

Computational Statistics

11/24/21

More sampling...

Recall Bayesian modeling:

- ① We model data X_1, \dots, X_n as having come from some distribution F with density $f = f(x_1, \dots, x_n; \theta_1, \dots, \theta_k)$ where $\theta_1, \dots, \theta_k$ are parameters defining the distribution. (which we want to estimate)
- ② We specify a prior distribution on $\vec{\theta}$, $\propto f(\vec{\theta})$, which reflects our beliefs (or lack thereof) on the value of $\vec{\theta}$.
- ③ We compute the posterior, conditional distribution, of $\vec{\theta} | \vec{X}$ as:

$$p(\vec{\theta} | \vec{X}) = p(\theta_1, \dots, \theta_k | X_1, \dots, X_n) = \frac{f(x_1, \dots, x_n | \theta_1, \dots, \theta_k) f(\vec{\theta})}{\int f(x_1, \dots, x_n | \theta_1, \dots, \theta_k) f(\vec{\theta}) d\vec{\theta}}$$

↑
posterior.

$\propto \underset{\substack{\uparrow \\ \text{likelihood}}}{\mathcal{L}(\vec{\theta})} f(\vec{\theta}) \underset{\substack{\uparrow \\ \text{prior}}}{f(\vec{\theta})}$

Bayes + MCMC

Recall the Metropolis - Hastings algorithm:

With $q(y|x)$ a proposal density, construct a M. chain for sampling from f as follows:

$$\text{set } X_{i+1} = \begin{cases} Y & \text{with probability } r(X_i, Y) \\ X_i & \text{with probability } 1 - r(X_i, Y) \end{cases}$$

$$\text{where } r(x, y) = \min \left\{ \frac{f(y)}{f(x)} \frac{q(x|y)}{q(y|x)}, 1 \right\}$$

If we wish to sample from the posterior $p(\theta|x)$ [Recall, x is the observed data], then we wish

$$p(\theta|x) = \frac{I(\theta) f(\theta)}{C(x)}$$

to construct a chain $\theta_1, \theta_2, \dots, \theta_i, \dots$ when

$$\theta_{i+1} = \begin{cases} Y & \text{with prob } r(\theta_i, Y) \\ \theta_i & \text{with prob } 1 - r(\theta_i, Y) \end{cases}$$

$$\text{and } r(\theta_i, Y) = \min \left\{ \frac{p(Y|x) q(\theta_i|Y)}{p(\theta_i|x) q(Y|\theta_i)}, 1 \right\}$$

$$= \min \left\{ \frac{I(Y) f(Y) / C(x)}{I(\theta_i) f(\theta_i) / C(x)} \frac{q(\theta_i|Y)}{q(Y|\theta_i)}, 1 \right\}$$

$$= \min \left\{ \frac{L(Y) f(Y)}{L(\theta_i) f(\theta_i)} \frac{q(\theta_i | Y)}{q(Y | \theta_i)}, 1 \right\}$$

Note: You do not need to compute the normalizing constant $C(x)$ in order to sample from p using MCMC!

Furthermore: If the prior f was uninformative, i.e. flat, $f(x) = C$ and therefore we have

$$r(\theta_i, Y) = \min \left\{ \frac{L(Y) q(\theta_i | Y)}{L(\theta_i) q(Y | \theta_i)}, 1 \right\}.$$

This works if θ is a scalar or vector $\vec{\theta}$.

In the case of $\vec{\theta} \in \mathbb{R}^k$, $\vec{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$ (NOTATION consistent with Efron & Tibshirani.)

We can use Gibbs Sampling.

Idea: Reduce the generation of $\vec{\theta}$ down to a series of univariate calculations.

Let $\vec{\theta}^{(k)} = (\theta_1, \dots, \theta_{k-1}, \theta_{k+1}, \dots, \theta_k) \in \mathbb{R}^{k-1}$ and

$g^{(k)}$ be conditional density of $\theta_k | \vec{\theta}^{(k)}, \vec{x}$ $\hat{=}$ the data

$$\theta_k | \vec{\theta}^{(k)}, \vec{x} \sim g^{(k)}(\theta_k | \vec{\theta}^{(k)}, \vec{x}).$$

Given $\vec{\theta}^{(0)}, \dots, \vec{\theta}^{(i)}$, to generate the components of $\vec{\theta}^{(i+1)}$ merely draw from

$$\theta_k^{(i+1)} \sim q_k(\theta_k | \vec{\theta}_{(-k)}^{(i)}, \vec{y}) \quad \text{for } k=1, \dots, K.$$

Then use the previous step's values of $\theta_j^{(i)}$ to condition and draw $\theta_k^{(i+1)}$.

(See Efron & Tibshirani §13.4 for worked example.)

Just a note: Imagine we have $x_1, \dots, x_n \sim \text{i.i.d. } N(\mu, \tau)$

and we wish to sample from the posterior for μ, τ .

We aim to generate a sequence of $\mu^{(1)}, \tau^{(1)}, \mu^{(2)}, \tau^{(2)}, \dots$

$$p(\mu, \tau | \vec{x}) \propto \mathcal{L}(\mu, \tau) f(\mu, \tau)$$

And to sample, we need

$$f(x, y) = f(x|y) f(y)$$

$$p(\mu | \tau, \vec{x}) = \frac{p(\mu, \tau | \vec{x})}{p(\tau)}$$

$$p(\tau | \mu, \vec{x}) = \frac{p(\mu, \tau | \vec{x})}{p(\mu)}$$

need to compute them,
or be smart about
choosing f as conjugate
prior, etc...

M-H is more general algorithm...

Monte Carlo Variance Reduction

For any MC method used to estimate

$$\mathbb{I} = \int f(x) dx, \quad \text{the variance scales as } \text{Var}(\hat{\mathbb{I}}) = \frac{C}{n}.$$

All one can do is make C smaller.

Analytic Transformations

Consider computing $P(X > 2)$ with $X \sim \text{Cauchy}$.

$$\begin{aligned} P(X > 2) = \mathbb{I} &= \int_2^{\infty} \frac{1}{\pi(1+x^2)} dx \\ &= \int_{-\infty}^{\infty} \mathbb{1}_{(x>2)} \frac{1}{\pi(1+x^2)} dx \end{aligned}$$

$$\Rightarrow \hat{\mathbb{I}} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{(x_i > 2)}$$

where x_i is a draw from Cauchy distribution.

Alternatively: $\mathbb{I} = \int_0^{1/2} \frac{1/y^2}{\pi(1+1/y^2)} dy$ (by change of variables)

Let $\tilde{\mathbb{I}} = \frac{1}{n} \sum_{i=1}^n \frac{1/y_i^2}{\pi(1+1/y_i^2)}$ where $y_i \sim U(0, 1/2)$

\Rightarrow one can show $\frac{\text{Var}(\tilde{\mathbb{I}})}{\text{Var}(\hat{\mathbb{I}})} \approx .001$

Analytic methods for reducing the variance are the best!

Importance Sampling

Any integral $I = \int f(x) dx$ can be rewritten as:

$$I = \int \frac{f(x)}{p(x)} p(x) dx$$

└ probability density $\int p(x) dx = 1$,
 $p \geq 0$.

Then $\hat{I} = \frac{1}{N} \sum_{j=1}^N \frac{f(x_j)}{p(x_j)}$ when $x_i \sim \text{iid}$ from p .

Intuitively, to compute $\int f$ we should sample where $|f|$ is large \rightarrow how should we choose p ? \hat{I} is clearly unbiased, we should minimize $\text{Var}(\hat{I})$.

$$\text{Var}(\hat{I}) = \frac{1}{N} \text{Var}\left(\frac{f(x)}{p(x)}\right) \text{ where } X \sim p.$$

$$\text{Var}\left(\frac{f(x)}{p(x)}\right) = \mathbb{E}\left(\frac{f^2(x)}{p^2(x)}\right) - \mathbb{E}\left(\frac{f(x)}{p(x)}\right)^2$$

$$\mathbb{E}\left(\frac{f(x)}{p(x)}\right)^2 = \left(\int \frac{f(x)}{p(x)} p(x) dx\right)^2 = \underbrace{\left(\int f(x) dx\right)^2}_{\text{does not depend on choice of } p!}$$

\Rightarrow \parallel

Consider the first term — by Jensen's Inequality:

$$\mathbb{E} \left(\frac{f^2(x)}{p(x)} \right) \geq \mathbb{E} \left(\frac{|f(x)|}{p(x)} \right)^2 \\ = \left(\int |f(x)| dx \right)^2$$

$$\Rightarrow \mathbb{E} \left(\frac{f^2(x)}{p^2(x)} \right) \geq \left(\int |f(x)| dx \right)^2$$

This bound can be achieved by choosing:

$$p(x) = \frac{|f(x)|}{\int |f(x)| dx}$$

$$\text{then } \mathbb{E} \left(\frac{f^2(x)}{p^2(x)} \right) = \mathbb{E} \left(\frac{f^2(x)}{|f(x)|^2} \left(\int |f(x)| dx \right)^2 \right)$$

$$= \left(\int |f(x)| dx \right)^2 \mathbb{E}(1) = \left(\int |f(x)| dx \right)^2$$

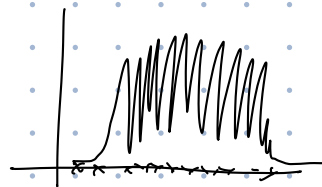
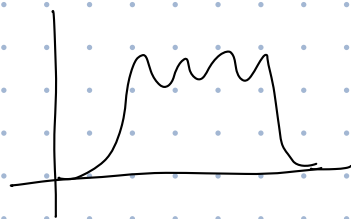
So choose p close to $|f|$, i.e. so that $\frac{|f|}{p} \approx 1$.

Density Estimation

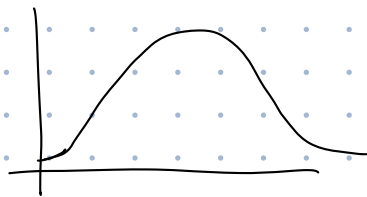
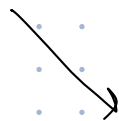
Setup: Observe data $X_1, \dots, X_n \sim F$, and the density is $f = F'$.

Goal: Estimate f using as few assumptions as possible.

\Rightarrow Still a smoothing problem:



undersmoothed estimate



oversmoothed

One possible measure of the error is the L^2 error:

$$\begin{aligned} \text{Loss} = L &= \int (\hat{f}(x) - f(x))^2 dx \\ &= \|\hat{f} - f\|^2 \\ &= \underbrace{\|\hat{f}\|^2 - 2(\hat{f}, f)}_{\mathcal{J}} + \|f\|^2 = \mathcal{J} + C. \end{aligned}$$

$$\begin{aligned} (\hat{f}, f) &= \text{inner product of } \hat{f} \text{ with } f \\ &= \int \hat{f}(x) f(x) dx \\ &= \int \hat{f}(x) dF(x) \\ &= E(\hat{f}(x)) \end{aligned}$$

Goal is to estimate \mathcal{J} .

As before, denote by $\hat{f}_{(-i)}$ the estimator obtained by leaving out x_i :

Def: CV estimate of the risk:

$$\underline{\underline{\hat{J} = \|\hat{f}\|^2 - \frac{2}{n} \sum \hat{f}_{(-i)}(x_i)}}$$

Histograms

Assume we are estimating f on $[0, 1]$, set $h = \frac{1}{m}$, then we have bins $B_1 = [0, h)$, $B_2 = [h, 2h)$, ..., $B_j = [(j-1)h, jh)$.

Denote by $Y_j = \# X_i$'s in bin j .

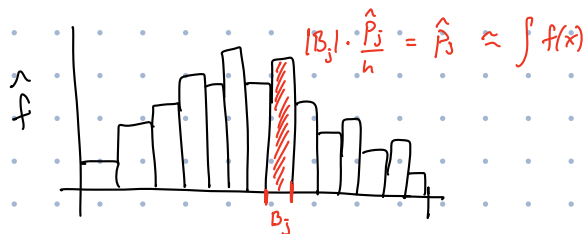
$\hat{p}_j = Y_j/n$. \leftarrow ^{estimate} probability of ending up in bin j

$p_j = \int_{B_j} f(x) dx$ \leftarrow true probability of landing in bin j .
 $= P(X \in B_j)$.

Histogram estimator: $\hat{f}(x) = \sum_{j=1}^m \frac{\hat{p}_j}{h} \mathbb{1}(x \in B_j)$.

\leftarrow Maybe a surprising factor?

Why not just \hat{p}_j ?



$$E(\hat{f}(x)) = \frac{E(\hat{p}_j)}{h} = \frac{p_j}{h} = \frac{1}{h} \int_{B_j} f(x) dx \approx \frac{1}{h} f(x) \cdot h = f(x).$$

Thm $E(\hat{f}(x)) = \frac{p_j}{h}$ for $x \in B_j$

$$\text{Var}(\hat{f}(x)) = \frac{p_j(1-p_j)}{n h^2}$$

and the risk can be computed as:

Then Assume that f' is "absolutely continuous" and

$\int (f')^2 < \infty$, then

$$R(\hat{f}, f) = \frac{h^2}{12} \int (f'(x))^2 dx + \frac{1}{nh} + o(h^2) + o\left(\frac{1}{n}\right)$$

and for fixed n , the minimum occurs at

$$h_* = \frac{1}{n^{1/3}} \left(\frac{6}{\int f'^2 dx} \right)^{1/3} \sim \frac{1}{n^{1/3}} \quad \text{and}$$

then $R(\hat{f}, f) \sim C \frac{1}{n^{2/3}}$.

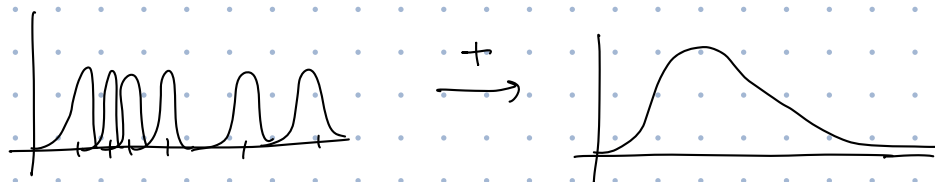
Remember: Since f is not known, minimize \hat{J} instead since it can actually be computed.

Kernel Density Estimator

If you had only one data point, x_i , what would you do?



Idea: Place a local kernel at each data point, and sum:



How wide should the kernel be?

Def: Kernel density estimator:

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x-x_i}{h}\right)$$

$$\int K = 1$$

$$\int x K = 0$$

$$\int x^2 K < \infty$$

Thm If f is continuous at x , and $h \rightarrow 0$, $nh \rightarrow \infty$,
then $\hat{f}(x) \xrightarrow{P} f(x)$:

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\hat{f}(x) - f(x)| > \varepsilon) = 0$$

Thm Let $R(x) = \mathbb{E}(|f(x) - \hat{f}(x)|^2)$ be the risk
at x . Then

$$R(x) = \frac{1}{4} \sigma_u^4 h^4 f''(x)^2 + \frac{f(x)}{nh} \int K^2(x) dx + o\left(\frac{1}{n}\right) + o(h^6).$$