# Importance Sampling

Problem: Estimate

$$E(\mathcal{O}) = \int \mathrm{d}x \; p(x)\mathcal{O}(x)$$

when $\mathcal{O}$ is only non-zero where $p(x)$ is very small.

What is the problem? Suppose

$$\int \mathrm{d}x \; p(x) \; \mathcal{I}_{\mathcal{O}>0} < 10^{-8}.$$

Then for numerical simulations, we have to draw $\approx 10^8$ samples from $p(x)$ to get even one sample where the expectation is non-zero.

Suppose

$$\gamma \equiv \int dx \; p(x) \; \mathcal{I}_{\mathcal{O}}.$$

Draw $N$ samples from $p(x)$. Then the estimate

$$\gamma_N = \frac{1}{N} \sum_i \mathcal{I}_{\mathcal{O}}(x_i)$$

has variance

$$Var(\gamma_N) = \frac{\gamma(1-\gamma)}{N}.$$

Confidence interval for $\gamma$ estimate (CLT!):

$$\gamma_N \pm z_{\alpha/2}\sqrt{Var(\gamma_N)}.$$

For example, if $\alpha = 0.01$, the confidence interval for $\gamma < 0.1$ is

$$2.576 \times \sqrt{Var(\gamma_N)} < 0.1\gamma_N.$$

How big must $N$ be?

$$N \approx \left(\frac{2.576}{0.1}\right)^2 \times \sqrt{\frac{1-\gamma}{\gamma}}.$$

If $\gamma \approx 10^{-6}$ we need almost $7 \times 10^8$ samples!

Basic idea of importance sampling

$$\gamma \equiv \int dx\ p(x)\ \mathcal{I}_{\mathcal{O}} = \int dx\ q(x)\ \frac{p(x)}{q(x)}\ \mathcal{I}_{\mathcal{O}} = E_q\Big(\frac{p(x)}{q(x)}\ \mathcal{I}_{\mathcal{O}}\Big).$$

Pick $q(x)$ so that $\mathcal{I}_{\mathcal{O}}(x) = 1$ is *not* rare in the distribution $q(x)$.

Evaluate the new observable $\mathcal{O}_q \equiv \frac{p(x)}{q(x)}\ \mathcal{I}_{\mathcal{O}}$.

Intuitively, maximize $q(x)$ where $\mathcal{I}_{\mathcal{O}}(x) > 0$.

*Any q* leads to the same expectation with more draws.
So pick $q$ to minimize the variance of $\mathcal{O}_q$.

$$\gamma \equiv \int \mathrm{d}x \; p(x) \; \mathcal{I}_{\mathcal{O}} = \int \mathrm{d}x \; q(x) \; \frac{p(x)}{q(x)} \; \mathcal{I}_{\mathcal{O}} = E_q\Big(\frac{p(x)}{q(x)} \; \mathcal{I}_{\mathcal{O}}\Big).$$

Pick $q(x)$ so that $\mathcal{I}_{\mathcal{O}}(x) = 1$ is *not* rare in the distribution $q(x)$.

Evaluate the new observable $\mathcal{O}_q \equiv \frac{p(x)}{q(x)} \; \mathcal{I}_{\mathcal{O}}$.

$$E_q\Big(\mathcal{O}_q^2\Big) = \int \mathrm{d}x \; q(x) \Big(\frac{p(x)}{q(x)}\Big)^2 \mathcal{I}_{\mathcal{O}}(x) = E_p(\mathcal{O}_q).$$

Does this always work?

NO!

$$p(x) = \lambda \exp(-\lambda x), \quad x \geq 0.$$

We want to estimate $\mathcal{I}_{x>y}$ for some large $y$.

Try
$$q(x) = \mu \exp(-\mu x)??$$

Then
$$E_p\left(\mathcal{I}_{x>y}\frac{p}{q}\right) = \frac{\lambda^2}{\mu}\int_y^\infty dx \ \exp(-(2\lambda - \mu)x)$$

$$= \frac{\lambda^2}{\mu(2\lambda - \mu)}\left(\exp(-(2\lambda - \mu)y - \exp(-(2\lambda - \mu)L)\right).$$

No minimum at any finite value of $\mu$.

Now, let's apply this IS measure redefinition to actually prove Cramér's theorem. Remember, I didn't actually prove Gärtner-Ellis, I just motivated it, and then cheated by using GE to prove Cramér.

## Cramér's Theorem, revisited

Remember Chebyshev. Redo the proof:

$$P(\frac{1}{n}\sum_i X_i > \mu + \epsilon) = P(\sum_i X_i - n\mu > n\epsilon) \geq \frac{E((\sum_i X_i - n\mu)^2)}{n^2\epsilon^2} = \frac{\sigma^2}{n\epsilon^2},$$

and the same argument holds for $P(-\frac{1}{n}\sum_i X_i > \mu + \epsilon)$.

Now, the same argument applies if we take $\psi$ to be any non-negative increasing function:

$$P(\frac{1}{n}\sum_i X_i > s) \leq P(\psi(\sum_i X_i) \geq \psi(ns)) \leq \frac{1}{\psi(ns)}E(\psi(\sum_i X_i)).$$

Take $\psi(s) \equiv \exp(\lambda s)$. Then

$$\frac{1}{n}P(\frac{1}{n}\sum_i X_i > s) \leq -\lambda s + \frac{1}{n}\sum_i \ln E(\exp(\lambda X_i)).$$

Recognize the cumulant generating function? So

$$-\frac{1}{n}P(\frac{1}{n}\sum_i X_i > s) \geq \sup(\lambda s - K_X(\lambda)).$$

# Cramér's Theorem: Lower bound

Define a new probability measure

$$q(X) \equiv p(X) \frac{\exp(\lambda X)}{M_X(\lambda)}.$$

This is still normalized:

$$E_q(1) = \frac{E_p(\exp(\lambda X))}{M_X(\lambda)} = 1.$$

Now

$$\frac{1}{n} \ln P(\frac{1}{n} \sum X_i \geq s) \geq \frac{1}{n} \ln P(s + \epsilon > \frac{1}{n} \sum X_i \geq s)$$

$$= \frac{1}{n} \ln E_q \left( \exp(-n\lambda X) M_X^n(\lambda) \mathcal{I}(s + \epsilon > \frac{1}{n} \sum X_i \geq s) \right)$$

$$\geq -\lambda(s + \epsilon) + K_X(\lambda) + \frac{1}{n} \ln P_q(s + \epsilon > \frac{1}{n} \sum X_i \geq s).$$

How do we choose $\lambda$ :

$$\lim_{\epsilon \downarrow 0} \ln P_q(s + \epsilon > \frac{1}{n} \sum X_i \geq s) \to 0?$$

Note $K'_X(\lambda) = E_q(X) = E(X \exp \lambda X) K_X(\lambda)^{-1}$,

$$K'_X(\lambda = 0) = \mu,$$

and

$$K'_X(\lambda \to \infty) = \text{maximum of support of } X.$$

So by the mean-value theorem there is some $\lambda(s) \in [0, \infty)$ such that

$$K'_X(\lambda(s)) = s + \frac{\epsilon}{2}.$$

So, finally,

$$-\frac{1}{n} \ln P(\frac{1}{n} \sum X_i \geq s) \leq (\lambda(s)s - K_X(\lambda(s))).$$

Let's use this probability 'tilting' to make ML maximization trivial.

# Mayer cluster expansions

For a single spin, taking values in $\{\pm 1\}$, the partition function is

$$Z = \frac{1}{2}\left[\exp(-\epsilon(w + b)) + \exp(-\epsilon(-w + b))\right]$$

where $N$ is the number of configurations and $\epsilon$ is the inverse temperature. Write $\exp(-\epsilon V_{ij}) = 1 - (1 - \exp(-\epsilon V_{ij})) \equiv 1 + \Delta_{ij}$. Then

$$Z = \frac{1}{N}\sum_{\text{configs}}\prod_{i<j}(1 + \Delta_{ij}).$$

Expand in $\Delta$ :

$$Z = 1 + \frac{1}{N}\sum_{\text{configs}}\left[\sum_{i<j}\Delta_{ij} + \sum_{i<j}\sum_{k<l}\Delta_{ij}\Delta_{kl} + \dots\right].$$

In usual statistical mechanics situations, this cluster expansion becomes a more or less geometric expansion because the constituents might be molecules so there are constraints on what configurations can contribute to each separate sum.

For us, things will be much simpler but it's good to know the general idea.

For a single spin, taking values in $\{\pm 1\}$, the partition function is

$$Z = \frac{1}{2}\big[\exp(-\epsilon w) + \exp(-\epsilon(-w))\big]$$

so we get

$$Z = \frac{1}{2}\big[(1 + \Delta_+) + (1 + \Delta_-)\big] = 1 + \frac{1}{2}(\Delta_+ + \Delta_-).$$

For $\epsilon$ small we can expand the $\Delta$ terms if we want:

$$Z \approx 1 - \frac{\epsilon}{2}(w - w) + \frac{\epsilon^2}{4}(w^2 + (-w)^2) = 1 + \frac{\epsilon^2}{2}(w^2).$$

Take-away message: In this limit $Z$ is trivial to calculate.

For two interacting spins, taking values in $\{\pm 1\}$, the partition function is

$$Z = \frac{1}{4} \sum_{\sigma_1, \sigma_2} \left[ \exp(-\epsilon w^i \sigma_i + w^3 \sigma_1 \sigma_2) \right].$$

Now let's be more organized about this expansion. Note (for any variable $\sigma = \pm 1$)

$$\exp(w\sigma) = \cosh w\sigma + \sinh w\sigma = \cosh w + \sigma \sinh w = \cosh w[1 + \sigma \tanh w].$$

So now the cosh terms don't depend on $\sigma$ values at all so they can come out of the sum over configurations.

$$Z = \frac{1}{4} \cosh(w^1) \cosh(w^2) \cosh(w^3) \sum_{\sigma_1, \sigma_2} (1 + \sigma_1 \tanh w^1)(1 + \sigma_2 \tanh w^2)(1 + \sigma_1 \sigma_2 \tanh w^3).$$

The only terms that will survive in the sum over configurations will be terms that have no linear $\sigma$ dependence on either of the $\sigma_i$, for example, $\tanh w^1 \tanh w^2 \tanh w^3$. For $\epsilon$ small, each $w$ comes with a factor of $\epsilon$ and $\tanh x \approx x$ for small $x$. Takeaway message: Even with interacting spins, $Z$ is very simple to calculate, especially for small $\epsilon$.

Let's go back to Lecture 1, where we saw

$$\frac{\partial D_{KL}(f||p)}{\partial \theta_i} = E(\mathcal{O}_i) - E_f(\mathcal{O}_i)$$

and we said that

$$E(\mathcal{O}_i) = \partial_{\theta_i} \ln Z(\theta)$$

was difficult to calculate because of $Z$ summing over all configurations.

Was I lying then or am I lying now??

The problem is that we have no idea of the magnitude of the interactions of the spins. In other words, you cannot change the temperature of a system if you have no idea what temperature it's at. All we have is some observed data.

Be concrete: RV: spins $\sigma_i \in \{\pm 1\}, i = 0, \ldots, N-1$.

Spin 'interactions': $\mathcal{O}_I \in \{\sigma_i, \sigma_i \sigma_j, \ldots\}$.

Want to assign probabilities for any realization of the RV:

$$P(\sigma|w) \equiv \frac{\exp\left(w^I \mathcal{O}_I(\sigma)\right)}{Z(w)},$$

and $ML$ maximization

$$\implies \boxed{w^* : P(\sigma_{\text{obs}}|w^*) = f(\sigma_{\text{obs}}) \equiv \frac{n(\sigma_{\text{obs}})}{\sum_{\text{obs}} n(\sigma_{\text{obs}})}.}$$

$$D_{KL}(f_{\text{obs}}||p) \equiv \sum_{\text{obs}} f_{\text{obs}} \ln\left(\frac{f_{\text{obs}}}{p_{\text{obs}}}\right)$$

# Where's the temperature?

Minimizing

$$D_{KL}(f_{\mathrm{obs}}\|p) \equiv \sum_{\mathrm{obs}} f_{\mathrm{obs}} \ln \left( \frac{f_{\mathrm{obs}}}{p_{\mathrm{obs}}} \right)$$

by gradient descent:

$$\delta w^I = \alpha \left[ E_f(\mathcal{O}_I) - E_p(\mathcal{O}_I) \right],$$

where (as before)

$$E_q(\mathcal{O}_I) \equiv \sum_{\mathrm{configs}\ \sigma_i = \pm 1} q(\sigma)\mathcal{O}_I(\sigma).$$

In the small $\epsilon$ limit,

$$Z(\epsilon w^I) = 1 + \frac{\epsilon^2}{2} \sum_I (w^I)^2 + \text{h.o. terms.}$$

So in this limit

$$E_{p_\epsilon}(\mathcal{O}_I) = \epsilon w^I,$$

where we define

$$p_\epsilon(\sigma) \equiv \frac{p(\sigma)p(\sigma)^{\epsilon-1}}{\sum_{\text{configs}}p(\text{config})p(\text{config})^{\epsilon-1}}.$$

How is this useful at all? What happened to ensuring that $f_{\text{obs}} = p_{\text{obs}}$?

Redefine

$$f_\epsilon(\text{obs}) = \frac{f(\text{obs})p(\text{obs})^{\epsilon-1})}{\sum_{\text{all observations}} f(\text{obs})p(\text{obs})^{\epsilon-1})}.$$

Now minimizing

$$D_{KL}(f_\epsilon||p_\epsilon)$$

still gives

$$f(\text{obs}) = p(\text{obs}),$$

but gradient descent is

$$\delta w^I \propto \left[E_{f_\epsilon}(\mathcal{O}_I) - E_{p_\epsilon}(\mathcal{O}_I)\right] = \left[E_{f_\epsilon}(\mathcal{O}_I) - \epsilon w^I\right].$$

Why is this useful??

The sum in $E_{f_\epsilon}$ is only over observations! NO sum over all configurations.

Hopfield solution: $\epsilon = 1 \implies w^I = E_f(\mathcal{O}_I)$

Analogy with damping term in a differential equation:

$$\frac{\mathrm{d}w^I}{\mathrm{d}t} = E_{f_\epsilon}(\mathcal{O}_I) - \epsilon w^I.$$

# Learning during Iteration



(a)

Mean squared error

0.010

0.005

0.000

Legend:
- $\varepsilon = 0.1$ (black dashed)
- $\varepsilon = 0.5$ (red solid)
- $\varepsilon = 0.8$ (blue dashed)

Iterations: 0, 25, 50, 75, 100

(e)

Mean squared error

0.010

0.005

0.000

Legend:
- $L = 5,000$ (black solid, circle)
- $L = 10,000$ (blue dashed, square)

$\varepsilon$: 0.00, 0.25, 0.50, 0.75, 1.00

M=20

$$MSE = \sqrt{\sum_I \left(w_I - w_I^{true}\right)^2}$$

$D_{KL}(\widetilde{f}||\widetilde{p})$

$$\langle E \rangle_f = \left\langle \sum_I w_I O_I(\sigma) \right\rangle_f$$

Large sample

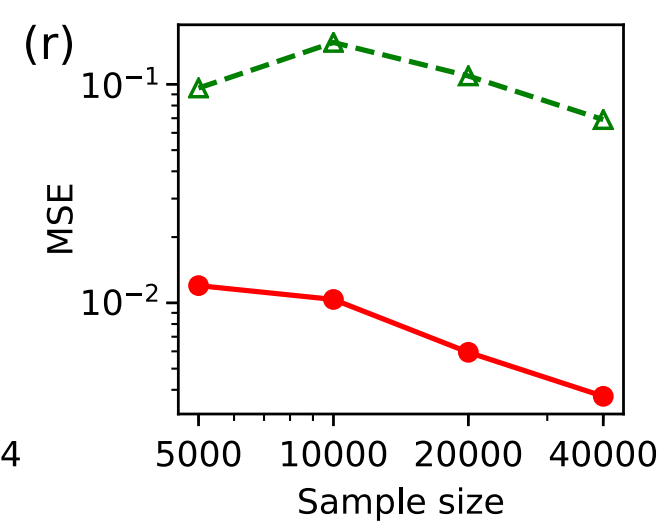Small sample

(a) HF
MLE
PLE
EM

weak coupling

$M = 20$

Inferred couplings

Actual couplings

(b)

(c)

MSE

Sample size

strong coupling

(d)

(e)

(f)

MSE

Sample size

(a) M=20

Inferred interactions vs Actual interactions
- Hopfield (○)
- Boltzmann (□)
- $\varepsilon$-machine (●)

(b)
Computational time vs System size
- Boltzmann (□)
- $\varepsilon$-machine (●)

(c) M=100

Inferred interactions vs Actual interactions
- Hopfield (○)
- $\varepsilon$-machine (●)
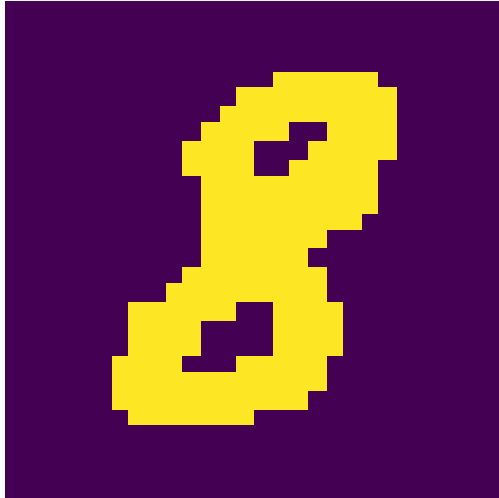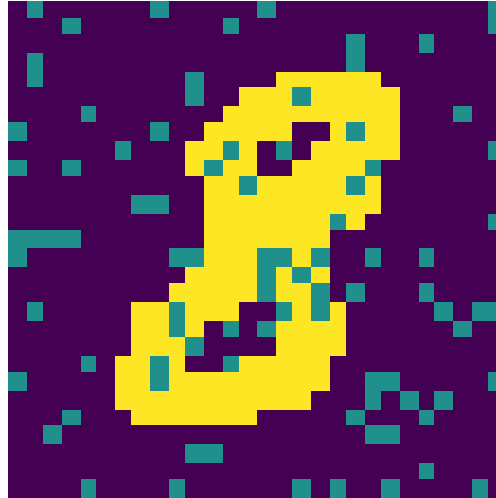
(d)
Histogram vs Energy per spin
- random (□)
- observed (●)

(a) original image    (b) noisy image    (c) recovered image