

# Model-based Survey Weighting Using Logistic Regression

Seonghun Lee

## Contents

<b>Abstract</b>	<b>1</b>
<b>1 Introduction</b>	<b>2</b>
<b>2 Hierarchical Gaussian Process Regression</b>	<b>3</b>
<b>3 Gaussian Processes</b>	<b>4</b>
3.1 Posterior prediction under gaussian process prior . . . . .	4
3.2 Computational cost reduction . . . . .	4
3.3 Connection to kernel ridge regression and crossed random effects model . . . . .	5
3.4 Uniform Convergence of the Posterior Predictive Mean . . . . .	7
<b>4 Connection to Survey Weighting</b>	<b>8</b>
<b>5 Results</b>	<b>9</b>
<b>6 Discussion</b>	<b>11</b>

*Abstract.* Survey weighting is a central tool for producing population-representative estimates in the presence of unequal sampling probabilities, nonresponse, and undercoverage. While classical design-based weights rely on known inclusion probabilities, modern surveys increasingly require model-based approaches to handle high-dimensional poststratification spaces and sparse cell counts. Building on prior work connecting logistic regression to equivalent survey weights, this paper develops a flexible Bayesian framework for survey weighting using hierarchical Gaussian process (HGP) regression. The proposed model decomposes the latent outcome-generating process into structured group-level effects and a smooth kernel-based component, allowing information to be shared across related poststratification cells while accommodating nonlinear and interaction effects. We derive closed-form expressions for posterior predictive means and covariances under monotonic link functions, and show how structured covariance representations yield substantial computational savings relative to standard Gaussian process models. A key theoretical contribution is a uniform convergence result for the posterior predictive mean under mild regularity conditions. Leveraging this result, we establish a formal connection between model-based poststratification and survey weighting by introducing *local equivalent unit weights*, defined via a first-order Taylor expansion of the poststratified estimator around the fitted latent mean surface. These weights generalize classical equivalent-weight representations for generalized linear models to the nonparametric Gaussian process setting, and admit simple closed forms for identity, logit, and probit links. Simulation studies demonstrate that the proposed HGP-based weights reduce bias and variance relative to logistic regression and kernel-based alternatives, particularly in sparse and highly structured populations. The results provide a principled and scalable approach to constructing survey weights that integrates modern Bayesian nonparametric modeling with classical design-based intuition.

## 1 Introduction

Survey weighting is a fundamental tool for correcting selection bias and achieving population-representative estimates, especially when data suffer from unequal sampling probabilities, nonresponse, or undercoverage. In our previous work, we explored a model-based framework for constructing base survey weights using logistic regression. That approach offered a structured alternative to traditional design-based weights, showing how predictive modeling could align with classical estimators like the Hájek estimator when appropriately normalized. However, logistic regression—while useful—assumes a linear structure on the log-odds scale and does not easily accommodate complex interactions or smooth variation across structured populations.

In this sequel, we introduce a more flexible model-based approach to survey weight construction using Hierarchical Gaussian Processes (HGPs). HGPs offer a principled way to model complex, nonlinear relationships across both continuous and discrete covariates, while borrowing strength across groups through hierarchical structure. This is particularly valuable in survey settings where small subgroup sizes or sparsely observed cells limit the stability of traditional poststratification or regression-based weighting methods.

A key computational innovation in our method is the use of structured covariance matrices that exploit the known grouping of the population. Rather than fitting a full Gaussian process across individual-level covariates—which can be computationally expensive due to the  $\mathcal{O}(n^3)$  scaling—we construct a covariance matrix over groups (e.g., poststratification cells from weighting variables). This reduces the dimensionality of the GP to the number of unique groups rather than individuals, yielding substantial computational savings while preserving flexibility in the modeled inclusion probabilities. In practice, this makes HGP-based weighting feasible even for large surveys with thousands of respondents and of poststratification cells.

We develop this approach within a fully Bayesian framework and compare the resulting population estimates to those from alternative models, and the resulting weights to both our previous model-based methods and classical design-based approaches. Our simulations demonstrate that the HGP-based method not only reduces bias in estimated outcomes and variance in estimated weights but also enhances robustness in the presence of data sparsity, providing more stable inference without requiring external raking or trimming. Additionally, while binary outcomes can limit a model’s flexibility and stability—especially in sparse cells—our method is designed to accommodate a broad range of outcome types, ensuring wider applicability across survey settings.

This paper contributes a flexible and scalable model-based framework for survey weighting, with a focus on structured populations. Our results support the use of HGPs for smoothing across cells and correcting for complex inclusion mechanisms—expanding the practical toolbox available for researchers seeking principled and efficient weighting strategies.

## 2 Hierarchical Gaussian Process Regression

Our goal is to estimate the population mean of an outcome variable  $y$  across structured subgroups, given a sample that sparsely covers the full poststratification space. Let the population be divided into  $H$  hierarchical groups (e.g., demographic strata for which population-level information is available) and  $K$  kernel groups (e.g., the full observed set of poststratification cells spanning all demographic combinations). In typical settings, we observe outcomes for only a small subset of the full set of poststratification cells (e.g., 1,100 observed cells out of 4,480 total), which leads to instability in direct cell-level estimation.

To overcome this sparsity, we propose a hierarchical Gaussian process (HGP) model that enables information sharing across related groups. Specifically, we assume that the latent outcome-generating process is composed of additive components: a structured group-level effect  $\mu_h$  and a smooth kernel-based effect  $\eta_k$  modeled by a Gaussian process.

Formally, for each observation  $i = 1, \dots, n$  in the sample, we assume the following model:

$$\mathbb{E}[Y \mid X] = g^{-1}(f(X)) + \mathbb{E}[\varepsilon], \quad \varepsilon \sim \mathcal{N}(0, \sigma^2),$$

where  $Y \in \mathbb{R}^n$ ,  $X \in \mathbb{R}^{n \times p}$ ,  $g(\cdot)$  is a monotonic link function, and the latent function is defined component-wise as

$$f_i = \mu_{h[i]} + \eta_{k[i]}, \quad \text{for } i = 1, \dots, n.$$

where  $\mu_{h[i]}$  is the fixed effect for group  $h[i] \in \{1, \dots, H\}$  and  $\eta_{k[i]}$  is the kernel-based effect for group  $k[i] \in \{1, \dots, K\}$ . For continuous outcomes, we set  $\sigma^2 > 0$ . For binary outcomes using a non-identity link (e.g., probit or logit), we set  $\sigma^2 = 0$  to reflect the deterministic link from the latent scale to the probability scale.

We place the following priors on the latent components:

$$\mu \sim \mathcal{N}(\boldsymbol{\mu}_{\text{prior}}, \Sigma_{\mu}), \quad \eta \sim \mathcal{GP}(0, K_{\eta}),$$

where  $\boldsymbol{\mu}_{\text{prior}} \in \mathbb{R}^H$  encodes known population-level means, and  $K_{\eta}$  is the Radial Basis Function (RBF) kernel defined over  $K$  structured groups. The squared exponential kernel is computed via:

$$[K_{\eta}]_{ij} = \alpha^2 \exp\left(-\frac{\|X_i - X_j\|^2}{2\rho^2}\right),$$

where  $K_{\eta} : \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}$  given that  $X_i$  and  $X_j$  are  $P$ -dimensional covariate vectors for groups  $i$  and  $j$ ,  $\alpha$  is the marginal standard deviation, and  $\rho$  is the length-scale parameter.

### 3 Gaussian Processes

#### 3.1 Posterior prediction under gaussian process prior

**Notation.** For convenience, we denote  $g(\mathbb{E}[Y | X])$  as  $g(\mu)$  throughout the remainder of this paper.

**Theorem 1** (Posterior Prediction in Gaussian Process with Monotonic Link). *Assume a monotonic link function  $g$  such that:*

$$g(\mu) \sim \mathcal{N}(\boldsymbol{\mu}_{\text{prior},h}, \boldsymbol{\Sigma}_{\mu}[h, h'] + \mathbf{K}_{\eta}[k, k'] + \sigma^2 \mathbf{I}_n),$$

then the posterior prediction of  $f_*$  given  $g(\mu)$  is:

$$f_* | g(\mu) \sim \mathcal{N}(\boldsymbol{\mu}_{f_*} + \boldsymbol{\Sigma}_{f_*y} \boldsymbol{\Sigma}_{yy}^{-1} (g(\mu) - \boldsymbol{\mu}_y), \boldsymbol{\Sigma}_{f_*f_*} - \boldsymbol{\Sigma}_{f_*y} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yf_*}),$$

which we denote as:

$$f_* | g(\mu) \sim \mathcal{N}(\mathbb{E}[f_* | g(\mu)], \text{Cov}(f_* | g(\mu))),$$

and hence:

$$g(\mathbf{y}_*) | g(\mu) \sim \mathcal{N}(\mathbb{E}[f_* | g(\mu)], \text{Cov}(f_* | g(\mu)) + \sigma^2 \mathbf{I}_{n_*}).$$

**Proof:** The joint distribution of  $\begin{bmatrix} g(\mu) \\ f_* \end{bmatrix}$  is:

$$\begin{bmatrix} g(\mu) \\ f_* \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}_{\text{prior},h} \\ \boldsymbol{\mu}_{\text{prior},h_*} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{\mu}[h, h'] + \mathbf{K}_{\eta}[k, k'] + \sigma^2 \mathbf{I}_n & \boldsymbol{\Sigma}_{\mu}[h, h_*] + \mathbf{K}_{\eta}[k, k_*] \\ \boldsymbol{\Sigma}_{\mu}[h_*, h] + \mathbf{K}_{\eta}[k_*, k] & \boldsymbol{\Sigma}_{\mu}[h_*, h_*] + \mathbf{K}_{\eta}[k_*, k_*] \end{bmatrix} \right),$$

which we denote as:

$$\begin{bmatrix} g(\mu) \\ f_* \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}_y \\ \boldsymbol{\mu}_{f_*} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yf_*} \\ \boldsymbol{\Sigma}_{f_*y} & \boldsymbol{\Sigma}_{f_*f_*} \end{bmatrix} \right).$$

Applying the standard formula for conditional distributions of multivariate Gaussians, we obtain:

$$f_* | g(\mu) \sim \mathcal{N}(\boldsymbol{\mu}_{f_*} + \boldsymbol{\Sigma}_{f_*y} \boldsymbol{\Sigma}_{yy}^{-1} (g(\mu) - \boldsymbol{\mu}_y), \boldsymbol{\Sigma}_{f_*f_*} - \boldsymbol{\Sigma}_{f_*y} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yf_*}).$$

Finally, because  $g(\mathbf{y}_*) = f_* + \boldsymbol{\epsilon}_*$  with  $\boldsymbol{\epsilon}_* \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_{n_*})$ , we conclude:

$$g(\mathbf{y}_*) | g(\mu) \sim \mathcal{N}(\mathbb{E}[f_* | g(\mu)], \text{Cov}(f_* | g(\mu)) + \sigma^2 \mathbf{I}_{n_*}).$$

#### 3.2 Computational cost reduction

**Theorem 2** (Computational Efficiency via Woodbury Identity). *The computational cost of inverting the  $n \times n$  covariance matrix  $\boldsymbol{\Sigma}_{\mathbf{y}\mathbf{y}}$  is reduced from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(n^2(H + K))$ , where  $H$  and  $K$  are the dimensions of the group-level covariates and  $H, K \ll n$ .*

**Proof:** Write the model as:

$$f + \boldsymbol{\epsilon} = Z_{\mu} \boldsymbol{\mu}_h + Z_{\eta} \boldsymbol{\eta}_k + \boldsymbol{\epsilon}$$

where  $Z_{\mu} \in \mathbb{R}^{n \times H}$  and  $Z_{\eta} \in \mathbb{R}^{n \times K}$  are group incidence matrices with entries

$$(Z_{\mu})_{ih} = \mathbf{1}\{h[i] = h\}, \quad (Z_{\eta})_{ik} = \mathbf{1}\{k[i] = k\},$$

and  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$ . Then, the covariance matrix is:

$$V = Z_{\mu} \boldsymbol{\Sigma}_{\mu} Z_{\mu}^{\top} + Z_{\eta} \mathbf{K}_{\eta} Z_{\eta}^{\top} + \sigma^2 \mathbf{I}_n$$

Using the Woodbury identity:

$$V^{-1} = \left( \sigma^2 I_n + U C U^\top \right)^{-1} = \sigma^{-2} I_n - \sigma^{-2} U \left( C^{-1} + \sigma^{-2} U^\top U \right)^{-1} U^\top \sigma^{-2}$$

with  $U = [Z_\mu \ Z_\eta] \in \mathbb{R}^{n \times (H+K)}$ ,  $C = \text{diag}(\Sigma_\mu, K_\eta)$ .

The major costs are:

- $U^\top U$ :  $\mathcal{O}(n^2(H+K))$
- $(C^{-1} + U^\top U)^{-1}$ :  $\mathcal{O}((H+K)^3)$

Full expression:  $\mathcal{O}(2n^2(H+K) + n(H+K)^2 + (H+K)^3)$ , which simplifies to  $\mathcal{O}(n^2(H+K))$ .

#### Algorithm for Computational Efficiency Using Cholesky Factorization

##### Step 1: Construct Covariance Matrices

Compute  $\Sigma_{yy}, \Sigma_{yf_*}, \Sigma_{f_*y}, \Sigma_{f_*f_*}$  using the estimated covariance matrices.

[The computational cost reduces from  $\mathcal{O}(n^3)$  to  $\mathcal{O}(n^2(H+K))$  and  $\mathcal{O}(n_*^3)$  to  $\mathcal{O}(n_*^2(H_*+K_*))$ .]

Continuous outcome:  $\Sigma_{yy} = \Sigma_{\mu, \text{full}} + K\eta_{\text{full}} + \sigma^2 I_n$

Binary outcome:  $\Sigma_{yy} = \Sigma_{\mu, \text{full}} + K\eta_{\text{full}} + \epsilon I_n$ , where  $\epsilon > 0$  ensures positive definiteness.

##### Step 2: Cholesky Decomposition & Posterior Mean

Compute  $L := \text{Cholesky}(\Sigma_{yy})$  and  $g(\mu)_c := g(\mu) - \mu_y = f - \mu_y$ .

Compute  $V_1 := \text{solve}(L^\top, \text{solve}(L, g(\mu)_c)) = (LL^\top)^{-1}g(\mu)_c$

Then,  $\mathbb{E}[f_* | g(\mu)] = \mu_{f_*} + \Sigma_{f_*y} V_1$

##### Step 3: Posterior Covariance

Compute  $V_2 := \text{solve}(L, \Sigma_{yf_*})$

Then,  $\text{Cov}(f_* | g(\mu)) = \Sigma_{f_*f_*} - V_2^\top V_2$

##### Step 4: Posterior Sampling

Identity link:

$$y_* | \mu \sim \mathcal{N}(\mathbb{E}[f_* | \mu], \text{Cov}(f_* | \mu) + \sigma^2 I_{n_*})$$

Binary outcome:

$$f_* | g(\mu) \sim \mathcal{N}(\mathbb{E}[f_* | g(\mu)], \text{Cov}(f_* | g(\mu)))$$

$$p_* \equiv g^{-1}(f_*)$$

$$y_* | f_* \sim \text{Bernoulli}(p_*)$$

### 3.3 Connection to kernel ridge regression and crossed random effects model

**Proposition 1** (Equivalence of KRR and GP Posterior Mean under Identity Link). *Consider our hierarchical Gaussian process model with an **identity link**  $g(x) = x$ , noise variance  $\sigma^2$ , and suppose the group-level prior covariance  $\Sigma_\mu \rightarrow 0$  (so that  $\mu_h \equiv 0$ ). Let*

$$\mathbf{K} = K_\eta[k, k'],$$

*the  $n \times n$  kernel matrix over the  $n$  observed groups, and let  $\mathbf{y} \in \mathbb{R}^n$  be the observed outcomes. Then the GP posterior mean at a new input  $x^*$  is*

$$\mathbb{E}[f(x^*) | \mathbf{y}] = \mathbf{k}(x^*, X)^\top (\mathbf{K} + \sigma^2 I_n)^{-1} \mathbf{y},$$

which coincides exactly with the Kernel Ridge Regression solution

$$\hat{f}(x^*) = \sum_{i=1}^n \alpha_i k(x_i, x^*) \quad \text{with} \quad \boldsymbol{\alpha} = (\mathbf{K} + \lambda I_n)^{-1} \mathbf{y} \quad \text{and} \quad \lambda = \sigma^2.$$

**Proof:** Under  $\Sigma_\mu \rightarrow 0$ , our latent function reduces to  $f = \eta \sim \mathcal{GP}(0, K_\eta)$ . With identity link and Gaussian noise,

$$\mathbf{y} \sim \mathcal{N}(0, \mathbf{K} + \sigma^2 I_n).$$

The standard GP posterior mean at  $x^*$  is

$$\mathbb{E}[f(x^*) \mid \mathbf{y}] = \mathbf{k}(x^*, X)^\top (\mathbf{K} + \sigma^2 I_n)^{-1} \mathbf{y},$$

where  $\mathbf{k}(x^*, X)_i = k_\eta(x^*, x_i)$ . Let  $\mathcal{H}_k$  be a Reproducing Kernel Hilbert Space (RKHS) associated with a positive-definite kernel function  $K(\cdot, \cdot)$ . Then, KRR solves

$$\hat{f} = \arg \min_{f \in \mathcal{H}_k} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{H}_k}^2$$

and admits the dual form

$$\hat{f}(x^*) = \mathbf{k}(x^*, X)^\top (\mathbf{K} + \lambda I_n)^{-1} \mathbf{y}.$$

Thus, setting  $\lambda = \sigma^2$  makes the GP posterior mean and the KRR prediction identical.

**Proposition 2.** *Without loss of generality, let the number of covariates  $p$  be 3. Let  $X \in \mathbb{R}^{K \times p}$  be the matrix of covariates for the  $K$  hierarchical groups, and set*

$$K_\eta = \tau^2 X X^\top.$$

*Assume the identity link  $g(x) = x$ , observation noise  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ , and a group-intercept prior  $\mu \sim \mathcal{N}(\mu_{\text{prior}}, \Sigma_\mu)$ . Let*

$$b_0 \sim \mathcal{N}(\mu_{\text{prior}}, \Sigma_\mu) \quad b_1 \sim \mathcal{N}(0, \tau^2 I_p)$$

*Then the hierarchical GP model*

$$y = \mu_h + \eta_k + \varepsilon, \quad \eta \sim \mathcal{N}(0, K_\eta),$$

*is equivalent to the crossed random effects model with a random intercept  $b_{0,ij}$  and a random-slope  $b_{1,ijk}$ :*

$$y_{ijk} = b_{0,ij} + X_{ijk} b_{1,ijk} + \varepsilon_{ijk}, \text{ where } i, j, k \text{ represent the levels of three covariates.}$$

**Proof:** By construction,

$$\eta \sim \mathcal{N}(0, K_\eta) = \mathcal{N}(0, \tau^2 X X^\top).$$

Introducing  $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \tau^2 I_p)$  and setting  $\eta = X \mathbf{u}$  yields

$$\text{Cov}(\eta) = X \text{Cov}(\mathbf{u}) X^\top = \tau^2 X X^\top = K_\eta, \quad \mathbb{E}[\eta] = \mathbf{0}.$$

Matching  $i, j$  to a group  $H$  and  $i, j, k$  to a group  $K$  gives

$$y_i = \mu_i + \eta_i + \varepsilon_i = b_{0,ij} + X_{ijk} b_{1,ijk} + \varepsilon_{ijk},$$

with the specified priors on  $b_0$  and  $b_1$ .

### 3.4 Uniform Convergence of the Posterior Predictive Mean

**Theorem 3.** Consider the HGP regression model with identity link:

$$y_i = f(x_i) + \varepsilon_i, \quad f_i = \mu_{h[i]} + \eta_{k[i]}, \quad \begin{cases} \mu \sim \mathcal{N}(\boldsymbol{\mu}_{\text{prior}}, \Sigma_\mu), \\ \eta \sim \mathcal{GP}(0, K_\eta), \end{cases} \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2),$$

where  $X \in \mathbb{R}^{n \times p}$  and  $K_\eta$  is the squared-exponential kernel over the  $K$  unique rows of  $X$ . Let

$$\hat{f}_n(x) = \mathbb{E}[\mu_{h(x)} + \eta_{k(x)} \mid y_{1:n}, X]$$

be the posterior predictive mean, and  $f_0(x)$  the true regression function of the same form.

Assume:

1.  $f_0$  lies in the RKHS of  $K_\eta$  with smoothness  $\alpha > 0$ .
2. The design rows  $\{X_{i,:}\}_{i=1}^n$  densely cover the compact covariate domain  $\mathcal{X}$ .
3. The prior covariance  $\Sigma_\mu$  has eigenvalues bounded away from zero.

Then there exists a rate

$$\varepsilon_n = n^{-\alpha/(2\alpha+1)} (\log n)^t$$

such that

$$\Pr\left(\sup_{x \in \mathcal{X}} |\hat{f}_n(x) - f_0(x)| > M \varepsilon_n \mid y_{1:n}, X\right) \xrightarrow{a.s.} 0$$

for every large  $M > 0$ .

**Proof:** We combine three key ingredients, adapted to the additive structure  $f = \mu + \eta$ :

**1. Posterior Contraction for  $(\mu, \eta)$ .** Show that the joint posterior concentrates in an  $\ell^\infty$ -ball around  $(\mu_0, \eta_0)$ , at rate  $\varepsilon_n$ , using standard RKHS-based concentration for  $\eta$  and Gaussian-prior concentration for  $\mu$ .

**2. Finite Dimensional Bernstein–von Mises.** We expand the process  $\eta$  using its first  $J_n$  eigenfunctions  $\phi_j(x)$  and corresponding coefficients  $\theta_j$ . We also include  $\mu$  as  $H$  additional finite parameters. The joint posterior distribution of  $\mu$  and the coefficients  $\theta_{1:J_n}$  is then approximated by a multivariate Bernstein–von Mises (BvM) theorem. This approximation holds for the parameter vector  $(\mu, \theta_{1:J_n})$  with dimension  $H + J_n = o(n)$ , yielding the following normal approximation:

$$(\mu, \theta_{1:J_n}) \mid y \approx \mathcal{N}(\widehat{(\mu, \theta)}, n^{-1}I),$$

where  $\widehat{(\mu, \theta)}$  represents the posterior mean of the parameters, and  $I$  is the identity matrix. This approximation is valid as  $n \rightarrow \infty$ , ensuring that the posterior for these parameters is concentrated around the true values with covariance scaling as  $n^{-1}$ .

**3. Sup-Norm Control.** Write

$$\hat{f}_n(x) - f_0(x) = \underbrace{\sum_{h=1}^H (\hat{\mu}_h - \mu_{0,h}) \mathbf{1}\{h(x) = h\}}_{O_P(n^{-1/2} H^{1/2})} + \underbrace{\sum_{j=1}^{J_n} (\hat{\theta}_j - \theta_{0j}) \phi_j(x)}_{O_P(n^{-1/2} J_n^{1/2})} + \underbrace{\sum_{j > J_n} \theta_{0j} \phi_j(x)}_{O(J_n^{-\alpha})}.$$

Choosing  $J_n \asymp n^{1/(2\alpha+1)}$  balances these terms so that  $\sup_x |\hat{f}_n - f_0| = O_P(\varepsilon_n)$ . A subsequence argument then upgrades to almost-sure convergence.

## 4 Connection to Survey Weighting

We now present a theorem that connects the regression function  $g^{-1}(f)$  to population-calibrated survey weights.

**Theorem 4** (Local equivalent unit weights). *Let  $z \equiv g(\mu) = g(E[Y \mid X]) \in \mathbb{R}^n$  denote the latent mean-scale vector, and let  $f_* \in \mathbb{R}^K$  denote the latent GP values at the  $K$  poststratification cells  $X_*$ . Under the setup of Theorem 3, the posterior predictive mean*

$$m_*(z) \equiv \mathbb{E}[f_* \mid z] = \mu_{f_*} + \Sigma_{f_*y} \Sigma_{yy}^{-1}(z - \mu_y)$$

*is affine in  $z$ .*

*Let  $N_k^{\text{pop}}$  be population cell counts,  $N = \sum_{k=1}^K N_k^{\text{pop}}$ , and  $N^{\text{pop}} = (N_1^{\text{pop}}, \dots, N_K^{\text{pop}})^\top$ . Define the model-based poststratified estimator*

$$T(z) \equiv \frac{1}{N} (N^{\text{pop}})^\top g^{-1}(m_*(z)),$$

*where  $g^{-1}$  is applied componentwise.*

*Fix the fitted value  $\hat{z}$ . Mirroring the exact equivalent-weight representation for generalized linear models, we replace global equality by first-order local equivalence. The local equivalent unit weights  $w \in \mathbb{R}^n$  are defined as the unique vector satisfying*

$$T(z) = T(\hat{z}) + \frac{1}{n} w^\top (y - \hat{y}) + o(\|z - \hat{z}\|), \quad z \rightarrow \hat{z}, \quad (1)$$

*where  $\hat{y}$  denotes the fitted outcomes under  $\hat{z}$ . Then the weights are given by*

$$\boxed{w = \frac{n}{N} \Sigma_{yy}^{-1} \Sigma_{yf_*} D_* N^{\text{pop}}}, \quad (2)$$

*with*

$$D_* = \text{diag}((g^{-1})'(m_{*,1}(\hat{z})), \dots, (g^{-1})'(m_{*,K}(\hat{z}))).$$

*Proof.* By Theorem 3, the plug-in surface  $\hat{z}$  converges uniformly to  $z = g(\mu)$ , so a first-order Taylor expansion of  $T(z)$  around  $\hat{z}$  is valid:

$$T(z) = T(\hat{z}) + \nabla T(\hat{z})^\top (z - \hat{z}) + o(\|z - \hat{z}\|).$$

The local equivalent unit weights are defined so that the linear term in this expansion can be represented as a weighted sample mean, yielding

$$\nabla T(\hat{z})^\top (z - \hat{z}) = \frac{1}{n} w^\top (y - \hat{y}).$$

Matching coefficients implies  $w = n \nabla T(\hat{z})$ . It remains to compute  $\nabla T(\hat{z})$ . Since  $m_*(z)$  is affine in  $z$ ,

$$\frac{\partial m_*}{\partial z} = \Sigma_{f_*y} \Sigma_{yy}^{-1}.$$

Applying the chain rule to  $T(z)$ , noting that  $g^{-1}$  acts componentwise, gives

$$\nabla T(\hat{z})^\top = \frac{1}{N} (N^{\text{pop}})^\top D_* \Sigma_{f_*y} \Sigma_{yy}^{-1}.$$

Transposing and multiplying by  $n$  yields the stated expression for  $w$  in (2).  $\square$



**Special cases.** Write

$$\eta_* \equiv m_*(\hat{z}) \in \mathbb{R}^K, \quad D_* = \text{diag}((g^{-1})'(\eta_{*,1}), \dots, (g^{-1})'(\eta_{*,K})).$$

Then Theorem 4 gives

$$w = \frac{n}{N} \Sigma_{yy}^{-1} \Sigma_{yf*} D_* N^{\text{pop}}.$$

- **Identity (linear).**  $g^{-1}(\eta) = \eta$ , so  $(g^{-1})'(\eta) = 1$  and  $D_* = I_K$ . Hence

$$w = \frac{n}{N} \Sigma_{yy}^{-1} \Sigma_{yf*} N^{\text{pop}}.$$

- **Logit.**  $g^{-1}(\eta) = \text{expit}(\eta) = \{1 + \exp(-\eta)\}^{-1}$ . Let  $p_* \equiv \text{expit}(\eta_*)$  (componentwise). Then  $(g^{-1})'(\eta) = \text{expit}(\eta)\{1 - \text{expit}(\eta)\}$ , so  $D_* = \text{diag}(p_{*,1}(1 - p_{*,1}), \dots, p_{*,K}(1 - p_{*,K}))$  and

$$w = \frac{n}{N} \Sigma_{yy}^{-1} \Sigma_{yf*} \text{diag}(p_* \odot (1 - p_*)) N^{\text{pop}},$$

where  $\odot$  denotes elementwise multiplication.

- **Probit.**  $g^{-1}(\eta) = \Phi(\eta)$ . Let  $\phi$  denote the standard normal pdf. Since  $(g^{-1})'(\eta) = \phi(\eta)$ , we have  $D_* = \text{diag}(\phi(\eta_{*,1}), \dots, \phi(\eta_{*,K}))$  and

$$w = \frac{n}{N} \Sigma_{yy}^{-1} \Sigma_{yf*} \text{diag}(\phi(\eta_*)) N^{\text{pop}}.$$

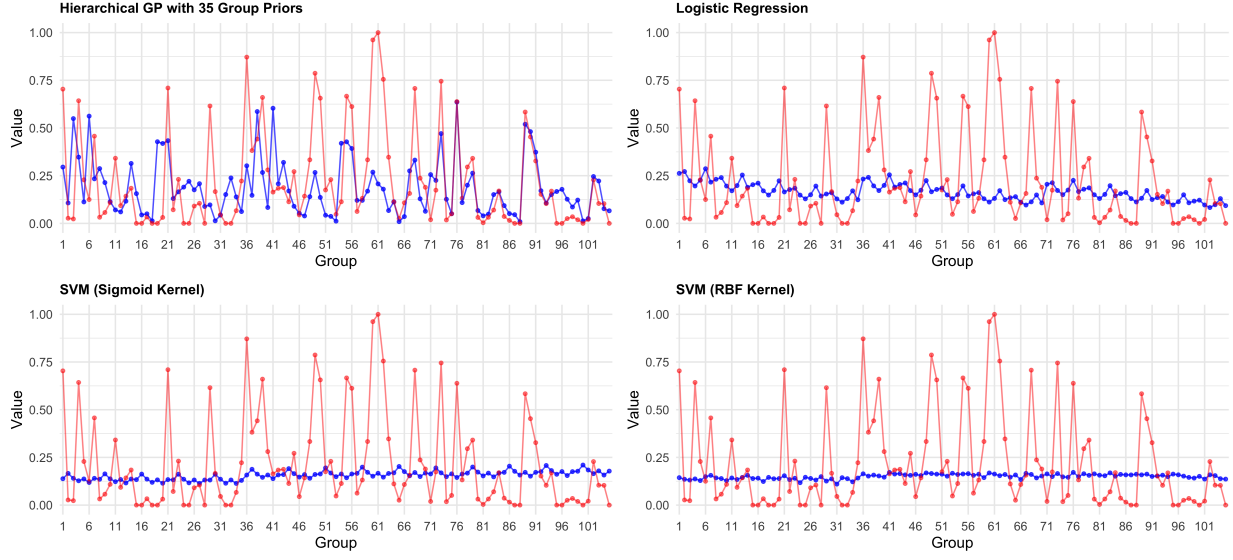
## 5 Results

In Table 1, we compare the Mean Squared Error (MSE) of four predictive models—Hierarchical Gaussian Process (HGP), Logistic Regression, Support Vector Machine with a Sigmoid kernel (SVM Sigmoid), and Support Vector Machine with a Radial Basis Function kernel (SVM RBF)—under two prior configurations: one with 7 group priors and another with 35 priors out of the 105 group priors. The MSE values quantify the discrepancy between each model’s predicted group rates and the actual population rates, with lower values indicating more accurate predictions.

Figures 1 and 2 visualize the predicted (blue) versus actual population (red) group rates under the 7-prior and 35-prior settings, respectively. These visualizations allow us to assess how closely each model’s predictions track the true group-level patterns. The analysis is conducted under intentionally sparse sampling conditions: the sample includes only 300 units from a total population of 10,000, and 70 out of the 105 possible population-defined groups are observed in the sample. This setup evaluates model performance in small-sample, partially observed group settings. Due to computational constraints, we chose to work with 105 total groups rather than the hundreds or thousands that may be encountered in large-scale population modeling. This choice balances realism with tractability while still providing a meaningful test of each model’s predictive accuracy.

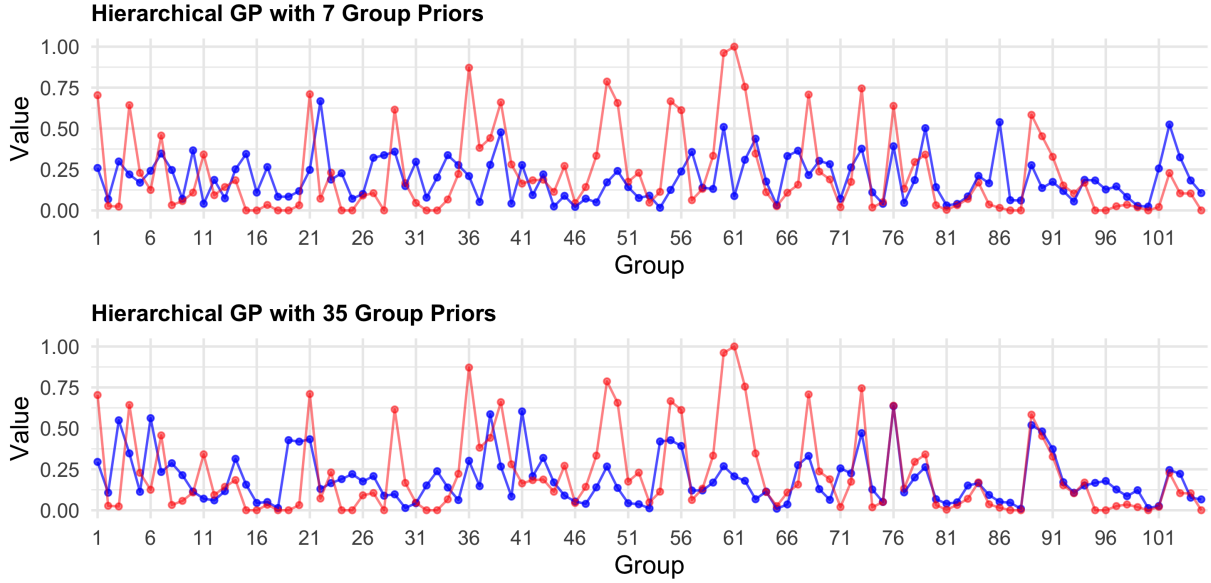
Priors	HGP	Logistic	SVM Sigmoid	SVM RBF
7 priors	0.062	0.066	0.068	0.065
35 priors	0.053	0.066	0.068	0.065

Table 1: MSE for Different Methods



*Blue = Predicted Rates, Red = Population Rates.*

Figure 1: Predicted Group Rates vs. Population Rates for 35 out of 105 priors.



*Blue = Predicted Rates, Red = Population Rates.*

Figure 2: HGP 7 Group Priors vs. 35 Group Priors.

We now present the results of the survey weights after applying Theorem 4. Table 2 compares the original population counts across cities with the corresponding weighted estimates derived from the adjusted survey weights. Table 3 summarizes the effective sample size after weighting, along with descriptive statistics of the normalized weights, highlighting their distribution and variability.

Table 2: Population Counts and Weighted Estimates by City

City	Pop. Count	Weighted Est.
1	1108	922.24
2	1221	985.26
3	1619	1959.19
4	1678	2540.39
5	1611	1537.32
6	1342	918.94
7	1421	1136.67

Table 3: Summary of Normalized Weights and Effective Sample Size

Statistic	Value
ESS	165 out of 300
Minimum	0.01125
1st Quartile	17.21502
Median	29.09463
Mean	34.36426
3rd Quartile	41.72785
Maximum	217.29819

## 6 Discussion

The data-generating process incorporated demographic covariates—education (3 levels), age group (5 levels), and city of residence (7 levels)—sampled to reflect empirical marginal distributions. Since we are dealing with weighting variables, all predictors are categorical. Ordinal variables, such as education, were treated numerically, while nominal variables, like city of residence, were encoded using one-hot encoding to appropriately capture their categorical nature. We simulated a binary outcome using both structured main effects and a smooth, nonparametric latent component to capture complex group heterogeneity. Structured effects were generated by applying linear coefficients to these covariates, while latent variation was introduced through a smooth, nonparametric function defined over 105 group cells formed by the interaction of city, education, and age. This latent function was drawn from a Gaussian process with a radial basis function (RBF) kernel, enabling flexible modeling of unobserved structure across cells. The final outcome probability was computed by applying a sigmoid transformation to the sum of the structured and latent components, resulting in probabilities that reflect both observed covariates and underlying heterogeneity.

Figure 1 and Table 1 highlight the clear advantage of the Hierarchical Gaussian Process (HGP) model in capturing latent structure across demographic subgroups, each defined by a unique combination of education, age group, and city of residence (yielding 105 total groups). Unlike logistic regression and support vector machines (SVMs), which only account for structured covariate effects and do not incorporate priors, the HGP flexibly adapts to variation across subgroups by leveraging prior information. This flexibility becomes increasingly effective as the amount of prior information grows: the HGP with 35 priors achieves the lowest mean squared error ( $MSE = 0.053$ ), outperforming all other models, whose MSEs remain comparatively high and flat (0.065–0.068) regardless of the number of priors. In Figure 1, the HGP demonstrates a stronger capacity to track the wide range of subgroup-level empirical outcome rates, particularly in areas where unobserved heterogeneity plays a dominant role. In contrast, the alternative models struggle to account for this broad distribution of subgroup signals, resulting in underfitting in regions with complex structure. These findings underscore the value of incorporating latent interaction patterns and prior structure, especially when observable covariates alone are insufficient to explain outcome variation. By modeling at the group level, our approach also reduces the dimensionality of the latent function and enables effective sharing of information across similar subgroups.

The target survey population consists of live births occurring in large U.S. cities with population over 200,000 between 1998 and 2000. This focus is motivated by the process of constructing base weights of the Future of Families and Child Wellbeing Study (FFCWS). The FFCWS sample is a stratified and multistage design with 4,898 children, oversampling births to unmarried mothers

at a ratio of 3 to 1 with the inclusion of a large number of Black, Hispanic, and low-income families. Follow-up interviews were conducted across seven waves, when children were approximately ages 1, 3, 5, 9, 15, and 22. In constructing the FFCWS weights, four demographic variables were used for poststratification, and geographical information was used to estimate population birth counts using the Centers for Disease Control and Prevention (CDC) annual natality data. These variables are mother’s marital status, race/ethnicity, age, education, and city of birth, resulting in 4,480 group cells. The size and complexity of the FFCWS sample and its cell structure impose substantial computational and storage demands. As such, we plan to apply our method to this target dataset in the near future.