DISCOUNT: Counting in Large Image Collections with Detector-Based Importance Sampling

Gustavo Perez Subhransu Maji* Daniel Sheldon*
Manning College of Information & Computer Sciences
University of Massachusetts, Amherst
{gperezsarabi, smaji, sheldon}@cs.umass.edu

Abstract

Many modern applications use computer vision to detect and count objects in massive image collections. However, when the detection task is very difficult or in the presence of domain shifts, the counts may be inaccurate even with significant investments in training data and model development. We propose DISCOUNTa detector-based importance sampling framework for counting in large image collections that integrates an imperfect detector with human-in-the-loop screening to produce unbiased estimates of counts. We propose techniques for solving counting problems over multiple spatial or temporal regions using a small number of screened samples and estimate confidence intervals. This enables end-users to stop screening when estimates are sufficiently accurate, which is often the goal in a scientific study. On the technical side we develop variance reduction techniques based on control variates and prove the (conditional) unbiasedness of the estimators. DISCOUNT leads to a 9-12× reduction in the labeling costs over naive screening for tasks we consider, such as counting birds in radar imagery or estimating damaged buildings in satellite imagery, and also surpasses alternative covariate-based screening approaches in efficiency.

1 Introduction

Many modern applications use computer vision to detect and count objects in massive image collections. For example, we are interested in applications that involve counting bird roosts in radar images and damaged buildings in satellite images. The image collections are too massive for humans to solve these tasks in the available time. Therefore, a common approach is to train a computer vision detection model and run it exhaustively on the images.

The task is interesting because the goal is not to generalize, but to achieve the scientific counting goal with sufficient accuracy for a *fixed* image collection. The best use of human effort is unclear: it could be used for model development, labeling training data, or even directly solving the counting task! A particular challenge occurs when the detection task is very difficult, so the accuracy of counts made on the entire collection is questionable even with huge investments in training data and model development. Some works resort to human screening of the detector outputs [1–3], which saves time compared to manual counting but is still very labor intensive.

These considerations motivate *statistical* approaches to counting. Instead of screening the detector outputs for all images, a human can "spot-check" some images to estimate accuracy, and, more importantly, use statistical techniques to obtain unbiased estimates of counts across unscreened images. In a related context, Meng et al. [4] proposed IS-count, which uses importance sampling

^{*}equal advising contribution

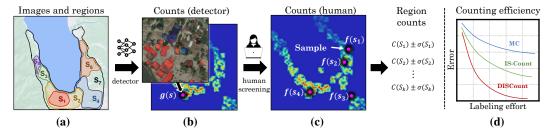


Figure 1: k-DISCOUNT uses detector-based importance sampling to screen counts and solve multiple counting problems. (a) Geographical regions where we want to estimate counts of damaged buildings. (b) Outputs of a damaged building detector on satellite imagery, which can be used to estimate counts g(s) for each tile (shows as dots). (c) Tiles selected for human screening to obtain true counts f(s), from which counts for all regions are joinly estimated by k-DISCOUNT. (d) Our experiments show that DISCOUNT outperforms naive (MC) and covariate-based sampling (IS-Count [4]).

to estimate total counts across a collection when (satellite) images are expensive to obtain by using spatial covariates to sample a subset of images.

We contribute counting methods for large image collections that build on IS-count in several ways. First, we work in a different model where images are freely available and it is possible to train a detector to run on all images, but the detector is not reliable enough for the final counting task, or its reliability is unknown. We contribute human-in-the-loop methods for count estimation using the detector to construct a proposal distribution, as seen in Fig. 1. Second, we consider solving multiple counting problems—for example, over disjoint or overlapping spatial or temporal regions—simultaneously, which is very common in practice. We contribute a novel sampling approach to obtain simultaneous estimates, prove their (conditional) unbiasedness, and show that the approach allocates samples to regions in a way that approximates the optimal allocation for minimizing variance. Third, we design confidence intervals, which are important practically to know how much human effort is needed. Fourth, we use variance reduction techniques based on control variates.

Our method produces unbiased estimates and confidence intervals with reduced error compared to covariate-based methods. In addition, the labeling effort is further reduced with DISCOUNT as we only have to verify detector predictions instead of producing annotations from scratch. On our tasks, DISCOUNT leads to a $9-12\times$ reduction in the labeling costs over naive screening and $6-8\times$ reduction over IS-Count. Finally, we show that solving multiple counting problems jointly can be done more efficiently than solving them separately, demonstrating a more efficient use of samples.

2 Related Work

Computer vision techniques have been deployed for counting in numerous applications where exhaustive human-labeling is expensive due to the sheer volume of imagery involved. This includes areas such as detecting animals in camera trap imagery [1, 5], counting buildings, cars, and other structures in satellite images [2, 6–8], species monitoring in citizen science platforms [5, 9], monitoring traffic in videos [10, 11], as well as various medicine, science and engineering applications. For many applications the cost associated with training an *accurate* model is considerably less than that of meticulously labeling the entire dataset. Even with a less accurate model, human-in-the-loop recognition strategies have been proposed to reduce annotation costs by integrating human validation with noisy predictions [12, 13].

Our approach is related to work in active learning [14] and semi-supervised learning [15], where the goal is to reduce human labeling effort to learn models that generalize on i.i.d. held out data. While these approaches reduce the cost of labels on training data, they often rely on large labeled test sets to estimate the performance of the model, which can be impractical. Active testing [16, 17] aims to reduce the cost of model evaluation by providing a statistical estimate of the performance using a small number of labeled examples. Unlike traditional learning where the goal is performance on held out data, the goal of active testing is to estimate performance on a *fixed* dataset. Similarly, our goal is to estimate the counts on a fixed dataset, but different from active testing we are interested in estimates of the true counts and not the model's performance. In particular, we want unbiased

estimates of counts even when the detector is unreliable. Importantly, since generalization is not the goal, overfitting to the dataset statistics may lead to more accurate estimates.

Statistical estimation has been widely used to conduct surveys (e.g., estimating population demographics, polling, etc.) [18]. In IS-Count [4], the authors propose an importance sampling approach to estimate counts in large image collections using humans-in-the-loop. They showed that one can count the number of buildings at the continental scale by sampling a small number of regions based on covariates such as population density and annotating those regions, thereby reducing the cost of obtaining high-resolution satellite imagery and human labels. However, for many applications the dataset is readily available, and running the detector is cost effective, but human screening is expensive. To address this, we propose using the detector to guide the screening process and demonstrate that this significantly reduces error rates in count estimation given a fixed amount of human effort. Furthermore, for some applications, screening the outputs of a detector can be significantly faster than to annotate from scratch, leading to additional savings.

An interesting question is what is the best way to utilize human screening effort to count on a dataset. For example, labels might be used to improve the detector, measure performance on the deployed dataset, or, as is the case in our work, to derive a statistical estimate of the counts. Our work is motivated by problems where improving the detector might require significant effort, but counts from the detector are correlated with true counts and can be used as a proposal distribution for sampling.

3 DISCOUNT: Detector-based IS-Count

Consider a counting problem in a discrete domain Ω (usually spatiotemporal) with elements $s \in \Omega$ that represent a single unit such as an image, grid cell, or day of year. For each s there is a ground truth "count" $f(s) \geq 0$, which can be any non-negative measurement, such as the number or total size of all objects in an image. A human can label the underlying images for any s to obtain s0.

Define $F(S) = \sum_{s \in S} f(s)$ to be the cumulative count for a region S. We wish to estimate the total counts $F(S_1), \ldots, F(S_k)$ for k different subsets $S_1, \ldots, S_k \subseteq \Omega$, or regions, while using human effort as efficiently as possible. The regions represent different geographic divisions or time ranges and may overlap — for example, in the roost detection problem we want to estimate cumulative counts of birds for each day of the year, while disaster-relief planners want to estimate building damage across different geographical units such as towns, counties, and states. Assume without loss of generality that $\bigcup_{i=1}^k S_i = \Omega$, otherwise the domain can be restricted so this is true.

We will next present our methods; derivations and proofs of all results are found in the appendix.

3.1 Single-Region Estimators

Consider first the problem of estimating the total count F(S) for a single region S. Meng et al. [4] studied this problem in the context of satellite imagery, with the goal of minimizing the cost of purchasing satellite images to obtain an accurate estimate.

Simple Monte Carlo [4] This is a baseline based on simple Monte Carlo sampling. Write $F(S) = \sum_{s \in S} f(s) = |S| \cdot \mathbb{E}_{s \sim \mathrm{Unif}(S)}[f(s)]$. Then the following estimator, which draws n random samples uniformly in S to estimate the total, is unbiased:

$$\hat{F}_{MC}(S) = |S| \cdot \frac{1}{n} \sum_{i=1}^{n} f(s_i), \quad s_i \sim \text{Unif}(S).$$

IS-Count [4] Meng et al. then proposed an estimator based on importance sampling [19]. Instead of sampling uniformly, the method samples from a proposal distribution q that is cheap to compute for all $s \in S$. For example, to count buildings in US satellite imagery, the proposal distribution could use maps of artificial light intensity, which are freely available. The importance sampling estimator is:

$$\hat{F}_{IS}(S) = \frac{1}{n} \sum_{i=1}^{n} \frac{f(s_i)}{q(s_i)}, \quad s_i \sim q.$$

DISCOUNT IS-count assumes *images* are costly to obtain, which motivates using external covariates for the proposal distribution. However, in many scientific tasks, the images are readily available, and the key cost is that of human supervision. In this case it is possible to train a detection model and run it on all images to produce an approximate count g(s) for each s. Define $G(S) = \sum_{s \in S} g(s)$ to be the approximate detector-based count for region s. We propose the *detector-based IS-count* ("DISCOUNT") estimator, which uses the proposal distribution proportional to s0 on region s1, i.e., with density s1 and s2 is s3 are constant. The importance-sampling estimator then specializes to:

$$\hat{F}_{\text{DIS}}(S) = G(S) \cdot \frac{1}{n} \sum_{i=1}^{n} \frac{f(s_i)}{g(s_i)}, \quad s_i \sim \bar{g}_S.$$

To interpret DISCOUNT, let $w_i = f(s_i)/g(s_i)$ be the ratio of the true count to the detector-based count for the ith sample s_i or (importance) weight. DISCOUNT reweights the detector-based total count G(S) by the average weight $\bar{w} = \frac{1}{n} \sum_{i=1}^n w_i$, which can be viewed as a correction factor based on the tendency to over- or under-count, on average, across all of S.

DISCOUNT is unbiased as long as $\bar{g}(s) > 0$ for all $s \in S$ such that f(s) > 0. Henceforth, we assume detector counts are pre-processed if needed so that g(s) > 0 for all relevant units, for example, by adding a small amount to each count.

3.2 k-DISCOUNT

We now return to the multiple region counting problem. A naive approach would be to run DISCOUNT separately for each region. However, this is suboptimal. First, it allocates samples equally to each region, regardless of their size or predicted count. Intuitively, we want to allocate more effort to regions with higher predicted counts. Second, if regions overlap it is wasteful to repeatedly draw samples from each one to solve the estimation problems separately.

k-DISCOUNT We propose estimators based on n samples drawn from all of Ω with probability proportional to g. Then, we can estimate F(S) for any region using only the samples from S. Specifically, the k-DISCOUNT estimator is

$$\hat{F}_{k\mathrm{DIS}}(S) = \begin{cases} G(S) \cdot \bar{w}(S) & n(S) > 0 \\ 0 & n(S) = 0 \end{cases}, \qquad s_i \sim \bar{g}_{\Omega},$$

where $n(S) = |\{i : s_i \in S\}|$ is the number of samples in region S and $\bar{w}(S) = \frac{1}{n(S)} \sum_{i:s_i \in S} w_i$ is the average importance weight for region S.

Claim 1. The k-DISCOUNT estimator $\hat{F}_{kDIS}(S)$ is conditionally unbiased given at least one sample in region S. That is, $\mathbb{E}[\hat{F}_{kDIS}(S) \mid n(S) > 0] = F(S)$.

The unconditional bias can also be analyzed (see Appendix). Overall, bias has negligible practical impact. It occurs only when the sample size n(S) is zero, which is an event that is both observable and has probability $(1 - p(S))^n$ that decays exponentially in n, where $p(S) = G(S)/G(\Omega)$.

In terms of variance, k-DISCOUNT behaves similarly to DISCOUNT run on each region S with sample size equal to $\mathbb{E}[n(S)] = np(S)$. To first order, both approaches have variance $\frac{G(S)^2 \cdot \sigma^2(S)}{np(S)}$ where $\sigma^2(S)$ is the importance-weight variance. In the case of *disjoint* regions, running DISCOUNT on each region is the same as *stratified importance sampling* across the regions, and the allocation of np(S) samples to region S is optimal in the following sense:

Claim 2. Suppose S_1, \ldots, S_k partition Ω and the importance weight variance $\sigma^2(S_i) = \sigma^2$ is constant across regions. Assume DISCOUNT is run on each region S_i with n_i samples. Given a total budget of n samples, the sample sizes that minimize $\sum_{i=1}^k \mathrm{Var}(\hat{F}_{DIS}(S_i))$ are given by $n_i = np(S_i) = nG(S_i)/G(\Omega)$.

The analysis uses reasoning similar to the *Neyman allocation* for stratified sampling [18], and shows that k-DISCOUNT approximates the optimal allocation of samples to (disjoint) regions under the stated assumptions. One key difference is that k-DISCOUNT draws samples from all of Ω and then assigns them to regions, which is called "post-stratification" in the sampling literature [18]. An exact

variance analysis in the Appendix reveals that, if the expected sample size np(S) for a region is very small, k-DISCOUNT may have up to 30% "excess" variance compared to stratification due to the random sample size, but the excess variance disappears quickly and both approaches have the same asymptotic variance. A second key difference to stratification is that regions can overlap; k-DISCOUNT's approach of sampling from all of Ω and then assigning samples to regions extends cleanly to this setting.

3.3 Control Variates

Control variates are functions h(s) whose integrals $H(S) = \sum_{s \in S} h(s)$ are known and can be combined with importance sampling using the following estimator:

$$\hat{F}_{k\mathrm{DIS}cv}(S) = \begin{cases} G(S) \cdot \bar{w}_h(S) + H(S) & n(S) > 0 \\ 0 & n(S) = 0 \end{cases}, \quad s_i \sim \bar{g}_{\Omega},$$

where $\bar{w}_h(S) = \frac{1}{n(S)} \sum_{i:s_i \in S} w_{h,i}$ and $w_{h,i} = (f(s_i) - h(s_i))/g(s_i)$. It is clear that $\hat{F}_{kDIScv}(S)$ has the same expectation as $\hat{F}_{kDIS(S)}$, but $\hat{F}_{kDIScv}(S)$ might have a lower variance under certain conditions (if f and h are sufficiently correlated [19]). For bird counting, estimated counts from previous years could be used as control variates as migration is periodic to improve count estimates (see experiments in § 4 for details).

3.4 Confidence intervals

Confidence intervals for k-DISCOUNT can be constructed in a way similar to standard importance sampling. For a region S, first estimate the importance weight variance $\sigma^2(S)$ as:

$$\hat{\sigma}^2(S) = \frac{1}{n(S)} \sum_{i: s_i \in S} \left(\frac{f(s_i)}{g(s_i)} - \frac{\hat{F}_{k \text{DIS}}(S)}{G(S)} \right)^2.$$

An approximate $1-\alpha$ confidence interval is then given by $\hat{F}_{k\text{DIS}}(S) \pm z_{\alpha/2} \cdot G(S) \cdot \hat{\sigma}(S) / \sqrt{n(S)}$, where z_{γ} is the $1-\gamma$ quantile of the standard normal distribution, e.g., $z_{0.025}=1.96$ for a 95% confidence interval. The theoretical justification is subtle due to scaling by the *random* sample size n(S). It is based on the following asymptotic result, proved in the Appendix.

Claim 3. The k-DISCOUNT estimator with scaling factor $G(S)\hat{\sigma}(S)/\sqrt{n(S)}$ is asymptotically normal, that is, the distribution of $\frac{\hat{F}_{kDIS}(S)-F(S)}{G(S)\cdot\hat{\sigma}(S)/\sqrt{n(S)}}$ converges to $\mathcal{N}(0,1)$ as $n\to\infty$.

In preliminary experiments we observed that for small expected sample sizes the importance weight variance $\sigma^2(S)$ can be underestimated leading to intervals that are too small — as an alternative, we propose a practical heuristic for smaller sample sizes where $\hat{\sigma}^2(\Omega)$ is used instead of $\hat{\sigma}^2(S)$; that is, all samples are used to estimate variability of importance weights for each region S.

4 Experimental Setup

In this section we describe the counting tasks and detection models (§ 4.1–4.2) and the evaluation metrics (§ 4.3) we will use to evaluate different counting methods. We focus on two applications: counting roosting birds in weather radar images and counting damaged buildings in satellite images of a region struck by a natural disaster.

4.1 Counting Roosting Birds from Weather Radar

Many species of birds and bats congregate in large numbers at nighttime or daytime roosting locations. Their departures from these "roosts" are often visible in weather radar, from which it's possible to estimate their numbers [20–22]. The US "NEXRAD" weather radar network [23] has collected data for 30 years from 143+ stations and provides an unprecedented opportunity to study long-term and wide-scale biological phenomenon such as roosts [24, 25]. However, the sheer volume of radar scans (>250M) prevents manual analysis and motivates computer vision approaches [26–28, 3].

Unfortunately, the best computer vision models [3, 28] for detecting roosts have average precision only around 50% and are not accurate enough for fully automated scientific analysis, despite using state-of-the-art methods such as Faster R-CNNs [29] and training on thousands of human annotations—the complexity of the task suggests substantial labeling and model development efforts would be needed to improve accuracy, and may be impractical.

Previous work [30, 31] used a roost detector combined with manual screening of the detections to analyze more than 600,000 radar scans spanning a dozen stations in the Great Lakes region of the US to reveal patterns of bird migration over two decades. The vetting of nearly 64,000 detections was orders of magnitude faster than manual labeling, yet still required a substantial 184 hours of manual effort. Scaling to the entire US network would require at least an order of magnitude more effort, thus motivating a statistical approach.

We use the exhaustively screened detections from the Great Lakes analysis in [30, 31] to systematically analyze the efficiency of sampling based counting. The data is organized into domains $\Omega^{\text{sta},yr}$ corresponding to 12 stations and 20 years (see Fig. 7 in Appendix B). Thus the domains are disjoint and treated separately. Counts are collected for each day s by running the detector using all radar scans for that day to detect and track roost signatures and then mapping detections to bird counts using the measured radar "reflectivity" within the tracks. For the approx-

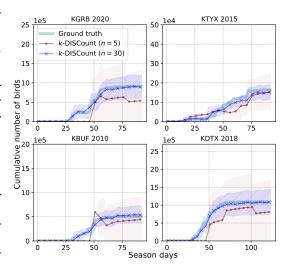


Figure 2: Count estimates with confidence intervals for two station years (i.e., KGRB 2020 and KBUF 2010) using different numbers of samples.

imate count g(s) we use the automatically detected tracks, while for the true count f(s) we use the manually screened and corrected tracks. For a single domain, i.e., each station-year, we divide a complete roosting season into temporal regions in three different scenarios: (1) estimating bird counts up to each day in the roosting season (i.e., regions are nested prefixes of days in the entire season), (2) the end of each quarter of (i.e., regions are nested prefixes of quarters in the entire season), and (3) estimating each quarter's count (each region is one quarter). We measure error using the fully-screened data and average errors across all domains and regions. Fig. 2 shows the counts and confidence intervals estimated using k-DISCOUNT for the first scenario on four station-years.

4.2 Counting Damaged Buildings from Satellite Images

Building damage assessment from satellite images [32, 33] is often used to plan humanitarian response after a natural disaster strikes. However, the performance of computer vision models degrades when applied to new regions and disaster types. Our approach can be used to quickly vet the data produced by the detector to correctly estimate counts in these scenarios.

We use the building damage detection model by [34], the winner of the xView2 challenge [35]. The model is based on U-Net [36] to detect buildings in the pre-disaster image, followed by a "siamese network" that incorporates at pre- and post-disaster images to estimate damage. The model is trained on the xBD dataset [37] that contains building and damage annotations spanning multiple geographical regions and disaster types (e.g., earthquake, hurricane, tsunami, etc.). While the dataset contains four levels of damage (i.e., 0: no-damage, 1: minor-damage, 2: major-damage, and 3: destroyed), in this work we combine all damage levels (i.e., classes 1-3) into a single "damage" class.

We consider the Palu Tsunami from 2018; the data consists of 113 high-resolution satellite images labeled with 31,394 buildings and their damage levels. We run the model on each tile s to estimate the number of damaged buildings g(s), while the ground-truth number of damaged buildings is used as f(s). Our goal is to estimate the cumulative damaged building count in sub-regions expanding from the area with the most damaged buildings as shown in Fig. 9 in the Appendix C. To define the sub-regions, we sort all m images by their distance from the epicenter (defined as the image tile with

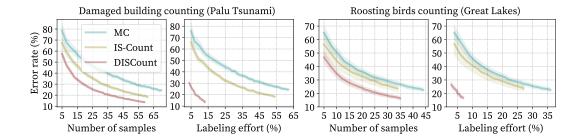


Figure 3: **Detector-based sampling.** Estimation error of damaged building counts in the Palu Tsunami region from the xBD dataset (left) and counting roosting birds from the Great Lakes radar stations in the US from NEXRAD data (right). We get lower error with DISCOUNT compared to IS-Count and simple Monte Carlo sampling (MC). The labeling effort is further reduced with DISCOUNT since the user is not required to label an image from scratch but only to verify outputs from the detector (See § 5 for details). The estimation errors are averaged over 1000 runs.

the most damaged buildings) and then divide into chunks or "annuli" A_1, \ldots, A_7 of size m/7. The task is to estimate the cumulative counts $S_j = \bigcup_{i=1}^j A_i$ of the first j chunks for j from 1 to 7.

4.3 Evaluation

We measure the fractional error between the *true* and the estimated counts averaged over all regions in a domain $S_1, \ldots, S_k \subseteq \Omega$ as:

$$\operatorname{Error}(\Omega) = \frac{1}{k} \sum_{i=1}^{k} \frac{|F(S_i) - \hat{F}(S_i)|}{F(\Omega)}.$$

For the bird counting task, for any given definition of regions within one station-year Ω (i.e., cumulative days or quarters defined in § 4.1) we report the error averaged across all station-years corresponding to 12 stations and ≈ 20 years. For the damaged building counting problem there is only a single domain corresponding to the Palu Tsunami region. In addition, we calculate the average confidence interval width normalized by $F(\Omega)$. We run 1000 trials and plot average metrics $\pm 1.96 \times$ std. error over the trials. We also evaluate confidence interval coverage, which is the fraction of confidence intervals that contain the true count over all domains, regions, and trials.

5 Results

In this section, we present the results comparing detector-based to covariate-based sampling. Also, we show reductions in labeling effort and demonstrate the advantages of estimating multiple counts jointly. Finally, we show confidence intervals and control variates results.

Detector-based sampling reduces error We first compare DISCOUNT (detector-based sampling) to IS-Count and simple Monte Carlo sampling for estimating $F(\Omega)$, that is, the total counts of birds in a complete roosting season for a given station year, or damaged buildings in the entire disaster region. Fig. 3 shows the error rate as a function of number of labeled samples (i.e., the number of distinct s_i sampled, since each s is labeled at most once). In the buildings application, a sample refers to an image tile of size 1024×1024 pixels, while for the birds a sample refers to a single day.

Using the detector directly without any screening results in high error rates — roughly 136% and 149% for estimating the total count for the damaged buildings and bird counting tasks respectively. Meng et al. [4] show the advantages of using importance sampling with screening to produce count estimates with base covariates as opposed to simple Monte Carlo sampling (MC vs. IS-Count). For the bird counting task, we construct a non-detector covariate $g_{\rm IS}$ by fitting a spline to f(s) with 10% of the days from an arbitrarily selected station-year pair (station KBUF in 2001). For the damaged building counting task, the covariate $g_{\rm IS}$ is the true count of all buildings (independent of the damage) obtained using the labels provided with the xBD dataset.

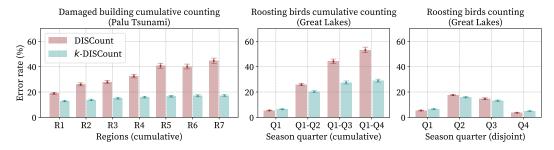


Figure 4: **Solving multiple counting problems jointly.** Estimation error of counting damaged buildings in the Palu Tsunami region from the xBD dataset (left) and counting roosting birds from the Great Lakes radar stations in the US from NEXRAD data (right). We compare solving the counting problems jointly (*k*-DISCOUNT) against solving the counting problems separately (DISCOUNT). We use 10 samples for both these tests. The estimation errors are averaged over 1000 runs.

Covariate-based sampling (IS-Count) leads to significant savings over simple Monte Carlo sampling (MC), but DISCOUNT provides further improvements. In particular, to obtain an error rate of 20% DISCOUNT requires $\approx 1.6\times$ fewer samples than IS-Count and $\approx 3\times$ fewer samples than MC for both counting problems.

Screening leads to a further reduction in labeling effort DISCOUNT alleviates the need for users to annotate an image from scratch, such as identifying an object and drawing a bounding box around it. Instead, users only need to verify the detector's output, which tends to be a quicker process. In a study by Su et al. [38] on the ImageNet dataset [39], the median time to draw a bounding-box was found to be 25.5 seconds, whereas verification took only 9.0 seconds (this matches the screening time of \approx 10s per bounding-box in [31, 30]). The right side of Fig. 3 presents earlier plots with the x-axis scaled based on labeling effort, computed as $100 \cdot c \cdot n/|\Omega|$, where n denotes the number of screened samples and $c \in [0,1]$ represents the fraction of time relative to labeling from scratch. For instance, the labeling effort is 100% when all elements must be labeled from scratch (c=1 and $n=|\Omega|$). For DISCOUNT, we estimate $c_{\text{DIS}}=9.0/(25.5+9.0)=0.26$, since annotating from scratch requires both drawing and verification, while screening requires only verification. To achieve the same 20% error rate, DISCOUNT requires $6\times$ less effort than IS-Count and $9\times$ less effort than MC for the bird counting task, and $8\times$ less effort than IS-Count and $12\times$ less effort than MC for building counting.

Multiple counts can be estimated efficiently (k-DISCOUNT) To solve multiple counting problems, we compared k-DISCOUNT to using DISCOUNT separately on each region. For bird counting, the task was to estimate four quarterly counts (cumulative or individual) as described in § 4.1. For k-DISCOUNT, we sampled n=40 days from the complete season to estimate the counts simultaneously. For DISCOUNT, we solved each of the four problems separately using n/4=10 samples per region for the same total number of samples. For building damage counting, the task was to estimate seven cumulative counts as described in § 4.2. For k-DISCOUNT, we used n=70 images sampled from the entire domain, while for DISCOUNT we used n/7=10 sampled images per region.

Fig. 4 shows that solving multiple counting problems jointly (*k*-DISCOUNT) is better than solving them separately (DISCOUNT). For the cumulative tasks, *k*-DISCOUNT makes much more effective use of samples from overlapping regions. For single-quarter bird counts, *k*-DISCOUNT has slightly higher error in Q1 and Q4 and lower errors in Q2 and Q3. This can be understood in terms of sample allocation: *k*-DISCOUNT allocates in proportion to predicted counts, which provides more samples and better accuracy in Q2-Q3, when many more roosts appear, and approximates the optimal allocation of Claim 2. DISCOUNT allocates samples equally, so has slightly lower error for the smaller Q1 and Q4 counts. In contrast, for building counting, *k*-DISCOUNT has lower error even for the smallest region R1, since this has the most damaged buildings and thus gets more samples than DISCOUNT. Fig. 5 (left) shows *k*-DISCOUNT outperforms simple Monte Carlo (adapted to multiple regions similarly to *k*-DISCOUNT) for estimating cumulative daily bird counts as in Fig. 2.

Confidence intervals We measure the width and coverage of the estimated confidence intervals (CIs) per number of samples for cumulative daily bird counting; see examples in Fig. 2. We compare the CIs of k-DISCOUNT, k-DISCOUNT-cv (control variates), k-DISCOUNT-cv- $\sigma(\Omega)$ (using all sam-

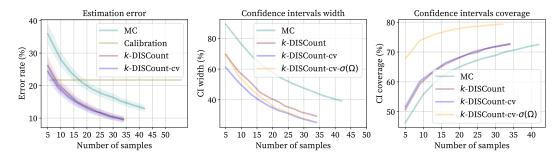


Figure 5: Control variates and confidence intervals on bird counting. We compare simple Monte Carlo (MC), calibration with isotonic regression, and variations of k-DISCOUNT that include control variates (-cv) and improved variance estimates $(-\sigma(\Omega))$. (left) Error rates using k-DISCOUNT are significantly smaller than MC and calibration. (middle) Confidence intervals' width. (right) Confidence intervals' coverage. The error and the confidence intervals' width are slightly reduced when control variates are used while maintaining the coverage. Furthermore, k-DISCOUNT-cv- $\sigma(\Omega)$ improves the coverage. The results are averaged over all station-years and over 1,000 runs.

ples to estimate variance), and simple Monte Carlo sampling in Fig. 5. When using control variates, the error rate and the CI width are slightly reduced while keeping the same coverage. CI coverage is lower than the nominal coverage (95%) for all methods, but increasing with sample size and substantially improved by k-DISCOUNT-cv- $\sigma(\Omega)$, which achieves up to $\approx 80\%$ coverage. Importance weight distributions can be heavily right-skewed and the variance easily underestimated [40].

DISCOUNT improves over a calibration baseline We implement a calibration baseline where the counts are estimated as $\hat{F}_{CAL}(S) = \sum_{s \in S} \hat{\phi}(g(s))$, where we learn an isotonic regression model $\hat{\phi}$ between the predicted and true counts trained for each station using 15 uniformly selected samples from one year from that station. Results are shown as the straight line in Fig. 5 (left). DISCOUNT outperforms calibration with less than 10 samples per station suggesting the difficulties in generalization across years using a simple calibration approach.

Control variates (k-**DISCOUNT-cv**) We perform experiments adding control variates to k-DISCOUNT in the roosting birds counting problem. We use the calibrated detector counts $\hat{\phi}(g(s))$ defined above as the control variate for each station year. Fig. 5 shows that control variates reduce the confidence interval width (middle: k-DISCOUNT vs. k-DISCOUNT-cv) without hurting coverage (right). In addition, the error of the estimate is reduced slightly, as shown in Fig. 5 (left). Note that this is achieved with a marginal increase in the labeling effort.

6 Discussion and Conclusion

We contribute methods for counting in large image collections with a detection model. When the task is complex and the detector is imperfect, allocating human effort to estimate the scientific result directly might be more efficient than improving the detector. For instance, performance gains from adding more training data may be marginal for a mature model. Our proposed solution produces accurate and unbiased estimates with a significant reduction in labeling costs from naive and covariate-based screening approaches. We demonstrate this in two real-world open problems where data screening is still necessary despite large investments in model development. Our approach is limited by the availability of a good detector, and confidence interval coverage is slightly low; possible improvements are to use bootstrapping or corrections based on importance-sampling diagnostics [40].

7 Acknowledgements

We thank Wenlong Zhao for the deployment of the roost detector, Maria Belotti, Yuting Deng, and our Colorado State University AeroEco Lab collaborators for providing the screened data of the Great Lakes radar stations, and Yunfei Luo for facilitating the building detections on the Palu Tsunami region. This work was supported by the National Science Foundation award #2017756.

References

- [1] Mohammad Sadegh Norouzzadeh, Anh Nguyen, Margaret Kosmala, Alexandra Swanson, Meredith S Palmer, Craig Packer, and Jeff Clune. Automatically identifying, counting, and describing wild animals in camera-trap images with deep learning. *Proceedings of the National Academy of Sciences*, 115(25):E5716–E5725, 2018.
- [2] W. Nurkarim and A.W. Wijayanto. Building footprint extraction and counting on very high-resolution satellite imagery using object detection deep learning framework. *Earth Sci Inform*, 2023.
- [3] Gustavo Perez, Wenlong Zhao, Zezhou Cheng, Maria Carolina T. D. Belotti, Yuting Deng, Victoria F. Simons, Elske Tielens, Jeffrey F. Kelly, Kyle G. Horton, Subhransu Maji, and Daniel Sheldon. Using spatio-temporal information in weather radar data to detect and track communal bird roosts. *bioRxiv*, 2022. doi: 10.1101/2022.10.28.513761. URL https://www.biorxiv.org/content/early/2022/10/31/2022.10.28.513761.
- [4] Chenlin Meng, Enci Liu, Willie Neiswanger, Jiaming Song, Marshall Burke, David Lobell, and Stefano Ermon. Is-count: Large-scale object counting from satellite images with covariate-based importance sampling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.
- [5] Devis Tuia, Benjamin Kellenberger, Sara Beery, Blair R Costelloe, Silvia Zuffi, Benjamin Risse, Alexander Mathis, Mackenzie W Mathis, Frank van Langevelde, Tilo Burghardt, et al. Perspectives in machine learning for wildlife conservation. *Nature communications*, 13(1):792, 2022.
- [6] Jeannine Cavender-Bares, Fabian D Schneider, Maria João Santos, Amanda Armstrong, Ana Carnaval, Kyla M Dahlin, Lola Fatoyinbo, George C Hurtt, David Schimel, Philip A Townsend, et al. Integrating remote sensing with ecology and evolution to advance biodiversity conservation. *Nature Ecology & Evolution*, 6(5):506–519, 2022.
- [7] Marshall Burke, Anne Driscoll, David B Lobell, and Stefano Ermon. Using satellite imagery to understand and promote sustainable development. *Science*, 371(6535):eabe8628, 2021.
- [8] Jens Leitloff, Stefan Hinz, and Uwe Stilla. Vehicle detection in very high resolution satellite images of city areas. *IEEE transactions on Geoscience and remote sensing*, 48(7):2795–2806, 2010.
- [9] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The Inaturalist species classification and detection dataset. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [10] Myounggyu Won. Intelligent traffic monitoring systems for vehicle classification: A survey. *IEEE Access*, 8:73340–73358, 2020.
- [11] Benjamin Coifman, David Beymer, Philip McLauchlan, and Jitendra Malik. A real-time computer vision system for vehicle tracking and traffic surveillance. *Transportation Research Part C: Emerging Technologies*, 6(4):271–288, 1998.
- [12] Steve Branson, Catherine Wah, Florian Schroff, Boris Babenko, Peter Welinder, Pietro Perona, and Serge Belongie. Visual recognition with humans in the loop. In *European Conference on Computer Vision (ECCV)*, 2010.
- [13] Catherine Wah, Grant Van Horn, Steve Branson, Subhransu Maji, Pietro Perona, and Serge Belongie. Similarity comparisons for interactive fine-grained categorization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [14] Burr Settles. Active learning literature survey. 2009.
- [15] Olivier Chapelle, Bernhard Scholkopf, and Alexander Zien. Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews]. IEEE Transactions on Neural Networks, 20(3):542–542, 2009.

- [16] Phuc Nguyen, Deva Ramanan, and Charless Fowlkes. Active testing: An efficient and robust framework for estimating accuracy. In *International Conference on Machine Learning (ICML)*, 2018.
- [17] Jannik Kossen, Sebastian Farquhar, Yarin Gal, and Tom Rainforth. Active testing: Sample-efficient model evaluation. In *International Conference on Machine Learning*, 2021.
- [18] William G Cochran. Sampling techniques. John Wiley & Sons, 1977.
- [19] Art B. Owen. Monte Carlo theory, methods and examples. 2013.
- [20] David W. Winkler. Roosts and migrations of swallows. Hornero, 21(2):85–97, 2006.
- [21] Jeffrey J Buler, Lori A Randall, Joseph P Fleskes, Wylie C Barrow Jr, Tianna Bogart, and Daria Kluver. Mapping wintering waterfowl distributions using weather surveillance radar. *PloS one*, 7(7):e41571, 2012.
- [22] Jason W Horn and Thomas H Kunz. Analyzing nexrad doppler radar images to assess nightly dispersal patterns and population trends in brazilian free-tailed bats (tadarida brasiliensis). *Integrative and Comparative Biology*, 48(1):24–39, 2008.
- [23] Next generation weather radar (NEXRAD). National Centers for Environmental Information (NCEI), Feb 2023. URL https://www.ncei.noaa.gov/products/radar/next-generation-weather-radar.
- [24] Kenneth V. Rosenberg, Adriaan M. Dokter, Peter J. Blancher, John R. Sauer, Adam C. Smith, Paul A. Smith, Jessica C. Stanton, Arvind Panjabi, Laura Helft, Michael Parr, and Peter P. Marra. Decline of the north american avifauna. *Science*, 366(6461):120–124, 2019. doi: 10.1126/science.aaw1313.
- [25] Francisco Sánchez-Bayo and Kris AG Wyckhuys. Worldwide decline of the entomofauna: A review of its drivers. *Biological conservation*, 232:8–27, 2019.
- [26] Carmen Chilson, Katherine Avery, Amy McGovern, Eli Bridge, Daniel Sheldon, and Jeffrey Kelly. Automated detection of bird roosts using NEXRAD radar data and convolutional neural networks. *Remote Sensing in Ecology and Conservation*, 2018.
- [27] Tsung-Yu Lin, Kevin Winner, Garrett Bernstein, Abhay Mittal, Adriaan M Dokter, Kyle G Horton, Cecilia Nilsson, Benjamin M Van Doren, Andrew Farnsworth, Frank A La Sorte, et al. Mistnet: Measuring historical bird migration in the us using archived weather radar data and convolutional neural networks. *Methods in Ecology and Evolution*, 10(11):1908–1922, 2019.
- [28] Z. Cheng, S. Gabriel, P. Bhambhani, D. Sheldon, S. Maji, A. V, and D. Winkler. Detecting and tracking communal bird roosts in weather radar data. *Association for the Advancement of Artificial Intelligence (AAAI)*, 2020.
- [29] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. Advances in neural information processing systems, 2015.
- [30] Maria Carolina TD Belotti, Yuting Deng, Wenlong Zhao, Victoria F Simons, Zezhou Cheng, Gustavo Perez, Elske Tielens, Subhransu Maji, Daniel Sheldon, Jeffrey F Kelly, et al. Long-term analysis of persistence and size of swallow and martin roosts in the us great lakes. *Remote Sensing in Ecology and Conservation*, 2023.
- [31] Yuting Deng, Maria Carolina TD Belotti, Wenlong Zhao, Zezhou Cheng, Gustavo Perez, Elske Tielens, Victoria F Simons, Daniel R Sheldon, Subhransu Maji, Jeffrey F Kelly, et al. Quantifying long-term phenological patterns of aerial insectivores roosting in the great lakes region using weather surveillance radar. *Global Change Biology*, 29(5):1407–1419, 2023.
- [32] KeumJi Kim and SeongHwan Yoon. Assessment of building damage risk by natural disasters in south korea using decision tree analysis. *Sustainability*, 10(4), 2018. ISSN 2071-1050. doi: 10.3390/su10041072. URL https://www.mdpi.com/2071-1050/10/4/1072.

- [33] Liwei Deng and Yue Wang. Post-disaster building damage assessment based on improved u-net. *Scientific Reports*, 12(1):2045–2322, 2022.
- [34] DIUx-xView. Diux-xview/xview2_first_place: 1st place solution for "xview2: Assess building damage" challenge. URL https://github.com/DIUx-xView/xView2_first_place.
- [35] The xview2 ai challenge. URL https://www.ibm.com/cloud/blog/the-xview2-ai-challenge.
- [36] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2015.
- [37] Ritwik Gupta, Richard Hosfelt, Sandra Sajeev, Nirav Patel, Bryce Goodman, Jigar Doshi, Eric Heim, Howie Choset, and Matthew Gaston. Creating xBD: A Dataset for Assessing Building Damage from Satellite Imagery. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2019.
- [38] Hao Su, Jia Deng, and Li Fei-Fei. Crowdsourcing annotations for visual object detection. In Workshops at the Twenty-Sixth AAAI Conference on Artificial Intelligence, 2012.
- [39] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2009.
- [40] Tim C Hesterberg. Estimates and confidence intervals for importance sampling sensitivity analysis. *Mathematical and computer modelling*, 23(8-9):79–85, 1996.
- [41] Marko Žnidarič. Asymptotic expansion for inverse moments of binomial and poisson distributions. *The Open Mathematics, Statistics and Probability Journal*, 1(1), 2009.
- [42] Aad W. van der Vaart and Jon A. Wellner. Weak Convergence and Empirical Processes: With Applications to Statistics. Springer, 1996.

A Derivations

A.1 IS-Count

Take p(s) = 1/|S| and $\tilde{f}(s) = |S|f(s)$, we want $\mathbb{E}_p[\tilde{f}(s)] = \sum_{s \in S} \frac{1}{|S|} |S|f(s) = F(S)$. Importance sampling with proposal q gives

$$F(S) = \mathbb{E}_p[\tilde{f}(s)]$$

$$= \mathbb{E}_q \left[\frac{p(s)}{q(s)} \tilde{f}(s) \right]$$

$$= \mathbb{E}_q \left[\frac{1/|S|}{q(s)} |S| f(s) \right]$$

$$= \mathbb{E}_q \left[\frac{f(s)}{q(s)} \right]$$

A.2 DISCount

Take $q = \bar{g}_S$ in IS-Count, then

$$F(S) = \mathbb{E}_q \left[\frac{f(s)}{q(s)} \right]$$
$$= \mathbb{E}_{\bar{g}_S} \left[\frac{f(s)}{g(s)/G(S)} \right]$$
$$= G(S) \cdot \mathbb{E}_{\bar{g}_S} \left[\frac{f(s)}{g(s)} \right]$$

A.3 k-DISCount

Proof of Claim 1. For any m > 0 we have

$$\mathbb{E}\left[\hat{F}_{k\mathrm{DIS}}(S) \mid n(S) = m\right] = \mathbb{E}\left[G(S) \cdot \frac{1}{m} \sum_{i=1}^{n} w_{i} \cdot \mathbb{I}[s_{i} \in S] \mid n(S) = m\right]$$

$$= G(S) \cdot \frac{1}{m} \sum_{i=1}^{n} \mathbb{E}\left[w_{i} \cdot \mathbb{I}[s_{i} \in S] \mid n(S) = m\right]$$

$$= G(S) \cdot \frac{1}{m} \sum_{i=1}^{n} \Pr[s_{i} \in S \mid n(S) = m] \cdot \mathbb{E}\left[w_{i} \mid s_{i} \in S, n(S) = m\right]$$

$$= G(S) \cdot \frac{1}{m} \sum_{i=1}^{n} \Pr[s_{i} \in S \mid n(S) = m] \cdot \mathbb{E}\left[w_{i} \mid s_{i} \in S\right]$$

$$= G(S) \cdot \frac{1}{m} \sum_{i=1}^{n} \frac{m}{n} \cdot \mathbb{E}\left[w_{i} \mid s_{i} \in S\right]$$

$$= G(S) \cdot \frac{1}{m} \sum_{i=1}^{n} \frac{m}{n} \cdot \frac{F(S)}{G(S)}$$

$$= F(S)$$

In the third line, we used the fact that $\mathbb{E}\left[h(X)\cdot\mathbb{I}[X\in A]\right]=\Pr[x\in A]\cdot\mathbb{E}[h(X)\mid X\in A]$ for any random variable X and event A (see Lemma 1 below). In the fourth line we used the fact that s_i is conditionally independent of n(S) given $s_i\in S$, since $n(S)=\mathbb{I}[s_i\in S]+\sum_{j\neq i}\mathbb{I}[s_j\in S]$ and the latter sum is independent of s_i . In the fifth line we used the fact that $\Pr[s_i\in S\mid n(S)=m]=\frac{m}{n}$ because $n(S)=\sum_{j=1}^n\mathbb{I}[s_j\in S]$ and the terms in the sum are exchangeable. In the sixth line we

computed the conditional expectation as follows using the fact that the conditional density of s_i given $s_i \in S$ is equal to $g(s_i)/G(S)$:

$$\mathbb{E}\left[w_i \mid s_i \in S\right] = \mathbb{E}\left[\frac{f(s_i)}{g(s_i)} \mid s_i \in S\right] = \sum_{s \in S} \frac{g(s_i)}{G(S)} \cdot \frac{f(s_i)}{g(s_i)} = \frac{1}{G(S)} \sum_{s \in S} f(s) = \frac{F(S)}{G(S)}.$$

The unconditional bias of k-DISCOUNT can also be analyzed:

Claim 4. Let $p(S) = \Pr[s_i \in S] = G(S)/G(\Omega)$. The bias of the k-DISCOUNT estimator is $\mathbb{E}[\hat{F}_{k\text{DIS}}(S)] - F(S) = -(1 - p(S))^n$.

In particular, bias decays exponentially with n and quickly becomes negligible, with magnitude at most ϵ for $n \ge \log(1/r)/\log(1/\epsilon)$ and r = 1 - p(S). Further, the bias is easily computable from the detector counts and therefore known prior to sampling, and the event that leads to a biased estimate (n(S) = 0) is observed after sampling. All these factors make bias a very minor concern.

Proof of Claim 4. Using Claim 1, we compute the unconditional expectation as

$$\mathbb{E}\left[\hat{F}_{k\mathrm{DIS}}(S)\right] = \Pr[n(S) = 0] \cdot \mathbb{E}\left[\hat{F}_{k\mathrm{DIS}}(S) \mid n(S) = 0\right] + \Pr[n(S) > 0] \cdot \mathbb{E}[\hat{F}_{k\mathrm{DIS}}(S) \mid n(S) > 0]$$

$$= \Pr[n(S) > 0] \cdot F(S)$$

$$= \left(1 - \left(1 - \frac{G(S)}{G(\Omega)}\right)^n\right) \cdot F(S).$$

In the final line, $1 - G(S)/G(\Omega)$ is probability that $s_i \notin S$ for a single i, and $(1 - G(S)/G(\Omega))^n = \Pr[n(S) = 0]$ is the probability that $s_i \notin S$ for all i. Rearranging gives the result. \square

Lemma 1. $\mathbb{E}\left[h(X) \cdot \mathbb{I}[X \in A]\right] = \Pr[x \in A] \cdot \mathbb{E}[h(X) \mid X \in A]$ for any random variable X and event A.

Proof. Observe

$$\begin{split} \mathbb{E}\left[h(X)\cdot\mathbb{I}[X\in A]\right] &= \sum_{x}\Pr[X=x]h(x)\,\mathbb{I}[x\in A] \\ &= \sum_{x}\Pr[X=x,X\in A]h(x) \\ &= \Pr[X\in A]\sum_{x}\Pr[X=x\mid X\in A]h(x) \\ &= \Pr[X\in A]\cdot\mathbb{E}[h(X)\mid X\in A]. \end{split}$$

A.4 Optimal allocation of samples for DISCOUNT to disjoint regions

Proof of Claim 2. The proof is similar to that of Theorem 5.6 in [18]. We prove the claim for k=2; the proof generalizes to larger k in an obvious way. The variance of DISCOUNT on S_i is

$$\operatorname{Var}(\hat{F}_{\operatorname{DIS}}(S_i)) = \frac{G(S_i)^2 \cdot \sigma^2(S_i)}{n_i} = \sigma^2 \frac{G(S_i)^2}{n_i}.$$

We want to minimize $\sum_{i} \operatorname{Var}(\hat{F}_{DIS}(S_i))$, which with k=2 is proportional to

$$V = \frac{G(S_1)^2}{n_1} + \frac{G(S_2)^2}{n_2}.$$

¹The k-DISCOUNT estimator can be debiased by dividing by $u=1-(1-p(S))^n<1$. However, this leads to higher overall error: if n(S)=0, the estimator is unchanged, and conditioned on the event n(S)>0 the estimator becomes biased and has higher variance by a factor of $1/u^2>1$.

By the Cauchy-Shwarz inequality, for any $n_1, n_2 > 0$,

$$Vn = \left(\frac{G(S_1)^2}{n_1} + \frac{G(S_2)^2}{n_2}\right)(n_1 + n_2) \ge (G(S_1) + G(S_2))^2.$$

If we substitute $n_i = G(S_i)/Z$ for any Z on the left of the inequality and simplify, we see the inequality becomes tight, so the minimum is achieved. We further require $\sum_i n_i = n$, so choose Z so

$$n_i = n \cdot \frac{G(S_i)}{G(\Omega)} = np(S_i).$$

A.5 k-DISCOUNT variance

Recall that $p(S) = \Pr[s_i \in S] = G(S)/G(\Omega)$ is the probability of a sample landing in S under the sampling distribution \bar{g}_{Ω} . Define

$$\sigma^{2}(S) = \operatorname{Var}(f(s_{i})/g(s_{i}) \mid s_{i} \in S) = \sum_{s \in S} \frac{g(s)}{G(S)} \cdot \left(\frac{f(s)}{g(s)} - \frac{F(S)}{G(S)}\right)^{2}.$$

to be the variance of the importance weight for $s_i \sim \bar{g}_S$.

Claim 5. Let r = 1 - p(S). The variance of the k-DISCOUNT estimator is given by

$$\operatorname{Var}(\hat{F}_{kDIS}(S)) = G(S)^2 \cdot \sigma^2(S) \cdot (1 - r^n) \cdot \mathbb{E}\left[\frac{1}{n(S)} \left| n(S) > 0 \right] + F(S)^2 \cdot r^n \cdot (1 - r^n).\right]$$

where
$$(1-r)^n \mathbb{E} [1/n(S) \mid n(S) > 0] = \sum_{j=1}^n (1/j) \cdot \text{Binomial}(j; n, p(S)).$$

The second term in the variance arise from the possibility that no samples land in S; it decays exponentially in n and is negligible compared to the first term. The first term can be compared to the variance $G(S)^2 \cdot \sigma^2(S) \cdot \frac{1}{m}$ of importance sampling with exactly m samples allocated to S and the proposal distribution \bar{g}_S , i.e., DISCOUNT. Because the sample size n(S) is random, the correct scaling factor for k-DISCOUNT is $(1-r^n) \mathbb{E}[1/n(S) \mid n(S)>0]$, which it turns out is asymptotically equivalent to 1/(np(S)), i.e., DISCOUNT with a sample size of $m=np(S)=\mathbb{E}[n(S)]$ —see Claim 6 below. We find that for a small expected sample size (around 4) there can be up to 30% "excess variance" due to the randomness in the number of samples (see Figure 6), but that this disappears quickly with larger expected sample size.

Claim 6. Let $\hat{F}_{kDIS,n}$ and $\hat{F}_{DIS,m}$ be the k-DISCOUNT and DISCOUNT estimators with sample sizes n and m, respectively. The asymptotic variance of k-DISCOUNT is given by

$$\lim_{n \to \infty} n \operatorname{Var}(\hat{F}_{kDIS,n}(S)) = G(S)^2 \cdot \sigma^2(S) / p(S).$$

This is asymptotically equivalent to DISCOUNT with sample size $m = \mathbb{E}[n(S)] = np(S)$. That is

$$\lim_{n\to\infty} \frac{\operatorname{Var}(\hat{F}_{kDIS,n}(S))}{\operatorname{Var}(\hat{F}_{DIS,\lceil np(S)\rceil}(S))} = 1.$$

Proof of Claim 5. By the law of total variance,

$$\operatorname{Var}(\hat{F}_{k\text{DIS}}(S)) = \mathbb{E}\left[\operatorname{Var}(\hat{F}_{k\text{DIS}}(S) \mid n(S))\right] + \operatorname{Var}\left(\mathbb{E}[\hat{F}_{k\text{DIS}}(S) \mid n(S)]\right). \tag{1}$$

We will treat each term in Eq. 1 separately. For the first term, from the definition of k-DISCOUNT we see

$$\operatorname{Var}(\hat{F}_{k\text{DIS}}(S) \mid n(S)) = \begin{cases} 0 & n(S) = 0\\ G(S)^2 \cdot \frac{\sigma^2(S)}{n(S)} & n(S) > 0 \end{cases}.$$

Therefore

$$\mathbb{E}\left[\operatorname{Var}(\hat{F}_{k\text{DIS}}(S) \mid n(S))\right] = G(S)^2 \cdot \Pr[n(S) > 0] \,\mathbb{E}\left[\frac{\sigma^2(S)}{n(S)} \,\middle|\, n(S) > 0\right]$$
$$= G(S)^2 \cdot \sigma^2(S) \cdot (1 - r^n) \cdot \mathbb{E}\left[\frac{1}{n(S)} \,\middle|\, n(S) > 0\right].$$

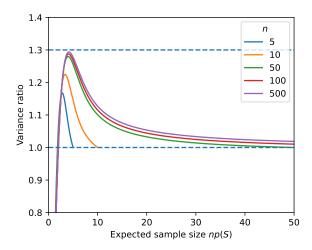


Figure 6: Ratio of variance of k-DISCOUNT with expected sample size np(S) to DISCOUNT with np(S) samples; uses formula from Claim 5 (first term only).

In the last line, we used the fact that $n(S) \sim \text{Binomial}(n, p(S))$, so $\Pr[n(S) > 0] = 1 - r^n$ where r = 1 - p(S). The summation for $(1 - r^n) \mathbb{E}[1/n(S) \mid n(S) > 0]$ follows from the same fact.

For the second term in Eq. (1), from the definition of k-DISCOUNT and conditional unbiasedness (Claim 1), we have

$$\mathbb{E}[\hat{F}_{k\mathrm{DIS}}(S) \mid n(S)] = \begin{cases} 0 & \text{if } n(S) = 0 \\ F(S) & \text{if } n(S) > 0 \end{cases}$$
$$= F(S) \cdot \mathrm{Bernoulli}(1 - r^n).$$

The variance is therefore

$$\operatorname{Var}\left(\mathbb{E}[\hat{F}_{k\text{DIS}}(S) \mid n(S)]\right) = F(S)^2 \cdot r^n \cdot (1 - r^n).$$

Putting the two terms together yields the result.

Proof of Claim 6. By Claim 5 we have

$$\lim_{n \to \infty} n \operatorname{Var}(\hat{F}_{k \text{DIS},n}(S)) = \lim_{n \to \infty} n \cdot G(S)^2 \cdot \sigma^2(S) \cdot (1 - r^n) \cdot \mathbb{E}\left[\frac{1}{n(S)} \left| n(S) > 0 \right| + \lim_{n \to \infty} n \cdot F(S)^2 \cdot r^n \cdot (1 - r^n). \right]$$
(2)

The second limit on the right side is zero, because $nr^n \to 0$ as $n \to \infty$ (recall that r < 1) and the other factors are bounded. We will show the first limit on the right side is equal to $G(S)^2 \cdot \sigma^2(S)/p(S)$, which will prove the first part of the result. The asymptotic expansion of [41] (Corollary 3) states that

$$(1-r^n)\mathbb{E}[1/n(S) \mid n(S) > 0] = \frac{1}{np(S)} + \mathcal{O}\left(\frac{1}{(np(S))^2}\right).$$

Using this expansion in the limit gives:

$$\lim_{n \to \infty} n \cdot G(S)^2 \cdot \sigma^2(S) \cdot (1 - r^n) \cdot \mathbb{E}\left[\frac{1}{n(S)} \left| n(S) > 0 \right] \right]$$

$$= \lim_{n \to \infty} G(S)^2 \cdot \sigma^2(S) / p(S) \cdot (1 + \mathcal{O}(1/n))$$

$$= G(S)^2 \cdot \sigma^2(S) / p(S)$$

The variance of DISCOUNT with sample size m is $Var(\hat{F}_{DIS,m}(S)) = G(S)^2 \cdot \sigma^2(S)/m$. Setting $m = \lceil np(S) \rceil$ and using the second to last line above we have

$$\lim_{n \to \infty} \frac{\operatorname{Var}(\hat{F}_{k\mathrm{DIS},n}(S))}{\operatorname{Var}(\hat{F}_{\mathrm{DIS},\lceil np(S)\rceil}(S))} = \lim_{n \to \infty} \frac{n \operatorname{Var}(\hat{F}_{k\mathrm{DIS},n}(S))}{n \operatorname{Var}(\hat{F}_{\mathrm{DIS},\lceil np(S)\rceil}(S))}$$
$$= \lim_{n \to \infty} \frac{G(S)^2 \cdot \sigma^2(S)/p(S) \cdot (1 + \mathcal{O}(1/n))}{n \cdot G(S)^2 \cdot \sigma^2(S)/\lceil np(S)\rceil}$$
$$= 1$$

A.6 Control Variates

Recall that with control variates the weight is redefined as

$$w_i = (f(s_i) - h(s_i))/g(s_i).$$

The expectation of the weight given $s_i \in S$ is

$$\mathbb{E}[w_i \mid s_i \in S] = \sum_{s \in S} \frac{g(s)}{G(S)} \frac{f(s) - h(s)}{g(s)} = \frac{F(S) - H(S)}{G(S)}$$

Therefore

$$\mathbb{E}[\bar{w}_{cv}(S) \mid n(S) > 0] = \frac{F(S) - H(S)}{G(S)}$$

Therefore

$$\mathbb{E}[\hat{F}_{k\text{DIS}cv}(S) \mid n(S) > 0] = G(S) \cdot \frac{F(S) - H(S)}{G(S)} + H(S) = F(S)$$

A.7 Confidence Intervals

Proof of Claim 3. Let w_1, w_2, \ldots be an iid sequence of importance weights for samples in S, i.e., $w_i = f(s_i)/g(s_i)$ for $s_i \sim \bar{g}_S$. Each weight w_i has mean F(S)/G(S) and variance $\sigma^2(S)$. Let $\bar{\omega}_n = \frac{1}{n} \sum_{i=1}^n w_i$. By the central limit theorem,

$$\sqrt{n}(\bar{\omega}_n - F(S)/G(S)) \xrightarrow{D} \mathcal{N}(0, \sigma^2(S))$$

Recall that $\hat{F}_{k\mathrm{DIS},n}(S) = G(S) \cdot \bar{w}_n(S)$ where $\bar{w}_n(S)$ is the average of the importance weights for samples that land in S when drawn from all of Ω (for clarity in the proof we add subscripts for sample size to all relevant quantities). It is easy to see that $\bar{w}_n(S)$ is equal in distribution to $\bar{\omega}_{n(S)}$ where $n(S) \sim \mathrm{Binomial}(n,p(S))$ and n(S) is independent of the sequence of importance weights — this follows from first choosing the number of samples that land in S and then choosing their locations conditioned on being in S. From Theorem 3.5.1 of [42] (with $N_n = n(S)$ and $c_n = np(S)$) it then follows that

$$\sqrt{n(S)} (\bar{w}_n(S) - F(S)/G(S)) \stackrel{D}{\longrightarrow} \mathcal{N}(0, \sigma^2(S))$$

Rearranging yields

$$\frac{\hat{F}_{k\text{DIS}}(S) - F(S)}{G(S)/\sqrt{n(S)}} \xrightarrow{D} \mathcal{N}(0, \sigma^2(S))$$

After dividing by $\hat{\sigma}_n(S)$, the result follows from Slutsky's lemma if $\hat{\sigma}_n^2(S) \xrightarrow{P} \sigma^2(S)$, which follows from a similar application of Theorem 3.5.1 of [42].

B Counting tasks

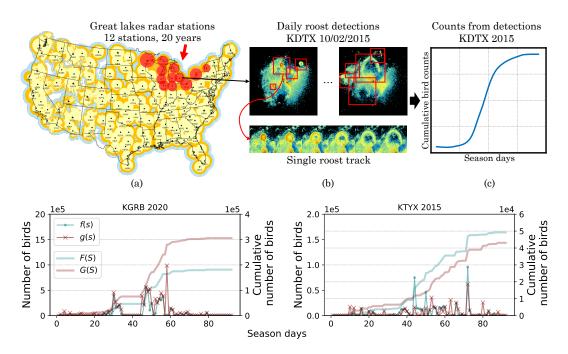


Figure 7: **Counting roosting birds in radar images.** (a) The US weather radar network has collected data for 30 years from 143+ stations and provides an unprecedented opportunity to study long-term and wide-scale biological phenomenon such as roosts. (b) Counts are collected for each day s by running the detector using all radar scans for that day to detect and track roost signatures and then mapping to bird counts using the radar "reflectivity" within the tracks. The figure shows two scans for the KDTX station (Detroit, MI) on the same day, along roost detections which appear as expanding rings. By tracking these detections across a day one can estimate the number of birds in each roost. (c) Cumulative bird counts in the complete roosting season by aggregating counts across all tracked roosts and days. (bottom) Examples of *true* bird counts (blue) and detector counts (red) during a roosting season for station KGRB 2020 and KTYX 2015.

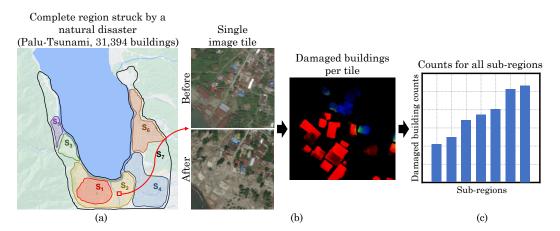


Figure 8: Counting damaged buildings in satellite images. Building damage assessment from satellite images [32, 33] is often used to plan humanitarian response after a natural disaster strikes. (a) We consider the Palu Tsunami from 2018; the data consists of 113 high-resolution satellite images labeled with 31,394 buildings and their damage levels. (b) Counts are collected per tile using beforeand after-disaster satellite images. Colors indicate different levels of damage (e.g., red: "destroyed"). (c) Damaged building counts per sub-region.

C Palu Tsunami regions

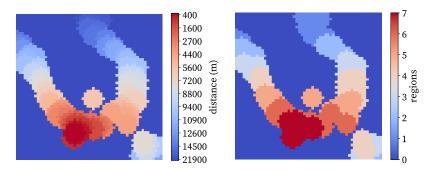


Figure 9: (left) Distance from the area with the most damaged buildings. (right) Regions defined for damaged building counting (See \S 4.2).