathcal{C}_1^{\text{pre}}$ has benign blow-up.

Assumption 14 (Benign Blow-up) Assume that $\mathcal{C}_1^{\text{pre}}$ has benign blow-up, that is to say, for the blow-up set $\mathcal{C}_1^{\text{pre}}(u) = \{v : g(v) = u\}$; $\|u - u^0\| \leq \epsilon$, we have

$$\|V(\mathcal{C}_1^{\text{pre}}(u)) - V(\mathcal{C}_1^{\text{pre}}(u^0))\| \leq c \epsilon;$$

where c denotes a constant independent of ϵ .

In the one-dimensional case \mathbb{R} , Assumption 14 easily holds. For the high-dimensional cases, such a bound usually requires that depends on the dimension.

Theorem 15 (Convergence Rate) Assume that the non-conformity scores in the calibration fold have no ties. Under Assumption 12, Assumption 13 and Assumption 14, we have that

$$\|V(\mathcal{C}_1^{\text{pre}}) - V(\mathcal{C}_1^{\text{pre}})\| \leq c L_1 L_2 \frac{1}{n}.$$

Proof. Firstly, as derived in Romano et al. (2019), when the non-conformity score in the calibration fold has no ties (the probability is zero), we have

$$P(g(v) = y) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(g(v) = y); \quad (10)$$

where v ; Q_1 denotes the surrogate feature, the trained feature, and the quantile value in Algorithm 3, respectively.

By Assumption 13 that the inverse quantile function is Lipschitz around Q_1 , we have

$$\|Q_1 - Q_1\| \leq L_2 \frac{1}{n}.$$

Therefore, for any $u \in \mathcal{C}_1^{\text{pre}}$, there exists $u^0 \in \mathcal{C}_1^{\text{pre}}$ such that

$$\|u - u^0\| \leq \epsilon, \quad \|g(v) - g(v^0)\| \leq L_1 \|v - v^0\| \leq L_1 L_2 \frac{1}{n}. \quad (11)$$

We note that bounding $\|g(v) - g(v^0)\|$ requires that the region of v, v^0 are both balls, and therefore one can select v^0 as the point with the smallest distance to v . Since the region $\mathcal{C}_1^{\text{pre}}$ has benign blow-up, we have that

$$|V(\mathcal{C}_1^{\text{pre}}) - V(\mathcal{C}_1^{\text{pre}})| \leq c L_1 L_2 \frac{1}{n}.$$

Besides, the following equation naturally holds due to Equation (10).

$$V(\mathcal{C}_1^{\text{pre}}) = V(\mathcal{C}_1^{\text{pre}}):$$

Therefore, we conclude with the following inequality,

$$\|V(\mathcal{C}_1^{\text{pre}}) - V(\mathcal{C}_1^{\text{pre}})\| \leq c L_1 L_2 \frac{1}{n}.$$

Therefore, as the sample size in the calibration fold goes to infinity, the length of the trained band converges to $V(\mathcal{C}_1^{\text{pre}})$.

□

Figure 5: The model architecture of the uni-dimensional and synthetic multi-dimensional target regression experiments. The dropout layers are omitted.

Figure 6: The model architecture of the semantic segmentation experiment.

B EXPERIMENTAL DETAILS

Section B.1 introduces the omitted experimental details. Section B.2 provide experimental evidence to validate cubic conditions. Section B.3 shows that Feature CP performs similarly to vanilla CP for untrained neural networks, validating that Feature CP works due to semantic information trained in feature space. Section B.4 introduces Feature CQR which applies feature-level techniques on CQR and Section B.5 reports the corresponding group coverage. Finally, Section B.9 provides other additional experiments omitted in the main text.

B.1 EXPERIMENTAL DETAILS

Model Architecture. The model architecture of the uni-dimensional and synthetic multi-dimensional target regression task is shown in Figure 5. The feature function and prediction head includes two linear layers, respectively. Moreover, the model architecture of the FCN used in the semantic segmentation experiment is shown in Figure 6, which follows the official implementation of PyTorch. The batch normalization and dropout layers are omitted in the figure. We use the ResNet50 backbone as default and take two convolution layers as We select the Layer4 output of ResNet50 as our surrogate feature.

Training protocols. In the unidimensional and synthetic dimensional target regression experiments, we randomly divide the dataset into training, calibration, and test sets with the proportion 0.8:0.1:0.1. As for the semantic segmentation experiment, because the labels of the pre-divided test set are not accessible, we re-split the training, calibration, and test sets randomly on the original training set of Cityscapes. We remove the class (unlabeled) from the labels during calibration and testing, and use the weighted mean square error as the training objective where the class weights are adopted from Paszke et al. (2016).

Band Estimation. In estimating the band length, we deploy Band Estimation based on the score calculated by Algorithm 2. In our experiment, we choose the number of Bands via Algorithm 2 via cross-validation.

Table 4: Comparison between CP and Feature CP.

METHOD DATASET	CP				FEATURE CP			
	COVERAGE		LENGTH		COVERAGE		LENGTH	
COMMUNITY	89:82	0:95	1:99	0:09	89:62	0:83	1:99	0:19
FACEBOOK1	90:11	0:33	3:17	0:59	90:07	0:32	2:08	0:17
FACEBOOK2	89:99	0:18	2:72	0:37	89:98	0:17	1:97	0:17
MEPS19	90:51	0:21	4:25	0:21	90:55	0:17	3:58	0:35
MEPS20	89:80	0:61	4:17	0:36	89:76	0:61	3:37	0:38
MEPS21	89:92	0:54	4:67	0:30	89:94	0:54	3:93	0:54
STAR	90:30	1:17	0:41	0:02	90:35	0:99	0:41	0:03
BIO	90:27	0:26	1:97	0:01	90:20	0:29	1:87	0:04
BLOG	90:08	0:27	3:32	0:38	90:06	0:33	2:81	0:40
BIKE	89:65	0:87	1:91	0:04	89:61	0:69	1:79	0:08

Randomness. We train each model ν times with different random seeds and report the mean and standard deviation value across all the runs as the experimental results (as shown in Figure 2 and Table 1).

Details of transforming segmentation classification problem into a regression task. The original semantic segmentation problem is to find the one-hot label whose size is $(C; W; H)$ via logistic regression, where C is the number of the classes, and W and H are the width and height of the image. We use Gaussian Blur to smooth the values in each channel. At this time, the smoothed label y ranges from 0 to 1. Then, we use the double log trick to convert the label space from $(1; 1)$, i.e., $y = \log(-\log(y))$. Finally, we use mean square error loss to train

Definition of weighted length. We formulate the weighted length as

$$\text{weighted length} = \frac{1}{|I|} \sum_{i \in I} \sum_{j \in J} w_i^{(j)} |C(X_i)_j|^{(i)};$$

where $w_i^{(j)}$ is the corresponding weight in each dimension. We remark that although the formulation of $w_i^{(j)}$ is usually sample-dependent, we omit the dependency of the sample and denote it by $w^{(j)}$ when the context is clear. We next show how to derive $w^{(j)}$ in practice.

Generally speaking, we hope that $w^{(j)}$ is large when being informative (i.e., in non-boundary regions). Therefore, for the j -th pixel after Gaussian Blur whose value $Y^{(j)} \in [0, 1]$, its corresponding weight is defined as

$$w^{(j)} = \frac{j^{2Y^{(j)}}}{W} \cdot 1j \in [0, 1];$$

where $W = \sum_j j^{2Y^{(j)}}$ and $1j$ is a scaling factor.

At a colloquial level, $w^{(j)}$ is close to 1 if $Y^{(j)}$ is close to 0 or 1. In this case, $Y^{(j)}$ being close to 0 or 1 means that the pixel is far from the boundary region. Therefore, the weight indicates the degree to which a pixel is being informative (not in object boundary regions).

Calibration details. During calibration, to get the best value for the number of steps, we take a subset (one-fifth) of the calibration set as the additional validation set. We calculate the non-conformity score on the rest of the calibration set with various values of M and then evaluate on the validation set to get the best one whose coverage is just over 0.95. The final trained surrogate feature \hat{y} is close to the true feature because $\frac{\|y - \hat{y}\|_2}{\|y\|_2} < 1\%$. In practice, the surrogate feature after optimization satisfies $\frac{\|y - \hat{y}\|_2}{\|y\|_2} < 1\%$.

Comparison between Feature CP and Vanilla CP. We next present the specific statistics of Figure 2 in Table 4 and Table 7.

Table 5: Validate cubic conditions.

SPACE	FEATURE SPACE		OUTPUT SPACE			
METRIC	$M_j Q_1$	$V_{D_{ca}}^f$	$V_{D_{ca}^j}^f$	$M[Q_1$	$H(V_{D_{ca}}^f)$	$H(V_{D_{ca}^j}^f)]$
COMMUNITY	0:1150	0:0290		0:8073	0:1450	
FACEBOOK1	0:2491	0:0391		1:8950	0:2058	
FACEBOOK2	0:2387	0:0960		1:7918	0:5220	
MEPS19	0:2403	0:0161		1:7511	0:1150	
MEPS20	0:2485	0:0571		1:7936	0:3532	
MEPS21	0:2528	0:0488		1:8686	0:3350	
STAR	0:0230	0:0025		0:1605	0:0123	
BIO	0:1051	0:0056		0:8509	0:0368	
BLOG	0:3537	0:1135		2:3769	0:5421	
BIKE	0:0921	0:0058		0:7759	0:0469	

Figure 7: Length v.s. cubic metric. Larger cubic metric implies better efficiency (shorter band).

B.2 CERTIFYING CUBIC CONDITIONS

In this section, we validate the cubic conditions. The most important component for the cubic condition is Condition 2, which claims that conducting the quantile step would not hurt much efficiency. We next provide experiment results in Table 5 on comparing the average distance between each sample to their quantile in feature space $M_j Q_1$, $V_{D_{ca}}^f$, $V_{D_{ca}^j}^f$ and in output space $M[Q_1$, $H(V_{D_{ca}}^f; D_{ca})$, $H(V_{D_{ca}^j}^f; D_{ca})]$. We here take $\alpha = 1$ for simplicity. The significant gap in Table 5 validates that the distance in feature space is significantly smaller than that in output space, although we did not consider the Lipschitz factor for computational simplicity.

Besides, we plot the relationship between the efficiency (band length) v.s. cubic metric in Figure 7. Specifically, cubic metric here represents the core statement in cubic condition (statement 2), which implies a metric form like $M_j Q_1$, $V_{D_{ca}}^f$, $V_{D_{ca}^j}^f$. The results are shown in Figure 7.

B.3 FEATURE CP WORKS DUE TO SEMANTIC INFORMATION IN FEATURE SPACE

Experiment results illustrate that feature-level techniques improve the efficiency of conformal prediction methods (e.g., Feature CP vs. CP, Feature CQR vs. CQR). We claim that exploiting the semantic information in feature space is the key to our algorithm. Different from most existing conformal prediction algorithms, which regard the base model as a black-box model, feature-level operations allow seeing the training process via the trained feature. This is novel and greatly broadens the scope of conformal prediction algorithms. For a well-trained base model, feature-level techniques improve efficiency by utilizing the powerful feature embedding abilities of well-trained neural networks.

In contrast, if the base model is untrained with random initialization (whose representation space does not have semantic meaning), Feature CP returns a similar band length as the baseline (see Table 6). This validates the hypothesis that Feature CP's success lies in leveraging the inductive bias of deep representation learning. Fortunately, realistic machine learning models usually contain meaningful information in the feature space, enabling Feature CP to perform well.

Table 6: Untrained base model comparison between conformal prediction and Feature CP. The base model is randomly initialized but not trained with the training fold. Experiment results show that Feature CP cannot outperform vanilla CP if the base model is not well-trained.

METHOD	VANILLA CP				FEATURE CP			
	COVERAGE		LENGTH		COVERAGE		LENGTH	
COMMUNITY	90:28	1:70	4.85	0:22	90:68	1:33	4:92	0:77
FACEBOOK1	90:15	0:15	3:42	0:25	90:16	0:12	3:20	0:50
FACEBOOK2	90:17	0:11	3:51	0:26	90:12	0:14	3:34	0:39
MEPS19	90:81	0:46	4:02	0:16	90:86	0:30	4:22	0:48
MEPS20	90:10	0:60	4:10	0:28	90:28	0:46	4:02	0:41
MEPS21	89:78	0:44	4:08	0:16	89:85	0:58	3:81	0:32
STAR	90:07	0:77	2:23	0:18	89:47	1:84	2:24	0:40
BIO	90:06	0:19	4:25	0:11	90:11	0:07	4:44	0:74
BLOG	90:13	0:34	2:41	0:15	90:16	0:26	2:58	0:49
BIKE	89:53	0:78	4:65	0:15	89:61	0:86	4:13	0:38

Algorithm 4 Feature Conformalized Quantile Regression (Feature CQR)

Require: Level α , dataset $\mathcal{D} = f(X_i; Y_i)_{i \in \mathcal{I}}$, test point X^0 ,

- 1: Randomly split the dataset into a training fold \mathcal{D}_{tr} , $(X_i; Y_i)_{i \in \mathcal{I}_{tr}}$ together with a calibration fold \mathcal{D}_{ca} , $(X_i; Y_i)_{i \in \mathcal{I}_{ca}}$;
- 2: Train a base machine learning model $f^{lo}(\cdot)$ and $f^{hi}(\cdot)$ using \mathcal{D}_{tr} to estimate the quantile of response Y_i , which returns $\hat{Y}_i^{lo}; \hat{Y}_i^{hi}$;
- 3: For each $i \in \mathcal{I}_{ca}$, calculate the index $s_i^{lo} = I(\hat{Y}_i^{lo} - Y_i)$ and $s_i^{hi} = I(\hat{Y}_i^{hi} - Y_i)$;
- 4: For each $i \in \mathcal{I}_{ca}$, calculate the non-conformity score $v_i^{lo} = v_i^{lo, c_1^0}$ where v_i^{lo} is derived on the lower bound function with Algorithm 2;
- 5: Calculate the $(1 - \alpha)$ -th quantile Q_1^{lo} of the distribution $\frac{1}{|\mathcal{I}_{ca}|} \sum_{i \in \mathcal{I}_{ca}} v_i^{lo} + 1$;
- 6: Apply Band Estimation on test data feature \mathbf{X}^0 with perturbation Q_1^{lo} and prediction head g^{lo} , which returns $[C_0^0; C_1^0]$;
- 7: Apply STEP 4-6 similarly with higher quantile, which returns $[C_0^1; C_1^1]$;
- 8: Derive $C_1^{fcqr}(X)$ based on Equation (12);

Ensure: $C_1^{fcqr}(X)$.

B.4 FEATURE CONFORMALIZED QUANTILE REGRESSION

In this section, we show that feature-level techniques are pretty general in that they can be applied to most of the existing conformal prediction algorithms. Specifically, We take Conformalized Quantile Regression (CQR, Romano et al. (2019)) as an example and propose Feature-level Conformalized Quantile Regression (Feature CQR). The core idea is similar to Feature CP (See Algorithm 3), where we conduct calibration steps in the feature space. We summarize the Feature CQR algorithm in Algorithm 4.

Similar to CQR, Algorithm 4 also considers the one-dimension case where $d=1$. We next discuss the steps in Algorithm 4. Firstly, different from Feature CP, Feature CQR follows the idea of CQR that the non-conformity score can be negative (see Step 4). Such negative scores help reduce the band length, which improves efficiency. This is achieved by the index calculated in Step 3. Generally, if the predicted value is larger than the true value $\hat{Y}_i > Y_i$, we need to adjust \hat{Y}_i^{lo} to be smaller, and vice versa. Step 8 follows the adjustment, where we summarize the criterion in Equation (12), given the two bands $[C_0^0; C_1^0]$ and $[C_0^1; C_1^1]$.

³Here, we set $\alpha = 1$ for simplicity

Table 7: Comparison between CQR and Feature CQR. Feature CQR achieves better efficiency while maintaining effectiveness.

METHOD	CQR				FEATURE CQR			
	COVERAGE		LENGTH		COVERAGE		LENGTH	
COMMUNITY	90:33	1:36	1:60	0:08	90:23	1:65	1.23	0:15
FACEBOOK1	89:94	0:23	1:15	0:04	92:00	0:20	1.00	0:04
FACEBOOK2	89:99	0:03	1:25	0:08	92:19	0:17	1.08	0:06
MEPS19	90:26	0:38	2:41	0:26	91:25	0:43	1.48	0:19
MEPS20	89:78	0:60	2:47	0:10	90:9	0:39	1.34	0:35
MEPS21	89:52	0:32	2:26	0:20	90:2	0:53	1.69	0:20
STAR	90:99	1:08	0:20	0:00	89:88	0:33	0.13	0:01
BIO	90:09	0:36	1:39	0:01	89:88	0:27	1.22	0:02
BLOG	90:15	0:15	1:47	0:05	91:49	0:26	0.89	0:03
BIKE	89:38	0:25	0:58	0:01	89:95	0:98	0.38	0:02

$$\begin{aligned}
 &\text{if } c_1^o < 0; c_1^{hi} < 0; \text{ return } C_1^{fcqr}(X) = [C_1^o; C_1^{hi}]; \\
 &\text{if } c_1^o < 0; c_1^{hi} > 0; \text{ return } C_1^{fcqr}(X) = [C_1^o; C_0^{hi}]; \\
 &\text{if } c_1^o > 0; c_1^{hi} < 0; \text{ return } C_1^{fcqr}(X) = [C_1^o; C_1^{hi}]; \\
 &\text{if } c_1^o > 0; c_1^{hi} > 0; \text{ return } C_1^{fcqr}(X) = [C_1^o; C_0^{hi}];
 \end{aligned}
 \tag{12}$$

Similar to Feature CP, we need a Band Estimation step to approximate the band length used in Step 6. One can change it into Band Detection if necessary. Different from Feature CP where Band Estimation always returns the upper bound of the band, Feature CQR can only approximate it. We conduct experiments to show that this approximation does not lose effectiveness since the coverage is always approximate to α . Besides, different from CQR, which considers adjusting the upper and lower with the same value, we adjust them separately, which is more flexible in practice (see Step 7).

We summarize the experiments result in Table 7. Feature CQR achieves better efficiency while maintaining effectiveness. Here we provide 90% confidence band using five repeated experiments with different random seeds.

B.5 GROUP COVERAGE FOR FEATURE CONFORMALIZED QUANTILE REGRESSION

This section introduces the group coverage returned by feature-level techniques, which implies the performance conditional coverage, namely $\mathbb{P}(\mathbb{Y} \in C(X) | X)$. Specifically, we split the test set into three groups according to their response values, and report the minimum coverage over each group.

We remark that the group coverage of feature-level conformal prediction stems from its vanilla version. That is to say, when the vanilla version has a satisfying group coverage, its feature-level version may also return a relatively satisfying group coverage. Therefore, we did not provide Feature CP here because vanilla CP cannot return a good group coverage.

We summarize the experiment results in Table 8. Although we did not provide a theoretical guarantee for group coverage, Feature CQR still outperforms vanilla CQR in various datasets in terms of group coverage. Among ten datasets, Feature CQR outperforms vanilla CQR in four datasets, and is comparable with vanilla CQR in five datasets. Although the advantage is not universal, improving group coverage via feature-level techniques is still possible.

We note that there is still one dataset where vanilla CQR outperforms Feature CQR. We attribute the possible failure reason of Feature CQR on the dataset FACEBOOK2 to the failure of base models. As stated in Section B.3, Feature CQR only works when the base model is well-trained. However, when grouping according to the returned values, it is possible that there exists one group that is not well-trained during the training process. This may cause the failure of Feature CQR on the dataset FACEBOOK2.

Table 8: Comparison between Feature CQR and CQR in terms of group coverage.

GROUP COVERAGE	FEATURE CQR	VANILLA CQR
COMMUNITY	76:05 5:13	78:77 4:17
FACEBOOK1	65:52 0:95	66:68 1:69
FACEBOOK2	65:66 1:41	70:78 1:39
MEPS19	76:67 2:17	71:26 1:23
MEPS20	77:70 0:90	71:26 3:20
MEPS21	74:71 2:36	70:74 1:83
STAR	84:62 2:77	82:20 5:72
BIO	84:80 1:05	80:03 1:48
BLOG	59:43 0:60	49:10 0:54
BIKE	81:07 1:65	78:22 2:44

Table 9: Ablation study of the number of layers l ($l = 0:1$) in unidimensional tasks, where the default setting is $l = 2; g = 2$.

DATASET	BIO				BIKE				BLOG			
	COVERAGE		LENGTH		COVERAGE		LENGTH		COVERAGE		LENGTH	
f : 2 g : 2	90:20	0:39	1:873	0:06	89:61	0:94	1:794	0:11	90:06	0:44	2:811	0:54
f : 3 g : 1	90:24	0:32	1:961	0:02	89:72	1:10	1:917	0:14	90:16	0:34	3:319	0:22
f : 1 g : 3	90:00	0:46	1:860	0:12	89:72	0:81	1:748	0:10	90:11	0:43	2:595	0:38

B.6 ROBUSTNESS OF SPLITTING POINT

As discussed in the main text, the empirical coverage of Feature CP is pretty robust to the splitting point. We show both experiments of small neural networks (Table 9) and large neural networks (Table 10). We remark that one can also apply standard cross-validation to find the best splitting point, which may return a shorter band.

B.7 FCP IN CLASSIFICATION

In this section, we show how to deploy FCP in classification problems. The basic ideas follow Algorithm 3, except for two points:

1. We use cross-entropy loss instead of MSE loss when calculating the non-conformity score.
2. We do not deploy LiPRA, but use a sampling method to return the empirical confidence band, which induces a looser bound.

For each sample in the calibration fold, we sample 100 samples in the confidence band of feature space, and return the corresponding predictions in output space to get the empirical confidence band. We summarize the experiment results in Table 2, where we follow the statistics of APS and RAPS in Angelopoulos et al. (2021b).

B.8 THE VALIDATION OF NON-CONFORMITY SCORE

Distribution for non-conformity score in calibration fold. We plot the distribution of calibration score in Figure 8. We plot each non-conformity score in the calibration fold. The distribution of non-conformity scores is smooth and single-peak in the real-world dataset, meaning that the proposed score is reasonable.

Besides, one may wonder what happens if the predicted function is close to the ground truth function, which may lead to a zero quantile of the non-conformity score. We conduct experiments to show that zero quantiles would not underestimate the empirical non-conformity score. Specifically, we create a dataset that the prediction target is given by $f(x) + \epsilon$, where f denotes a three-layer neural network, $\epsilon = 0$ with probability 0.901 and ϵ follows a standard Gaussian with probability 0.099.

Table 10: Ablation study of the number of layers γ ($\gamma = 0:1$) in large neural networks on ImageNet.

SPLITTING POINT (PREDICTION HEAD)	COVERAGE		LENGTH	
THE LAST LAYER	0:900	0:0018	3:32	0:10
LAST TWO LAYERS	0:901	0:0047	3:26	0:03
LAST THREE LAYERS	0:900	0:0029	3:30	0:01

Figure 8: The distribution of the calibration score for the segmentation task. This indicates that the definition of our non-conformity score is a proper one.

Here the 0:901 is to let the 90% quantile of the non-conformity score on a calibration set to be 0. Besides, we make the predicted model satisfy $\mathbb{P}(s \leq \gamma) = \gamma$. We summarize the results ($\gamma = 0:1$) as follows (we repeated experiments with different γ). Here s is generated with a uniform distribution between 0 and 1. Note that the zero-length is implied by zero quantiles.

However, in practice it is impossible for the predicted model to be therefore exact zero-quantile is impossible. We then summarize the results when the model is obtained via training on a training fold. The non-zero length implies that the quantile is non-zero.

B.9 ADDITIONAL EXPERIMENT RESULTS

This section provides more experiment results omitted in the main text.

Visualization for the segmentation problem. We also provide more visualization results for the segmentation problem in Figure 9.

Table 11: Zero quantile does not effects the empirical coverage ().

TRIAL	COVERAGE	LENGTH
1	90.63	0
2	90.08	0
3	89.72	0
4	89.66	0
5	90.21	0

Table 12: Zero quantile does not effects the empirical coverage ().

TRIAL	COVERAGE	LENGTH
1	90.55	0.1274
2	89.72	0.1105
3	89.75	0.1120
4	89.31	0.1202
5	90.10	0.1245

Figure 9: More visualization results for the Cityscapes segmentation task.