

Our goal is to estimate expectations of functions

$$E_{p(x)} [f(x)] \approx \sum_i \frac{f(x_i)}{n} \text{ with } x_i \sim p(x) \text{ independent}$$

using monte carlo estimators.

This is unbiased

$$E \left[\sum_i \frac{f(x_i)}{n} \right] = \sum_i \frac{1}{n} E[f(x)] = E[f(x)]$$

$$\text{Var} \left(\sum_i \frac{f(x_i)}{n} \right) = \frac{\text{Var}(f(x))}{n} = \frac{\sigma^2}{n}$$

We can reduce the variance by increasing n . Can we do better?

(Caveat: We have to balance the cost of the extra work with the cost of drawing extra samples.)

Two motivating examples

① Reinforcement learning - policy gradients

Take an action according to $\pi(a)$ get stochastic reward $R(a)$, want to maximize

$$E_{\pi} [R(a)] \text{ but } \nabla E_{\pi} [R(a)] = E_{\pi} [R(a) \nabla \log \pi(a)]$$

is high variance. We may be limited in the number of samples.

We can use control variates to reduce the variance of the estimator.

② Suppose you are responsible for Google Home metrics. You collect millions of interactions with the device and want to know what % are "satisfactory". You can evaluate the utterances by having someone listen to each interaction, but this is expensive. With a limited budget, what can you do?

Antithetics

$$\mu = E[f(x)]$$

$$f(x_i) = \mu + (f(x_i) - \mu) = \mu + \epsilon_i \Rightarrow \sum_i \frac{1}{n} f(x_i) = \mu + \sum_i \frac{\epsilon_i}{n}$$

Each ϵ_i has variance σ^2 and the average has variance $\frac{\sigma^2}{n}$ due to independence. If we sample in a coupled fashion we might get even more cancellation.

$$X \xrightarrow{\text{involution}} \tilde{X} \text{ s.t. } f(x) \text{ and } f(\tilde{x})$$

Somehow have opposite errors

For example

$$X \sim N(\mu, \sigma^2) \text{ then } \tilde{X} = \mu - (x - \mu) = 2\mu - x$$

The effectiveness depends on f

$$\begin{aligned} \text{Var}(\hat{\mu}_{anti}) &= \text{Var}\left(\frac{1}{N} \sum_i \frac{f(x_i) + f(\tilde{x}_i)}{2}\right) = \frac{1}{N^2} \frac{\text{Var}(f(x_i)) + \text{Var}(f(\tilde{x}_i)) + 2\text{Cov}(f(x_i), f(\tilde{x}_i))}{2} \\ &= \frac{1}{N^2} (N\sigma^2 + \text{Cov}(f(x), f(\tilde{x}))) = \frac{1}{N} \sigma^2 (1 + \rho) \end{aligned}$$

So it is good if $f(x)$ & $f(\tilde{x})$ are anti-correlated.

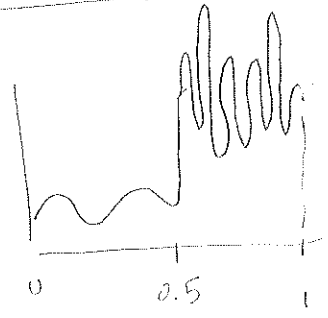
(Q1)
$$f(x) = \frac{f(x) + f(\tilde{x})}{2} + \frac{f(x) - f(\tilde{x})}{2} = f_E(x) + f_O(x)$$

Show that f_E and f_O are orthogonal and that

$$\text{Var}(\hat{\mu}_{anti}) = 2 \frac{\text{Var}(f_E(x))}{n} \text{ whereas } \text{Var}(\hat{\mu}) = \frac{\text{Var}(f_E) + \text{Var}(f_O)}{n}$$

Sampling x & \tilde{x} can sometimes be much cheaper than sampling x twice.

Stratification



might estimate the function in two pieces and possibly weight the right piece more

Partition space D_1, \dots, D_J w/ $P(X \in D_j) = w_j$, $P_j(x) = P(X \in D_j)$

Then sample n_j samples $X_{ij} \sim P_j$ and form

$$\hat{\mu}_{\text{strat}} = \sum_{j=1}^J \frac{w_j}{n_j} \sum_{i=1}^{n_j} f(X_{ij})$$

Check: $\hat{\mu}_{\text{strat}}$ is unbiased

$$\text{Var}(\hat{\mu}_{\text{strat}}) = \text{Var}\left(\sum_{j=1}^J \frac{w_j}{n_j} \sum_{i=1}^{n_j} f(X_{ij})\right) = \sum_{j=1}^J \frac{w_j^2}{n_j^2} n_j \sigma_j^2 = \sum_{j=1}^J \frac{w_j^2}{n_j} \sigma_j^2$$

Proportional allocation says $n_j = w_j n$, so

$$\text{Var}(\hat{\mu}_{\text{prop}}) = \sum_{j=1}^J \frac{w_j}{n} \sigma_j^2$$

Using the law of variance

$$\text{Var}(f(x)) = E[\text{Var}(f(x|y))] + \text{Var}(E[f(x|y)])$$

$$\sigma^2 = \underbrace{\sum_{j=1}^J w_j \sigma_j^2}_{\text{within strata}} + \underbrace{\sum_{j=1}^J w_j (\mu_j - \mu)^2}_{\text{between strata}}$$

So $\text{Var}(\hat{\mu}_{\text{prop}}) \leq \text{Var}(\hat{\mu})$ allowing us to drop the between strata variance term

Optimal allocation is possible, but tricky because σ_j^2 is unknown

Q1 for the metrics problem, how could you use stratified sampling? Make additional reasonable assumptions.

Q2 Can you construct $f = f_B + f_w$ similar to the antithetics where f_B contains the between strata variance and f_w contains the within strata variance

Common Random Numbers

(4)

Suppose we want

$$E[f(x) - g(x)] \text{ for } f \text{ and } g \text{ related}$$

eg., $f(x) = h(x, \theta)$ and $g(x) = h(x, \tilde{\theta})$

Then $E[f(x)] - E[g(x)] = E[f(x) - g(x)]$

Which is better?

$$\text{Var}(\hat{\mu}_{\text{common}}) = \frac{\text{Var}(f(x)) + \text{Var}(g(x)) - 2\text{cov}(f(x), g(x))}{n}$$

$$\text{Var}(\hat{\mu}_{\text{indep}}) = \frac{\text{Var}(f(x)) + \text{Var}(g(x))}{n}$$

If f and g are closely related, this is great. Also useful for

$$E[f(x) - f(\tilde{x})] \text{ if we can reparameterize from a}$$

common source of randomness.

eg. $X \sim N(\mu, \sigma^2)$ then $Z \sim N(0, 1)$ $X = \mu + \sigma Z$
 $\tilde{X} \sim N(\tilde{\mu}, \tilde{\sigma}^2)$ $\tilde{X} = \tilde{\mu} + \tilde{\sigma} Z$

For example,

$$\nabla_{\theta} E_{\pi(a, \theta)} [R(a)] \approx \frac{E_{\pi(a, \theta + \varepsilon)} [R(a)] - E_{\pi(a, \theta - \varepsilon)} [R(a)]}{2\varepsilon}$$

can get a much lower variance estimator w/ common random numbers.

Conditioning

If we want

$$E[f(x,y)] = E_x \left[\underbrace{E_y[f(x,y)|x]}_{h(x)} \right] = E_x[h(x)]$$

Now

$$\text{Var}(h(x)) = \text{Var}_x(f(x)) - E_x[\text{Var}_y(f(x,y)|x)] \geq \text{Var}_x(f(x))$$

Caution: If $h(x)$ is expensive, then may not be useful.

① Suppose $\nabla E[f(h)]$ $h = H(z)$ $z = \log \alpha + \log u - \log(1-u)$
 $u \sim U(0,1)$

Logistic random variable

Now

$$E_h[f(h) \nabla \log p(h)] = \nabla E_h[f(h)] = \nabla E_z[f(H(z))] = E_z[f(H(z)) \nabla \log p(z)]$$

which is better?

② Suppose $\pi(a) = \pi(a_1) \pi(a_2|a_1) \dots \pi(a_n|a_{1:n-1})$

Compute an estimator for

$$\nabla E_\pi[\log \pi]$$

That does 1 step of conditioning - i.e. local conditional entropy.

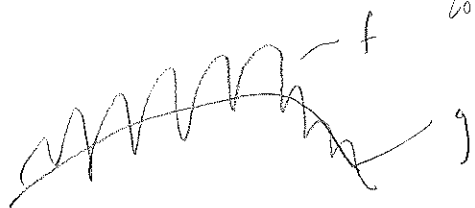
Control Variables - allowing us to use closed form solutions

(6)

We note

$$E[f(x)] = E[\underbrace{f(x) - g(x)}_{\text{complex residual}}] + \underbrace{E[g(x)]}_{\text{helpful}}$$

if this is tractable, this can be helpful



$$\hat{\mu}_{cv} = \frac{1}{n} \sum_{i=1}^n f(x_i) - \eta(g(x_i) - \bar{g}) + \bar{g}$$

check: unbiased

$$\text{Var}(\hat{\mu}_{cv}) = \frac{1}{n} (\text{Var}(f) + \eta^2 \text{Var}(g) - 2\eta \text{Cov}(f, g))$$

$$\text{So } \eta_{\text{opt}} = \frac{\text{Cov}(f, g)}{\text{Var}(g)}$$

$$\text{Thus } \text{Var}(\hat{\mu}_{cv}) = \frac{1}{n} \left(\text{Var}(f) + \frac{\text{Cov}(f, g)^2 \text{Var}(g)}{\text{Var}(g)^2} - 2 \frac{\text{Cov}(f, g)^2}{\text{Var}(g)} \right)$$

$$= \frac{1}{n} \left(1 - \frac{\text{Cov}(f, g)^2}{\text{Var}(g) \text{Var}(f)} \right) \text{Var}(f)$$

So any cv is helpful if we can compute η optimally.

Note that

$$\text{Var}(\hat{\mu}_{cv}) = E \left[\left(\frac{1}{n} \sum_{i=1}^n f(x_i) - \eta(g(x_i) - \bar{g}) \right)^2 \right] - E \left[\frac{1}{n} \sum_{i=1}^n f(x_i) \right]^2 \rightarrow \text{constant in } \eta$$

So the optimal η minimizes a least squares problem. Can estimate η_{opt} by least squares (or leave one out least squares for unbiasedness).

(Ex 1) In the problem

$$E_{\pi} [R(a) \nabla \log \pi]$$

$\nabla \log \pi$ is a control variate because $E[\nabla \log \pi] = \nabla E_{\pi} [1] = 0$

So $E_{\pi} [(R(a) - b) \nabla \log \pi]$ is unbiased, b could be estimated from minibatch statistics.

Q Prop

Policy gradients in continuous control

$$E_s E_{\pi} [\hat{Q} \nabla \log \pi] \quad \text{this is the key quantity to compute (on policy)}$$

\hat{Q} = sample of discounted returns.

$$\textcircled{1} \quad Q^\pi(a, s) = r(a, s) + \gamma E_\pi [Q^\pi(a', s)]$$

can be estimated off policy $\Rightarrow Q_w(a, s)$ parameterized critic

$\textcircled{2}$ How can we use $Q_w(a, s)$ to reduce variance?

MuProp showed any linear function of a can be used

$$E_{a \sim \pi} [(ma + b) \nabla \log \pi(a)] = m E_{a \sim \pi} [a \nabla \log \pi(a)] = m \nabla E_{a \sim \pi} [a]$$

as long as this is tractable
eg. Gaussian policy

$$m = Q'_w(\bar{a}, s)$$

So, we get

$$E_s \left[E_{a \sim \pi} \left[(\hat{Q} - \eta(Q_w(\bar{a}, s) + (a - \bar{a})Q'_w(\bar{a}, s))) \nabla \log \pi \right] + \eta Q'_w(\bar{a}, s) \nabla E_\pi [a] \right]$$

Alternatively

$$E_s \left[E_{a \sim \pi} \left[(Q(a, s) - \eta(Q_w(a, s))) \nabla \log \pi \right] + \eta \nabla E_\pi [Q_w(a, s)] \right]$$

tractable if π is reparameterizable

$\textcircled{3}$ Can you use a linear approximation to reduce the variance further?

$$\nabla E_\pi [Q_w(a, s)] = \nabla E_\xi [Q_w(a(\xi, \theta), s)]$$

$$= E_\xi [\nabla Q_w(a(\xi, \theta), s)]$$

$$= E_\xi [Q'_w(a, s) |_{a=a(\xi, \theta)} \nabla a(\xi, \theta)]$$

$$= Q'_w(\bar{a}, s) \nabla a(\xi, \theta)$$

$\eta(s)$ can be learned online w/ least squares.

⑧

Ⓔ How can you use a model-based system to reduce the variance of the model free policy gradient?

Rebar

h is discrete. Want to compute

$$E [f(h) \nabla \log p(h)]$$

$$= \nabla E_p [f(h)] = \nabla \left(E_h [f(h) - \eta E_{z|h} [f(\sigma_t(z))]] + \eta E_z [f(\sigma_t(z))] \right)$$

$$= E_h [(f(h) - \eta E_{z|h} [f(\sigma_t(z))]) \nabla \log p(h) - \eta \nabla E_{z|h} [f(\sigma_t(z))]] + \eta \nabla E_z [f(\sigma_t(z))]$$

reparameterizable

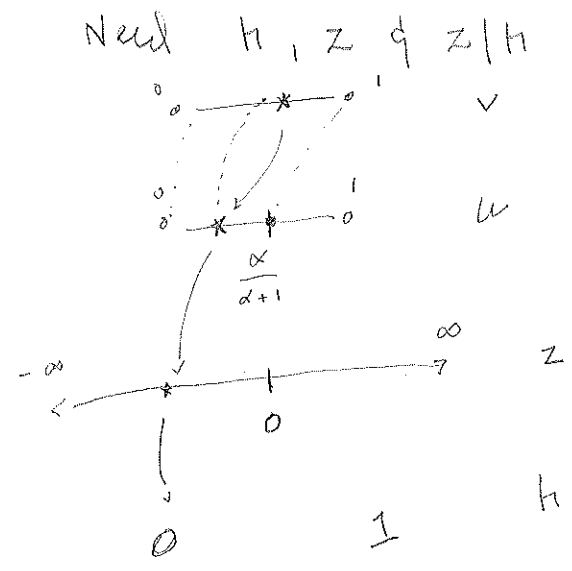
$$Z = \log \alpha + \log U - \log(1-U) \quad U \sim \text{Unif}(0,1)$$

$$h = H(z) \quad \text{so that } h \sim \text{Bern} \left(\frac{\alpha}{\alpha+1} \right)$$

σ_t is a tempered sigmoid that approximates $H(z)$

$$h = H(z) \approx \sigma \left(\frac{z}{T} \right) = \sigma_t(z)$$

How can we use common random numbers?



(21) Why is this okay?