
Adaptive Stratified Active Statistical Inference

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Active Statistical Inference (ASI) leverages machine learning predictions to guide
2 label acquisition and improve statistical inference under limited labeling budgets,
3 but lacks the nuanced understanding of model reliability varying across the data
4 manifold, leading to inefficient budget allocation. To address this, we propose
5 Adaptive Stratified Active Statistical Inference (AdaStrat-ASI), a framework that
6 replaces ASI’s global sampling rule with stratum-specific local policies, designed
7 via a model-guided ‘scouting’ phase. We provide theoretical results showing that
8 AdaStrat-ASI achieves strictly lower asymptotic variance compared to ASI while
9 preserving the inferential guarantees of ASI. We verify our theoretical findings
10 through empirical results on real-world datasets demonstrate that AdaStrat-ASI
11 yields tighter confidence intervals than existing baselines under the same labeling
12 budget.

13 1 Introduction

14 High-quality labels data are often a bottleneck in data-driven research. Consequently, practitioners
15 are increasingly adopting machine learning predictions as inexpensive proxies for missing labels, in
16 applications ranging from remote sensing [28, 58, 51], proteomics [29], and electoral systems [59],
17 enabling large-scale measurement possible in settings where exhaustive human annotation would be
18 prohibitively expensive. Yet prediction is not inference: even highly accurate predictions do not, by
19 themselves, provide statistical guarantees required for decision-making. Confidence intervals and
20 hypothesis tests still require access to ground truth labels.

21 Active Statistical Inference (ASI) [65] addresses this gap by using predictions not as a substitute for
22 labels but as a guide for acquiring them. It samples labels according to a single *global* uncertainty rule
23 and then corrects the estimator using an inverse propensity weighting. This yields statistically valid
24 inference for a broad class of M-estimation targets under tight budgets. The limitation is that ASI
25 asks a global question of a local object. Its sampling rule normalizes uncertainty over the entire data
26 pool, so all points compete for the same labeling budget on a single global scale. This is natural when
27 predictive reliability is roughly homogeneous. In practice, this is seldom the case – a model may be
28 accurate on one subpopulation, biased in another, calibrated in one region and confidently wrong in
29 another [4, 17]. LLM annotators [34, 64, 49, 56] exhibit the same pattern: they may perform well
30 on well-specified or easily verifiable tasks [35, 12], yet become unreliable on open-ended, weakly
31 grounded or demographically sensitive inputs [57, 26].

32 *Can ASI be made more efficient by exploiting region-specific reliability differential without*
33 *sacrificing validity?*

34 We answer this question with AdaStrat-ASI, an adaptive stratified framework for active statistical
35 inference. The key idea is to make reliability local: given a partition of data into strata – defined
36 by domain knowledge, observed covariate groupings, or model-induced structure – AdaStrat-ASI
37 replaces ASI’s single global sampling rule with stratum-specific local policies. While labels are still

38 acquired actively within each stratum, the total budget allocation is no longer spread according to one
39 global uncertainty rule. Instead, AdaStrat-ASI allocates more labels to strata where the predictor is
40 locally less useful for inference, and fewer labels to strata where the predictor is locally stable. In an
41 oracle design, these label budgets would be assigned to strata based according to their contribution
42 to asymptotic variance of the final estimator. But this requires access to true labels to calculate.
43 AdaStrat-ASI resolves this cold-start problem through a label-free scouting stage. Before querying
44 any labels, it uses the black-box predictor on the unlabelled data to estimate a proxy difficulty score
45 for each stratum, measuring the local variability of proxy loss gradients which serve as a model-based
46 signal for how much the stratum may affect the final estimator. This results in a Neyman-style
47 allocation (Section 4.1). After this allocation step, a local ASI sampling rule is applied within
48 each stratum. Thus, our method separates two decisions which are conflated in global ASI: *which*
49 *points* to label within a region, and *how much budget* that region should receive. We summarize our
50 contributions as follows:

- 51 • **Method.** We propose Adaptive Stratified Active Statistical Inference (AdaStrat-ASI), a stratified
52 extension of ASI that replaces one global sampling policy with stratum-specific local policies
53 under a fixed labeling budget to improve statistical efficiency.
- 54 • **Validity.** We prove that our method is statistically valid: under standard regularity conditions, it is
55 asymptotically normal and supports valid confidence intervals (Theorem 5.3).
- 56 • **Efficiency.** We characterize an oracle optimal stratum allocation (Proposition 5.4), relate stratum
57 difficulty to predictor error (Lemma 5.5), and provide a sufficient condition under which AdaStrat-
58 ASI provably attains a strictly smaller asymptotic parameter covariance than ASI (Theorem 5.6),
59 clarifying when stratification yields a principled efficiency gain.
- 60 • **Empirical validation.** Our results empirically validate our theoretical claims, demonstrating
61 AdaStrat-ASI’s narrower confidence intervals under identical labeling budgets compared against
62 standard baselines on synthetic and real-world datasets. The source code for this project will be
63 made publicly available upon acceptance of the manuscript.

64 **Technical Overview.** While stratification is standard, our focus is *ML-integrated inference*. Any
65 inference framework that treats the predictor as a globally homogeneous oracle is theoretically
66 misaligned with reality [55, 5]. AdaStrat-ASI resolves this misalignment by treating the latent
67 structure of model performance in strata as the central theme in the sampling design by proposing
68 a novel estimator (Eq. 3). We provide the first theoretical analysis within the ASI framework that
69 explicitly links model error to inference efficiency—showing how it enters the asymptotic covariance
70 (Lemma C.3 in Appendix) and drives the Neyman-style oracle label allocation across strata (Lemma
71 5.5), offering a nuanced understanding of *why* and *where* ML-guided inference succeeds or fails.
72 Further, since optimal label allocation depends on unknown labels, we introduce a novel, *label-free*
73 *scouting* mechanism that uses cheaply available proxy predictions $f(X)$ to estimate per-stratum
74 “difficulty” and allocate a fixed labeling budget *without peeking* at the ground truth, effectively solving
75 the cold-start problem for stratified active learning. Taken together, stratification is the mathematically
76 natural response to the heterogeneity inherent in black-box predictors and is a principled route to
77 improved statistical efficiency (Theorem 5.6).

78 2 Related Works

79 **ML-enabled Data Annotation and Inference.** A burgeoning body of literature exists on inference
80 from adaptively collected data, encompassing works such as Kato et al. [31], Zhang et al. [61], Cook
81 et al. [14], Lin et al. [37]. Additionally, adaptive experimental design, including approaches like Hahn
82 et al. [23], Chandak et al. [8], has garnered significant attention. Under the regime of black-box
83 predictions being well calibrated, ASI aligns most closely with our pursuits. It positions itself uniquely
84 by adaptively leveraging modern, black-box machine learning uncertainties for data collection while
85 adhering to budget constraints for statistical inference. ASI shares similarities with prediction-
86 powered inference (PPI) and other recent works on inference with ML predictions [2, 3, 44, 42].
87 Notably, ASI considers the same class of inferential targets studied in Angelopoulos et al. [2, 3], Zrníc
88 and Candès [66]. Recently, Li et al. propose Robust Sampling in ASI to mitigate the risk of unreliable
89 uncertainty estimates by interpolating between uniform and uncertainty-based sampling, ensuring
90 performance no worse than uniform sampling under model misspecification. In parallel, LLMs have
91 shown great potential in handling text-annotation tasks without prior task-specific training, sometimes
92 even outperforming crowd workers [21, 63, 38, 10, 33]. Inspired by active learning’s role in LLM
93 Annotation [62, 41], Gligoric et al. extended ASI to leverage LLM-provided confidence to guide

94 targeted human annotation and produce statistically valid estimates and confidence intervals under
 95 the confidence-driven inference (CCI) framework.

96 **Stratified Sampling and Adaptivity.** Standard variance reduction techniques like Stratified Sampling
 97 can be viewed as a statistical ‘divide-and-conquer’ strategy: rather than estimating parameters over a
 98 complex, heterogeneous data manifold, one partitions the space into homogeneous strata and solves
 99 simpler local estimation problems [13]. Although related in goal, Fisch et al. combined the variance
 100 reduction benefits of Stratified Sampling for PPI. However AdaStrat-ASI is fundamentally different
 101 in approach and idea. AdaStrat-ASI, is an extension of ASI i.e. our goal is to strategically collect
 102 data to yield powerful inference, leveraging known structure in data for variance reduction. Further,
 103 the efficiency of method hinges on label allocation inspired by Neyman allocation, which dictates
 104 that the optimal sample distribution is proportional to true variance [45]. However, the internal
 105 variance for the task of statistical inference is an unknown a priori. Our work addresses this cold start
 106 problem by creating an approximation of the optimality criterion using a two-stage adaptive design,
 107 reminiscent of Explore-then-Commit (EtC) strategies in bandit literature [20] or stratified Monte
 108 Carlo methods [6](discussed in detail in Section 4).

109 **Region-based Active Learning.** Zhang and Chaudhuri demonstrated that in agnostic settings,
 110 absolute disagreement [24] is overly conservative. They introduced confidence-rated predictors to
 111 study disagreement-based learning in sub-regions. This idea was further extended to contextual
 112 bandits through the lens of “Sub-Regions of Disagreement” [43]: treating disagreement as a binary,
 113 global property masks opportunities for exploitation; instead, the input space should be decomposed
 114 into sub-regions where the optimal policy is locally identifiable, effectively isolating the “hard”
 115 regions of the space. Our approach also shares conceptual lineage with Region-Based Active
 116 Learning (RBAL) [15], which adapts the sampling strategy based on the difficulty of learning within
 117 specific sub-regions of the feature space. Similar to RBAL, AdaStrat-ASI acknowledges that model
 118 performance is not uniform; however, while RBAL focuses on optimizing a hypothesis h during
 119 training, our work focuses on the statistical inference of population parameters.

120 3 Preliminaries

121 **Problem Setup.** We observe unlabeled samples $X_1, \dots, X_n \sim \mathbb{P}_X$, with corresponding labels Y_i
 122 initially unobserved, and aim to estimate a population parameter θ^* defined as the minimizer of an
 123 expected loss:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{(X,Y) \sim \mathbb{P}} [\ell_{\theta}(X, Y)],$$

124 for a convex loss function ℓ_{θ} . Here, $(X, Y) \sim \mathbb{P} = \mathbb{P}_X \times \mathbb{P}_{Y|X}$ denotes a generic feature-label pair
 125 drawn from the underlying data distribution. This formulation captures a broad class of M -estimation
 126 problems, including mean estimation, quantiles, and regression coefficients.

127 **Active Statistical Estimator.** We assume access to an off-the-shelf predictive model $f(X)$ for
 128 labels $Y \in \mathbb{R}$ given covariates $X \in \mathcal{X}$. We design a *sampling rule* $\pi : \mathcal{X} \rightarrow [0, 1]$ and collect label
 129 Y_i with probability $\pi(X_i)$. Our goal is to perform inference under a strict labeling budget ($n_b \ll n$).
 130 To overcome this, we use $f(X)$ as a strategic guide to select the most informative points to label.
 131 We denote the number of empirically collected labels as n_{lab} . The intuition is to model $\pi(\cdot)$ to
 132 have a high(low) value for instances where f is uncertain(certain). At a technical level, the rule is
 133 derived from a measure of model-output uncertainty $u(x)$ (computed from f). The definition of
 134 uncertainty in classification and regression is taken from Section 4 of Zrnić and Candès [65]. Let
 135 $\xi_i \sim \text{Bernoulli}(\pi(X_i))$ indicate whether point i is labeled, so that $n_{\text{lab}} = \sum_{i=1}^n \xi_i$. To stay under a
 136 budget n_b , we set $\pi(x) = \frac{u(x)}{\mathbb{E}[u(X)]} \cdot \frac{n_b}{n}$, which yields $\mathbb{E}[n_{\text{lab}}] = \mathbb{E}[\pi(X)] \cdot n \leq n_b$.

137 The *active estimator* i.e. Equation 1 from Zrnić and Candès [65] is defined as $\hat{\theta}^{\pi} = \arg \min_{\theta} L^{\pi}(\theta)$,
 138 where our *active empirical risk* is as follows,

$$L^{\pi}(\theta) := \frac{1}{n} \sum_{i=1}^n \left[\ell_{\theta,i}^f + \left(\ell_{\theta,i} - \ell_{\theta,i}^f \right) \frac{\xi_i}{\pi(X_i)} \right] \quad (1)$$

139 To simplify notation, let $L(\theta) = \mathbb{E}[\ell_{\theta}(X, Y)]$, $\ell_{\theta,i} = \ell_{\theta}(X_i, Y_i)$, and $\ell_{\theta,i}^f = \ell_{\theta}(X_i, f(X_i))$. It is
 140 evident that the estimator exemplifies augmented inverse propensity weighting (AIPW) estimation

141 [50]. Further, when the sampling rule is uniform, i.e., $\pi(x) = \frac{n_b}{n}$ for all x , the estimator imitates
 142 the prediction-powered estimator [2]. Classical baselines such as $\hat{\theta}^{\text{noML}} := \arg \min_{\theta} \frac{1}{n_b} \sum_{i=1}^{n_b} \ell_{\theta, i}$
 143 utilize pure labels.

144 In many data-efficient learning and inference settings, a complementary goal is to ensure that
 145 labeled data are well-distributed across the underlying data manifold [7, 52, 16]. Motivated by this
 146 perspective, we quantify the geometric spread of the labeled set via the diameter of the induced
 147 manifold. For a labeled dataset $\mathcal{X}_{\text{label}} \subset \mathbb{R}^d$ we define diameter of manifold $\text{Diam}(\cdot)$ as,

$$\text{Diam}(\mathcal{X}_{\text{lab}}) = \max_{\substack{x_i, x_j \in \mathcal{X}_{\text{lab}} \\ i \neq j}} \|x_i - x_j\| \quad (2)$$

148 The diameter captures the maximal geometric separation among labeled points and serves as a proxy
 149 for the spread of the labeled data manifold.

150 **Stratification and Model Performance.** In accordance with previous works [19], we assume that
 151 the input space X is partitioned in advance into K non-empty, mutually exclusive, and exhaustive
 152 strata $\mathcal{A} = (\mathcal{A}_1, \dots, \mathcal{A}_K)$, where K is a finite integer. Let $n_k = |\mathcal{A}_k|$ denote the size of stratum
 153 k , such that $\sum_{k=1}^K n_k = n$. Let w_k represent the probability mass of stratum k in the finite
 154 population, where $w_k = \mathbb{P}(X \in \mathcal{A}_k)$. Naturally, w_k becomes $\frac{n_k}{n}$. A major premise of our work is
 155 that the performance of f is not uniform across the input space [55, 25, 5]. To formalize this, let
 156 $\mathbb{E}[(f(X) - Y)^2]$ be the expected MSE (Mean Squared Error) measuring $f(X)$'s fidelity to Y . While
 157 natural in regression settings, there exist extensions on MSE in classification settings (e.g., Hamming
 158 distance for binary; 0–1 loss for nominal multiclass; MSE for ordinal labels). We characterize the
 159 model's reliability within each stratum k by an upper bound on its expected error, denoted by ϵ_k . This
 160 means, $\mathbb{E}[(f(X) - Y)^2 \mid X \in \mathcal{A}_k] \leq \epsilon_k$. We further define $\epsilon^* = \max_k \epsilon_k$ as the worst-case error
 161 across all strata. In our framework, a lower ϵ_k implies $f(X)$ is a stronger proxy for Y in stratum k .

162 4 Main Algorithm

163 AdaStrat-ASI is designed to approximate the *oracle* way of spending a fixed labeling budget across
 164 strata. If we knew, for each stratum \mathcal{A}_k , how much labeling within it would reduce the variance
 165 of the final M-estimator, then the optimal design would follow a Neyman-style principle: allocate
 166 more labels to strata that are (i) prevalent (larger w_k) and (ii) *hard* for the predictor in the sense of
 167 contributing more to the estimator's asymptotic covariance. Formally, the oracle allocation depends on
 168 the stratum-specific, label-dependent quantities (e.g. true error, and resulting covariance contribution),
 169 and is therefore unavailable at design time. Our key observation is that these oracle difficulty terms
 170 admit a reasonable *label-free proxy* that can be estimated using only the black-box predictor $f(\cdot)$. We
 171 define a stratum difficulty score $\hat{\sigma}_k^f$ using the proxy Hessian and gradient covariance computer from
 172 $f(\cdot)$ on unlabeled data, and then allocate the main labeling budget according to $n_{b_k} \propto w_k \cdot \sqrt{\hat{\sigma}_k^f}$
 173 which mirrors the oracle Neyman allocation but requires no access to Y . Algorithm 1 implements
 174 this idea via a short *scouting* phase for estimating $\hat{\sigma}_k^f$ followed by *attacking* phase for sampling
 175 and labeling according to the resulting stratum budgets, and finally computing an ASI-style (Eq. 1)
 176 estimator. Additional details on AdaStrat-ASI are described as follows.

177 4.1 Scouting Stage

178 In this phase, we leverage the availability of unlabeled data and the black-box model $f(x)$ to
 179 estimate the geometric properties of the loss landscape without querying any human labels. Let
 180 $\hat{\theta} := \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell_{\theta}(X_i, f(X_i))$. First, we compute the proxy Hessian \hat{H}_f using the entire
 181 unlabeled dataset at $\hat{\theta}$. This matrix captures the global steepness of the proxy loss surface, serving as
 182 a substitute for the unknown true Hessian. Then, for each stratum k , compute the covariance matrix
 183 \hat{V}_k^f of the loss gradients, using $f(x)$, at $\hat{\theta}$. This measures how much the loss gradient vector fluctuates
 184 within that specific sub-population, solely based on the model's perspective. Finally, we combine
 185 the global curvature and local variance to compute the *difficulty* score $\hat{\sigma}_k^f = \text{Tr}(\hat{H}_f^{-1} \hat{V}_k^f \hat{H}_f^{-1})$. This
 186 scalar quantifies how hard it is to minimize the loss in stratum k relative to the global objective. As a

187 result, we get a set of scores $\{\hat{\sigma}_1^f, \dots, \hat{\sigma}_K^f\}$ that dictate the allocation of the labeling budget in the
 188 attacking stage.

189 4.2 Attacking Stage

190 Let n_{b_k} denote the number of labels allocated to stratum k , such that $\sum_{k=1}^K n_{b_k} = n_b$. We now
 191 distribute the labeling budget n_b to the areas identified as most critical using difficulty scores
 192 derived in the scouting stage: strata with higher weighted difficulty ($n_{b_k} \propto w_k \sqrt{\hat{\sigma}_k^f}$) receive more
 193 samples. In practice, n_{b_k} is realized by drawing $(n_{b_1}, \dots, n_{b_K}) \sim \text{Multinomial}(n_b, \vec{z})$, where
 194 $z_k = \frac{w_k \sqrt{\hat{\sigma}_k^f}}{\sum_{j=1}^K w_j \sqrt{\hat{\sigma}_j^f}}$. This construction guarantees $\mathbb{E}[\sum_{k=1}^K n_{b_k}] = n_b$.

195 Once the budget is distributed, we locally apply ASI to each stratum. Specifically, we observe
 196 i.i.d. unlabeled pre-partitioned covariates $\mathcal{A}_{k1:n_k}$. We consider a budget-scaled family of sampling
 197 rules per stratum $\pi_{\eta_k}(x) = \eta_k u(x)$, $\eta_k \in \mathcal{H} \subseteq \mathbb{R}^+$, and η_k is chosen to match the stratum
 198 label budget. We set $\hat{\eta}_k = \max\{\eta_k \in \mathcal{H} : \eta_k \sum_{i=1}^{n_k} u(X_i) \leq n_{b_k}\}$. This gives rise to
 199 $\vec{\eta} := [\eta_1, \eta_2, \dots, \eta_k]$ and $\hat{\vec{\eta}} := [\hat{\eta}_1, \hat{\eta}_2, \dots, \hat{\eta}_k]$. The later deploys $\pi_{\vec{\eta}}$, to ensure the expected number
 200 of collected labels does not exceed the overall budget. We denote $\hat{\theta}^{\vec{\eta}} \equiv \hat{\theta}^{\pi_{\vec{\eta}}}$. Optionally, to
 201 guarantee stability and some level of acceptance, i.e. to lower bound our acceptance probabilities,
 202 we rely on τ -mixing, i.e. essentially interpolation of our sampling rule and uniform sampling. Thus,
 203 $\pi_k^{(\tau)}(x) := \frac{n_{b_k}}{n_k} \left[(1 - \tau) \frac{u(x)}{\mathbb{E}_{X \sim \mathcal{A}_k}[u(X)]} + \tau \right]$, where $\tau \in (0, 1)$. We query the oracle for labels Y_i
 204 based on the randomized decisions $\xi_{i,k} \sim \text{Bern}(\pi_k^{(\tau)}(X_i))$. It is worth noting, in the extreme situation
 205 of $[\hat{\sigma}_1^f, \dots, \hat{\sigma}_K^f] = \vec{0}$ computed from Section 4.1, we use zero human labels for all strata ($n_{b_k} = 0$ for
 206 all k) and purely rely on model output for computing the risk.

207 4.3 Final Estimation

208 Thus, our *stratified active estimator* is defined as $\hat{\theta}^{\pi_{\vec{\eta}}} = \arg \min_{\theta} L^{\pi_{\vec{\eta}}}(\theta)$, where a simple stratum-
 209 wise rearrangement of Eq. (1) gives us our unbiased *stratified active empirical risk*

$$L^{\pi_{\vec{\eta}}}(\theta) = \frac{1}{n} \sum_{k=1}^K \left[\sum_{i \in \mathcal{A}_k} \left[\ell_{\theta,i}^f + \frac{\xi_{i,k}}{\pi_k^{(\tau)}(X_i)} (\ell_{\theta,i} - \ell_{\theta,i}^f) \right] \right] \quad (3)$$

210 A short calculation show us that our method stays under budget; $\mathbb{E}[n_{\text{lab}}] = \mathbb{E} \left[\sum_{k=1}^K \mathbb{E}[\sum_{i \in \mathcal{A}_k} \xi_{i,k} \mid n_{b_k}, \{X_i\}] \right] = \mathbb{E} \left[\sum_{k=1}^K \sum_{i \in \mathcal{A}_k} \pi_k(X) \right] = \mathbb{E} \left[\sum_{k=1}^K n_{b_k} \right] \leq n_b$. On comparing Eq. (3) to Eq. (1),
 211 the only change comes from using a distinct sampling rule for every k , which in-turn depends on the
 212 label allocation from the previous scouting stage.
 213

214 We provide some insights into data stratification for AdaStrat-ASI in the Appendix in Section E.4.

215 5 Theoretical Analysis

216 Like ASI, our setup also requires standard, mild smoothness assumptions on the loss ℓ_{θ} , articulated
 217 as follows:

218 **Assumption 5.1.** The loss ℓ is smooth if:

- 219 1. $\ell_{\theta}(x, y)$ is differentiable at θ^* for all (x, y) ;
- 220 2. ℓ_{θ} is locally Lipschitz around θ^* : there is a neighborhood of θ^* such that $\ell_{\theta}(x, y)$ is $C(x, y)$ -
 221 Lipschitz and $\ell_{\theta}(x, f(x))$ is $C(x)$ -Lipschitz in θ , where $\mathbb{E}[C(X, Y)^2] < \infty, \mathbb{E}[C(X)^2] < \infty$;
- 222 3. $L(\theta) = \mathbb{E}[\ell_{\theta}(X, Y)]$ and $L^f(\theta) = \mathbb{E}[\ell_{\theta}(X, f(X))]$ are double-differentiable, and $H_{\theta^*} =$
 223 $\nabla^2 \mathbb{E}[\ell_{\theta^*}(X, Y)] = \nabla^2 L(\theta^*) \succ 0$;

224 **Assumption 5.2.** Additionally, we assume the gradient of the loss function $\nabla \ell_{\theta}(x, y)$ is stratumwise
 225 L_y -Lipschitz continuous with respect to the second argument y i.e. for $x \in \mathcal{A}_k, \|\nabla \ell_{\theta}(x, y) -$
 226 $\nabla \ell_{\theta}(x, y')\| \leq L_y(x) \|y - y'\|$ where $\mathbb{E}[L_y(X)^2 | \mathcal{A}_k] < \infty$.

227 In this section, we begin by proving that our AdaStrat-ASI produces valid confidence intervals
 228 and our estimator asymptotically normal. We then shed some light on the design of the Scouting
 229 Stage (Section 4.1) of Algorithm 1 in the Appendix B. We provide justification for relying on model
 230 predictions for label allocation. Finally, we show under mild conditions, our estimator is guaranteed
 231 to be more statistically efficient than the standard ASI estimator, resulting in a smaller asymptotic
 232 covariance.

233 In ASI it is established that the empirically chosen budget parameter $\hat{\eta}$ converges to the true optimal
 234 parameter η^* using Claim A.1 in Zrnić and Candès [65]. In AdaStrat-ASI, we run independent budget
 235 optimizations within each stratum k . We show simultaneous convergence happens for all strata (i.e.
 236 $\mathbb{P}(\hat{\eta}_k \neq \eta_k^*) \rightarrow 0 \quad \forall k$) in Claim C.1 in Appendix.

Theorem 5.3. *Assume the loss is smooth (Ass. 5.1). Assume that for each stratum k , there exists an optimal budget parameter $\eta_k^* \in \mathcal{H}$ such that $\mathbb{P}(\hat{\eta}_k \neq \eta_k^*) \rightarrow 0$. Then, if $\hat{\theta}^{\hat{\eta}^*} \xrightarrow{p} \theta^*$, we have:*

$$\sqrt{n}(\hat{\theta}^{\hat{\eta}} - \theta^*) \xrightarrow{d} N(0, \Sigma_{\hat{\eta}^*}), \text{ where}$$

$$\Sigma_{\hat{\eta}^*} = H_{\theta^*}^{-1} \left[\sum_{k=1}^K w_k \text{Var} \left(\nabla \ell_{\theta^*}(X, f(X)) + \frac{\xi^{\eta_k^*}}{\pi_{\eta_k^*}(X)} (\nabla \ell_{\theta^*}(X, Y) - \nabla \ell_{\theta^*}(X, f(X))) \mid X \in \mathcal{A}_k \right) \right] H_{\theta^*}^{-1},$$

and $\xi^{\eta_k^*} \sim \text{Bern}(\pi_{\eta_k^*}(X))$. Consequently, for any consistent plug-in estimate $\hat{\Sigma} \xrightarrow{p} \Sigma_{\hat{\eta}^*}$, the confidence interval $\mathcal{C}_\alpha = (\hat{\theta}_j^{\hat{\eta}} \pm z_{1-\alpha/2} \sqrt{\frac{\hat{\Sigma}_{jj}}{n}})$ is a valid $(1 - \alpha)$ -asymptotic confidence interval for θ_j^* :

$$\lim_{n \rightarrow \infty} \mathbb{P}(\theta_j^* \in \mathcal{C}_\alpha) = 1 - \alpha.$$

237 A key challenge lies in the allocation of n_b labels across K strata. Let R_k^* represent the asymptotic
 238 weighted variance of the loss residuals within stratum k ; $R_k^* := \text{Tr}(H_{\theta^*}^{-1} W_k^* H_{\theta^*}^{-1})$. W_k^* is the
 239 Weighted Covariance of the Residual Gradients in stratum k : $W_k^* = \mathbb{E}[\frac{1}{s_k(X)} (\nabla \ell_{\theta^*}(X, Y) -$
 240 $\nabla \ell_{\theta^*}(X, f(X))) (\cdot)^\top | \mathcal{A}_k]$. Here, $s_k(X) = \frac{u(X)}{\mathbb{E}_{X \sim \mathcal{A}_k}[u(X)]}$ is the relative uncertainty-sampling score.

241 **Proposition 5.4.** *Given a fixed budget n_b and fixed stratum characteristics $\{w_k, R_k^*\}_{k=1}^K$, such that*
 242 *$w_k > 0$, $R_k^* > 0$, the allocation distribution vector $\{n_{b_1}, \dots, n_{b_K}\}$ that minimizes AdaStrat-ASI's*
 243 *asymptotic variance $\Sigma_{\hat{\eta}^*}$ (Theorem 5.3) subject to $\sum_{k=1}^K n_{b_k} = n_b$, $n_{b_k} > 0$, $\forall k$ is unique and*
 244 *satisfies: $n_{b_k}^* = n_b \cdot p_k^*$, where $p_k^* = \frac{w_k \sqrt{R_k^*}}{\sum_{j=1}^K w_j \sqrt{R_j^*}}$*

245 When deeply inspected, on abstracting away H_{θ^*} and $s_k(X)$ from R_k^* , we notice it is affected by W_k^* .
 246 W_k^* further at a basal level decomposes as, $\text{Cov}(\nabla \ell_{\theta^*}(X, Y) | \mathcal{A}_k) + \text{Cov}(\nabla \ell_{\theta^*}(X, f(X)) | \mathcal{A}_k) -$
 247 $2\text{Cov}(\nabla \ell_{\theta^*}(X, Y), \nabla \ell_{\theta^*}(X, f(X)) | \mathcal{A}_k)$. We notice R_k^* is small when the model perfectly tracks
 248 the variations in Y . Conversely, surviving variance is high when the model and true gradients'
 249 are misaligned (third term ≈ 0). The aforementioned allocation strategy prioritizes stratum where
 250 such *cancellations* fail. Intuitively, if $Y \approx \hat{Y}$ in stratum k , $\sqrt{R_k^*}$ decreases, implying spending
 251 less labeling budget there and relying more on the predictor. While this characterizes the optimal
 252 allocation for AdaStrat-ASI's estimator's variance reduction, the quantity R_k^* depends on both Y and
 253 \hat{Y} and cannot be estimated without labeled data, making it unsuitable for zero-label scouting. This
 254 limitation is not exclusive to AdaStrat-ASI, but applies to any ML-enabled IPW estimator.

Lemma 5.5. *Assuming the loss function is smooth (Ass. 5.1), loss gradient is Lipschitz in y (Ass. 5.2) for any stratum k , and sampling policy enforces a strict lower bound (τ -mixing) on all selection probabilities. Then the square root of asymptotic weighted variance of the loss residuals $\sqrt{R_k^*}$ is bounded by the stratum-specific absolute model error $\sqrt{\epsilon_k}$. i.e.,*

$$\sqrt{R_k^*} = \mathcal{O}(\sqrt{\epsilon_k})$$

255 Crucially, the inaccessibility of R_k^* does not undermine the validity of the algorithm. While R_k^* is
 256 a latent oracle quantity requiring access to true labels, Lemma 5.5 shows that it is upper bounded
 257 by the model error ϵ_k . This observation does not render the value of R_k^* irrelevant; rather, it shows
 258 that its magnitude is controlled by a measurable notion of model fidelity. In regimes where ϵ_k is
 259 small—whether assessed through validation data or prior domain knowledge an oracle allocation

260 would assign negligible labeling budget to that stratum. In fact, in the limit of a perfect local model
 261 ($\epsilon_k = 0$), the optimal allocation requires zero human labels for the stratum. Importantly, this collapse
 262 is not a smooth limit but a consequence of the discrete nature of budget allocation: once empirical
 263 difficulty estimates vanish across all strata ($\epsilon^* = 0$), no human labels are required.

264 **Justification for using f during Scouting Stage.** To obtain a practical and analyzable allocation rule,
 265 we rely on cheap, model predictions as discussed in Section 4.1. In a stratified regime, $\hat{\sigma}_k^f$ effectively
 266 measures the stability of the model’s decision-making within a partition. $\hat{\sigma}_k^f$ is inherently designed
 267 to be small in stable, homogeneous regions where f behaves consistently. In such stable regions, a
 268 reasonably calibrated model is also likely to have consistently low residual variance. Similarly in
 269 a heterogenous strata, high proxy gradient variance indicates that the stratum lies along complex
 270 model decision boundaries leading to unstable predictions. These are precisely the regions where
 271 the underlying task is difficult, implying that $f(x)$ is prone to errors and R_k^* is large. We notice,
 272 $\hat{\sigma}_k^f$ and R_k^* are affected similarly under different conditions. While σ_k^f is not a labeled estimate of
 273 oracle difficulty R_k^* , it serves as a ranking signal for reallocating a fixed budget toward strata where
 274 the predictor’s behavior is less stable. Similar gradient based variance heuristics have been used to
 275 capture fit difficulty in coreset selection [40] and impact of data on learning [1, 30, 54, 53], aligning
 276 well in practice.

277 Further, we empirically validate this proxy on the California Housing Dataset (introduced later in Section E.2.1) by comparing the oracle allocation \mathbf{p}^* (computed with labels) to our proxy allocation $\hat{\mathbf{p}}$, finding close agreement (Mean Absolute Error (MAE) of 0.0050, Figure 1). Moreover, for budget parity with competing baselines, we reserve the entire human labeling budget for the estimation stage; thus scouting must be zero-label, making cheap model predictions the natural source of information before labeling. To formalize AdaStrat-ASI’s efficiency gains, we view ASI’s asymptotic variance, through the lens of stratification. This is purely an analytical decomposition, as ASI does not algorithmically rely on stratified sampling. Recall, from Eq. (1) and Eq. (3), ASI and AdaStrat-ASI, differ in the use of sampling rule. The former relies on a global policy for the entire data manifold, however the latter relies on k local policies, one for each stratum. We also know, AdaStrat-ASI enforces fixed local allocation $n_{b_k}^{\text{AdaStrat-ASI}}$ (In previous sections, we refer to this simply as n_{b_k}).

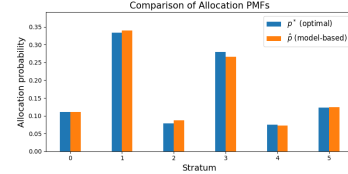


Figure 1: Allocation Comparison

292 But in ASI, $\pi(x) = \frac{n_b}{n} \cdot \frac{u(x)}{\mathbb{E}[u[X]]} \propto u(x)$. This means, high uncertainty points pull budget toward
 293 them strongly and low uncertainty points have a weak pull. Since $\pi(x) \propto u(x)$, this becomes:
 294 $\mathbb{E}[n_{b_k}^{\text{ASI}}] \propto \sum_{x \in A_k} u(x)$, where $n_{b_k}^{\text{ASI}}$ denotes the random number of samples that fall into stratum k
 295 after ASI is ran. From this, we can clearly infer: in ASI, the implicit local budget allocation depends
 296 on the total uncertainty volume in that stratum. Let $\gamma_k = w_k^2 R_k^*$ denote the weighted oracle difficulty
 297 of stratum k .

298 **Theorem 5.6.** Assume the conditions for asymptotic normality of the active M -estimator hold (
 299 Theorem 5.3). Let Σ_{ASI} and $\Sigma_{\text{AdaStrat-ASI}}^1$ be the asymptotic covariance matrices for the ASI and
 300 AdaStrat-ASI estimators respectively. Then Variance reduction is guaranteed if weighted sum of label
 301 reciprocals for ASI is more than AdaStrat-ASI i.e. if,

$$\sum_{k=1}^K \frac{\gamma_k}{\mathbb{E}[n_{b_k}^{\text{ASI}}]} > \sum_{k=1}^K \frac{\gamma_k}{n_{b_k}^{\text{AdaStrat-ASI}}} \quad (4)$$

then,

$$\text{Tr}(\Sigma_{\text{ASI}}) > \text{Tr}(\Sigma_{\text{AdaStrat-ASI}})$$

302 Theorem 5.6 intuitively shows us, AdaStrat-ASI guarantees variance reduction by *stealing* samples
 303 from easy strata and giving them to hard strata. We further empirically study the difference in
 304 labelling distribution in Appendix E.9. One notes, if γ_k is constant for all k , then Eq. (4), is never
 305 true; the weighted distributions become the same. This aligns with AdaStrat-ASI’s mechanism to
 306 exploit the heterogeneity across strata to improve efficiency.

¹We use these notations instead of standard Σ_* and $\Sigma_{\hat{\pi}_*}$, for ease of readability and direct comparison.

307 The gap in Eq. 4 is large when the model’s uncertainty (which drives ASI’s allocation $\mathbb{E}[n_{b_k}^{\text{ASI}}]$) is
 308 misaligned with the model’s true error, which drives AdaStrat-ASI’s allocation. For example, if
 309 a model is “confidently wrong” in a specific stratum, ASI will starve it of samples, inflating the
 310 variance. On the other hand, AdaStrat-ASI will detect high pseudo-gradient fluctuations and will
 311 assign an informed budget $n_{b_k}^{\text{AdaStrat-ASI}}$, ensuring robust estimation.

312 It is also worth noting Eq. 4 is always true if n_{b_k} is the global minimizer of the true objective
 313 function as shown in Proposition 5.4. Thus the inequality holds if and only if the proxy Allocation
 314 is closer to the Oracle Allocation (see Figure 1). But more importantly, we wish to emphasize the
 315 mathematical safety of our scouting heuristic : because AdaStrat-ASI strictly decouples the Scouting
 316 stage from the Attacking stage, a failure of the proxy model is entirely benign to the statistical validity
 317 of the estimator (Theorem 5.3). That is, even if the scouting phase generates a highly suboptimal
 318 allocation vector, the resulting confidence intervals remain statistically valid. The heuristic only
 319 governs efficiency (variance reduction; Theorem 5.6), never validity. We provide more details in
 320 Appendix D.

321 6 Experiments

322 **Motivational Synthetic Example.** To isolate the effect of
 323 stratification on variance reduction, we construct a highly controlled
 324 1D logistic regression setting. We draw $x \sim \mathcal{U}_{[-6,6]}$
 325 and generate labels from a ground-truth model $y \mid x \sim$
 326 Bernoulli($\sigma(\theta_1 \cdot x + \theta_0)$), where $(\theta_1, \theta_0) = (1, 0)$. A logistic
 327 regression model is fit on 5000 such samples, yielding an acqui-
 328 sition model $\hat{f}(x) = \sigma(\hat{\theta}_1 \cdot x + \hat{\theta}_0)$. We consider the task of esti-
 329 mating the population mean $\mathbb{E}[Y]$ under a deliberately structured
 330 test distribution consisting of two Gaussian clusters: $\mathcal{C}_1 : x \sim$
 331 $\mathcal{N}(8, 0.1^2)$, $y = 1$, $\mathcal{C}_2 : x \sim \mathcal{N}(-8, 0.1^2)$, $y = 0$,
 332 each containing 5000 points. By construction, the true mean
 333 equals 0.5. Both clusters are symmetrically located far from
 334 the decision boundary at $x = 0$, and therefore lie in regions of
 335 near-deterministic predictions under \hat{f} (Figure 3).

336 We notice both clusters lie in regions of low model uncertainty
 337 and are symmetrically positioned with respect to the decision
 338 boundary. Consequently, model uncertainty based sampling pro-
 339 vides no additional signal for prioritization, causing ASI to ef-
 340 fectively reduce to uniform sampling. For stratification, we rely
 341 on with the data-generating structure. As a result, the empirical
 342 variance estimates within each stratum collapse to zero across
 343 trials (Figure 4). Although contrived, this example highlights
 344 the mechanism by which stratification improves efficiency: when strata align with regions of high
 345 black-box model certainty and low outcome variability, stratified ASI can dramatically reduce both
 346 empirical variance(and in extreme cases labeling cost). In realistic settings, such perfect alignment is
 347 unlikely; nevertheless, this toy example serves to isolate and illustrate the maximal benefit achievable
 348 through stratification.

349 **Real world Datasets.** We provide experimental details in the Appendix (Section E.1).

350 **AlphaFold Dataset.** Inspired by Angelopoulos et al. [2] and Zrnić and Candès [65], we incorpo-
 351 rate the AlphaFold dataset($n = 10,802$) sourced from the `ppi_py` package² as an application
 352 to proteomics. Rather than estimating the odds ratio of phosphorylation as in Angelopoulos
 353 et al. [2], Zrnić and Candès [65], we shift the target of inference to the *mean disorder rate* i.e.
 354 $\theta^* = \mathbb{E}[Y] = \mathbb{P}(\text{residue is disordered})$, to estimate the fraction of PTM-modified residues that are
 355 truly intrinsically disordered.

356 From Figure 5, we see AdaStrat-ASI offering the largest reduction in CI width(approximately
 357 30%(ASI), 32%(Uniform), and 64%(Classical) reduction). AdaStrat-ASI’s samples also yield

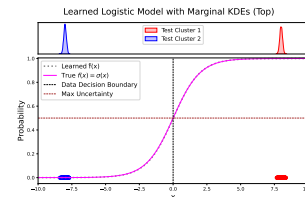


Figure 2: Toy Setup

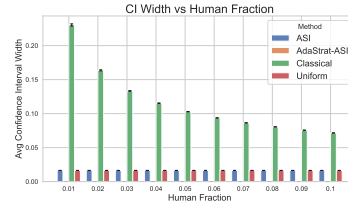


Figure 3: CI width

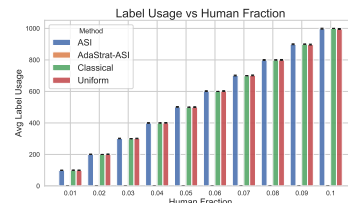


Figure 4: Label Usage

²Github Repository: https://github.com/aangelopoulos/ppi_py/tree/main

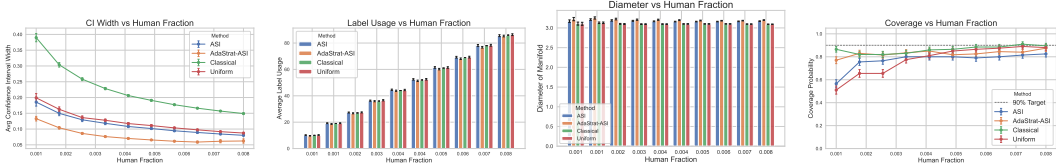


Figure 5: AlphaFold Dataset

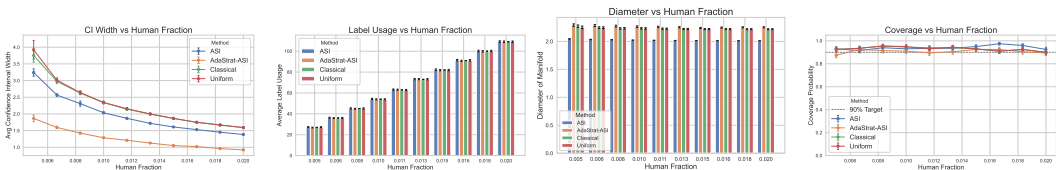


Figure 6: Politeness Dataset

358 the largest diameter. Interestingly, other ML-powered methods struggle with coverage at tight
 359 budgets ($n_b < 0.004$). By contrast, AdaStrat-ASI decomposes the uncertainty into stratified sub-
 360 regions. This local estimation allows us to more accurately capture the residual variance within each
 361 stratum, effectively shielding the CI from becoming biased.

362 **Politeness Dataset.** We use the dataset³ used by Gligoric et al. to gauge politeness of online
 363 text requests posted on Stack Exchange and Wikipedia ($n = 5,480$). In this estimation task, θ^*
 364 corresponds to the logistic regression coefficient β_{hedge} measuring the impact of a linguistic feature
 365 such as hedging on the perceived politeness, $\text{logit}(P(H_{\text{polite}} = 1 | X_{\text{hedge}})) = \beta_{\text{hedge}} X_{\text{hedge}} + \beta_0$,
 366 where $X_{\text{hedge}} = 1$ indicates the presence of the hedge marker and $H_{\text{polite}} = 1$ indicates annotation
 367 as polite.

368 From Figure 6, all methods achieve coverage close to the nominal 90% level across budgets, indicating
 369 that inferential validity is preserved throughout. Among the ML-enabled approaches, AdaStrat-ASI
 370 most consistently maintains the best fidelity. Label usage is nearly identical across all methods.
 371 In particular, AdaStrat-ASI yields a large reduction in confidence interval width relative to other
 372 baselines across all human fractions, with the most pronounced gains at smaller budgets; up to 54%
 373 and 45.67% in case of Classical and ASI respectively. AdaStrat-ASI exhibits a modest but consistent
 374 increase in diameter indicating broader capture of the data manifold.

375 These results collectively highlight that, even when uncertainty alone is informative enough to
 376 maintain coverage, stratification provides an additional efficiency gain by promoting efficiency
 377 (reduced CI). Also, stratifying data first and then applying ASI within each stratum, catapults the
 378 budget sampling, causing broader labeled data manifold spread.

379 7 Conclusion

380 **Limitations and Discussion.** A potential pitfall of AdaStrat-ASI, inherited from ASI, is the
 381 reliance on well-calibrated estimates of model uncertainty. A poorly-calibrated model may also
 382 lead to suboptimal allocation during scouting. On datasets with strata having lack of inter-stratum
 383 heterogeneity, AdaStrat-ASI may not be able to reduce CI width(Theorem 5.6). Finally, like ASI, our
 384 theoretical guarantees rest on the assumption that data is i.i.d.

385 We propose Adaptive Stratified Active Statistical Inference (AdaStrat-ASI), a novel inference frame-
 386 work that directs labeling effort toward regions with the greatest impact on the statistical efficiency.
 387 By leveraging model uncertainty locally, AdaStrat-ASI achieves more powerful inferences with better
 388 spread of labeled data. We provide theoretical guarantees for its asymptotic validity and demonstrate
 389 its practical effectiveness across datasets. We believe this work establishes a new principle for
 390 efficient data collection, showing that latent structure in data of the statistical problem is as crucial as
 391 the uncertainty of the predictive model.

³Github Repository: <https://github.com/kristinagligoric/confidence-driven-inference/tree/main/datasets>

392 **References**

- 393 [1] C. Agarwal, D. D’souza, and S. Hooker. Estimating example difficulty using variance of
394 gradients, 2022. URL <https://arxiv.org/abs/2008.11600>.
- 395 [2] A. N. Angelopoulos, S. Bates, C. Fannjiang, M. I. Jordan, and T. Zrnic. Prediction-powered
396 inference, 2023. URL <https://arxiv.org/abs/2301.09633>.
- 397 [3] A. N. Angelopoulos, J. C. Duchi, and T. Zrnic. Ppi++: Efficient prediction-powered inference.
398 2024. URL <https://arxiv.org/abs/2311.01453>.
- 399 [4] J. Angwin, J. Larson, S. Mattu, and L. Kirchner. Machine bias: There’s soft-
400 ware used across the country to predict future criminals. and it’s biased against
401 blacks. ProPublica, May 2016. URL [https://www.propublica.org/article/
402 machine-bias-risk-assessments-in-criminal-sentencing](https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing). Accessed 2025-12-13.
- 403 [5] S. L. Blodgett, S. Barocas, H. Daumé III, and H. Wallach. Language (technology) is power: A
404 critical survey of “bias” in NLP. In D. Jurafsky, J. Chai, N. Schlueter, and J. Tetreault, editors,
405 *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages
406 5454–5476, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/
407 2020.acl-main.485. URL <https://aclanthology.org/2020.acl-main.485/>.
- 408 [6] A. Carpentier, R. Munos, and A. Antos. Adaptive strategy for stratified monte carlo sampling.
409 *Journal of Machine Learning Research*, 16(68):2231–2271, 2015. URL [http://jmlr.org/
410 papers/v16/carpentier15a.html](http://jmlr.org/papers/v16/carpentier15a.html).
- 411 [7] K. Chaloner and I. Verdinelli. Bayesian experimental design: A review. *Statistical Science*, 1995.
412 URL [https://projecteuclid.org/journals/statistical-science/volume-10/
413 issue-3/Bayesian-Experimental-Design-A-Review/10.1214/ss/1177009939](https://projecteuclid.org/journals/statistical-science/volume-10/issue-3/Bayesian-Experimental-Design-A-Review/10.1214/ss/1177009939).
414 full.
- 415 [8] Y. Chandak, S. Shankar, V. Syrgkanis, and E. Brunskill. Adaptive instrument design for indirect
416 experiments, 2023. URL <https://arxiv.org/abs/2312.02438>.
- 417 [9] T. Chen and C. Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the
418 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, page
419 785–794. ACM, Aug. 2016. doi: 10.1145/2939672.2939785. URL [http://dx.doi.org/10.
420 1145/2939672.2939785](http://dx.doi.org/10.1145/2939672.2939785).
- 421 [10] C.-H. Chiang and H.-y. Lee. Can large language models be an alternative to human evaluations?
422 In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting
423 of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15607–15631,
424 Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.
425 acl-long.870. URL <https://aclanthology.org/2023.acl-long.870/>.
- 426 [11] J. A. Clayton and F. S. Collins. Policy: Nih to balance sex in cell and animal studies. *Nature*,
427 509(7500):282–283, 2014. URL <https://www.nature.com/articles/509282a>.
- 428 [12] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek,
429 J. Hilton, R. Nakano, C. Hesse, and J. Schulman. Training verifiers to solve math word
430 problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
- 431 [13] W. Cochran. *Sampling Techniques*. Wiley Series in Probability and Statistics. Wiley, 1977.
432 ISBN 9780471162407. URL <https://books.google.com/books?id=8Y4QQAIAAJ>.
- 433 [14] L. N. Cook, A. Mishler, A. Ramdas, and T. Zrnić. Semiparametric efficient inference from adap-
434 tively collected data. In *Proceedings of the Conference on Learning Theory and Representations*,
435 2024. URL <https://openreview.net/forum?id=mfsGIZpwi0>.
- 436 [15] C. Cortes, G. DeSalvo, C. Gentile, M. Mohri, and N. Zhang. Region-based active learning.
437 In K. Chaudhuri and M. Sugiyama, editors, *Proceedings of the Twenty-Second International
438 Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine
439 Learning Research*, pages 2801–2809. PMLR, 16–18 Apr 2019. URL [https://proceedings.
440 mlr.press/v89/cortes19a.html](https://proceedings.mlr.press/v89/cortes19a.html).

- 441 [16] S. Dasgupta and D. Hsu. Hierarchical sampling for active learning. In *Proceedings of the*
442 *25th International Conference on Machine Learning*, ICML '08, page 208–215, New York,
443 NY, USA, 2008. Association for Computing Machinery. ISBN 9781605582054. doi: 10.1145/
444 1390156.1390183. URL <https://doi.org/10.1145/1390156.1390183>.
- 445 [17] J. Dressel and H. Farid. The accuracy, fairness, and limits of predicting recidivism. *Science*
446 *Advances*, 4(1):aao5580, 2018. doi: 10.1126/sciadv.aao5580. URL <https://www.science.org/doi/10.1126/sciadv.aao5580>.
- 448 [18] C. Dwork and A. Roth. The algorithmic foundations of differential privacy. *Found. Trends*
449 *Theor. Comput. Sci.*, 9(3–4):211–407, Aug. 2014. ISSN 1551-305X. doi: 10.1561/04000000042.
450 URL <https://doi.org/10.1561/04000000042>.
- 451 [19] A. Fisch, J. Maynez, R. A. Hofer, B. Dhingra, A. Globerson, and W. W. Co-
452 hen. Stratified prediction-powered inference for effective hybrid evaluation of language
453 models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tom-
454 czak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*,
455 volume 37, pages 111489–111514. Curran Associates, Inc., 2024. doi: 10.52202/
456 079017-3541. URL [https://proceedings.neurips.cc/paper_files/paper/2024/
457 file/c9fcd02e6445c7dfbad6986abee53d0d-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/c9fcd02e6445c7dfbad6986abee53d0d-Paper-Conference.pdf).
- 458 [20] A. Garivier, E. Kaufmann, and T. Lattimore. On explore-then-commit strategies, 2016. URL
459 <https://arxiv.org/abs/1605.08988>.
- 460 [21] F. Gilardi, M. Alizadeh, and M. Kubli. ChatGPT outperforms crowd workers for text-annotation
461 tasks. *Proceedings of the National Academy of Sciences*, 120(30):e2305016120, 2023. URL
462 <https://www.pnas.org/doi/10.1073/pnas.2305016120>.
- 463 [22] K. Gligoric, T. Zrníc, C. Lee, E. Candes, and D. Jurafsky. Can unconfident LLM annota-
464 tions be used for confident conclusions? In L. Chiruzzo, A. Ritter, and L. Wang, editors,
465 *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Asso-*
466 *ciation for Computational Linguistics: Human Language Technologies (Volume 1: Long*
467 *Papers)*, pages 3514–3533, Albuquerque, New Mexico, Apr. 2025. Association for Computa-
468 tional Linguistics. ISBN 979-8-89176-189-6. doi: 10.18653/v1/2025.naacl-long.179. URL
469 <https://aclanthology.org/2025.naacl-long.179/>.
- 470 [23] J. Hahn, K. Hirano, and D. Karlan. Adaptive experimental design using the propensity score.
471 *Journal of Business & Economic Statistics*, 29(1):96–108, 2011. doi: 10.1198/jbes.2009.08161.
472 URL <https://doi.org/10.1198/jbes.2009.08161>.
- 473 [24] S. Hanneke. Theory of disagreement-based active learning. *Found. Trends Mach. Learn.*,
474 7(2–3):131–309, June 2014. ISSN 1935-8237. doi: 10.1561/22000000037. URL <https://doi.org/10.1561/22000000037>.
- 476 [25] D. Hovy and A. Søgaard. Tagging performance correlates with author age. In C. Zong and
477 M. Strube, editors, *Proceedings of the 53rd Annual Meeting of the Association for Computational*
478 *Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume*
479 *2: Short Papers)*, pages 483–488, Beijing, China, July 2015. Association for Computational
480 Linguistics. doi: 10.3115/v1/P15-2079. URL <https://aclanthology.org/P15-2079/>.
- 481 [26] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin, and
482 T. Liu. A survey on hallucination in large language models: Principles, taxonomy, challenges,
483 and open questions. *ACM Transactions on Information Systems*, 43(2):1–55, Jan. 2025. ISSN
484 1558-2868. doi: 10.1145/3703155. URL <http://dx.doi.org/10.1145/3703155>.
- 485 [27] G. W. Imbens and D. B. Rubin. *Causal Inference in Statistics, Social, and Biomedical Sciences*.
486 Cambridge University Press, 2015. URL <https://doi.org/10.1017/CB09781139025751>.
- 487 [28] N. Jean, M. Burke, M. Xie, W. M. Davis, D. B. Lobell, and S. Ermon. Combining satellite
488 imagery and machine learning to predict poverty. *Science*, 353(6301):790–794, 2016. URL
489 <https://www.science.org/doi/abs/10.1126/science.aaf7894>.

- 490 [29] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvu-
491 nakool, R. Bates, A. Zidek, A. Potapenko, et al. Highly accurate protein structure predic-
492 tion with alphafold. *Nature*, 596(7873):583–589, 2021. URL [https://doi.org/10.1038/](https://doi.org/10.1038/s41586-021-03819-2)
493 s41586-021-03819-2.
- 494 [30] A. Katharopoulos and F. Fleuret. Not all samples are created equal: Deep learning with
495 importance sampling. In J. Dy and A. Krause, editors, *Proceedings of the 35th International*
496 *Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*,
497 pages 2525–2534. PMLR, 10–15 Jul 2018. URL [https://proceedings.mlr.press/v80/](https://proceedings.mlr.press/v80/katharopoulos18a.html)
498 katharopoulos18a.html.
- 499 [31] M. Kato, T. Ishihara, J. Honda, and Y. Narita. Efficient adaptive experimental design for average
500 treatment effect estimation, 2025. URL <https://arxiv.org/abs/2002.05308>.
- 501 [32] M. Kearns and A. Roth. *The Ethical Algorithm: The Science of Socially Aware Algorithm Design*.
502 Oxford University Press, 2019. URL <https://dl.acm.org/doi/10.5555/3379082>.
- 503 [33] S. Kim, J. Shin, Y. Cho, J. Jang, S. Longpre, H. Lee, S. Yun, S. Shin, S. Kim, J. Thorne, and
504 M. Seo. Prometheus: Inducing fine-grained evaluation capability in language models, 2024.
505 URL <https://arxiv.org/abs/2310.08491>.
- 506 [34] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa. Large language models are zero-shot
507 reasoners, 2023. URL <https://arxiv.org/abs/2205.11916>.
- 508 [35] A. Lewkowycz, A. Andreassen, D. Dohan, E. Dyer, H. Michalewski, V. Ramasesh, A. Slone,
509 C. Anil, I. Schlag, T. Gutman-Solo, Y. Wu, B. Neyshabur, G. Gur-Ari, and V. Misra. Solving
510 quantitative reasoning problems with language models, 2022. URL [https://arxiv.org/](https://arxiv.org/abs/2206.14858)
511 abs/2206.14858.
- 512 [36] P. Li, T. Zrnic, and E. Candès. Robust sampling for active statistical inference, 2025. URL
513 <https://arxiv.org/abs/2511.08991>.
- 514 [37] L. Lin, K. Khamaru, and M. J. Wainwright. Semi-parametric inference based on adaptively
515 collected data, 2025. URL <https://arxiv.org/abs/2303.02534>.
- 516 [38] Y. Liu, D. Iter, Y. Xu, S. Wang, R. Xu, and C. Zhu. G-Eval: NLG Evaluation using Gpt-4 with
517 Better Human Alignment. In H. Bouamor, J. Pino, and K. Bali, editors, *Proceedings of the*
518 *2023 Conference on Empirical Methods in Natural Language Processing*, pages 2511–2522,
519 Singapore, Dec. 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.
520 emnlp-main.153. URL <https://aclanthology.org/2023.emnlp-main.153/>.
- 521 [39] S. Lohr. *Sampling: Design and Analysis*. Chapman & Hall/CRC Texts in Statistical Science.
522 CRC Press, 2021. ISBN 9781000478235. URL [https://books.google.com/books?id=](https://books.google.com/books?id=DahGAAAAQBAJ)
523 DahGAAAAQBAJ.
- 524 [40] S. H. A. Mahmood, M. Yin, and R. Khanna. On the support vector effect in dnns: Rethinking
525 data selection and attribution. In *Proceedings of the 31st ACM SIGKDD Conference on*
526 *Knowledge Discovery and Data Mining V.1*, KDD '25, page 1020–1031, New York, NY, USA,
527 2025. Association for Computing Machinery. ISBN 9798400712456. doi: 10.1145/3690624.
528 3709295. URL <https://doi.org/10.1145/3690624.3709295>.
- 529 [41] K. Margatina, G. Vernikos, L. Barrault, and N. Aletras. Active learning by acquiring contrastive
530 examples. In M.-F. Moens, X. Huang, L. Specia, and S. W.-t. Yih, editors, *Proceedings of*
531 *the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 650–663,
532 Online and Punta Cana, Dominican Republic, Nov. 2021. Association for Computational
533 Linguistics. doi: 10.18653/v1/2021.emnlp-main.51. URL [https://aclanthology.org/](https://aclanthology.org/2021.emnlp-main.51/)
534 2021.emnlp-main.51/.
- 535 [42] J. Miao, X. Miao, Y. Wu, J. Zhao, and Q. Lu. Assumption-lean and data-adaptive post-prediction
536 inference, 2024. URL <https://arxiv.org/abs/2311.14220>.
- 537 [43] P. R. Mohanty. Beyond Disagreement-based Learning for Contextual Bandits. 7 2023.
538 doi: 10.25394/PGS.23739957.v1. URL [https://hammer.purdue.edu/articles/thesis/](https://hammer.purdue.edu/articles/thesis/Beyond_Disagreement-based_Learning_for_Contextual_Bandits/23739957)
539 Beyond_Disagreement-based_Learning_for_Contextual_Bandits/23739957.

- 540 [44] K. Motwani and D. Witten. Revisiting inference after prediction, 2024. URL <https://arxiv.org/abs/2306.13746>.
541
- 542 [45] J. Neyman. On the two different aspects of the representative method: The method of stratified
543 sampling and the method of purposive selection. *Journal of the Royal Statistical Society*, 97(4):
544 558–625, 1934. ISSN 09528385. URL <http://www.jstor.org/stable/2342192>.
- 545 [46] C. Ober, D. A. Loisel, and Y. Gilad. Sex-specific genetic architecture of human disease. *Nature Reviews Genetics*, 9(12):911–922, 2008. URL <https://www.nature.com/articles/nrg2415>.
546
547
- 548 [47] OpenAI. Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.
- 549 [48] R. K. Pace and R. Barry. Sparse spatial autoregressions. *Statistics & Probability Letters*, 33(3):
550 291–297, 1997. URL [https://doi.org/10.1016/S0167-7152\(96\)00140-X](https://doi.org/10.1016/S0167-7152(96)00140-X).
- 551 [49] J. S. Park, J. O’Brien, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein. Generative agents:
552 Interactive simulacra of human behavior. In *Proceedings of the 36th Annual ACM Symposium*
553 *on User Interface Software and Technology*, UIST ’23, New York, NY, USA, 2023. Association
554 for Computing Machinery. ISBN 9798400701320. doi: 10.1145/3586183.3606763. URL
555 <https://doi.org/10.1145/3586183.3606763>.
- 556 [50] J. M. Robins, A. Rotnitzky, and L. P. Zhao. Estimation of regression coefficients when some
557 regressors are not always observed. *Journal of the American Statistical Association*, 89(427):
558 846–866, 1994. doi: 10.1080/01621459.1994.10476818. URL [https://doi.org/10.1080/](https://doi.org/10.1080/01621459.1994.10476818)
559 [01621459.1994.10476818](https://doi.org/10.1080/01621459.1994.10476818).
- 560 [51] E. Rolf, J. Proctor, T. Carleton, I. Bolliger, V. Shankar, M. Ishihara, B. Recht, and S. Hsiang.
561 A generalizable and accessible approach to machine learning with global satellite im-
562 agery. *Nature Communications*, 12(1):4392, 2021. URL [https://doi.org/10.1038/](https://doi.org/10.1038/s41467-021-24638-z)
563 [s41467-021-24638-z](https://doi.org/10.1038/s41467-021-24638-z).
- 564 [52] O. Sener and S. Savarese. Active learning for convolutional neural networks: A core-set
565 approach. In *International Conference on Learning Representations*, 2018. URL [https://](https://openreview.net/forum?id=H1aIuk-RW)
566 openreview.net/forum?id=H1aIuk-RW.
- 567 [53] B. Settles, M. Craven, and S. Ray. Multiple-instance active learning. In *Proceedings of the 21st*
568 *International Conference on Neural Information Processing Systems*, NIPS’07, page 1289–1296,
569 Red Hook, NY, USA, 2007. Curran Associates Inc. ISBN 9781605603520.
- 570 [54] M. Toneva, A. Sordoni, R. T. des Combes, A. Trischler, Y. Bengio, and G. J. Gordon. An
571 empirical study of example forgetting during deep neural network learning, 2019. URL
572 <https://arxiv.org/abs/1812.05159>.
- 573 [55] A. Torralba and A. A. Efros. Unbiased look at dataset bias. In *Proceedings of the 2011 IEEE*
574 *Conference on Computer Vision and Pattern Recognition*, CVPR ’11, page 1521–1528, USA,
575 2011. IEEE Computer Society. ISBN 9781457703942. doi: 10.1109/CVPR.2011.5995347.
576 URL <https://doi.org/10.1109/CVPR.2011.5995347>.
- 577 [56] J. T.-Z. Wei, F. Zufall, and R. Jia. Operationalizing content moderation "accuracy" in the digital
578 services act, 2024. URL <https://arxiv.org/abs/2305.09601>.
- 579 [57] L. Weidinger, J. Uesato, M. Rauh, C. Griffin, P.-S. Huang, J. Mellor, A. Glaese, M. Cheng,
580 B. Balle, A. Kasirzadeh, C. Biles, S. Brown, Z. Kenton, W. Hawkins, T. Stepleton, A. Birhane,
581 L. A. Hendricks, L. Rimell, W. Isaac, J. Haas, S. Legassick, G. Irving, and I. Gabriel. Taxonomy
582 of risks posed by language models. In *Proceedings of the 2022 ACM Conference on Fairness,*
583 *Accountability, and Transparency*, FAccT ’22, page 214–229, New York, NY, USA, 2022.
584 Association for Computing Machinery. ISBN 9781450393522. doi: 10.1145/3531146.3533088.
585 URL <https://doi.org/10.1145/3531146.3533088>.
- 586 [58] M. Xie, N. Jean, M. Burke, D. Lobell, and S. Ermon. Transfer learning from deep features for
587 remote sensing and poverty mapping. In *Proceedings of the AAAI Conference on Artificial Intel-*
588 *ligence*, volume 30, 2016. URL <https://dl.acm.org/doi/10.5555/3016387.3016457>.

- 589 [59] M. Zdun. Machine politics: How america casts and counts its
590 votes. *Reuters*, 2022. URL [https://www.reuters.com/world/us/
591 us-midterm-elections-how-america-casts-counts-its-votes-2022-11-02/](https://www.reuters.com/world/us/us-midterm-elections-how-america-casts-counts-its-votes-2022-11-02/).
- 592 [60] C. Zhang and K. Chaudhuri. Beyond disagreement-based agnostic active learning, 2014. URL
593 <https://arxiv.org/abs/1407.2657>.
- 594 [61] K. W. Zhang, L. Janson, and S. Murphy. Statistical inference with m-estimators on adaptively
595 collected data. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances
596 in Neural Information Processing Systems*, 2021. URL [https://openreview.net/forum?
597 id=TJ0Qw_vMlAj](https://openreview.net/forum?id=TJ0Qw_vMlAj).
- 598 [62] R. Zhang, Y. Li, Y. Ma, M. Zhou, and L. Zou. LLMaAA: Making large language models
599 as active annotators. In *The 2023 Conference on Empirical Methods in Natural Language
600 Processing*, 2023. URL <https://openreview.net/forum?id=B6Gdg7u04y>.
- 601 [63] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing,
602 H. Zhang, J. E. Gonzalez, and I. Stoica. Judging LLM-as-a-judge with MT-bench and chatbot
603 arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and
604 Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=uccHPGDlao>.
- 605 [64] C. Ziems, W. Held, O. Shaikh, J. Chen, Z. Zhang, and D. Yang. Can large language models
606 transform computational social science? *Computational Linguistics*, 50(1):237–291, Mar. 2024.
607 doi: 10.1162/coli_a_00502. URL <https://aclanthology.org/2024.c1-1.8/>.
- 608 [65] T. Zrnić and E. J. Candès. Active statistical inference. In *Proceedings of the 41st International
609 Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, 2024.
610 URL <https://proceedings.mlr.press/v235/zrnic24a.html>. PMLR v235.
- 611 [66] T. Zrnić and E. J. Candès. Cross-prediction-powered inference. *Proceedings of the National
612 Academy of Sciences*, 121(15):e2322083121, 2024. doi: 10.1073/pnas.2322083121. URL
613 <https://www.pnas.org/doi/abs/10.1073/pnas.2322083121>.

615 **A Impact Statement**

616 This work advances principled statistical inference under limited labeling budgets by improving
 617 use of labeled data and variance efficiency through stratified active sampling. Beyond efficiency
 618 gains, the stratified framework enables explicit control over how labeling budgets are allocated across
 619 predefined subgroups, allowing practitioners to enforce balance or representational constraints when
 620 required. This capability is particularly relevant in settings such as survey research, social science, and
 621 AI-assisted data annotation, where subgroup considerations may be mandated for reliability or policy
 622 reasons. In the bigger picture, when applied responsibly, such efficiency gains may reduce costs,
 623 respondent burden, and reliance on large-scale labeling efforts. As a methodological contribution, the
 624 societal impact is indirect and depends on responsible downstream use.

625 **B Algorithm****Algorithm 1** Adaptive Stratified Active Statistical Inference

Input: Unlabeled data partitions $\mathcal{A}_k = \{X_{ik}\}_{i=1}^{n_k}$ for $k \in [K]$; total labeling budget n_b ; predictive model f ; loss $\ell_\theta(x, y)$; proxy parameter using f $\hat{\theta}$; error level α ; strata weight $w_k = \frac{n_k}{n}$.

— **Phase 1: Scouting** —

1: Compute the global proxy Hessian $\hat{H}_f := \frac{1}{n} \sum_{i=1}^n \nabla^2 \ell_{\hat{\theta}}(X_i, f(X_i))$

2:

for each stratum $k \in [K]$ **do**

2.1: Compute the proxy Covariance Matrix $\hat{V}_k^f := \text{Cov}(\nabla \ell_{\hat{\theta}}(X, f(X)) | X \in \mathcal{A}_k)$

2.2 Assemble $\hat{\Sigma}_k := \hat{H}_f^{-1} \hat{V}_k^f \hat{H}_f^{-1}$

2.3: Compute $\hat{\sigma}_k^f := \text{Tr}(\hat{\Sigma}_k)$

end for

— **Phase 2: Attacking** —

3: For each k , allocate the main budget using *proxy budget allocation weight* based on the

scouting results i.e. $(n_{b_1}, \dots, n_{b_K}) \sim \text{Multinomial}(n_b, \vec{z})$, where $z_k = \frac{w_k \sqrt{\hat{\sigma}_k^f}}{\sum_{j=1}^K w_j \sqrt{\hat{\sigma}_j^f}}$

4:

for each stratum $k \in [K]$ **do**

4.1: Compute sampling probabilities for all $x \in \mathcal{A}_k$: $\pi_k(x) = \frac{u(x) \cdot \hat{\eta}_k}{\mathbb{E}_{X \sim \mathcal{A}_k}[u(X)]}$, $\hat{\eta}_k = \frac{n_{b_k}}{n_k}$.

4.2 Optionally: Select $\tau \in (0, 1)$ and choose sampling rule $\pi_k^{(\tau)}(x) = (1 - \tau) \cdot \pi_k(x) + \tau \cdot \pi_k^{\text{unif}}$, where $\pi_k^{\text{unif}} = \frac{n_{b_k}}{n_k}$.

4.3 Sample labeling decisions: $\xi_{i,k} \sim \text{Bern}(\pi_k^{(\tau)}(X_i))$, $i \in [n_k]$

4.4 Collect labels $\{Y_i : \xi_{i,k} = 1\}$

end for

— **Phase 3: Final Estimation** —

5: Compute stratified active estimator $(\hat{\theta}^{\pi_{\vec{\eta}}})$ (minimizer of Eq. (3))

626 **C Proofs**

Claim C.1 (Consistency of $\hat{\eta}$ for AdaStrat-ASI). *Suppose that for each stratum $k \in [K]$, the local budget fraction converges: $n_{b_k}/n_k \xrightarrow{P} q_k \in (0, 1)$. If \mathcal{H} is discrete and for every k there is no $\eta \in \mathcal{H}$ such that $\eta \mathbb{E}[u(X) | X \in \mathcal{A}_k] = q_k$ exactly, then there exists a unique optimal vector $\vec{\eta}^* = (\eta_1^*, \dots, \eta_K^*) \in \mathcal{H}^K$ such that:*

$$\mathbb{P}(\vec{\eta} \neq \vec{\eta}^*) \rightarrow 0$$

Proof. Proof of Claim C.1. For any fixed stratum k , the estimated parameter $\hat{\eta}_k$ is defined as:

$$\hat{\eta}_k = \max \left\{ \eta \in \mathcal{H} : \eta \frac{1}{n_k} \sum_{i \in \mathcal{A}_k} u(X_i) \leq \frac{n_{b_k}}{n_k} \right\}.$$

This condition is equivalent to:

$$\eta \leq \frac{n_{b_k}/n_k}{\frac{1}{n_k} \sum_{i \in \mathcal{A}_k} u(X_i)}.$$

627 In the numerator (by assumption), $\frac{n_{b_k}}{n_k} \xrightarrow{p} q_k$. In the denominator, by Law of Large Numbers
 628 (within stratum k), the empirical average converges in probability to the population expectation:
 629 $\frac{1}{n_k} \sum_{i \in \mathcal{A}_k} u(X_i) \xrightarrow{p} \mathbb{E}[u(X) \mid X \in \mathcal{A}_k]$.

Applying the Continuous Mapping Theorem,

$$\frac{n_{b_k}/n_k}{\frac{1}{n_k} \sum_{i \in \mathcal{A}_k} u(X_i)} \xrightarrow{p} \frac{q_k}{\mathbb{E}[u(X) \mid X \in \mathcal{A}_k]}.$$

Let the gap to the nearest incorrect parameter be:

$$\varepsilon_k = \min_{\eta \in \mathcal{H}} \left| \eta - \frac{q_k}{\mathbb{E}[u(X) \mid \mathcal{A}_k]} \right|.$$

By the assumption that no η exactly equals the target ratio, we have $\varepsilon_k > 0$. Define the true optimal parameter for stratum k as:

$$\eta_k^* = \max \{ \eta \in \mathcal{H} : \eta \mathbb{E}[u(X) \mid \mathcal{A}_k] \leq q_k \}.$$

On the event E_k that the empirical ratio $\frac{n_{b_k}/n_k}{\frac{1}{n_k} \sum_{i \in \mathcal{A}_k} u(X_i)}$ is within distance ε_k of the population ratio, we must have $\hat{\eta}_k = \eta_k^*$. Convergence in probability implies $\mathbb{P}(E_k^c) \rightarrow 0$ as $n_k \rightarrow \infty$. Since K is finite, we apply a union bound over all strata:

$$\mathbb{P}(\hat{\eta} \neq \bar{\eta}^*) = \mathbb{P} \left(\bigcup_{k=1}^K \{ \hat{\eta}_k \neq \eta_k^* \} \right) \leq \sum_{k=1}^K \mathbb{P}(\hat{\eta}_k \neq \eta_k^*) \rightarrow 0.$$

630 Thus, the estimated parameter vector converges exactly to the optimal vector with probability
 631 approaching 1. \square

632 *Proof for Theorem 5.3.* We assume the smoothness conditions of Assumption 5.1 hold. Let $n_k =$
 633 $\sum_{i=1}^n \mathbb{I}(X_i \in A_k)$ be the random number of unlabeled samples falling into stratum k . Because the
 634 total pool of n samples is drawn i.i.d. from the marginal distribution P_X , $n_k \sim \text{Binomial}(n, w_k)$. By
 635 the Weak Law of Large Numbers, $n_k/n \xrightarrow{p} w_k$.

The AdaStrat-ASI estimator $\hat{\theta}^{\hat{\eta}}$ is defined as the minimizer of the post-stratified empirical risk (Equation 3):

$$L_n^{\hat{\eta}}(\theta) := \sum_{k=1}^K w_k \left(\frac{1}{n_k} \sum_{i \in A_k} L_{\theta,i}^{\hat{\eta}_k} \right)$$

636 where $L_{\theta,i}^{\eta_k} = \ell_{\theta}(X_i, f(X_i)) + (\ell_{\theta}(X_i, Y_i) - \ell_{\theta}(X_i, f(X_i))) \frac{\xi_i}{\pi_{\eta_k}(X_i)}$.

Let $\vec{\eta} = (\eta_1, \dots, \eta_K) \in \mathcal{H}^K$ parameterize the piecewise sampling policy $\pi_{\vec{\eta}}(X_i) = \sum_{k=1}^K \mathbb{I}(X_i \in A_k) \pi_{\eta_k}(X_i)$. For a given $\vec{\eta}$, we define the single-sample surrogate loss as $L_{\theta,i}^{\vec{\eta}} = \ell_{\theta}(X_i, f(X_i)) + (\ell_{\theta}(X_i, Y_i) - \ell_{\theta}(X_i, f(X_i))) \frac{\xi_i}{\pi_{\vec{\eta}}(X_i)}$. We define the stratified empirical expectation and empirical process operators over the full sample of size n :

$$\mathbb{E}_n[g(L_{\theta}^{\vec{\eta}})] := \frac{1}{n} \sum_{k=1}^K \sum_{i \in A_k} g(L_{\theta,i}^{\vec{\eta}}) \quad ; \quad \mathbb{G}_n[g(L_{\theta}^{\vec{\eta}})] := \frac{1}{\sqrt{n}} \sum_{k=1}^K \sum_{i \in A_k} \left(g(L_{\theta,i}^{\vec{\eta}}) - \mathbb{E}[g(L_{\theta,i}^{\vec{\eta}}) \mid A_k] \right)$$

637 Notice that $\mathbb{E}_n[L_{\hat{\eta}}^*]$ exactly matches our stratified empirical risk $L^{\pi_{\hat{\eta}}}(\theta)$ from Eq. (3).

From C.1, we show uniform consistency of the vector-valued Split ($\hat{\eta}$). Because $\bar{\eta}^*$ is fixed, we can apply the standard empirical process bounds to the global risk. By the differentiability and local Lipschitzness of the loss, for any $h_n = O_P(1)$ we have stochastic equicontinuity:

$$\mathbb{G}_n \left[\sqrt{n}(L_{\theta^*+h_n/\sqrt{n}}^* - L_{\theta^*}^*) - h_n^\top \nabla L_{\theta^*}^* \right] \xrightarrow{P} 0$$

By definition, this is equivalent to:

$$n\mathbb{E}_n[L_{\theta^*+h_n/\sqrt{n}}^* - L_{\theta^*}^*] = n(L(\theta^* + h_n/\sqrt{n}) - L(\theta^*)) + h_n^\top \mathbb{G}_n[\nabla L_{\theta^*}^*] + o_P(1)$$

A second-order Taylor expansion around the population minimizer θ^* yields:

$$n\mathbb{E}_n[L_{\theta^*+h_n/\sqrt{n}}^* - L_{\theta^*}^*] = \frac{1}{2}h_n^\top H_{\theta^*} h_n + h_n^\top \mathbb{G}_n[\nabla L_{\theta^*}^*] + o_P(1)$$

Crucially, because $\mathbb{P}(\hat{\eta} \neq \bar{\eta}^*) \rightarrow 0$, we can substitute the empirical estimate $\hat{\eta}$ without altering the asymptotic behavior:

$$n\mathbb{E}_n[L_{\hat{\eta}+h_n/\sqrt{n}}^* - L_{\hat{\eta}}^*] = \frac{1}{2}h_n^\top H_{\theta^*} h_n + h_n^\top \mathbb{G}_n[\nabla L_{\theta^*}^*] + o_P(1)$$

We evaluate the previous display at $h_n = \hat{h}_n := \sqrt{n}(\hat{\theta}^{\hat{\eta}} - \theta^*)$ (which is $O_P(1)$ by consistency) and at $h_n = \tilde{h}_n := -H_{\theta^*}^{-1} \mathbb{G}_n[\nabla L_{\theta^*}^*]$. Following the algebraic manipulation from ASI:

$$\frac{1}{2} \left(\sqrt{n}(\hat{\theta}^{\hat{\eta}} - \theta^*) - \tilde{h}_n \right)^\top H_{\theta^*} \left(\sqrt{n}(\hat{\theta}^{\hat{\eta}} - \theta^*) - \tilde{h}_n \right) + o_P(1) \leq 0$$

638 Since $H_{\theta^*} \succ 0$, it must be that $\sqrt{n}(\hat{\theta}^{\hat{\eta}} - \theta^*) - \tilde{h}_n \xrightarrow{P} 0$.

639 We now analyze the gradient vector $\sqrt{n}\nabla L_n^*(\theta^*)$, accounting for the random n_k :

$$\sqrt{n}\nabla L_n^*(\theta^*) = \sum_{k=1}^K w_k \sqrt{n} \left(\frac{1}{n_k} \sum_{i \in A_k} \nabla L_{\theta^*,i}^{\eta_k^*} \right)$$

640 Let $\mu_k = \mathbb{E}[\nabla L_{\theta^*,i}^{\eta_k^*} | A_k]$. We rewrite the scaled gradient by multiplying and dividing by $\sqrt{n_k}$:

$$\sqrt{n}\nabla L_n^*(\theta^*) = \sum_{k=1}^K w_k \sqrt{\frac{n}{n_k}} \left(\frac{1}{\sqrt{n_k}} \sum_{i \in A_k} (\nabla L_{\theta^*,i}^{\eta_k^*} - \mu_k) \right)$$

On expanding,

$$\sqrt{n}\nabla L_n^*(\theta^*) = \sum_{k=1}^K w_k \sqrt{\frac{n}{n_k}} \left(\frac{1}{\sqrt{n_k}} \sum_{i \in A_k} (\nabla L_{\theta^*,i}^{\eta_k^*} - \mu_k) \right) + \sqrt{n} \sum_{k=1}^K w_k \mu_k$$

Because the estimator deterministically multiplies by the true population weights w_k , this sum exactly perfectly matches the population first-order condition $\sum_{k=1}^K w_k \mu_k = 0$. Therefore, the entire drift term vanishes algebraically before any variance is taken:

$$\sqrt{n} \sum_{k=1}^K w_k \mu_k = \sqrt{n}(0) = 0$$

641 Because this deterministic drift is forced to zero for every finite sample n , it cannot contribute
642 to the asymptotic variance. The variance of the remaining sum depends strictly on the centered
643 within-stratum gradients.

Now, by the standard Central Limit Theorem within each stratum,

$$Z_k := \frac{1}{\sqrt{n_k}} \sum_{i \in A_k} (\nabla L_{\theta^*, i}^{\eta_k^*} - \mu_k) \xrightarrow{d} \mathcal{N}(0, V_k)$$

644 , where $V_k = \text{Var}(\nabla L_{\theta^*, i}^{\eta_k^*} \mid X \in A_k)$.

645 By Slutsky's theorem, since $\sqrt{n/n_k} \xrightarrow{p} 1/\sqrt{w_k}$, we have $w_k \sqrt{n/n_k} Z_k \xrightarrow{d} \sqrt{w_k} \mathcal{N}(0, V_k) \equiv$
 646 $\mathcal{N}(0, w_k V_k)$.

Because the samples are independent across strata, the sum converges to a normal distribution with the sum of the covariances:

$$\sqrt{n} \nabla L_n^{\eta^*}(\theta^*) \xrightarrow{d} \mathcal{N}\left(0, \sum_{k=1}^K w_k V_k\right)$$

Applying this to the Taylor expansion yields $\sqrt{n}(\hat{\theta}^{\hat{\eta}} - \theta^*) \xrightarrow{d} \mathcal{N}(0, \Sigma_{\hat{\eta}^*})$, where the covariance is exactly the weighted sum of the within-stratum variances:

$$\Sigma_{\hat{\eta}^*} = H_{\theta^*}^{-1} \left(\sum_{k=1}^K w_k \text{Var} \left(\nabla \ell_{\theta^*}(X, f(X)) + \frac{\xi_{\eta_k^*}}{\pi_{\eta_k^*}(X)} \Delta \mid A_k \right) \right) H_{\theta^*}^{-1}$$

647 where $\Delta = \nabla \ell_{\theta^*}(X, Y) - \nabla \ell_{\theta^*}(X, f(X))$. The final statement follows by Slutsky's theorem. \square

648 **Lemma C.2.** *Assume standard M-estimation regularity from Theorem 5.3. In Algorithm 1 the*
 649 *AdaStrat-ASI estimator's asymptotic variance's trace admits the decomposition*

$$\text{Tr}(\Sigma_{\hat{\eta}^*}) = C + n \sum_{k=1}^K \frac{w_k^2}{n_{b_k}} R_k^*, \quad (5)$$

650 where C does not depend on $(n_{b_1}, \dots, n_{b_K})$, and

$$R_k^* := \text{Tr}(H_{\theta^*}^{-1} W_k^* H_{\theta^*}^{-1}), \quad W_k^* := \mathbb{E} \left[\frac{1}{s_k(X)} \Delta \Delta^\top \mid X \in A_k \right]. \quad (6)$$

651 , where,

$$\Delta := \nabla \ell_{\theta^*}(X, Y) - \nabla \ell_{\theta^*}(X, f(X)).$$

Proof. We directly start from result of Theorem 5.3

$$\Sigma_{\hat{\eta}^*} = H_{\theta^*}^{-1} \left(\sum_{k=1}^K w_k V_k \right) H_{\theta^*}^{-1}$$

We now decompose V_k . By the law of total variance conditioned on $X \in A_k$, and using $\xi_i \sim \text{Bernoulli}(\pi_{\eta_k^*}(X_i))$, we separate the variance of the true gradient from the variance induced by active sampling:

$$V_k = \text{Var}(\nabla \ell_{\theta^*}(X, Y) \mid A_k) + \mathbb{E} \left[\left(\frac{1}{\pi_{\eta_k^*}(X)} - 1 \right) \Delta \Delta^\top \mid A_k \right]$$

By our active sampling design, the policy inside stratum k is $\pi_{\eta_k^*}(X) = \frac{n_{b,k}}{n_k} s_k(X)$. Substituting this yields:

$$V_k = \text{Var}(\nabla \ell_{\theta^*}(X, Y) \mid A_k) - \mathbb{E}[\Delta \Delta^\top \mid A_k] + \frac{n_k}{n_{b,k}} \mathbb{E} \left[\frac{1}{s_k(X)} \Delta \Delta^\top \mid A_k \right]$$

Asymptotically, $n_k/n \rightarrow w_k$, so the scaling factor becomes $\frac{n_k}{n_{b,k}} \rightarrow \frac{nw_k}{n_{b,k}}$. The limiting within-stratum variance becomes:

$$V_k = \tilde{C}_k + \frac{nw_k}{n_{b,k}} W_k^*$$

where $\tilde{C}_k = \text{Var}(\nabla \ell_{\theta^*}(X, Y) \mid A_k) - \mathbb{E}[\Delta \Delta^\top \mid A_k]$ is the irreducible variance independent of the active labeling budget. Finally, we substitute this back into the global covariance matrix and apply the trace operator $\text{Tr}(\cdot)$:

$$\text{Tr}(\Sigma_{\tilde{\eta}^*}) = \sum_{k=1}^K w_k \text{Tr}(H_{\theta^*}^{-1} \tilde{C}_k H_{\theta^*}^{-1}) + n \sum_{k=1}^K \frac{w_k^2}{n_{b,k}} \text{Tr}(H_{\theta^*}^{-1} W_k^* H_{\theta^*}^{-1})$$

Let $C = \sum_{k=1}^K w_k \text{Tr}(H_{\theta^*}^{-1} \tilde{C}_k H_{\theta^*}^{-1})$. Substituting the definition $R_k^* := \text{Tr}(H_{\theta^*}^{-1} W_k^* H_{\theta^*}^{-1})$, we obtain the final decomposition:

$$\text{Tr}(\Sigma_{\tilde{\eta}^*}) = C + n \sum_{k=1}^K \frac{w_k^2}{n_{b,k}} R_k^*$$

652 This completes the proof.

653

□

654 Proof for Proposition 5.4.

655 *Proof.* Under Lemma C.2, we formulate the constrained optimization problem using the method
656 of Lagrange multipliers. The Lagrangian \mathcal{L} is defined with respect to the decision variables n_{b_k} :

657 $\mathcal{L}(n_{b_1}, \dots, n_{b_K}, \lambda) = n \sum_{k=1}^K \frac{w_k^2 R_k^*}{n_{b_k}} + \lambda \left(\sum_{k=1}^K n_{b_k} - n_b \right)$ where λ is the Lagrange multiplier
658 enforcing the budget constraint. To find the critical points, we take the partial derivative with
659 respect to each n_{b_k} and set it to zero: $\frac{\partial \mathcal{L}}{\partial n_{b_k}} = -\frac{n w_k^2 R_k^*}{n_{b_k}^2} + \lambda = 0$ Rearranging the terms to isolate

660 the decision variable n_{b_k} : $n_{b_k}^2 = \frac{n w_k^2 R_k^*}{\lambda} \implies n_{b_k} = \sqrt{n} \frac{w_k \sqrt{R_k^*}}{\sqrt{\lambda}}$ This relation holds for all k .

661 We determine the constant λ by summing n_{b_k} over all strata and applying the budget constraint:

662 $\sum_{k=1}^K n_{b_k} = \sum_{k=1}^K \sqrt{n} \frac{w_k \sqrt{R_k^*}}{\sqrt{\lambda}} = \frac{\sqrt{n}}{\sqrt{\lambda}} \sum_{k=1}^K w_k \sqrt{R_k^*} = n_b$. Solving for the scaling factor $\frac{\sqrt{n}}{\sqrt{\lambda}}$:
663 $\frac{\sqrt{n}}{\sqrt{\lambda}} = \frac{n_b}{\sum_{j=1}^K w_j \sqrt{R_j^*}}$. Finally, substituting this back into the expression for n_{b_k} yields the optimal

664 allocation rule: $n_{b_k}^* = w_k \sqrt{R_k^*} \left(\frac{n_b}{\sum_{j=1}^K w_j \sqrt{R_j^*}} \right)$

665 One further notices, the objective is strictly convex in n_{b_k} over $n_{b_k} > 0$ because each term has a
666 second derivative $\left(\frac{2n w_k R_k^*}{n_{b_k}^3} > 0 \right)$. So any critical point is the unique global minimizer.

667

□

668 *Proof.* Proof for Lemma 5.5.

669 By definition, $\epsilon_k := \mathbb{E}[(Y - f(X))^2 \mid X \in \mathcal{A}_k]$.

670 Recall our notations,

671 $H_{\theta^*} := \mathbb{E}[\nabla^2 \ell_{\theta^*}(X, Y)]$, $\Delta := \nabla \ell_{\theta^*}(X, Y) - \nabla \ell_{\theta^*}(X, f(X))$, $W_k^* := \mathbb{E} \left[\frac{1}{s_k(X)} \Delta \Delta^\top \mid X \in \mathcal{A}_k \right]$,

672 $R_k^* := \text{Tr}(H_{\theta^*}^{-1} W_k^* H_{\theta^*}^{-1})$

673 We know, $\pi_k(x) = \frac{n_{b_k}}{n_k} s_k(x)$. On τ -mixing,

$$\begin{aligned} \pi_k^\tau(x) &= \frac{n_{b_k}}{n_k} s_k(x)(1 - \tau) + \frac{n_{b_k}}{n_k} \tau \\ &= \frac{n_{b_k}}{n_k} (s_k(x)(1 - \tau) + \tau) \\ &= \frac{n_{b_k}}{n_k} (s_k^\tau(x)) \end{aligned}$$

674 This implies $\forall x$,

$$s_k^\tau(x) > \tau \tag{7}$$

Also, from Assumption 5.2, gradient is Lipschitz in the target y : there exists a $L_y(X)$ with $\mathbb{E}[L_y(X)^2 | \mathcal{A}_k] \leq L_{y,k}^2 < \infty$ such that for all y, y' ,

$$\|\nabla \ell_{\theta^*}(X, y) - \nabla \ell_{\theta^*}(X, y')\| \leq L_y(X) |y - y'|$$

675 Thus, by Lipschitz-in- y , $\|\Delta\| \leq L_y(X) |Y - f(X)|$.

676 Squaring and taking stratumwise expectation,

$$\begin{aligned} \mathbb{E}[\|\Delta\|^2 | \mathcal{A}_k] &\leq \mathbb{E}[L_y(X)^2 (Y - f(X))^2 | \mathcal{A}_k] \\ &\leq L_{y,k}^2 \mathbb{E}[(Y - f(X))^2 | \mathcal{A}_k] \\ &= L_{y,k}^2 \epsilon_k \end{aligned} \tag{8}$$

677 Now,

$$\begin{aligned} \|W_k^*\|_{\text{op}} &\leq \mathbb{E}\left[\frac{1}{s_k^{\tau}(X)} \|\Delta \Delta^{\top}\|_{\text{op}} | \mathcal{A}_k\right] \\ &= \mathbb{E}\left[\frac{1}{s_k^{\tau}(X)} \|\Delta\|^2 | \mathcal{A}_k\right] \\ &\leq \frac{1}{\tau} \mathbb{E}[\|\Delta\|^2 | \mathcal{A}_k] \\ &\leq \frac{L_{y,k}^2}{\tau} \epsilon_k. \end{aligned}$$

678 Third line and fourth line are true because of Eq. 7 and Eq. 8 respectively.

679 Finally, assuming $\|H_{\theta^*}^{-1}\|_F < \infty$,

$$R_k^* = \text{Tr}(H_{\theta^*}^{-1} W_k^* H_{\theta^*}^{-1}) \leq \|H_{\theta^*}^{-1}\|_F^2 \|W_k^*\|_{\text{op}} \leq \frac{\|H_{\theta^*}^{-1}\|_F^2}{\tau} L_{y,k}^2 \epsilon_k$$

Since R_k^* and ϵ_k are non negative, for all k

$$\sqrt{R_k^*} = \mathcal{O}(\sqrt{\epsilon_k})$$

680

□

681 *Proof.* Proof for Theorem 5.6.

682 The idea of strata exists explicitly and implicitly in AdaStrat-ASI and ASI respectively; AdaStrat-ASI
683 uses pre-stratification(3), whereas we analyze ASI under a post-stratification((1)) lens under the
684 same strata composition.

Let $n_k^{\text{ASI}} := \sum_{i=1}^n \mathbf{1}\{X_i \in A_k\}$ be the random number of samples that fall into stratum k , such that $\sum n_k^{\text{ASI}} = n$. Let $\hat{w}_k = \frac{n_k^{\text{ASI}}}{n}$ be the empirical weight. Because the total pool is drawn i.i.d., the vector of counts follows a Multinomial distribution: $(n_1^{\text{ASI}}, \dots, n_K^{\text{ASI}}) \sim \text{Multinomial}(n, (w_1, \dots, w_K))$. By the properties of the Multinomial distribution, the empirical weights \hat{w}_k fluctuate around the true weights w_k . By the Central Limit Theorem, the scaled fluctuations converge to a Normal distribution:

$$\sqrt{n}(\hat{w} - w) \xrightarrow{d} \mathcal{N}(0, \Sigma_{\text{multinomial}})$$

685 where the covariance matrix of the weights has entries $\text{Cov}(\hat{w}_j, \hat{w}_k) = w_j(1 - w_j)$ if $j = k$, and
686 $-w_j w_k$ if $j \neq k$.

687 Let $\mu_k = \mathbb{E}[\nabla L_{\theta^*, i}^k | A_k]$ be the true mean gradient in stratum k . The population first-order optimality
688 condition guarantees that exactly $\sum_{k=1}^K w_k \mu_k = 0$.

Standard ASI (Equation 1) computes the unweighted global average. We scale it by \sqrt{n} to find its asymptotic distribution:

$$\sqrt{n} \nabla L_n^{\text{ASI}}(\theta^*) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \nabla L_i$$

We can group this sum into the K strata, and center the inner sums around μ_k :

$$\sqrt{n}\nabla L_n^{\text{ASI}}(\theta^*) = \sum_{k=1}^K \frac{1}{\sqrt{n}} \sum_{i \in A_k} (\nabla L_i - \mu_k) + \sum_{k=1}^K \frac{1}{\sqrt{n}} \sum_{i \in A_k} \mu_k$$

On closely inspecting the second term, the sum $\sum_{i \in A_k} \mu_k$ is just $n_k^{\text{ASI}} \mu_k$.

$$\text{Second Term} = \sum_{k=1}^K \frac{n_k^{\text{ASI}}}{\sqrt{n}} \mu_k = \sqrt{n} \sum_{k=1}^K \left(\frac{n_k^{\text{ASI}}}{n} \right) \mu_k = \sqrt{n} \sum_{k=1}^K \hat{w}_k \mu_k$$

Because $\sum_{k=1}^K w_k \mu_k = 0$, we can subtract it from this term without changing its value:

$$\text{Second Term} = \sqrt{n} \sum_{k=1}^K (\hat{w}_k - w_k) \mu_k$$

This term represents the error caused by the random sample sizes (n_k). Because $\sqrt{n}(\hat{w} - w)$ converges to the Multinomial covariance matrix, the variance of this linear combination is:

$$\text{Var} \left(\sum_{k=1}^K \sqrt{n}(\hat{w}_k - w_k) \mu_k \right) = \sum_{k=1}^K w_k \mu_k \mu_k^\top - \left(\sum_{k=1}^K w_k \mu_k \right) \left(\sum_{k=1}^K w_k \mu_k \right)^\top$$

Because the second bracket is strictly 0 (the first-order condition), the variance of the second term is exactly:

$$\sum_{k=1}^K w_k \mu_k \mu_k^\top \equiv \Delta_{\text{Bias}} > 0$$

689 Thus, ASI's empirical gradient contains the term $\sqrt{n} \sum \hat{w}_k \mu_k$. Because \hat{w}_k is a random variable, this
 690 term fluctuates, injecting exactly Δ_{Bias} into the final asymptotic covariance matrix. By forcing the
 691 estimator to use w_k instead of the empirical \hat{w}_k , AdaStrat-ASI algebraically destroys the multinomial
 692 fluctuation term before the variance is taken.

693 Now, performing a symmetrical k -stratum decomposition analysis on ASI's asymptotic variance as
 694 seen in Proof of Lemma C.2,

$$\begin{aligned} \text{Tr}(\Sigma_{\text{ASI}}) &= \Delta_{\text{Bias}} + T \left(\sum_{k=1}^K w_k (C_k + n_k \mathbb{E} \left[\frac{1}{n_{b_k}^{\text{ASI}}} \right] \mathbb{E} \left[\frac{\Delta \Delta^\top}{s_k(X)} \mid Z = k \right] \right) \\ &= \Delta_{\text{Bias}} + C + n \sum_{k=1}^K w_k^2 R_k^* \mathbb{E} \left[\frac{1}{n_{b_k}^{\text{ASI}}} \right] \end{aligned} \quad (9)$$

695 Subtracting Eq. (5) from Eq. (9),

$$\begin{aligned} \text{Gap} &= \text{Tr}(\Sigma_{\text{ASI}}) - \text{Tr}(\Sigma_{\text{AdaStrat-ASI}}) \\ &= \underbrace{\Delta_{\text{Bias}}}_{\text{Stability Gap(Term 1)}} + n \cdot \underbrace{\sum_{k=1}^K w_k^2 R_k^* \left(\mathbb{E} \left[\frac{1}{n_{b_k}^{\text{ASI}}} \right] - \frac{1}{n_{b_k}^{\text{AdaStrat-ASI}}} \right)}_{\text{Efficiency Gap(Term 2)}} \end{aligned}$$

696 The gap shown above reveals that AdaStrat-ASI reduces asymptotic variance through two distinct
 697 mechanisms: *stability* (via stratification) and *efficiency* (via allocation). We analyze the conditions
 698 under which the gap $\text{Tr}(\Sigma_{\text{ASI}}) - \text{Tr}(\Sigma_{\text{AdaStrat-ASI}})$ is strictly positive.

699 **1. Term 1:** The first term compares reduction of between-stratum variance. Since Δ_{Bias} is the trace
 700 of a quadratic form involving a covariance matrix, it is strictly non-negative.

2. Term 2: The second term compares the allocation strategies. By using Jensen's inequality, on convex function $f(x) = 1/x$, we have:

$$\sum_{k=1}^K \gamma_k \mathbb{E} \left[\frac{1}{n_{b_k}^{\text{ASI}}} \right] \geq \underbrace{\sum_{k=1}^K \frac{\gamma_k}{\mathbb{E}[n_{b_k}^{\text{ASI}}]}}_{\text{(Fixed Average)}}$$

701 The above inequality implies that removing the randomness of the sample counts (fixing $n_{b_k}^{\text{ASI}}$ to its
702 mean) reduces variance.

703 But it is given that,

$$\sum_{k=1}^K \frac{\gamma_k}{\mathbb{E}[n_{b_k}^{\text{ASI}}]} \geq \sum_{k=1}^K \frac{\gamma_k}{n_{b_k}^{\text{AdaStrat-ASI}}}$$

This implies,

$$\text{Tr}(\Sigma_{\text{ASI}}) > \text{Tr}(\Sigma_{\text{AdaStrat-ASI}})$$

704 Hence proved. \square

705 **Lemma C.3.** Assuming non-degenerate difficulty scores $0 < \hat{\sigma}_k^f < \infty \quad \forall k$. In Algorithm 1, as
706 the model becomes perfect ($\epsilon^* \rightarrow 0$), the asymptotic variance of AdaStrat-ASI converges to the
707 irreducible variance C at a rate of $\mathcal{O}(\epsilon^*)$.

Proof.

$$\begin{aligned} \text{Tr}(\Sigma_{\bar{\eta}^*}) &= C + n \sum_{k=1}^K \frac{w_k^2}{n_{b_k}} R_k^* \\ &= C + n \sum_{k=1}^K \frac{w_k^2}{n_b \left(\frac{w_k \sqrt{\hat{\sigma}_k^f}}{C'} \right)} R_k^* \quad \left(C' := \sum_{j=1}^K w_j \sqrt{\hat{\sigma}_j^f} \right) \\ &= C + \frac{n C'}{n_b} \sum_{k=1}^K w_k \frac{R_k^*}{\sqrt{\hat{\sigma}_k^f}} \\ &= C + \mathcal{O} \left(\sum_{k=1}^K w_k \frac{\epsilon_k}{\sqrt{\hat{\sigma}_k^f}} \right). \\ &= C + \mathcal{O}(\epsilon^*) \end{aligned}$$

708 The first line follows from Lemma C.2. The second line is true asymptotically.⁴ The fourth line
709 follows from Lemma 5.5.

710 This also confirms, the variance is bounded by the weighted average of the Noise-to-Signal ratio. \square

711 D Validity of Inference

Our AdaStrat-ASI estimator (Eq.(3)) depends on the attacking-phase sampling probabilities $\pi_k(x)$ only through their realized values in the denominator. When we write:

$$\pi_k(x) = \frac{n_{b_k}}{n_k} \cdot \frac{u(x)}{\mathbb{E}[u(X) | A_k]}$$

712 the budget fraction n_{b_k}/n_k enters as a deterministic scaling constant in the asymptotic limit (by
713 assumption: $n_{b_k}/n_k \rightarrow q_k \in (0, 1)$; see Claim C.1). Our proof for Theorem 5.3 does not depend

⁴Since the budget allocation follows a multinomial distribution, $n_{b_k}/n_b \rightarrow z_k$ almost surely in expectation. Thus, for the asymptotic variance analysis, we treat the allocation as deterministic: $n_{b_k} = n_b z_k$.

714 whether this q_k comes from Oracle Neyman allocation ($\propto \sqrt{R_k^*}$), informed empirical assistance(
715 i.e. Scouting Stage), random allocation (z_k arbitrary) or failed scouting (any convergent sequence);
716 the structure of the asymptotic variance formula remains identical. Fundamentally we care about
717 CLT holding within each fixed strata: $n_k \rightarrow \infty$ (consequence of $n \rightarrow \infty$) and that the sampling
718 probabilities $\pi_{\eta_k^*}(X_i)$ are bounded away from zero. We only care about the ratio $\sqrt{n/n_k} \rightarrow 1/\sqrt{w_k}$ (
719 by Slutsky’s Theorem) based on the stratum prevalence, not the budget allocation mechanism. Overall,
720 the numerical value of q_k changes, affecting efficiency, not validity of our method.

721 E Experiments

722 E.1 Training Details

723 For AlphaFold Dataset, each observation corresponds to a post-translationally modified protein
724 residue, where the binary response $Y_i \in \{0, 1\}$ indicates whether residue i belongs to an exper-
725 imentally validated intrinsically disordered region (IDR). Crucially, AlphaFold2’s pLDDT-based
726 disorder predictions $\hat{Y}_i \in [0, 1]$ are provided directly by `ppi_py`, requiring no additional proxy
727 model training. However for the other real-world datasets, we follow a consistent sample-splitting
728 protocol as in Zrníć and Candès [65]. Specifically, 50% of the data is used to learn models that guide
729 annotation allocation—including the predictive model in ASI(XGBoost [9]) and the auxiliary error
730 model used to calibrate LLM(GPT-4O [47]) confidence while statistical inference is performed
731 exclusively on the remaining half. To clarify, for the latter, the LLM itself is not trained on the data,
732 the error model that maps verbalized confidence scores to squared annotation error is learned on
733 the first half and evaluated only on the held-out half(CCI; Section 3.2 from Gligoric et al. [22]).
734 This separation preserves the inferential validity guarantees of ASI. For simplicity sake, we report
735 confidence-intervals(CIs) and coverage for the first parameter with $\alpha = 0.1$. We repeat our analysis
736 for 500 trials. All experiments were run on a single RTX A6000.

737 E.2 Additional Real Data Experiment

738 E.2.1 California Housing Dataset

739 For this setup, we use the California Housing dataset ($n = 20640$) [48] with the goal to predict the
740 median house value for California districts(expressed in hundreds of thousands of dollars), using
741 demography, location, and general information, while minimizing over squared loss.

742 From Figure 7, all methods achieve coverage close to the nominal 90% level, confirming inferential
743 validity. Notably, among ML-enabled approaches, AdaStrat-ASI consistently maintains best fidelity
744 to target coverage. Here, label usage is nearly identical across methods for a given human fraction,
745 ruling out differences in annotation cost as the primary driver of performance gaps.

746 Despite this, AdaStrat-ASI consistently achieves the smallest confidence interval width across all
747 budgets. Relative to ASI, this corresponds to an approximate 10%~14% reduction in CI width across
748 most human fractions, with larger gains at smaller budgets. Moreover, AdaStrat-ASI selects labeled
749 points with the largest diameter (from Eq. (2)), i.e., the widest spread across the data manifold.

750 E.3 Runtime Overhead

751 AdaStrat-ASI adds a small computational overhead primarily due to the the zero-label scouting phase
752 compared to ASI. For instance on the Housing Dataset, ASI takes 7.336 ± 0.369 ms and AdaStrat-ASI
753 takes 8.565 ± 0.189 ms.

754 E.4 Guidelines on stratifying data

755 The effectiveness of stratified inference critically depends on how strata are defined. In this section
756 we shed light on practical ways to identify strata in the inference data. As is standard in stratified
757 sampling, strata with very small sample sizes are merged with neighboring or similar strata to avoid
758 unstable variance estimates and degenerate allocations [39, 13].

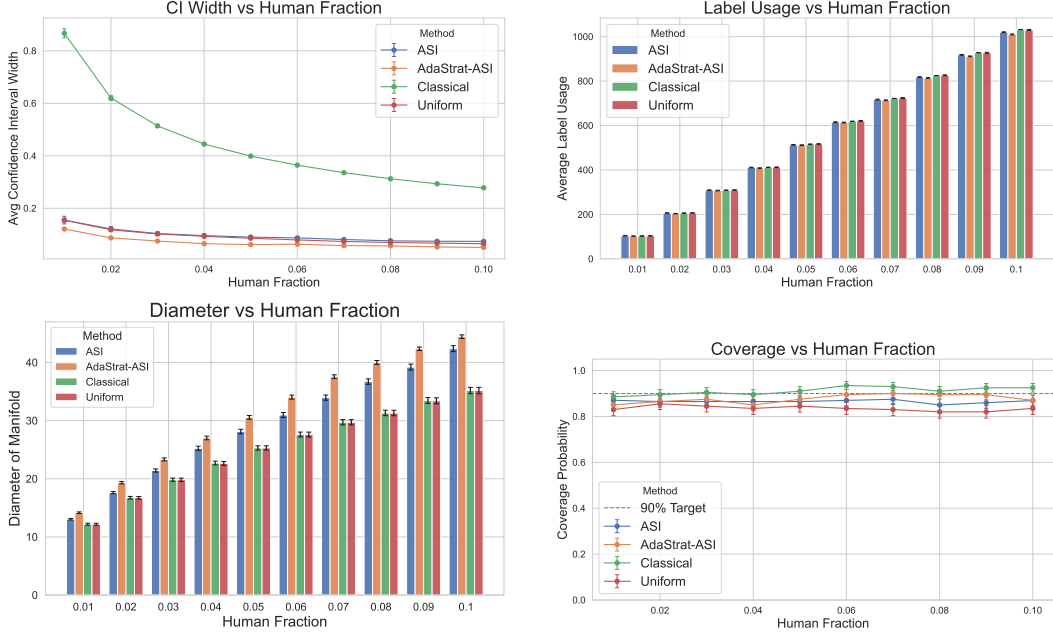


Figure 7: Housing Dataset

759 E.4.1 Domain Knowledge

760 In many real-world applications the choice of stratifying variable is not uniquely prescribed. Instead,
 761 it often hinges on practitioner intuition and domain expertise. Practitioners may select stratification
 762 variables based on prior empirical findings, exploratory analyses, or contextual understanding of the
 763 data-generating process. Such judgment-driven stratification leverages known scientific, biological,
 764 or social structure to partition the data into subpopulations that are expected to exhibit relatively
 765 homogeneous response behavior. For instance, in biomedical and clinical studies, stratifying by
 766 sex, age group, tissue type, or disease subtype is often essential, as these factors are known to
 767 induce systematic differences in outcomes and variability [11, 46]. Similarly, in social science and
 768 policy analysis, stratification by demographic or geographic categories (e.g., race, gender, region) is
 769 routinely employed to account for heterogeneity in behavior and responses [27]. In our framework,
 770 these domain-driven strata provide a natural and effective starting point.

771 For stratification in the California Dataset (Section E.2.1), we use the categorical variable
 772 `ocean_proximity`. Figure 8 shows the comparison between normalized trivial weight-based alloca-
 773 tion and weighted model-based allocation across strata for the California Housing dataset.

774 First, our proxy distribution is markedly non-uniform, indicating that the black-box model identifies
 775 systematic differences across strata and is non-trivially informative. Such structure validates the
 776 use of proxy-based scouting in AdaStrat-ASI: the allocation reflects meaningful variation in local
 777 difficulty, which is subsequently exploited during the attacking stage to reallocate a fixed labeling
 778 budget more effectively. Second, our proxy allocation differs from population-proportional sampling
 779 by a total variation distance of 0.093, indicating a considerable redistribution of sampling mass across
 780 four strata.

781 We take this opportunity to underline the importance of measuring the difficulty scores ($\hat{\sigma}_k^f$) in the
 782 Scouting Stage. Instead of allowing the scouting phase to determine allocation weights i.e. Figure 1,
 783 we use an allocation vector purely based on strata size i.e. Figure 8(Left). We use the former allocation
 784 method in Section E.2.1. Comparing results in Figure 9 to the results presented in Section E.2.1,
 785 there is not much change in width, label-usage, and diameter. However, with purely weight-based
 786 allocation, coverage is severely impacted.

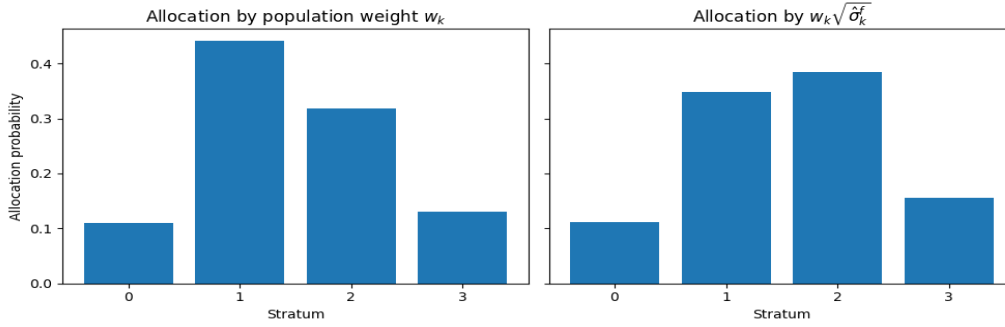


Figure 8: Influence of difficulty scores in Label Allocation

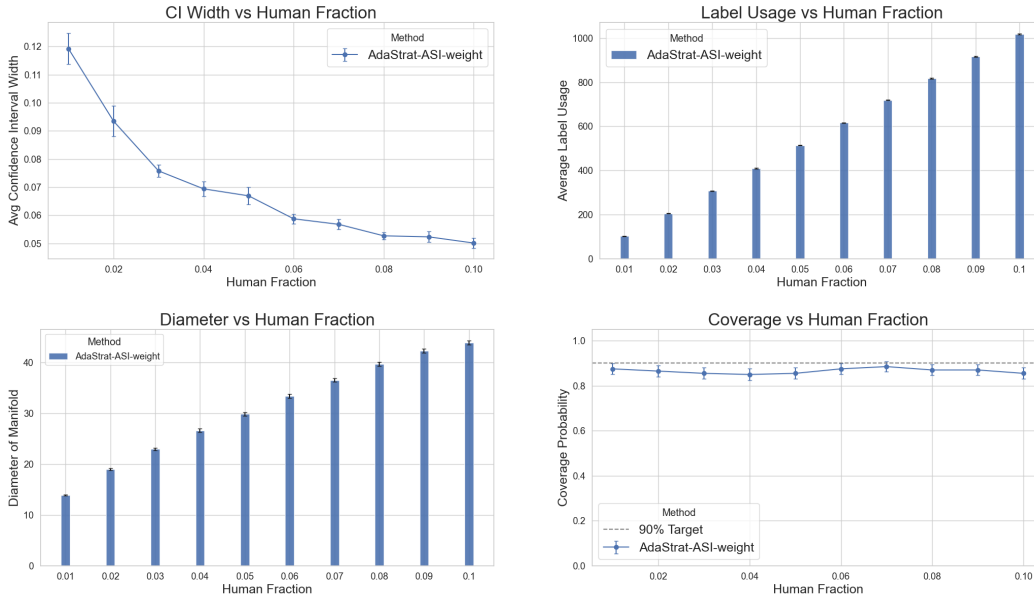


Figure 9: Impact of pure weight-based label allocation in Housing Dataset

787 E.4.2 Black Box Model based Stratification

788 In some applications, domain knowledge is unavailable or difficult to exploit, rendering intuition-
 789 driven stratification infeasible. This commonly arises when practitioners are provided with
 790 anonymized or abstracted datasets, where features lack semantic interpretation due to privacy, propri-
 791 etary constraints, or institutional policies [18, 32].

792 To address this scenario, we adopt a model-based stratification strategy that leverages cheap proxy
 793 information provided by a black-box predictive model f , which is already central to the ASI frame-
 794 work. Concretely, we stratify the unlabeled pool by jointly considering the model’s prediction’s
 795 $f(x)$ and uncertainty $u(x)$, yielding partitions that reflect both the model’s belief and its surety. Our
 796 stratification approach is loosely inspired by Section 5.2 of Fisch et al. [19].

797 Intuitively, regions where the model is highly confident are treated as coarse strata because the
 798 model’s predictions are expected to be reliable proxies for the true response, and additional labels
 799 are less informative. Conversely, regions of high uncertainty are partitioned more finely, benefiting
 800 from targeted labeling. This construction reinforces our ideology: if a strata is problematic (i.e., high
 801 uncertainty), chunk it into smaller fragments to minimize its contribution to global variance (w_k). A
 802 similar argument applies to constructing larger strata where the model’s predictions are reliable.

803 In practice, this can be implemented via grid-based stratification over the joint space of $(f(x), u(x))$.
 804 Breakpoints along each axis may be chosen using simple heuristics such as quantiles, histogram

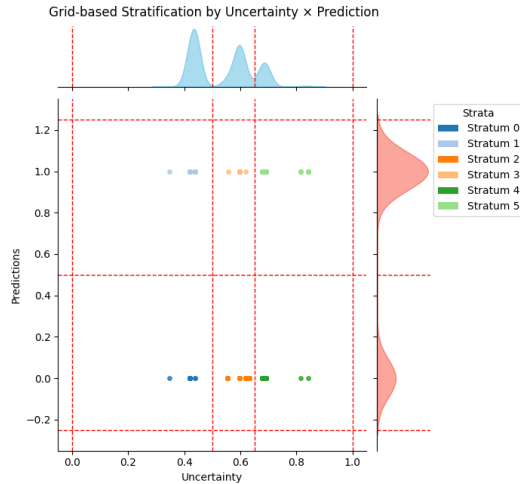


Figure 10: Stratification in Politeness Dataset

805 modes, or change-point detection, resulting in rectangular strata whose sizes shrink as uncertainty
 806 increases. Figure 10 illustrates this procedure on the Politeness dataset (Section 6).

807 E.5 Impact of Model Sensitivity on Label Allocation

808 To gauge impact on label allocation on $f(x)$ we extend our synthetic 1D example from Section 6.
 809 The setup remains exactly identical i.e. in Stratum 2, the predictor stays perfect: $f(x) = 0$, except in
 810 Stratum 1, we deliberately make the predictions noisy. $f(x) \sim 1 + \mathcal{N}(0, \sigma^2)$.

811 We run a small experiment for a fixed budget of 500 labels and $\sigma^2 = 1$. Because our scouting phase
 812 realizes Stratum 2 is perfectly solved ($f(x) = \hat{Y} = 0$), $\hat{\sigma}_2^f$ comes out to be 0, same as in the previous
 813 setting. Consequently, 0 budget labels were assigned to Stratum 2. However, in the case of Stratum 1,
 814 because of the artificial noise injection, $\hat{\sigma}_1^f \neq 0$ anymore (0.9708 to be precise). This causes Stratum
 815 1 to be assigned all 500 labels.

816 This setup perfectly demonstrates the impact on model sensitivity on the scouting stage. Thus, the
 817 proxy variance $\hat{\sigma}_1^f$ correctly detected the instability, and AdaStrat-ASI dynamically shifted the entire
 818 labeling budget to Stratum 1.

819 E.6 Random Stratification Ablation

820 As formalized in our variance reduction guarantee (Theorem 5.6), our method strictly requires the
 821 strata to capture latent heterogeneity. While we suggest ways to stratify data E.4, if the chosen strata
 822 are uninformative, AdaStrat-ASI’s performance can be severely affected. the budget allocation ($\hat{\sigma}_k^f$)
 823 approaches uniformity.

824 To empirically demonstrate this, we run an ablation experiment where we construct \mathcal{A} using a purely
 825 random split of the data into K equally sized strata (for simplicity), ignoring domain knowledge and
 826 model confidence entirely. We set $K = 6$ for direct comparison, on the Politeness Dataset (Section
 827 6).

828 From Figure 11 (Right), the proxy difficulty scores across these random strata are nearly identical. In
 829 this case, $n_{b_k} \propto w_k \approx \frac{1}{K}$. The allocation weights that serve input to the multinomial distribution
 830 follow an almost uniform distribution ($\approx \frac{1}{6}$). Stratification was clearly not *useful* here.

831 From Figure 14 compared to the other baselines discussed in Section E.4 and illustrated in Figure
 832 6, label usage and diameter for random allocation are quite similar. Compared to ASI, we still
 833 reduce CI-width. Interestingly, per our experimental design γ_k for all strata is roughly the same,
 834 which is why the efficiency gap is not exploitable. But, the reduction in CI widths comes from the

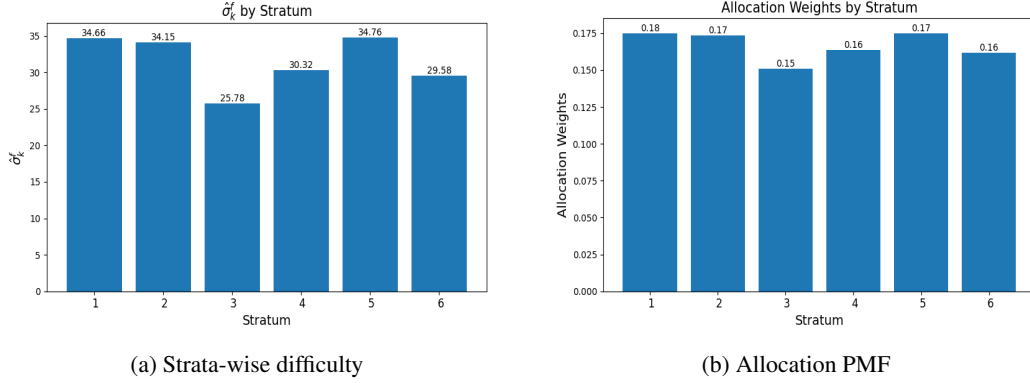


Figure 11: Impact of Random Stratification on Budget Allocation

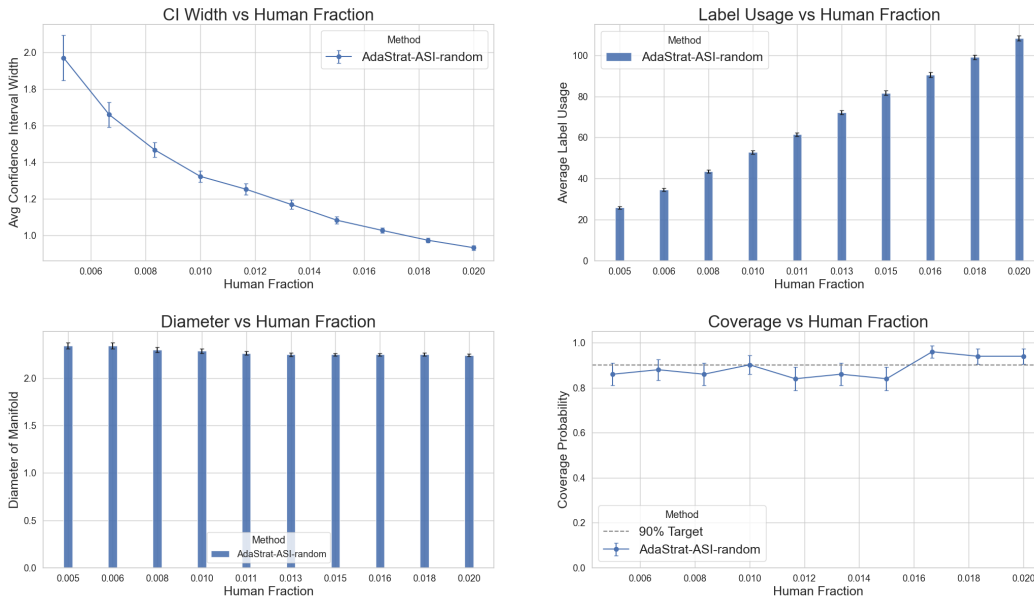


Figure 12: Impact of Random Stratification on Politeness Dataset

835 post-stratification structure i.e. Δ_{Bias} . Thus, random stratification kills the Efficiency Gap (budget is
 836 flat), but preserves the Stability Gap (structural post-stratification).

837 Compared to Informed Stratification, the width is almost identical, except for a slight difference in
 838 the stretch range of 0.01 – 0.013, where informed Stratification(ours) provides a minor improvement.
 839 However, random stratification falls short in terms of coverage.

840 **Poor Coverage.** Based on our random stratification setup, each Strata 1-6 contains a random mix of
 841 easy points and hard points. Now, when the Attacking Stage begins, inside every stratum, it samples
 842 points using active sampling policy $\pi_k(x) \propto u(x)$. Because strata are a random mix of everything,
 843 it contains some points with massive uncertainty $u(x)$ and many points with tiny uncertainty. The
 844 algorithm will heavily over-sample the high-uncertainty points and largely ignore the low-uncertainty
 845 points inside that bucket.

846 Since Equation (3) follows AIPW, we divide by the sampling probability: $\frac{1}{\pi_k(x)}$. Since the budget
 847 allocation is even, the algorithm in expectation is designed to exhaust their respective respective
 848 strata budgets. This design makes it highly likely that low-uncertainty points (where $\pi_k(x)$ is very
 849 small) are randomly drawn. Naturally, those point gets a massive weight in the final average. Because
 850 strata are a random mix of everything, the normalization factor $\mathbb{E}[u(X) | A_k]$ is essentially the global

851 average uncertainty. The probabilities $\pi_k(x)$ stretch from very near 1 to very near 0. This creates
 852 extreme propensity weights.

853 When an estimator has extreme propensity weights, under limited labeling budget, the asymptotic
 854 Central Limit Theorem takes a very long time to kick in. The empirical variance estimate used to
 855 construct the CIs becomes highly unstable and downward-biased. Hence, the algorithm calculates
 856 tight CIs, but because the weights are so erratic, the actual estimates bounce around wildly outside
 857 that CI.

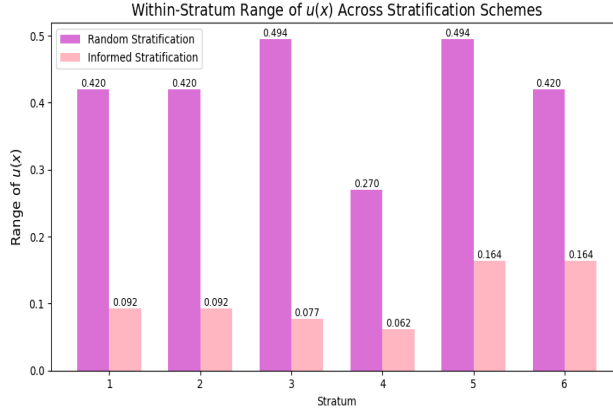


Figure 13: Range of Uncertainty across Strata

858 Figure 13 compares the within-stratum range of $u(x)$ under random and informed stratification. Since
 859 the two stratification procedures induce different partitions of the data, the stratum indices are not
 860 directly comparable across methods; the plot is intended only to contrast overall within-stratum
 861 homogeneity. The difference in performance goes on to show that informed stratification, in fact,
 862 facilitates the emergence on CLT within each stratum.

863 E.7 Random Allocation Ablation

We retain the informed stratification setup from Section 6. However, instead of allowing the scouting
 phase to determine allocation weights(Figure 1), we manually inject a purely randomly generated
 allocation vector:

$$\vec{z} = [0.1734, 0.0821, 0.2567, 0.1349, 0.0913, 0.2616]$$

864 These weights are generated independently of the data (summing to 1.0 via normalization) and
 865 represent a “failed” or “uninformed” scouting phase that assigns budget capriciously. This tests
 866 whether the method is robust to misspecified allocation—i.e., when budget is concentrated in arbitrary
 867 strata rather than those with proxy variance.

868 Compared to results illustrated in Figure 6, random budget allocation performs similarly in terms
 869 of label usage and diameter against other baselines. Specifically, compared to ASI, random budget
 870 allocated outperforms in terms of CI Width. This can be attributed to our pre-stratification mechanism(
 871 (3)), enforcing $\Delta_{\text{Bias}} > 0$. This is almost identical to what we see in the case of random strata
 872 allocation (Section E.6). However, compared to informed stratification, CI-width and coverage have
 873 noticeably different results.

874 We notice wider confidence intervals across all budgets, reflecting the suboptimal allocation violating
 875 the Neyman condition of Theorem 5.6. This is followed by severe undercoverage ($< 80\%$) at low
 876 budgets ($n_b < 0.012$), driven by budget starvation in high-variance strata preventing within-stratum
 877 CLT convergence. Here stabilization occurs at $87\% - 89\%$ coverage once all strata accumulate
 878 sufficient samples ($n_b \geq 0.012$) allowing CLT to activate. So, coverage eventually gets close, but
 879 finite-sample corrections from the arbitrary allocation still permanently degrades it. In fact, the
 880 coverage would asymptotically approach 90% as $n_b \rightarrow \infty$, but the convergence is governed by
 881 the most starved stratum, requiring budgets beyond our experimental range to fully eliminate the
 882 calibration gap.

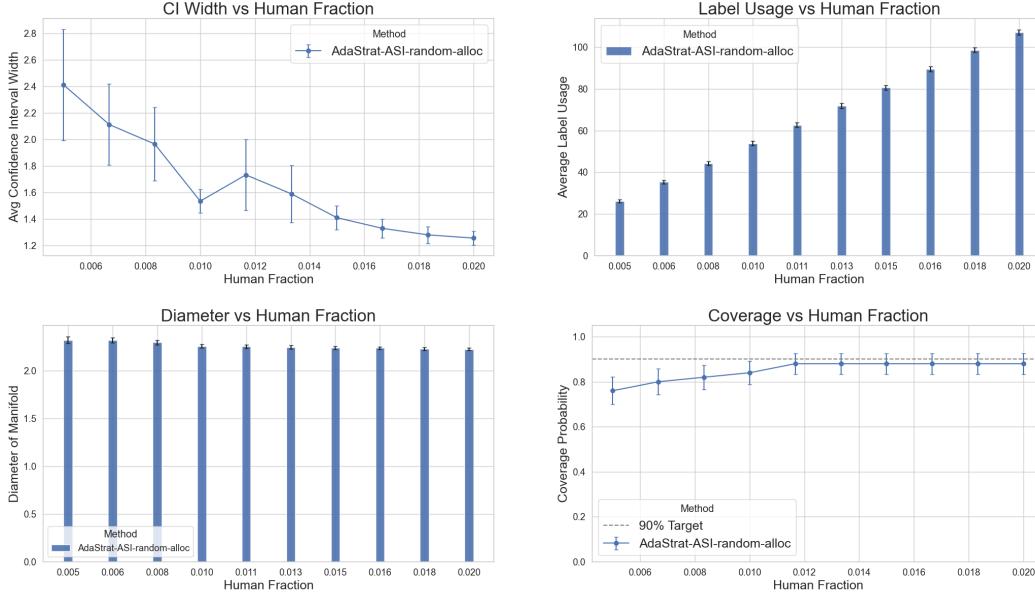


Figure 14: Impact of Random Label Allocation on Politeness Dataset

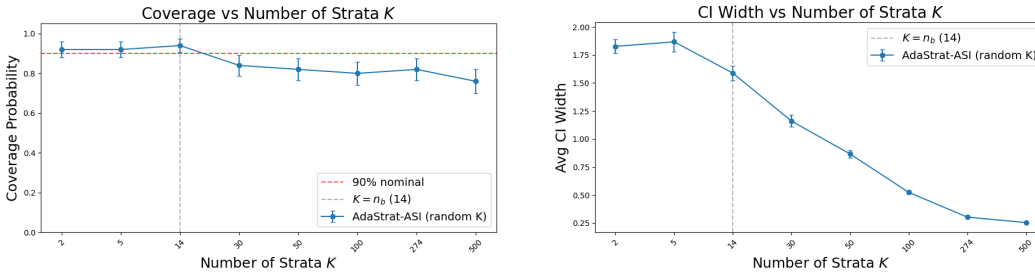


Figure 15: Impact of Rampant Stratification on Politeness Dataset

883 This confirms that stratification *and* allocation are complementary necessities. Even with strata that
 884 perfectly capture model heterogeneity, arbitrary budget allocation destroys the variance reduction
 885 guarantee. The scouting phase is not merely a convenience but a critical component for valid, efficient
 886 inference. So while practitioners *can* choose to skip the Scouting Stage to save on runtime(Section
 887 E.3), the performance becomes severely suboptimal.

888 E.8 Effect of K on AdaStrat-ASI

889 In this section we perform an ablation on K i.e. our strata count. We borrow the random stratification
 890 setup from Appendix E.6 , where we create K strata via random partitioning of the pool, varying
 891 $K \in \{2, 5, 14, 30, 50, 100, 274, 500\}$ on the Politeness dataset at a fixed labelling budget of 0.05%
 892 (14 labelled points in total). For each value of K we report empirical coverage and CI width. As
 893 shown in Figure 15, CI width and Coverage decrease as K increases. Specifically, coverage remains
 894 near the nominal level $1 - \alpha = 0.90$ for small to moderate $K \approx 14$, however, where the *expected*
 895 number of labelled points per stratum, $\frac{nb}{K}$ drops below one, coverage degrades sharply and CI width
 896 becomes unreliable. This behaviour is expected: the within-stratum sandwich variance estimator
 897 relies on a consistent estimate of V_k , which requires sufficiently many labelled observations per
 898 stratum for the central limit theorem to hold locally. Hence

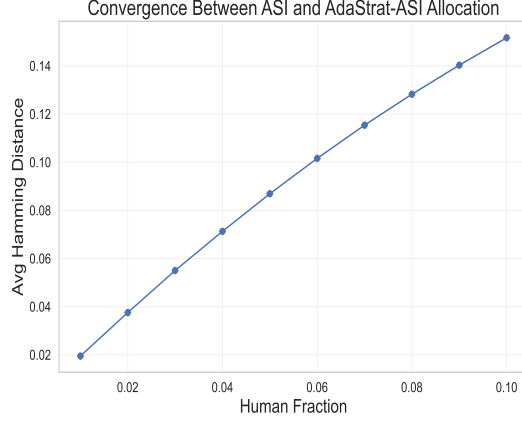


Figure 16: Difference in Label Distributions of ASI and AdaStrat-ASI

899 **E.9 Does stratification(scouting stage) make any difference?**

900 AdaStrat-ASI is not designed to reduce labeling cost or to conserve budget. By construction, it
 901 operates under a fixed labeling budget n_b and, in expectation, expends the same number of labels as
 902 standard ASI. This follows directly from the multinomial budget allocation in the attacking stage and
 903 is consistently observed empirically(Section 6). The objective of AdaStrat-ASI is instead to alter *how*
 904 a fixed labeling budget is allocated across the data space, rather than *how much* labeling is performed.

905 To gauge this impact, for the California Housing Dataset(Section E.2.1) for a fixed total labeling
 906 budget n_b , we define two random binary vectors over the inference pool of size n : $\xi^{\text{ASI}} \in \{0, 1\}^n$,
 907 and $\xi^{\text{AdaStrat}} \in \{0, 1\}^n$ where, $\sum_i \xi_i^{\text{ASI}} = \sum_i \xi_i^{\text{AdaStrat}} = n_b$, but the locations of the ones differ due
 908 to stratified reallocation.

We then compute

$$\bar{H}(\xi^{\text{ASI}}, \xi^{\text{AdaStrat}}) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{\xi_i^{\text{ASI}} \neq \xi_i^{\text{AdaStrat}}\}$$

909 and report its average over trials as shown in Figure 16.

- 910 • Low Budget: Both methods select very few points. Most entries in both vectors are zero.
 911 Hence average Hamming distance is close to zero.
- 912 • High Budget: The Hamming distance increases because stratification actively reshapes the
 913 allocation pattern under the same total budget.

914 Though not shown, at large budgets $n_b \rightarrow n$, both vectors will converge to the all-ones vector. Any
 915 allocation rule becomes identical $\bar{H} \rightarrow 0$. So the curve must eventually come back down.

916 In theory, the true shape is unimodal (low \rightarrow high \rightarrow zero). This analysis highlights that the primary
 917 impact of stratification lies in redistributing labels across strata when labels are scarce, which is
 918 precisely the regime where allocation decisions matter most due to practical constraints.

919 Thus, the variance reduction (Theorem 5.6) is not just theoretical; it stems from a concrete shift in
 920 which specific data points are labeled. The gap represents the *correction* AdaStrat-ASI applies to
 921 preventing over-spending on popular but easy strata.

We formalize this by modeling the unnormalized Hamming distance $H(\beta)$ as a quadratic:

$$H(\beta) = A\beta^2 + B\beta + C$$

922 where $\beta = \frac{n_b}{n}$.

923 We know, $H(0) = 0$. This implies, $C = 0$. We also know, $H(1) = 0$. This implies $A = -B$.

924 Recall, $H(\xi^{\text{ASI}}, \xi^{\text{AdaStrat}}) = H(\xi^{\text{ASI}}, \xi^{\text{AdaStrat}}) = |\xi^{\text{ASI}} \cup \xi^{\text{AdaStrat}}| - |\xi^{\text{ASI}} \cap \xi^{\text{AdaStrat}}| = 2n_b - 2|\xi^{\text{ASI}} \cap$
 925 $\xi^{\text{AdaStrat}}|$

926 For ease of notation, let $I = |\xi^{\text{ASI}} \cap \xi^{\text{AdaStrat}}|$.

927 Let ρ be the correlation between the selection policies.

928 **Case 1: Perfect Correlation ($\rho = 1$).** If the policies are identical, they pick the exact same top- n_b
929 items. Thus, $I = n_b$.

930 **Case 2: Random/Uncorrelated ($\rho = 0$).** If one policy is random relative to the other, the intersection
931 is just the probability of picking the same item twice: $P(\text{pick}) \times P(\text{pick}) = \beta \times \beta$. Thus, $I =$
932 $n \cdot \beta^2 = n_b \cdot \beta$.

933 We model the actual intersection as a weighted average of Perfect Overlap ($\rho = 1$) and Random
934 Overlap, ($\rho = 0$) weighted by the correlation coefficient ρ :

$$I(\beta) = \rho \cdot (n_b) + (1 - \rho) \cdot (n_b \beta)$$

935 On replacing and simplifying we get, $H(\beta) = 2n(1 - \rho)(\beta - \beta^2)$

936 Naturally the normalized Hamming Distance becomes,

$$\bar{H}(\beta) = (-\beta^2) \cdot 2(1 - \rho) + (\beta) \cdot 2(1 - \rho)$$

Magnitude: $\bar{H}(\beta)$'s value is maximized for β at $\beta = 0.5$ i.e.

$$\bar{H}(\beta) = \frac{1 - \rho}{2}$$

Similarly, $\bar{H}(\beta)$'s value is maximized for ρ at $\rho = 0$ i.e.

$$\bar{H}(\beta) = 2\beta(1 - \beta)$$

937 These findings jointly indicate that $0 \leq \bar{H}(\beta) \leq 0.5 \forall \beta, \rho$.

938 **Interpretation:** The amplitude A is directly influenced by $(1 - \rho)$, which represents the ‘decor-
939 relation’ between ASI and AdaStrat-ASI. Consequently, when ρ approaches 0, indicating a high
940 degree of decorrelation, the Hamming Distance over low label budgets exhibits a more pronounced
941 upward curvature, resulting in a higher Hamming Distance. In practical terms, stratifying and then
942 sampling reduces the correlation (ρ) because we concentrate on sampling from k partitions within
943 the n -dimensional data vector, rather than the entire data vector directly.

944 Conversely, when the global (ASI) and local (AdaStrat-ASI) sampling rules converge, decorrelation
945 approaches 0. This flattening of the curve ($\bar{H}(\beta) \rightarrow 0$) signifies that there is no actual change in the
946 labeling distributions, which aligns with intuitive reasoning.