# Model-assisted domain estimation of postfire tree regeneration in the western US using nearest neighbor techniques

**David L. R. Affleck**[a] **and George C. Gaines, III** [b]

[a]Department of Forest Management, University of Montana, Missoula, MT 59812, USA; [b]USDA Forest Service, Rocky Mountain Research Station Forest Inventory & Analysis, Missoula, MT 59808, USA

Corresponding author: **George C. Gaines III** (email: george.gaines@usda.gov)

## Abstract

Many nations administer national forest inventory programs for unbiased estimation of forest attributes over broad spatial and temporal extents. However, management and conservation decisions often demand reliable estimates for finer spatiotemporal domains. In the western US, wildfire activity is expanding and postfire regeneration must contend with a warmer, drier climate. We evaluate the potential of K nearest neighbor (KNN) strategies for estimation of stocking across postfire measurements of Forest Inventory & Analysis plots in 11 western US states, and subsequently for model-assisted (MA) estimation of stocking over domains defined by aggregations of burned areas within individual states and 4-year periods. In particular, we develop and evaluate a form of constrained KNN that allows for unbiased MA domain estimation under simple random sampling by drawing only on measurements external to a domain of interest. KNN strategies based on geographically, radiometrically, and climatically proximate measurements are found to provide more accurate estimates of stocking at the plot level than domain means. Applying the selected external KNN strategy also reduced standard errors of MA domain estimates by 16% over direct domain estimators, but bias correction introduces substantial variability over synthetic estimates. Further applications of the external constraint imposed on KNN are discussed.

**Key words:** small area estimation, forest inventory, forest regeneration, wildland fire, model-assisted estimation

## Introduction

The size, frequency, and severity of wildfires in the western United States increased dramatically in recent decades under a warmer, drier climate (Parks and Abatzoglou 2020), with many western forests burning more now than at any other point in the past millennium (Higuera et al. 2021). Climate and other factors limiting seedling recruitment increasingly lead to postfire regeneration failures and conversion of previously forested areas to non-forest cover types (Stevens-Rumann et al. 2017; Young et al. 2019), fueling broad concern for the future of western forests. While a growing body of literature documents changing fire regimes and assesses the biophysical mechanisms governing postfire regeneration, a need persists for precise, high-resolution estimates of postfire forest density in the western US.

Many nations maintain national forest inventory (NFI) programs that facilitate inference on forest traits over broad spatial and temporal extents, yet precise estimates for spatiotemporal subsets of forest populations, or domains, are of increasing interest throughout the world. The US Forest Service (USFS) Forest Inventory & Analysis (FIA) program administers an NFI comprised of a ground plot network that permits design-based inference on postfire forest conditions

(Bechtold and Patterson 2005). This NFI's coverage in time and space, however, is too coarse to deliver precise estimates of postfire forest attributes in burned area domains of management interest using traditional estimation techniques (Gaines and Affleck 2021). Small area estimation (SAE) techniques, documented in detail by Rao and Molina (2015), offer a possible solution. SAE techniques borrow extra-domain sample information to increase effective sample sizes for domains requiring more precise estimates. Such domains are commonly referred to as "small areas" though they may actually extend over large spatiotemporal extents.

Multiple approaches to SAE have been developed for forest inventory and monitoring applications. Gaines and Affleck (2021) borrowed extra-domain sample data in time and space for indirect domain estimation. However, most techniques also leverage statistical relationships between ground observations of vegetation traits and auxiliary variables. Some relate ground observations of target forest attributes to auxiliary data at the observational unit (e.g., field plot) level (e.g., Breidenbach and Astrup 2012), while others exploit relationships aggregated to the area (i.e., domain) level (e.g., Coulston et al. 2021). Some specify a functional association between response and explanatory variables via regression-based tech-

niques (e.g. Hill et al. 2018). Others specify an explanatory variable space from which to draw observations via K nearest neighbors (KNN) techniques (e.g. Bell et al. 2022). Haakana et al. (2020) explore poststratified estimation techniques in an SAE context. Across these studies, inference on population parameters proceed from design-based (e.g. Breidenbach and Astrup 2012; Baffetta et al. 2009), model-based (Ver Planck et al. 2018; Coulston et al. 2021), and hybrid (Magnussen et al. 2014) inferential paradigms.

Hill et al. (2018) compared regression estimators in a design-based framework, while Breidenbach and Astrup (2012) compared regression estimators in both design- and model-based frameworks. In both cases, regression estimators increased the precision of domain estimates relative to simple or weighted sample means by wide margins. Coulston et al. (2021) compared a poststratified domain estimator employed in a design-based inferential framework with empirical best linear unbiased predictor (EBLUP) estimators based on area-level models, relating forest removal estimates to sawmill survey and Landsat-based tree cover loss data.

Baffetta et al. (2009) compared direct expansion domain estimators with model-assisted (MA) domain estimators employing linear or KNN models of timber volume. The models were specified at the unit level and drew on spectral data from satellite images as auxiliary data. The MA estimators consisted of domain sums of pixel-level estimates (i.e., synthetic estimators) plus sample-based bias correction terms. Estimators of this form will be more precise than domain sample means provided that variability in the differences between observed and estimated values is smaller than the variability in the observed values themselves (Breidt and Opsomer 2017). In other words, such estimators increase precision of domain estimates (or, indeed, of population-level estimates) if the unit-level estimation approach closely approximates the actual values of the attribute of interest so that the bias correction component is small. Asymptotically, these estimators are unbiased and amenable to variance approximation (Särndal et al. 2003, ch. 6; Breidt and Opsomer 2017). McRoberts et al. (2022) trace the nomenclature and differentiate among alternative forms of MA estimators; McConville et al. (2020) provide an overview of MA estimation in the context of forest inventory.

Assessing the accuracy of small area estimators can be complex. Direct domain estimators, which only use domain sample observations, are generally unbiased but often unstable. Indirect estimators that borrow extra-domain sample data to increase precision incur bias. For design-based inference, estimators of the bias and mean squared error (MSE) of indirect domain estimators must draw on the estimated variance of the direct domain estimator (see González and Waksberg 1973; Marker 1995; Rao and Molina 2015). This estimated variance is, in turn, dependent on the domain sample—the often insufficient size of which leads to high instability. Gaines and Affleck (2021) implemented new and existing estimators of the design-MSE of indirect domain estimators of postfire tree regeneration in the western US, but found these to be unreliable (highly variable and frequently negative). They concluded that only the variance of their indirect estimators could be feasibly estimated, yielding incomplete accuracy assessments of inherently biased estimators. This, in turn, motivated the development of unbiased domain estimators.

McRoberts (2012) evaluated alternative KNN-based approaches to SAE of tree stem volume in a model-based inferential framework. Under a model-based approach to small area inference, domain differences are described with an explicit probabilistic model, and estimators are considered unbiased if model diagnostics suggest the model to be correctly specified. This not only highlights a need for approaches to model validation (see, e.g., McRoberts 2012), but also provides a means of derivation of model-based MSE estimators.
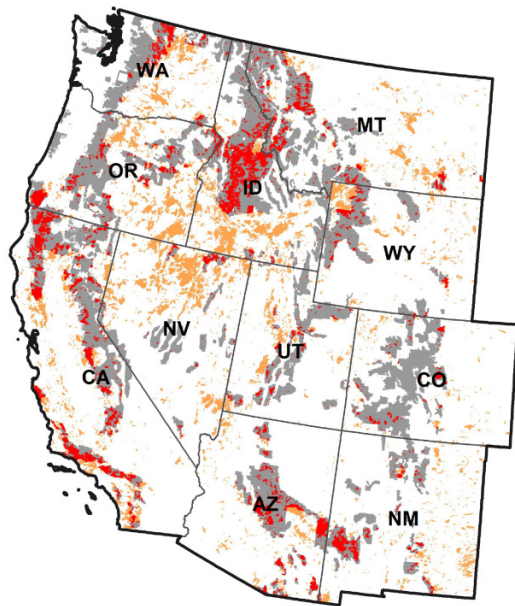
An advantage of the design-based inferential framework is the potential to pursue inherently design-unbiased or approximately design-unbiased MA small area estimators, precluding the need for bias estimation. Accuracy assessments are then permitted to focus entirely on well-documented variance estimation techniques (e.g., Särndal et al. 2003; Baffetta et al. 2009; Breidenbach and Astrup 2012; Hill et al. 2018).

Mandallaz et al. (2013) developed domain regression estimators with partially exhaustive auxiliary data. They included an additional small area indicator variable in their small area estimator to ensure a zero-mean residual property, which simplified variance estimation. Pertinently to this work, they describe a condition of externality, i.e. that the regression coefficients used to obtain unit-level estimates for a given domain were not estimated using the inventory data from that domain (and instead, perhaps, using data from a previous inventory). They further apply this condition to derive tractable forms for the small area regression estimator and its variance estimator. They cite empirical evidence (Mandallaz 2012) that this condition can be relaxed for large samples and thus fit their regression coefficients using observations from within and beyond the domain of interest.

When synthetic estimators draw on KNN rather than regression strategies, it is unclear whether the size of the complete sample will have a bearing on the dependence of unit-level estimates on observations from within the same domain. However, externality can be achieved by formulating synthetic domain estimates using only data from outside the domain of interest. We apply this idea to develop a novel KNN-based domain estimator. It can be viewed as a domain-level difference estimator (Särndal et al. 2003, p. 221) constructed from proxy values for the attribute of interest that draw only on sample observations in other domains. We show this estimator to be unbiased under simple random sampling (SRS), given a nonzero domain sample size and the availability of extra-domain data. The desirability of design-unbiased estimators is especially high for official governmental statistical programs driven by probabilistic sample data, including many NFI programs.

One objective of this research was to determine the accuracy with which postfire stocking in forested plots across the western US could be estimated using KNN methods with mapped topographic, climatic, and radiometric data products. The second objective was to evaluate whether precise MA estimators of stocking proportions for spatiotemporal domains within the western US could be developed using KNN-based assisting models. Of particular interest was whether

**Fig. 1.** Study area spanning the contiguous western USA (US Census Bureau 2023). Areas identified as burned by MTBS between 1984 and 2018 (MTBS 2023) are shown in red where they overlap USFS NFS lands and in orange on all other ownerships; unburned NFS lands (USFS 2023) are shown in gray.



domain estimators drawing only on extra-domain data for KNN estimation could be used without appreciable loss of precision. The spatiotemporal scope of the study, KNN estimation strategies, and domain estimators are described below. Results concerning the application of these estimators to burned area domains are then discussed, along with their advantages and limitations relative to synthetic and model-based estimators.

## Methodology

### Application scope

Our focus is on forested areas within the 11 contiguous western US states (Fig. 1) that burned within the period 2001–2018. This is a diverse landscape largely dominated by coniferous species but comprising subalpine forests along coastal and interior cordillera, as well as temperate rainforests and dry, woodland ecosystems (Perry et al. 2022). Forests are maintained under diverse ownerships across this region but, as described further below, the National Forest System (NFS; Fig. 1) lands managed by the USFS define the domains of interest in this study given the imperative to maintain forest cover across these areas. Both wildfire and managed fire are important disturbance agents and forest management tools in this region, although the extent, frequency, and severity of fire has varied over time and across the landscape. Historic fire perimeters (see Fig. 1) were obtained from the Monitoring Trends in Burn Severity (MTBS; Eidenshink et al. 2007) program, which maps all fires greater than 404 ha occurring across western states dating back to 1984.

The USFS FIA program provides repeat field measures of vegetation condition across all lands in the western US states. Since 2001, this program has been implemented as an unaligned equal-probability sample. Plots are located at random across a hexagonal tessellation with a density of approximately 1 location per 2400 ha, but hexagonal cells (and thus plots) are divided into 10 interpenetrating panels remeasured on decadal intervals (Bechtold and Patterson 2005). This annualized FIA program was implemented first in Arizona, California, and Oregon, but by 2011 was active in all 11 western states (albeit with interruptions in some years; see Burrill et al. 2021, Appendix J). On an FIA plot, trees above 12.7 cm diameter at breast height (dbh; 1.37 m) are tallied within four spatially disjoint subplots (7.32 m radius), each of which contains one 2.07 m radius microplot that is searched for saplings and seedlings below 12.7 cm dbh (Bechtold and Patterson 2005). Plots are located on all lands, but here we focus on those containing at least one forested condition (Burrill et al. 2021, p. 2–38). Additionally, we focus on subplot measurements lying within MTBS fire perimeters, taken at least 2 years postfire, and made between the years 2001 and 2018. In total, this amounted to to 5660 measurements at 4297 unique FIA non-intensified plot locations (1355 locations were measured twice; four were measured 3 times).

The attribute of central interest for estimation was postfire tree stocking, defined here as a binary variable taking the value 1 where density exceeded 740 trees ha$^{-1}$ and 0 otherwise. The threshold of 740 trees ha$^{-1}$ corresponds to the presence of at least one seedling per FIA microplot. However, determination of stocking was made at the plot level using counts of live trees of all sizes (weighted according to their appropriate expansion factors) identified on all burned subplots. Averaging observed stocking over subplots, let $y_i = y(\mathbf{w}_i, t_i)$ represent the stocking for the $i^{th}$ plot measurement, obtained in year $t_i$ at a plot centered at coordinate location $\mathbf{w}_i$. Actual tree density was also considered, but was highly skewed and very frequently 0, complicating inference on this scale (see Gaines and Affleck 2021).

True geographic coordinates were provided by the USFS for burned FIA plots and used to obtain climatic and radiometric attributes available through Google Earth Engine (Table 1; Gorelick et al. 2017). Monthly climatic attributes with a spatial resolution of approximately 4.6 km were drawn from the TerraClimate dataset (Abatzoglou et al. 2018) and averaged over 1984–2019. Annual means of these attributes used all monthly data; summer means used June-August data. Vegetation cover metrics for each year were obtained at a spatial resolution of 30 m from the Rangeland Analysis Platform (RAP) dataset (Allred et al. 2021). Based on initial analyses of the relationships between FIA plot tree density and RAP-estimated tree cover, the tree cover variable was square-root transformed (see Gaines 2022). In each case, burned subplots were intersected with source raster data and extracted attributes then averaged to the FIA plot level. Plot eastings and northings were obtained after projecting coordinates to an equidistant conic basis (standard parallels 33°N and 45°N; datum NAD83).

As detailed below, KNN estimation was evaluated as a tool to leverage the predictor variables of Table 1 for estimation

**Table 1.** Candidate predictor variables, sources, and units.

| Predictor | Units | Definition |
|---|---|---|
| EAST | m | Location easting in equidistant conic projection |
| NRTH | m | Location northing in equidistant conic projection |
| MLAG[1] | yr | Time interval between fire and measurement |
| TREE[2] | $\%^{\frac{1}{2}}$ | Percent tree cover (square-root transformed) |
| FOGR[2] | % | Percent cover of forbs and grasses |
| SHRB[2] | % | Percent shrub cover |
| TMXY[3] | °C | Mean annual maximum temperature |
| TMXS[3] | °C | Mean summer maximum temperature |
| AETY[3] | mm | Annual actual evapotranspiration |
| AETS[3] | mm | Summer actual evapotranspiration |
| DEFY[3] | mm | Annual water deficit |
| DEFS[3] | mm | Summer water deficit |

[1] From MTBS (Eidenshink et al. 2007).
[2] From the RAP (Allred et al. 2021).
[3] From TerraClimate (Abatzoglou et al. 2018).

**Table 2.** Characteristics of burned area domains by time period; summaries are for domains with at least one plot measurement.

| Burn period | Number of domains | Postfire lag (yr) min | Postfire lag (yr) max | Area (ha) min | Area (ha) median | Area (ha) max | $n_d$ min | $n_d$ median | $n_d$ max |
|---|---|---|---|---|---|---|---|---|---|
| 1984–1987 | 57 | 17 | 31 | 518 | 44 551 | 379 919 | 1 | 2.0 | 16 |
| 1988–1991 | 83 | 13 | 27 | 4317 | 73 020 | 197 744 | 1 | 2.0 | 12 |
| 1992–1995 | 80 | 9 | 23 | 9916 | 65 531 | 324 020 | 1 | 2.0 | 11 |
| 1996–1999 | 85 | 5 | 19 | 3857 | 53 407 | 331 863 | 1 | 2.0 | 10 |
| 2000–2003 | 125 | 2 | 15 | 29 115 | 172 507 | 431 524 | 1 | 6.0 | 24 |
| 2004–2007 | 87 | 2 | 11 | 8763 | 148 643 | 699 877 | 1 | 5.0 | 33 |
| 2008–2011 | 48 | 2 | 7 | 11 415 | 64 141 | 556 363 | 1 | 3.5 | 25 |
| 2012–2015 | 22 | 2 | 3 | 46 131 | 221 478 | 705 360 | 1 | 6.0 | 30 |
| All | 587 | 2 | 31 | 518 | 89 095 | 705 360 | 1 | 3.0 | 33 |

of mean stocking at a domain level. The domains $\mathcal{D}_d$ ($d = 1$, 2, …) of interest here are spatiotemporal in nature, defined both by the spatial extent of fires occurring in the western US states in specific years and by a particular length of time over which this extent has been allowed to develop without subsequent fire. For specificity, consider first a base domain $\mathcal{M}(2000, 10)$ corresponding to the 2010 condition of all NFS lands in Montana that burned in 2000, less any areas that reburned between 2001 and 2010. Note that this base domain is defined by two temporal parameters: the burn year (2000) and the lag-to-assessment-year (10 years). The domains of interest $\mathcal{D}_d$ consist of aggregations of base domains like $\mathcal{M}(2000, 10)$, such as all Montana NFS lands that burned in the 4-year span 2000–2003 and that experienced 10 years of regrowth without subsequent fire (e.g., $\{\mathcal{M}(2000, 10), \mathcal{M}(2001, 10), \mathcal{M}(2002, 10), \mathcal{M}(2003, 10)\}$). All the domains we consider below take such a form and so can be referenced by a state, a 4-year span of burn years, and a postfire measurement lag. The domains are confined to NFS lands and to MTBS-mapped fires. Characteristics of collections of such domains as grouped by burn period are summarized in Table 2. We emphasize that all lands constitute the focal population for FIA and that the number of measurements taken within any subsequently delineated

spatiotemporal domain is a random variable with respect to the sampling design.

## KNN model development

Domain estimation is supported by KNN-estimated stocking surfaces formulated from the attributes summarized in Table 1. To identify the best subsets of those attributes and the number of neighbors ($K$) to draw on, KNN was first implemented at the plot level. Adapting the notation of Baffetta et al. (2009), a corresponding KNN estimate of tree stocking $\widehat{y}_i(s) = \widehat{y}(\mathbf{w}_i, t_i, s)$ for plot measurement $i$ is

$$(1) \qquad \widehat{y}_i(s) = \frac{1}{K} \sum_{k=1}^{K} y_{H(i,k,s)}$$

where $H(i, k, s)$ is a random variable returning the label of the $k^{\text{th}}$-nearest measurement of tree stocking within sample set $s$ to measurement $i$. As described further below, the sample set $s$ may include all available measurements or only measurements from beyond a domain of interest. In the former case, referred to below as unconstrained KNN, we will denote estimates as $\widehat{y}_i(s_a)$, indicating that all measurements in dataset—including any internal to the domain in which $y_i$

was observed—were considered as potential neighbors. In the latter case, subsequently referred to as external KNN, we will denote estimates as $\widehat{y}_i(s_{-d})$ to reinforce that only measurements external to the domain $d$ of the focal location were considered as potential neighbors.

Proximity for the purposes of KNN estimation was defined over spaces spanned by subsets of the predictors given in Table 1. With $\mathbf{x}_i = \mathbf{x}(\mathbf{w}_i, t_i)$ being the vector of selected predictor values for plot measurement $i$, the squared distance to measurement $j$ was defined as

$$(2) \qquad d_{\mathbf{x}}(i, j) = (\mathbf{x}_i - \mathbf{x}_j)^T \Sigma_{s_a}^{-1} (\mathbf{x}_i - \mathbf{x}_j)$$

where $\Sigma_{s_a}$ is a diagonal matrix containing the variances of the $\mathbf{x}$-attributes as obtained over the complete sample of FIA measurements $s_a$.

The KNN estimator (1) is a simple average of KNNs. The value for $K$, as well as the number and identities of predictors, were selected using a cross-validation approach. Specifically, the set of measurements $s_a$ was randomly divided into 10 folds and KNN estimates produced for each plot measurement using only data from the other 9 folds (necessarily precluding a measurement from being one of its own neighbors). This was carried out for all combinations of predictors (from $p = 1$ predictor to $p = 12$ predictors) and for values of $K$ ranging from 1 to 50. It was also carried out separately for unconstrained and external KNN, with the latter implemented such that only measurements from beyond the domain of the focal measurement ($s_{-d}$) were allowed as candidate neighbors. More specifically, in this case proximity was assessed only for measurements outside of the domain of the focal measurement and outside of the cross-validation fold of the focal measurement; measurements within the same fold or within the same domain as the focal measurement (or even within the spatial extent of that domain) were assigned arbitrarily large distances.

An overall MSE was obtained for each predictor combination and value of $K$, as well as a corresponding standard error (SE) from the variation across folds. MSE was calculated only using observations within domains, i.e.,

$$(3) \qquad \text{MSE} = \frac{\sum_d \sum_{i \in s_d} [y_i - \widehat{y}_i(s)]^2}{\sum_d n_d}$$

where $s_d$ is the subset of measurements falling within domain $\mathcal{D}_d$ and $n_d$ is the size of that subsample. MSE was restricted to domain data even though measurements not belonging to any domains (i.e., falling within burned areas off of NFS lands) were considered as potential neighbors. The strategy (predictor combination and value for $K$) with the lowest MSE was then identified, as well as the simplest strategy (fewest predictors and lowest value of $K$) with MSE lying within one SE of the lowest MSE. The latter was selected as the KNN strategy for domain estimation.

The selected KNN strategy yielded stocking estimates in the unit interval. Rounding these allowed for estimation of classification (stocked vs. non-stocked) accuracies, both overall and by condition. Again, only observations within domains

contributed to accuracy estimates:

$$(4) \qquad \text{CA} = \frac{\sum_d \sum_{i \in s_d} C(y_i, \widehat{y}_i(s))}{\sum_d n_d}$$

where $C(y_i, \widehat{y}_i(s))$ is an indicator variable taking the value 1 if the rounded value of $\widehat{y}_i(s)$ equals $y_i$ and 0 otherwise.

## Domain estimation

Mean stocking of a domain is defined as

$$(5) \qquad \mu_{y,d} = \frac{1}{A_d} \int_{\mathcal{D}_d} y(\mathbf{w}, t_d) \, \mathrm{d}\mathbf{w}$$

where $A_d$ is the area of $\mathcal{D}_d$ and $y(\mathbf{w}, t_d)$ is stocking at location $\mathbf{w}$ in the postfire assessment year $t_d$ associated with that location in $\mathcal{D}_d$. Since stocking requires an explicit spatial support, we take the continuous variable $y(\mathbf{w}, t_d)$ to have the same support as an FIA plot. That is, an FIA plot centered at $\mathbf{w}$ will yield an observation of $y(\mathbf{w}, t_d)$ whether $\mathbf{w}$ is in the interior or near the edge of a burned area domain.

While the number of FIA plots distributed across the western US is large, the number of plots lying within any area burned in a particular interval and measured at a specific time-since-fire is generally small. This follows from the fact that the nominal spatiotemporal frequency of the FIA network is approximately one plot measurement per 24000 ha·yr, which we assume uniform across the western states. As a result, direct estimates (sensu Rao and Molina 2015) of domain means that rely only on measurements taken over the domain spatial extent and at the appropriate time-since-fire (i.e., the domain sample $s_d$) are expected to be insufficiently precise. That is, for an equal-intensity design like FIA, the direct estimator of the Hájek form

$$(6) \qquad \bar{y}_d = \frac{1}{n_d} \sum_{i \in s_d} y_i$$

will have high variance owing to small and variable domain sample sizes ($n_d$). Provided $n_d > 0$, it will be unbiased under SRS (Särndal et al. 2003, p. 396) but only approximately unbiased for other equal intensity designs (like the FIA design) and then only where the expected value of $n_d$ is large. Thus, we turn to MA estimators that use data from beyond the domain of interest to support inferences about the $\mu_{y,d}$. To do so, define the synthetic domain estimator as

$$(7) \qquad \mu_{\tilde{y},d}(s) = \frac{1}{A_d} \int_{\mathcal{D}_d} \tilde{y}(\mathbf{w}, t_d, s) \, \mathrm{d}\mathbf{w}$$

where $\tilde{y}(\mathbf{w}, t_d, s)$ is a KNN-estimated surface spanning $\mathcal{D}_d$. This surface is estimated in accordance with eqs. 1 and 2, using the selected $K$ and predictor vector. Though continuous, in practice this surface is pixel-wise constant since all predictor variables share the same value within a given pixel. As a mean of KNN estimates spanning $\mathcal{D}_d$, $\mu_{\tilde{y},d}(s)$ is a function of the plot-level values of stocking, including some potentially lying outside of $\mathcal{D}_d$. As it does not depend in any way on $n_d$, it

may be quite precise but is generally design-biased. However, it is possible to estimate its bias from the domain sample $s_d$ as

(8) $\quad \bar{e}_d = \dfrac{1}{n_d} \sum_{i \in s_d} e_i(s) = \dfrac{1}{n_d} \sum_{i \in s_d} [y_i - \tilde{y}_i(s)]$

where the $e_i(s)$ correspond to discrepancies between observed stocking and the synthetic KNN surface, with $\tilde{y}_i(s) = \tilde{y}(\mathbf{w}_i, t_i, s)$ obtained by intersecting the footprint of plot $i$ with this surface.

The bias correction term (8) suggests augmenting the synthetic estimator (7) to obtain

(9) $\quad \widehat{\mu}_{y,d} = \mu_{\tilde{y},d}(s) + \bar{e}_d$

and below we distinguish between two cases. In the first, unconstrained KNN is applied using the complete set of FIA plot measurements ($s_a$). Thus, the selection of plot locations within the domain determines the point set over which the bias correction term is calculated and affects the probabilistic properties of the KNN surface for $\mathcal{D}_d$. In this case, with outcomes internal to the domain of interest being available for KNN estimation, we refer to the two-part estimator (9) as an unconstrained domain-level MA estimator. The second case obtains where external KNN is applied using only measurements external to the domain of interest ($s_{-d}$). This provides a bias correction term where the $\tilde{y}_i(s_{-d})$ are independent of $s_d$, leading to a form of (9) that we refer to as an external domain MA estimator. As developed in the Appendix, this external MA estimation strategy is motivated by its unbiasedness for $\mu_{y,d}$ under SRS. We note, however, that like direct estimator (6) unbiasedness does not necessarily carry over to spatially structured sampling designs (such as the FIA design) owing to the randomness of $n_d$.

Variance expressions for MA estimators of the form $\widehat{\mu}_{y,d}$ generally rely on approximations, improving for large $n_d$, that reduce to functions of the expected values of the discrepancies $e_i$ (see e.g., Baffetta et al. 2009). However, we show in the Appendix that under SRS the variance of the external MA estimator of $\widehat{\mu}_{y,d}$ can be written as

(10) $\quad \mathsf{V}\left[\widehat{\mu}_{y,d}\right] = \mathsf{E}\left[\mathsf{V}\left[\bar{e}_d \mid n_d, s_{-d}\right]\right]$

provided only that $n_d > 0$. That is, the variance of the external MA estimator is the conditional variance of the bias correction term, averaged over possible domain sample sizes ($n_d$) and sample selections external to the domain of interest ($s_{-d}$). This then suggests the simple variance estimator

(11) $\quad \widehat{\mathsf{V}}\left[\widehat{\mu}_{y,d}\right] = \dfrac{1}{n_d(n_d - 1)} \sum_{i \in s_d} [e_i - \bar{e}_d]^2$

with $e_i = e_i(s_{-d})$ and $\bar{e}_d = \frac{1}{n_d} \sum_i e_i$. Under SRS, this unbiasedly estimates the conditional variance of the external MA estimator (see Appendix). The same result does not obtain for the unconstrained MA estimator owing to dependence of the

$\tilde{y}_i(s_a)$ on $s_d$. Nonetheless, this variance estimator has been applied in this context (Baffetta et al. 2009; Bell et al. 2022) and is used for the unconstrained MA estimator below (with $e_i = e_i(s_a)$).

The variance of the direct domain estimator (6) can be written as

(12) $\quad \mathsf{V}[\bar{y}_d] = \mathsf{E}\left[\dfrac{\sigma_{y,d}^2}{n_d}\right]$

for $n_d > 0$, where $\sigma_{y,d}^2$ is the variance of $y$ within domain $\mathcal{D}_d$. Based on eq. 12, the SE of $\bar{y}_d$ can be estimated from

$$\widehat{\mathsf{V}}[\bar{y}_d] = \dfrac{1}{n_d(n_d - 1)} \sum_{i \in s_d} (y_i - \bar{y}_d)^2 = \dfrac{1}{n_d} \widehat{\sigma}_{y,d}^2$$

provided $n_d > 1$. Given this form for the estimated variance of $\bar{y}_d$, we can interpret the efficiency of an alternative estimator of the domain mean in terms of a proportionate increase or decrease to the domain sample size used for direct estimation. Specifically, if an alternative estimator has SE equal to $b \times \frac{\widehat{\sigma}_{y,d}}{\sqrt{n_d}} = \frac{\widehat{\sigma}_{y,d}}{\sqrt{b^{-2} n_d}}$, then its use can be interpreted as equivalent to a $b^{-2}$ factor increase or decrease in $n_d$.

Finally, we note that Särndal et al. (2003, p. 224) establish conditions under which the variance of a difference estimator of a population mean will be smaller than that of an expansion estimator. Under SRS, this results when the correlation $r_{\hat{y},y}$ between MA estimates and actual values of $y$ exceeds half the ratio of their standard deviations—that is, when $r_{\hat{y},y} > 0.5 \frac{\sigma_{\hat{y}}}{\sigma_y}$. Below, we evaluate this inequality using sample-based estimates of the correlation and standard deviations.
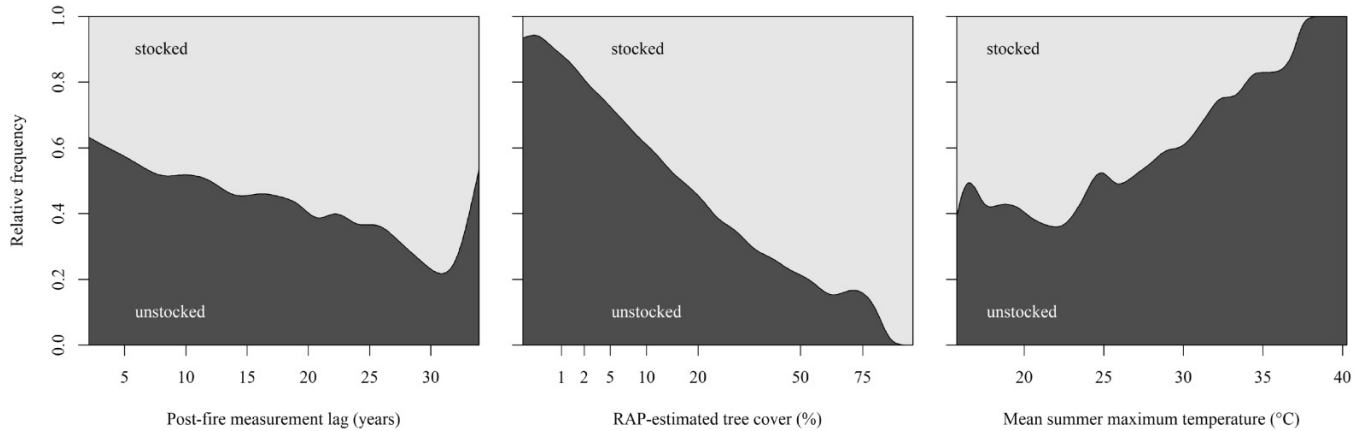
## Results

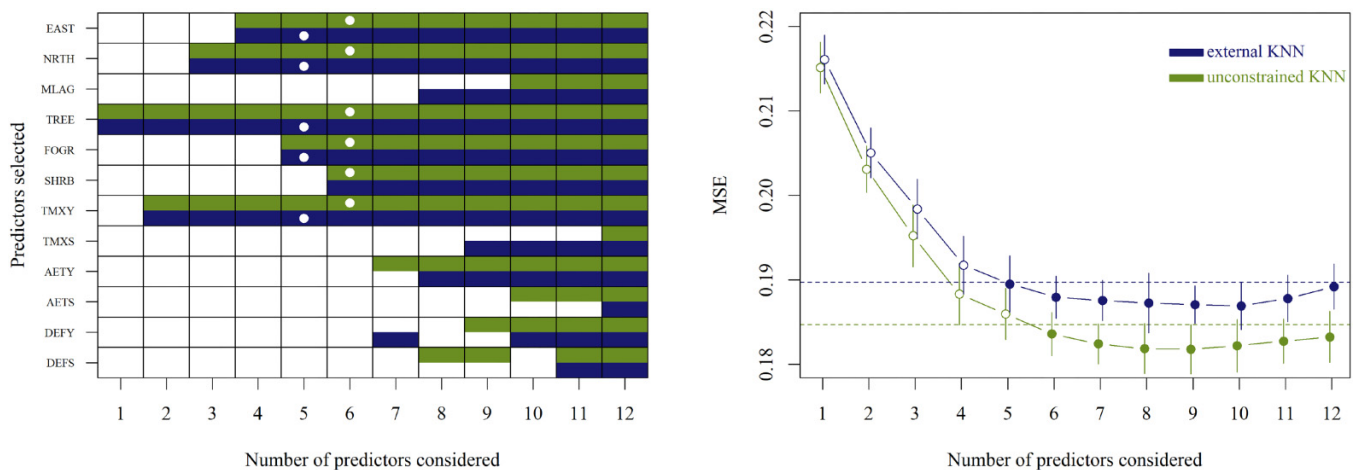### Plot-level KNN estimation

Extreme variability was observed in densities of seedlings, saplings, and mature trees across the 5660 postfire FIA measurements, but the binary stocking variable used here averaged 0.49 (or 49%) and varied over a relatively narrow range across US states (0.46–0.61). Nevada, where mean stocking was only 0.12, provided an exception. With respect to the auxiliary variables of Table 1, key patterns included increasing stocking rates with time-since-fire (MLAG) and RAP-estimated tree cover (TREE), as well as decreased stocking with increased mean annual temperature (TMXY; Fig. 2). For geographic reference, measurements with the highest values of TMXY were found in southern New Mexico and Arizona, and in the central valley of California.

All 12 predictors in Table 1 were made available for KNN estimation, but lowest MSEs were achieved using only 9 (unconstrained KNN) or 10 (external KNN; Fig. 3). Furthermore, MSEs within one cross-validated SE of the minimum could be achieved using only five (external KNN) or six (unconstrained KNN) predictors. The left panel of Fig. 3 shows the best com-

**Fig. 2.** Conditional density of stocking relative to time-since-fire (MLAG; left); estimated tree cover on square-root scale (TREE; center); and mean summer maximum temperature (TMXS; right).



**Fig. 3.** Best predictor sets (left) and corresponding MSE (±1 SE; right) for KNN estimation as a function of the number of predictors considered; results for unconstrained KNN are shown in green, those for external KNN in blue. In the left panel, the selected predictors (and *p*) are identified with white dots; in the right panel, MSEs within one SE of the lowest MSE (dashed lines) are drawn with filled symbols.
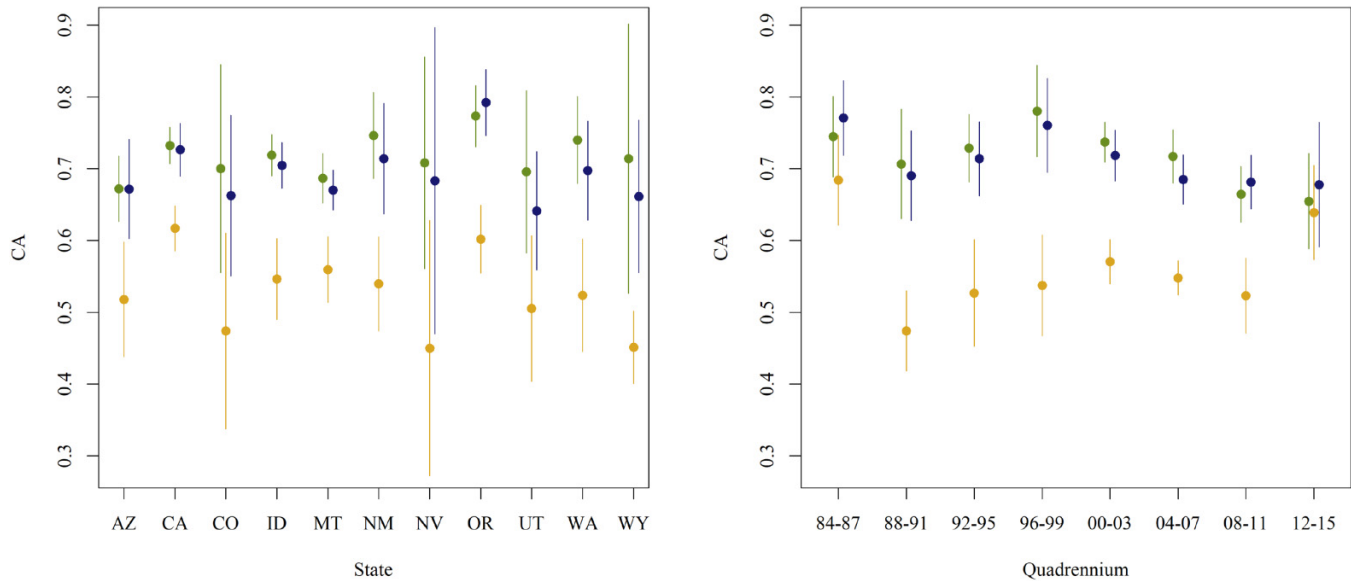


binations of predictors for each KNN strategy and number of predictors considered. For both KNN strategies, the best single predictor variable was TREE and the best combination of two predictors was TREE and TMXY. However, when considering only combinations with MSE within one SE of the minimum, the geographic coordinates (EAST, NRTH) were consistently among the selected predictor set, as was TREE and one of the maximum temperature variables (TMXY or TMXS). Figure 3 does not show the values of *K* minimizing the MSEs of the various KNN strategies, but these ranged from 18 to 45 with larger values generally associated with smaller number of predictors. For both the selected $p = 5$ external KNN and $p = 6$ unconstrained KNN strategies, the smallest value of *K* maintaining an MSE within one SE of the minimum was 29. However, in all cases the best unconstrained KNN strategies achieved lower MSE than the best external KNN strategies (Fig. 3 right panel). For reference, the MSE associated with the use of the simple domain means (6) was 0.280 (±0.003 SE), substantially larger than any of those for KNN.

Figure 4 shows the classification accuracies of the selected KNN strategies and of the direct domain estimators. Classification accuracies of the two KNN strategies consistently fell within ±2 SEs of one another, but both were generally much higher than those of the direct domain estimators. The exception to this was for measurements falling in the earliest (1984–1987) and most recent (2012–2015) burn periods, which were also associated with the longest and shortest post-fire measurement lags, respectively.

Simple linear correlations between stocking and KNN-estimated stocking were 0.51 and 0.49 for the unconstrained and external strategies, respectively. At the same time, the ratios of standard deviations of estimated-to-observed stocking ($\hat{\sigma}_{\hat{y}}/\hat{\sigma}_y$) were approximately 0.5 in both cases. Thus, $\hat{r}_{y,\hat{y}} \gg 0.5\hat{\sigma}_{\hat{y}}/\hat{\sigma}_y$ for both unconstrained and external KNN, suggesting that difference estimation based on KNN-derived surfaces would yield improved estimates of stocking at the population level.

**Fig. 4.** Classification accuracies (±2 SEs) of selected KNN strategies across US states (left) and burn periods (right). Results for unconstrained KNN are shown in green, for external KNN in blue, and for direct estimation in orange.



## Domain-level estimation

Despite differences in MSEs and predictor variables, unconstrained and external MA estimates of domain means were highly correlated (Pearson correlation 0.96), as were their estimated SEs (correlation 0.93). As such, below we focus primarily on external MA domain estimates. These were linearly related to the direct domain estimates (Fig. 5; correlation 0.90), with greatest differentiation at the extremes ($\bar{y}_d = 0$ or 1). The latter resulted when domain plot measurements were all 0 or all 1, and occurred primarily when domain sample sizes were small ($n_d = 1$ or 2). In such instances, the external MA estimator frequently resulted in domain estimates outside the unit interval.

The distribution of bias correction terms (8) for the external MA domain estimators are shown in Fig. 6. These bias correction terms are highly variable for small domain sample sizes (esp. for $n_d \leq 5$). However, as $n_d$ increases beyond approximately 10 measurements, the distribution shows diminishing variability around a value close to 0. This suggests that the external KNN synthetic estimator is approximately unbiased at the domain level, despite drawing on $K = 29$ neighbors across the five-predictor design space. Examining these bias correction terms for domains within states (not shown), the distributions were again approximately symmetric around 0. The exception was for burned-area domains in UT, where the distribution was shifted to positive values, indicating that the selected KNN strategy underestimated stocking in these domains. These results were also observed for MA domain estimators based on unconstrained KNN (not shown).

Estimated SEs of the MA and direct domain estimates are plotted in Fig. 7. Differences were slight for domains with samples sizes larger than approximately 20, but are more apparent for more sparsely sampled domains. Fitting a weighted (by $n_d$) linear regression yielded a slope of 0.84 (SE 0.023), indicating the estimated SEs of the external MA

estimator were on average 16% smaller than those of the direct estimates. Alternatively, this apparent increase in precision of the external MA estimates over the direct estimates is equivalent to a 42% (i.e., $b^{-2} = 0.84^{-2} = 1.42$) increase in domain sample size.
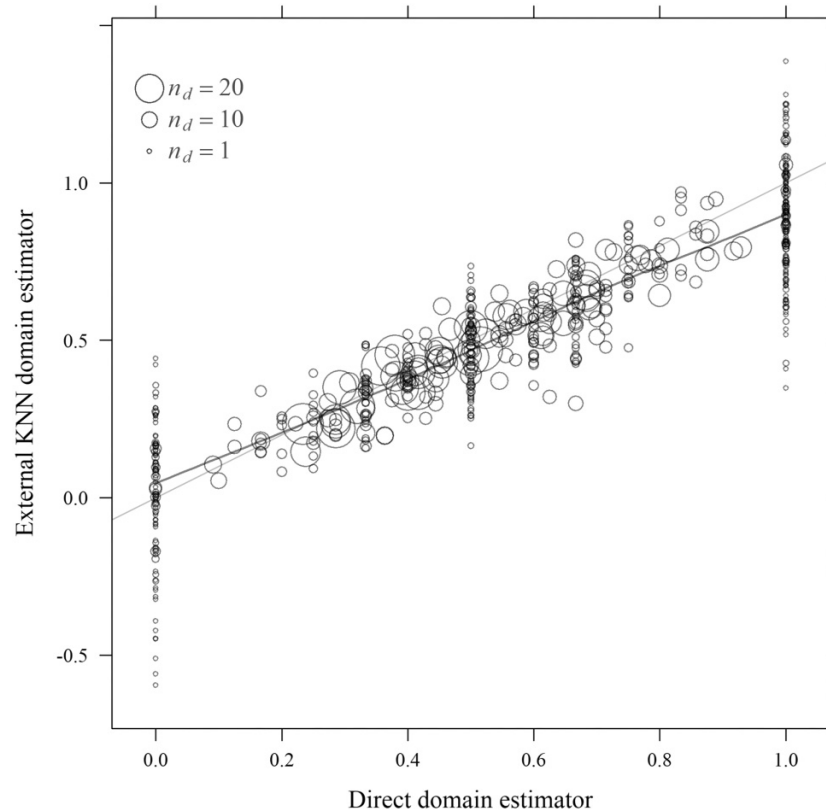
A subset of domain-level estimates are shown in Fig. 8. Specifically, this figure shows estimates of mean stocking 3–15 years after fire on NFS lands that burned between 2000 and 2003 in MT. Direct estimates vary considerably year-over-year, owing in part to the fact that measurements from distinct locations are used. That is, the 17 sample locations used to estimate stocking 12 years postfire are distinct from the 17 locations used to estimate stocking 13 years postfire. The synthetic estimator based on external KNN is much smoother and suggests increased stocking over time. The external MA domain estimates build on the smooth synthetic estimates but take on additional variability owing to the bias correction terms. The ±2 SE intervals drawn around the domain estimates generally widen as sample size decreases. Intervals for the external MA estimates are not uniformly narrower than for the direct estimates, but are narrower on average and more so for smaller $n_d$.

## Discussion

Across the FIA measurement data, KNN yielded favorable estimation results from intuitively defined neighborhoods. Both the unconstrained and external KNN strategies achieved MSEs that were appreciably lower than using direct domain estimates, and classification accuracies that were appreciably higher (Fig. 4). The neighborhood dimensions also spanned a rational basis for estimation (Fig. 3 left), defining nearest measurements to be close geographically (EAST, NRTH), radiometrically (TREE, and FOGR or SHRB), and climatically (TMXY or TMXS, and/or AETY). Put differently, a geographically proxi-

**Fig. 5.** Model assisted (MA) versus direct estimates of domain means with symbol scales according to domain sample size; MA estimates are based on the selected external KNN strategy.
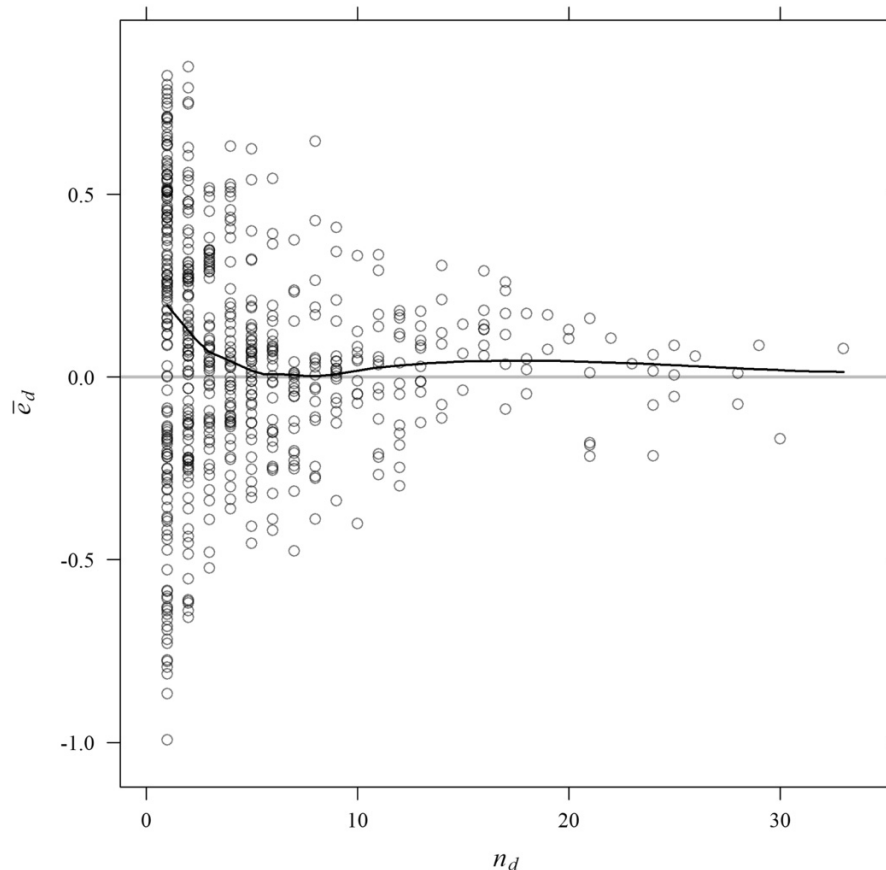


mate measurement may not rank highly as a neighbor if it has a very different inferred tree cover (e.g., perhaps due to a different time-since-fire) or a very different climatic regime (e.g., owing to complex terrain). These outcomes notwithstanding, no claim is made that these are the best combinations of predictors for estimating postfire stocking. Indeed, our use and selection of geographic position variables is perhaps an indication of the importance of spatially varying processes that could be described by other attributes. In particular, the harmonic predictors derived from Landsat time series and used by Bell et al. (2022) could be useful in this postdisturbance context. However, the determination of a best set of predictors was not a central focus of this study, and our methods could readily be extended to other predictor sets.

With the number of potential predictors limited to 12 (Table 1), a full search over the space of predictor combinations and a wide range of $K$ was feasible. McRoberts (2012) emphasized the importance of predictor variable selection in KNN, owing in part to the fact that MSE can increase as the number of predictors (and the dimensionality of the neighborhood) increases. This was observed in our analyses (Fig. 3 right), where cross-validated MSEs were lowest when the numbers of predictors were 9 or 10 (depending on the KNN strategy). Using cross-validation, our selected predictors sets were reduced further down to five or six (Fig. 3 left). This process also resulted in selection of $K = 29$ neighbors for both KNN strategies. Such a large value for $K$ can be expected to yield precise but biased estimates: precise owing to the

large number of neighbors being averaged over and biased owing to the larger distances/conditions from which neighbors are drawn (Hastie et al., 2009, §13.3). Using a canonical correlation-based distance metric, Bell et al. (2022) arrived at a value of $K$ of similar magnitude ($K = 28$), though their application involved estimation of aggregate tree biomass by species. Given our focus on a single forest attribute, we did not consider using canonical correlation for neighborhood determination. Also, for computational reasons, we opted for a variance-weighted Euclidean distance function (2) (similar to a Mahalanobis distance but with diagonal variance matrix) and a simple unweighted average of the KNNs. McRoberts (2012) reviews and evaluates other KNN distance functions and neighbor-weighting strategies.

Modifying the KNN procedure to draw only from an external sample proved computationally (and algorithmically) simple. The idea is similar to cross-validation in that data from each domain of interest are withheld and an assisting model (KNN-based) is then developed from the remaining data. Unlike cross-validation, however, there is no subsequent aggregation of results across withheld partitions because each domain is an entity of interest.

The external KNN procedure resulted in neighborhoods of reduced dimension and slightly increased MSE (Fig. 3). The latter was not surprising given that observations from within the same domain as a focal measurement can generally be expected to have experienced the most similar disturbance conditions (and pre- and postdisturbance conditions). However,

**Fig. 6.** Bias correction terms (8) used in the external KNN MA estimator versus domain sample size.
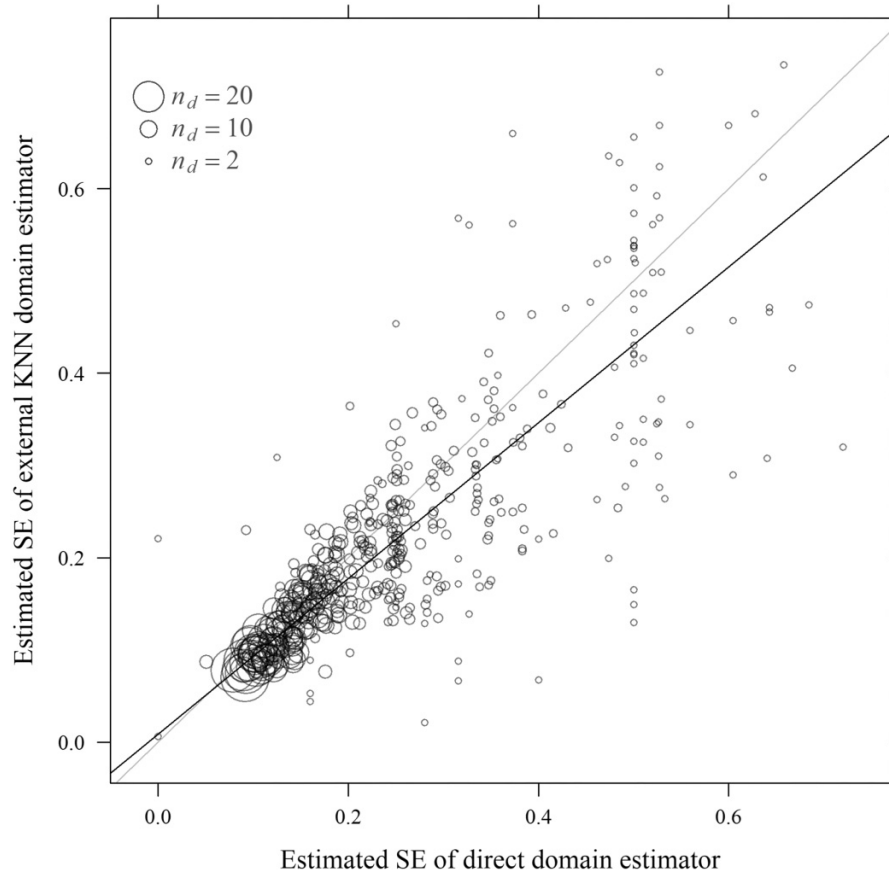


with typically small domain samples from which to draw (Table 2), the unconstrained KNN strategy was nearly always forced to draw on data from beyond the domain of the focal measurement. Thus, the impacts of the neighborhood constraint on the nature and performance of the KNN strategies was relatively small (Fig. 4).

In this study, domains were tessellated with approximately 30 m × 30 m pixels, derived from projection of the RAP rasters (themselves based on Landsat sensor resolution; Allred et al. 2021). In contrast, an FIA plot is a dispersed cluster of four circular subplots (each with nested microplots) covering a combined area of 673 m$^2$. For KNN model development and subsequently for KNN surface bias estimation, we assigned values to plot measurements through spatial overlays of burned subplots and weighted averages of intersected pixel values. Other studies have approached the problem differently. For example, McRoberts (2012) and Bell et al. (2022) used only central FIA subplot observations and indexed these against spatial layers using only that subplot's coordinate location. Yet short of squaring plot dimensions with pixel dimensions—infeasible in practice—there is no clear method of reconciling differences in spatial scales. A consequence of this is that we cannot conclude that the KNN strategy (predictor set and K) selected at the plot-scale exhibits similar performance when applied in domain estimation. However, at the domain estimation stage we use KNN only for synthetic esti-

mation, and differences in scale do not impact our ability to evaluate the performance of external MA domain estimates via the bias correction term (8).

The spatial and temporal structure of the FIA design also affects the properties of our domain estimators. In particular, owing to the spatial dispersion of plot locations across hexagonal cells and to the temporal dispersion of measurements over interpenetrating panels (Bechtold and Patterson 2005), the joint inclusion density (sensu Cordy 1993) of any two points on the landscape is not strictly positive. Ignoring this and applying variance estimators derived for SRS likely leads to variance overestimation (Gregoire and Valentine 2007, p. 55) even as the dispersion could be quite advantageous from a variance reduction standpoint—it precludes measurement of spatiotemporally proximate locations where the resources of interest are likely more similar (see Stevens 1997). The same design structures also barred us from generalizing the unbiasedness of our external MA estimator. Specifically, our approach of conditioning on the observed domain sample size breaks down with the FIA design because juxtaposition of hexagonal cells and domain boundaries can then be informative with regards to conditional sampling intensity (see Appendix). Simulated sampling can be used to evaluating the magnitude of bias (if any) in our external MA estimator when applied in a stratified design like that used by the FIA and is the subject of an ongoing research. We further

**Fig. 7.** Estimated SEs of the direct and external KNN MA domain estimates with symbols scaled according to domain sample size; trend (black line) is the $n_d$-weighted linear regression. Only shown are results for domains with $n_d > 1$.



conjecture that modifying the bias correction term (8) to a Horvitz–Thompson form (i.e., $A_d^{-1}\sum_i \frac{e_i(s_{-d})}{\pi_i}$ with $\pi_i$ being the inclusion density) would yield an unbiased estimator. As discussed by Särndal et al. (2003, pp. 391–393), however, this can be expected to increase estimator variance and thus we did not pursue the approach here. Finally, the FIA design elements also impacted how we implemented our external KNN approach. Specifically, the external sample $s_d$ was made up of measurements from outside the spatial footprint of the domain of interest $d$ and did not include measurements lying within that footprint but obtained outside the temporal scope of the domain. For example, a 2005 measurement lying within an area burned in 2000 could not contribute to a direct estimate of 2010 stocking, but neither was it allowed to be a neighbor for external KNN estimation of 2010 stocking. This was done because the inclusion of such a measurement in the sample—$n$ for an off-year—is informative about the conditional inclusion density of measurements in the year of interest.
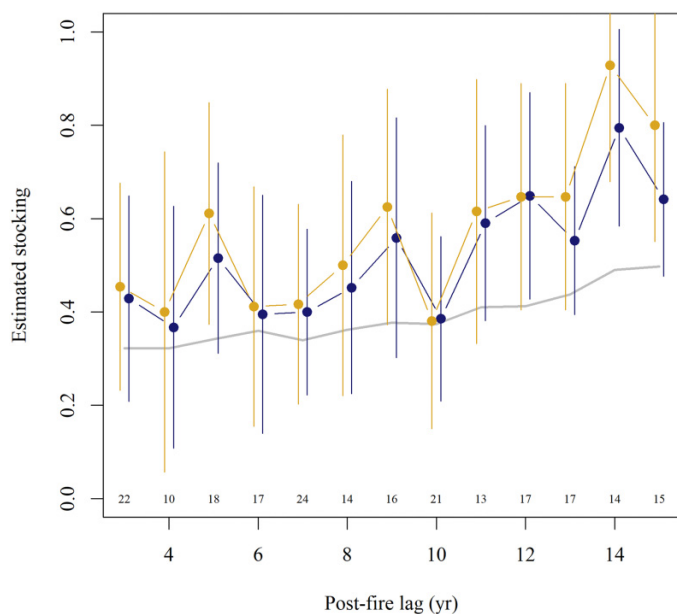
Aggregating KNN estimates to the domain level, the unconstrained and external MA estimates were very similar. As noted above, this is likely attributable to the fact that both KNN strategies drew heavily on extra-domain measurements ($K = 29$ neighbors while the median domain sample

size was only 3; Table 2). Estimated SEs were slightly smaller for the unconstrained MA estimator, though it also carried a design bias. That bias is likely small but is difficult to estimate (Gaines and Affleck 2021). Indeed, that difficulty is what motivated the development and evaluation of the external KNN approach. Under SRS and if $n_d > 0$, it provides unbiased estimates of domain means and a variance estimator that does not rely on linear approximation (see Appendix).

External MA estimates were similar also to the direct estimates for many of the domains studied here. The exception was where the domain sample size dropped (toward 1), the direct estimate took a value of 0 or 1, and the observed variability in stocking collapsed to 0 (Fig. 5). Where that occurred, the external MA estimator sometimes departed substantially from the direct estimate and could take values less than 0 or greater than 1. The latter is not a result of the external KNN strategy—it occurred also with the unconstrained KNN domain estimates and, in general, will be possible for any MA estimator employing an additive bias correction term (9).

The external MA estimator did yield apparent reductions in SEs of 16% over direct estimates (Fig. 7). As noted, this can be interpreted as equivalent to a 42% increase in $n_d$. Nonetheless, as illustrated in Fig. 8, substantial uncertainty remains. This is in part owing to high variability in the bias correction el-

**Fig. 8.** Estimated stocking in domains defined by different postfire assessment lags over NFS lands that burned between 2000 and 2003 in MT. Synthetic estimates from external KNN estimates shown in gray, MA estimates in blue, and direct estimates in gold. Intervals span $\pm 2$ SE; domain sample sizes ($n_d$) are printed above the $x$-axis.



ement for small $n_d$, particularly $n_d < 10$ (Fig. 6). As suggested by Bell et al. (2022), it may be beneficial to rely entirely on synthetic estimation where $n_d$ falls below a certain threshold. Doing so incurs a design bias of a magnitude difficult to estimate, but which is likely smaller than the gains in precision obtained by dropping $\bar{e}_d$. In fact, that the mean bias correction tends to 0 as $n_d$ increases (Fig. 6) suggests that the bias of the external KNN synthetic estimator is low. Of course, relying entirely on synthetic estimation is also necessary where $n_d = 0$. Bell et al. (2022) and McRoberts (2012) describe a model-based inferential approach for KNN domain estimation and evaluate methods for estimating (model-based) variance. The latter can be computationally demanding but likely feasible for domains of the size considered here (Table 2).

For design-based inference, the KNN synthetic estimator contributes to the uncertainty of the external MA estimator (9) in a quite different way. The synthetic estimator effectively establishes a new baseline (departing from $\bar{y}_d$) around which variability is reckoned. Here, as in difference estimation, we have conditioned on the synthetic estimation surface and focused entirely on the within-domain dispersion around it (cf. eqs. (10) and (11)). The cross-validation techniques for MA variance estimation discussed by McConville et al. (2020) could potentially capture variation induced by the inherent randomness of $s_{-d}$. However, with the external MA domain estimator it is reasonable to condition on the synthetic estimator because the factors affecting its variation are largely distinct from (external to) the design-induced sampling variability within the domain of interest.

Finally, the external KNN approach developed here can be extended to other forms of synthetic estimation. In particular, it can be applied with regression to provide an external generalized regression estimator (Gaines 2022). Like the approach developed by Mandallaz (2013), this would necessitate estimation of distinct regression coefficients for each domain of interest. However, it would depart from their strategy by excluding the data from the domain of interest (thus ensuring externality) rather than by allowing for domain-specific coefficients. Also, while Baffetta et al. (2009) adopt the term "empirical difference estimator" for what is equivalent to our unconstrained MA estimator, we would argue that this term is better applied to the external MA estimator. This is because it is only in the latter case that the proxy function is derived from sources outside the domain of interest. Of course, when considering simultaneous inference for multiple domains there may be a need to account for cross-domain dependence in synthetic estimation surfaces. However, this is not central to our interest as the method was developed primarily for unbiased estimation of stocking within specific (and possibly singular) domains. Whereas many model-based approaches to SAE require a collection of domains to be defined in order to estimate the mean of any single domain, the approaches studied here are applicable where there exists only a single domain of interest, provided a sample of larger scope exists.

## Acknowledgements

## Article information

### History dates

### Copyright

### Data availability
MTBS wildfire perimeter vectors are available for download at https://mtbs.gov/. FIA tabular data are available for download at https://apps.fs.usda.gov/fia/datamart/datamart.html. Contact USFS FIA for information regarding FIA sample location coordinates.

## Author information

### Author ORCIDs
George C. Gaines, III https://orcid.org/0000-0002-4341-6802

### Author contributions
Conceptualization: DA, GG
Data curation: DA, GG
Formal analysis: DA, GG
Funding acquisition: DA
Investigation: DA, GG
Methodology: DA, GG
Project administration: DA, GG
Resources: DA
Supervision: DA
Validation: DA, GG
Visualization: DA, GG
Writing – original draft: DA, GG
Writing – review & editing: DA, GG

### Competing interests
There are no competing interests that bias or might be seen to bias this work.

## References

Abatzoglou, J.T., Dobrowski, S.Z., Parks, S.A., and Hegewisch, K.C. 2018. Terraclimate, a high-resolution global dataset of monthly climate and climatic water balance from 1958–2015. Sci. Data. **5**: 170191. doi:10.1038/sdata.2017.191.

Allred, B.W., Bestelmeyer, B.T., Boyd, C.S., Brown, C., Davies, K.W., Duniway, M.C., et al. 2021. Improving Landsat predictions of rangeland fractional cover with multitask learning and uncertainty. Meth. Ecol. Evol. **12**: 841–849. doi:10.1111/2041-210X.13564.

Baffetta, F., Fattorini, L., Franceschi, S., and Corona, P. 2009. Design-based approach to k-nearest neighbours technique for coupling field and remotely sensed data in forest surveys. Remote Sens. Environ. **113**: 463–475. doi:10.1016/j.rse.2008.06.014.

Bechtold, W.A., and Patterson, P.L. 2005. The enhanced forest inventory and analysis program-national sampling design and estimation procedures. Gen. Tech. Rep. SRS-80, USDA Forest Service, Southern Research Station, Asheville, NC.

Bell, D.M., Wilson, B.T., Werstak, C.E., Jr., Oswalt, C.M., and Perry, C.H. 2022. Examining k-nearest neighbor small area estimation across scales using national forest inventory data. Front. Forests Glob. Change **5**: 763422. doi:10.3389/ffgc.2022.763422.

Breidenbach, J., and Astrup, R. 2012. Small area estimation of forest attributes in the Norwegian national forest inventory. Eur. J. Forest Res. **131**: 1255–1267. doi:10.1007/s10342-012-0596-7.

Breidt, F.J., and Opsomer, J.D. 2017. Model-assisted survey estimation with modern prediction techniques. Stat. Sci. **32**: 190–205. doi:10.1214/16-STS589.

Burrill, E.A., DiTommaso, A.M., Turner, J.A., Pugh, S.A., Christiansen, G., Perry, C.J., and Conkling, B.L. 2021. The forest inventory and analysis database: database description and user guide version 9.0 for phase 2. USDA forest service manual. Available from https://www.fia.fs.usda.gov/library/database-documentation/current/ver90/FIADB%20User%20Guide%20P2-_9-0-_final.pdf [accessed February 2022].

Cordy, C.B., 1993. An extension of the Horvitz-Thompson theorem to point sampling from a continuous universe. Stat. Prob. Lett. **18**: 353–362. doi:10.1016/0167-7152(93)90028-H.

Coulston, J.W., Green, P.C., Radtke, P.J., Prisley, S.P., Brooks, E.B., Thomas, V.A., et al. 2021. Enhancing the precision of broad-scale forestland removals estimates with small area estimation techniques. Forestry: An Int. J. Forest Res. **94**: 427–441. doi:10.1093/forestry/cpaa045.

Eidenshink, J., Schwind, B., Brewer, K., Zhu, Z., Quayle, B., and Howard, S. 2007. A project for monitoring trends in burn severity. Fire Ecol. **3**: 3–21. doi:10.4996/fireecology.0301003.

Gaines, G.C., III, 2022. Small area estimation of postfire tree density in the western united states using an annualized forest inventory. Ph.D. dissertation, University of Montana, Missoula, MT.

Gaines, G.C., III, and Affleck, D.L.R., 2021. Small area estimation of postfire tree density using continuous forest inventory data. Front. Forests Glob. Change **4**: 761509. doi:10.3389/ffgc.2021.761509.

González, M.E., and Waksberg, J. 1973. Estimation of the error of synthetic estimates. Paper presented at First Meeting of the International Association of Survey Statisticians, Vienna, Austria, 18-25 August, 1973.

Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., and Moore, R. 2017. Google Earth Engine: planetary-scale geospatial analysis for everyone. Remote Sen. Environ. **202**: 18–27. doi:10.1016/j.rse.2017.06.031.

Gregoire, T.G., and Valentine, H.T. 2007. Sampling strategies for natural resources and the environment. Chapman and Hall/CRC, Boca Raton, FL. doi:10.1201/9780203498880.

Haakana, H., Heikkinen, J., Katila, M., and Kangas, A. 2020. Precision of exogenous post-stratification in small-area estimation based on a continuous national forest inventory. Can. J. Forest Res. **50**: 359–370. doi:10.1139/cjfr-2019-0139.

Hastie, T., Tibshirani, R., and Friedman, J. 2009. The elements of statistical learning: Data mining, inference, and prediction. Springer Science & Business Media, Berlin.

Higuera, P.E., Shuman, B.N., and Wolf, K.D. 2021. Rocky Mountain subalpine forests now burning more than any time in recent millennia. Proc. Natl. Acad. Sci. U. S. A. **118**: 1–5. doi:10.1073/pnas.2103135118.

Hill, A., Mandallaz, D., and Langshausen, J. 2018. A double-sampling extension of the German national forest inventory for design-based small area estimation on forest district levels. Remote Sen. **10**: 1052. doi:10.3390/rs10071052.

Magnussen, S., Næsset, E., and Gobakken, T. 2014. An estimator of variance for two-stage ratio regression estimators. Forest Sci. **60**: 663–676. doi:10.5849/forsci.12-163.

Mandallaz, D. 2012. Design-based properties of small-area estimators in forest inventory with two-phase sampling. ETH Zurich, Department of Environmental Systems Science, Technical Report. Available from e-collection.library.ethz.ch [accessed June 2021].

Mandallaz, D. 2013. Design-based properties of some small-area estimators in forest inventory with two-phase sampling. Can. J. Forest Res. **43**: 441–449. doi:10.1139/cjfr-2012-0381.

Mandallaz, D., Breschan, J., and Hill, A. 2013. New regression estimators in forest inventories with two-phase sampling and partially exhaustive information: a design-based Monte Carlo approach with applications to small-area estimation. Can. J. Forest Res. **43**: 1023–1031. doi:10.1139/cjfr-2013-0181.

Marker, D.A. 1995. Small area estimation: a Bayesian perspective. Ph.D. dissertation, University of Michigan, Ann Arbor, MI.

McConville, K.S., Moisen, G.G., and Frescino, T.S. 2020. A tutorial on model-assisted estimation with application to forest inventory. Forests, **11**: 244. doi:10.3390/f11020244.

McRoberts, R.E. 2012. Estimating forest attribute parameters for small areas using nearest neighbors techniques. Forest Ecol. Manag. **272**: 3–12. doi:10.1016/j.foreco.2011.06.039.

McRoberts, R.E., Næsset, E., Heikkinen, J., Chen, Q., Strimbu, V., Esteban, J., et al. 2022. On the model-assisted regression estimators using remotely sensed auxiliary data. Remote Sen. Environ. **281**: 113168. doi:10.1016/j.rse.2022.113168.

MTBS. 2023. Monitoring Trends in Burn Severity (MTBS) burned areas boundaries for 1984-2021. MTBS Project, USDA Forest Service and US Geological Survey. doi:10.5066/P9IED7RZ.

Parks, S.A., and Abatzoglou, J.T. 2020. Warmer and drier fire seasons contribute to increases in area burned at high severity in western US forests from 1985 to 2017. Geophys. Res. Lett. **47**: e2020GL089858. doi:10.1029/2020GL089858.

Perry, C.H., Finco, M.V., and Wilson, B.T. 2022. Forest Atlas of the United States. FS-1172. USDA Forest Service, Washington, DC.

Rao, J.N.K., and Molina, I. 2015. Small area estimation, 2nd ed. Wiley, New York.

Särndal, C.E., Swensson, B., and Wretman, J. 2003. Model assisted survey sampling. Springer Science & Business Media, Berlin.

Stevens, D.L., Jr. 1997. Variable density grid-based sampling designs for continuous spatial populations. Environmetrics, **8**: 167–195. doi:10.1002/(SICI)1099-095X(199705)8:3⟨167::AID-ENV239⟩3.0.CO;2-D.

Stevens-Rumann, C.S., Kemp, K.B., Higuera, P.E., Harvey, B.J., Rother, M.T., Donato, D.C., et al. 2017. Evidence for declining forest resilience to wildfires under climate change. Ecol. Lett. **21**: 243–252.

US Census Bureau. 2023. Cartographic boundary files—States. Census Bureau, US Department of Commerce. Available from https://www2.census.gov/geo/tiger/GENZ2018/shp/cb_2018_us_state_500k.zip [accessed January 2023].

USFS. 2023. National Forest System land units. Geospatial Technology and Applications Center, USDA Forest Service (USFS). Available from https://data.fs.usda.gov/geodata/edw/edw_resources/shp/S_USA.NFSLandUnit.zip/ [accessed January 2023].

Ver Planck, N.R., Finley, A.O., Kershaw, J.A., Jr, Weiskittel, A.R., and Kress, M.C. 2018. Hierarchical Bayesian models for small area estimation of forest variables using lidar. Remote Sen. Environ. **204**: 287–295. doi:10.1016/j.rse.2017.10.024.

Young, D.J.N., Werner, C.M., Welch, K.R., Young, T.P., Safford, H.D., and Latimer, A.M. 2019. Post-fire forest regeneration shows limited climate tracking and potential for drought-induced type conversion. Ecology, **100**: e02571. doi:10.1002/ecy.2571.

# Appendix A

Results concerning the mean and variance of the MA domain estimator given by (9) can be derived under SRS, but consider first the wider class of equal probability sampling designs and the condition $n_d > 0$. The latter is a precondition for the application of the MA domain estimator (9) under any form of KNN (or regression) estimation. Under these conditions

$$(A1) \quad E\left(\bar{e}_d\right) = E\left[E\left(\bar{e}_d \mid n_d\right)\right],$$

$$= E\left[E\left(\frac{1}{n_d}\sum_{i=1}^{n_d} y_i \mid n_d\right)\right] - E\left[E\left(\frac{1}{n_d}\sum_{i=1}^{n_d}\tilde{y}_i\left(s\right) \mid n_d\right)\right]$$

$$= E\left[\frac{1}{n_d}\sum_{i=1}^{n_d} E\left(y_i \mid n_d\right)\right] - E\left[\frac{1}{n_d}\sum_{i=1}^{n_d} E\left(\tilde{y}_i\left(s\right) \mid n_d\right)\right]$$

Suppressing the temporal dependence of $y_i$ and $\tilde{y}_i(s)$ for ease of exposition, the first of the inner expectations of (A1) can be written as

$$(A2) \quad E\left(y_i \mid n_d\right) = E\left(y\left(\mathbf{w}_i\right) \mid n_d\right)$$

$$= \int_{\mathcal{D}_d} y\left(\mathbf{w}\right) f\left(\mathbf{w}|n_d\right) \, \mathrm{d}\mathbf{w}$$

where $f(\mathbf{w}|n_d)$ is the conditional probability density function for location $\mathbf{w}$. In general, $f(\mathbf{w}|n_d)$ can vary with $\mathbf{w}$. In particular, under the FIA design the dispersion of sample points effected by the hexagonal tessellation means that a realized domain sample size can reduce $f(\mathbf{w}|n_d)$ to 0 for points in hexagons intersected by the domain bound-

aries. Yet under SRS, $f\left(\mathbf{w}|n_d\right) = f\left(\mathbf{w}\right) = \frac{1}{A_d}$ and thus (A2) simplifies to

$$(A3) \quad E\left(y_i \mid n_d\right) = \frac{1}{A_d}\int_{\mathcal{D}_d} y\left(\mathbf{w}\right) \, \mathrm{d}\mathbf{w} = \mu_{y,d}$$

The second inner expectation of (A1), $E\left(\tilde{y}_i\left(s\right) \mid n_d\right) = E\left(\tilde{y}\left(\mathbf{w}_i, s\right) \mid n_d\right)$ is not analytically tractable for $s = s_a$. The reason is that $\tilde{y}\left(\mathbf{w}_i, s_a\right)$ is a random function not only of the evaluation location $\mathbf{w}_i$, but also of all other sample locations through the nonlinear random function $H(i, k, s_a)$ of eq. (1). However, if $s = s_d$ then

$$(A4) \quad E\left(\tilde{y}\left(\mathbf{w}_i, s_{-d}\right) \mid n_d\right) = E\left[E\left(\tilde{y}\left(\mathbf{w}_i, s_{-d}\right) \mid n_d, s_{-d}\right) \mid n_d\right]$$

$$= E\left[\int_{\mathcal{D}_d} \tilde{y}\left(\mathbf{w}, s_{-d}\right) f\left(\mathbf{w}|n_d, s_{-d}\right) \, \mathrm{d}\mathbf{w} \mid n_d\right]$$

In general, $f(\mathbf{w}|n_d, s_{-d})$ can again vary with $\mathbf{w}$ but under SRS $f\left(\mathbf{w}|n_d, s_{-d}\right) = f\left(\mathbf{w}\right) = \frac{1}{A_d}$ such that

$$(A5) \quad E\left(\tilde{y}\left(\mathbf{w}_i, s_{-d}\right) \mid n_d\right) = E\left[\frac{1}{A_d}\int_{\mathcal{D}_d} \tilde{y}\left(\mathbf{w}, s_{-d}\right) \, \mathrm{d}\mathbf{w} \mid n_d\right]$$

$$= E\left[\mu_{\tilde{y},d}\left(s_{-d}\right) \mid n_d\right]$$

Substituting (A3) and (A5) into (A1) allow for the result

$$(A6) \quad E\left(\bar{e}_d\right) = E\left[\frac{1}{n_d}\sum_{i=1}^{n_d}\mu_{y,d}\right] - E\left[\frac{1}{n_d}\sum_{i=1}^{n_d} E\left[\mu_{\tilde{y},d}\left(s_{-d}\right) \mid n_d\right]\right]$$

$$= \mu_{y,d} - E\left[\mu_{\tilde{y},d}\left(s_{-d}\right)\right]$$

As such, under SRS the external KNN MA estimator (9) has expectation

$$(A7) \quad E\left(\widehat{\mu}_{y,d}\right) = E\left[\mu_{\tilde{y},d}\left(s_{-d}\right) + \bar{e}_d\right]$$

$$= E\left[\mu_{\tilde{y},d}\left(s_{-d}\right)\right] + \mu_{y,d} - E\left[\mu_{\tilde{y},d}\left(s_{-d}\right)\right]$$

$$= \mu_{y,d}$$

That is, the external KNN MA domain estimator is design-unbiased under SRS given $n_d > 0$. Unfortunately, this result does not generalize to the wider class of equal probability designs. In particular, it does not carry over to the FIA design where conditional (on $n_d$) spatial sampling intensities are not necessarily equal across arbitrary domains owing to the design's underlying hexagonal spatial structure. More importantly in the context of this research, result (A7) does not generalize to the unconstrained KNN approach that draws on the complete sample ($s_a$). Where KNN uses the full sample or any subsample that includes elements of $s_d$, $\tilde{y}_i(s)$ is a random nonlinear function of the sample locations (cf. Baffetta et al. 2009) and we cannot solve (A3) analytically.

Carrying forward the above noted conditions, the variance of the external KNN MA estimator can be written as

$$(A8) \quad V\left(\widehat{\mu}_{y,d}\right) = E\left[V\left(\widehat{\mu}_{y,d} \mid n_d, s_{-d}\right)\right] + V\left[E\left(\widehat{\mu}_{y,d} \mid n_d, s_{-d}\right)\right]$$

Focusing on the inner expectation of (A8) and using the above results (A3) and (A5)

$$(A9) \quad E\left(\widehat{\mu}_{y,d} \mid n_d, s_{-d}\right) = E\left(\mu_{\tilde{y},d}\left(s_{-d}\right) \mid n_d, s_{-d}\right) + E\left(\frac{1}{n_d}\sum_{i=1}^{n_d} y_i \mid n_d, s_{-d}\right) - E\left(\frac{1}{n_d}\sum_{i=1}^{n_d}\tilde{y}_i\left(s_{-d}\right) \mid n_d, s_{-d}\right)$$

$$= \mu_{\tilde{y},d}\left(s_{-d}\right) + \frac{1}{n_d}\sum_{i=1}^{n_d} E\left(y_i \mid n_d\right) - \frac{1}{n_d}\sum_{i=1}^{n_d} E\left(\tilde{y}_i\left(s_{-d}\right) \mid n_d, s_{-d}\right)$$

$$= \mu_{\tilde{y},d}\left(s_{-d}\right) + \frac{1}{n_d}\sum_{i=1}^{n_d}\mu_{y,d} - \frac{1}{n_d}\sum_{i=1}^{n_d}\mu_{\tilde{y},d}\left(s_{-d}\right) = \mu_{y,d}$$

Therefore the full variances reduces to

$$(A10) \quad V\left(\widehat{\mu}_{y,d}\right) = E\left[V\left(\mu_{\tilde{y},d}\left(s_{-d}\right) + \bar{e}_d \mid n_d, s_{-d}\right)\right] + V\left[\mu_{y,d}\right]$$

$$= E\left[V\left(\bar{e}_d \mid n_d, s_{-d}\right)\right]$$

That is, the design variance of the external KNN MA domain estimator is simply the conditional variance of the bias correction term, averaged over all possible partitions of the population/domain sample ($n_d$ and $s_{-d}$). We can estimate $V\left(\widehat{\mu}_{y,d}\right)$ as

$$(A11) \quad \widehat{V}\left(\widehat{\mu}_{y,d}\right) = \frac{1}{n_d\left(n_d - 1\right)}\sum_{i\in s_d}\left[e_i\left(s_{-d}\right) - \bar{e}_d\right]^2$$

though this does not recognize the additional variation induced by variability in $n_d$ and $s_{-d}$, nor that from development of the KNN model.