

Scaling issues w/ rejection sampling

e.g. from Learning in Graphical Models

We need an envelope q ^(or 'proposal') to match the entire target distribution f .

Even if we have q s.t. $\sup \{f/q\} \leq K < \infty$

Finding K may be v. hard

K may be v. large \Rightarrow alg. may be very inefficient

Toy E.g. Sample M -dim. Gaussian w/ mean $\underline{0}$.
Want $X \sim N_M(\underline{0}, \sigma_f^2 I) \leftarrow f$ (target)
Say q is M -dim. Gaussian w/ mean $\underline{0}$.
'Propose' $Y \sim N_M(\underline{0}, \sigma_q^2 I) \leftarrow q$ (proposal)

Suppose $\sigma_q = 1.01 \sigma_f$ so proposal(q) is heavier-tailed than target(f).

Smallest K satisfying $\sup_x \frac{f(x)}{q(x)} \leq K$

$$\sup_x \frac{f(x)}{q(x)} = \frac{\overset{\text{maximized at origin}}{(2\pi\sigma_f^2)^{-M/2}}}{(2\pi\sigma_q^2)^{-M/2}} = \left(\frac{\sigma_q}{\sigma_f}\right)^M = \exp\{M \log(1.01)\}$$

If $M = 1000$, $K \approx 20,000$
 \Rightarrow acc. rate $< 1/20,000$

K grows exponentially w/ M .

[Of course, we would sample these rvs individually in practice \because they are independent.]

Aside : generating multivariate normals $N_p(\underline{\mu}, \Sigma)$.

Recall : generating univariate normal : $x \sim N(0, 1)$

Then, $\sigma x + \mu \sim N(\mu, \sigma^2)$

Multivariate normal : $\underline{x} \sim N_p(\underline{0}, I_{NN})$ (p i.i.d. $N(0, 1)$)

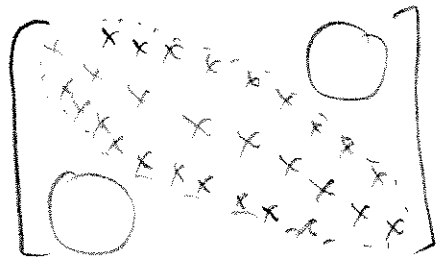
Find C , 'Choleski factor' of Σ st. $CC^T = \Sigma$.

$$C \underline{x} \sim N_p(0, C I C^T) = N_p(0, \Sigma)$$

$$C \underline{x} + \underline{\mu} \sim N_p(\underline{\mu}, \Sigma).$$

Choleski factoring/decomposition is computationally expensive, of order $p^3/3$ flops (floating point operations).

For large p : if matrix has special structure, use it. For e.g. band matrix



Common in inverse var. of Markov random field

models
E.g. AR-1, $\Sigma^{-1} = \begin{bmatrix} \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times & \times & \times & \times & \times & \times \end{bmatrix}$

Choleski takes $p b_w^2$ flops where b_w = bandwidth

= upper b_w (# non-zero diag above diagonal)
+ lower " (" " " below)
+ 1

mvnorm

MC 15(a)

M.C. s. errors
Assessing accuracy of estimates
 $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f$ $\hat{\mu}_n = \frac{\sum_{i=1}^n g(X_i)}{n}$

From CLT before:

$$\hat{\mu}_n \approx N(\mu, \sigma_g^2/n) \text{ for large } n.$$

$$\text{where } \sigma_g^2 = \text{Var}_f g(x)$$

Easy to estimate σ_g^2 : compute sample variance
based on $g(X_1), \dots, g(X_n)$.

$$\hat{\sigma}_g^2 = \frac{1}{n-1} \sum_{i=1}^n (g(X_i) - \hat{\mu}_n)^2$$

Along w/ our estimates, report our estimated M.C. s. error
 $= \hat{\sigma}_g / \sqrt{n}$.

Standard, sensible frequentist statistics!

Importance Sampling

Not a method for generating samples from f .

Instead estimate expectations w.r.t. f by reweighting samples from another distr. g (importance function).

Basic idea: $\mu = \int g(x) f(x) dx = E_f\{g(x)\}$

Before: $\hat{\mu}_n = \frac{\sum_{i=1}^n g(X_i)}{n}$ when $X_1, \dots, X_n \stackrel{iid}{\sim} f$

Now if $\mu = \int \left\{ g(x) \frac{f(x)}{g(x)} \right\} g(x) dx$ ~~(Assume $g(x) > 0 \Rightarrow f(x) \geq 0$)~~
if $g(x) > 0$ whenever $g(x)f(x) \neq 0$

Imp. sampling estimate, $\hat{\mu}_n = \frac{\sum_{i=1}^n g(Y_i) \frac{f(Y_i)}{g(Y_i)}}{n}$ when $Y_1, \dots, Y_n \stackrel{iid}{\sim} g$

$\hat{\mu}_n \rightarrow \mu$ (as before by SLLN)
if $E_f|g| < \infty$

Problem: we are assuming normalizing constant is known.

More general ratio imp. sampling estimate: assuming $g(x) > 0$ whenever $f(x) > 0$

Note that ~~E_f~~ $E_g \left\{ \frac{f(x)}{g(x)} \right\} = \int \frac{f(x)}{g(x)} g(x) dx = 1$.

So, $\mu = \frac{E_g \left\{ g(x) \frac{f(x)}{g(x)} \right\}}{E_g \left\{ \frac{f(x)}{g(x)} \right\}}$
 $= \frac{E_g \left\{ g(x) \frac{h(x)}{g(x)} \right\}}{E_g \left\{ \frac{h(x)}{g(x)} \right\}}$

where $\frac{h(x)}{c} = f(x)$

Suggesting the estimator: $Y_i \stackrel{iid}{\sim} g$

$\tilde{\mu}_n = \frac{\sum_{i=1}^n g(Y_i) \frac{h(Y_i)}{g(Y_i)}}{\sum_{i=1}^n \frac{h(Y_i)}{g(Y_i)}} \rightarrow \mu$

SLLN + a version of Slutsky's Thm.

$\tilde{\mu}_n$ has a small bias (unlike $\hat{\mu}_n$) but it has a lower MSE (JLW pg. 33)

Note that
$$\tilde{\mu}_n = \sum_{i=1}^n g(Y_i) \underbrace{\left[\frac{h(Y_i)/q(Y_i)}{\sum_{i=1}^n h(Y_i)/q(Y_i)} \right]}_{\text{'normalized' weights}}$$

$$= \sum_{i=1}^n g(Y_i) \tilde{w}(Y_i)$$

where
$$\tilde{w}(Y_i) = \frac{h(Y_i)/q(Y_i)}{\sum_{i=1}^n h(Y_i)/q(Y_i)}$$

In practice, $\tilde{w}(Y_i)$ is numerically more stable than unnormalized version.

Note that numerator and denominator can be estimated using different sets of samples.

$$\mu = \frac{E_{q_1} \left(g(x) \frac{h(x)}{q_1(x)} \right)}{E_{q_2} \left(\frac{h(x)}{q_2(x)} \right)}$$

as long as $q_1(x) > 0$ whenever $g(x)h(x) \neq 0$
and $q_2(x) > 0$ " $h(x) > 0$

Can estimate μ from $Y_1, \dots, Y_n \stackrel{iid}{\sim} q_1$
 $W_1, \dots, W_n \stackrel{iid}{\sim} q_2$

Useful when $g(x)h(x) \neq 0$ on a very different region than where $h(x) > 0$. e.g. when estimating tail probabilities.

Another use of imp. sampling:

Can use single sample (set of samples) to compute expectations w.r.t. many different distr. by reweighting samples differently each time.

Suppose we have a parametric family $\{f_\theta: \theta \in \Theta\}$ and we want

$$\mu(\theta) = E_{f_\theta}(g(x)) = \int g(x) f_\theta(x) dx \text{ for } \theta \in \Theta$$

E.g. $f_\theta(x)$ is normal density and $\theta = (\mu, \sigma^2)$
And, we want $E(x^2)$ for different combinations of μ, σ^2 .

How this may be useful: If we want to find

θ that maximizes $E_{f_\theta}(g(x))$.

Note: also useful for Monte Carlo maximum likelihood.

Naive M.C. assuming we can sample from $f_\theta \forall \theta \in \Theta$.

$X_1^{(\theta_1)}, \dots, X_{n_1}^{(\theta_1)} \sim f_{\theta_1}$ to estimate $E_{f_{\theta_1}}(g(x))$

$X_1^{(\theta_k)}, \dots, X_{n_k}^{(\theta_k)} \sim f_{\theta_k} \dots E_{f_{\theta_k}}(g(x))$

Too many samples to generate especially if k is large, Θ is highly multivariate.

However, suppose we find q s.t.

$$f_\theta(x) > 0 \Rightarrow q(x) > 0 \quad \forall x, \forall \theta \in \Theta$$

Then, simulate one set of samples

$$X_1, \dots, X_n \sim q.$$

Use importance sampling

$$\text{Estimate of } E_{f_\theta}(g(x)) = \sum_{i=1}^n g(x_i) w_\theta(x_i)$$

$$\text{where } w_\theta(x_i) = \frac{h_\theta(x_i) / q(x_i)}{\sum_{i=1}^n \frac{h_\theta(x_i)}{q(x_i)}}$$

~~where~~ and $\frac{h_\theta(x)}{c(\theta)} = f_\theta(x)$

Note: $c(\theta)$ is a function of θ so really a normalizing function.

Estimate $E_{f_{\theta_i}}(g(x))$ for any $\theta_i \in \Theta$ using single set of samples.

Very efficient (but not due to variance reduction.)
 q ~~need~~ ^{should} not be optimal for a single f_{θ_i} ,
but rather should be adequate $\forall \theta \in \Theta$.

Difficult to find good q in general. (may not be able to get away w/ single q in practice for hard problems.)

Observation:

- (1) Importance sampling is useful for estimating normalizing constants. (This is implicitly used in ratio imp. sampling.)

We know $h(x)/c = f(x)$

Given $h(x)$, how can we estimate c ?

$$\therefore E_q \left\{ \frac{f(x)}{q(x)} \right\} = 1$$

We have $E_q \left\{ \frac{h(x)}{q(x)} \right\} = c$ normalizing constant is written here as an expectation

Can use Monte Carlo to estimate expectation.
Useful for Monte Carlo max. liked.

- (2) General purpose tool for integration.

Suppose you want to find $\int_a^b \psi(x) dx$

If you have q s.t. $q(x) > 0$ when $\psi(x) \neq 0$,

$$\int_a^b \psi(x) dx = \int_a^b \left[\frac{\psi(x)}{q(x)} \right] q(x) dx = E_q \left\{ \frac{\psi(x)}{q(x)} \right\}.$$

Rare event problems: importance sampling can be very useful.

E.g. want $P(Z > 4.5)$ where $Z \sim N(0,1)$.

Naive M.C.: $Z_1, \dots, Z_n \stackrel{iid}{\sim} N(0,1)$

$$\hat{\mu}_n = \sum_{i=1}^n \mathbb{I}(Z_i > 4.5) / n$$

Even for $n = 100,000$, $\hat{\mu}_n$ usually 0 ($\because P(Z > 4.5)$ is small)

Instead use $q = \text{shifted Exp}(4.5, 1)$

$Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{shifted Exp}(4.5, 1)$

Shifted Exp: $\text{Expon}(\text{scale}=1)$ shifted right to $\neq 4.5$



$$pdf = \frac{e^{-(x-4.5)}}{\int_{4.5}^{\infty} e^{-x} dx}$$

My Experiment

Importance sampling:
With $n = 10,000$

$$\hat{\mu}_n = \frac{\sum_{i=1}^n \mathbb{I}(Y_i > 4.5) \frac{f(Y_i)}{q(Y_i)}}{n}$$

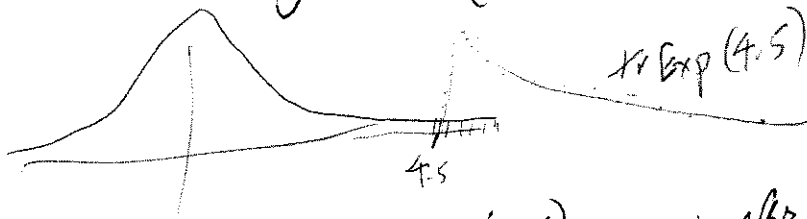
gives 3.35×10^{-6}

$n = 100,000$ gives 3.38×10^{-6}

True value (R's pnorm function): 3.3976×10^{-6}

Case study: $P(\bar{X} > 4.5)$

$$Z \sim N(0,1)$$



Simple
imp. est.

$$\begin{aligned} q &= \text{tr Exp}(4.5) \\ q &= \text{tr Exp}(2) \\ q &= \text{tr E}(5) \end{aligned}$$

works well
" "

does not work

} support ≤ 4.5

} support > 4.5

Ratio est.

$$\begin{aligned} q &= \text{tr E}(4.5) \Rightarrow \tilde{\mu} = 1 \\ q &= \text{tr E}(2) \Rightarrow \tilde{\mu} = 0.001 \\ q &= \text{tr E}(5) \Rightarrow \tilde{\mu} = 1 \end{aligned}$$

does
not
work

does
not
work

- Q. 1. How can we get ratio est. to work?
- Q. 2. What are conditions under which
- (1) simple imp. est work?
- (2) ratio imp. est work?

Ratio est. works if
~~Ratio est.~~ does not " "

$$\begin{aligned} q &= N(4.5, 1) \\ q &= \text{tr Exp}(1) \end{aligned}$$

(Need 1 million samples
to get good estimate)
even 10 million
samples

1 million samples, $\tilde{\mu} = 6.9 \times 10^{-6}$

FOR MYSELF (NOT CLASS) MC 20 extra

Notes: Naïve imp. sampling works as long as q is s.t. $g(x)f(x) \neq 0 \Rightarrow g(x) > 0$ a.s.

Less stringent than ratio imp. sampling

condition: q is s.t. $f(x) > 0 \Rightarrow g(x) > 0$ a.s.

For e.g. for our toy problem $q = N(4.5, 1)$ works

but $q = \text{SkExp}(4.5)$ does not work (when using ratio imp. sampling.)

Soln: Can use different importance fns. for numerator and denominator.

Tail probability problems are common: e.g. hypothesis testing. (p-values).

Estimating probabilities of rare events. e.g. physics, astronomy

Importance sampling : M.C. standard errors

Estimate is useless w/o estimate of its variability

Simple imp. sampling : easy!

$$\begin{aligned}\text{Var}(\hat{M}_n) &= \text{Var} \left\{ \frac{1}{n} \sum_{i=1}^n g(x_i) \frac{f(x_i)}{q(x_i)} \right\} \\ &= \frac{1}{n} \text{Var}_q \left\{ g(x) \frac{f(x)}{q(x)} \right\} \quad (\because x_1, \dots, x_n \stackrel{\text{iid}}{\sim} q)\end{aligned}$$

\therefore we have sample $x_1, \dots, x_n \stackrel{\text{iid}}{\sim} q$,

estimate $\text{Var}_q \left\{ g(x) \frac{f(x)}{q(x)} \right\}$ by sample variance
of $g(x_1) \frac{f(x_1)}{q(x_1)}, \dots, g(x_n) \frac{f(x_n)}{q(x_n)}$, say $\hat{\sigma}^2$

$$\text{Est. M.C. s. error} = \frac{\hat{\sigma}}{\sqrt{n}}$$

Note: No guarantee that ~~above~~ ^{M.C.} s. error is finite

$$\text{If } E_q \left\{ g^2(x) \frac{f^2(x)}{q^2(x)} \right\} < \infty \quad \text{and} \quad E_q \left\{ \frac{f^2(x)}{q^2(x)} \right\} < \infty$$

then s. error is finite.

Implication: q w/ lighter tails than f are not appropriate since variance will be infinite for many

for g .

$$\left(\text{Second condition } E_q \left\{ \frac{f^2(x)}{q^2(x)} \right\} = \int \frac{f^2(x)}{q(x)} dx = \int \frac{f(x)}{q(x)} f(x) dx \leq \int K f(x) dx \right. \\ \left. = K < \infty \text{ if resampling condition is satisfied} \right)$$

However, more important to consider ratio imp. estimator.
 To rigorously quantify M.C. s-error, need
 C.L.T.

Imp. sampling CLT: If $E_g \left\{ \frac{f^2(x)}{g^2(x)} \right\} < \infty$ and
 $E_g \left\{ g^2(x) \frac{f^2(x)}{g^2(x)} \right\} < \infty$ then (1) the CLT holds for the
 ratio imp. sampling estimate and (2) the method of
 moments estimate of its asymptotic variance is
 consistent.

Multivariate CLT: If $X_1, \dots, X_n \stackrel{iid}{\sim} g$, $w(x_i) = \frac{f(x_i)}{g(x_i)}$,

$$\sqrt{n} \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n g(x_i) w(x_i) - \mu \\ \frac{1}{n} \sum_{i=1}^n w(x_i) - 1 \end{pmatrix} \xrightarrow{d} N(0, \Sigma_{2 \times 2}).$$

To derive variance of ratio estimator, appeal to
 delta-method where if $Q(u, v)$ is real-valued, then

$$\sqrt{n} \left(Q \left(\frac{1}{n} \sum_{i=1}^n g(x_i) w(x_i), \frac{1}{n} \sum_{i=1}^n w(x_i) \right) - Q(\mu, 1) \right) \\ \xrightarrow{d} N(0, \underbrace{Q'(\mu, 1)^T \Sigma Q'(\mu, 1)}_{\sigma^2, \text{asymptotic variance}})$$

where Q' is a 2×1 matrix of partial derivatives.

~~For~~

For ratio estimation $Q(a, b) = a/b$,

$$Q'(a, b) = \left(\frac{1}{b}, -\frac{a}{b^2} \right)^T$$

$$Q'(\mu, 1) = (1, -\mu)^T$$

$$\Sigma = \begin{bmatrix} \text{Var}(a) & \text{Cov}(a, b) \\ \text{Cov}(a, b) & \text{Var}(b) \end{bmatrix}$$

$$\text{where } a = \frac{1}{n} \sum_{i=1}^n g(x_i) w(x_i), \quad b = \frac{1}{n} \sum_{i=1}^n w(x_i)$$

$$\text{Thus, } \sqrt{n} (\tilde{\mu}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$$

$$\text{where } \sigma^2 = Q'(\mu, 1)^T \Sigma Q'(\mu, 1).$$

If we obtain consistent estimator of σ^2 , such as method of moments est., can appeal to Slutsky's Thm.

Simplifying, we get

$$\hat{\sigma}^2 = n \left(\frac{\sum_{i=1}^n g(x_i) \tilde{w}(x_i)}{\sum_{i=1}^n \tilde{w}(x_i)} \right)^2 \left[\frac{\sum_{i=1}^n g^2(x_i) \tilde{w}^2(x_i)}{\left(\sum_{i=1}^n g(x_i) \tilde{w}(x_i) \right)^2} + \frac{\sum_{i=1}^n \tilde{w}^2(x_i)}{\left(\sum_{i=1}^n \tilde{w}(x_i) \right)^2} - 2 \frac{\sum_{i=1}^n g(x_i) \tilde{w}^2(x_i)}{\left(\sum_{i=1}^n g(x_i) \tilde{w}(x_i) \right) \sum_{i=1}^n \tilde{w}(x_i)} \right]$$

Comparison to rejection sampling

1. Is a generalization of rejection sampling, and is always at least as efficient \nwarrow (Y. Chen, 2006)
when using χ^2 distances
2. Do not require

$$\sup_{x \in \Omega} \frac{f(x)}{g(x)} \leq K \quad \text{for some } K < \infty$$

However, this is very desirable: if above holds and $E_g \{g^2(x)\} < \infty$ CLT holds and m.o.m. variance estimates are consistent.

3. Do not need to find K above.
4. May not want g to match f well; may be more important to match gf , for e.g.
5. Variance estimates are more unreliable.
6. Numerical stability issues: need to exponentiate $\log(h/g)$ to calculate weights. No need for this in rejection sampling (can do everything in log-scale).
7. Harder to estimate density. Use cdf at many pts.

$$F(x) = \int \mathbb{I}(X \leq x) f(x) dx = E_f \{ \mathbb{I}(X \leq x) \}.$$

Imp. sampling estimates can be unstable

Suppose $q(x)$ is v. small for $x \in S$. where $g(x)h(x)$ is very large.

If $x^* \in S$ is obtained, it gets huge weight $\left(\frac{g(x)h(x)}{q(x)}\right)$

$\Rightarrow \tilde{M}$ changes drastically. Improvement E. Toudes 2011, Truncated importance sampling

M.C. s.error: Even after many samples drawn, may not observe $x^* \in S$ ($\because q(x^*)$ is small).

Hence, may never see x^* s.t. $g(x^*) \frac{h(x^*)}{q(x^*)}$ is huge,

so s.error will be heavily underestimated.

Worse/unstable M.C. estimates may imply worse/unreliable s.error estimates. (Worst possible situation: bad estimates, but we think we have good estimates!)

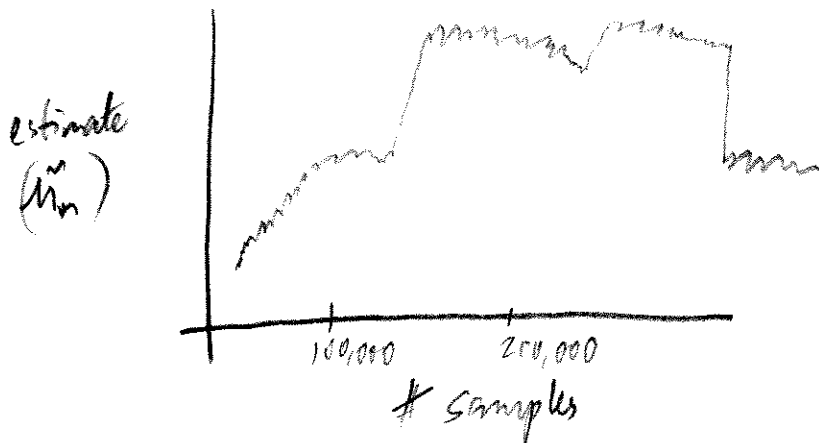
Importance sampling:

How many samples are necessary (what is N ?)
for good estimates?

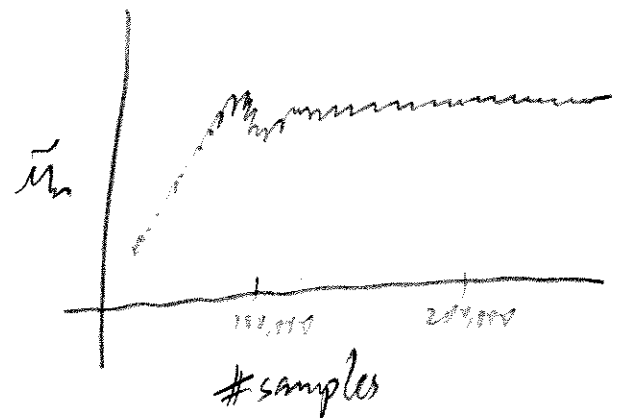
Best approach: M.C. s.errors: when errors are small enough, stop sampling.

Problem: M.C. s.errors are not always reliable so it is best to use M.C. s.errors + some simple heuristics ('reality checks').

Trace plots:



(a)



(b)

(a) Obviously poorly converging estimate

(b) Apparently more stable " (but cannot be 100% sure).

Plots: can tell you if estimate is obviously unstable but cannot tell you when estimator (and M.C. s.e.s) are reliable. (jump could happen after MC 28 million iteration for plot (b).)

Techniques for finding importance functions

Some approaches (not mutually exclusive):

1. Exponential tilting see Ross
2. Defensive importance fn. (T. Hesterberg '95 Technometrics)

Use a mixture: $q(x) = p q_1(x) + (1-p) q_2(x)$ $p \in (0,1)$
w/ $q_1(x)$ matching $g(x)f(x)$ well but $q_2(x)$ is heavy-tailed. Set p close to 1; this ensures variances are finite.

3. Laplace approx: (Tierney '89, Tierney & Kadane '86)

Assume $\log f$ admits Taylor expansion about its mode.

Let $l(x) + a = \log f(x)$ for some constant a .

Let \hat{x} satisfy $l'(\hat{x}) = 0$ l' = vector of 1st derivatives of l .

and let $l''(x)$ be matrix of 2nd derivatives of l .

Then,
$$M = \int g(x) f(x) dx = \int g(x) e^{l(x)+a} dx = e^a \int g(x) e^{l(x)} dx$$
$$\approx e^a \int g(x) \exp(l(\hat{x}) + (x-\hat{x})^T l''(\hat{x}) (x-\hat{x})/2) dx \quad (*)$$

(∵ $l'(\hat{x}) = 0$)

