

1 Variance Reduction (continued)

1.3 Stratification

The idea in stratified sampling is to split up the domain \mathcal{D} of \mathbf{X} into separate regions, take a sample of points from each such region, and combine the results to estimate $E(f(\mathbf{X}))$. Intuitively, if each region gets its fair share of points then we should get a better answer. We might be able to do better still by oversampling within the important strata and under-sampling those in which f is nearly constant.

Our goal is to estimate $\mu = \int_{\mathcal{D}} f(\mathbf{x})p(\mathbf{x})d\mathbf{x}$. We partition \mathcal{D} into mutually exclusive and exhaustive regions \mathcal{D}_j , for $j = 1, \dots, J$. These regions are the strata. We write $\omega_j = P(\mathbf{X} \in \mathcal{D}_j)$ and to avoid trivial issues, we assume $\omega_j > 0$ for all j . Next, let $p_j(\mathbf{x}) = \omega_j^{-1}p(\mathbf{x})1_{\mathbf{x} \in \mathcal{D}_j}$, the conditional density of \mathbf{X} given that $\mathbf{X} \in \mathcal{D}_j$. To use stratified sampling, we must know the probabilities ω_j of the strata, and we must also know how to sample $\mathbf{X} \sim p_j$ for $j = 1, \dots, J$. These conditions are quite reasonable. When we are defining strata, we naturally prefer ones we can sample from.

Let $\mathbf{X}_{ij} \stackrel{i.i.d.}{\sim} p_j$ for $i = 1, \dots, n_j$ and $j = 1, \dots, J$ be sampled independently. The stratified sampling estimate of μ is

$$\hat{\mu}_{\text{strat}} = \sum_{j=1}^J \frac{\omega_j}{n_j} \sum_{i=1}^{n_j} f(\mathbf{X}_{ij}).$$

We choose $n_j > 0$ so that $\hat{\mu}_{\text{strat}}$ is properly defined. Unless otherwise specified, we make sure that $n_j > 2$, which will allow the variance estimate below to be applied.

Note that the stratified estimator of μ is unbiased as

$$\begin{aligned} E(\hat{\mu}_{\text{strat}}) &= \sum_{j=1}^J \omega_j E\left(\frac{1}{n_j} \sum_{i=1}^{n_j} f(\mathbf{X}_{ij})\right) \\ &= \sum_{j=1}^J \omega_j \int_{\mathcal{D}_j} f(\mathbf{x})p_j(\mathbf{x})d\mathbf{x} \\ &= \sum_{j=1}^J \int_{\mathcal{D}_j} f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &= \int_{\mathcal{D}} f(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &= \mu. \end{aligned}$$

Let $\mu_j = \int_{\mathcal{D}_j} f(\mathbf{x})p_j(\mathbf{x})d\mathbf{x}$ and $\sigma_j^2 = \int_{\mathcal{D}_j} (f(\mathbf{x}) - \mu_j)^2 p_j(\mathbf{x})d\mathbf{x}$ be the j th stratum mean and variance, respective. The variance of the stratified sampling estimate is

$$\text{Var}(\hat{\mu}_{\text{start}}) = \sum_{j=1}^J \omega_j^2 \frac{\sigma_j^2}{n_j}.$$

An immediate consequence is that $Var(\widehat{\mu}_{\text{start}}) = 0$ for integrands f that are constant within strata \mathcal{D}_j for all $j = 1, 2, \dots, J$.

For error estimation, we can proceed as follows. Denoting $f(X_{ij})$ by Y_{ij} , we have

$$\widehat{Var}(\widehat{\mu}_{\text{start}}) = \sum_{j=1}^J \omega_j^2 \frac{s_j^2}{n_j},$$

where

$$\widehat{\mu}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} Y_{ij} \quad \text{and} \quad s_j^2 = \frac{1}{n_j - 1} \sum_{i=1}^{n_j} (Y_{ij} - \widehat{\mu}_j)^2.$$

A CLT based 99% confidence interval for μ is

$$\widehat{\mu}_{\text{start}} \mp 2.58 \sqrt{\widehat{Var}(\widehat{\mu}_{\text{start}})}.$$

A natural choice for stratum sample sizes is proportional allocation, $n_j = n\omega_j$. In our analysis, we will suppose that all the n_j are integers. We can usually choose n and D_j to make this so, or else accept small non-proportionalities due to rounding. For proportional allocation, $\widehat{\mu}_{\text{strat}}$ reduces to the ordinary sample mean

$$\widehat{\mu}_{\text{strat}} = \frac{1}{n} \sum_{j=1}^J \sum_{i=1}^{n_j} f(\mathbf{X}_{ij}),$$

and $Var(\widehat{\mu}_{\text{start}})$ becomes

$$Var(\widehat{\mu}_{\text{strat}}) = \frac{1}{n} \sum_{j=1}^J \omega_j \sigma_j^2.$$

Note that the variance of $f(\mathbf{X})$ can be decomposed as follows:

$$\begin{aligned} \sigma^2 &= Var(f(\mathbf{X})) \\ &= E(Var(f(\mathbf{X})|\mathbf{X} \in \mathcal{D}_j)) + Var(E(f(\mathbf{X})|\mathbf{X} \in \mathcal{D}_j)) \\ &= \sum_{j=1}^J \omega_j \sigma_j^2 + \sum_{j=1}^J \omega_j (\mu_j - \mu)^2. \end{aligned}$$

Therefore,

$$Var(\widehat{\mu}_{\text{strat}}) = \frac{1}{n} \sum_{j=1}^J \omega_j \sigma_j^2 \leq \frac{\sigma^2}{n}.$$

It shows that stratified sampling with proportional allocation cannot have larger variance than simple Monte Carlo sampling. A good stratification scheme is one that reduces the within-stratum variance.

1.4 Conditioning

Sometimes we can do part of the problem in closed form, and then do the rest of it by Monte Carlo. For example, suppose that we want to find $\mu = \int_0^1 \int_0^1 f(x, y) dx dy$, where $f(x, y) = e^{yg(x)}$. It is easy to integrate out y for fixed x , yielding $h(x) = (e^{g(x)} - 1)/g(x)$. Then we have a one dimensional problem, which may be simpler to handle. In general, suppose that $\mathbf{X} \in \mathbb{R}^k$ and $\mathbf{Y} \in \mathbb{R}^{d-k}$ are random vectors and that we want to estimate $E(f(\mathbf{X}, \mathbf{Y}))$. The natural estimate is $\hat{\mu} = (1/n) \sum_{i=1}^n f(\mathbf{X}_i, \mathbf{Y}_i)$, where $(\mathbf{X}_i, \mathbf{Y}_i) \in \mathbb{R}^d$ are independent samples from the joint distribution of (\mathbf{X}, \mathbf{Y}) . Now, let $h(\mathbf{x}) = E(f(\mathbf{X}, \mathbf{Y})|\mathbf{X} = \mathbf{x})$. We might also estimate μ by

$$\hat{\mu}_{\text{cond}} = \frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i)$$

where \mathbf{X}_i are independently sampled from the distribution of \mathbf{X} . The justification for the method is that $E(f(\mathbf{X}, \mathbf{Y})) = E(E(f(\mathbf{X}, \mathbf{Y})|\mathbf{X})) = E(h(\mathbf{X}))$. The method is called conditioning, or conditional Monte Carlo. The main requirement for conditioning is that we must be able to compute $h(\cdot)$. Note that

$$\text{Var}(\hat{\mu}_{\text{cond}}) = \frac{1}{n} \text{Var}(h(\mathbf{X})) = \frac{1}{n} \text{Var}(E(f(\mathbf{X}, \mathbf{Y})|\mathbf{X})) \leq \text{Var}(\hat{\mu}),$$

as

$$\text{Var}(f(\mathbf{X}, \mathbf{Y})) = E(\text{Var}(f(\mathbf{X}, \mathbf{Y})|\mathbf{X})) + \text{Var}(E(f(\mathbf{X}, \mathbf{Y})|\mathbf{X})).$$

Conditioning is a special case of derandomization. The function $f(\mathbf{X}, \mathbf{Y})$ has two sources of randomness, \mathbf{X} and \mathbf{Y} . For any given \mathbf{x} and random \mathbf{Y} we replace the random value $f(\mathbf{x}, \mathbf{Y})$ by its expectation $h(\mathbf{x})$, removing one of the two sources of randomness.

1.5 Control Variates

Suppose that we want to find $\mu = E(f(\mathbf{X}))$ and that we know the value $\theta = E(h(\mathbf{X}))$, where $h(\cdot)$ is ‘similar’ to $f(\cdot)$. The precise meaning of ‘similar’ depends on the problem. For a value $\beta \in \mathbb{R}$, the regression estimator of μ is given by

$$\hat{\mu}_{\beta} = \frac{1}{n} \sum_{i=1}^n (f(\mathbf{X}_i) - \beta h(\mathbf{X}_i)) + \beta \theta = \hat{\mu} + \beta (\theta - \hat{\theta}),$$

where $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(\mathbf{X}_i)$ and $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n h(\mathbf{X}_i)$. The variance of the regression estimator is

$$\text{Var}(\hat{\mu}_{\beta}) = \frac{1}{n} (\text{Var}(f(\mathbf{X})) - 2\beta \text{Cov}(f(\mathbf{X}), h(\mathbf{X})) + \beta^2 \text{Var}(h(\mathbf{X}))).$$

Let β_{opt} denote the best value of β in the sense that $\text{Var}(\hat{\mu}_{\beta})$ is minimum at $\beta = \beta_{\text{opt}}$. Then

$$\beta_{\text{opt}} = \frac{\rho\sigma}{\sigma_h} \quad \text{and} \quad \text{Var}(\hat{\mu}_{\beta_{\text{opt}}}) = \frac{\sigma^2}{n} (1 - \rho^2),$$

where $\rho = \text{Corr}(f(\mathbf{X}), h(\mathbf{X}))$, $\sigma^2 = \text{Var}(f(\mathbf{X}))$ and $\sigma_h^2 = \text{Var}(h(\mathbf{X}))$. In the regression estimator, any control variate that correlates with f is helpful, even one that correlates negatively. In

practice we do not know β_{opt} and hence we estimate it by

$$\hat{\beta} = \frac{\sum_{i=1}^n (f(\mathbf{X}_i) - \hat{\mu}) (h(\mathbf{X}_i) - \hat{\theta})}{\sum_{i=1}^n (h(\mathbf{X}_i) - \hat{\theta})^2},$$

In general $E(\hat{\mu}_{\hat{\beta}}) \neq \mu$, but this bias is usually small. The estimated variance of $\hat{\mu}_{\hat{\beta}}$ is

$$\widehat{Var}(\hat{\mu}_{\hat{\beta}}) = \frac{1}{n^2} \sum_{i=1}^n \left(f(\mathbf{X}_i) - \hat{\mu}_{\hat{\beta}} - \hat{\beta} (h(\mathbf{X}_i) - \hat{\theta}) \right)^2.$$

A 99% confidence interval is $\hat{\mu}_{\hat{\beta}} \pm 2.58 \sqrt{\widehat{Var}(\hat{\mu}_{\hat{\beta}})}$.