

Truthfulness of Decision-Theoretic Calibration Measures

Mingda Qiao

Massachusetts Institute of Technology

MINGDA.QIAO.CS@GMAIL.COM

Eric Zhao

University of California, Berkeley

ERIC.ZH@BERKELEY.EDU

Editors: Nika Haghtalab and Ankur Moitra

Abstract

Calibration measures quantify how much a forecaster’s predictions violates calibration, which requires that forecasts are unbiased conditioning on the forecasted probabilities. Two important desiderata for a calibration measure are its *decision-theoretic implications* [KPLST23] (i.e., downstream decision-makers that best-respond to the forecasts are always no-regret) and its *truthfulness* [HQYZ24] (i.e., a forecaster approximately minimizes error by always reporting the true probabilities). Existing measures satisfy at most one of the properties, but not both.

We introduce a new calibration measure termed *subsampled step calibration*, $\text{StepCE}^{\text{sub}}$, that is both decision-theoretic and truthful. In particular, on any product distribution, $\text{StepCE}^{\text{sub}}$ is truthful up to an $O(1)$ factor whereas prior decision-theoretic calibration measures suffer from an $e^{-\Omega(T)}\text{-}\Omega(\sqrt{T})$ truthfulness gap. Moreover, in any smoothed setting where the conditional probability of each event is perturbed by a noise of magnitude $c > 0$, $\text{StepCE}^{\text{sub}}$ is truthful up to an $O(\sqrt{\log(1/c)})$ factor, while prior decision-theoretic measures have an $e^{-\Omega(T)}\text{-}\Omega(T^{1/3})$ truthfulness gap. We also prove a general impossibility result for truthful decision-theoretic forecasting: any complete and decision-theoretic calibration measure must be discontinuous and non-truthful in the non-smoothed setting.

Keywords: calibration, truthfulness, decision-theoretic calibration, online learning

1. Introduction

Probabilistic forecasts play a central role in data-driven decision-making across broad application domains including finance, meteorology, and medicine [MW84, DF83, WM68, JOKOM12, KSB21, VCV15, BF⁺02, CAT16]. One of the greatest forms of utility provided by high-quality forecasting is that it enables downstream agents to confidently base their decision-making on the forecasts without any other knowledge of the future. For example, a weather station’s forecast of the probability of rain in the evening provides utility by informing individuals that may be debating whether to bring an umbrella to dinner. A related and widely studied requirement of forecasting is *calibration* [Bri50, Daw82, FV98], which requires that predicted probabilities align with long-run empirical frequencies of events. Calibration requires, for example, that it rains 70% of the days where the weather station forecasts a 70% chance of rain. Importantly, any downstream agent that bases their rational decision-making on perfectly calibrated forecasts will not incur any positive regret [FH21].

While perfect calibration is generally unachievable, there exist a number of calibration measures, such as expected calibration error (ECE) [FV98], smooth calibration error (SCE) [KF08], and U-Calibration (UCal) [KPLST23], which formalize a notion of approximate calibration and quantify deviations from perfect calibration. U-Calibration is a calibration measure of particular significance for decision-making applications as it is defined as the worst regret that a rational agent

can incur by blindly following a forecaster [KPLST23]. In contrast, most other calibration measures, such as smooth calibration error, are not “decision-theoretic” in that they do not provide guarantees for the regret of downstream agents.

Because the Bayes optimal classifier is perfectly calibrated, it seems natural to view calibration as incentivizing a forecaster to produce predictions that are consistent with their beliefs. This is not the case: forecasters that know the future are incentivized by most calibration measures to produce *non-truthful* predictions [FH21, QV21, HQYZ24]. For example, a forecaster may publicly forecast a 50% chance of rain even if they know for certain that there is a 100% chance of rain. To this end, [HQYZ24] proposed a set of desiderata for a calibration measure that includes, among common sense requirements like completeness (correct predictions have low error) and soundness (incorrect predictions have high error), a notion of *truthfulness*: a calibration measure should not penalize forecasters that know the future for predicting the true probabilities of events. [HQYZ24] formalizes truthfulness by defining a calibration measure having a truthfulness gap as the asymptotic separation between the expected value of a calibration measure on a truthful forecaster and on a strategic forecaster, when both know exactly the probability with which future events will occur. They also show that there exists a simple modification of smooth calibration error that is sound, complete, and truthful: compute smooth calibration error over randomly *subsampled* timesteps rather than the entire time horizon. However, the resulting calibration measure, like smooth calibration error, is not decision-theoretic in that it provides no meaningful guarantees for downstream agents.

In contrast, the U-Calibration measure is decision-theoretic but not truthful. This means that minimizing the expected worst-case regret of a downstream agent requires intentionally misrepresenting one’s knowledge of future events. This is in stark contrast to the maximization of downstream utilities, which always incentivizes an aligned forecaster to predict consistently with their beliefs [Bri50]. The literature leaves unresolved whether there exists any decision-theoretic calibration measure that both provides no-regret guarantees for downstream agents and satisfies the usual desiderata of truthfulness, soundness and completeness.

There are reasons to believe such a “best of all worlds” calibration measure is not possible. First, the technique used by [HQYZ24] to derive a truthful calibration measure from the Smooth Calibration Error measure [KF08] does not appear to suffice when applied to the U-Calibration measure. Second, the best responses of downstream agents are typically discontinuous in the forecasts they are given, and discontinuities are intimately connected with non-truthfulness [HQYZ24]. This raises the questions:

Is there a calibration measure that both provides decision-theoretic guarantees and incentivizes honest forecasting?

Is there a fundamental conflict between minimizing the regret of downstream agents and being truthful?

In this work, we show that the answer to both questions can be “Yes”: there is a fundamental conflict between truthfulness and decision-theoretic guarantees—but not with smoothed analysis, under which we can design a calibration measure that is the best of all worlds.

1.1. Overview of Results

U-Calibration is far from truthful. We identify two sources of non-truthfulness in U-Calibration that the subsampling technique of [HQYZ24] does not remedy. The first source arises from the

discontinuity of U-Calibration error. The second source is an incentive for forecasters to hedge their predictions, i.e., exaggerate their uncertainty, and materially contributes to the non-truthfulness of U-Calibration even under smoothed analysis.

As a result, even with subsampling and smoothed analysis, U-Calibration exhibits a *truthfulness gap*. We say a calibration measure has an α - β truthfulness gap, which we define formally in (4), if it gives a truthful forecaster an error of $\geq \beta$ but the optimal strategic forecaster’s error is below α .

Propositions 6, 9, and 10, Informally. *Both the U-Calibration measure UCal and its subsampled variant UCal^{sub} suffer from an $O(\sqrt{T})$ - $\Omega(T)$ truthfulness gap due to the discontinuity of UCal, and an $e^{-\Omega(T)}$ - $\Omega(\text{poly}(T))$ truthfulness gap due to the hedging incentives of UCal.*

We also prove a general impossibility result that suggests the non-truthfulness of U-Calibration is, to a degree, unavoidable. We later show that this result can be softened with smoothed analysis.

Proposition 7, Informally. *For any calibration measure, at least one of the following must be true:*

- *It is not complete: consistently forecasting a 50% chance of heads given a sequence of T fair coins does not yield an $O(\sqrt{T})$ error.*
- *It is not decision-theoretic: it does not always upper bound the regret of downstream agents.*
- *It is not truthful: there is an $O(\sqrt{T})$ - $\Omega(T)$ truthfulness gap.*

Step calibration error. We introduce *step calibration*: a sound, complete, and decision-theoretic calibration measure that provides no-regret guarantees for all downstream agents. Given a sequence of events $x_1, \dots, x_T \in \{0, 1\}$ and predictions $p_1, \dots, p_T \in [0, 1]$, the step calibration error is defined:

$$\text{stepCE}(x, p) := \sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \alpha] \right|.$$

Step calibration is equivalent, up to a constant factor, to a variant of V-Calibration that uses a slightly different baseline to disincentivize hedging behavior by penalizing excessively conservative probabilistic forecasts (Fact 11). In addition to step calibration being a complete and sound calibration measure, we also demonstrate an algorithm that achieves an $\tilde{O}(\sqrt{T})$ step calibration error for the adversarial prediction setting.¹

Proposition 12 and Theorem 23, Informally. *The step calibration error is sound, complete, and decision-theoretic. Moreover, there is a forecasting algorithm that guarantees an expected step calibration error of $O(\sqrt{T \log T})$, even if the events are adversarially and adaptively chosen.*

Truthfulness under smoothed analysis. We show that—under smoothed analysis—the impossibility of simultaneously providing decision-theoretic guarantees and truthfulness largely disappears.

1. Concurrent works by [OKK25] and [RSB⁺25] reached a closely related definition of decision-theoretic calibration, though interestingly these works are motivated by—in contrast to truthfulness—the study of omniprediction and testable calibration measures.

At a high-level, smoothing negates an adversary’s ability to exploit the inherently discontinuous nature of a downstream agent’s decision-making. Importantly, we show that the “subsamped” variant of step calibration,

$$\text{stepCE}^{\text{sub}}(x, p) := \mathbb{E}_{S \sim \text{Unif}(2^{[T]})} \left[\sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \alpha \wedge t \in S] \right| \right],$$

is truthful under smoothed analysis.

We say that a calibration measure is (α, β) -truthful gap if, given any prior distribution over the sequence of events, the error incurred by the truthful forecaster is upper bounded by the optimal strategic forecaster’s error, up to a factor of α and an additive term of β ; see Equation (4) for a formal definition.

Theorem 16, Informally. *Subsampled step calibration error is $(O(\sqrt{\log(1/c)}), \text{polylog}(T/c))$ -truthful when each conditional probability is drawn from a distribution with density $\leq 1/c$.*

This $O(\sqrt{\log(1/c)})$ factor is tight for $\text{stepCE}^{\text{sub}}$. For non-smoothed product distributions, we can obtain a stronger truthfulness result showing that $\text{stepCE}^{\text{sub}}$ is $(O(1), 0)$ -truthful (Proposition 22). $\text{stepCE}^{\text{sub}}$ also retains the desiderata of stepCE , and is decision-theoretic, sound, and complete, and admits an $\tilde{O}(\sqrt{T})$ algorithm in the adversarial setting. This is because, as we show in Lemma 24, $\frac{1}{2}\text{stepCE}(x, p) \leq \text{stepCE}^{\text{sub}}(x, p) \leq \frac{1}{2}\text{stepCE}(x, p) + O(\sqrt{T})$.

1.2. Related Work

Most closely related to our work are the previous studies of sequential binary calibration with respect to various calibration measures. The seminal work of [FV98] showed that asymptotic calibration can be achieved even for adversarially chosen events. Implicit in their paper is a sublinear rate of $O(T^{2/3})$ on the ECE incurred by the forecaster; a more detailed proof was given by [Har22]. On the lower bound side, an $\Omega(\sqrt{T})$ bound is trivial and the first non-trivial lower bound of $\Omega(T^{0.528})$ was shown by [QV21]. A recent breakthrough of [DDF⁺24] improved the upper bound to $O(T^{2/3-\varepsilon})$ for some constant $\varepsilon > 0$, and gave the best known lower bound of $\Omega(T^{0.54389})$.

The analogous question has been studied for other calibration measures, including smooth calibration [KF08, QZ24], U-Calibration [KPLST23], distance from calibration [QZ24, ACRS25], calibration decision loss [HW24], and subsampled smooth calibration [HQYZ24]. These calibration measures relax the ECE in different ways, so that the forecaster can achieve a faster rate of $\tilde{O}(\sqrt{T})$, circumventing the super- \sqrt{T} lower bounds for the ECE.

The systematic study of calibration measures was initiated by [BGHN23], who focused on the offline setting and proposed the *distance from calibration* as a ground truth. Their work identified calibration measures that are continuous and consistent (i.e., being polynomially related to the distance from calibration). Subsequent work studied calibration measures that satisfy other natural axioms, including being decision-theoretic [KPLST23, NRRX23, RS24, HW24] and being truthful [HQYZ24]. Remarkably, the distance from calibration, while being a natural measure, is neither decision-theoretic nor truthful in the sequential setting. We also note that the proper calibration and cutoff calibration measures introduced by concurrent works [OKK25] and [RSB⁺25] in their study of testable calibration measures match our definition of *step calibration* up to constant factors.

The truthfulness of calibration measures is also intimately connected to the study of multi-calibration [HJKRR18], which designs calibration measures that incentivize forecasters to “truthfully” predict the true probabilities of each feature by evaluating calibration error over many feature

subsets—a practice similar to [HQYZ24]’s use of subsampling to enforce truthfulness and algorithmically linked to Blackwell approachability [Bla56, Fos99, Har22, HJZ23]. Truthfulness can also be viewed as enforcing an online, multi-timestep notion of outcome indistinguishability [DKR⁺21].

Smoothed analysis was introduced by Spielman and Teng [ST04, ST09] for analyzing the “typical” runtime of the simplex method. [KST09] introduced a smoothed analysis model for supervised learning, in which the data distribution is a product distribution over $\{0, 1\}^n$ with marginal probabilities randomly perturbed. This circumvents hard instances that are specific for the uniform distribution but easily learnable under the perturbed distribution. Subsequent work studied tensor decomposition [BCMV14] and decision tree learning [BDM20, BLQT21] in similar setups.

More closely related to our work are the smoothed analysis for online learning introduced by [RST11]. A recent series of work extends this setting to adaptive adversaries, showing that online learning against a smoothed adversary is not much harder than learning in the offline (batch) setup [HRS20, HRS24, BDGR22, HHSY22, BS22, BP23, BSR23, BST23]. In these models, the smoothed analysis limits the adversary’s ability of concentrating the probability mass at a “hard region” in the instance space. As a result, the learner may circumvent the canonical hard instance of threshold functions, which is easily learnable in the offline setting, and cannot be learned in an online setting without smoothing. Our work applies the smoothed analysis to avoid the large truthfulness gap in the non-smoothed setting following the same intuition.

2. Preliminaries

2.1. Sequential Prediction and Calibration

Sequential prediction. In the basic (non-smoothed) prediction setup, a sequence of events $x \in \{0, 1\}^T$ is sampled from a distribution \mathcal{D} . At each time step $t \in [T]$, the forecaster makes a prediction $p_t \in [0, 1]$, after which x_t is revealed. Formally, a deterministic forecaster is a function $\mathcal{A} : \bigcup_{t=1}^T \{0, 1\}^{t-1} \rightarrow [0, 1]$, where $\mathcal{A}(b_1, b_2, \dots, b_{t-1})$ specifies the forecaster’s prediction at step t upon observing $x_{1:(t-1)} = b_{1:(t-1)}$. We will write $(x, p) \sim (\mathcal{D}, \mathcal{A})$ to denote sampling events $x \in \{0, 1\}^T$ and predictions $p \in [0, 1]^T$ from the joint distribution naturally induced by distribution \mathcal{D} and forecaster \mathcal{A} , i.e., by sampling $x \sim \mathcal{D}$ and setting $p_t = \mathcal{A}(x_1, x_2, \dots, x_{t-1})$ for each $t \in [T]$. We could have defined the forecaster to be randomized or a function of both the outcomes $x_{1:(t-1)}$ and its own predictions $p_{1:(t-1)}$, but restricting to deterministic functions of the outcomes $x_{1:(t-1)}$ comes without loss of generality.

The smoothed setting. The prediction setting above can be equivalently viewed as the nature specifying the conditional probability of $x_t = 1$ given x_1, x_2, \dots, x_{t-1} . We will also consider a *smoothed* setting, where each conditional probability is perturbed by a noise of magnitude $c > 0$. Formally, the nature specifies a mapping $\mathcal{P} : \bigcup_{t=1}^T \{0, 1\}^{t-1} \mapsto \Delta_c$, where Δ_c is the family of distributions over $[0, 1]$ with densities bounded by $1/c$ everywhere. At each step t , the nature realizes $p_t^* \sim \mathcal{P}(x_1, x_2, \dots, x_{t-1})$ and credibly reveals the value of p_t^* to the forecaster.² The forecaster predicts p_t , and the event x_t is sampled from Bernoulli(p_t^*) and revealed. Formally, the forecaster’s prediction p_t is a function of both $x_{1:(t-1)}$ and p_t^* , i.e., $\mathcal{A} : \bigcup_{t=1}^T (\{0, 1\}^{t-1} \times [0, 1]) \rightarrow [0, 1]$.

2. In practice, this corresponds to the forecaster acquiring certain side information about x_t , thus changing their belief of $p_t^* = \Pr[x_t = 1 \mid \text{observations}]$. The smoothness assumption would then correspond to assumptions on the side information, which might ensure that the distribution of p_t^* is not too spiky.

Note that the non-smoothed setting—in which the nature specifies a fixed distribution \mathcal{D} over $\{0, 1\}^T$ —can be viewed as a “0-smoothed” setting where each $\mathcal{P}(b_1, b_2, \dots, b_{t-1})$ is the degenerate distribution at value $\Pr_{x \sim \mathcal{D}} [x_t = 1 \mid x_{1:(t-1)} = b_{1:(t-1)}]$. Furthermore, as in the non-smoothed setting, the nature \mathcal{P} and the forecaster \mathcal{A} naturally induce a joint distribution over the triple (x, p^*, p) , denoted by sampling $(x, p^*, p) \sim (\mathcal{P}, \mathcal{A})$ (or a subset thereof) in the rest of the paper.

One way to interpret the smoothed setting is that by limiting Nature’s precision, we limit its ability to select pathological edge cases. We note that this is different from assuming away the possibility that Nature introduces rare events, which the smoothed setting still allows for.

Calibration measures. A calibration measure $\text{CM}_T : \{0, 1\}^T \times [0, 1]^T \rightarrow [0, T]$ quantifies the quality of a forecaster’s prediction. We omit the subscript T when it is clear from context. The expected penalty incurred by forecaster \mathcal{A} on distribution \mathcal{D} is defined as $\text{err}_{\text{CM}}(\mathcal{D}, \mathcal{A}) := \mathbb{E}_{(x,p) \sim (\mathcal{D}, \mathcal{A})} [\text{CM}(x, p)]$. For the smoothed setting, we define $\text{err}_{\text{CM}}(\mathcal{P}, \mathcal{A}) := \mathbb{E}_{(x,p) \sim (\mathcal{P}, \mathcal{A})} [\text{CM}(x, p)]$.

One would naturally expect a calibration measure to be both *complete* (accurate predictions lead to a small penalty) and *sound* (inaccurate predictions receive a large penalty). We adopt a variant of the definition in [HQYZ24]. In the following, $\mathbf{1}_T$ denotes the T -dimensional all-1 vector.

Definition 1 (Completeness and soundness [HQYZ24]) *A calibration measure CM is complete if: (1) For any $x \in \{0, 1\}^T$, predicting the events x gives $\text{CM}_T(x, x) = 0$; (2) For any $\alpha \in [0, 1]$, predicting the constant probability α gives $\mathbb{E}_{x_1, \dots, x_T \sim \text{Bernoulli}(\alpha)} [\text{CM}_T(x, \alpha \cdot \mathbf{1}_T)] = o_\alpha(T)$; (3) For any $p \in [0, 1]^T$ that is perfectly calibrated with respect to $x \in \{0, 1\}^T$, $\text{CM}_T(x, p) = o(T)$. The calibration measure is sound if: (1) For any $x \in \{0, 1\}^T$, $\text{CM}_T(x, \mathbf{1}_T - x) = \Omega(T)$; (2) For any $\alpha, \beta \in [0, 1]$ such that $\alpha \neq \beta$, $\mathbb{E}_{x_1, \dots, x_T \sim \text{Bernoulli}(\alpha)} [\text{CM}_T(x, \beta \cdot \mathbf{1}_T)] = \Omega_{\alpha, \beta}(T)$. Here, $o_\alpha(\cdot)$ and $\Omega_{\alpha, \beta}(\cdot)$ hide constant factors that depend on the subscripted parameters.*

We strengthened the definition of completeness in [HQYZ24] by adding a third constraint—perfectly calibrated predictions should receive a low (sublinear) penalty. To the best of our knowledge, all calibration measures satisfy this condition, with most satisfying this condition with $\text{CM}(x, p) = 0$ while the SSCE introduced by [HQYZ24] satisfies this with $\text{CM}(x, p) = O(\sqrt{T})$.

2.2. Decision-Theoretic Calibration

Consider a decision-making setting where an agent acts on the basis of the forecaster’s predictions. Formally, consider a repeated game where, at each round t , an agent chooses an action $a_t \in \mathcal{A}$ informed by the forecaster’s prediction p_t . The agent’s utility at time t is a function $u : \mathcal{A} \times \{0, 1\} \rightarrow [-1, 1]$ of its action a_t and the event x_t . In this setup, the agent assumes that p_t is an accurate forecast of the probability that x_t occurs and thus selects a_t to maximize $\mathbb{E}_{x \sim \text{Bernoulli}(p_t)} [u(a_t, x)]$. One benefit of calibrated forecasts is that agents can make decisions according to the forecasts and be no-regret [FH21]. That is, the agent’s expected cumulative utility $\sum_t u(a_t, x_t)$ when given calibrated forecasts $p_{1:T}$ will never be worse than when given a *base rate forecaster* which predicts $p_t = \frac{1}{T} \sum_{t=1}^T x_t$ for all t .

Providing forecasts with no-regret guarantees for agents with any utility is equivalent to providing forecasts that satisfy the no-regret property with respect to any (piecewise linear) proper scoring rule [KPLST23]. A (bounded) scoring rule is a function $S : \{0, 1\} \times [0, 1] \rightarrow [-1, 1]$, where $S(x, p)$ denotes the loss incurred by the forecaster, when it predicts value $p \in [0, 1]$ on an outcome that turns out to be $x \in \{0, 1\}$. A scoring rule S is *proper* if, for any $\alpha \in [0, 1]$, the function $p \mapsto \mathbb{E}_{x \sim \text{Bernoulli}(\alpha)} [S(x, p)]$ is minimized at $p = \alpha$, i.e., predicting the true probability minimizes

the expected loss. The U-Calibration error of [KPLST23] quantifies the mis-calibration using the worst possible external regret:

$$\text{UCal}(x, p) := \sup_S \left[\sum_{t=1}^T S(x_t, p_t) - \inf_{\beta \in [0,1]} \sum_{t=1}^T S(x_t, \beta) \right], \quad (1)$$

where the supremum is taken over all proper scoring rules S .

More generally, we refer to any calibration measure CM as being *decision-theoretic* if, for all events x and predictions p , the calibration measure is lower bounded by U-Calibration up to a universal constant factor: $\text{CM}(x, p) \geq \Omega(1) \cdot \text{UCal}(x, p)$. A decision-theoretic calibration measure upper bounds the external regret of any agent that acts on the forecaster's predictions.

V-Calibration. We will work with a calibration measure known as V-Calibration [KPLST23] that is more technically convenient. V-Calibration is a modification of U-Calibration obtained from limiting the supremum in its definition (Equation (1)) to a narrow class of scoring rules of the form

$$S_\alpha(x, p) := (\alpha - x) \cdot \text{sgn}(p - \alpha) \quad (2)$$

for any $\alpha \in [0, 1]$. Formally, the V-Calibration error is defined as

$$\text{VCal}(x, p) := \sup_{\alpha, \beta \in [0,1]} \left[\sum_{t=1}^T S_\alpha(x_t, p_t) - \sum_{t=1}^T S_\alpha(x_t, \beta) \right]. \quad (3)$$

Despite its simpler form, V-Calibration is equivalent to U-Calibration up to constant factors:

Lemma 2 (Theorem 8 of [KPLST23]) *For any $x \in \{0, 1\}^T$ and $p \in [0, 1]^T$, it holds that*

$$\frac{1}{2} \text{UCal}(x, p) \leq \text{VCal}(x, p) \leq \text{UCal}(x, p).$$

We can rewrite the V-Calibration measure into an alternative form without the scoring rules. We prove the following proposition in Section C.1.

Proposition 3 *The V-Calibration error takes the alternative form*

$$\text{VCal}(x, p) = 2 \cdot \sup_{\alpha \in [0,1]} \max\{X_-^{(\alpha)} - \alpha N_-^{(\alpha)}, \alpha N_+^{(\alpha)} - X_+^{(\alpha)}\},$$

where $N_-^{(\alpha)} := \sum_{t=1}^T \mathbb{1}[p_t < \alpha]$, $N_+^{(\alpha)} := \sum_{t=1}^T \mathbb{1}[p_t > \alpha]$, $X_-^{(\alpha)} := \sum_{t=1}^T x_t \cdot \mathbb{1}[p_t < \alpha]$, and $X_+^{(\alpha)} := \sum_{t=1}^T x_t \cdot \mathbb{1}[p_t > \alpha]$.

2.3. Truthfulness

Calibration measures are often seen as measuring how close a forecaster's predictions are to the true probabilities that events occur. However, even if one knows the exact probability that an event x_t will occur, calibration does not necessarily incentivize one to predict truthfully [HQYZ24]. Consider the truthful forecaster for the non-smoothed setting specified by $\mathcal{D} \in \Delta(\{0, 1\}^T)$,

$$\mathcal{A}^{\text{truthful}}(\mathcal{D})(b_1, b_2, \dots, b_{t-1}) := \Pr_{x \sim \mathcal{D}} [x_t = 1 \mid x_{1:(t-1)} = b_{1:(t-1)}],$$

which can be argued to be the only forecaster that makes the “right” predictions on distribution \mathcal{D} . Given a reasonable calibration measure CM , one might expect the error of the truthful forecaster to be close to the optimal error $\text{OPT}_{\text{CM}}(\mathcal{D}) := \inf_{\mathcal{A}} \text{err}_{\text{CM}}(\mathcal{D}, \mathcal{A})$, where \mathcal{A} ranges over all deterministic forecasters. This property is known as *truthfulness* [HQYZ24], where we say that a calibration measure CM is (α, β) -truthful if, for every $\mathcal{D} \in \Delta(\{0, 1\}^T)$,

$$\text{err}_{\text{CM}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D})) \leq \alpha \cdot \text{OPT}_{\text{CM}}(\mathcal{D}) + \beta. \quad (4)$$

Conversely, CM is said to have an α - β truthfulness gap if, for some distribution \mathcal{D} , $\text{OPT}_{\text{CM}}(\mathcal{D}) \leq \alpha$ and $\text{err}_{\text{CM}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D})) \geq \beta$.

In smoothed settings where the conditional probabilities are sampled according to \mathcal{P} and revealed to the forecaster, the truthful forecaster, $\mathcal{A}^{\text{truthful}}$, simply maps $(b_{1:(t-1)}, p_t^*)$ to p_t^* for any $t \in [T]$ and $(b_{1:(t-1)}, p_t^*) \in \{0, 1\}^{t-1} \times [0, 1]$. We define $\text{OPT}_{\text{CM}}(\mathcal{P}) := \inf_{\mathcal{A}} \text{err}_{\text{CM}}(\mathcal{P}, \mathcal{A})$, and a CM as being (α, β) -truthful if $\text{err}_{\text{CM}}(\mathcal{P}, \mathcal{A}^{\text{truthful}}) \leq \alpha \cdot \text{OPT}_{\text{CM}}(\mathcal{P}) + \beta$. We similarly define the α - β truthfulness gap for smoothed settings.

U-Calibration is known to not be a truthful calibration measure [HQYZ24]. This might be counterintuitive since, by definition, truthful forecasting minimizes the expected penalty for each *fixed* proper scoring rule. However, after taking a supremum over all proper scoring rules, the truthful forecaster ceases to be optimal for the resulting measure.

Proposition 4 (Proposition A.3 of [HQYZ24]) *The U-Calibration error has an $O(1)$ - $\Omega(\sqrt{T})$ truthfulness gap.*

One example of a truthful calibration measure is the *Subsampled Smooth Calibration Error* (SSCE) introduced by [HQYZ24]. SSCE is a variant of the *smooth calibration error* calibration measure introduced by [KF08]: $\text{smCE}(x, p) := \sup_{f \in \mathcal{F}} \sum_{t=1}^T f(p_t)(x_t - p_t)$, where \mathcal{F} is the family of 1-Lipschitz functions from $[0, 1]$ to $[-1, 1]$. SSCE is defined by subsampling a subset of the time horizon, and evaluating the Smooth Calibration Error on it. Formally, letting $\text{Unif}(S)$ denote the uniform distribution over a finite set S and $x|_S$ denote the $|S|$ -dimensional vector formed by the entries of x indexed by S :

$$\text{SSCE}(x, p) := \mathbb{E}_{S \sim \text{Unif}(2^{[T]})} [\text{smCE}(x|_S, p|_S)] = \mathbb{E}_{y \sim \text{Unif}(\{0, 1\}^T)} \left[\sup_{f \in \mathcal{F}} \sum_{t=1}^T y_t \cdot f(p_t) \cdot (x_t - p_t) \right]. \quad (5)$$

In light of Proposition 7, since SSCE is complete and truthful, it cannot be decision-theoretic.

3. Technical Overview

3.1. Non-truthfulness of U-Calibration and Its Variants

The non-truthfulness of the U-Calibration error (Proposition 4) comes from the incentive for a dishonest forecaster to “patch up” their previous mis-calibration. Specifically, [HQYZ24] considered a length- T sequence that consists of $T/2$ independent random bits followed by $T/2$ ones. In this case, the truthful forecaster predicts $1/2$ in the first half, and typically incurs an $\Omega(\sqrt{T})$ bias on those bits. This further translates into an $\Omega(\sqrt{T})$ U-Calibration error. On the other hand, a strategic forecaster may deliberately predict a biased value of $5/8$ on the first half, and continue predicting

5/8 on the second half until the bias is close to 0.³ The resulting U-Calibration error can then be bounded by $O(1)$ in expectation.

This $O(1)$ - $\Omega(\sqrt{T})$ truthfulness gap, however, would vanish once we apply the *subsampling* technique of [HQYZ24]. The subsampled version, UCal^{sub} , evaluates the U-Calibration error on a random subset of the horizon. This introduces a $\Theta(\sqrt{T})$ error in the resulting penalty, so the strategic forecaster at best outperforms the truthful one by a constant factor. One might naturally wonder whether UCal^{sub} is truthful in general. Unfortunately, as we outline below, there exist two additional “failure modes” of the U-Calibration error that cannot be remedied by subsampling alone.

Example 1: Non-truthfulness due to discontinuity. We start by noting that the U-Calibration error, $\text{UCal}(x, p)$, is not continuous in p . Suppose that, for some small $\varepsilon > 0$, we have

$$(x_t, p_t) = \begin{cases} (1, 1/2 - \varepsilon), & t \leq T/2, \\ (0, 1/2 + \varepsilon), & t > T/2. \end{cases}$$

Note that p is almost calibrated: $\tilde{p} = (1/2, 1/2, \dots, 1/2)$ is entry-wise close to p , and perfectly calibrated with respect to x , which implies $\text{UCal}(x, \tilde{p}) = 0$. However, $\text{UCal}(x, p)$ is much larger: Consider the equivalent formulation of the V-Calibration error in Proposition 3 and take $\alpha = 1/2$. There are $N_- = T/2$ steps on which $p_t < \alpha$, and the outcomes on those steps sum up to $X_- = T/2$. By Lemma 2, we have $\text{UCal}(x, p) \geq \text{VCal}(x, p) \geq 2(X_- - \alpha N_-) = \Omega(T)$. Taking $\varepsilon \rightarrow 0^+$ gives triples (x, p, \tilde{p}) such that $\|p - \tilde{p}\|_\infty \rightarrow 0$ but $\text{UCal}(x, p) = \Omega(T)$ and $\text{UCal}(x, \tilde{p}) = 0$ are far away.

This implies that any complete and decision-theoretic calibration measure CM must be discontinuous. For the triple (x, p, \tilde{p}) constructed as above, completeness (Definition 1) implies $\text{CM}(x, \tilde{p}) = o(T)$ while being decision-theoretic requires $\text{CM}(x, p) \geq \Omega(1) \cdot \text{UCal}(x, p) = \Omega(T)$. Taking $\varepsilon \rightarrow 0^+$ shows that CM is discontinuous.

The example above does not immediately give the $O(\sqrt{T})$ - $\Omega(T)$ truthfulness gap in Proposition 6. Towards showing that the truthful forecaster incurs an $\Omega(T)$ U-Calibration error, we need to design a sequence of true probabilities $p_1^*, p_2^*, \dots, p_T^* \approx 1/2$ and a threshold $\alpha \in [0, 1]$, such that $x_t = 1$ whenever $p_t^* < \alpha$ and $x_t = 0$ whenever $p_t^* > \alpha$. This is *very* unlikely to happen if each x_t is independently sampled from $\text{Bernoulli}(p_t^*)$.

However, when nature picks each p_t^* based on the previous outcomes $x_{1:(t-1)}$, the hoped-for property *can* be guaranteed via a simple binary search. At step $t = 1$, the nature starts with an interval $[l_1, r_1] = [1/2 - \varepsilon, 1/2 + \varepsilon]$ and picks $p_1^* = (l_1 + r_1)/2$ as the middle point. After realizing $x_1 \sim \text{Bernoulli}(p_1^*)$, if $x_1 = 1$, the nature updates $[l_2, r_2] \leftarrow [p_1^*, r_1]$; otherwise, $[l_2, r_2] \leftarrow [l_1, p_1^*]$. If the nature repeats this T steps, we can verify that, for $\alpha := (l_{T+1} + r_{T+1})/2$ and every $t \in [T]$: (1) $p_t^* < \alpha$ implies $x_t = 1$; (2) $p_t^* > \alpha$ implies $x_t = 0$. Then, a similar argument shows that truthful forecasting leads to $\text{UCal}(x, p^*) = \Omega(T)$. Furthermore, this argument is robust to subsampling, i.e., we also have $\text{UCal}^{\text{sub}}(x, p^*) = \Omega(T)$. In contrast, a strategic forecaster might choose to predict $p_t = 1/2$ at every step. Since each p_t^* is in $[1/2 - \varepsilon, 1/2 + \varepsilon]$, as long as $\varepsilon = O(1/\sqrt{T})$, $\sum_{t=1}^T x_t$ would concentrate around $T/2 \pm O(\varepsilon T + \sqrt{T}) = T/2 \pm O(\sqrt{T})$. This shows that $\text{OPT}_{\text{UCal}}, \text{OPT}_{\text{UCal}^{\text{sub}}} = O(\sqrt{T})$, and thus the $O(\sqrt{T})$ - $\Omega(T)$ truthfulness gap.

We also further generalize this example to show that no calibration measure can be complete, decision-theoretic and non-trivially truthful simultaneously.

3. This happens with high probability, as the mean of the first half concentrates around $1/2 < 5/8$, while the mean of the entire sequence concentrates around $3/4 > 5/8$.

Example 2: Non-truthfulness due to hedging. In the previous example, it is crucial that nature specifies the conditional expectation of each bit (p_t^*) adaptively and with arbitrary precision. Nevertheless, we show that UCal can still have a large truthfulness gap, even if the events are drawn from a product distribution, the marginal probabilities of which are, in turn, drawn from smooth distributions. At a high level, this is because the regret minimization for downstream agents incentivizes *hedging behaviors*, where forecasters benefit from exaggerating the uncertainty of future events.

We start with a simple, non-smoothed setting: For each $t \in [T]$, we set $p_t^* = 1/5$ if $t \leq T/2$, and $p_t^* = 4/5$ if $t > T/2$. Each x_t is independently sampled from $\text{Bernoulli}(p_t^*)$. As in Example 1, truthful prediction typically leads to an $\Omega(\sqrt{T})$ bias at predicted values $1/5$ and $4/5$ each, and results in an $\Omega(\sqrt{T})$ U-Calibration error. In contrast, the forecaster can significantly lower its penalty by predicting $p_t = 2/5$ at $t \leq T/2$ and $p_t = 3/5$ at $t > T/2$ instead. In light of Lemma 2 and Proposition 3, it suffices to upper bound the value of $\max\{X_- - \alpha N_-, \alpha N_+ - X_+\}$ for different values of α . The worst cases are when $\alpha \rightarrow (2/5)^+$ and $\alpha \rightarrow (3/5)^-$. In the former case, we have $N_- = T/2$ while X_- concentrates around $(T/2) \cdot (1/5) = T/10$, and is typically (except with probability $e^{-\Omega(T)}$) smaller than $\alpha N_- = T/5$. Similarly, when $\alpha \approx 3/5$, we have $N_+ = T/2$ and X_+ concentrates around $(T/2) \cdot (4/5) = 2T/5$, and is extremely unlikely to be below $\alpha N_+ = 3T/10$. This establishes the $e^{-\Omega(T)}\text{-}\Omega(\sqrt{T})$ truthfulness gap for UCal, and the same construction works for the subsampled version UCal^{sub} as well.

In the c -smoothed setting for some small constant c , instead of setting each p_t^* to $1/5$ or $4/5$, we draw each p_t^* independently and uniformly from either $[1/5 - c, 1/5 + c]$ or $[4/5 - c, 4/5 + c]$. Here, a complication is that we cannot “catch” a high U-Calibration error by naïvely setting $\alpha = 1/5 + c$ (and applying Lemma 2 and Proposition 3). This is because we would end up with $N_- = T/2$ and X_- concentrating around $(T/2) \cdot (1/5) = T/10$, which is *lower* than $\alpha N_- = (1/5 + c) \cdot (T/2) = T/10 + \Omega(T)$. Instead, we pick $\alpha = 1/5 - (1 - \gamma)c$ for some $\gamma > 0$ to be chosen. Since a uniform sample from $[1/5 - c, 1/5 + c]$ falls into $[1/5 - c, \alpha]$ with probability $(\gamma c)/2c = \gamma/2$, we expect $N_- = \Theta(\gamma T)$. Furthermore, conditioning on that $p_t^* \in [1/5 - c, \alpha]$, x_t has an expectation of $\frac{(1/5 - c) + \alpha}{2} = \alpha - \gamma c/2$. This shows that X_- concentrates around $(\alpha - \gamma c/2) \cdot N_- = \alpha N_- - O(\gamma N_-)$ up to a typical deviation of $\sqrt{N_-} = \Theta(\sqrt{\gamma T})$. If we set $\gamma = \Theta(T^{-1/3})$ appropriately, the deviation would be $\Theta(\sqrt{\gamma T}) = \Theta(T^{1/3})$, dominating the $O(\gamma N_-) = O(\gamma^2 T) = O(T^{1/3})$ bias. This leads to an $\Omega(T^{1/3})$ penalty in both UCal and UCal^{sub} . In contrast, the dishonest forecasts (that take value either $2/5$ or $3/5$) incur an $e^{-\Omega(T)}$ penalty under either UCal or UCal^{sub} . This shows that neither UCal and UCal^{sub} can be truthful with sub-poly(T) parameters, even in the $\Omega(1)$ -smoothed setting.

3.2. Truthfulness of Subsampled Step Calibration

Showing that $\text{stepCE}^{\text{sub}}$ is truthful involves two steps: lower bounding the optimal penalty that can be achieved by a (possibly dishonest) forecaster, and upper bounding the penalty incurred by the truthful forecaster. The first part follows from a result of [HQYZ24]: regardless of the forecasting algorithm \mathcal{A} , it holds that

$$\mathbb{E}_{(x,p) \sim (\mathcal{D}, \mathcal{A}), y \sim \text{Unif}(\{0,1\}^T)} \left[\left| \sum_{t=1}^T y_t \cdot (x_t - p_t) \right| \right] \geq \Omega(\mathbb{E}[\gamma(\text{Var}_T)]), \quad (6)$$

where $\gamma(x) = \begin{cases} x, & x \leq 1, \\ \sqrt{x}, & x > 1 \end{cases}$ and the random variable Var_T is defined as $\text{Var}_T := \sum_{t=1}^T p_t^*(1 - p_t^*)$. Since the left-hand side of Equation (6) is a lower bound on $\text{stepCE}^{\text{sub}}(x, p)$, we also have $\text{OPT}_{\text{stepCE}^{\text{sub}}}(\mathcal{D}) \geq \Omega(\mathbb{E}[\gamma(\text{Var}_T)])$.

For simplicity, we assume in this section that Var_T is always $\Omega(T)$ (e.g., when every conditional probability is bounded away from 0 and 1), and focus on upper bounding $\text{err}_{\text{stepCE}^{\text{sub}}}(\mathcal{D}, \mathcal{A}^{\text{truthful}})$ by $\tilde{O}(\sqrt{T})$. The general case that Var_T can be much lower than T can be handled via a doubling trick similar to the technique of [HQYZ24]. Also, we will upper bound stepCE instead of $\text{stepCE}^{\text{sub}}$, as all the analyses would naturally generalize to the subsampled version.

Warm-up #1: Product distributions. We start with the special case that \mathcal{D} is a product distribution, i.e., $p^* \in [0, 1]^T$ is fixed, and each x_t is independently sampled from $\text{Bernoulli}(p_t^*)$. Without loss of generality, $p_1^* < p_2^* < \dots < p_T^*$. Then, the step calibration error can be written as:

$$\text{stepCE}(x, p^*) = \max_{t \in [T]} \left| \sum_{i=1}^t (x_i - p_i^*) \right|.$$

The above is simply the maximum deviation $\max_{t \in [T]} |X_t|$ of a random walk $(X_t)_{t=0}^T$ defined as $X_0 = 0$ and $X_t = X_{t-1} + (x_t - p_t^*)$. Naïvely controlling its expectation via Hoeffding's inequality and a union bound would give an $O(\sqrt{T \log T})$ upper bound. We can shave the logarithmic factor using Kolmogorov's inequality, which gives

$$\Pr \left[\max_{t \in [T]} |X_t| \geq \tau \right] \leq \frac{\mathbb{E}[X_T^2]}{\tau^2} \leq \frac{T}{4\tau^2}.$$

Integrating this tail bound would then give $\mathbb{E}[\max_{t \in [T]} |X_t|] = O(\sqrt{T})$.

Warm-up #2: Truthfulness up to an $O(\sqrt{\log(T/c)})$ factor. Unfortunately, the analysis above does not immediately generalize to non-product distributions, as the conditional probabilities of p_1^*, \dots, p_T^* are random and may not have a fixed ordering. One might resort to a “covering + union bound” argument, but the family of step functions does not admit a finite covering (in the ℓ_∞ sense).

Fortunately, in the smoothed setting in which each p_t^* is randomly drawn from a c -smoothed distribution, a simple discretization argument would suffice. Let $\varepsilon := c/T^2$ and consider an ε -net of the interval $[0, 1]$: $V_\varepsilon := \{0, \varepsilon, 2\varepsilon, \dots, 1\}$. We will relax $\text{stepCE}(x, p^*)$ into the following:

$$\max_{\alpha \in V_\varepsilon} \left| \sum_{t=1}^T (x_t - p_t^*) \cdot \mathbb{1}[p_t^* \in [0, \alpha]] \right|. \quad (7)$$

Suppose that the supremum over $\alpha \in [0, 1]$ in $\text{stepCE}(x, p^*)$ is achieved by some $\alpha^* \in [i\varepsilon, (i+1)\varepsilon]$. Then, the difference between the values of $\sum_{t=1}^T (x_t - p_t^*) \cdot \mathbb{1}[p_t^* \in [0, \alpha]]$ at α^* versus at $\alpha' = i\varepsilon$ is at most

$$\sum_{t=1}^T \mathbb{1}[p_t^* \in (\alpha', \alpha^*]] \leq \sum_{t=1}^T \mathbb{1}[p_t^* \in [i\varepsilon, (i+1)\varepsilon]],$$

the number of steps on which p_t^* falls into the interval $[i\varepsilon, (i+1)\varepsilon]$ of length ε . Since each p_t^* is drawn from a distribution with density at most $1/c$, $\Pr[p_t^* \in [i\varepsilon, (i+1)\varepsilon]] \leq \varepsilon/c = 1/T^2$. It then follows that the effect of replacing $[0, 1]$ with the ε -net V_ε is negligible.

It remains to upper bound the expectation of Equation (7). Since for each $\alpha \in V_\varepsilon$, $\sum_{t=1}^T (x_t - p_t^*) \cdot \mathbb{1}[p_t^* \in [0, \alpha]]$ is the outcome of a T -step martingale, applying Hoeffding's inequality with a union bound over V_ε gives an upper bound of $O(\sqrt{T \log |V_\varepsilon|}) = O(\sqrt{T \log(T/c)})$. Extending this to $\text{stepCE}^{\text{sub}}$ shows that $\text{stepCE}^{\text{sub}}$ is truthful up to an $O(\sqrt{\log(T/c)})$ factor.

Removing the $\text{polylog}(T)$ factor. We further tighten the multiplicative factor to $\sqrt{\log(1/c)}$ via a chaining argument. Our technique amounts to controlling the maximum deviation $\max_{t \in [T]} |X_t|$ in a martingale $(X_t)_{t=0}^T$ without applying Kolmogorov's inequality, and rather applies a more ‘‘combinatorial’’ analysis. This analysis turns out to be generalizable to non-product distributions.

As a warm-up, we revisit the toy problem below:

Max-deviation of random walk: Consider the random walk $(X_t)_{t=0}^T$ where $X_0 = 0$ and $X_t = X_{t-1} \pm 1$ with equal probability. Prove that $\mathbb{E} [\max_{t \in [T]} |X_t|] = O(\sqrt{T})$.

While above would follow from Kolmogorov's inequality and an integration, here is a different proof: We consider $\approx \log_2 T$ ‘‘levels’’ of random variables. The zeroth level consists of only $X_T - X_0$. The first level contains $X_T - X_{T/2}$ and $X_{T/2} - X_0$. The second level contains $X_T - X_{3T/4}$, $X_{3T/4} - X_{2T/4}$, \dots . In general, the i -th level divides the horizon into 2^i blocks of length $T/2^i$, and considers the displacement within each block. Then, we note that each X_t can be written as the sum of at most $\approx \log_2 T$ terms, at most one from each level. It follows that $\mathbb{E} [\max_{t \in [T]} |X_t|]$ is at most $\sum_{i=0}^{\log_2 T} Y_i$, where Y_i is the expectation of the maximum absolute value among level i . Since level i contains 2^i terms, each of which is a sum of $T/2^i$ independent samples from $\text{Unif}(\{\pm 1\})$, Hoeffding's inequality with a union bound gives $Y_i = O(\sqrt{(T/2^i) \log 2^i})$. Summing over all i shows $\mathbb{E} [\max_{t \in [T]} |X_t|] = O(\sqrt{T})$.

Here is how we apply the above to the analysis of the c -smoothed setting. We consider two discretization of $[0, 1]$: $V_c := \{0, c, 2c, \dots, 1\}$ and $V_\varepsilon := \{0, \varepsilon, 2\varepsilon, \dots, 1\}$ for $\varepsilon = c/T^2$. As argued earlier, replacing the interval $[0, 1]$ in stepCE with V_ε comes with a negligible error. If we could further replace V_ε with V_c , we would be done: controlling the maximum over V_c only involves a union bound over $|V_c| = O(1/c)$ martingales, and leads to an $O(\sqrt{T \log(1/c)})$ upper bound. Thus, it remains to control the error when we replace V_ε with V_c . We divide the interval $[0, 1]$ into sub-intervals of length $c/2, c/4, c/8, \dots, \varepsilon = c/T^2$. For each i between 1 and $O(\log T)$, the i -th level of the division consists of $2^i/c$ intervals of length $c/2^i$. For each level i and $j \in [2^i/c]$, we consider:

$$\sum_{t=1}^T (x_t - p_t^*) \cdot \mathbb{1}[(j-1) \cdot (c/2^i) \leq p_t^* \leq j \cdot (c/2^i)],$$

which is the outcome of a T -step martingale. Furthermore, since each p_t^* is sampled from a distribution with density $\leq 1/c$, it falls into the length- $(c/2^i)$ interval with probability at most 2^{-i} . Therefore, the contribution of the i -th level to the step calibration error is upper bounded by $\sqrt{(2^{-i}T) \cdot \log(2^i/c)}$. Summing over all i gives the desired upper bound of $O(\sqrt{T \log(1/c)})$.

3.3. Minimize Step Calibration in the Adversarial Setup

We sketch the proof of the Theorem 23 by giving a simple non-constructive argument for the $O(\sqrt{T \log T})$ error rate; we derive an explicit and efficient algorithm in the actual proof.

We apply the minimax argument of [Har22] for minimizing the ℓ_1 calibration error (also called the ECE) in an adversarial prediction setting. First, we restrict the forecaster so that its prediction is

always a multiple of $1/\sqrt{T}$. Then, we note that both the adversary and the forecaster have finitely many deterministic strategies—each deterministic strategy of the adversary (resp. forecaster) maps the history (all the previous outcomes and predictions) to the next outcome (resp. prediction). By the minimax theorem, it suffices to show that, against any given, possibly randomized strategy of the adversary, the forecaster can achieve an $O(\sqrt{T \log T})$ error with respect to stepCE.

In this scenario, at each step t , the forecaster can compute the conditional probability $p_t^* = \Pr [x_t = 1 \mid x_{1:(t-1)}, p_{1:(t-1)}]$ using the adversary’s strategy. Then, the forecaster predicts p_t obtained by rounding p_t^* to the nearest multiple of $1/\sqrt{T}$. To control the resulting step calibration error, we note that there are only $O(\sqrt{T})$ values of α that need to be considered (namely, the multiples of $1/\sqrt{T}$). For each fixed α , we can bound $\left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \alpha] \right|$ by the sum of two terms: one involving $(x_t - p_t^*)$ and another involving $(p_t^* - p_t)$. Since $|p_t^* - p_t| \leq 1/\sqrt{T}$ for all t , the latter term is always $O(\sqrt{T})$. For the former, we apply a union bound over the $O(\sqrt{T})$ values of α . The resulting bound would scale as $O(\sqrt{T \log T})$. While this argument does not give an explicit algorithm, we can also cast step calibration minimization as a Blackwell approachability problem [Bla56, Fos99] and obtain an explicit algorithm using min-max game dynamics [Har22, HJZ23].

Acknowledgments

This work is supported by the Google Research Fellowship and the National Science Foundation Graduate Research Fellowship Program under Grant No. DGE 2146752. The authors thank Kunhe Yang for valuable discussions on the project.

References

- [ACRS25] Eshwar Ram Arunachaleswaran, Natalie Collina, Aaron Roth, and Mirah Shi. An elementary predictor obtaining distance to calibration. In *Symposium on Discrete Algorithms (SODA)*, pages 1366–1370, 2025.
- [BCM14] Aditya Bhaskara, Moses Charikar, Ankur Moitra, and Aravindan Vijayaraghavan. Smoothed analysis of tensor decompositions. In *Symposium on Theory of Computing (STOC)*, pages 594–603, 2014.
- [BDGR22] Adam Block, Yuval Dagan, Noah Golowich, and Alexander Rakhlin. Smoothed online learning is as easy as statistical learning. In *Conference on Learning Theory (COLT)*, pages 1716–1786, 2022.
- [BDM20] Alon Brutzkus, Amit Daniely, and Eran Malach. ID3 learns juntas for smoothed product distributions. In *Conference on Learning Theory (COLT)*, pages 902–915, 2020.
- [BF⁺02] Henri Berestycki, Igor Florent, et al. Asymptotics and calibration of local volatility models. *Quantitative finance*, 2(1):61, 2002.
- [BGHN23] Jarosław Błasiok, Parikshit Gopalan, Lunjia Hu, and Preetum Nakkiran. A unifying theory of distance from calibration. In *Symposium on Theory of Computing (STOC)*, pages 1727–1740, 2023.

- [Bla56] David Blackwell. An analog of the minimax theorem for vector payoffs. *Pacific Journal of Mathematics*, 6(1):1 – 8, 1956. Publisher: Pacific Journal of Mathematics, A Non-profit Corporation.
- [BLQT21] Guy Blanc, Jane Lange, Mingda Qiao, and Li-Yang Tan. Decision tree heuristics can fail, even in the smoothed setting. In *International Conference on Randomization and Computation (RANDOM)*, pages 45:1–45:16, 2021.
- [BP23] Adam Block and Yury Polyanskiy. The sample complexity of approximate rejection sampling with applications to smoothed online learning. In *Conference on Learning Theory (COLT)*, pages 228–273, 2023.
- [Bri50] Glenn W. Brier. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78(1):1–3, 1950.
- [BS22] Adam Block and Max Simchowitz. Efficient and near-optimal smoothed online learning for generalized linear functions. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 7477–7489, 2022.
- [BSR23] Adam Block, Max Simchowitz, and Alexander Rakhlin. Oracle-efficient smoothed online learning for piecewise continuous decision making. In *Conference on Learning Theory (COLT)*, pages 1618–1665, 2023.
- [BST23] Adam Block, Max Simchowitz, and Russ Tedrake. Smoothed online learning for prediction in piecewise affine systems. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 41663–41674, 2023.
- [CAT16] Cynthia S Crowson, Elizabeth J Atkinson, and Terry M Therneau. Assessing calibration of prognostic risk scores. *Statistical methods in medical research*, 25(4):1692–1706, 2016.
- [Daw82] A. P. Dawid. The well-calibrated bayesian. *Journal of the American Statistical Association*, 77(379):605–610, 1982.
- [DDF⁺24] Yuval Dagan, Constantinos Daskalakis, Maxwell Fishelson, Noah Golowich, Robert Kleinberg, and Princewill Okoroafor. Breaking the $T^{2/3}$ barrier for sequential calibration. *arXiv preprint arXiv:2406.13668v3*, 2024.
- [DF83] Morris H DeGroot and Stephen E Fienberg. The comparison and evaluation of forecasters. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 32(1-2):12–22, 1983.
- [DKR⁺21] Cynthia Dwork, Michael P Kim, Omer Reingold, Guy N Rothblum, and Gal Yona. Outcome indistinguishability. In *Symposium on Theory of Computing (STOC)*, pages 1095–1108, 2021.
- [FH21] Dean P. Foster and Sergiu Hart. Forecast hedging and calibration. *Journal of Political Economy*, 129(12):3447–3490, 2021.

- [Fos99] Dean P. Foster. A proof of calibration via blackwell’s approachability theorem. *Games and Economic Behavior*, 29(1-2):73–78, 1999.
- [FV98] Dean P. Foster and Rakesh V. Vohra. Asymptotic calibration. *Biometrika*, 85(2):379–390, 1998.
- [Har22] Sergiu Hart. Calibrated forecasts: The minimax proof. *arXiv preprint arXiv:2209.05863*, 2022.
- [HHSY22] Nika Haghtalab, Yanjun Han, Abhishek Shetty, and Kunhe Yang. Oracle-efficient on-line learning for smoothed adversaries. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 4072–4084, 2022.
- [HJKRR18] Ursula Hébert-Johnson, Michael Kim, Omer Reingold, and Guy Rothblum. Multicalibration: Calibration for the (computationally-identifiable) masses. In *International Conference on Machine Learning (ICML)*, pages 1939–1948, 2018.
- [HJZ23] Nika Haghtalab, Michael Jordan, and Eric Zhao. A unifying perspective on multicalibration: Game dynamics for multi-objective learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 72464–72506, 2023.
- [HQYZ24] Nika Haghtalab, Mingda Qiao, Kunhe Yang, and Eric Zhao. Truthfulness of calibration measures. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 117237–117290, 2024.
- [HRS20] Nika Haghtalab, Tim Roughgarden, and Abhishek Shetty. Smoothed analysis of on-line and differentially private learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 9203–9215, 2020.
- [HRS24] Nika Haghtalab, Tim Roughgarden, and Abhishek Shetty. Smoothed analysis with adaptive adversaries. *Journal of the ACM (JACM)*, 71(3):1–34, 2024.
- [HW24] Lunjia Hu and Yifan Wu. Predict to minimize swap regret for all payoff-bounded tasks. In *Foundations of Computer Science (FOCS)*, pages 244–263, 2024.
- [JOKOM12] Xiaoqian Jiang, Melanie Osl, Jihoon Kim, and Lucila Ohno-Machado. Calibrating predictive model estimates to support personalized medicine. *Journal of the American Medical Informatics Association*, 19(2):263–274, 2012.
- [KF08] Sham M. Kakade and Dean P. Foster. Deterministic calibration and nash equilibrium. *Journal of Computer and System Sciences*, 74(1):115–130, 2008.
- [KPLST23] Robert Kleinberg, Renato Paes Leme, Jon Schneider, and Yifeng Teng. U-calibration: Forecasting for an unknown agent. In *Conference on Learning Theory (COLT)*, pages 5143–5145, 2023.
- [KSB21] Benjamin Kompa, Jasper Snoek, and Andrew L Beam. Second opinion needed: communicating uncertainty in medical machine learning. *NPJ Digital Medicine*, 4(1):4, 2021.

- [KST09] Adam Tauman Kalai, Alex Samorodnitsky, and Shang-Hua Teng. Learning and smoothed analysis. In *Foundations of Computer Science (FOCS)*, pages 395–404, 2009.
- [MW84] Allan H Murphy and Robert L Winkler. Probability forecasting in meteorology. *Journal of the American Statistical Association*, 79(387):489–500, 1984.
- [NRRX23] Georgy Noarov, Ramya Ramalingam, Aaron Roth, and Stephan Xie. High-dimensional unbiased prediction for sequential decision making. In *OPT 2023: Optimization for Machine Learning*, 2023.
- [OKK25] Princewill Okoroafor, Robert Kleinberg, and Michael P. Kim. Near-optimal algorithms for omniprediction, 2025.
- [QV21] Mingda Qiao and Gregory Valiant. Stronger calibration lower bounds via sidestepping. In *Symposium on Theory of Computing (STOC)*, pages 456–466, 2021.
- [QZ24] Mingda Qiao and Letian Zheng. On the distance from calibration in sequential prediction. In *Conference on Learning Theory (COLT)*, pages 4307–4357, 2024.
- [RS24] Aaron Roth and Mirah Shi. Forecasting for swap regret for all downstream agents. In *Economics and Computation (EC)*, pages 466–488, 2024.
- [RSB⁺25] Raphael Rossellini, Jake A. Soloff, Rina Foygel Barber, Zhimei Ren, and Rebecca Willett. Can a calibration metric be both testable and actionable?, 2025.
- [RST11] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning: Stochastic, constrained, and smoothed adversaries. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1764–1772, 2011.
- [ST04] Daniel A Spielman and Shang-Hua Teng. Smoothed analysis of algorithms: Why the simplex algorithm usually takes polynomial time. *Journal of the ACM (JACM)*, 51(3):385–463, 2004.
- [ST09] Daniel A Spielman and Shang-Hua Teng. Smoothed analysis: an attempt to explain the behavior of algorithms in practice. *Communications of the ACM*, 52(10):76–84, 2009.
- [VCV15] Ben Van Calster and Andrew J Vickers. Calibration of risk prediction models: impact on decision-analytic performance. *Medical decision making*, 35(2):162–169, 2015.
- [WM68] Robert L. Winkler and Allan H. Murphy. “Good” probability assessors. *Journal of Applied Meteorology and Climatology*, 7(5):751–758, 1968.

Appendix A. Non-truthfulness of U-Calibration

Previous observations of non-truthfulness in calibration measures centered around a specific source of non-truthfulness: an incentive that calibration measures provide to forecasters to “cancel out” previous errors in their forecast by intentionally mispredicting future events. This form of non-truthfulness can be remedied by randomly subsampling which timesteps are included in the calibration measure computation [HQYZ24].

In this section, we identify two new and qualitatively distinct sources of non-truthfulness in the U-Calibration measure, neither of which can be remedied with subsampling alone. The first source of non-truthfulness arises from the inherently discontinuous nature of the U-Calibration measure and is largely inevitable: one can show that any reasonable calibration measure that provides decision-theoretic guarantees must also be non-truthful due to discontinuity. However, this form of non-truthfulness requires an adversary to precisely choose the event distribution \mathcal{D} and, as we will later see, largely disappears under smoothed analysis. The second source of non-truthfulness arises from the incentive that U-Calibration provides to a forecaster to *hedge* their predictions, i.e., to exaggerate their uncertainty in their predictions. This form of non-truthfulness remains even in the smoothed setting.

A.1. Non-truthfulness from Discontinuity

The U-Calibration measure is discontinuous in a forecaster’s prediction. From a decision-theoretic perspective, this is because an agent’s mapping from the forecaster’s prediction to an action is usually discontinuous: a marginal change in the probability of an event occurring may result in an agent switching actions and perhaps incurring significantly higher or lower regret. The following proposition describes one such case.

Proposition 5 *For any $\varepsilon \in (0, 1/2)$ and even number $T \in \mathbb{Z}$, there is a sequence of events $x_{1:T}$ and predictions $p_{1:T}$ such that changing $p_{1:T}$ to a similar alternative set of predictions $\tilde{p}_{1:T}$ where $\|p_{1:T} - \tilde{p}_{1:T}\|_\infty \leq \varepsilon$ increases the U-Calibration measure from $\text{UCal}(x_{1:T}, p_{1:T}) = 0$ to $\text{UCal}(x_{1:T}, \tilde{p}_{1:T}) = \Omega(T)$.*

Our proof of this proposition also shows that *no* calibration measure can be complete, continuous, and decision-theoretic simultaneously.

Proof We define the events $x_{1:T}$ and original predictions $p_{1:T}$ as $(x_t, p_t) = (1, \frac{1}{2} - \varepsilon)$ for the first half of timesteps $t \in [T/2]$ and $(x_t, p_t) = (0, \frac{1}{2} + \varepsilon)$ for the second half of timesteps $t > T/2$. Recall the equivalent form of the V-Calibration error from Proposition 3. For $\alpha = 1/2$, there are $N_-^{(\alpha)} = T/2$ steps on which $p_t < \alpha$, and the events on those steps sum up to $X_-^{(\alpha)} = T/2$. It follows that

$$\text{VCal}(x, p) \geq 2 \left(X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)} \right) = \Omega(T).$$

This further implies $\text{UCal}(x, p) = \Omega(T)$ due to the equivalence of UCal and VCal (Lemma 2). On the other hand, the alternative predictions $\tilde{p}_{1:T} = \frac{1}{2} \cdot \mathbf{1}_T$ would guarantee $\text{UCal}(x, p) = 0$ and $\|p - \tilde{p}\|_\infty \leq \varepsilon$. ■

The discontinuity of the U-Calibration measure provides a source of non-truthfulness that cannot be avoided with the subsampling technique of [HQYZ24]. The following proposition demonstrates

an $O(\sqrt{T})$ - $\Omega(T)$ truthfulness gap for U-Calibration, as well as for its subsampled variant UCal^{sub} :

$$\text{UCal}^{\text{sub}}(x_{1:T}, p_{1:T}) = \mathbb{E}_{S \sim \text{Unif}(2^{[T]})} [\text{UCal}(x|_S, p|_S)]. \quad (8)$$

Proposition 6 *Both the U-Calibration measure UCal and its subsampled variant UCal^{sub} suffer from an $O(\sqrt{T})$ - $\Omega(T)$ truthfulness gap.*

Proof We will analyze the truthfulness gaps of V-Calibration VCal and its subsampled version VCal^{sub} ; the proposition would then follow from Lemma 2.

The distribution of events. Let $\varepsilon \in (0, 1/4)$ be sufficiently small such that $\varepsilon = O(1/\sqrt{T})$. We now construct a distribution \mathcal{D} such that the conditional probability p_t^* at every timestep t is guaranteed to fall into the interval $[\frac{1}{2} - \varepsilon, \frac{1}{2} + \varepsilon]$. Let us define \mathcal{D} as fixing $p_1^* = \frac{1}{2}$ and, for all $t \in [T]$, letting

$$p_{t+1}^* = \begin{cases} p_t^* + \frac{\varepsilon}{2^t} & \text{if } x_t = 1, \\ p_t^* - \frac{\varepsilon}{2^t} & \text{if } x_t = 0. \end{cases}$$

That is, distribution \mathcal{D} adversarially sets the probabilities p_t^* in a binary search fashion based on the realizations of the preceding events $x_{1:t-1}$.

Truthful forecasts give a linear penalty. By our adversarial construction of \mathcal{D} , regardless of the realization of (x, p^*) , the following holds for $\alpha^* := p_{T+1}^*$ and every $t \in [T]$: (1) $p_t^* \neq \alpha^*$; (2) $p_t^* < \alpha^*$ implies $x_t = 1$; and (3) $p_t^* > \alpha^*$ implies $x_t = 0$.

Towards lower bounding $\text{VCal}(x, p^*)$, we consider the equivalent formulation of the V-Calibration error in Proposition 3 at $\alpha = \alpha^*$. Our construction guarantees $X_-^{(\alpha)} = N_-^{(\alpha)}$ and $X_+^{(\alpha)} = 0$. Since $\alpha \in [1/2 - \varepsilon, 1/2 + \varepsilon] \subseteq [1/4, 3/4]$, we have

$$X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)} = (1 - \alpha) \cdot N_-^{(\alpha)} \geq N_-^{(\alpha)} / 4$$

and

$$\alpha \cdot N_+^{(\alpha)} - X_+^{(\alpha)} = \alpha \cdot N_+^{(\alpha)} \geq N_+^{(\alpha)} / 4.$$

It follows that

$$\max \left\{ X_-^{(\alpha)} - \alpha N_-^{(\alpha)}, \alpha N_+^{(\alpha)} - X_+^{(\alpha)} \right\} \geq \max \left\{ N_-^{(\alpha)}, N_+^{(\alpha)} \right\} / 4 \geq T/8.$$

By Proposition 3, $\text{VCal}(x, p^*) \geq \Omega(T)$ always holds, which in turn gives $\text{err}_{\text{VCal}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D})) = \Omega(T)$.

To lower bound $\text{err}_{\text{VCal}^{\text{sub}}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D}))$, we note that the same argument as above gives

$$\text{VCal}(x|_S, p^*|_S) \geq \Omega(|S|)$$

for every $S \subseteq [T]$. Taking an expectation over $S \sim \text{Unif}(2^{[T]})$ gives $\text{err}_{\text{VCal}^{\text{sub}}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D})) = \Omega(T)$.

Dishonest forecasts with an $O(\sqrt{T})$ penalty. On the other hand, constantly predicting $1/2$ gives an $O(\sqrt{T})$ error with respect to both VCal and VCal^{sub} . To see this, we apply [Proposition 3](#):

$$\text{VCal}(x, \tfrac{1}{2} \cdot \mathbf{1}_T) = 2 \cdot \sup_{\alpha \in [0,1]} \max \left\{ X_-^{(\alpha)} - \alpha N_-^{(\alpha)}, \alpha N_+^{(\alpha)} - X_+^{(\alpha)} \right\}.$$

Since all predictions take value $1/2$, we have

$$X_-^{(\alpha)} - \alpha N_-^{(\alpha)} = \mathbb{1}[\alpha > 1/2] \cdot \left(\sum_{t=1}^T x_t - \alpha T \right) \leq \left| \sum_{t=1}^T x_t - \frac{T}{2} \right|$$

and

$$\alpha N_+^{(\alpha)} - X_+^{(\alpha)} = \mathbb{1}[\alpha < 1/2] \cdot \left(\alpha T - \sum_{t=1}^T x_t \right) \leq \left| \sum_{t=1}^T x_t - \frac{T}{2} \right|.$$

Therefore, we have

$$\text{VCal}(x, \tfrac{1}{2} \cdot \mathbf{1}_T) \leq 2 \left| \sum_{t=1}^T (x_t - 1/2) \right| \leq 2 \left| \sum_{t=1}^T (x_t - p_t^*) \right| + 2 \left| \sum_{t=1}^T (p_t^* - 1/2) \right|.$$

Since $p_t^* \in [1/2 - \varepsilon, 1/2 + \varepsilon]$ holds for every $t \in [T]$, the second term above is always at most $2\varepsilon T = O(\sqrt{T})$. The first term is the deviation of a T -step martingale with bounded differences, and is thus bounded by $O(\sqrt{T})$ in expectation. Therefore, we have

$$\text{OPT}_{\text{VCal}}(\mathcal{D}) \leq \mathbb{E}_{x \sim \mathcal{D}} [\text{VCal}(x, \tfrac{1}{2} \cdot \mathbf{1}_T)] \leq O(\sqrt{T}).$$

Again, the argument above can be easily extended to show that

$$\mathbb{E}_{x \sim \mathcal{D}} [\text{VCal}(x|_S, \tfrac{1}{2} \cdot \mathbf{1}_{|S|})] \leq O(\sqrt{|S|}) \leq O(\sqrt{T})$$

holds for every fixed $S \subseteq [T]$. Taking an expectation over $S \sim \text{Unif}(2^{[T]})$ gives

$$\text{OPT}_{\text{VCal}^{\text{sub}}}(\mathcal{D}) \leq \mathbb{E}_{x \sim \mathcal{D}, S \sim \text{Unif}(2^{[T]})} [\text{VCal}(x|_S, \tfrac{1}{2} \cdot \mathbf{1}_{|S|})] \leq O(\sqrt{T}).$$

This establishes the $O(\sqrt{T})$ - $\Omega(T)$ truthfulness gaps for VCal and VCal^{sub} and finishes the proof. ■

We can generalize the proof of [Proposition 6](#) beyond the U-Calibration measure to *any* decision-theoretic calibration measure that satisfies the very weak condition of completeness ([Definition 1](#)). Recall that a calibration measure is decision-theoretic if it upper bounds the external regret of an agent that acts on the forecaster's predictions (or equivalently upper bounds U-Calibration) and a calibration measure is complete if it is low for a base-rate forecaster when all events occur with the same constant probability. This generalization, stated formally in following proposition, implies that—without smoothed analysis—the requirement of a calibration measure providing decision-theoretic guarantees is directly at odds with that of being truthful.

Proposition 7 Consider any complete decision-theoretic calibration measure CM_T . Suppose that its completeness guarantee for $\alpha = \frac{1}{2}$ takes the rate of f , i.e.,

$$\mathbb{E}_{x_1, \dots, x_T \sim \text{Bernoulli}(1/2)} [\text{CM}_T(x, \frac{1}{2} \cdot \mathbf{1}_T)] = O(f(T)).$$

Then, CM_T has a truthfulness gap of $O(f(T)) \cdot \Omega(T)$.

Proof As in the proof of Proposition 6, we will define the distribution \mathcal{D} by setting $p_1^* = \frac{1}{2}$ and

$$p_{t+1}^* = \begin{cases} p_t^* + \frac{\varepsilon}{2^t} & \text{if } x_t = 1, \\ p_t^* - \frac{\varepsilon}{2^t} & \text{if } x_t = 0 \end{cases}$$

for some small $\varepsilon = O(f(T)/T^2)$. We therefore have from Proposition 6 that the truthful forecaster's U-Calibration measure is lower bounded by $\mathbb{E}_{x \sim \mathcal{D}} [\text{UCal}(x, p^*)] \geq \Omega(T)$. Because CM_T is a decision-theoretic calibration measure, it must upper bound the U-Calibration measure, meaning that

$$\text{err}_{\text{CM}_T}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D})) = \mathbb{E}_{x \sim \mathcal{D}} [\text{CM}_T(x, p^*)] \geq \Omega \left(\mathbb{E}_{x \sim \mathcal{D}} [\text{UCal}(x, p^*)] \right) \geq \Omega(T).$$

To upper bound the calibration measure for the non-truthful forecaster, we first observe that the construction of \mathcal{D} guarantees $|p_t^* - 1/2| \leq \varepsilon$ for all $t \in [T]$. Thus, the total variation distance between \mathcal{D} and $\text{Unif}(\{0, 1\}^T)$ is $O(\varepsilon T)$. Since $\text{CM}_T(\cdot, \cdot)$ takes value in $[0, T]$, by our choice of $\varepsilon = O(f(T)/T^2)$ and the completeness of CM, we have

$$\mathbb{E}_{x \sim \mathcal{D}} [\text{CM}_T(x, 1/2 \cdot \mathbf{1}_T)] \leq \mathbb{E}_{x \sim \text{Unif}(\{0, 1\}^T)} [\text{CM}_T(x, 1/2 \cdot \mathbf{1}_T)] + T \cdot O(\varepsilon T) = O(f(T)).$$

■

However, the existence of this conflict—and the proof of Proposition 7—hinges on the assumption that an adversary has the ability to choose the event distribution \mathcal{D} , in particular each conditional probability p_t^* , in an arbitrarily precise way so as to exploit discontinuity. As we will later see, under smoothed analysis where an adversary has limited precision in choosing p_t^* , this inevitability of non-truthfulness for decision-theoretic measures largely disappears.

A.2. Non-truthfulness from Hedging

We now demonstrate a second source of non-truthfulness in the U-Calibration measure that arises from the incentivization of *hedging*: forecasters can reduce their expected U-Calibration measure by portraying their beliefs as being more uncertain than they truly are. This form of non-truthfulness does not depend on the existence of an adversary that is able to choose the distribution \mathcal{D} , or equivalently $p_{1:T}^*$, with high precision. It also cannot be remedied with the technique of randomly sub-sampling timesteps.

One might expect that hedging, by introducing bias to a forecaster's predictions, should increase the U-Calibration measure. However, we can construct a setting, stated formally in the following proposition, where hedging results in an $\Omega(T)$ bias but a U-Calibration measure of 0. Moreover, we can design this example to be asymmetric where the forecaster only hedges its predictions in one direction, e.g., if it believes the event will likely occur.

Proposition 8 *There is a sequence of events $x_{1:T}$ and predictions $p_{1:T}$ with zero U-Calibration measure $\text{UCal}(x, p) = 0$ but linear bias $\left| \sum_{t=1}^T (x_t - p_t) \right| = \Omega(T)$.*

Proof We will construct a simple example for $T = 2$, which can be repeated an arbitrary number of times to attain the claim for arbitrary T . Consider the events $x = (0, 1)$ and forecasts $p = (0, 3/4)$. To show that $\text{UCal}(x, p) = 0$, by Lemma 2 and proposition 3, it suffices to prove that the following holds for every $\alpha \in [0, 1]$:

$$\max \left\{ X_-^{(\alpha)} - \alpha N_-^{(\alpha)}, \alpha N_+^{(\alpha)} - X_+^{(\alpha)} \right\} \leq 0,$$

where $N_-^{(\alpha)}$ (resp., $N_+^{(\alpha)}$) denotes the number of timesteps on which the prediction is strictly below (resp., above) α , and $X_-^{(\alpha)}$ (resp., $X_+^{(\alpha)}$) denotes the sum of the events on those steps. This can be done via the following cases analysis:

- When $\alpha = 0$, we have $N_- = X_- = 0$ and $N_+ = X_+ = 1$. This gives $X_- - \alpha N_- = 0$ and $\alpha N_+ - X_+ = -1 < 0$.
- When $\alpha \in (0, 3/4)$, we have $N_- = 1$, $X_- = 0$ and $N_+ = X_+ = 1$. This gives $X_- - \alpha N_- = -\alpha < 0$ and $\alpha N_+ - X_+ = \alpha - 1 < 0$.
- When $\alpha = 3/4$, we have $N_- = 1$, $X_- = 0$, and $N_+ = X_+ = 0$. This gives $X_- - \alpha N_- = -\alpha < 0$ and $\alpha N_+ - X_+ = 0$.
- When $\alpha \in (3/4, 1]$, we have $N_- = 2$, $X_- = 1$, and $N_+ = X_+ = 0$. This gives $X_- - \alpha N_- = 1 - 2\alpha < 0$ and $\alpha N_+ - X_+ = 0$.

Therefore, we have $\text{UCal}(x, p) = 0$. On the other hand, the total bias of the predictions is $|(x_1 + x_2) - (p_1 + p_2)| = 1/4$. Repeated $T/2$ times, this gives a bias of $T/8 = \Omega(T)$ and a U-Calibration measure of 0. \blacksquare

This example demonstrates that the bias of hedging predictions does not negatively affect the U-Calibration measure. On the other hand, hedging often provides an explicit advantage. For example, hedging an honest prediction of $p_t^* = 1$ down to $p_t = 3/4$ may not incur any cost, but may instead be of benefit for a previous or future timestep t' where the event outcome is high variance with $p_{t'}^* = 3/4$ and the forecaster seeks to “dilute” the variance. In this way, we can construct an $e^{-\Omega(T)}$ - $\Omega(\sqrt{T})$ truthfulness gap for the U-Calibration measure and its subsampled variant. Later, we will extend this construction to the smoothed setting.

Proposition 9 *Both the U-Calibration measure UCal and its subsampled version UCal^{sub} have an $e^{-\Omega(T)}$ - $\Omega(\sqrt{T})$ truthfulness gap.*

Proof Again, we analyze V-Calibration rather than U-Calibration, i.e., we will establish the truthfulness gap of VCal and VCal^{sub} , and the proposition would then follow from Lemma 2.

Let T be an even number. Let $p^* = (\frac{1}{5}, \frac{1}{5}, \dots, \frac{1}{5}, \frac{4}{5}, \frac{4}{5}, \dots, \frac{4}{5})$ be the vector with $\frac{T}{2}$ copies of $\frac{1}{5}$ and $\frac{4}{5}$ each and \mathcal{D} be the product distribution $\prod_{t=1}^T \text{Bernoulli}(p_t^*)$. Similarly, let p be the alternative “non-truthful” prediction $(\frac{2}{5}, \frac{2}{5}, \dots, \frac{2}{5}, \frac{3}{5}, \frac{3}{5}, \dots, \frac{3}{5})$ where again both $\frac{2}{5}$ and $\frac{3}{5}$ appear exactly $\frac{T}{2}$ times. We will show that, under both VCal and VCal^{sub} , the truthful forecaster $\mathcal{A}^{\text{truthful}}(\mathcal{D})$ that predicts according to p^* receives an $\Omega(\sqrt{T})$ penalty, whereas predicting according to p leads to an $e^{-\Omega(T)}$ penalty.

Truthful forecasts lead to an $\Omega(\sqrt{T})$ penalty. We start by showing that $\text{err}_{\text{VCal}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D}))$ and $\text{err}_{\text{VCal}^{\text{sub}}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D}))$ are both lower bounded by $\Omega(\sqrt{T})$. Towards applying [Proposition 3](#), we fix $\alpha = 1/5 + \varepsilon$ for an arbitrarily small $\varepsilon > 0$, and aim to lower bound the quantity

$$X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)} = X_-^{(\alpha)} - (1/5 + \varepsilon) \cdot \frac{T}{2},$$

where $N_-^{(\alpha)} = T/2$ is the number of timesteps on which the prediction is strictly smaller than $\alpha = 1/5 + \varepsilon$, and $X_-^{(\alpha)} = \sum_{t=1}^{T/2} x_t$ is the sum of the $T/2$ events on those steps.

By definition of \mathcal{D} , $X_-^{(\alpha)}$ follows the distribution $\text{Binomial}(T/2, 1/5)$. Then, it holds with probability $\Omega(1)$ that $X_-^{(\alpha)} \geq (T/2) \cdot (1/5) + \Omega(\sqrt{T}) = T/10 + \Omega(\sqrt{T})$, i.e., $X_-^{(\alpha)}$ exceeds its mean by $\Omega(\sqrt{T})$. Conditioning on this event, by [Proposition 3](#), we have

$$\text{VCal}(x, p^*) \geq X_-^{(\alpha)} - (1/5 + \varepsilon) \cdot \frac{T}{2} \geq \Omega(\sqrt{T}) - \varepsilon T.$$

This implies that

$$\text{err}_{\text{VCal}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D})) = \mathbb{E}_{x \sim \mathcal{D}} [\text{VCal}(x, p^*)] \geq \Omega(1) \cdot [\Omega(\sqrt{T}) - \varepsilon T].$$

Taking $\varepsilon \rightarrow 0^+$ proves $\text{err}_{\text{VCal}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D})) = \Omega(\sqrt{T})$.

For VCal^{sub} , we note that the argument above shows that, for any fixed $S \subseteq [T]$,

$$\mathbb{E}_{x \sim \mathcal{D}} [\text{VCal}(x|_S, p^*|_S)] \geq \Omega\left(\sqrt{|S \cap [T/2]|}\right).$$

Then, over the randomness in $S \sim \text{Unif}(2^{[T]})$, $|S \cap [T/2]|$ follows $\text{Binomial}(T/2, 1/2)$. It follows that $\mathbb{E}_{S \sim \text{Unif}(2^{[T]})} \left[\sqrt{|S \cap [T/2]|} \right] = \Omega(\sqrt{T})$, and

$$\text{err}_{\text{VCal}^{\text{sub}}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D})) = \mathbb{E}_{x \sim \mathcal{D}, S \subseteq [T]} [\text{VCal}(x|_S, p^*|_S)] = \Omega(\sqrt{T}).$$

Dishonest forecasts lead to an exponentially small penalty. Suppose that the forecaster strategically predicts according to $p = (\frac{2}{5}, \frac{2}{5}, \dots, \frac{2}{5}, \frac{3}{5}, \frac{3}{5}, \dots, \frac{3}{5})$ instead. Again, we start with the analysis for VCal . Let $X_1 := \sum_{t=1}^{T/2} x_t$ and $X_2 := \sum_{t=T/2+1}^T x_t$ denote the sums of the first and second halves of the event sequence, respectively. Note that over the randomness in $x \sim \mathcal{D}$, X_1 follows $\text{Binomial}(T/2, 1/5)$ and X_2 follows $\text{Binomial}(T/2, 4/5)$. By a Chernoff bound, with probability at least $1 - e^{-\Omega(T)}$, the following three inequalities hold simultaneously:

$$X_1 \leq \frac{T}{5}, \quad X_2 \geq \frac{3T}{10}, \quad \frac{2T}{5} \leq X_1 + X_2 \leq \frac{3T}{5}. \quad (9)$$

It remains to show that the inequalities above together imply $\text{VCal}(x, p) = 0$. Assuming this, since the V-Calibration error is at most T , we would then have $\mathbb{E}_{x \sim \mathcal{D}} [\text{VCal}(x, p)] \leq e^{-\Omega(T)} \cdot T = e^{-\Omega(T)}$ as desired. By [Proposition 3](#), it suffices to show that the inequalities in (9) together imply that, for every $\alpha \in [0, 1]$,

$$\max \left\{ X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)}, \alpha \cdot N_+^{(\alpha)} - X_+^{(\alpha)} \right\} \leq 0.$$

This can be done via the following case analysis:

- When $\alpha \in [0, 2/5)$, we have $N_-^{(\alpha)} = X_-^{(\alpha)} = 0$, $N_+^{(\alpha)} = T$ and $X_+^{(\alpha)} = X_1 + X_2$. It follows that

$$X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)} = 0 \quad \text{and} \quad \alpha \cdot N_+^{(\alpha)} - X_+^{(\alpha)} = \alpha T - (X_1 + X_2) \leq \frac{2}{5}T - \frac{2T}{5} = 0.$$

- When $\alpha = 2/5$, we have $N_-^{(\alpha)} = X_-^{(\alpha)} = 0$, $N_+^{(\alpha)} = T/2$ and $X_+^{(\alpha)} = X_2$. It follows that

$$X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)} = 0 \quad \text{and} \quad \alpha \cdot N_+^{(\alpha)} - X_+^{(\alpha)} = \alpha T/2 - X_2 \leq \frac{T}{5} - \frac{3T}{10} < 0.$$

- When $\alpha \in (2/5, 3/5)$, we have $N_-^{(\alpha)} = N_+^{(\alpha)} = T/2$, $X_-^{(\alpha)} = X_1$, and $X_+^{(\alpha)} = X_2$. It follows that

$$X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)} = X_1 - \alpha \cdot \frac{T}{2} \leq \frac{T}{5} - \frac{2}{5} \cdot \frac{T}{2} = 0$$

and

$$\alpha \cdot N_+^{(\alpha)} - X_+^{(\alpha)} = \alpha T/2 - X_2 \leq \frac{3}{5} \cdot \frac{T}{2} - \frac{3T}{10} = 0.$$

- The remaining cases that $\alpha = 3/5$ and $\alpha \in (3/5, 1]$ hold by symmetry.

Therefore, we conclude that

$$\text{OPT}_{\text{VCal}}(\mathcal{D}) \leq \mathbb{E}_{x \sim \mathcal{D}} [\text{VCal}(x, p)] \leq e^{-\Omega(T)}.$$

Towards upper bounding $\text{OPT}_{\text{VCal}^{\text{sub}}}(\mathcal{D})$, we analyze the quantity

$$\mathbb{E}_{x \sim \mathcal{D}} [\text{VCal}^{\text{sub}}(x, p)] = \mathbb{E}_{x \sim \mathcal{D}, S \sim \text{Unif}(2^{[T]})} [\text{VCal}(x|_S, p|_S)].$$

As in the analysis for VCal, we will identify a high-probability event (over the randomness in both x and S) that: (1) happens with probability $1 - e^{-\Omega(T)}$; (2) implies $\text{VCal}(x|_S, p|_S) = 0$. It would then immediately follow that $\mathbb{E}_{x \sim \mathcal{D}} [\text{VCal}^{\text{sub}}(x, p)] = e^{-\Omega(T)}$.

Let $N_1 := |S \cap \{1, 2, \dots, T/2\}|$ (resp., $N_2 := |S \cap \{T/2 + 1, T/2 + 2, \dots, T\}|$) denote the number of timesteps among the first half (resp., the second half) of the sequence that get subsampled into S . Let $X_1 := \sum_{t \in S} x_t \cdot \mathbb{1}[t \leq T/2]$ and $X_2 := \sum_{t \in S} x_t \cdot \mathbb{1}[t > T/2]$ denote the sum of the events on those steps, respectively.

Note that each of N_1 , N_2 , X_1 and X_2 can be written as a sum of $\Omega(T)$ independent Bernoulli random variables. Furthermore, over the randomness in both x and S , we have $\mathbb{E}[N_1] = \mathbb{E}[N_2] = T/4$, $\mathbb{E}[X_1] = T/20$ and $\mathbb{E}[X_2] = T/5$. Let $\varepsilon := 1/40$ be a small constant. By a Chernoff bound and the union bound, the following conditions hold simultaneously with probability $1 - e^{-\Omega(T)}$:

- $N_1, N_2 \in [T/4 - \varepsilon T, T/4 + \varepsilon T]$.
- $X_1 \leq T/20 + \varepsilon T$ and $X_2 \geq T/5 - \varepsilon T$.
- $X_1 + X_2 \in [T/4 - \varepsilon T, T/4 + \varepsilon T]$.

It remains to show that, assuming the conditions above, we have $\text{VCal}(x|_S, p|_S) = 0$. Towards applying [Proposition 3](#), we analyze the quantity

$$\max \left\{ X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)}, \alpha \cdot N_+^{(\alpha)} - X_+^{(\alpha)} \right\}$$

for $\alpha \in [0, 1]$ with respect to events $x|_S$ and predictions $p|_S$:

- When $\alpha \in [0, 2/5)$, we have

$$N_-^{(\alpha)} = 0, X_-^{(\alpha)} = 0, N_+^{(\alpha)} = N_1 + N_2, X_+^{(\alpha)} = X_1 + X_2.$$

It follows that

$$X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)} = 0$$

and

$$\alpha \cdot N_+^{(\alpha)} - X_+^{(\alpha)} = \alpha \cdot (N_1 + N_2) - (X_1 + X_2) \leq \frac{2}{5} \cdot \left(\frac{T}{2} + 2\varepsilon T \right) - \left(\frac{T}{4} - \varepsilon T \right) = -\frac{T}{20} + \frac{9}{5}\varepsilon T < 0.$$

- When $\alpha = 2/5$, we have

$$N_-^{(\alpha)} = 0, X_-^{(\alpha)} = 0, N_+^{(\alpha)} = N_2, X_+^{(\alpha)} = X_2.$$

It follows that

$$X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)} = 0$$

and

$$\alpha \cdot N_+^{(\alpha)} - X_+^{(\alpha)} = \frac{2}{5} \cdot N_2 - X_2 \leq \frac{2}{5} \cdot \left(\frac{T}{4} + \varepsilon T \right) - \left(\frac{T}{5} - \varepsilon T \right) = -\frac{T}{10} + \frac{7}{5}\varepsilon T < 0.$$

- When $\alpha \in (2/5, 3/5)$, we have

$$N_-^{(\alpha)} = N_1, X_-^{(\alpha)} = X_1, N_+^{(\alpha)} = N_2, X_+^{(\alpha)} = X_2.$$

It follows that

$$X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)} = X_1 - \alpha \cdot N_1 \leq \frac{T}{20} + \varepsilon T - \frac{2}{5} \cdot \left(\frac{T}{4} - \varepsilon T \right) = -\frac{T}{20} + \frac{7}{5}\varepsilon T < 0.$$

and

$$\alpha \cdot N_+^{(\alpha)} - X_+^{(\alpha)} = \alpha \cdot N_2 - X_2 \leq \frac{3}{5} \cdot \left(\frac{T}{4} + \varepsilon T \right) - \left(\frac{T}{5} - \varepsilon T \right) = -\frac{T}{20} + \frac{8}{5}\varepsilon T < 0.$$

- Finally, by symmetry, we have an upper bound of 0 in the cases that $\alpha = 3/5$ and $\alpha \in (3/5, 1]$.

Therefore, we conclude that

$$\text{OPT}_{\text{VCal}^{\text{sub}}}(\mathcal{D}) \leq \mathbb{E}_{x \sim \mathcal{D}} [\text{VCal}^{\text{sub}}(x, p)] \leq e^{-\Omega(T)}.$$

■

By adding noises to the marginal probabilities, we can show that an $e^{-\Omega(T)}$ -poly(T) truthfulness gap remains even if we assume a smoothed setting where an adversary cannot precisely choose the conditional probabilities $p_{1:T}^*$. Recall that, in the c -smoothed setting, the adversary is forced to sample each conditional probability p_t^* from a distribution with density upper bounded by $1/c$.

Proposition 10 *Both the U-Calibration measure UCal and its subsampled version UCal^{sub} have an $e^{-\Omega(T)}$ - $\Omega(T^{1/3})$ truthfulness gap, even when the event distribution is a product distribution with marginal probabilities perturbed by independent noises uniformly sampled from $[-c, c]$ for some constant $c \in (0, 1/5)$.*

Intuitively, the constructions of both [Propositions 9](#) and [10](#) involve incentivizing forecasters to exaggerate the uncertainty of their forecasts. Hedging one's forecasts allows any and all downstream agents to avoid paying for the variance of the events and attain zero regret with high probability. Notably, this form of non-truthfulness that arises when studying decision-theoretic calibration measures is qualitatively different from the non-truthfulness shown by [\[HQYZ24\]](#), which is more concerned with incentivizing forecasters to dishonestly “patch up” previous erroneous forecasts than with incentivizing forecasters to proactively hedge for (potentially future) uncertainty.

Proof Let T be an even number. We construct a c -smoothed prior \mathcal{P} by sampling each event probability p_t^* in the first half of timesteps (i.e., $t \leq T/2$) from the uniform distribution over $[\frac{1}{5} - c, \frac{1}{5} + c]$, and similarly sampling each p_t^* independently and uniformly from $[\frac{4}{5} - c, \frac{4}{5} + c]$ for $t \geq T/2 + 1$. More formally, for every $t \in [T]$ and $b_{1:(t-1)} \in \{0, 1\}^{t-1}$, we have

$$\mathcal{P}(b_{1:(t-1)}) = \begin{cases} \text{Unif}([\frac{1}{5} - c, \frac{1}{5} + c]), & t \leq T/2, \\ \text{Unif}([\frac{4}{5} - c, \frac{4}{5} + c]), & t \geq T/2 + 1. \end{cases}$$

Dishonest forecasts with low penalties. Let $p = (\frac{2}{5}, \frac{2}{5}, \dots, \frac{2}{5}, \frac{3}{5}, \frac{3}{5}, \dots, \frac{3}{5})$ be the vector with $T/2$ copies of $\frac{2}{5}$ and $\frac{3}{5}$ and represent the alternative “non-truthful” prediction. The analysis of the non-truthful forecaster remains largely the same as in the proof of [Proposition 9](#). This is because, by drawing $x \in \{0, 1\}^T$ according to \mathcal{P} , the marginal distribution of x remains the same as earlier, i.e., $x \sim \mathcal{D} := \prod_{t=1}^T \text{Bernoulli}(\bar{p}_t)$ where $\bar{p} = (\frac{1}{5}, \frac{1}{5}, \dots, \frac{1}{5}, \frac{4}{5}, \frac{4}{5}, \dots, \frac{4}{5})$. It follows that

$$\text{OPT}_{\text{UCal}}(\mathcal{P}) \leq \mathbb{E}_{x \sim \mathcal{P}} [\text{UCal}(x, p)] = \mathbb{E}_{x \sim \mathcal{D}} [\text{UCal}(x, p)] \leq e^{-\Omega(T)},$$

and $\text{OPT}_{\text{UCal}^{\text{sub}}}(\mathcal{P}) \leq e^{-\Omega(T)}$ by the same token.

Truthful forecasts give a high U-Calibration error. It remains to analyze the truthful forecaster by lower bounding the expectation of $\text{UCal}(x, p^*)$ and $\text{UCal}^{\text{sub}}(x, p^*)$. In light of [Lemma 2](#) and [proposition 3](#), we let $\alpha := \frac{1}{5} - (1 - 2\gamma)c \in [\frac{1}{5} - c, \frac{1}{5} + c]$ for some $\gamma \in [0, 1]$ to be chosen later, and focus on lower bounding the quantity

$$X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)},$$

where $N_-^{(\alpha)} = \sum_{t=1}^T \mathbb{1}[p_t^* < \alpha]$ is the number of steps on which the prediction is strictly below α , and $X_-^{(\alpha)} = \sum_{t=1}^T x_t \cdot \mathbb{1}[p_t^* < \alpha]$ is the sum of the events on those steps.

By our construction of \mathcal{P} , for each $t \in \{1, 2, \dots, T/2\}$, p_t^* is drawn independently and uniformly from $[\frac{1}{5} - c, \frac{1}{5} + c]$. It follows that

$$\Pr[p_t^* < \alpha] = \frac{\alpha - (\frac{1}{5} - c)}{2c} = \frac{2\gamma c}{2c} = \gamma,$$

and $N_-^{(\alpha)}$ follows $\text{Binomial}(T/2, \gamma)$. Furthermore, conditioning on that $p_t^* < \alpha$, p_t^* is uniformly distributed over $[\frac{1}{5} - c, \alpha]$, which implies that the conditional probability of $x_t = 1$ is $(\frac{1}{5} - c + \alpha)/2 = \frac{1}{5} - (1 - \gamma)c$. Therefore, conditioning on the value of $N_-^{(\alpha)}$, $X_-^{(\alpha)}$ follows $\text{Binomial}(N_-^{(\alpha)}, \frac{1}{5} - (1 - \gamma)c)$.

By a Chernoff bound, as long as $\gamma = \Omega(1/T)$, $N_-^{(\alpha)} \in [\gamma T/4, \gamma T]$ holds with probability $\Omega(1)$. Furthermore, conditioning on the realization of $N_-^{(\alpha)} \in [\gamma T/4, \gamma T]$, $X_-^{(\alpha)}$ has a conditional expectation of $[\frac{1}{5} - (1 - \gamma)c] \cdot N_-^{(\alpha)}$ and a conditional variance of

$$N_-^{(\alpha)} \cdot \left[\frac{1}{5} - (1 - \gamma)c \right] \cdot \left[\frac{4}{5} + (1 - \gamma)c \right] \geq \frac{\gamma T}{4} \cdot \left(\frac{1}{5} - c \right) \cdot \frac{4}{5} \geq \Omega(\gamma T),$$

where the last step applies $c < 1/5$. By a central limit theorem, it holds with probability at least $\Omega(1)$ that

$$X_-^{(\alpha)} \geq \left[\frac{1}{5} - (1 - \gamma)c \right] \cdot N_-^{(\alpha)} + 2\sqrt{\gamma T}. \quad (10)$$

Assuming that both $N_-^{(\alpha)} \in [\gamma T/4, \gamma T]$ and (10) hold, for $\gamma = T^{-1/3}$, we have

$$\begin{aligned} X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)} &\geq \left[\frac{1}{5} - (1 - \gamma)c \right] \cdot N_-^{(\alpha)} - \left[\frac{1}{5} - (1 - 2\gamma)c \right] \cdot N_-^{(\alpha)} + 2\sqrt{\gamma T} \\ &= 2\sqrt{\gamma T} - \gamma c \cdot N_-^{(\alpha)} \\ &\geq 2\sqrt{\gamma T} - \gamma c \cdot \gamma T && (N_-^{(\alpha)} \leq \gamma T) \\ &\geq T^{1/3}. && (\gamma = T^{-1/3}, c \leq 1) \end{aligned}$$

Therefore, it holds with probability $\Omega(1)$ that $X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)} \geq T^{1/3}$. By [Lemma 2](#) and [proposition 3](#), we have

$$\text{err}_{\text{UCal}}(\mathcal{P}, \mathcal{A}^{\text{truthful}}) = \mathbb{E}_{(x, p^*) \sim \mathcal{P}} [\text{UCal}(x, p^*)] \geq 2 \mathbb{E}_{(x, p^*) \sim \mathcal{P}} \left[\max\{X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)}, 0\} \right] = \Omega(T^{1/3}).$$

Truthful forecasts give a high subsampled U-Calibration error. The analysis for UCal^{sub} is almost the same. Consider a random subset $S \sim \text{Unif}(2^{[T]})$. To lower bound $\text{UCal}(x|_S, p^*|_S)$, we examine the quantity

$$X_-^{(\alpha)} - \alpha \cdot N_-^{(\alpha)},$$

where $\alpha = \frac{1}{5} - (1 - 2\gamma)c$, $N_-^{(\alpha)} = \sum_{t=1}^T \mathbb{1}[t \in S \wedge p_t^* < \alpha]$, and $X_-^{(\alpha)} = \sum_{t=1}^T x_t \cdot \mathbb{1}[t \in S \wedge p_t^* < \alpha]$.

For each $t \in \{1, 2, \dots, T/2\}$, the events $t \in S$ and $p_t^* < \alpha$ are independent and happen with probabilities $1/2$ and γ , respectively. Therefore, $N_-^{(\alpha)}$ follows $\text{Binomial}(T/2, \gamma/2)$. Moreover, given $N_-^{(\alpha)}$, the conditional distribution of $X_-^{(\alpha)}$ is still $\text{Binomial}(N_-^{(\alpha)}, \frac{1}{5} - (1 - \gamma)c)$. Then, the rest of the analysis goes through by considering the typical realization of $N_-^{(\alpha)} \in [\gamma T/8, \gamma T/2]$. ■

Appendix B. Step Calibration

In light of [Proposition 3](#), V-Calibration ([Equation \(3\)](#)) corresponds to testing the predictions on intervals of form $[0, \alpha]$ and $(\alpha, 1]$: If, among the steps where the prediction is $< \alpha$ (resp., $> \alpha$), the actual fraction of ones is significantly higher (resp., lower) than α , we know that the predictions must be far from calibration. Naturally, we would get a stronger measure if, in the above comparison, we replace α with the actual average of the predictions, and take an absolute value to penalize both over- and under-estimation.

Formally, we consider the following calibration measure, which we term the *step calibration error*:

$$\text{stepCE}(x, p) := \sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \alpha] \right|.$$

Compared with the smooth calibration error of [\[KF08\]](#), the step calibration error replaces the family of Lipschitz functions with the family of “step functions”: $\{p \mapsto \mathbb{1}[p \leq \alpha] : \alpha \in [0, 1]\}$.

The definition above is robust in the sense that, if we replace the step functions with the union of a constant number of intervals, the resulting error only increases by a constant factor. The definition above also has a decision-theoretic interpretation: If we replace the benchmark in V-Calibration (i.e., $\inf_{\beta \in [0, 1]} \sum_{t=1}^T S_\alpha(x_t, \beta)$) with the sum

$$\sum_{t=1}^T \mathbb{E}_{x'_t \sim \text{Bernoulli}(p_t)} [S_\alpha(x'_t, p_t)] = \sum_{t=1}^T (p_t - \alpha) \cdot \text{sgn}(\alpha - p_t)$$

and add an additional absolute value, we get

$$\sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - \alpha) \cdot \text{sgn}(\alpha - p_t) - \sum_{t=1}^T (p_t - \alpha) \cdot \text{sgn}(\alpha - p_t) \right| = \sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha - p_t) \right|,$$

which has a similar form to stepCE, except that the step function is replaced with the sign function. This definition can be interpreted as follows: Assuming that the bits were indeed draw from $\text{Bernoulli}(p_1)$ through $\text{Bernoulli}(p_T)$, we expect a loss of $\mathbb{E}_{x'_t \sim \text{Bernoulli}(p_t)} [S_\alpha(x'_t, p_t)]$ at each step t . If our actual loss, $\sum_{t=1}^T S_\alpha(x_t, p_t)$, is significantly higher or lower, we are certain that the predictions cannot be calibrated.

We prove in [Section C.2](#) that the step calibration error is equivalent to the above variant of V-Calibration up to a constant factor.

Fact 11 For any $x \in \{0, 1\}^T$ and $p \in [0, 1]^T$,

$$\frac{1}{3} \sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha - p_t) \right| \leq \text{stepCE}(x, p) \leq \sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha - p_t) \right|.$$

B.1. Step Calibration Error is Complete, Sound and Decision-Theoretic

We show that the step calibration error is complete and sound in the sense of [Definition 1](#). Furthermore, it is decision-theoretic, i.e., stepCE always upper bounds the U-Calibration error up to a constant factor.

Proposition 12 *The step calibration error, stepCE , is complete and sound. Moreover, for any $x \in \{0, 1\}^T$ and $p \in [0, 1]^T$, it holds that*

$$\text{stepCE}(x, p) \geq \frac{1}{8} \text{UCal}(x, p).$$

Proof We start by showing that stepCE is decision-theoretic. Let A be a shorthand for

$$\sup_{\alpha \in [0, 1]} \max \left\{ X_-^{(\alpha)} - \alpha N_-^{(\alpha)}, \alpha N_+^{(\alpha)} - X_+^{(\alpha)} \right\}.$$

By [Lemma 2](#) and [proposition 3](#), we have

$$A = \frac{1}{2} \text{VCal}(x, p) \geq \frac{1}{4} \text{UCal}(x, p),$$

so it suffices to prove $\text{stepCE}(x, p) \geq A/2$.

Step calibration is decision-theoretic. Note that, over all possible values of $\alpha \in [0, 1]$, the term $\max \left\{ X_-^{(\alpha)} - \alpha N_-^{(\alpha)}, \alpha N_+^{(\alpha)} - X_+^{(\alpha)} \right\}$ only takes finitely many (concretely, $O(T)$) different values. Thus, the supremum A can be achieved at some $\alpha^* \in [0, 1]$, i.e.,

$$A = \max \left\{ X_-^{(\alpha^*)} - \alpha^* N_-^{(\alpha^*)}, \alpha^* N_+^{(\alpha^*)} - X_+^{(\alpha^*)} \right\}.$$

We consider the following two cases, depending on which of the two terms above is larger.

Case 1: $A = X_-^{(\alpha^*)} - \alpha^* N_-^{(\alpha^*)}$. If $\alpha^* = 0$, we have $N_-^{(\alpha^*)} = X_-^{(\alpha^*)} = 0$. Then, $A = 0$, and $\text{stepCE}(x, p) \geq A/2$ vacuously holds. If $\alpha^* > 0$, we can find $\beta \in [0, \alpha^*)$ such that $(\beta, \alpha^*) \cap \{p_1, p_2, \dots, p_T\} = \emptyset$. Then, for every $t \in [T]$, $p_t < \alpha^*$ holds if and only if $p_t \leq \beta$. It follows that

$$\begin{aligned} \text{stepCE}(x, p) &\geq \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \beta] \\ &= \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t < \alpha^*] && (p_t < \alpha^* \iff p_t \leq \beta) \\ &\geq \sum_{t=1}^T (x_t - \alpha^*) \cdot \mathbb{1}[p_t < \alpha^*] && (p_t < \alpha^* \implies x_t - p_t \geq x_t - \alpha^*) \\ &= X_-^{(\alpha^*)} - \alpha^* N_-^{(\alpha^*)} = A, \end{aligned}$$

where the first step relaxes the supremum in the definition of stepCE to a fixed term at β , and removes the absolute value.

Case 2: $A = \alpha^* N_+^{(\alpha^*)} - X_+^{(\alpha^*)}$. In this case, we note that

$$\sum_{t=1}^T (p_t - x_t) \cdot \mathbb{1}[p_t > \alpha^*] \geq \sum_{t=1}^T (\alpha^* - x_t) \cdot \mathbb{1}[p_t > \alpha^*] = \alpha^* N_+^{(\alpha^*)} - X_+^{(\alpha^*)} = A.$$

Since $\mathbb{1}[p_t > \alpha^*] = \mathbb{1}[p_t \leq 1] - \mathbb{1}[p_t \leq \alpha^*]$ holds for every $t \in [T]$, we have $B_1 - B_{\alpha^*} = A$, where

$$B_1 = \sum_{t=1}^T (p_t - x_t) \cdot \mathbb{1}[p_t \leq 1] \quad \text{and} \quad B_{\alpha^*} = \sum_{t=1}^T (p_t - x_t) \cdot \mathbb{1}[p_t \leq \alpha^*].$$

Then, either B_1 or B_{α^*} must have an absolute value of at least $A/2$. It follows that

$$\text{stepCE}(x, p) \geq \max_{\beta \in \{\alpha^*, 1\}} \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \beta] \right| = \max\{|B_{\alpha^*}|, |B_1|\} \geq A/2.$$

Completeness and soundness. To verify the completeness, we note that stepCE is always upper bounded by the expected calibration error (ECE) defined as

$$\text{ECE}(x, p) := \sum_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t = \alpha] \right|.$$

This is because, for every $\alpha \in [0, 1]$, it holds that

$$\left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \alpha] \right| \leq \sum_{\alpha' \in [0, \alpha]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t = \alpha'] \right| \leq \text{ECE}(x, p).$$

Then, since the ECE is complete (e.g., [HQYZ24, Table 2]), stepCE is also complete.

Similarly, since stepCE is lower bounded by UCal up to a constant factor, stepCE inherits the soundness of the U-Calibration error (e.g., [HQYZ24, Table 2]). \blacksquare

B.2. Step Calibration Error is Truthful After Subsampling

As a warm-up, we start by proving that the subsampled version of step calibration, $\text{stepCE}^{\text{sub}}$, is (α, β) -truthful for some parameters $\alpha, \beta = \text{polylog}(T/c)$ in the c -smoothed setting. Later, we will improve the parameter α to $O(\sqrt{\log(1/c)})$ and also prove its tightness (for the truthfulness of $\text{stepCE}^{\text{sub}}$).

Recall that in the c -smoothed setting, the conditional expectation of each x_t (given $x_{1:(t-1)}$) is sampled from a distribution with density upper bounded by $1/c$ (e.g., the uniform distribution over an interval of length c). Formally, the setting is described by a function $\mathcal{P} : \bigcup_{t=1}^T \{0, 1\}^{t-1} \rightarrow \Delta_c([0, 1])$, where for each $t \in [T]$ and $(x_1, x_2, \dots, x_{t-1}) \in \{0, 1\}^{t-1}$, $\mathcal{P}(x_1, x_2, \dots, x_{t-1})$ specifies a distribution over $[0, 1]$ —with a density upper bounded by c —from which the conditional expectation of $x_t | x_1, x_2, \dots, x_{t-1}$ is drawn. The sequence $x \in \{0, 1\}^T$ is determined sequentially as follows: For each $t \in [T]$, we draw $p_t^* \sim \mathcal{P}(x_1, x_2, \dots, x_{t-1})$ and then draw $x_t \sim \text{Bernoulli}(p_t^*)$. Furthermore, the true conditional probability for the event $x_t = 1$, p_t^* , is observable by the forecaster. Thus, the truthful forecaster always predicts p_t^* at time t .

The $c = \Omega(1)$ regime. Towards showing that $\text{stepCE}^{\text{sub}}$ is truthful, we will lower bound the optimal error that can be achieved on \mathcal{P} (namely, $\text{OPT}_{\text{stepCE}^{\text{sub}}}(\mathcal{P})$) and upper bound the error incurred by the truthful forecaster (namely, $\text{err}_{\text{stepCE}^{\text{sub}}}(\mathcal{P}, \mathcal{A}^{\text{truthful}})$). When $c = \Omega(1)$ is a fixed constant, we can easily obtain the following lower bound:

$$\text{OPT}_{\text{stepCE}^{\text{sub}}}(\mathcal{P}) = \Omega(\sqrt{T}),$$

i.e., every forecaster must incur an $\Omega(\sqrt{T})$ penalty measured by $\text{stepCE}^{\text{sub}}$ in expectation. The expectation above is over the randomness in the realized outcomes x as well as the forecaster that generates the predictions p . To see this, we note that

$$\text{stepCE}(x, p) := \sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \in [0, \alpha]] \right| \geq \left| \sum_{t=1}^T (x_t - p_t) \right|,$$

and it follows that

$$\text{stepCE}^{\text{sub}}(x, p) \geq_{y \sim \text{Unif}(\{0, 1\}^T)} \mathbb{E} \left[\left| \sum_{t=1}^T y_t \cdot (x_t - p_t) \right| \right].$$

Then, the results of [HQYZ24] lower bound the right-hand side above by $\mathbb{E}[\gamma(\text{Var}_T)]$, where

$$\text{Var}_T := \sum_{t=1}^T p_t^*(1 - p_t^*)$$

is the *realized variance* up to time T , and $\gamma(x) := x \cdot \mathbb{1}[x \leq 1] + \sqrt{x} \cdot \mathbb{1}[x > 1]$. In the $c = \Omega(1)$ regime, we can show that $\text{Var}_T \geq \Omega(T)$ with high probability, which implies the $\Omega(\sqrt{T})$ lower bound.

It remains to show that, when $c = \Omega(1)$, truthful forecasts lead to an $O(\sqrt{T})$ error with respect to the $\text{stepCE}^{\text{sub}}$ measure. As before, it suffices to show this for stepCE , as the argument should extend to the subsampled version easily. Via an easy “covering + union bound” argument, we can prove an upper bound of

$$\text{err}_{\text{stepCE}^{\text{sub}}}(\mathcal{P}, \mathcal{A}^{\text{truthful}}) = \mathbb{E}_{(x, p^*) \sim \mathcal{P}} [\text{stepCE}^{\text{sub}}(x, p^*)] = O\left(\sqrt{T \log(T/c)}\right),$$

almost matching $\text{OPT}_{\text{stepCE}^{\text{sub}}} = \Omega(\sqrt{T})$ in the $c = \Omega(1)$ regime.

The small- c regime. The argument above, unfortunately, does not directly apply to the case where c is small. This is because, in the worst case, Var_T can be as low as $\Theta(cT)$ (e.g., when each conditional distribution follows the uniform distribution over $[0, c]$). As a result, we can at best lower bound $\text{OPT}_{\text{stepCE}^{\text{sub}}}$ by $\Omega(\sqrt{cT})$. Then, our upper bound on the $\text{stepCE}^{\text{sub}}$ incurred by the truthful forecaster— $\mathbb{E}[\text{stepCE}^{\text{sub}}(x, p^*)] = O(\sqrt{T \log(T/c)})$ —would be higher by a factor of $\sqrt{1/c}$.

Instead, we will replace both bounds—the lower bound on $\text{OPT}_{\text{stepCE}^{\text{sub}}}(\mathcal{P})$ and the upper bound on $\text{err}_{\text{stepCE}^{\text{sub}}}(\mathcal{P}, \mathcal{A}^{\text{truthful}})$ —with ones that depend on the realized variance of the distribution, i.e.,

$$\text{Var}_T := \sum_{t=1}^T p_t^*(1 - p_t^*).$$

It should be noted that the “right” way of defining Var_T should be using p_t^* , rather than using the mean of the distribution $\mathcal{P}(x_1, x_2, \dots, x_{t-1})$. This is because revealing the value of $p_t^* \sim \mathcal{P}(x_{1:(t-1)})$ might significantly decrease the remaining variance in x_t .⁴

The following lemma is implicit in [HQYZ24]:

4. Consider the case that $\mathcal{P}(x_{1:(t-1)})$ is symmetric around $1/2$ and puts most of the probability mass near 0 and 1.

Lemma 13 (Theorem 6.5 of [HQYZ24]) *For any choice of \mathcal{P} , it holds that*

$$\text{OPT}_{\text{stepCE}^{\text{sub}}}(\mathcal{P}) = \Omega \left(\mathbb{E}_{\mathcal{P}} [\gamma(\text{Var}_T)] \right),$$

where $\gamma(x) := x \cdot \mathbb{1}[x \leq 1] + \sqrt{x} \cdot \mathbb{1}[x > 1]$.

Proof Note that

$$\text{stepCE}(x, p) = \sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \in [0, \alpha]] \right| \geq \left| \sum_{t=1}^T (x_t - p_t) \right|.$$

It follows that

$$\text{stepCE}^{\text{sub}}(x, p) \geq \mathbb{E}_{y \sim \text{Unif}(\{0, 1\}^T)} \left[\left| \sum_{t=1}^T y_t \cdot (x_t - p_t) \right| \right],$$

so the rest of the proof follows from [HQYZ24] (which relaxes the SSCE to the same expression as above). \blacksquare

We can show that the $\text{stepCE}^{\text{sub}}$ incurred by the truthful forecaster nearly matches the above, up to a multiplicative factor and an additive term of $\text{polylog}(T/c)$:

Lemma 14 *For any \mathcal{P} that is c -smoothed, we have*

$$\text{err}_{\text{stepCE}^{\text{sub}}}(\mathcal{P}, \mathcal{A}^{\text{truthful}}) \leq O \left(\mathbb{E} [\gamma(\text{Var}_T)] \cdot \sqrt{\log(T/c)} + \log^2(T/c) \right),$$

where $\gamma(x) := x \cdot \mathbb{1}[x \leq 1] + \sqrt{x} \cdot \mathbb{1}[x > 1]$.

Combining Lemmas 13 and 14 shows that $\text{stepCE}^{\text{sub}}$ is (α, β) -truthful for $\alpha = O(\sqrt{\log(T/c)})$ and $\beta = O(\log^2(T/c))$ in any c -smoothed setting.

Proof [Proof sketch] In the definition of stepCE , we replace $\alpha \in [0, 1]$ in the supremum with a (c/T^2) -net of $[0, 1]$. Note that the change in the calibration measure is bounded by the maximum number of p^* s that fall into the same length- (c/T^2) bin. This can be shown to be $O(1)$ in expectation and will not be dominating.

For each fixed choice of α (out of the $O(T^2/c)$ choices), the quantity $\sum_{t=1}^T (x_t - p_t^*) \cdot \mathbb{1}[p_t^* \in [0, \alpha]]$ induces a martingale. Then, we apply the same doubling trick as in [HQYZ24] to bound its deviation from zero. For the block in which the realized variance is roughly 2^k ($k = 0, 1, 2, \dots, O(\log T)$), we take a union bound over $O(T^2/c)$ martingales, each with $\leq T$ increments bounded between ± 1 and a total realized variance $\leq 2^k$. Freedman's inequality and the union bound show that the maximum deviation is at most $O \left(\sqrt{2^k \log(T/c)} + \log(T/c) \right)$ in expectation. Summing over k would then give an upper bound of

$$\mathbb{E} \left[\sqrt{\text{Var}_T} \right] \cdot O(\sqrt{\log(T/c)}) + O(\log T) \cdot O(\log(T/c)).$$

Finally, we note that $\sqrt{x} \leq \gamma(x) + 1$ holds for every $x \geq 0$, so the upper bound above can be further relaxed to $\mathbb{E} [\gamma(\text{Var}_T)] \cdot O(\sqrt{\log(T/c)}) + O(\log^2(T/c))$. \blacksquare

B.3. A Tighter Truthfulness Guarantee

We give a more refined analysis that improves the $\sqrt{\log(T/c)}$ multiplicative factor to $\sqrt{\log(1/c)}$ by implementing a more involved covering strategy.

We start by noting that an $\Omega(\sqrt{\log(1/c)})$ factor is unavoidable; this follows from a “binary search” construction similar to the one for the $O(\sqrt{T})$ - $\Omega(T)$ truthfulness gap of U-Calibration (Proposition 6).

Proposition 15 *For any $c \in [\Theta(2^{-T}), \Theta(1)]$, in the c -smoothed setting, the step calibration error measure has an $O(\sqrt{cT}+1)$ - $\Omega(\sqrt{T \log(1/c)})$ truthfulness gap and the subsampled step calibration error measure has an $O(\sqrt{T})$ - $\Omega(\sqrt{T \log(1/c)})$ truthfulness gap.*

Proof Without loss of generality, we assume that $c \leq \frac{1}{16}$ and T is divisible by $5 \lfloor \log_2 \frac{1}{8c} \rfloor$.

Construction of \mathcal{P} . We construct the c -smoothed setting \mathcal{P} by specifying the distribution from which each p_t^* is drawn. We divide the timesteps $[T]$ into $k+1$ epochs T_1, \dots, T_k, T_{k+1} , where $k = \lfloor \log_2 \frac{1}{8c} \rfloor$, epoch $T_{k+1} = \{\frac{T}{5} + 1, \frac{T}{5} + 2, \dots, T\}$, and epoch $T_i = \{\frac{T(i-1)}{5k} + 1, \frac{T(i-1)}{5k} + 2, \dots, \frac{Ti}{5k}\}$ for all $i \in [k]$. This division is well-defined by our assumption that T is divisible by $5k$.

For each $i \in [k]$, every timestep t in epoch T_i will have the same distribution of p_t^* : p_t^* is uniformly distributed over $[w_i - \frac{c}{2}, w_i + \frac{c}{2}]$, where w_i is defined in terms of the realized events of the previous epoch $\{x_t\}_{t \in T_{i-1}}$: We set $w_1 = \frac{1}{2}$ and, for every $i \geq 2$,

$$w_i = \begin{cases} w_{i-1} + \frac{1}{2^{i+2}} & \text{if } \mu_{i-1} \geq w_{i-1}, \\ w_{i-1} - \frac{1}{2^{i+2}} & \text{if } \mu_{i-1} < w_{i-1}, \end{cases}$$

where we use $\mu_i = \frac{1}{|T_i|} \sum_{t \in T_i} x_t$ to denote the average outcome in epoch T_i . This guarantees that $|w_i - w_{i'}| \geq 2^{-(k+3)} \geq c$ for all $i, i' \in [k+1]$ where $i \neq i'$ and that every p_t^* falls into the interval $[1/4, 3/4]$. This in turn ensures $\alpha^* = w_{k+1}$ satisfies, for every $i \in [k]$: (1) $w_i \notin [\alpha^* - c, \alpha^* + c]$, (2) $w_i < \alpha^*$ implies $\mu_i \geq w_i$, and (3) $w_i > \alpha^*$ implies $\mu_i < w_i$.

For the last epoch, which spans the last $4T/5$ timesteps, we define p_t^* to be sampled uniformly from $[0, c]$ for timesteps $t \in [T/5 + 1, 3T/5]$ and p_t^* to be sampled uniformly from $[1 - c, 1]$ for timesteps $t \in [3T/5 + 1, T]$.

Truthful forecaster. We will analyze the subsampled step calibration error of the truthful forecaster, with the (non-subsampled) step calibration error bound following identically. We first lower bound $\text{stepCE}^{\text{sub}}$ by fixing $\alpha = \alpha^*$ and removing the absolute value:

$$\begin{aligned} \text{stepCE}^{\text{sub}}(x, p^*) &= \mathbb{E}_{S \sim \text{Unif}(2^{[T]})} \left[\sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t^*) \cdot \mathbb{1}[p_t^* \in [0, \alpha] \wedge t \in S] \right| \right] \\ &\geq \mathbb{E}_{S \sim \text{Unif}(2^{[T]})} \left[\sum_{t=1}^T (x_t - p_t^*) \cdot \mathbb{1}[p_t^* \in [0, \alpha^*] \wedge t \in S] \right]. \end{aligned}$$

We next bound in expectation the right-hand summand for each epoch T_i individually. Fixing such an epoch i and realization of previous outcomes $x'_{1:\max T_{i-1}}$ where $\max T_{i-1}$ is the last timestep of

epoch $i - 1$, we can lower bound

$$\begin{aligned}
 & \mathbb{E}_{\substack{S \sim \text{Unif}(2^{[T]}) \\ (x, p^*) \sim \mathcal{P}}} \left[\sum_{t \in T_i} (x_t - p_t^*) \cdot \mathbb{1}[p_t^* \in [0, \alpha^*] \wedge t \in S] \mid x_{1:\max T_{i-1}} = x'_{1:\max T_{i-1}} \right] \\
 &= \frac{1}{2} \mathbb{E}_{(x, p^*) \sim \mathcal{P}} \left[\sum_{t \in T_i} (x_t - w_i) \cdot \mathbb{1}[w_i \in [0, \alpha^*]] \mid x_{1:\max T_{i-1}} = x'_{1:\max T_{i-1}} \right] \\
 &= \frac{1}{2} \mathbb{E}_{(x, p^*) \sim \mathcal{P}} \left[\mathbb{E}_{X \sim \text{Binomial}(|T_i|, w_i)} [\max\{X - |T_i| \cdot w_i, 0\}] \mid x_{1:\max T_{i-1}} = x'_{1:\max T_{i-1}} \right] \\
 &\geq \frac{1}{2} \min_{p \in [1/4, 3/4]} \mathbb{E}_{X \sim \text{Binomial}(|T_i|, p)} [\max\{X - |T_i| \cdot p, 0\}].
 \end{aligned}$$

In the above, the first step applies $w_i \notin [\alpha^* - c, \alpha^* + c]$ and marginalizes out S . The second step holds since $w_i \in [0, \alpha^*]$ if and only if $\mu_i \geq w_i$, which is equivalent to $\sum_{t \in T_i} x_t \geq |T_i| \cdot w_i$. The last inequality follows from $w_i \in [1/4, 3/4]$. To bound the last expectation, we first observe that the variance σ^2 of $X \sim \text{Binomial}(|T_i|, p)$ is at least $\sigma^2 \geq \frac{3}{16}|T_i|$. By the Berry-Esseen central limit theorem, we have $\Pr[X - |T_i| \cdot p \leq \sigma] \leq \Phi(1) + O(1/\sqrt{|T_i|})$. Thus,

$$\mathbb{E}_X [\max\{X - |T_i| \cdot p, 0\}] \geq \sigma \cdot \Pr[X - |T_i| \cdot p \geq \sigma] \geq \Omega(\sigma) \geq \Omega(\sqrt{|T_i|}).$$

We can thus bound the step calibration error in the first k epochs by

$$\mathbb{E}_{(x, p^*) \sim \mathcal{P}, S \sim \text{Unif}(2^{[T]})} \left[\sum_{i=1}^k \sum_{t \in T_i} (x_t - p_t^*) \cdot \mathbb{1}[p_t^* \in [0, \alpha^*] \wedge t \in S] \right] \geq \Omega\left(\sum_{i=1}^k \sqrt{|T_i|}\right) = \Omega(\sqrt{Tk}).$$

By definition of \mathcal{P} , in the last epoch $k + 1$, $p_t^* \in [0, \alpha^*]$ for the first $2T/5$ timesteps and $p_t^* \notin [0, \alpha^*]$ for the last $2T/5$ timesteps. Thus,

$$\begin{aligned}
 & \mathbb{E}_{(x, p^*) \sim \mathcal{P}, S \sim \text{Unif}(2^{[T]})} \left[\sum_{t \in T_{k+1}} (x_t - p_t^*) \cdot \mathbb{1}[p_t^* \in [0, \alpha^*] \wedge t \in S] \right] \\
 &= \frac{1}{2} \mathbb{E}_{(x, p^*) \sim \mathcal{P}} \left[\sum_{t=T/5+1}^{3T/5} (x_t - c/2) \right] = 0
 \end{aligned}$$

Combining the epochs, we obtain a subsampled step calibration error of at least

$$\text{err}_{\text{stepCE}^{\text{sub}}}(\mathcal{P}, \mathcal{A}^{\text{truthful}}) = \mathbb{E}_{(x, p^*) \sim \mathcal{P}} [\text{stepCE}^{\text{sub}}(x, p^*)] \geq \Omega(\sqrt{Tk}) = \Omega(\sqrt{T \log(1/c)}).$$

We can easily verify that the same analysis goes through without subsampling, i.e., we also have $\text{err}_{\text{stepCE}}(\mathcal{P}, \mathcal{A}^{\text{truthful}}) = \Omega(\sqrt{T \log(1/c)})$.

Non-truthful forecaster. We now turn to upper bounding the step calibration error of a dishonest forecaster. Suppose that the dishonest forecaster, denoted by \mathcal{A} , predicts $p_t = \frac{1}{2}$ for the first $T/5$ timesteps and computes $\Delta = \sum_{t=1}^{T/5} x_t - \frac{T}{10}$, which denotes the deviation of realized outcomes from our prediction. Note that $|\Delta| \leq T/10$. If $\Delta < 0$, $p_t = \frac{1}{2}$ was an overestimate. The forecaster \mathcal{A} then predicts $p_t = \frac{c}{2}$ for timesteps $t \in [T/5 + 1, 3T/5]$. For $t \in [3T/5 + 1, T]$, \mathcal{A} predicts $\frac{1}{2}$ until the bias at $\frac{1}{2}$ becomes non-negative. Formally, let $T' := \max \left\{ t \in [T] \mid \Delta + \sum_{\tau=3T/5+1}^{t-1} (x_{\tau-1} - \frac{1}{2}) < 0 \right\}$. Forecaster \mathcal{A} predicts $p_t = \frac{1}{2}$ for $t \in [3T/5 + 1, T']$ and $p_t = 1 - \frac{c}{2}$ for $t \in [T' + 1, T]$. That is, we intentionally underestimate the true $p_t^* = 1 - \frac{c}{2}$ by guessing $p_t = \frac{1}{2}$ to cancel out the bias from the first k epochs. Note that this implies $\Delta + \sum_{t=3T/5+1}^T (x_{t-1} - \frac{1}{2}) \leq \frac{1}{2}$. Let us condition on the event E that $T' \leq 9T/10$ (and $\Delta < 0$), which occurs with probability at least $1 - \exp(-\Omega(T))$ by Hoeffding's inequality as the complementary event would require that at least $T/10$ of the timesteps $t \in [3T/5 + 1, 9T/10]$ result in $x_t = 0$. Then,

$$\begin{aligned}
& \mathbb{E}_{(x,p) \sim (\mathcal{P}, \mathcal{A})} [\text{stepCE}(x, p) \mid \Delta < 0] \\
&= \mathbb{E}_{(x,p) \sim (\mathcal{P}, \mathcal{A})} \left[\sup_{\alpha \in [0,1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \in [0, \alpha]] \right| \mid \Delta < 0 \right] \\
&\leq \mathbb{E}_{(x,p) \sim (\mathcal{P}, \mathcal{A})} \left[\left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t = \frac{1}{2}] \right| \mid E \right] + \mathbb{E}_{(x,p) \sim (\mathcal{P}, \mathcal{A})} \left[\left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t = \frac{c}{2}] \right| \mid E \right] \\
&\quad + \mathbb{E}_{(x,p) \sim (\mathcal{P}, \mathcal{A})} \left[\left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t = 1 - \frac{c}{2}] \right| \mid E \right] + T \cdot \exp(-\Omega(T)) \\
&\leq \frac{1}{2} + \mathbb{E}_{X \sim \text{Binomial}(2T/5, \frac{c}{2})} \left[\left| X - \frac{cT}{5} \right| \right] + \mathbb{E}_{X \sim \text{Binomial}(T-T', \frac{c}{2})} \left[\left| X - \frac{c(T-T')}{2} \right| \right] + o(1) \\
&\leq O(\sqrt{cT} + 1).
\end{aligned}$$

The case where $\Delta \geq 0$ follows symmetrically. Defining $T' := \max \{ t \in [3T/5] \mid \Delta + \sum_{\tau=T/5+1}^{t-1} (x_{\tau-1} - \frac{1}{2}) < 0 \}$, we fix predictions of $p_t = 1 - \frac{c}{2}$ for timesteps $t \in [3T/5 + 1, T]$, $p_t = \frac{1}{2}$ for $t \in [T/5 + 1, T']$, and $p_t = \frac{c}{2}$ for $t \in [T' + 1, 3T/5]$. The event E that $T' \leq 5T/10$ (and $\Delta > 0$) occurs with probability at least $1 - \exp(-\Omega(T))$ by Hoeffding's inequality as the complementary event would require that at least $T/10$ of the timesteps $t \in [T/5 + 1, 5T/10]$ result in $x_t = 1$. Then,

$$\mathbb{E}_{(x,p) \sim (\mathcal{P}, \mathcal{A})} [\text{stepCE}(x, p) \mid \Delta \geq 0] \leq O(\sqrt{cT} + 1).$$

Therefore, we have

$$\text{OPT}_{\text{stepCE}}(\mathcal{P}) \leq \mathbb{E}_{(x,p) \sim (\mathcal{P}, \mathcal{A})} [\text{stepCE}(x, p)] \leq O(\sqrt{cT} + 1).$$

Applying the inequality $\text{stepCE}^{\text{sub}}(x, p) \leq \frac{1}{2} \text{stepCE}(x, p) + O(\sqrt{T})$ from [Lemma 24](#) gives

$$\text{OPT}_{\text{stepCE}^{\text{sub}}}(\mathcal{P}) \leq O(\sqrt{T}).$$

Since both $\text{err}_{\text{stepCE}}(\mathcal{P}, \mathcal{A}^{\text{truthful}})$ and $\text{err}_{\text{stepCE}}(\mathcal{P}, \mathcal{A}^{\text{truthful}})$ are lower bounded by $\Omega(\sqrt{T \log(1/c)})$, there is an $O(\sqrt{cT} + 1) - \Omega(\sqrt{T \log(1/c)})$ truthfulness gap for the step calibration error and an $O(\sqrt{T}) - \Omega(\sqrt{T \log(1/c)})$ truthfulness gap for the subsampled step calibration error. \blacksquare

We now turn to proving the tighter upper bound for the truthful forecaster's error.

Theorem 16 *For any $T \geq 2$, smoothness parameter $c \in (0, 1]$, and c -smoothed setting specified by \mathcal{P} , the expected subsampled step calibration error of truthfully predicting conditional probabilities $p_{1:T}^*$, $\text{err}_{\text{stepCE}^{\text{sub}}}(\mathcal{P}, \mathcal{A}^{\text{truthful}})$, is upper bounded by*

$$\begin{aligned} & \left(547\sqrt{\ln(1/c)} + 1161 \right) \mathbb{E} \left[\sqrt{\text{Var}_T} \right] + (\log_2 T + 3)[16 \log_2(T) \log_2(T/c) + 15] + 8\sqrt{2 \ln(1/c)} + 17\sqrt{2} \\ & \leq O \left(\sqrt{\log(1/c)} \right) \cdot \mathbb{E} \left[\sqrt{\text{Var}_T} \right] + O \left(\log^2(T) \log(T/c) \right). \end{aligned}$$

Combined with the lower bound $\text{OPT}_{\text{stepCE}^{\text{sub}}}(\mathcal{P}) \geq \Omega(\mathbb{E}_{\mathcal{P}}[\gamma(\text{Var}_T)])$ from Lemma 13 and the inequality $\sqrt{x} \leq \gamma(x) + 1$, the theorem above shows that $\text{stepCE}^{\text{sub}}$ is $(O(\sqrt{\log(1/c)}), \text{polylog}(T/c))$ -truthful in c -smoothed settings.

Proof For convenience, given $\alpha \in [0, 1]$ and $\mathbf{y} \in \{0, 1\}^T$, we define the martingale

$$M_t(\alpha, \mathbf{y}) := \sum_{s=1}^t y_s \cdot \mathbb{1}[p_s^* \leq \alpha] \cdot (x_s - p_s^*),$$

adapted to the filtration $(\mathcal{F}_t)_{t \in [T]}$ generated by sampling $p_t^* \sim \mathcal{P}(x_1, x_2, \dots, x_{t-1})$ and $x_t \sim \text{Bernoulli}(p_t^*)$ for each $t = 1, 2, \dots, T$. We can verify that $M_t(\alpha, \mathbf{y})$ is a martingale by observing that, conditioned on any realization of $x_{1:(s-1)} = x'_{1:(s-1)}$, we indeed have

$$\begin{aligned} & \mathbb{E}_{p_s^*, x_s} \left[y_s \cdot \mathbb{1}[p_s^* \leq \alpha] \cdot (x_s - p_s^*) \mid x_{1:(s-1)} = x'_{1:(s-1)} \right] \\ &= y_s \cdot \mathbb{E}_{p_s^* \sim \mathcal{P}(x'_{1:(s-1)})} \left[\mathbb{1}[p_s^* \leq \alpha] \cdot \mathbb{E}_{x_s \sim \text{Bernoulli}(p_s^*)} [x_s - p_s^*] \right] \\ &= 0. \end{aligned}$$

We can thus equivalently write our main claim as:

$$\begin{aligned} & \mathbb{E}_{\substack{(x, p^*) \sim \mathcal{P} \\ \mathbf{y} \sim \text{Unif}(\{0, 1\}^T)}} \left[\sup_{\alpha \in [0, 1]} |M_T(\alpha, \mathbf{y})| \right] \\ & \leq \left(547\sqrt{\ln(1/c)} + 1161 \right) \mathbb{E} \left[\sqrt{\text{Var}_T} \right] + (\log_2 T + 3)[16 \log_2(T) \log_2(T/c) + 15] + 8\sqrt{2 \ln(1/c)} + 17\sqrt{2}. \end{aligned}$$

Dividing into epochs. For any realization of $x_{1:T} \sim \mathcal{P}$, we will divide the timesteps $[T]$ into epochs in the same style as [HQYZ24]: Let the k -th epoch be the shortest period (following the $(k-1)$ -th epoch) whose realized variance is roughly at least 2^{k-1} . Formally, consider sequence $\tau_0, \tau_1, \dots \in \mathbb{Z} \cup \{\infty\}$, where each τ_k denotes the last step of the k -th epoch and is defined as $\tau_0 = 0$ and

$$\tau_k := \min\{t \in [\tau_{k-1} + 1, T] \mid \text{Var}_t - \text{Var}_{\tau_{k-1}} + (t - \tau_{k-1}) \cdot c \geq 2^{k-1}\} \cup \{\infty\}.$$

We will also write $I_k := \{\tau_{k-1} + 1, \tau_{k-1} + 2, \dots, \min\{T, \tau_k\}\}$ to denote the timesteps in epoch k .

We first recall the following useful facts about this epoch division from [HQYZ24]. The realized variance in epoch k , $\text{Var}_{\tau_k} - \text{Var}_{\tau_{k-1}}$, is at most $2^{k-1} + 1/4 \leq 2^k$. For the k -th epoch to be complete (i.e., $\tau_k < \infty$), a necessary condition is that $\text{Var}_T + cT \geq 2^{k-1}$. In particular, since $\text{Var}_T + cT \leq T/4 + T < 2T$, the $(\lceil \log_2 T \rceil + 2)$ -th epoch is never complete, i.e., $\tau_{\lceil \log_2 T \rceil + 2} = \infty$. The random variable $\mathbb{1}[t \in I_k]$ for an epoch k is measurable by $x_{1:(t-1)}$. We also note that the k -th epoch cannot last more than $\lceil 2^{k-1}/c \rceil$ steps: If $\tau_{k-1} + \lceil 2^{k-1}/c \rceil = t \leq T$, we would have $\text{Var}_t - \text{Var}_{\tau_{k-1}} + (t - \tau_{k-1}) \cdot c \geq \lceil 2^{k-1}/c \rceil \cdot c \geq 2^{k-1}$, which implies $\tau_k \leq \tau_{k-1} + \lceil 2^{k-1}/c \rceil$.

Next, we observe that $\mathbb{E}[\sqrt{\text{Var}_T}] = \Omega(\sqrt{cT})$ holds in every c -smoothed setting. Recall that each p_t^* is sampled from a distribution with density bounded by $1/c$. Therefore, with probability at least $1 - (1/c) \cdot (c/4) = 3/4$, we have $p_t^* \in [c/8, 1 - c/8]$, which implies $p_t^*(1 - p_t^*) \geq (c/8) \cdot (1 - c/8) \geq \frac{7c}{64}$. Therefore, with probability at least $1/2$, $p_t^*(1 - p_t^*) \geq \frac{7c}{64}$ holds for at least $T/2$ values of $t \in [T]$, and it follows that

$$\mathbb{E}[\sqrt{\text{Var}_T}] \geq \frac{1}{2} \cdot \sqrt{\frac{7c}{64} \cdot \frac{T}{2}} \geq \frac{\sqrt{7cT}}{16\sqrt{2}}.$$

We also want to bound the exponentially weighted sum

$$\sum_{k=2}^{\lceil \log_2 T \rceil + 2} \sqrt{2^k} \Pr[\tau_{k-1} < \infty].$$

To this end, let $j \in \mathbb{Z}$ be the random variable defined such that $\text{Var}_T + cT \in [2^j, 2^{j+1})$. We observe that $\sqrt{\text{Var}_T + cT} \geq \sqrt{2^j}$. Also recall that, for any $k \geq 2$, $\tau_{k-1} < \infty$ holds only if $\text{Var}_T + cT \geq 2^{k-2}$, which in turn holds only if $k - 2 \leq j$. We can thus bound

$$\begin{aligned} \sum_{k=2}^{\lceil \log_2 T \rceil + 2} \sqrt{2^k} \Pr[\tau_{k-1} < \infty] &\leq \mathbb{E} \left[\sum_{k=2}^{\lceil \log_2 T \rceil + 2} \sqrt{2^k} \cdot \mathbb{1}[\text{Var}_T + cT \geq 2^{k-2}] \right] \\ &\leq \mathbb{E} \left[\sum_{k=2}^{j+2} \sqrt{2^k} \right] \\ &\leq (4 + 2\sqrt{2}) \mathbb{E} [2^{j/2}] \\ &\leq (4 + 2\sqrt{2}) \mathbb{E} [\sqrt{\text{Var}_T + cT}] \\ &\leq 20(2 + \sqrt{2}) \mathbb{E} [\sqrt{\text{Var}_T}]. \end{aligned} \tag{11}$$

The last step applies $\mathbb{E}[\sqrt{\text{Var}_T}] \geq \frac{\sqrt{7cT}}{16\sqrt{2}}$, which gives

$$\mathbb{E}[\sqrt{\text{Var}_T + cT}] \leq \mathbb{E}[\sqrt{\text{Var}_T}] + \sqrt{cT} \leq \mathbb{E}[\sqrt{\text{Var}_T}] \cdot \left(1 + \frac{16\sqrt{2}}{\sqrt{7}}\right) \leq 10 \mathbb{E}[\sqrt{\text{Var}_T}].$$

Dividing step functions into rounded segments. Let $V_\varepsilon = \{0, \varepsilon, 2\varepsilon, \dots, \lfloor 1/\varepsilon \rfloor \varepsilon\}$ denote the multiples of ε in $[0, 1]$. Note that $|V_\varepsilon| = \lfloor 1/\varepsilon \rfloor + 1 \leq 2/\varepsilon$ elements. Let us fix an epoch k , condition

on the epoch being reached (i.e., $\tau_{k-1} < \infty$) and define the following restriction of the martingale $M_t(\alpha, \mathbf{y})$ to timesteps lying in epoch k :

$$M_t(\alpha, \mathbf{y}, k) := \sum_{s=1}^t y_s \cdot \mathbb{1}[p_s^* \leq \alpha] \cdot (x_s - p_s^*) \cdot \mathbb{1}[s \in I_k].$$

Recall that we want to bound the supremum of this martingale over different values of the step threshold α . Let us fix a step threshold $\alpha \in [0, 1]$ for now. Fixing some $\varepsilon^* > 0$ and integer $m \geq 1$ which we will specify later, we can define $w_0 = \lfloor \alpha / \varepsilon^* \rfloor \cdot \varepsilon^*$ as a rounding down of α onto the grid V_{ε^*} of resolution ε^* . Then, for all $i \in [m]$, we define recursively $w_i = w_{i-1} + 2^{-i} \varepsilon^*$ if $w_{i-1} + 2^{-i} \varepsilon^* \leq \alpha$ and $w_i = w_{i-1}$ otherwise. Equivalently, each w_i is the rounding of α down to the nearest multiple of $2^{-i} \varepsilon^*$. Note that

$$\{[0, w_0]\} \cup \{(w_{i-1}, w_i] : i \in [m]\}$$

forms a partition of $[0, w_m]$, so we have the decomposition

$$\mathbb{1}[x \leq \alpha] = \mathbb{1}[x \in [0, w_0]] + \sum_{i=1}^m \mathbb{1}[x \in (w_{i-1}, w_i]] + \mathbb{1}[x \in (w_m, \alpha)].$$

Thus,

$$\begin{aligned} |M_t(\alpha, \mathbf{y}, k)| &= \left| \sum_{s=1}^t y_s \cdot \left(\mathbb{1}[p_s^* \leq w_0] + \sum_{i=1}^m \mathbb{1}[p_s^* \in (w_{i-1}, w_i]] + \mathbb{1}[p_s^* \in (w_m, \alpha)] \right) \cdot (x_s - p_s^*) \cdot \mathbb{1}[s \in I_k] \right| \\ &\leq \left| \sum_{s=1}^t y_s \cdot \mathbb{1}[p_s^* \in (w_m, \alpha)] \cdot (x_s - p_s^*) \cdot \mathbb{1}[s \in I_k] \right| \\ &\quad + \left| \sum_{s=1}^t y_s \cdot \mathbb{1}[p_s^* \leq w_0] \cdot (x_s - p_s^*) \cdot \mathbb{1}[s \in I_k] \right| \\ &\quad + \sum_{i=1}^m \left| \sum_{s=1}^t y_s \cdot \mathbb{1}[p_s^* \in (w_{i-1}, w_i]] \cdot (x_s - p_s^*) \cdot \mathbb{1}[s \in I_k] \right|. \end{aligned}$$

We can simplify this expression by noting that the first term above is upper bounded by

$$\sum_{s=1}^t \mathbb{1}[p_s^* \in (w_m, w_m + 2^{-m} \cdot \varepsilon^*)],$$

since $\alpha < w_m + 2^{-m} \cdot \varepsilon^*$ while $y_s \in \{0, 1\}$ and $|x_s - p_s^*| \leq 1$ hold for every s . Let us also, with some abuse of notation, define for every interval $S \subset [0, 1]$:

$$M_t(S, \mathbf{y}, k) := \sum_{s=1}^t y_s \cdot \mathbb{1}[p_s^* \in S] \cdot (x_s - p_s^*) \cdot \mathbb{1}[s \in I_k].$$

Finally, we recall that $w_i \in V_{2^{-i} \cdot \varepsilon^*}$ for every $i \in \{0, 1, \dots, m\}$. Putting these together, we can simplify our upper bound on $|M_T(\alpha, \mathbf{y}, k)|$ to:

$$\begin{aligned} \sup_{\alpha \in [0,1]} |M_T(\alpha, \mathbf{y}, k)| &\leq \max_{w_0 \in V_{\varepsilon^*}} |M_T(w_0, \mathbf{y}, k)| + \sum_{i=1}^m \max_{w_i \in V_{2^{-i} \cdot \varepsilon^*}} |M_T((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k)| \\ &\quad + \max_{w_m \in V_{2^{-m} \cdot \varepsilon^*}} \sum_{t=1}^T \mathbb{1} [p_t^* \in (w_m, w_m + 2^{-m} \cdot \varepsilon^*)]. \end{aligned} \quad (12)$$

Bounding each summand. We first note the following lemmas, which we will prove in the sequel.

Lemma 17 *For any epoch k ,*

$$\mathbb{E} \left[\max_{w_0 \in V_{\varepsilon^*}} |M_T(w_0, \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] \leq \sqrt{2^{k+1} \ln(4/\varepsilon^*)} + \sqrt{2^{k-1} \pi} + 2 \ln(4/\varepsilon^*) + 2.$$

Lemma 18 *For any epoch k and level $i \in [m]$,*

$$\begin{aligned} &\mathbb{E} \left[\max_{w_i \in V_{2^{-i} \cdot \varepsilon^*}} |M_T((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] \\ &\leq \sqrt{(\varepsilon^* \cdot 2^{k-i+1}/c^2) \ln(2^{i+2}/\varepsilon^*)} + \sqrt{\pi \cdot \varepsilon^* \cdot 2^{k-i-1}/c^2} + 2 \ln(2^{i+2}/\varepsilon^*) + 2. \end{aligned}$$

Lemma 19 *For any epoch k ,*

$$\mathbb{E} \left[\max_{w_m \in V_{2^{-m} \cdot \varepsilon^*}} \sum_{t=1}^T \mathbb{1} [p_t^* \in (w_m, w_m + 2^{-m} \cdot \varepsilon^*)] \right] \leq 2T \cdot \frac{\varepsilon^*}{2^m c} + 3(\ln(2^{m+1}/\varepsilon^*) + 1).$$

We now set the parameters ε^* and m as

$$\varepsilon^* := c^2, \quad m := \lfloor \log_2 T \rfloor,$$

and substitute them into the above lemmas for the following simplified bound for each epoch k .

Lemma 20 *The subsampled step calibration error in epoch k is upper bounded by*

$$\mathbb{E} \left[\sup_{\alpha \in [0,1]} M_t(\alpha, \mathbf{y}, k) \mid \tau_{k-1} < \infty \right] \leq 8\sqrt{2^k \ln(1/c)} + 17\sqrt{2^k} + 16 \log_2(T) \log_2(T/c) + 15.$$

Summing Lemma 20 over all epochs, we have

$$\begin{aligned} &\mathbb{E} \left[\sup_{\alpha \in [0,1]} M_t(\alpha, \mathbf{y}) \right] \\ &\leq \sum_{k=1}^{\lfloor \log_2 T \rfloor + 2} \Pr[\tau_{k-1} < \infty] \mathbb{E} \left[\sup_{\alpha \in [0,1]} M_t(\alpha, \mathbf{y}, k) \mid \tau_{k-1} < \infty \right] \\ &\leq \sum_{k=1}^{\lfloor \log_2 T \rfloor + 2} \Pr[\tau_{k-1} < \infty] \left[8\sqrt{2^k \ln(1/c)} + 17\sqrt{2^k} + 16 \log_2(T) \log_2(T/c) + 15 \right] \\ &\leq \sum_{k=1}^{\lfloor \log_2 T \rfloor + 2} \Pr[\tau_{k-1} < \infty] \left[\sqrt{2^k} \left(8\sqrt{\ln(1/c)} + 17 \right) + 16 \log_2(T) \log_2(T/c) + 15 \right]. \end{aligned}$$

Applying Equation (11) gives the upper bound

$$\begin{aligned}
 & \mathbb{E} \left[\sup_{\alpha \in [0,1]} M_t(\alpha, \mathbf{y}) \right] \\
 & \leq \left(8\sqrt{\ln(1/c)} + 17 \right) \cdot \left[\sqrt{2} + \sum_{k=2}^{\lceil \log_2 T \rceil + 2} \Pr[\tau_{k-1} < \infty] \cdot \sqrt{2^k} \right] + (\lceil \log_2 T \rceil + 2) \cdot [16 \log_2(T) \log_2(T/c) + 15] \\
 & \leq \left(547\sqrt{\ln(1/c)} + 1161 \right) \mathbb{E} \left[\sqrt{\text{Var}_T} \right] + (\log_2 T + 3)[16 \log_2(T) \log_2(T/c) + 15] + 8\sqrt{2 \ln(1/c)} + 17\sqrt{2}.
 \end{aligned}$$

■

Remaining proofs. We now prove Lemma 17, Lemma 18, Lemma 19, and Lemma 20.

Lemma 17 For any epoch k ,

$$\mathbb{E} \left[\max_{w_0 \in V_{\varepsilon^*}} |M_T(w_0, \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] \leq \sqrt{2^{k+1} \ln(4/\varepsilon^*)} + \sqrt{2^{k-1}\pi} + 2 \ln(4/\varepsilon^*) + 2.$$

Proof We first note that the realized variance of $M_t(w_0, \mathbf{y}, k)$, for any particular choice of $w_0 \in V_{\varepsilon^*}$, is bounded by the realized variance within the k -th epoch, which is in turn upper bounded by 2^k . Thus, for any fixed choice of w_0 and any $p \in (0, 1)$, Freedman's inequality gives that, with probability at least $1 - p$,

$$|M_T(w_0, \mathbf{y}, k)| \leq \sqrt{2 \cdot 2^k \cdot \ln(2/p)} + 2 \ln(2/p).$$

Applying a union bound over the $|V_{\varepsilon^*}| \leq 2/\varepsilon^*$ choices of $w_0 \in V_{\varepsilon^*}$, we have with probability $1 - p$:

$$\max_{w_0 \in V_{\varepsilon^*}} |M_T(w_0, \mathbf{y}, k)| \leq \sqrt{2^{k+1} \ln((4/\varepsilon^*)/p)} + 2 \ln((4/\varepsilon^*)/p).$$

We can then use the standard inequality $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ (for $a, b \geq 0$) to relax the above into

$$\max_{w_0 \in V_{\varepsilon^*}} |M_T(w_0, \mathbf{y}, k)| \leq \sqrt{2^{k+1} \ln(4/\varepsilon^*)} + \sqrt{2^{k+1} \ln(1/p)} + 2 \ln(4/\varepsilon^*) + 2 \ln(1/p).$$

Taking the layer-cake representation, we have

$$\begin{aligned}
 & \mathbb{E} \left[\max_{w_0 \in V_{\varepsilon^*}} |M_t(w_0, \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] \\
 & \leq \int_0^1 \left[\sqrt{2^{k+1} \ln(4/\varepsilon^*)} + \sqrt{2^{k+1} \ln(1/p)} + 2 \ln(4/\varepsilon^*) + 2 \ln(1/p) \right] dp \\
 & = \sqrt{2^{k+1} \ln(4/\varepsilon^*)} + 2 \ln(4/\varepsilon^*) + \sqrt{2^{k+1}} \int_0^1 \sqrt{\ln(1/p)} dp + 2 \int_0^1 \ln(1/p) dp \\
 & = \sqrt{2^{k+1} \ln(4/\varepsilon^*)} + 2 \ln(4/\varepsilon^*) + \sqrt{2^{k-1}\pi} + 2,
 \end{aligned}$$

where the last equality uses the identities $\int_0^1 \sqrt{\ln(1/x)} dx = \frac{\sqrt{\pi}}{2}$ and $\int_0^1 \ln(1/x) dx = 1$. ■

Lemma 18 For any epoch k and level $i \in [m]$,

$$\begin{aligned} & \mathbb{E} \left[\max_{w_i \in V_{2^{-i} \cdot \varepsilon^*}} |M_T((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] \\ & \leq \sqrt{(\varepsilon^* \cdot 2^{k-i+1}/c^2) \ln(2^{i+2}/\varepsilon^*)} + \sqrt{\pi \cdot \varepsilon^* \cdot 2^{k-i-1}/c^2} + 2 \ln(2^{i+2}/\varepsilon^*) + 2. \end{aligned}$$

Proof Let us fix a $w_i \in V_{2^{-i} \cdot \varepsilon^*}$. We will now study segments of length $2^{-i} \cdot \varepsilon^*$. We first bound the realized variance of the martingale $M_t((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k)$.

Fact 21 Let $v_t := M_t((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k) - M_{t-1}((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k)$ denote the t -th term of the martingale. We have that

$$\sum_{s=1}^T \mathbb{E} [|v_s|^2 \mid x_{1:(s-1)}] \leq \varepsilon^* 2^{k-i}/c^2.$$

Proof Because the prior \mathcal{P} is smooth, the probability that the conditional probability lies in any short interval is small. More specifically,

$$\Pr [p_s^* \in (w_i, w_i + 2^{-i} \cdot \varepsilon^*) \mid x_{1:(s-1)}] \leq 2^{-i} \varepsilon^*/c$$

holds for any $x_{1:(s-1)} \in \{0, 1\}^{s-1}$. Also, let t' denote the first timestep of the k -th epoch and observe that the event $\mathbb{1}[t' = t]$ is measurable with $x_{1:t-1}$. Further recalling that the maximum length of each epoch is $\lceil 2^{k-1}/c \rceil$, we then have

$$\begin{aligned} & \sum_{s=1}^T \mathbb{E}_{x \sim \mathcal{P}} [|v_s|^2 \mid x_{1:(s-1)}] \\ &= \sum_{s=1}^T \mathbb{E}_{(x, p^*) \sim \mathcal{P}} [y_s \cdot \mathbb{1}[p_s^* \in (w_i, w_i + 2^{-i} \cdot \varepsilon^*)] \cdot (x_s - p_s^*)^2 \cdot \mathbb{1}[s \in I_k] \mid x_{1:(s-1)}] \\ &\leq \sum_{s=1}^T \mathbb{E}_{(x, p^*) \sim \mathcal{P}} [\mathbb{1}[p_s^* \in (w_i, w_i + 2^{-i} \cdot \varepsilon^*)] \cdot \mathbb{1}[s \in I_k] \mid x_{1:(s-1)}] \\ &\leq \sum_{s=1}^T \mathbb{E}_{x_{1:(s-1)}} \left[\mathbb{1}[t' = s] \sum_{\tau=0}^{\lceil 2^{k-1}/c \rceil - 1} \mathbb{E}_{(x, p^*) \sim \mathcal{P}} [\mathbb{1}[p_{s+\tau}^* \in (w_i, w_i + 2^{-i} \cdot \varepsilon^*)] \mid x_{1:(s-1)}] \right] \\ &\leq \sum_{s=1}^T \mathbb{E}_{x_{1:(s-1)}} \left[\mathbb{1}[t' = s] \sum_{\tau=0}^{\lceil 2^{k-1}/c \rceil - 1} 2^{-i} \varepsilon^*/c \right] \\ &\leq \lceil 2^{k-1}/c \rceil \cdot 2^{-i} \varepsilon^*/c \\ &\leq 2^k/c \cdot 2^{-i} \varepsilon^*/c = \varepsilon^* 2^{k-i}/c^2. \end{aligned}$$

■

Freedman's inequality gives for $p \in (0, 1)$, with probability at least $1 - p$,

$$|M_t((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k)| \leq \sqrt{2(\varepsilon^* \cdot 2^{k-i}/c^2) \ln(2/p)} + 2 \ln(2/p).$$

Taking a union bound on $V_{2^{-i}, \varepsilon^*}$, with probability at least $1 - p$,

$$\begin{aligned} & \max_{w_i \in V_{2^{-i}, \varepsilon^*}} |M_t((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k)| \\ & \leq \sqrt{(\varepsilon^* \cdot 2^{k-i+1}/c^2) \ln(4 \cdot (2^i/\varepsilon^*)/p)} + 2 \ln(4 \cdot (2^i/\varepsilon^*)/p). \end{aligned}$$

We then take the layer cake representation as before

$$\begin{aligned} & \mathbb{E} \left[\max_{w_i \in V_{2^{-i}, \varepsilon^*}} |M_t((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] \\ & \leq \int_0^1 \sqrt{(\varepsilon^* \cdot 2^{k-i+1}/c^2) \ln(2^{i+2}/\varepsilon^*)} + \sqrt{(\varepsilon^* \cdot 2^{k-i+1}/c^2) \ln(1/p)} + 2 \ln(4 \cdot (2^i/\varepsilon^*)/p) \, dp \\ & \leq \sqrt{(\varepsilon^* \cdot 2^{k-i+1}/c^2) \ln(2^{i+2}/\varepsilon^*)} + 2 \ln(2^{i+2}/\varepsilon^*) + \int_0^1 \sqrt{(\varepsilon^* \cdot 2^{k-i+1}/c^2) \ln(1/p)} + 2 \ln(1/p) \, dp \\ & = \sqrt{(\varepsilon^* \cdot 2^{k-i+1}/c^2) \ln(2^{i+2}/\varepsilon^*)} + 2 \ln(2^{i+2}/\varepsilon^*) + \sqrt{\pi \cdot \varepsilon^* \cdot 2^{k-i-1}/c^2} + 2. \end{aligned}$$

■

Lemma 19 For any epoch k ,

$$\mathbb{E} \left[\max_{w_m \in V_{2^{-m}, \varepsilon^*}} \sum_{t=1}^T \mathbb{1}[p_t^* \in (w_m, w_m + 2^{-m} \cdot \varepsilon^*)] \right] \leq 2T \cdot \frac{\varepsilon^*}{2^m c} + 3(\ln(2^{m+1}/\varepsilon^*) + 1).$$

Proof Let $X := \max_{w_m \in V_{2^{-m}, \varepsilon^*}} \sum_{t=1}^T \mathbb{1}[p_t^* \in (w_m, w_m + 2^{-m} \cdot \varepsilon^*)]$ denote the random variable of interest. Since \mathcal{P} is c -smoothed, for each fixed $w_m \in V_{2^{-m}, \varepsilon^*}$ and every $t \in [T]$, it holds that

$$\Pr[p_t^* \in (w_m, w_m + 2^{-m} \varepsilon^*) \mid x_{1:(t-1)}] \leq 2^{-m} \varepsilon^* / c.$$

Therefore,

$$\sum_{t=1}^T \mathbb{1}[p_t^* \in (w_m, w_m + 2^{-m} \varepsilon^*)]$$

is stochastically dominated by a binomial random variable that follows $\text{Binomial}(T, 2^{-m} \varepsilon^* / c)$. By the multiplicative Chernoff bound, we have

$$\Pr \left[\sum_{t=1}^T \mathbb{1}[p_t^* \in (w_m, w_m + 2^{-m} \varepsilon^*)] \geq (1 + \delta) \mu \right] \leq \exp \left(-\frac{\delta^2}{2 + \delta} \mu \right) \leq \exp \left(-\frac{\delta \mu}{3} \right),$$

where $\delta \geq 1$ and $\mu = T \cdot \frac{2^{-m} \varepsilon^*}{c}$.

By the union bound, for every $\delta \geq 1$, we have

$$\Pr[X \geq (1 + \delta)\mu] \leq |V_{2^{-m} \cdot \varepsilon^*}| \cdot \exp\left(-\frac{\delta\mu}{3}\right) \leq \lceil 2^m / \varepsilon^* \rceil \cdot \exp\left(-\frac{\delta\mu}{3}\right). \quad (13)$$

Therefore, we have

$$\begin{aligned} \mathbb{E}[X] &= \int_0^{+\infty} \Pr[X \geq \tau] \, d\tau \\ &\leq 2\mu + \int_{2\mu}^{+\infty} \Pr[X \geq \tau] \, d\tau \\ &\leq 2\mu + \mu \int_1^{+\infty} \Pr[X \geq (1 + \tau)\mu] \, d\tau. \end{aligned}$$

Plugging Equation (13) into the integral above gives

$$\int_1^{+\infty} \Pr[X \geq (1 + \tau)\mu] \, d\tau \leq \int_0^{+\infty} \min\{\lceil 2^m / \varepsilon^* \rceil \cdot e^{-\mu\tau/3}, 1\} \, d\tau = \frac{\ln \lceil 2^m / \varepsilon^* \rceil + 1}{\mu/3},$$

where the second step applies the identity

$$\int_0^{+\infty} \min\{ne^{-ax}, 1\} \, dx = \int_0^{(\ln n)/a} 1 \, dx + \int_{(\ln n)/a}^{+\infty} ne^{-ax} \, dx = \frac{\ln n}{a} + \frac{n}{a} \cdot e^{-(a \ln n)/a} = \frac{\ln n + 1}{a}.$$

Therefore, we conclude that

$$\mathbb{E}[X] \leq 2\mu + 3(\ln \lceil 2^m / \varepsilon^* \rceil + 1) = 2T \cdot \frac{\varepsilon^*}{2^m c} + 3(\ln \lceil 2^m / \varepsilon^* \rceil + 1) \leq 2T \cdot \frac{\varepsilon^*}{2^m c} + 3(\ln(2^{m+1} / \varepsilon^*) + 1).$$

■

Lemma 20 *The subsampled step calibration error in epoch k is upper bounded by*

$$\mathbb{E} \left[\sup_{\alpha \in [0, 1]} M_t(\alpha, \mathbf{y}, k) \mid \tau_{k-1} < \infty \right] \leq 8\sqrt{2^k \ln(1/c)} + 17\sqrt{2^k} + 16 \log_2(T) \log_2(T/c) + 15.$$

Proof By our choice of $\varepsilon^* = c^2$ and $m = \lfloor \log_2 T \rfloor$, we have $1/\varepsilon^* \leq 1/c^2$, $2^m \leq T < 2^{m+1}$. Filling these into Lemma 17, we have

$$\begin{aligned} &\mathbb{E} \left[\max_{w_0 \in V_{\varepsilon^*}} |M_t(w_0, \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] \\ &\leq \sqrt{2^{k+1} \ln(4/c^2)} + \sqrt{2^{k-1} \pi} + 2 \ln(4/c^2) + 2 \\ &\leq 2\sqrt{2^k \ln(1/c)} + \sqrt{2 \ln 4} \cdot \sqrt{2^k} + \sqrt{\pi/2} \cdot \sqrt{2^k} + 4 \ln(1/c) + (2 \ln 4 + 2) \\ &\leq 2\sqrt{2^k \ln(1/c)} + 3\sqrt{2^k} + 4 \ln(1/c) + 5. \end{aligned} \quad (14)$$

The bound in [Lemma 19](#) reduces to

$$\begin{aligned} \mathbb{E} \left[\max_{w_m \in V_{2^{-m}, \varepsilon^*}} \sum_{t=1}^T \mathbb{1} [p_t^* \in (w_m, w_m + 2^{-m} \cdot \varepsilon^*)] \right] &\leq 2T \cdot \frac{c^2}{2^m c} + 3 (\ln(2^{m+1}/c^2) + 1) \\ &\leq 4c + 3 \ln(2T/c^2) + 3 \\ &\leq 10 + 6 \ln(T/c). \end{aligned} \quad (15)$$

Similarly, plugging these constants into [Lemma 18](#) and summing over $i \in [m]$ gives

$$\begin{aligned} &\sum_{i=1}^m \mathbb{E} \left[\max_{w_i \in V_{2^{-i}, \varepsilon^*}} |M_t((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] \\ &\leq \sum_{i=1}^m \left(\sqrt{2^{k-i+1} \ln(2^{i+2}/c^2)} + \sqrt{2^{k-i-1} \pi} + 2 \ln(2^{i+2}/c^2) + 2 \right). \end{aligned} \quad (16)$$

To further simplify the right-hand side of [Equation \(16\)](#), we note that

$$\begin{aligned} \sum_{i=1}^m \sqrt{2^{k-i+1} \ln(2^{i+2}/c^2)} &= \sum_{i=1}^m \sqrt{2^k \ln(4/c^2) \cdot 2^{1-i} + 2^k \cdot 2^{1-i} \ln 2^i} \\ &\leq \sqrt{2^k \ln(4/c^2)} \cdot \sum_{i=1}^{+\infty} \sqrt{2^{1-i}} + \sqrt{2^k} \cdot \sum_{i=1}^{+\infty} \sqrt{2^{1-i} \ln 2^i} \\ &\leq (2 + \sqrt{2}) \sqrt{2^k \ln(4/c^2)} + 5\sqrt{2^k}, \end{aligned}$$

where the last inequality uses the geometric series $\sum_{n=1}^{\infty} \sqrt{2^{1-n}} = 2 + \sqrt{2}$ and $\sum_{n=1}^{\infty} \sqrt{2^{1-n} \ln 2^n} \approx 4.88 < 5$. We can similarly bound

$$\sum_{i=1}^m \sqrt{2^{k-i-1} \pi} \leq (1 + 1/\sqrt{2}) \sqrt{\pi 2^k}$$

and

$$\sum_{i=1}^m 2 \ln(2^{i+2}/c^2) = 2m \ln(4/c^2) + m(m+1) \ln 2$$

using the arithmetic series $\sum_{n=1}^m n = \frac{m(m+1)}{2}$. Putting these together, we have the following simplified expression for [Equation \(16\)](#):

$$\begin{aligned} &\sum_{i=1}^m \mathbb{E} \left[\max_{w_i \in V_{2^{-i}, \varepsilon^*}} |M_t((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] \\ &\leq (2 + \sqrt{2}) \sqrt{2^k \ln(4/c^2)} + (5 + (1 + 1/\sqrt{2})\sqrt{\pi}) \cdot \sqrt{2^k} + 2m \ln(4/c^2) + m(m+1) \ln 2 + 2m \\ &\leq 4\sqrt{2^k \ln(4/c^2)} + 9\sqrt{2^k} + 2m \ln(4/c^2) + 4m^2. \end{aligned} \quad (17)$$

We can then sum Equation (14), Equation (15) and Equation (17) to obtain the following simplification of Equation (12):

$$\begin{aligned}
& \mathbb{E} \left[\sup_{\alpha \in [0,1]} |M_t(\alpha, \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] \\
& \leq \mathbb{E} \left[\max_{w_0 \in V_{\varepsilon^*}} |M_t(w_0, \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] + \sum_{i=1}^m \mathbb{E} \left[\max_{w_i \in V_{2^{-i} \cdot \varepsilon^*}} |M_t((w_i, w_i + 2^{-i} \cdot \varepsilon^*), \mathbf{y}, k)| \mid \tau_{k-1} < \infty \right] \\
& \quad + \mathbb{E} \left[\max_{w_m \in V_{2^{-m} \cdot \varepsilon^*}} \sum_{t=1}^T \mathbb{1} [p_t^* \in (w_m, w_m + 2^{-m} \cdot \varepsilon^*)] \right] \\
& \leq \left(2\sqrt{2^k \ln(1/c)} + 3\sqrt{2^k} + 4\ln(1/c) + 5 \right) + \left(4\sqrt{2^k \ln(4/c^2)} + 9\sqrt{2^k} + 2m \ln(4/c^2) + 4m^2 \right) \\
& \quad + (10 + 6 \ln(T/c)) \\
& \leq 8\sqrt{2^k \ln(1/c)} + 17\sqrt{2^k} + 8m^2 + 4m \ln(1/c) + 6 \ln(T/c) + 4 \ln(1/c) + 15.
\end{aligned}$$

The last step above applies

$$\sqrt{2^k \ln(4/c^2)} = \sqrt{2^k \ln 4 + 2 \cdot 2^k \ln(1/c)} \leq \sqrt{2} \cdot \sqrt{2^k \ln(1/c)} + \sqrt{\ln 4} \cdot \sqrt{2^k}$$

and

$$2m \ln(4/c^2) = 4m \ln 2 + 4m \ln(1/c) \leq 4m^2 + 4m \ln(1/c).$$

Finally, the lemma follows from

$$8m^2 + 6 \ln(T/c) \leq 8(\log_2 T)^2 + 6 \log_2(T) \log_2(1/c) \leq 8 \log_2(T) \log_2(T/c)$$

and

$$4m \ln(1/c) + 4 \ln(1/c) \leq 8m \ln(1/c) \leq 8 \log_2(T) \log_2(T/c).$$

■

B.4. $(O(1), 0)$ -Truthfulness on Product Distributions

For product distributions, the subsampled step calibration error, $\text{stepCE}^{\text{sub}}$, enjoys a stronger $(O(1), 0)$ -truthfulness guarantee, even in the non-smoothed setting. Recall from Proposition 9 that the subsampled U-Calibration error, in contrast, can have an $e^{-\Omega(T)} \cdot O(\sqrt{T})$ truthfulness gap on product distributions.

Proposition 22 *On every product distribution \mathcal{D} over $\{0, 1\}^T$, it holds that*

$$\text{err}_{\text{stepCE}^{\text{sub}}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D})) \leq O(\text{OPT}_{\text{stepCE}^{\text{sub}}}(\mathcal{D})),$$

where the $O(\cdot)$ notation hides a universal constant that does not depend on \mathcal{D} .

Proof Suppose that $\mathcal{D} = \prod_{t=1}^T \text{Bernoulli}(p_t^*)$ for some $p^* \in [0, 1]^T$. By Lemma 13, we have

$$\text{OPT}_{\text{stepCE}^{\text{sub}}}(\mathcal{D}) = \Omega(\gamma(\text{Var}_T)),$$

where $\text{Var}_T := \sum_{t=1}^T p_t^*(1 - p_t^*)$ and $\gamma(x) := x \cdot \mathbb{1}[x \leq 1] + \sqrt{x} \cdot \mathbb{1}[x > 1]$.

It remains to show that the truthful forecaster $\mathcal{A}^{\text{truthful}}(\mathcal{D})$, which predicts $p_t = p_t^*$ at every step t , satisfies

$$\text{err}_{\text{stepCE}^{\text{sub}}}(\mathcal{D}, \mathcal{A}^{\text{truthful}}(\mathcal{D})) = \mathbb{E}_{x \sim \mathcal{D}} [\text{stepCE}^{\text{sub}}(x, p^*)] = O(\gamma(\text{Var}_T)).$$

We consider the following two cases, depending on whether Var_T is below or above 1.

Case 1: $\text{Var}_T \leq 1$. We note that, for every $y \in \{0, 1\}^T$ and $\alpha \in [0, 1]$,

$$\left| \sum_{t=1}^T y_t \cdot (x_t - p_t) \cdot \mathbb{1}[p_t \in [0, \alpha]] \right| \leq \sum_{t=1}^T |x_t - p_t|.$$

It follows that $\text{stepCE}^{\text{sub}}(x, p) \leq \sum_{t=1}^T |x_t - p_t|$, and

$$\mathbb{E}_{x \sim \mathcal{D}} [\text{stepCE}^{\text{sub}}(x, p^*)] \leq \sum_{t=1}^T \mathbb{E}_{x_t \sim \text{Bernoulli}(p_t^*)} [|x_t - p_t^*|] = \sum_{t=1}^T 2p_t^*(1 - p_t^*) = 2\gamma(\text{Var}_T).$$

Case 2: $\text{Var}_T > 1$. Without loss of generality, we assume that $p_1^* \leq p_2^* \leq \dots \leq p_T^*$, as the behavior of the truthful forecaster and the resulting $\text{stepCE}^{\text{sub}}$ are invariant up to the reordering of timesteps. For fixed $x, y \in \{0, 1\}^T$, we have

$$\sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T y_t \cdot (x_t - p_t^*) \cdot \mathbb{1}[p_t^* \in [0, \alpha]] \right| \leq \max_{t \in [T]} \left| \sum_{i=1}^t y_i \cdot (x_i - p_i^*) \right|.$$

Taking an expectation over $x \sim \mathcal{D}$ and $y \sim \text{Unif}(\{0, 1\}^T)$ shows that

$$\mathbb{E}_{x \sim \mathcal{D}} [\text{stepCE}^{\text{sub}}(x, p_t^*)] \leq \mathbb{E}_{x \sim \mathcal{D}, y \sim \text{Unif}(\{0, 1\}^T)} \left[\max_{t \in [T]} \left| \sum_{i=1}^t y_i \cdot (x_i - p_i^*) \right| \right].$$

Over the randomness in (x, y) , consider the random walk $(X_t)_{t=0}^T$ defined as $X_t := \sum_{i=1}^t y_i(x_i - p_i^*)$. Then, $(X_t)_{t=0}^T$ forms a martingale in which the increment at time t has a variance of $p_t^*(1 - p_t^*)/2$. The total variance is then $\mathbb{E}[X_T^2] = \sum_{t=1}^T p_t^*(1 - p_t^*)/2 = \text{Var}_T/2$. Then, by Kolmogorov's inequality,

$$\Pr \left[\max_{t \in [T]} |X_t| \geq \tau \right] \leq \frac{\mathbb{E}[X_T^2]}{\tau^2} = \frac{\text{Var}_T}{2\tau^2}$$

holds for every $\tau > 0$. It follows that

$$\begin{aligned} \mathbb{E} \left[\max_{t \in [T]} |X_t| \right] &= \int_0^{+\infty} \Pr \left[\max_{t \in [T]} |X_t| \geq \tau \right] d\tau \\ &\leq \int_0^{+\infty} \min \left\{ \frac{\text{Var}_T}{2\tau^2}, 1 \right\} d\tau = O(\sqrt{\text{Var}_T}) = O(\gamma(\text{Var}_T)). \end{aligned}$$

■

B.5. An $O(\sqrt{T})$ Step Calibration Error Algorithm

Theorem 23 *Algorithm 1 guarantees an expected step calibration error of $O(\sqrt{T \log T})$, even when the T events are adversarially and adaptively chosen. Moreover, the forecasts made by the algorithm each randomize over at most two probabilities.*

By the inequality $\text{stepCE}^{\text{sub}}(x, p) \leq \frac{1}{2} \text{stepCE}(x, p) + O(\sqrt{T})$ (Lemma 24), the same algorithm guarantees an $O(\sqrt{T \log T})$ expected error with respect to the subsampled version $\text{stepCE}^{\text{sub}}$ as well.

Algorithm 1 Prediction Algorithm with Hedge

Require: Number of buckets $k \geq 2$, time horizon T

- 1: Initialize weight w_1 to be uniform over $\{\pm 1\} \times [k]$
 - 2: **for** $t = 1$ to T **do**
 - 3: **if** $\mathbb{E}_{(\sigma, i) \sim w_t} \left[\sigma \cdot \mathbb{1} \left[\frac{j-1}{k-1} \leq \frac{i-1}{k-1} \right] \right] \geq 0$ for all $j \in [k]$ **then**
 - 4: Predict $p_t = 1$
 - 5: **else if** $\mathbb{E}_{(\sigma, i) \sim w_t} \left[\sigma \cdot \mathbb{1} \left[\frac{j-1}{k-1} \leq \frac{i-1}{k-1} \right] \right] \leq 0$ for all $j \in [k]$ **then**
 - 6: Predict $p_t = 0$
 - 7: **else**
 - 8: Find $j \in [k-1]$ such that

$$\mathbb{E}_{(\sigma, i) \sim w_t} \left[\sigma \cdot \mathbb{1} \left[\frac{j-1}{k-1} \leq \frac{i-1}{k-1} \right] \right] \cdot \mathbb{E}_{(\sigma, i) \sim w_t} \left[\sigma \cdot \mathbb{1} \left[\frac{j}{k-1} \leq \frac{i-1}{k-1} \right] \right] \leq 0$$
 - 9: Find $q \in [0, 1]$ such that

$$q \cdot \mathbb{E}_{(\sigma, i) \sim w_t} \left[\sigma \cdot \mathbb{1} \left[\frac{j-1}{k-1} \leq \frac{i-1}{k-1} \right] \right] + (1-q) \cdot \mathbb{E}_{(\sigma, i) \sim w_t} \left[\sigma \cdot \mathbb{1} \left[\frac{j}{k-1} \leq \frac{i-1}{k-1} \right] \right] = 0$$
 - 10: Predict $p_t = \frac{j-1}{k-1}$ with probability q and $p_t = \frac{j}{k-1}$ with probability $1-q$
 - 11: **end if**
 - 12: Observe $x_t \in \{0, 1\}$ from Nature
 - 13: Compute w_{t+1} by applying the Hedge algorithm to the cost functions $(c_{x_\tau, p_\tau})_{\tau \in [t]}$, where

$$c_{x_\tau, p_\tau}(\sigma, i) = 1 - \frac{1}{2} \sigma \cdot \mathbb{1} \left[p_\tau \leq \frac{i-1}{k-1} \right] \cdot (x_\tau - p_\tau).$$
 - 14: **end for**
-

Proof Given sign $\sigma \in \{\pm 1\}$ and bucket $i \in [k]$, we define $\ell_{\sigma, i} : \{0, 1\} \times [0, 1] \rightarrow [-1, 1]$ to map from a realization $x \in \{0, 1\}$ and prediction $p \in [0, 1]$ to a loss

$$\ell_{\sigma, i}(x, p) = \sigma \cdot \mathbb{1} \left[p \leq \frac{i-1}{k-1} \right] \cdot (x - p).$$

Then, on events $x_{1:T} \in \{0, 1\}^T$ and predictions $p_{1:T} \in \left\{ \frac{i-1}{k-1} : i \in [k] \right\}^T$, the step calibration error can be written as

$$\text{stepCE}(x_{1:T}, p_{1:T}) = \max_{\sigma^* \in \{\pm 1\}, i^* \in [k]} \sum_{t=1}^T \ell_{\sigma^*, i^*}(x_t, p_t).$$

Furthermore, the cost function $c_{x,p} : \{\pm 1\} \times [k] \mapsto [0, 1]$ used by the Hedge algorithm in [Algorithm 1](#) is simply

$$c_{x,p}(\sigma, i) = 1 - \frac{1}{2} \cdot \ell_{\sigma,i}(x, p).$$

Guarantees of Hedge. Consider the sequence of distributions $w_{1:T} \in \Delta(\{\pm 1\} \times [k])$ given by applying Hedge to the cost functions $(c_{x_t, p_t})_{t \in [T]}$:

$$w_t = \text{Hedge}(c_{x_1, p_1}, c_{x_2, p_2}, \dots, c_{x_{t-1}, p_{t-1}}).$$

Hedge guarantees that, even though (x_t, p_t) is allowed to depend on $w_{1:t}$ at each step $t \in [T]$, we still have

$$\min_{\sigma^* \in \{\pm 1\}, i^* \in [k]} \sum_{t=1}^T c_{x_t, p_t}(\sigma^*, i^*) \geq \sum_{t=1}^T \mathbb{E}_{(\sigma, i) \sim w_t} [c_{x_t, p_t}(\sigma, i)] - O\left(\sqrt{T \log k}\right).$$

Equivalently,

$$\max_{\sigma^* \in \{\pm 1\}, i^* \in [k]} \sum_{t=1}^T \ell_{\sigma^*, i^*}(x_t, p_t) \leq \sum_{t=1}^T \mathbb{E}_{(\sigma, i) \sim w_t} [\ell_{\sigma, i}(x_t, p_t)] + O\left(\sqrt{T \log k}\right).$$

Therefore, to upper bound $\text{stepCE}(x, p)$, it remains to control the loss $\mathbb{E}_{(\sigma, i) \sim w_t} [\ell_{\sigma, i}(x_t, p_t)]$ at each step.

Control the per-step loss. Fix a timestep $t \in [T]$. We condition on the realization of $x_{1:(t-1)}$ and $p_{1:(t-1)}$, and shorthand w for $w_t \in \Delta(\{\pm 1\} \times [k])$. For a fixed prediction distribution $\mathbf{p} \in \Delta(\{\frac{i-1}{k-1} : i \in [k]\})$, we can write

$$\begin{aligned} \max_{x \in \{0, 1\}} \mathbb{E}_{p \sim \mathbf{p}} \left[\mathbb{E}_{(\sigma, i) \sim w} [\ell_{\sigma, i}(x, p)] \right] &= \max_{x \in \{0, 1\}} \mathbb{E}_{(\sigma, i) \sim w} \left[\sigma \cdot \mathbb{E}_{p \sim \mathbf{p}} \left[\mathbb{1} \left[p \leq \frac{i-1}{k-1} \right] \cdot (x - p) \right] \right] \\ &= \max_{x \in \{0, 1\}} \mathbb{E}_{p \sim \mathbf{p}} [x \cdot C_p] - \mathbb{E}_{p \sim \mathbf{p}} [p \cdot C_p] \\ &= \max \left\{ \mathbb{E}_{p \sim \mathbf{p}} [C_p], 0 \right\} - \mathbb{E}_{p \sim \mathbf{p}} [p \cdot C_p], \end{aligned}$$

where $C_p := \mathbb{E}_{(\sigma, i) \sim w} \left[\sigma \cdot \mathbb{1} \left[p \leq \frac{i-1}{k-1} \right] \right]$.

Suppose that $C_p \geq 0$ holds for all $p \in \{\frac{i-1}{k-1} : i \in [k]\}$. Then, we can let \mathbf{p} be the degenerate distribution at 1 and get

$$\max_{x \in \{0, 1\}} \mathbb{E}_{p \sim \mathbf{p}} \left[\mathbb{E}_{(\sigma, i) \sim w} [\ell_{\sigma, i}(x, p)] \right] = \max\{C_1, 0\} - 1 \cdot C_1 = 0.$$

Similarly, if $C_p \leq 0$ holds for every $p \in \{\frac{i-1}{k-1} : i \in [k]\}$, we can set \mathbf{p} to be deterministically 0 and get

$$\max_{x \in \{0, 1\}} \mathbb{E}_{p \sim \mathbf{p}} \left[\mathbb{E}_{(\sigma, i) \sim w} [\ell_{\sigma, i}(x, p)] \right] = \max\{C_0, 0\} - 0 \cdot C_0 = 0.$$

If neither holds, there must exist $p_1, p_2 \in \left\{ \frac{i-1}{k-1} : i \in [k] \right\}$ such that $p_2 - p_1 = \frac{1}{k-1}$ and $C_{p_1} \cdot C_{p_2} \leq 0$. Then, there also exists $q \in [0, 1]$ such that

$$q \cdot C_{p_1} + (1 - q) \cdot C_{p_2} = 0.$$

We accordingly let \mathbf{p} take value p_1 with probability q and take value p_2 with probability $1 - q$. This choice ensures $\mathbb{E}_{\mathbf{p} \sim \mathbf{p}} [C_p] = q \cdot C_{p_1} + (1 - q) \cdot C_{p_2} = 0$, which further implies

$$\begin{aligned} \max_{x \in \{0,1\}} \mathbb{E}_{\mathbf{p} \sim \mathbf{p}} \left[\mathbb{E}_{(\sigma,i) \sim w} [\ell_{\sigma,i}(x, \mathbf{p})] \right] &= \max_{\mathbf{p} \sim \mathbf{p}} \left\{ \mathbb{E}_{\mathbf{p} \sim \mathbf{p}} [C_p], 0 \right\} - \mathbb{E}_{\mathbf{p} \sim \mathbf{p}} [p \cdot C_p] \\ &= 0 - q \cdot (p_1 C_{p_1}) - (1 - q) \cdot (p_2 C_{p_2}) \\ &= -p_1 (q C_{p_1} + (1 - q) C_{p_2}) - (p_2 - p_1) (1 - q) C_{p_2} \\ &= -p_1 \cdot \mathbb{E}_{\mathbf{p} \sim \mathbf{p}} [C_p] - (p_2 - p_1) (1 - q) C_{p_2} \\ &\leq \frac{1}{k-1}, \end{aligned}$$

with the last inequality following from $\mathbb{E}_{\mathbf{p} \sim \mathbf{p}} [C_p] = 0$, $p_2 - p_1 = \frac{1}{k-1}$ and $|C_{p_2}| \leq 1$.

Note that our construction of \mathbf{p} coincides with the random choice of p_t at each timestep t in [Algorithm 1](#). Therefore, it holds for every $t \in [T]$ that

$$\mathbb{E}_{(\sigma,i) \sim w_t} [\ell_{\sigma,i}(x_t, p_t)] = \mathbb{E}_{w_t} \left[\mathbb{E}_{(\sigma,i) \sim w_t} [\ell_{\sigma,i}(x_t, p_t) \mid w_t] \right] \leq \mathbb{E}_{w_t} \left[\max_{x' \in \{0,1\}} \mathbb{E}_{\substack{(\sigma,i) \sim w_t \\ p_t \sim \mathbf{p}_t}} [\ell_{\sigma,i}(x', p_t)] \right] \leq \frac{1}{k-1}.$$

We conclude that

$$\begin{aligned} \mathbb{E} [\text{stepCE}(x, p)] &= \mathbb{E} \left[\max_{\sigma^* \in \{\pm 1\}, i^* \in [k]} \sum_{t=1}^T \ell_{\sigma^*, i^*}(x_t, p_t) \right] \\ &\leq \sum_{t=1}^T \mathbb{E}_{(\sigma,i) \sim w_t} [\ell_{\sigma,i}(x_t, p_t)] + O(\sqrt{T \log k}) \\ &\leq \frac{T}{k-1} + O(\sqrt{T \log k}). \end{aligned}$$

Choosing $k = T$ gives the $O(\sqrt{T \log T})$ bound. ■

Appendix C. Basic Facts

In this section, we prove the equivalent formulation of the V-calibration error [\[KPLST23\]](#) ([Proposition 3](#)), establish the decision-theoretic interpretation of step calibration ([Fact 11](#)), and show that $\text{stepCE}^{\text{sub}}$ is close to stepCE in general ([Lemma 24](#)).

C.1. Proof of Proposition 3

Proposition 3 *The V-Calibration error takes the alternative form*

$$\text{VCal}(x, p) = 2 \cdot \sup_{\alpha \in [0, 1]} \max\{X_-^{(\alpha)} - \alpha N_-^{(\alpha)}, \alpha N_+^{(\alpha)} - X_+^{(\alpha)}\},$$

where $N_-^{(\alpha)} := \sum_{t=1}^T \mathbb{1}[p_t < \alpha]$, $N_+^{(\alpha)} := \sum_{t=1}^T \mathbb{1}[p_t > \alpha]$, $X_-^{(\alpha)} := \sum_{t=1}^T x_t \cdot \mathbb{1}[p_t < \alpha]$, and $X_+^{(\alpha)} := \sum_{t=1}^T x_t \cdot \mathbb{1}[p_t > \alpha]$.

Proof Recall that $S_\alpha(x, p) := (\alpha - x) \cdot \text{sgn}(p - \alpha)$ and that the V-Calibration error is given by

$$\begin{aligned} \text{VCal}(x, p) &= \sup_{\alpha, \beta \in [0, 1]} \left[\sum_{t=1}^T S_\alpha(x_t, p_t) - \sum_{t=1}^T S_\alpha(x_t, \beta) \right] \\ &= \sup_{\alpha \in [0, 1]} \left[\sum_{t=1}^T S_\alpha(x_t, p_t) - \sum_{t=1}^T S_\alpha(x_t, \mu) \right] \\ &= \sup_{\alpha \in [0, 1]} \left[\sum_{t=1}^T (x_t - \alpha) \cdot \text{sgn}(\alpha - p_t) - \sum_{t=1}^T (x_t - \alpha) \cdot \text{sgn}(\alpha - \mu) \right], \end{aligned}$$

where the optimal choice of β is always $\beta = \mu := \frac{1}{T} \sum x_t$, since S_α is proper.

For $\alpha \in [0, 1]$, we introduce the shorthands

$$f(\alpha) := \sum_{t=1}^T (x_t - \alpha) \cdot \text{sgn}(\alpha - p_t) - \sum_{t=1}^T (x_t - \alpha) \cdot \text{sgn}(\alpha - \mu)$$

and

$$g(\alpha) := \max\{X_-^{(\alpha)} - \alpha N_-^{(\alpha)}, \alpha N_+^{(\alpha)} - X_+^{(\alpha)}\}.$$

Towards proving $\sup_{\alpha \in [0, 1]} f(\alpha) = 2 \sup_{\alpha \in [0, 1]} g(\alpha)$, we will first show that, for $S := \{p_1, p_2, \dots, p_T, \mu\}$,

$$\sup_{\alpha \in [0, 1]} f(\alpha) = \sup_{\alpha \in [0, 1] \setminus S} f(\alpha) \quad \text{and} \quad \sup_{\alpha \in [0, 1]} g(\alpha) = \sup_{\alpha \in [0, 1] \setminus S} g(\alpha).$$

In other words, ignoring the case that α coincides with a prediction p_t or the overall average μ does not change either supremum. Then, we will show that, for every $\alpha \in [0, 1] \setminus S$, we indeed have $f(\alpha) = 2g(\alpha)$. The desired identity would then follow from

$$\sup_{\alpha \in [0, 1]} f(\alpha) = \sup_{\alpha \in [0, 1] \setminus S} f(\alpha) = 2 \sup_{\alpha \in [0, 1] \setminus S} g(\alpha) = 2 \sup_{\alpha \in [0, 1]} g(\alpha).$$

Ignore atypical values for f . First, we focus on the value of $f(\alpha_0)$ for some $\alpha_0 \in S = \{p_1, p_2, \dots, p_T, \mu\}$. We claim that

$$f(\alpha_0) = \frac{1}{2} \left[\lim_{\alpha \rightarrow \alpha_0^-} f(\alpha) + \lim_{\alpha \rightarrow \alpha_0^+} f(\alpha) \right], \quad (18)$$

i.e., $f(\alpha_0)$ is equal to the average of the left-sided and right-sided limits of f at α_0 . Assuming Equation (18), if $\alpha_0 \in (0, 1)$, when α approaches α_0 from one of the two sides, the limit of $f(\alpha)$ is at least $f(\alpha_0)$. Then, excluding α_0 does not decrease the supremum of $f(\alpha)$.

It remains to deal with the case that $\alpha_0 \in \{0, 1\}$. When $\alpha_0 = 0$, we still have Equation (18); the issue is that the left-sided limit (as $\alpha \rightarrow 0^-$) does not contribute to the supremum $\sup_{\alpha \in [0,1]} f(\alpha)$. The previous argument would still go through if we could further show that

$$\lim_{\alpha \rightarrow 0^-} f(\alpha) \leq \lim_{\alpha \rightarrow 0^+} f(\alpha), \quad (19)$$

since (18) and (19) together imply that $\lim_{\alpha \rightarrow 0^+} f(\alpha) \geq f(0)$, so that we can safely ignore the $\alpha_0 = 0$ case. A symmetric argument would deal with the $\alpha_0 = 1$ case via showing $\lim_{\alpha \rightarrow 1^-} f(\alpha) \geq \lim_{\alpha \rightarrow 1^+} f(\alpha)$.

To verify Equation (18), we consider the term $(x_t - \alpha) \cdot \text{sgn}(\alpha - p_t)$ in $f(\alpha)$. When $\alpha_0 \neq p_t$, the term is continuous at α_0 , and contributes equally to both sides of (18). When $\alpha_0 = p_t$, we have

$$\lim_{\alpha \rightarrow \alpha_0^-} (x_t - \alpha) \cdot \text{sgn}(\alpha - p_t) = -(x_t - \alpha) \quad \text{and} \quad \lim_{\alpha \rightarrow \alpha_0^+} (x_t - \alpha) \cdot \text{sgn}(\alpha - p_t) = x_t - \alpha.$$

The average of the two limits is exactly 0, which is equal to $(x_t - \alpha_0) \cdot \text{sgn}(\alpha_0 - p_t)$ since $\alpha_0 = p_t$. By the same token, each term $-(x_t - \alpha) \cdot \text{sgn}(\alpha - \mu)$ also contributes equally to both sides of (18).

To verify Equation (19), we again note that each term $(x_t - \alpha) \cdot \text{sgn}(\alpha - p_t)$ contributes to both sides equally if $p_t \neq 0$. When $p_t = 0$, we have

$$\lim_{\alpha \rightarrow 0^-} (x_t - \alpha) \cdot \text{sgn}(\alpha - p_t) = -x_t \leq 0 \quad \text{and} \quad \lim_{\alpha \rightarrow 0^+} (x_t - \alpha) \cdot \text{sgn}(\alpha - p_t) = x_t \geq 0.$$

For the term $-(x_t - \alpha) \cdot \text{sgn}(\alpha - \mu)$, again, it suffices to verify the case that $\mu = 0$, where

$$\lim_{\alpha \rightarrow 0^-} [-(x_t - \alpha) \cdot \text{sgn}(\alpha - \mu)] = x_t = 0 \quad \text{and} \quad \lim_{\alpha \rightarrow 0^+} [-(x_t - \alpha) \cdot \text{sgn}(\alpha - \mu)] = x_t = 0.$$

In the above, we use the fact that $\mu = 0$ implies that $x_t = 0$ for every $t \in [T]$. This proves Equation (19) and allows us to ignore the $\alpha_0 = 0$ case (and, by symmetry, the $\alpha_0 = 1$ case).

Ignore atypical values for g . Again, we start with the easier case that $\alpha_0 \in S \cap (0, 1)$. In this case, there exists $\delta > 0$ such that: (1) $\alpha_0 - \delta, \alpha_0 + \delta \in [0, 1]$; (2) $[\alpha_0 - \delta, \alpha_0 + \delta] \setminus \{\alpha_0\}$ contains no elements in $S = \{p_1, p_2, \dots, p_T, \mu\}$. Concretely, we can choose

$$\delta = \min \left(\left\{ \frac{1}{2} |\alpha_0 - \beta| : \beta \in S \setminus \{\alpha_0\} \right\} \cup \{\alpha_0, 1 - \alpha_0\} \right) > 0.$$

Then, recalling that $g(\alpha) = \max \{X_-^{(\alpha)} - \alpha N_-^{(\alpha)}, \alpha N_+^{(\alpha)} - X_+^{(\alpha)}\}$, we have

$$X_-^{(\alpha_0)} - \alpha_0 N_-^{(\alpha_0)} = X_-^{(\alpha_0 - \delta)} - (\alpha_0 - \delta) N_-^{(\alpha_0 - \delta)} \leq g(\alpha_0 - \delta).$$

The first step above holds since our choice of δ guarantees $p_t \notin [\alpha_0 - \delta, \alpha_0]$ for each $t \in [T]$, which further implies $p_t < \alpha_0 \iff p_t < \alpha_0 - \delta$. By the same token, we have

$$\alpha_0 N_+^{(\alpha_0)} - X_+^{(\alpha_0)} = (\alpha_0 + \delta) N_+^{(\alpha_0 + \delta)} - X_+^{(\alpha_0 + \delta)} \leq g(\alpha_0 + \delta).$$

Therefore, we have

$$g(\alpha_0) \leq \max\{g(\alpha_0 - \delta), g(\alpha_0 + \delta)\}$$

and $\alpha_0 - \delta, \alpha_0 + \delta \in [0, 1] \setminus S$. This shows that ignoring the case that $\alpha_0 \in S \cap (0, 1)$ does not affect the supremum of $g(\alpha)$.

It remains to deal with the case that $\alpha_0 \in \{0, 1\}$. When $\alpha_0 = 0$, we have

$$g(0) = \max \left\{ X_-^{(0)} - 0 \cdot N_-^{(0)}, 0 \cdot N_+^{(0)} - X_+^{(0)} \right\} = 0,$$

since $X_-^{(0)} = 0$ and $X_+^{(0)} \geq 0$. By symmetry, we also have $g(1) = 0$. Therefore, ignoring these two values of g does not affect the supremum.

Analysis for typical α . It remains to show that $f(\alpha) = 2g(\alpha)$ holds for every $\alpha \in [0, 1] \setminus S$. Note that, in this case, the factors $\text{sgn}(\alpha - p_t)$ and $\text{sgn}(\alpha - \mu)$ do not take value 0 in $f(\alpha)$. When $\alpha < \mu$, $f(\alpha)$ is given by

$$\begin{aligned} f(\alpha) &= \sum_{t \in [T]: p_t < \alpha} (x_t - \alpha) - \sum_{t \in [T]: p_t > \alpha} (x_t - \alpha) + \sum_{t=1}^T (x_t - \alpha) \\ &= \left(X_-^{(\alpha)} - \alpha N_-^{(\alpha)} \right) - \left(X_+^{(\alpha)} - \alpha N_+^{(\alpha)} \right) + \left(X_-^{(\alpha)} + X_+^{(\alpha)} - \alpha T \right) \\ &= 2 \left(X_-^{(\alpha)} - \alpha N_-^{(\alpha)} \right). \end{aligned}$$

Furthermore, since $\alpha < \mu$, we have

$$\alpha N_-^{(\alpha)} + \alpha N_+^{(\alpha)} = \alpha T < \mu T = X_-^{(\alpha)} + X_+^{(\alpha)},$$

which implies $X_-^{(\alpha)} - \alpha N_-^{(\alpha)} > \alpha N_+^{(\alpha)} - X_+^{(\alpha)}$. It follows that

$$f(\alpha) = 2 \left(X_-^{(\alpha)} - \alpha N_-^{(\alpha)} \right) = 2 \max \left\{ X_-^{(\alpha)} - \alpha N_-^{(\alpha)}, \alpha N_+^{(\alpha)} - X_+^{(\alpha)} \right\} = 2g(\alpha).$$

Similarly, when $\alpha > \mu$, we have

$$f(\alpha) = \left(X_-^{(\alpha)} - \alpha N_-^{(\alpha)} \right) - \left(X_+^{(\alpha)} - \alpha N_+^{(\alpha)} \right) - \left(X_-^{(\alpha)} + X_+^{(\alpha)} - \alpha T \right) = 2 \left(\alpha N_+^{(\alpha)} - X_+^{(\alpha)} \right)$$

and

$$\alpha N_+^{(\alpha)} - X_+^{(\alpha)} > X_-^{(\alpha)} - \alpha N_-^{(\alpha)},$$

which also imply $f(\alpha) = 2g(\alpha)$. This shows that $f(\alpha) = 2g(\alpha)$ holds for every $\alpha \in [0, 1] \setminus S$ and completes the proof. \blacksquare

C.2. Proof of Fact 11

Fact 11 For any $x \in \{0, 1\}^T$ and $p \in [0, 1]^T$,

$$\frac{1}{3} \sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha - p_t) \right| \leq \text{stepCE}(x, p) \leq \sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha - p_t) \right|.$$

Proof Fix $x \in \{0, 1\}^T$ and $p \in [0, 1]^T$. Before proving the inequalities, we first show that

$$\left| \sum_{t=1}^T (x_t - p_t) \right| \leq \sup_{\alpha \in [0, 1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha - p_t) \right|, \quad (20)$$

i.e., the total bias is upper bounded by the variant of V-Calibration error.

Upper bound the total bias. If $\sum_{t=1}^T (x_t - p_t) \geq 0$, we have

$$\sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(1 - p_t) = \sum_{t=1}^T (x_t - p_t) - \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t = 1] \geq \sum_{t=1}^T (x_t - p_t),$$

where the last step holds since $(x_t - p_t) \cdot \mathbb{1}[p_t = 1]$ takes value 0 when $p_t \neq 1$, and takes value $x_t - 1 \leq 0$ when $p_t = 1$. It follows that

$$\left| \sum_{t=1}^T (x_t - p_t) \right| = \sum_{t=1}^T (x_t - p_t) \leq \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(1 - p_t) \leq \sup_{\alpha \in [0,1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha - p_t) \right|.$$

Similarly, when $\sum_{t=1}^T (x_t - p_t) < 0$, we have

$$\sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(0 - p_t) = - \sum_{t=1}^T (x_t - p_t) + \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t = 0] \geq - \sum_{t=1}^T (x_t - p_t),$$

which implies

$$\left| \sum_{t=1}^T (x_t - p_t) \right| = - \sum_{t=1}^T (x_t - p_t) \leq \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(0 - p_t) \leq \sup_{\alpha \in [0,1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha - p_t) \right|.$$

The upper bound part. Now we upper bound $\text{stepCE}(x, p)$. It suffices to show that

$$\left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \alpha] \right| \leq \sup_{\beta \in [0,1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\beta - p_t) \right|$$

holds for every $\alpha \in [0, 1]$.

When $\alpha = 1$, the above is exactly given by Equation (20). When $\alpha < 1$, we can always find $\alpha' > \alpha$ such that $\{p_1, p_2, \dots, p_T\} \cap (\alpha, \alpha'] = \emptyset$. Then, for every $t \in [T]$, we have $\mathbb{1}[p_t \leq \alpha] = \mathbb{1}[p_t \leq \alpha']$. Furthermore, since $\alpha' \notin \{p_1, p_2, \dots, p_T\}$, we have $\mathbb{1}[p_t \leq \alpha'] = \mathbb{1}[0 \leq \alpha' - p_t] = \frac{1}{2}(1 + \text{sgn}(\alpha' - p_t))$. It follows that

$$\begin{aligned} \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \alpha] \right| &= \left| \sum_{t=1}^T (x_t - p_t) \cdot \frac{1}{2} (1 + \text{sgn}(\alpha' - p_t)) \right| \\ &\leq \frac{1}{2} \left| \sum_{t=1}^T (x_t - p_t) \right| + \frac{1}{2} \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha' - p_t) \right| \\ &\leq \sup_{\beta \in [0,1]} \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\beta - p_t) \right|. \end{aligned} \quad (\text{Equation (20)})$$

This proves the upper bound on $\text{stepCE}(x, p)$.

The lower bound part. For the other direction, it suffices to prove that

$$\left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha - p_t) \right| \leq 3\text{stepCE}(x, p)$$

holds for every $\alpha \in [0, 1]$. Note that

$$\text{sgn}(\alpha - p_t) = \mathbb{1}[p_t \leq \alpha] - \mathbb{1}[p_t \geq \alpha] = \mathbb{1}[p_t \leq \alpha] - \mathbb{1}[p_t \leq 1] + \mathbb{1}[p_t < \alpha].$$

If $\alpha = 0$, the last indicator $\mathbb{1}[p_t < \alpha]$ is always zero, and we have

$$\left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha - p_t) \right| \leq \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \alpha] \right| + \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq 1] \right| \leq 2\text{stepCE}(x, p).$$

When $\alpha > 0$, we can always find $\alpha' < \alpha$ such that $\{p_1, p_2, \dots, p_T\} \cap (\alpha', \alpha) = \emptyset$. Then, $\mathbb{1}[p_t < \alpha] = \mathbb{1}[p_t \leq \alpha']$ holds for every $t \in [T]$, and it follows that

$$\begin{aligned} & \left| \sum_{t=1}^T (x_t - p_t) \cdot \text{sgn}(\alpha - p_t) \right| \\ &= \left| \sum_{t=1}^T (x_t - p_t) \cdot (\mathbb{1}[p_t \leq \alpha] - \mathbb{1}[p_t \leq 1] + \mathbb{1}[p_t \leq \alpha']) \right| \\ &\leq \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \alpha] \right| + \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq 1] \right| + \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \leq \alpha'] \right| \\ &\leq 3\text{stepCE}(x, p). \end{aligned}$$

This proves the lower bound on $\text{stepCE}(x, p)$. ■

C.3. Step Calibration and Its Subsampled Variant

Lemma 24 For any $x \in \{0, 1\}^T$ and $p \in [0, 1]^T$,

$$\frac{1}{2}\text{stepCE}(x, p) \leq \text{stepCE}^{\text{sub}}(x, p) \leq \frac{1}{2}\text{stepCE}(x, p) + O(\sqrt{T}).$$

Proof Fix $x \in \{0, 1\}^T$ and $p \in [0, 1]^T$. For the lower bound part, suppose that the supremum in $\text{stepCE}(x, p)$ is achieved at α^* ,⁵ i.e.,

$$\text{stepCE}(x, p) = \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \in [0, \alpha^*]] \right|.$$

5. α^* is well defined, as the term in the supremum takes at most $T + 1$ different values over all $\alpha \in [0, 1]$.

Then, we have

$$\begin{aligned}
 \text{stepCE}^{\text{sub}}(x, p) &\geq \mathbb{E}_{y \sim \text{Unif}(\{0,1\}^T)} \left[\sum_{t=1}^T y_t \cdot (x_t - p_t) \cdot \mathbb{1}[p_t \in [0, \alpha^*]] \right] \\
 &\geq \left| \mathbb{E}_{y \sim \text{Unif}(\{0,1\}^T)} \left[\sum_{t=1}^T y_t \cdot (x_t - p_t) \cdot \mathbb{1}[p_t \in [0, \alpha^*]] \right] \right| \quad (\text{convexity of } x \mapsto |x|) \\
 &= \frac{1}{2} \left| \sum_{t=1}^T (x_t - p_t) \cdot \mathbb{1}[p_t \in [0, \alpha^*]] \right| = \frac{1}{2} \text{stepCE}(x, p).
 \end{aligned}$$

For the upper bound, we assume without loss of generality that $p_1 \leq p_2 \leq \dots \leq p_T$, since both stepCE and $\text{stepCE}^{\text{sub}}$ are invariant to the reordering of the entries (in x and p simultaneously). Let $S := \{t \in [T-1] : p_t < p_{t+1}\} \cup \{T\}$. For each $t \in \{0, 1, \dots, T\}$, define $A_t := \sum_{i=1}^t (x_i - p_i)$. Over the randomness in $y \sim \text{Unif}(\{0, 1\}^T)$, define

$$X_t := \sum_{i=1}^t (y_i - 1/2) \cdot (x_i - p_i).$$

Then, $\text{stepCE}(x, p)$ and $\text{stepCE}^{\text{sub}}(x, p)$ can be simplified into

$$\text{stepCE}(x, p) = \max_{t \in S} |A_t|,$$

and

$$\begin{aligned}
 \text{stepCE}^{\text{sub}}(x, p) &= \mathbb{E}_y \left[\max_{t \in S} |A_t/2 + X_t| \right] \\
 &\leq \frac{1}{2} \max_{t \in S} |A_t| + \mathbb{E}_y \left[\max_{t \in [T]} |X_t| \right] \\
 &= \frac{1}{2} \text{stepCE}(x, p) + \mathbb{E}_y \left[\max_{t \in [T]} |X_t| \right].
 \end{aligned}$$

Therefore, it remains to control the term $\mathbb{E}_y [\max_{t \in [T]} |X_t|]$ by $O(\sqrt{T})$. Note that $(X_t)_{t=0}^T$ is a martingale in which each displacement $(X_t - X_{t-1}) \mid X_{t-1}$ has a variance of $(x_t - p_t)^2/4 \leq 1/4$. Kolmogorov's inequality implies that, for every $\tau > 0$,

$$\Pr \left[\max_{t \in [T]} |X_t| \geq \tau \right] \leq \frac{T/4}{\tau^2}.$$

It follows that

$$\begin{aligned}
 \mathbb{E} \left[\max_{t \in [T]} |X_t| \right] &= \int_0^{+\infty} \Pr \left[\max_{t \in [T]} |X_t| \geq \tau \right] d\tau \\
 &\leq \int_0^{+\infty} \min \left\{ \frac{T}{4\tau^2}, 1 \right\} d\tau = O(\sqrt{T}).
 \end{aligned}$$

■