

741 A Proof of Proposition 1

742 The goal of Proposition 1 is to bound the approximation error
743 of the linear surrogate model on the fairness constraint man-
744 ifold. Specifically, we show that with high probability, the
745 linear proxy value $\mathbf{w}_\phi^\top \mathbf{h}$ remains close to the reference value
746 ϕ_{proxy} when \mathbf{h} is sampled from the level set H_ϕ .

747 *Proof.* The proof relies on establishing an integral identity
748 along a path connecting the reference point \mathbf{h}_ϕ to a target
749 point \mathbf{h} on the manifold, and then bounding the error us-
750 ing statistical properties of the domain geometry and gradient
751 misalignment.

752 **Step 1: Establishing the Gradient Integral Identity.** Let
753 $H_\phi = \{\mathbf{h} \in \mathcal{H} : \varphi(\mathbf{h}) = \phi\}$ be the constraint set defined in
754 the proposition. Consider an arbitrary distribution shift vector
755 $\mathbf{h} \in H_\phi$ and the fixed reference point $\mathbf{h}_\phi \in H_\phi$. We define a
756 straight-line path $l(t)$ connecting them:

$$l(t) = \mathbf{h}_\phi + t(\mathbf{h} - \mathbf{h}_\phi), \quad t \in [0, 1].$$

757 According to the *Fundamental Theorem of Calculus for Line*
758 *Integrals*, the difference in the fairness evaluation φ between
759 the endpoints can be expressed as the integral of the gradient
760 along the path:

$$\begin{aligned} \varphi(\mathbf{h}) - \varphi(\mathbf{h}_\phi) &= \int_0^1 \nabla \varphi(l(t))^\top \frac{d}{dt} l(t) dt \\ &= \int_0^1 \nabla \varphi(l(t))^\top (\mathbf{h} - \mathbf{h}_\phi) dt. \end{aligned}$$

761 Since both \mathbf{h} and \mathbf{h}_ϕ lie on the same level set H_ϕ , we have
762 $\varphi(\mathbf{h}) = \phi$ and $\varphi(\mathbf{h}_\phi) = \phi$. Consequently, the left-hand side
763 is exactly zero, yielding the identity:

$$0 = \int_0^1 \nabla \varphi(l(t))^\top (\mathbf{h} - \mathbf{h}_\phi) dt. \quad (10)$$

764 **Step 2: Introducing the Linear Surrogate and Error**
765 **Bound.** We define the approximation error of the linear sur-
766rogate \mathbf{w}_ϕ relative to the proxy reference $\phi_{\text{proxy}} := \mathbf{w}_\phi^\top \mathbf{h}_\phi$
767 as:

$$Y(\mathbf{h}) := |\mathbf{w}_\phi^\top \mathbf{h} - \phi_{\text{proxy}}| = |\mathbf{w}_\phi^\top (\mathbf{h} - \mathbf{h}_\phi)|.$$

768 Using the property that $\int_0^1 dt = 1$, we rewrite the linear term
769 $\mathbf{w}_\phi^\top (\mathbf{h} - \mathbf{h}_\phi)$ as an integral $\int_0^1 \mathbf{w}_\phi^\top (\mathbf{h} - \mathbf{h}_\phi) dt$. Subtracting the
770 gradient identity (10) (which equals 0) from this expression
771 allows us to isolate the misalignment between the surrogate
772 and the true gradient:

$$\begin{aligned} \mathbf{w}_\phi^\top (\mathbf{h} - \mathbf{h}_\phi) - 0 &= \int_0^1 \mathbf{w}_\phi^\top (\mathbf{h} - \mathbf{h}_\phi) dt \\ &\quad - \int_0^1 \nabla \varphi(l(t))^\top (\mathbf{h} - \mathbf{h}_\phi) dt \\ &= \int_0^1 (\mathbf{w}_\phi - \nabla \varphi(l(t)))^\top (\mathbf{h} - \mathbf{h}_\phi) dt. \end{aligned}$$

Taking the absolute value and applying the *Cauchy-Schwarz* 773
inequality to the integrand: 774

$$\begin{aligned} Y(\mathbf{h}) &= \left| \int_0^1 (\mathbf{w}_\phi - \nabla \varphi(l(t)))^\top (\mathbf{h} - \mathbf{h}_\phi) dt \right| \\ &\leq \int_0^1 \|\mathbf{w}_\phi - \nabla \varphi(l(t))\| \cdot \|\mathbf{h} - \mathbf{h}_\phi\| dt. \end{aligned} \quad (11)$$

Since the term $\|\mathbf{h} - \mathbf{h}_\phi\|$ is independent of t , we can separate 775
the geometric distance from the gradient deviation: 776

$$Y(\mathbf{h}) \leq \|\mathbf{h} - \mathbf{h}_\phi\| \cdot \underbrace{\int_0^1 \|\mathbf{w}_\phi - \nabla \varphi(\mathbf{h}_\phi + t(\mathbf{h} - \mathbf{h}_\phi))\| dt}_{\text{Gradient Misalignment along path}}. \quad (12)$$

Step 3: Bounding the Expected Error. We now consider 777
the expectation over the distribution of \mathbf{h} on the manifold H_ϕ . 778
Taking the expectation $\mathbb{E}_{\mathbf{h} \in H_\phi}$ on both sides of (12) and ap- 779
plying the *Cauchy-Schwarz inequality for random variables* 780
($\mathbb{E}[XY] \leq \sqrt{\mathbb{E}[X^2]} \sqrt{\mathbb{E}[Y^2]}$). Let the path integral residual 781
be denoted by $R = \int_0^1 \|\mathbf{w}_\phi - \nabla \varphi(l(t))\| dt$. Then we have: 782

$$\mathbb{E}[Y(\mathbf{h})] \leq \sqrt{\mathbb{E}[\|\mathbf{h} - \mathbf{h}_\phi\|^2]} \sqrt{\mathbb{E}[R^2]}.$$

Comparing this with the definitions given in Proposition 1, 783
we identify the two terms on the right-hand side as A and B : 784

- $A = \sqrt{\mathbb{E}[\|\mathbf{h} - \mathbf{h}_\phi\|^2]}$ represents the effective diameter. 785
- $B = \sqrt{\mathbb{E}\left[\left(\int_0^1 \|\mathbf{w}_\phi - \nabla \varphi(\mathbf{h}_\phi + t(\mathbf{h} - \mathbf{h}_\phi))\| dt\right)^2\right]}$ 786
quantifies the gradient misalignment. 787

Thus, the bound on the expected error is: 788

$$\mathbb{E}[Y(\mathbf{h})] \leq A \cdot B. \quad (13)$$

Step 4: Deriving the Probabilistic Concentration Bound. 789
Finally, we apply *Markov's Inequality* to the non-negative 790
random variable $Y(\mathbf{h})$. For any threshold $\delta > 0$: 791

$$P(Y(\mathbf{h}) \geq \delta) \leq \frac{\mathbb{E}[Y(\mathbf{h})]}{\delta} \leq \frac{A \cdot B}{\delta}.$$

Equivalently, the probability that the error remains within δ 792
is: 793

$$P(Y(\mathbf{h}) \leq \delta) \geq 1 - \frac{A \cdot B}{\delta}.$$

To ensure a confidence level of at least $1 - \beta$ (where $\beta \in$ 794
 $(0, 1)$), we set the failure probability bound equal to β , which 795
implies $\frac{A \cdot B}{\delta} = \beta$, or $\delta = \frac{A \cdot B}{\beta}$. Substituting this δ back yields 796
the final result: 797

$$P\left(|\mathbf{w}_\phi^\top \mathbf{h} - \phi_{\text{proxy}}| \leq \frac{A \cdot B}{\beta}\right) \geq 1 - \beta. \quad \square \quad 798$$

799 B Details of the Experiment in RQ1

In this section, we provide a detailed description of the exper- 800
imental setup for RQ1, designed to verify the estimation accu- 801
racy of CoRe against an analytically solvable ground truth. 802

803 **B.1 Dataset Construction**

804 To calculate the exact consistency radius, we constructed a
 805 synthetic environment based on a trinomial distribution vari-
 806 able $X \in \{0, 1, 2\}$. The dataset was generated through the
 807 following steps:

- 808 1. **Performance metric function (M):** We randomly sam-
 809 pled 10 performance vectors $\mathbf{m} \in [0, 1]^3$, where each
 810 component m_i represents the performance of the model
 811 on the data category i .
- 812 2. **Audit Distributions (P):** We randomly sampled 10
 813 probability vectors $\mathbf{p} \in \mathbb{R}^3$ from a Dirichlet distribution,
 814 ensuring $\sum p_i = 1$. This represents the initial audit data
 815 distribution.

816 The Cartesian product of these sets yields 100 distinct au-
 817 diting configurations (pairs of M and P).

818 **B.2 Feasible Threshold Selection**

819 A critical parameter in our problem is the fairness tolerance
 820 threshold ϵ . Randomly choosing ϵ is problematic because cer-
 821 tain fairness levels may be mathematically impossible for a
 822 given model M , regardless of the distribution shift.

823 To ensure valid and comprehensive testing, we first derived
 824 the **feasible range** of fairness values for each model. The
 825 theoretical boundaries of achievable fairness φ on the prob-
 826 ability simplex are determined by the roots of the following
 827 quadratic equation:

$$\frac{-b \pm \sqrt{b^2 - 4ac}}{2a}, \quad (14)$$

828 where the coefficients depend on the model performance vec-
 829 tor \mathbf{m} : $a = -1$, $b = m_0 + m_1 - 2m_2$, and $c = (m_0 -$
 830 $m_2)(m_2 - m_1)$.

831 For each of the 100 configurations, we calculated this
 832 feasible range and then uniformly sampled 10 values of ϵ
 833 within it. This sampling strategy generates a total of 1,000
 834 test instances, covering the full spectrum of auditing sce-
 835 narios—from loose constraints (easy to satisfy) to extremely
 836 tight constraints (borderline feasibility).

837 **B.3 Microscopic Analysis: Accuracy Decay on
 838 Corner Cases**

839 While the main text reports the aggregated error metrics, we
 840 provide here a fine-grained analysis of how the estimation be-
 841 havior changes with the strictness of the fairness constraint.

842 Figure 5 visualizes the relationship between the ground-
 843 truth radius (Solid Line) and the radius estimated by CoRe
 844 (Points) for a representative (M, P) pair. The x-axis repre-
 845 sents the chosen tolerance ϵ .

846 **Observation:** The results show a clear pattern: CoRe
 847 achieves high precision in the central region of the feasible
 848 range. However, the estimation error tends to increase as ϵ
 849 approaches the boundaries (indicated by the shaded regions
 850 in the figure).

851 **Analysis:** This phenomenon reveals the geometric char-
 852 acteristics of the consistency radius problem. The boundary
 853 regions correspond to extreme fairness values that are only
 854 achievable by highly skewed probability distributions. In

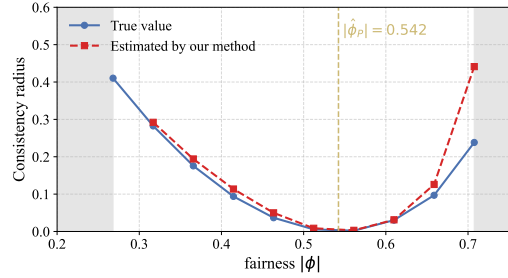


Figure 5: Visual comparison of the Ground-Truth Consistency Ra-
 dius vs. CoRe Estimated Radius for a single configuration (M, P)
 across the feasible ϵ spectrum. The shaded regions denote the
 boundaries of the feasible fairness range.

855 these "corner cases," the geometry of the fairness constraint
 856 manifold becomes highly curved. Since CoRe relies on a lin-
 857 ear surrogate vector \mathbf{w}^* to approximate this non-convex con-
 858 straint locally, the approximation error naturally grows when
 859 the local curvature is sharp.

860 Nevertheless, this analysis confirms that for the vast major-
 861 ity of practical auditing scenarios (the non-extreme regions),
 862 our convex relaxation method provides a reliable and tight
 863 estimation. The deviations at the boundaries also explain the
 864 outliers observed in the aggregated error statistics.

865 **C Details of the Experiment in RQ2**

866 This section details the model architectures, training proto-
 867 cols, and fairness intervention strategies for the three real-
 868 world datasets used in RQ2.

869 **C.1 Dataset and Model Settings**

870 **MovieLens-1M (Interaction Data).** The dataset contains
 871 approximately 6,000 users and 4,000 movies. We prepro-
 872 cessed the data by removing users and items with fewer than
 873 five interactions and re-indexing the remaining IDs. We em-
 874 ployed a leave-one-out splitting strategy: for each user, the
 875 most recent interaction was used for testing, while the pre-
 876 ceding interactions formed the training set.

877 We trained a LightGCN model, a state-of-the-art collabo-
 878 rative filtering method. The model consists of user and item
 879 embedding layers with a dimension of 64. Following the stan-
 880 dard LightGCN implementation, we simplified the graph con-
 881 volution by removing feature transformation and nonlinear
 882 activation. The model was trained using the Adam optimizer
 883 with a learning rate of 0.005 and a batch size of 1024. We
 884 used the Bayesian Personalized Ranking (BPR) loss, sam-
 885 pling one negative item for each positive instance. Early
 886 stopping was triggered if the validation NDCG@10 did not
 887 improve for 5 consecutive epochs.

888 **Adult (Structured Data).** We trained a Multilayer
 889 Perceptron (MLP) on the Adult dataset to predict in-
 890 come levels ($\leq 50K$ vs. $> 50K$). Missing values were
 891 removed, and categorical features were encoded using
 892 LabelEncoder, while numerical features were standard-
 893 ized using StandardScaler. The dataset was split into

894 80% training and 20% testing, with 20% of the training data
895 reserved for validation.

896 The MLP consists of two hidden layers with 64 and
897 32 units, respectively, using ReLU activation and Dropout
898 (rate=0.2). The model was trained using the Adam optimizer
899 (learning rate 1×10^{-3}) and binary cross-entropy loss. We
900 applied early stopping with a patience of 5 epochs based on
901 validation loss.

902 **CelebA (Image Data).** For the image classification task,
903 we used the CelebA dataset to predict the "Smiling" attribute.
904 We center-cropped the images to 178×178 and resized them
905 to 128×128 . The pixel values were normalized to the range
906 $[0, 1]$. We utilized a ResNet-18 architecture. To simulate a
907 standard auditing scenario where models are often trained
908 from scratch on specific datasets, we trained the ResNet-18
909 without pre-trained weights. The model was optimized using
910 Adam with a learning rate of 1×10^{-4} and a batch size of 64
911 for 20 epochs. The sensitive attribute was defined as "Male"
912 (Gender).

913 D Details of the Experiment for Addressing 914 practitioners' Queries

915 In this section, we provide the detailed experimental setup
916 and data generation processes for the exploratory analyses
917 presented in Section 6.

918 D.1 Simulation Setup for Radius and Fairness 919 Correlation

920 To analyze the relationship between the consistency radius
921 and the difference between the fairness evaluation result
922 $|\varphi(D)|$ and the threshold ϵ , we conducted a simulation study
923 using synthetic data.

924 **Data Generation.** We randomly constructed 3,000 syn-
925 thetic datasets. For each dataset, we generated a metric vector
926 \mathbf{m} and a group attribute vector \mathbf{g} as follows:

- 927 • **Sample Size:** Each synthetic dataset consists of $N =$
928 100 samples.
- 929 • **Metric Values (\mathbf{m}):** The metric values (e.g., predic-
930 tion scores or losses) were sampled independently from
931 a Uniform distribution over the interval $[0, 1]$, i.e., $m_i \sim$
932 $\mathcal{U}(0, 1)$.
- 933 • **Group Attribute (\mathbf{g}):** We initialized binary group la-
934 bels $\{0, 1\}$ randomly. To eliminate the confounding ef-
935 fect of group imbalance, we enforced a strictly balanced
936 group ratio. Specifically, we adjusted the counts such
937 that exactly 50 samples belonged to the protected group
938 ($g = 0$) and 50 samples to the non-protected group
939 ($g = 1$).

940 **Correlation Analysis.** For each dataset, we computed the
941 consistency radius given a fixed fairness threshold ϵ (e.g.,
942 $\epsilon = 0$ or $\epsilon = 0.1$). We then calculated the absolute difference
943 $|\hat{\phi}_P| - \epsilon$ and quantified its correlation with the consistency
944 radius using the Kendall rank correlation coefficient (τ), uti-
945 lizing the `scipy.stats` library.

946 D.2 Simulation Setup for Dataset Merging

947 To investigate the impact of merging audit datasets on the
948 consistency radius, we analyzed pairs of datasets based on
949 the assumption of identical group proportions.

950 **Dataset Construction.** We constructed 3,000 pairs of
951 datasets $(\mathcal{D}_1, \mathcal{D}_2)$. Each dataset in the pair was generated in-
952 dependently using the same procedure described above ($N =$
953 100, $\mathbf{m} \sim \mathcal{U}(0, 1)$, and strictly balanced groups). Since both
954 \mathcal{D}_1 and \mathcal{D}_2 are strictly balanced (50% group 0, 50% group
955 1), merging them naturally satisfies the requirement that the
956 group proportions remain consistent across the individual and
957 merged datasets.

958 **Evaluation.** For each pair, we created a merged dataset
959 $\mathcal{D}_{\text{concat}} = \mathcal{D}_1 \cup \mathcal{D}_2$ by concatenating the metric and group
960 vectors. We then computed the consistency radius for \mathcal{D}_1 ,
961 \mathcal{D}_2 , and $\mathcal{D}_{\text{concat}}$ separately under the same threshold ϵ . We
962 compared the radius of the merged dataset against the mini-
963 mum and maximum radii of the individual datasets to observe
964 the trends shown in Figure 4b.

965 E CoRe Algorithm Implementation Details

966 The CoRe algorithm was implemented in Python. For the
967 inner-level convex optimization problem (finding the optimal
968 shift \mathbf{h} under a linear constraint), we utilized the CVXPY li-
969 brary with the ECOS solver, which is efficient for handling
970 entropy-based objectives and linear constraints.

971 The binary search in the outer loop was configured with a
972 convergence tolerance of $\text{tol} = 1 \times 10^{-4}$ and a maximum
973 of 30 iterations. The regularization parameter λ for stabiliz-
974 ing group scales was set to 1.0 by default. All experiments
975 were conducted on a Linux server equipped with an Intel
976 Xeon Gold 6240 CPU and an NVIDIA GeForce RTX 3090
977 GPU. The estimation of the consistency radius typically con-
978 verges within seconds for audit datasets of size $n = 1,000$,
979 confirming the computational efficiency of the method.