

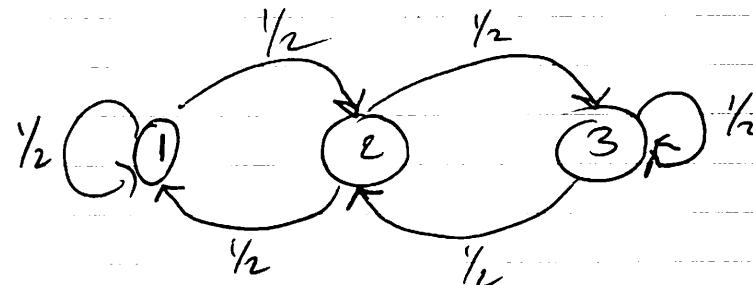
Markov Chams

\Rightarrow A Markov Chain is a stochastic process $\{X_i\}$ which satisfies

$$P(X_i | X_{i-1}, \dots, X_0) = P(X_i | X_{i-1})$$

so that the current state of the chain X_i only depends on the previous value X_{i-1} .

\Rightarrow The possible values $\{X_i\}$ can take is called the state space



$$P = \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 1/2 \\ 1/2 & 1/2 & 1/2 \end{bmatrix} \quad \text{Transition matrix (Kernel)}$$

$$P(X_n=j | X_{n-1}=i) = P_{ij}$$

Suppose we start the chain with initial dist. $(1, 0, 0)$, what is distribution after n iterations?

$$(0, 0, 1)$$

Given transition matrix P

$$P(X_n=j | X_0=i) = (P^n)_{ij}$$

where $P^n = \underbrace{P \times P \times P \times \dots \times P}_{n \text{ times}}$

$$\pi_n = \pi_0 P^n = \underbrace{\pi_0 P \times P \times P \times \dots \times P}_{n \text{ times}}$$

Let π^* be a vector (distribution) such that

$$\pi_* P = \pi_*$$

Call π_* a stationary distribution and a MC
stationary or said to be stationary if it
reaches ~~has~~ this distribution.

Baire limit theorem says that under some
conditions $\|\pi_* - \pi_n\| \rightarrow 0$ as $n \rightarrow \infty$.

No matter how we start the chain (π_0) π_n
approaches π_* .

Assumptions:

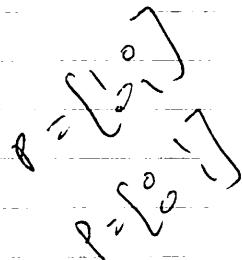
① Irreducible: $P_{ij}^n > 0$ for some n and

All pairs i, j are accessible from one another
for every i, j

② Aperiodic: the chain does not make deterministic visits to a subset of the state space.

$$\text{ex: } X_0 = 0 \quad X_n = X_{n-1} + \varepsilon_n$$

$$\varepsilon_n = \begin{cases} 1 & \frac{1}{2} \\ -1 & \frac{1}{2} \end{cases}$$



For n even, the chain hits subset of even #s
For n odd " " " " " odd #s

③ Positive Recurrent: Every state is visited i.o.
The waiting time between visits from ~~state~~ any 2 states
is finite

④ Stationary dist exists

Q

Let X_0, X_1, \dots be an irreducible, openable, Markov chain w/stationary dist. π .

Let $X_0 \sim \pi_0$, some arbitrary distribution. Then

$$\pi_n(i) \rightarrow \pi(i) \quad \forall i \text{ as } n \rightarrow \infty$$

irreducible = all states communicate w/each other

Time Reversiblity

A Markov Chain is Time reversible if

$$(X_0, X_1, \dots, X_n) \stackrel{D}{=} (X_n, X_{n-1}, \dots, X_0)$$

If $\{X_n\}$ is TR, then

$$(X_0, X_1) \stackrel{D}{=} (X_1, X_0)$$

$$\Rightarrow X_0 \stackrel{D}{=} X_1$$

$$\Rightarrow \pi_1 = \pi_0$$

Since $\pi_1 = \pi_0 P$, the initial dist. is stationary

Let $\pi = \pi_0$. TR chains have

$$\pi_i P_{ij} = \pi_j P_{ji}$$

g. from $i \rightarrow j$ g. from $j \rightarrow i$

Note That

$$\text{TP}(X_n = j \mid X_{n-1} = i) = \frac{\text{TP}(X_n = j, X_{n-1} = i)}{\pi_i}$$

$$\text{TR} = \frac{\text{TP}(X_n = i, X_{n-1} = j)}{\pi_j}$$

$$= \text{TP}(X_{n-1} = j \mid X_n = i)$$

$$\Rightarrow P(i \rightarrow j) = P(j \rightarrow i) \quad \text{"flux" of } i \rightarrow j = j \rightarrow i$$

Main Ideas:

- ① We want to sample some complicated density π .
- ② We know that certain Markov chains will converge to a stationary distribution (if exists)
- ③ How to construct a MC s.t. the sequence of values $\{X_n\}$ converges to a target distribution π ?

MCMC \Rightarrow

We know that if a MC with transition matrix P or Kernel $K(x, y)$ is time reversible wrt π . Then π must be the stationary dist of the MC.

Given the chain we start at some point and run it until convergence

Problem: We cannot start at the stationary distribution!

Metropolis-Hastings

Let $q(y|x)$ be a transition density from which we can easily simulate.

Let $X_0 = x$

$$\frac{p(x|y)}{p(x)}$$

① Simulate a candidate $y \sim q(y|x)$

② Let $\alpha(y|x) = \min\left\{\frac{\pi(y) q(x|y)}{\pi(x) q(y|x)}, 1\right\}$

not necessary

③ Simulate $U \sim \text{unif}(0,1)$ ~~and accept~~

If $U \leq \alpha(y|x)$ then set $X_{n+1} = y$.

Otherwise set $X_{n+1} = x$.

✓

Why does this work? We need to show
~~first~~ let $K(y|x)$ be the transition density of the Metropolized chain.

We need to show

$$\pi(y) K(x|y) = \pi(x) K(y|x)$$

$$K(y|x) = \alpha(y|x) q(y|x), \text{ so we need}$$

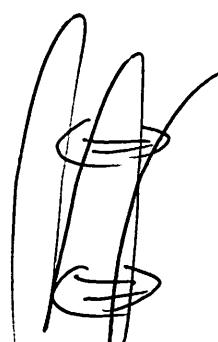
$$\alpha(y|x) q(y|x) \pi(x) = \alpha(x|y) q(x|y) \pi(y)$$

$$\min\left(\frac{\pi(y) q(x|y)}{\pi(x) q(y|x)}, 1\right) q(y|x) \pi(x) = \min\left(\frac{\pi(x) q(y|x)}{\pi(y) q(x|y)}, 1\right) q(x|y) \pi(y)$$

$$\min\left(\pi(x) q(y|x), 1\right) q(y|x) \pi(x) = \min\left(\pi(x) q(y|x), 1\right) \pi(y)$$

Observe

Constants
of prop. cancel
out



TOPIC: Metropolis Hastings Time-Reversal

FILE UNDER:

PAGE:

$$K(y|x) = \alpha(y|x) q(y|x)$$

$$+ \mathbb{1}\{y=x\} \left[1 - \int \alpha(s|x) q(s|x) ds \right]$$

(2)

$$(1) \text{ Show } K(y|x) \pi(x) = K(x|y) \pi(y)$$

$$\Leftrightarrow \alpha(y|x) q(y|x) \pi(x) = \alpha(x|y) q(x|y) \pi(y)$$

$$\Leftrightarrow \min\left(\frac{\pi(y)}{\pi(x)} \frac{q(x|y)}{q(y|x)}, 1\right) q(y|x) \pi(x) \\ = \min\left(\frac{\pi(x)}{\pi(y)} \frac{q(y|x)}{q(x|y)}, 1\right) q(x|y) \pi(y)$$

$$\cancel{\int_a q(y|x) + \int_b q(y|x) ds} \Rightarrow \min\left(\pi(y) q(x|y), q(y|x) \pi(x)\right) = \min\left(\pi(x) q(y|x), q(x|y) \pi(y)\right)$$

$$(2) \quad \cancel{\mathbb{1}\{y=x\}} r(x) = \mathbb{1}\{y=x\} \left[1 - \int \alpha(s|x) q(s|x) ds \right]$$

$$\cancel{s_x(y) r(x) \pi(x)} = s_y(x) r(y) \pi(y)$$

over set of points where $y=x$. Obsrvns!

Variations on M-H

Random Walk

Let $q(y|x)$ be defined by

$$y = x + \varepsilon$$

where $\varepsilon \sim g$ and g is symmetric about 0.

Note $q(y|x) = g(\varepsilon)$

$$\left. \begin{aligned} q(x|y) &\stackrel{||}{=} g(-\varepsilon) \end{aligned} \right\} \cancel{= g(\varepsilon)}$$

So the Metropolis acceptance prob is

$$\alpha(y|x) = \min \left\{ \frac{\pi(y) q(x|y)}{\pi(x) q(y|x)}, 1 \right\}$$

$$= \min \left\{ \frac{\pi(y)}{\pi(x)}, 1 \right\}$$

Start at $x_n = x$.

① Simulate $\varepsilon \sim g$. Let $y = x + \varepsilon$

② Simulate $u \sim \text{unif}(0,1)$.

③ If $u \leq \alpha(y|x_n) = \min \left(\frac{\pi(y)}{\pi(x_n)}, 1 \right)$

then $x_{n+1} = y$, else $x_{n+1} = x_n$.

Ex. Generating Normal R.V.s

Let g be $\text{Unif}[-8, 8]$ dist.

① Simulate $\varepsilon \sim U(-8, 8)$

② ~~Re~~ Simulate $u \sim \text{Unif}(0, 1)$

③ Accept $y = x + \varepsilon$ if

$$u \leq \min \left\{ \frac{\pi(y)}{\pi(x)}, 1 \right\}$$

Independence Metropolis Algorithm (Tierney)

Define $q(y|x) = q(y)$.

So

$$\alpha(y|x) = \min \left(\frac{\pi(y)q(x)}{\pi(x)q(y)}, 1 \right)$$

$$= \min \left(\frac{w(y)}{w(x)}, 1 \right)$$

(\hookrightarrow importance weights)

① Simulate $y \sim q(y)$

② Simulate $u \sim \text{Unif}(0, 1)$

③ Accept y if $u \leq \min \left(\frac{w(y)}{w(x)}, 1 \right)$

Independence Metropolis Sampler seems to work well in same situations as rejection sampling.

If $C = \sup_x \frac{\pi(x)}{g(x)} < \infty$, then we have

$$\|\pi_n - \pi\| < k\rho^n \text{ for } 0 < \rho < 1$$

Rate of convergence ρ depends on C being small (close to 1).

Gibbs Sampler (2 variables)

Let $\pi(x, y)$ be your target density.

Let $Z = (x, y)$ so that we want to generate a MC $\{Z_n\}$ where π is the stationary dist. of the chain. The full conditionals of π are

$$p(y|x) = \frac{\pi(x, y)}{\pi(x)} \propto \pi(x, y)$$

$$p(x|y) = \frac{\pi(x, y)}{\pi(y)} \propto \pi(x, y)$$

Let $Z_n = (x_n, y_n)$. The Gibbs sampler obtains

Z_{n+1} by

① Simulate $X_{n+1} \sim p(x|y_n)$

② Simulate $Y_{n+1} \sim p(y|x_{n+1})$

$Z_{n+1} = (x_{n+1}, y_{n+1})$

~~For MCs, need~~

time it sample
space is product

~~Irreducibility — positive probability of visiting any set of measure > 0~~

~~Aperiodic — no regularly scheduled visits~~

~~Positive Harris Recurrence~~

\hookrightarrow Any set with $\pi(A) > 0$ is visited infinitely often

\Rightarrow Imply $X_n | X_0 = x \rightarrow \pi$

\downarrow
~~stationary dist~~

Slice Sampling

Given state x_n , generate a

auxiliary variable

$y_{n+1} \sim \text{Unif}[0, \pi(x_n)]$

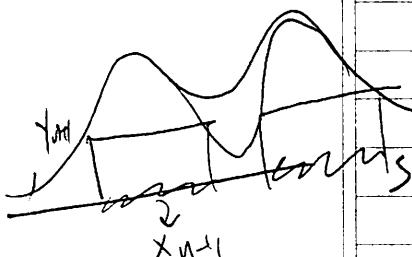
Given y_{n+1} , generate

$x_{n+1} \sim \text{Unif}\{x : \pi(x) \geq y_{n+1}\}$

Note that

$$\pi(x) = \int \mathbb{1}\{0 \leq y \leq \pi(x)\} dy$$

so flat



$$f(x, y) = \mathbb{1}\{0 \leq y \leq \pi(x)\} \text{ or}$$

a joint density

Gibbs Sampling (we want to sample from some posterior; or other complicated density)

For K -dim distributions, only need $K_{\text{conditional}}^{\text{full}}$ dists, not $K(K-1)$ conditionals i.e. K w/ DA.

$p(x, y, z) \rightarrow$ target density
 \downarrow

$$p(x|y, z)$$

$$p(y|x, z)$$

$$p(z|x, y)$$

Given the n th iteration (x_n, y_n, z_n)

① Sample $X_{n+1} \sim p(x|y_n, z_n)$

② Sample $Y_{n+1} \sim p(y|X_{n+1}, z_n)$

③ Sample $Z_{n+1} \sim p(z|X_{n+1}, Y_{n+1})$

$$[x_n, y_n, z_n] \xrightarrow{D} [x, y, z]$$

$$[x_n] \xrightarrow{D} [x]$$

$$[y_n] \xrightarrow{D} [y]$$

$$[z_n] \xrightarrow{D} [z]$$

Get
Geman +
Geman

$$\frac{1}{N} \sum_{n=1}^N f(x_n, y_n, z_n) \rightarrow \mathbb{E}[f(x, y, z)]$$

$$K(z_{n+1} | z_n) = p(y_{n+1} | x_{n+1}) p(x_{n+1} | y_n)$$

(no x_n)

~~Defn
Joint
Conditional
Marginal~~

Ex: Simulate $N(\mu, \Sigma) = N(\mu_1, \mu_2), \begin{pmatrix} \sigma_1^2 & \rho \\ \rho & \sigma_2^2 \end{pmatrix}$

$$\text{Set } z_n = (x_n, y_n)$$

- ① Simulate $x_{n+1} \sim N\left(\mu_1 + \frac{\rho(\mu_2 - \mu_1)}{1-\rho}, \sigma_1^2(1-\rho)\right)$
- ② Simulate $y_{n+1} \sim N\left(\mu_2 + \frac{\rho(\mu_1 - \mu_2)}{1-\rho}, \sigma_2^2(1-\rho)\right)$

Ex: Let $y_1, \dots, y_n \sim N(\mu, \tau^{-1})$

Suppose we put priors (indep.).

$$\mu \sim N(0, W^{-1})$$

$$\tau \sim \text{Gamma}(\alpha, \beta)$$

No conjugacy, so we need Gibbs sampling to explore posterior $\mu, \tau | y_1, \dots, y_n$.

$$p(\mu, \tau, y) \propto L(\mu, \tau | y) p(\mu) p(\tau)$$

$$p(\mu | \tau, y) \propto L(\mu, \tau | y) p(\mu)$$

$$p(\tau | \mu, y) \propto L(\mu, \tau | y) p(\tau)$$

$$p(\mu, \tau, \gamma) = \left(\frac{\tau}{2\pi}\right)^{n/2} \exp\left(-\frac{\tau}{2} \sum (y_i - \mu)^2\right)$$

$$\times \left(\frac{w}{2\pi}\right)^{1/2} \exp\left(-\frac{w}{2} \mu^2\right)$$

$$\times \frac{\tau^{\alpha-1}}{\Gamma(\alpha)} \beta^\alpha \exp(-\tau\beta)$$

$$p(\tau | \mu, \gamma) \propto \tau^{n/2} \exp\left(-\frac{\tau}{2} \sum (y_i - \mu)^2\right)$$

$$\times \tau^{\alpha-1} \exp(-\tau\beta)$$

$$= \tau^{\alpha-1 + n/2} \exp\left(-\tau\left(\beta + \frac{1}{2} \sum (y_i - \mu)^2\right)\right)$$

$$= \text{Gamma}\left(\alpha + \frac{n}{2}, \beta + \frac{1}{2} \sum (y_i - \mu)^2\right)$$

$$p(\mu | \tau, \gamma) \propto \exp\left(-\frac{\tau}{2} \sum (y_i - \mu)^2\right)$$

$$\times \exp\left(-\frac{w}{2} \mu^2\right)$$

$$= \exp\left(-\frac{(n\tau + w)}{2} \mu^2 - (\sum y_i) \tau \mu\right)$$

$$= N\left(\frac{(n\tau)}{(n\tau + w)} \sum y_i, \frac{1}{n\tau + w}\right)$$

Gibbs sampler iterates between those densities

$$p(x, y^2, \mu, \sigma^2) = \prod_{i=1}^n p(y_i | \mu, \sigma^2)$$

~~x, y, μ~~

Relationship b/w Gibbs sampling + Metropolis - Hastings

In M-H we have proposal $q(Y|X)$ and acceptance prob.

$$\alpha(Y|X) = \min\left(\frac{\pi(Y) q(X|Y)}{\pi(X) q(Y|X)}, 1\right)$$

~~One can do
The Gibbs Sampling~~

Single-component M-H

Let $X^{(t)} = (X_1^{(t)}, X_2^{(t)}, \dots, X_K^{(t)})$ be a K -vector representing our quantities of interest (parameters) at iteration t .

Let $X_{-i}^{(t)} = (X_1^{(t)}, X_2^{(t)}, \dots, X_{i-1}^{(t)}, X_{i+1}^{(t)}, \dots, X_K^{(t)})$

SCMH, at iteration t updates the i th component $X_i^{(t)}$ via:

(1) Sample $Y_i \sim q_i(Y_i | X_{-i}, X_i)$ for component i

(2) Let $\alpha(Y_i | X_{-i}^{(t)}) = \min\left(\frac{\pi(Y_i | X_{-i}) q(X_i | Y_i, X_{-i})}{\pi(X_i | X_{-i}) q(Y_i | X_i, X_{-i})}, 1\right)$

(3) Accept Y_i for component i w/prob $\alpha(Y_i | X_{-i}^{(t)})$

Repeat K times.

$\pi(X_i | X_{-i})$ is a full conditional distribution

In Gibbs sampling

$$q(y_i | X_i, X_{-i}) = \pi(y_i | X_{-i})$$

Thus gives us

$$\alpha(y_i | X_i, X_{-i}) = \min\left(\frac{\pi(y_i | X_{-i}) \pi(x_i | x_i)}{\pi(x_i | X_{-i}) \pi(y_i | X_{-i})}, 1\right)$$

$$= 1$$

So Gibbs sampling is like SCMA but it always accepts

General M-H is very exploratory, maybe too much.

Blocking

Hybrid Gibbs Sampler

Sometimes it is not possible to sample directly from a full conditional dist.

One can use a hybrid GS in that case.

Suppose we have 2-variable problem w/ (X, Y) and we can easily sample $p(X|Y)$ but not $p(Y|X)$. Then for ~~$\mathbb{P}(Y_n | X_n)$~~ $y_n | x_n$

① Simulate $\hat{y} \sim q(\hat{y} | x_n)$, draw $u \sim \text{unif}(0, 1)$

② Accept y if $u \leq \min\left(\frac{p(\hat{y}|x_n) q(y_n|\hat{y})}{p(y_n|x_n) q(\hat{y}|x_n)}, 1\right)$

We just need to simulate one value and move on for things to work

236 hkt

Hybrid GS

$$Y_{ij} \sim \text{Bernoulli}(p_{ij})$$

$$\text{logit } p_{ij} = \alpha_i + \beta_j$$

~~$\alpha_i \sim N(\mu, \sigma^2)$~~

$$\alpha_i \sim N(\mu, \sigma^2)$$

~~$\beta_j \sim N(0, D)$~~

$$\beta_j \sim N(0, D)$$

$$\gamma^* \sim \text{IG}(a, b)$$

$$p(y, \alpha, \mu, \sigma^2) \propto \left[\prod_{i=1}^n \left[\prod_{j=1}^J p(Y_{ij} | \alpha_i) \right] p(\alpha_i | \mu, \sigma^2) \right] \pi(\mu) \pi(\sigma^2)$$

$$(1) p(\alpha_i | m) \propto \left[\prod_{j=1}^J p(Y_{ij} | \alpha_i) \right] p(\alpha_i | \mu, \sigma^2) \quad i = 1, \dots, n$$

$$(2) p(\mu | m) \propto \left[\prod_{i=1}^n p(\alpha_i | \mu, \sigma^2) \right] \pi(\mu)$$

$$= N\left(\frac{D}{D+\sigma^2} \bar{\alpha}, \frac{\sigma^2}{D+\sigma^2} D\right)$$

$$(3) p(\sigma^2 | m) \propto \left[\prod_{i=1}^n p(\alpha_i | \mu, \sigma^2) \right] \pi(\sigma^2)$$

$$\text{IG}\left(a + \frac{n}{2}, b + \frac{1}{2} \sum (\alpha_i - \mu)^2\right)$$

Reparameterization in Metropolis

Suppose we want to sample from a full conditional $p(\lambda | u)$ where $\lambda \in (0, 1)$. Random walk is tricky b/c of edge/boundary cases. We can

$$\text{Let } z = \log f(\lambda) = \log \frac{\lambda}{1-\lambda}.$$

Then propose values of $z \in (-\infty, \infty)$, ~~and~~ accept/reject, and transform back.

① propose $z^* \sim g(z | z_n)$, $z_n = \log f(\lambda_n)$

$$② \alpha(z | z_n) = \min \left\{ \frac{p(\log f^{-1}(z^*) | u)}{p(\log f^{-1}(z_n) | u)} \frac{g(z_n | z^*)}{g(z^* | z_n)} \frac{|\mathcal{J}(z^*)|}{|\mathcal{J}(z_n)|} \right\}$$

where $|\mathcal{J}(z)|$ is the determinant of Jacobian of transformation from $u \mapsto z$ (i.e. $\log f^{-1}$)

$$|\mathcal{J}(z)| = \frac{1}{|\mathcal{J}(u)|} \rightarrow z \mapsto u$$

Hamiltonian MCMC

The energy in a system is defined by the Hamiltonian,

$$H(q, p) = U(q) + K(p)$$

potential energy

Kinetic energy

The partial derivatives of $H(q, p)$ determine how q and p change over time t ,

$$\frac{dq_i}{dt} = \frac{\partial H}{\partial p_i} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Hamilton's equations}$$

$$\frac{dp_i}{dt} = -\frac{\partial H}{\partial q_i} \quad \text{for } i=1, \dots, d = \text{# dimensions}$$

\Rightarrow Solving these differential equations gives you $q(t)$ and $p(t)$ which allow you to describe how the system changes from time t to time $t+s$.

Properties of Hamiltonian dynamics

- ① Reversible \rightarrow preserve stationary distribution
- ② Volume preservation \rightarrow no Jacobian
- ③ Conservation of Hamiltonian \rightarrow acceptance prob = 1

Basic Procedure for HMC

- ① Substitute momentum

leapfrog Method - moving around the system

Assume $H(q, p) = U(q) + K(p)$

and $K(p) = \frac{1}{2} \sum_{i=1}^d \frac{p_i^2}{m_i}$ (Gaussian w/ diagonal var).

For step size ϵ , the basic idea is

$$p_i(t+\epsilon) = p_i(t) + \epsilon \frac{dp_i(t)}{dt} = p_i(t) - \epsilon \frac{\partial U}{\partial q_i}(q(t))$$

$$q_i(t+\epsilon) = q_i(t) + \epsilon \frac{dq_i(t)}{dt} = q_i(t) + \epsilon \frac{p_i(t)}{m_i}$$

Basic 1st order approx doesn't work well because
of error propagation

leapfrog Method says

$$\textcircled{1} \quad p_i(t + \varepsilon/2) = p_i(t) - \frac{\varepsilon}{2} \frac{\partial U}{\partial q_i}(q(t))$$

$$\textcircled{2} \quad q_i(t + \varepsilon) = q_i(t) + \frac{p_i(t + \varepsilon/2)}{m_i}$$

$$\textcircled{3} \quad p_i(t + \varepsilon) = p_i(t + \varepsilon/2) - \frac{\varepsilon}{2} \frac{\partial U}{\partial q_i}(q(t + \varepsilon))$$

~~HMC Algorithm~~

HMC Algorithm

\textcircled{1} Simulate p from

\textcircled{2} Run leapfrog for L steps with step ε to generate
proposal

\textcircled{3} Accept / Reject

For statistical applications $q(\tau)$ is vector of parameters that "evolves" over time and $p(\tau)$ is arbitrary. Usually, assume

$$K(p) = \frac{1}{2} p' M^{-1} p \quad \text{where} \quad M = \begin{bmatrix} m_{11} & 0 \\ 0 & m_{22} \end{bmatrix}$$

$$U(q) = -\log \pi(q) = -\log p(i) L(q|D)$$

$$\Rightarrow H(q, p) = U(q) + K(p) = -\log \pi(q) + \frac{1}{2} p' M^{-1} p$$

\Rightarrow 1. Need $\frac{\partial U}{\partial q_i}$ the gradient of U to run the Leapsong algorithm.

Symbolic differentiation?

deriv, D.

HMC per stat : At time t , we have $q(t)$

① Draw $p \sim N(0, M)$

② Make half step

$$p^{(\epsilon)} = p - \frac{\epsilon}{2} \frac{\partial U^*(q(t))}{\partial q}$$

③ Take L steps, $j=1, \dots, L-1$

$$q_i^{(t+\epsilon)} = q_i(t) + \epsilon p_i(t)$$

$$p_i^{(t+\epsilon)} = p_i(t) - \epsilon \frac{\partial U}{\partial q_i}(q(t))$$

④ $p = p - \frac{\epsilon}{2} \frac{\partial U}{\partial q}(q(t))$

⑤ $\alpha = \min \left\{ \frac{\exp(-H(q^*, p^*))}{\exp(-H(q, p))}, 1 \right\}$

(proposal or symmetric)

(P)

Check autocorrelation of your chain (marginally) using acf. Goal chains should have decreasing acf. Faster decrease indicates fast mixing.

$$\text{acf}(1) = \frac{\frac{1}{n} \sum (x_t - \bar{x})(x_{t+1} - \bar{x})}{\frac{1}{n} \sum (x_t - \bar{x})^2}$$

(P)

Computing clusters allow easy implementation of multiple chains. So why not?



Simulated Annealing

Suppose we want to minimize a function

$h(\theta)$, where θ is a vector of parameters, $\theta \in S$.

The basic idea is to simulate a MC with the stationary dist.

$$\pi_T(\theta) \propto \exp(-h(\theta)/T)$$

The parameter T is called the "temperature".

Markov Chains for "simulation", annealing (cooling) for minimization.

Let $S^* = \{\theta^* \in S : h(\theta^*) = \min_{\theta} h(\theta)\}$

the set of global minimizers

Define a distribution π^* on S^* so that

$$\pi^*(\theta) \propto 1 \quad \begin{cases} \text{for } \theta \in S^* \text{ and } 0 \\ \text{otherwise} \end{cases}$$

MC for simulation
"Annealing"

~~This (As $T \downarrow 0$) we have $\pi_T \rightarrow \pi^*$~~

~~$\pi_T(\theta) \rightarrow \pi(\theta)$ for s~~

FACT:

~~This~~ As $T \downarrow 0$, we have $\pi_T \rightarrow \pi^*$

Obviously, we'd like to simulate from π^* to get the global minimum, but we can't.

But for a fixed $T > 0$, we can simulate from π_T .

Use Metropolis algorithm with a symmetric proposal density, i.e. ~~$q(\theta_t | \theta_{t-1}) = q(\theta_{t-1} | \theta_t)$~~

$$q(\theta_n | \theta_{n-1}) = q(\theta_{n-1} | \theta_n)$$

Suppose we are in state n , θ_n ,

① Simulate $\theta \sim q(\theta | \theta_n)$

② Simulate $u \sim \text{Unif}(0, 1)$

③ Set $\theta_{n+1} = \theta$ if $u \leq \min\left(\frac{\exp(-h(\theta)/T)}{\exp(-h(\theta_n)/T)}, 1\right)$

$$= \min\left(\exp\left(-\underbrace{(h(\theta) - h(\theta_n))}_{\Delta h}\right)/T, 1\right)$$

Otherwise set $\theta_{n+1} = \theta_n$

④ Decrease T .

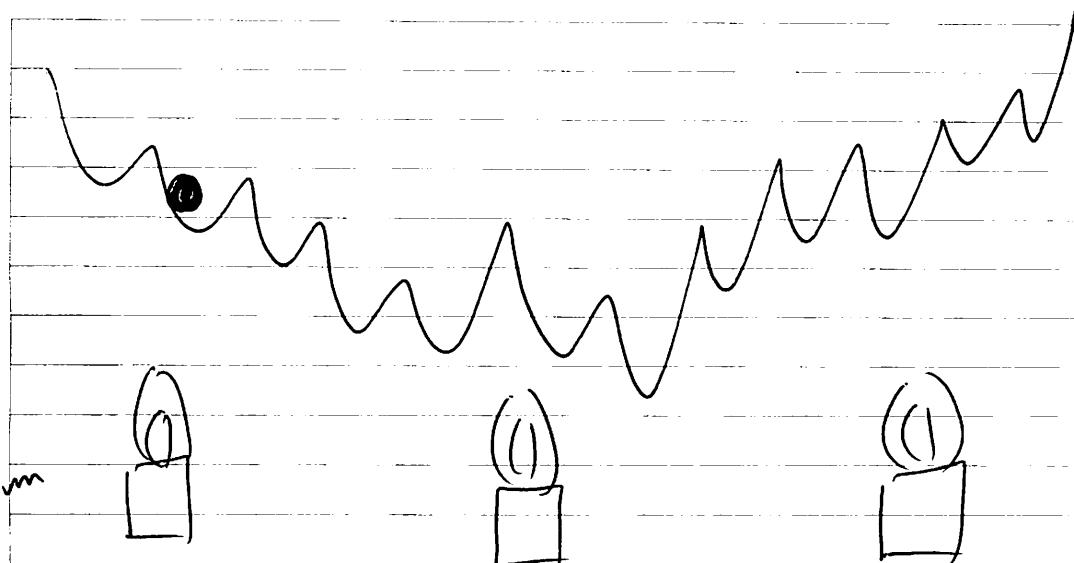
If ~~$\Delta h < 0$~~ , we always go
If $\Delta h > 0$, we go w/prob $\exp(-\Delta h/T)$

Hi temp
↳ ball moves all over

low temp
↳ ball stays in minimum

down hill

up hill



If $\Delta h < 0$, we always go

If $\Delta h > 0$, we go w/prob $\exp(-\Delta h/T)$

When $T \rightarrow$ closer to 0, it is less likely we will go up hill.

Can't cool too fast or will get stuck.

Need a cooling schedule like

$$T_n = \frac{q}{\log(n+b)} \quad (\text{for global convergence})$$

Anything else is "too fast" and will get stuck in a local mode.

T_n : If

$$\textcircled{1} \quad T_{n+1} < T_n \quad \forall n$$

$$\textcircled{2} \quad T_n \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

$$\textcircled{3} \quad \frac{T_n}{T_0} = \gamma / \log(n+b)$$

$$\text{Then } \|\pi_{T_n} - \pi^*\| \rightarrow 0$$

~~if we use a Metropolis step at L -th iterations, then we have SCMH.~~

~~Other possibilities are rejection sampling~~

K-Variable Gibbs sampler

let $(x_1^{(n)}, x_2^{(n)}, \dots, x_k^{(n)})$ be the state at ~~the~~ step n . The Gibbs sampler obtains the next state $(n+1)$ by

$$x_1^{(n+1)} \sim x_1 | x_2^{(n)}, x_3^{(n)}, \dots, x_k^{(n)}$$

$$x_2^{(n+1)} \sim x_2 | x_1^{(n+1)}, x_3^{(n)}, \dots, x_k^{(n)}$$

:

$$x_k^{(n+1)} \sim x_k | x_1^{(n+1)}, \dots, x_{k-1}^{(n+1)}$$

{Outline's
not!}



Notes on Gibbs Sampling



- ① It is sometimes useful to update groups of variables at a time (i.e. vector update rather than univariate). This is called "block Gibbs".
- ② Startly values are a guess. Can try frequentist values (means, medians, modes) or sample from priors. Maybe multiple starts
- ③ Good to use a burn in period before ~~taking averages~~ / storing values

- (4) Monitoring convergence is central
- look at marginal trace plots
 - try multiple chains with widely dispersed starting values
 - Monitor the acceptance rate ($\sim 30\%$ is good)
 - for > 2 dimensions, $\sim 50\%$ for 1, 2-D.
 - too low acceptance \rightarrow (slow convergence)
 - too high acceptance \rightarrow (slow mixing)

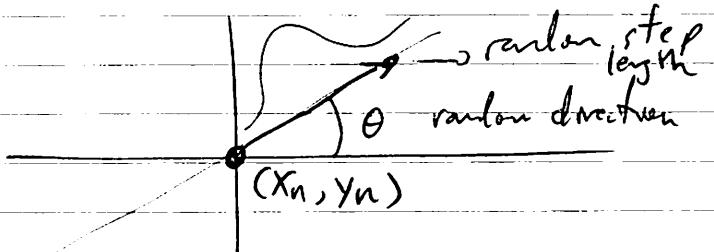
- (5) Order of updating need not be sequential. One can do 1, 2, 3, 4, 3, 2, 1, 2, 3 ...
- or randomly choose a component ("random scan")
 - randomly permute the order of updating
(1, 2, 3, 4), (1 3 4 2), (4 3 2 1) ...

- (6) Suppose X is multivariate and X_n is the current state. Choose a random direction e_n . Sample a scalar r from the density

$$p(r) \propto \pi(X_n + r e_n)$$

Set $X_{n+1} = X_n + r e_n$

In 2-D



- searches in more directions (than just along axes)
- how to sample from $p(r)$?
- using a Metropolis step results in random w/r Metropolis

(7) Gibbs sampler w/ improper priors is bad because it will run even if posterior is improper (no warning)

$$\text{ex } X \sim N(\theta, \sigma^2)$$

$$p(\theta | \sigma) = \frac{1}{\sigma^2}$$

$$\begin{aligned} & \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} p(x | \theta, \sigma) p(\theta | \sigma) d\theta dx = \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}\sigma^3} e^{-\frac{1}{2\sigma^2}(x-\theta)^2} \frac{1}{\sigma^2} d\theta dx \\ &= \int_{-\infty}^{+\infty} \frac{1}{\sigma^2} d\theta = \infty \end{aligned}$$

$\Rightarrow p(\theta | \sigma)$ is not a valid prior

But

$$\begin{aligned} p(\theta | \sigma, x) &\propto e^{-\frac{1}{2\sigma^2}(x-\theta)^2} = N(x, \sigma^2) \\ p(\sigma | \theta, x) &\propto \frac{1}{\sigma^3} e^{-(\theta-x)^2/2\sigma^2} \\ &= \frac{1}{\sqrt{\pi}} \text{ where } c \\ c &\sim \text{Exp}(\frac{(\theta-x)^2}{2}) \end{aligned}$$

Gibbs Sampler is easy to run!

\hookrightarrow This chain is "null recurrent": There is a stationary/mvariant function, but it is not a density.



(8) Your Gibbs Sampler may converge but that isn't always a good thing

Mantorny Convergence (Goldman + Rubin)

Start J
chains

let $x_j^{(0)}, x_j^{(1)}, \dots$ be the j^{th} Mankov Chain

(1) Discard $x_j^{(0)}, x_j^{(1)}, \dots, x_j^{(D-1)}$ values (burn in)

Use values $\underbrace{x_j^{(0)}, x_j^{(1)}, \dots, x_j^{(D-1)}}_{L \text{ values}}, \dots, x_j^{(D+L-1)}$

$$x_j^{(1)}, \dots, x_j^{(L)}$$

(2) Calculate

$$\bar{x}_j = \frac{1}{L} \sum_{t=1}^L x_j^{(t)} \quad (\text{chain mean})$$

$$\bar{x}_* = \frac{1}{J} \sum_{j=1}^J \bar{x}_j \quad (\text{grand mean})$$

$$B = \frac{1}{J-1} \sum_{j=1}^J (\bar{x}_j - \bar{x}_*)^2 \quad (\text{between chain variance})$$

~~$s_j^2 = \frac{1}{L-1} \sum_{t=1}^L (x_j^{(t)} - \bar{x}_j)^2$~~ (within chain variance)

$$W = \frac{1}{J} \sum_{j=1}^J s_j^2$$

(3) let $R = \frac{W + B}{W}$

$$\sqrt{R} < 1.2 \text{ acceptable?}$$

MCMC Diagnostics



You need to babysit your Gibbs sampler.
Look at output!



Try exhausting all iid sampling options



For estimation, one can use batch means.
Suppose we have chain x_1, x_2, x_3, \dots and we want $\mathbb{E}h(x)$. Then we have

$$h(x_1), h(x_2), h(x_3), \dots, h(x_K), h(x_{K+1}), \dots$$

Batch 1

$$\text{let } b_1 = \frac{1}{K} \sum_{i=1}^K h(x_i)$$

$$b_2 = \frac{1}{K} \sum_{i=K+1}^{2K} h(x_i)$$

Thus for
apart are
not only independent

Ergodic Theorem says

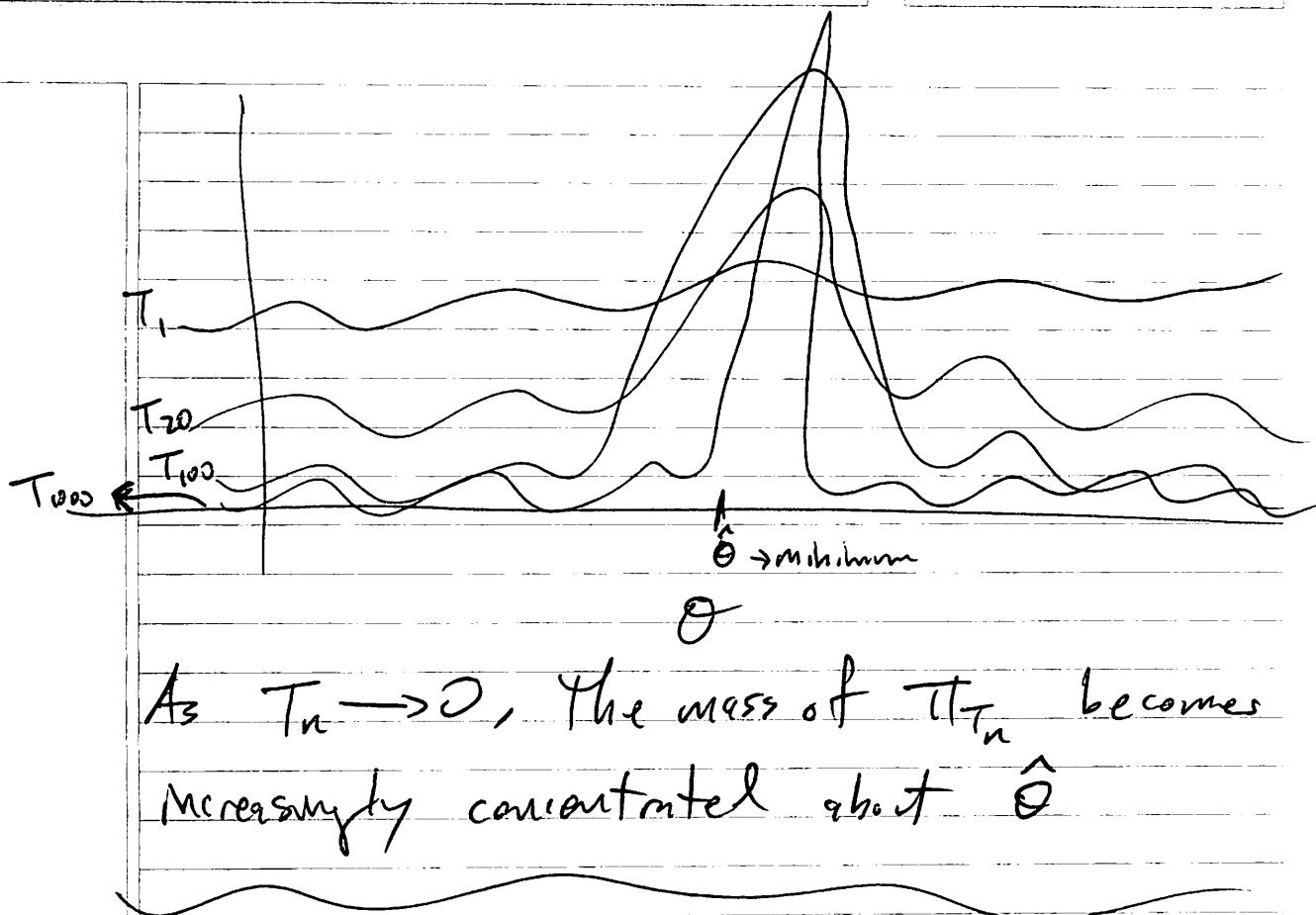
$$\bar{b} = \frac{1}{m} \sum_{i=1}^m b_i = \frac{1}{KM} \sum_{i=1}^{KM} h(x_i) \rightarrow \mathbb{E}[h(x)]$$

of K
batchs

The batches should be roughly independent
(by ergodic theorem) so for large m

$$\sqrt{m} \left(\frac{\bar{b} - \mathbb{E}h(x)}{s} \right) \rightarrow N(0, 1)$$

where $s^2 = \frac{1}{m} \sum (b_i - \bar{b})^2$, the sample variance
of b_i 's.



Perfect Sampling (brief intro)

We know that under certain conditions, a MC will converge to a stationary distribution.
 \hookrightarrow after time $n = \infty$

But we don't know when!

We know a burn-in is a good idea, but how much?

When can we declare "convergence"?

Propp + Wilson (1996) discovered that perfect random samples can be obtained in finite (but stochastic) time.

\Rightarrow Coupling from the past algorithm

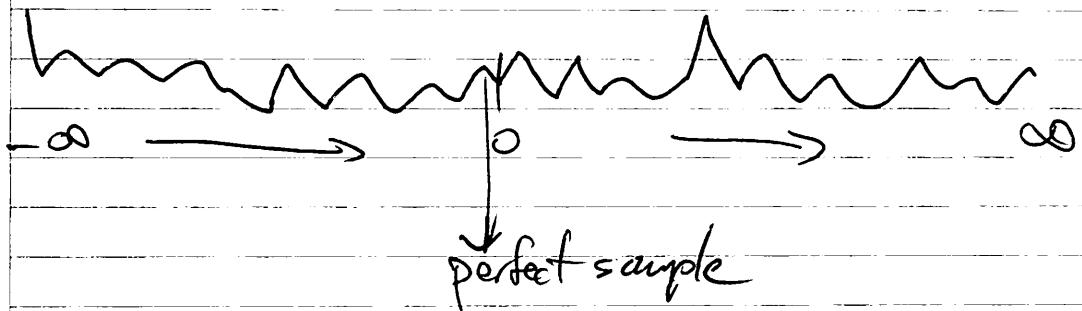
Basic idea
of CFTP

①

If we start at $n=0$, then a sample from state $n=\infty$ will be a "perfect" sample from π , the stationary distribution.

②

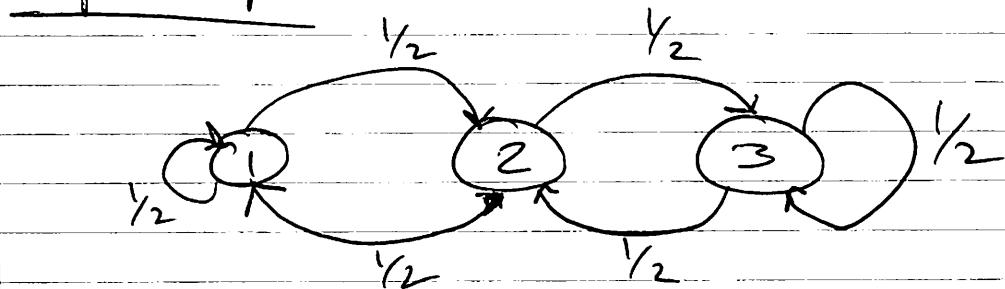
So why don't we start the chain at state $n=-\infty$, and ~~take a sample~~ use the value drawn at state $n=0$?



Both chains run for ∞ amount of time so must be stationary.

Simple example

Discrete
state space



Transition matrix

$$P = \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 0 & 1/2 \\ 0 & 1/2 & 1/2 \end{bmatrix}$$

Stationary dist: $\pi = (1/3, 1/3, 1/3)$

The stationary dist. is uniform on all 3 states.

How do we simulate this Markov Chain?

Suppose ~~X_n~~ X_n is the value in state n .

Draw a $U \sim \text{uniform}(0, 1)$.

If $X_n = 1$ then goto $\begin{cases} X_{n+1} = 1 & u \leq \frac{1}{2} \\ X_{n+1} = 2 & u > \frac{1}{2} \end{cases}$

If $X_n = 2$ then goto $\begin{cases} X_{n+1} = 1 & u \leq \frac{1}{2} \\ X_{n+1} = 3 & u > \frac{1}{2} \end{cases}$

If $X_n = 3$ then goto $\begin{cases} X_{n+1} = 2 & u \leq \frac{1}{2} \\ X_{n+1} = 3 & u > \frac{1}{2} \end{cases}$

Imagine running the chain in this way starting at $n = -\infty$, and ending at $n = 0$.

Can we emulate this process?

Notice:

① We can generate $U_n, U_{n+1}, U_{n+2}, \dots$ because they are independent of everything.

② $X_{n+1} = \phi(U_n, X_n)$

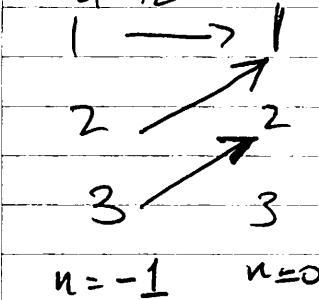
CFTP Algorithm for 1 perfect sample

~~Suppose we draw $U_1 \leq \frac{1}{2}$~~

Start 3 parallel MCs each starting in state 1, 2, or 3 respectively.

Suppose we draw $U_1 \leq \frac{1}{2}$ and start at $n = -1$

$$U_1 \leq \frac{1}{2}$$

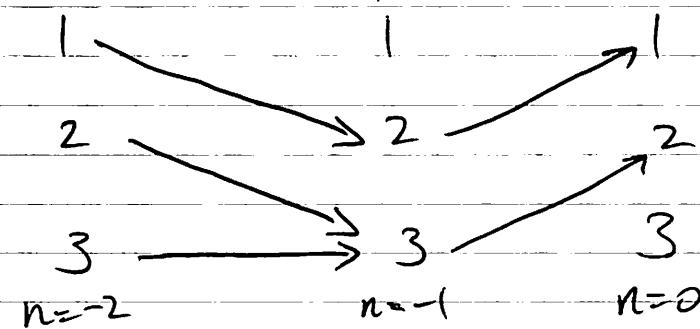


If we start the chains at $n = -1$ and draw $U_1 \leq \frac{1}{2}$, then we are at ~~state 1 or 2~~ state 1 or 2.
(not 3)

Suppose we generate $U_{-2} > \frac{1}{2}$ and start at $n = -2$.

$$U_{-2} > \frac{1}{2}$$

$$U_{-1} \leq \frac{1}{2} \text{ (fixed)}$$



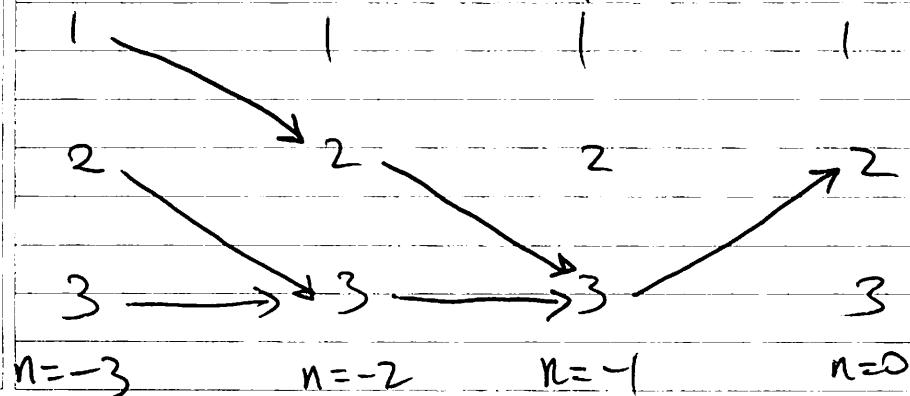
We "lost" the second chain, or the 2nd and 3rd chains coalesced.

Now suppose generate $U_{-3} > \frac{1}{2}$ and start at $n = -3$.

$$U_{-3} > \frac{1}{2}$$

$$U_{-2} > \frac{1}{2}$$

$$U_1 \leq \frac{1}{2}$$



When we start from $n = -3$, all of the chains coalesce at state $n = -1$.

~~BB Negg~~

\Rightarrow Regardless of what happens b/w ~~times~~ $n = -\infty$ and $n = -3$ (we don't know), we know that at ~~some~~ ^{time} $n = -1$ the will all coalesce at ~~at~~ state ③.

\Rightarrow Regardless of starting value, at the $n = 0$, we will be at state ②.

\Rightarrow ② is a perfect sample from stationary dist.

\Rightarrow Coalescence "breaks" dependence on the starting value ~~exactly~~

Note: In this chain, the arrows don't cross, so we only need to keep track of "top" chain and "bottom" chain

① Why can't run forward until copying?
↳ only copies at ① and ③

② Why can't run until copying and then a few more?
How many more?

Image Analysis (Simple)

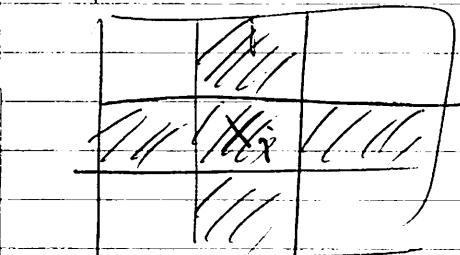
Let $X = (X_1, X_2, \dots, X_n)$ be an image (tree)
where X_i is the value for pixel i .

Suppose X_i is a r.v. such that

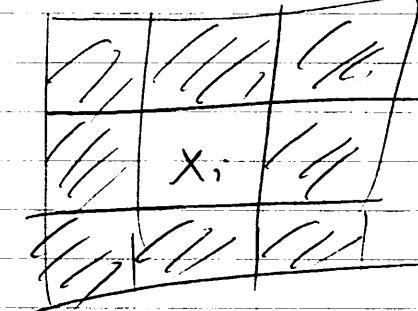
$$X_i | X_{S_i} \stackrel{D}{=} X_i | X_{S_i}$$

where S_i is a neighbourhood of x_i

first order S_i



2nd order S_i



$$\text{So } p(x_i | X_{S_i}) = p(x_i | X_{S_i})$$

Let $X_i = 0$ or 1 both with positive probability (binary image).

\Rightarrow for a 40×40 image, the state space is of size 2^{1600} !

Uniquely characterize joint distributions of X
(Hammersley-Clifford)

Let y_i be a noisy version of x_i
and assume that

$$p(y_1 \rightarrow y_n | x_1 \rightarrow x_n) = \prod_{i=1}^n p(y_i | x_i)$$

e.g. white noise process

\rightarrow likelihood for $x_1 \rightarrow x_n$ given $y_1 \rightarrow y_n$

We want to sample $p(x|y) \propto p(y|x) p(x)$

One class of priors for x is

$$p(x) \propto \exp\left(\sum_{i \neq j} \phi(x_i - x_j)\right)$$

where ϕ is symmetric about 0.

Using a first order neighbourhood structure for x_i

we can Gibbs sample each pixel

$$x_i^{(t+1)} | x_{-i}^{(t)}, y \sim p(x_i | x_{S_i}^{(t)}, y)$$

full condition

$$\propto p(y | x_{S_i}^{(t)}, x_i^{(t)}) p(x_i^{(t)}, x_i^{(t)})$$

$$p(y_i | x_i^{(t)}) \cancel{p(x_i^{(t)})} p(x_i^{(t)} | x_{-i}^{(t)}) \cancel{p(x_{-i}^{(t)})}$$

$$p(x_i^{(t)} | x_{S_i}^{(t)})$$

Ex: Binary Images

Given observed values $y_1 \rightarrow y_n \in \{0, 1\}$, we want to estimate true image $x_1 \rightarrow x_n \in \{0, 1\}$

$\otimes x_i = 1$ indicates presence of a species

$y_i = 1$ observed a species

$$\text{Let } p(y|x) \propto \exp\left(\alpha \sum \mathbb{1}\{y_i = x_i\}\right)$$

Species

{ Use pairwise difference prior for x ,

$$p(x) \propto \exp\left(\beta \sum_{i \sim j}^n \mathbb{1}\{x_i = x_j\}\right)$$

$$p(x_i^{(t+1)} = 1 | x_{-i}^{(t)}, y) \propto p(y|x_i) p(x_i|x_{-i})$$

$$\propto \exp\left(\alpha \mathbb{1}\{y_i = x_i\}\right) \cancel{\exp\left(\beta \sum_{i \sim j} \mathbb{1}\{x_i = x_j\}\right)}$$

$$\times \exp\left(\beta \sum_{i \sim j} \mathbb{1}\{x_i = x_j\}\right)$$

$$= \exp\left(\alpha \mathbb{1}\{y_i = x_i\} + \beta \sum_{j \in S_{x_i}} \mathbb{1}\{x_j = x_i\}\right)$$

$$\frac{\exp(\alpha \mathbb{1}\{y_i = 0\} + \beta \sum_{j \in S_{x_i}} \mathbb{1}\{x_j = 0\}) + \exp(\alpha \mathbb{1}\{y_i = 1\} + \beta \sum_{j \in S_{x_i}} \mathbb{1}\{x_j = 1\})}{\exp(\alpha \mathbb{1}\{y_i = 0\} + \beta \sum_{j \in S_{x_i}} \mathbb{1}\{x_j = 0\}) + \exp(\alpha \mathbb{1}\{y_i = 1\} + \beta \sum_{j \in S_{x_i}} \mathbb{1}\{x_j = 1\})}$$

$$= \left[1 + \exp\left(\alpha (\mathbb{1}\{y_i = 0\} - \mathbb{1}\{x_i = 1\}) + \beta \sum_{j \in S_{x_i}} (\mathbb{1}\{x_j = 0\} - \mathbb{1}\{x_j = 1\})\right) \right]^{-1}$$

Floating Point Numbers

$$n = f \times 10^e \rightarrow \text{exponent}$$

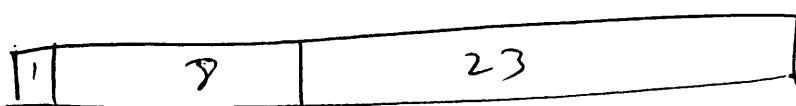
↓
fraction/Mantissa

$$3.14 = 3.14 \times 10^0$$

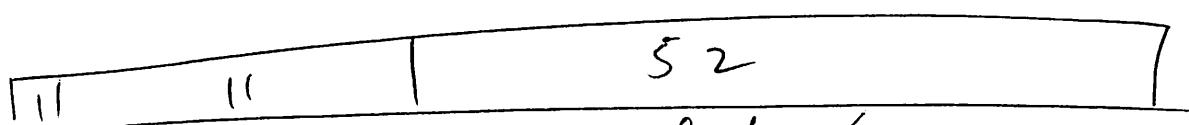
$$1941 = 1.941 \times 10^3$$

Late 1970s

IEEE 754
Standard



Single prec.



Double prec.

Sign exponent
 ↓
 1 = negative
 0 = positive
 excess 127
 excess 1023

fraction/
significant

$$\sqrt{2} \times \sqrt{2} = ?$$

$$a \leftarrow \sqrt{2}$$

$$a * a = 2 \quad \text{False}$$

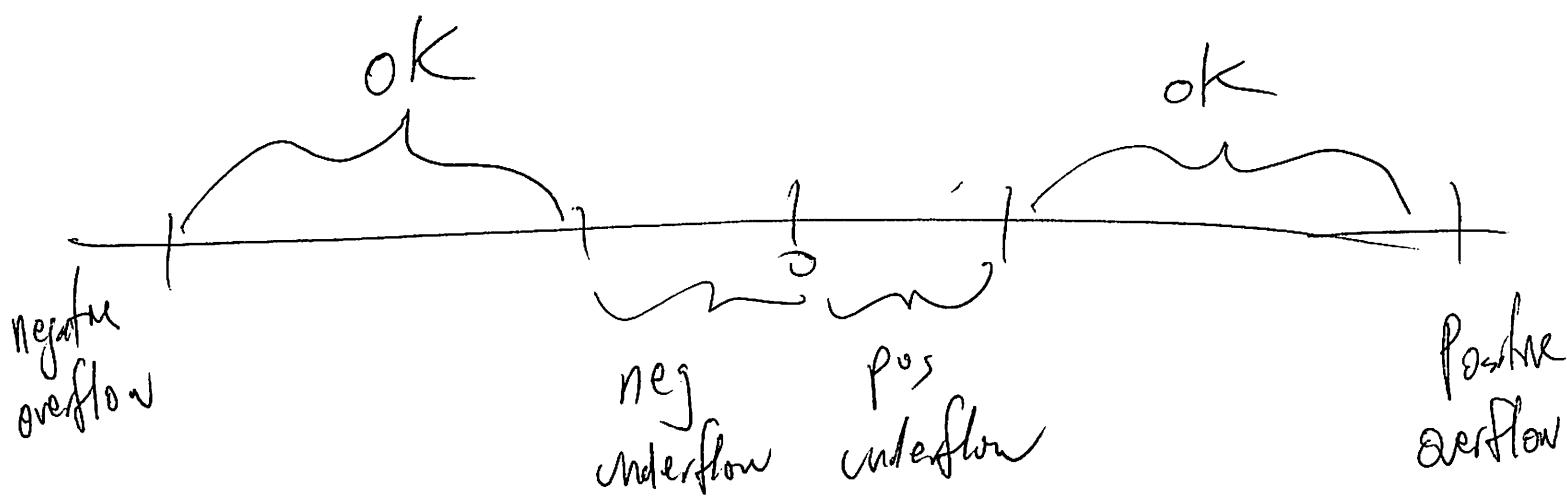
Inf = exp \Rightarrow all 1s
 Fraction = 0

$$\times 10^{-16}$$

rounding

NaN = exp \Rightarrow all 1s
 Fraction = any non-zero pattern

	single	double
smallest	2^{-126}	2^{-1022}
largest	2^{128}	2^{1024}



Linear Models:

$$Y = f(X) = \sum_{j=1}^P \beta_j X_j = \beta' X$$

Additive Model

$$f(X) = f(X_1, \dots, X_p) = \sum_{j=1}^P s_j(x_j)$$

More General:

$$f(X) = \sum_{k=1}^M g_m \underbrace{(w_m' X)}_{V_m = \text{derived feature}} \quad \begin{array}{l} \text{Projection} \\ \text{Pursuit} \end{array}$$

Linear Regression: $M=1$, $g_m(x) = \underline{\underline{w}}' X$, w_m unknown

Additive model: $M=p$, g_m smooth, $w_m = (0, 0, 1, \dots)$

PPR = Additive model + non-coordinate axis directions

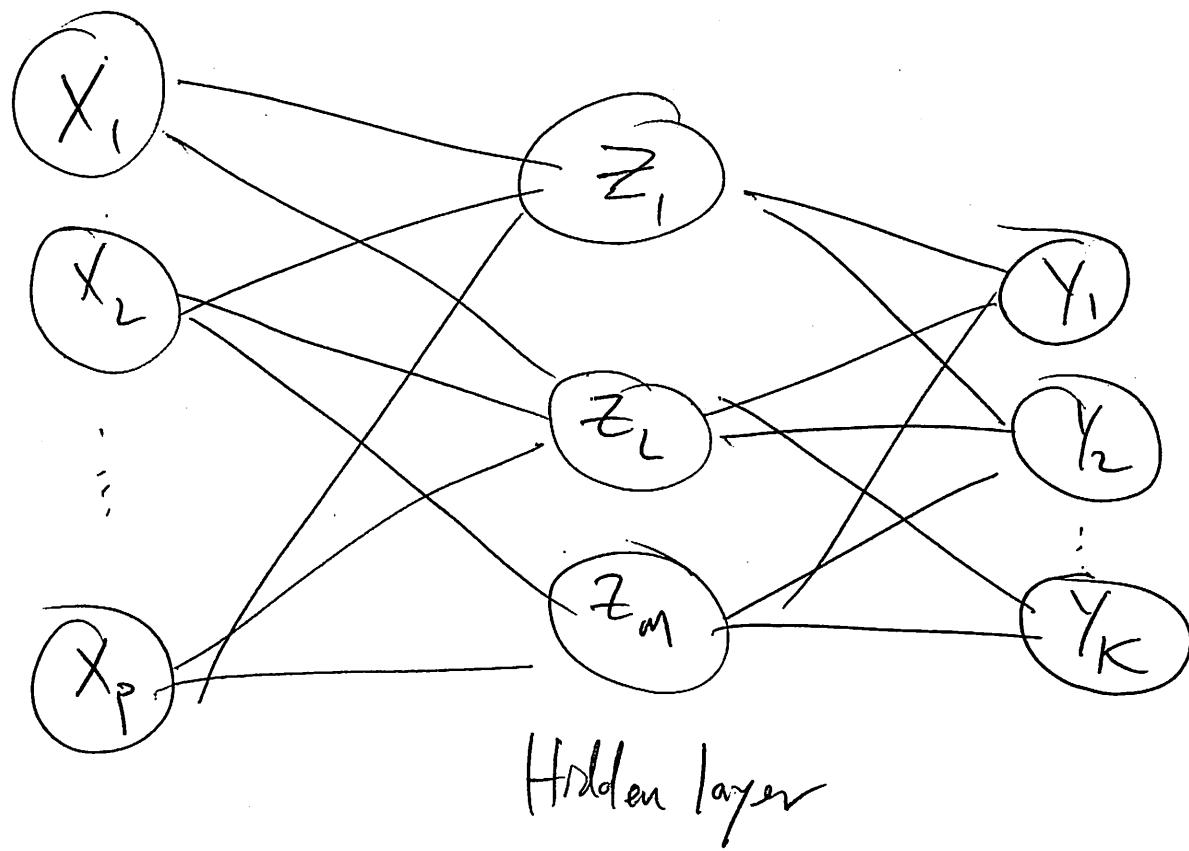
PCR: $M < p$, $g_m(x) = \beta_m X$, w_m = PC directions
only dep on X

Neural Network: $M = \text{large}$; w_m unknown

w/ 1 hidden layer

$$g_m(x) = \beta_m \sigma(x_{0m} + \|x_m\|_2 x)$$

$$\sigma(x) = \frac{1}{1+e^{-x}} \quad (\text{logistic})$$



Fitting General projection pursuit: estimate w_m and g_m

$$\cancel{f(w'x)} \approx g(w_0'x) + g'(w_0'x)(w - w_0)$$

$$f(x_i) = \sum_{m=1}^M g_m(w_m' x_i)$$

$$\text{Minimize } \sum_{i=1}^N (y_i - f(x_i))^2 = \sum_{i=1}^N \left(y_i - \sum_{m=1}^M g_m(w_m' x_i) \right)^2$$

For fitting g_m, w_m

$$\textcircled{1} \text{ Let } r_i = y_i - \sum_{k \neq m} g_k(w_k' x_i)$$

$$\textcircled{2} \text{ Approximate } g_m(w' x_i) \doteq g(w_0' x_i) + g'(w_0' x_i)(w - w_0)' x_i$$

$$\textcircled{3} \sum_{i=1}^N (r_i - g_m(w' x_i))^2 \doteq \sum_{i=1}^N \left(r_i - g(w_0' x_i) - g'(w_0' x_i)(w - w_0)' x_i \right)^2$$

$$\cancel{\sum_{i=1}^N g'(w_0' x_i)^2 \left[\left(w_0' x_i + \frac{r_i - g(w_0' x_i)}{g'(w_0' x_i)} \right) - w' x_i \right]^2}$$

$$\sum_{i=1}^N g'(w_0' x_i)^2 \left[\left(w_0' x_i + \frac{r_i - g(w_0' x_i)}{g'(w_0' x_i)} \right) - w' x_i \right]^2$$

$\textcircled{4}$ Linear regression

④ Given w_m , estimate g_m with
Smooth of r_i on $g_m(w_m'x_i)$ or
Simple linear regression.

⑤ Go back to ① until convergence
strategy:

- ① Start with $M = 1$
- ② ~~Use~~ Build M new terms in forward stage-wise manner (Fix previous terms)
- ③ Use CV to estimate OOS error
- ④ Increase M until error stops decreasing

~~Neural Networks Fitting~~

$$z_m = \sigma(\alpha_{0m} + \alpha_m' X) \quad m=1, \dots, M$$

Regression: $f(x) = \beta_0 + \beta_1 z_m, \quad z = (z_1, \dots, z_M)$

Classification: $f(x) = \frac{e}{1+e}$

Neural Networks

$$Y_1 \rightarrow Y_K$$

$$Z_m = \sigma(a_{0m} + \alpha_m^T X) , \quad m=1 \rightarrow M$$

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

$$T_K = \beta_{0K} + \beta_K^T Z , \quad Z = (Z_1 \rightarrow Z_M)$$

$$T = (T_1 \rightarrow T_K)$$

Regression : $f_K(x) = g_K(T_K) = \beta_{0K} + \beta_K^T Z = \beta_{0K} + \sum_{m=1}^{Mx} \beta_{km} \sigma(\)$

Classification: $f_K(x) = \frac{e^{T_K}}{\sum_{l=1}^K e^{T_l}} = \underbrace{\text{Softmax}}$

y_i

$$f_K(x) = \frac{e^{T_i}}{e^{T_1} + e^{T_0}} = \underbrace{\text{function}}$$

Error Function

~~$$f_1 = \frac{1-f}{2} R(\beta, \alpha) = \sum_{K=1}^K \sum_{i=1}^N (y_{ik} - f_k(x_i))^2$$~~

$$R(\beta, \alpha) = \sum_{i=1}^N \sum_{K=1}^K (y_{ik} - f_k(x_i))^2$$

~~Recall~~

Fitting NN via back propagation:

Recall $X_i = i^{\text{th}} \text{ input}, i \in \{1, \dots, N\}$

$$z_{im} = \sigma(\alpha_m + \beta_m' X_i)$$

$$z_i = (z_{i1}, \dots, z_{iM})$$

$$z = \begin{bmatrix} z_{i1} \\ \vdots \\ z_{iM} \end{bmatrix}$$

The error function $R(\beta, \alpha)$, ~~is not gradient~~

$$R(\beta, \alpha) = \sum_{i=1}^N R_i = \sum_{i=1}^N \left[\sum_{k=1}^K \underbrace{(y_{ik} - f_k(x_i))^2}_{\text{prac}} \right]$$

$$f_k(x_i) = g_k(\beta_k' z_i)$$

$$\frac{\partial R_i}{\partial \beta_{km}} = \underbrace{-2(y_{ik} - f_k(x_i)) g'_k(\beta_k' z_i)}_{\delta_{ki}} z_{im}$$

$K=1 \dots K$
 $M=1 \dots M$

$$\frac{\partial R_i}{\partial \alpha_{ml}} = \sum_{k=1}^K -2(y_{ik} - f_k(x_i)) g'_k(\beta_k' z_i) \beta_{km} \sigma'(\alpha_m' X_i) x_{il}$$

$l=1 \dots p$

δ_{mi}

Gradient descent update 13

$$\beta_{km}^{(t+1)} = \beta_{km}^{(t)} - \gamma_t \sum_{i=1}^N \frac{\partial R_i}{\partial \beta_{km}^{(t)}}$$

$$\alpha_{ml}^{(t+1)} = \alpha_{ml}^{(t)} - \gamma_t \sum_{i=1}^N \frac{\partial R_i}{\partial \alpha_{ml}^{(t)}}$$

γ_t is learning rate (step / cyth)

~~($\gamma_r = \gamma_t$ for backpropagation)~~

~~($\text{Pr}(Y=1)$)~~

~~1 - p~~
If we write

$$f_{im} = \sigma(\alpha_{0m} + \alpha_n x_i)$$

$$\frac{\partial R_i}{\partial \beta_{km}} = S_{ki} z_{im}$$

$$\frac{\partial R_i}{\partial \alpha_{ml}} = S_{mi} x_{il}$$

$$S_{ui} = \sum_{k=1}^K \delta_{ki} \beta_{km} \sigma'(\alpha_m' x_i)$$

$$= \sigma'(\alpha_m' x_i) \cdot \sum_{k=1}^K \delta_{ki} \beta_{km}$$

Back propagation
equations

$$\sum_{m_i}^M \sum_{k=1}^K \delta_{k,i} f_{km} \circ (\cancel{\alpha_{k,i}}) \beta_{ml}$$

~~Forward pass:~~

- ① Fix weights, compute $\hat{f}_k(x_i)$
- ② Compute errors $\delta_{k,i}$ and then $\sum_m \delta_{mi}$
- ③ Compute gradients $\frac{\partial R_i}{\partial \beta_{km}}$, $\frac{\partial R_i}{\partial \alpha_{me}}$
- ④ Update β_{km} , α_{me}
- ⑤ Go to ①

Learning rate γ_t can be fixed

Online learning

$$\beta_{km}^{(t+1)} = \beta_{km}^{(t)} - \gamma_t \frac{\partial R_i}{\partial \beta_{km}^{(t)}} \quad \text{for each } i=1, \dots, N$$

$$\alpha_{me}^{(t+1)} = \alpha_{me}^{(t)} - \gamma_t \frac{\partial R_i}{\partial \alpha_{me}^{(t)}}$$

For online learning γ_t must satisfy

$$\gamma_t \rightarrow 0 \quad t \rightarrow \infty$$

$$\sum \gamma_t = \infty$$

$$\sum \gamma_t^2 < \infty$$

$\gamma_t = \gamma_t$ works.

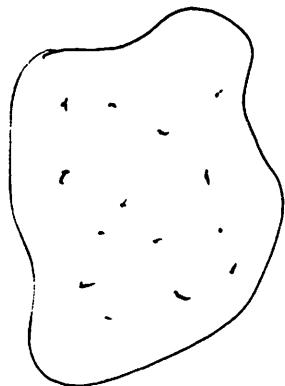
Penalization

① Stop early

② $J(f, \alpha) = \sum_{km} \beta_{km}^2 + \sum_{ml} \alpha_{ml}^2$

Let $R(f, \alpha) = R(\beta, \alpha) + \lambda J(\beta, \alpha)$

Ridge
penalty



Spatial Modelling

$$Y(s) = \mu(s) + w(s) + \varepsilon(s)$$

fixed effect → measurement error?
 microscale noise?

$$\mu(s) = X^T(s)\beta \rightarrow \begin{array}{l} \text{latent Gaussian} \\ \text{Spatial process} \\ \text{meaning?} \end{array}$$

$$w(s) \sim N(0, \sigma^2 H(\phi))$$

$$\varepsilon(s) \sim N(0, \tau^2)$$

s || range

↳ nugget effect

What's meaning
of $w(s)$?

$$H_{ij} = \text{Cov}(w(s_i), w(s_j))$$

$$= \rho(s_i - s_j; \phi) \quad \text{anisotropic}$$

$$= \rho(\|s_i - s_j\|; \phi) \quad \text{isotropic}$$

e.g. $\rho(s_i - s_j, \phi) = \exp(-\phi \|s_i - s_j\|)$

$$\text{let } Y = (Y(s_1), Y(s_2), \dots, Y(s_n))$$

$$\Theta = (\beta, \sigma^2, \phi, \tau^2)$$

$$Y|\theta, w \sim N(X\beta + w, \tau^2 I)$$

$$w | \sigma^2, \phi \sim N(0, \sigma^2 H(\phi))$$

$$\begin{matrix} \sigma^2 \\ \phi \\ \beta \\ \tau^2 \end{matrix} \sim \text{Priors}$$

Get posteriors

~~$\theta(\theta, \phi)$~~

$$p(y, \theta, w) \propto p(y|\theta, w) p(w|\theta) p(\theta)$$

$p(\alpha^y) p(\tau^z) p(\phi) p(\beta)$

w

$$p(\beta|w) \propto p(y|\beta, \tau^z, w) p(\beta)$$

informative

$$N(x_\beta + w, \tau^z I) N(0, D)$$

$$= \mathcal{N}(x_{\beta} + \frac{x_{\alpha^y} - x_{\beta}}{\tau^z}, (I - S)D)$$

$S = D[\tau^z I + D]^{-1}$

$$p(\tau^z|w) \propto p(y|\beta, \tau^z, w) p(\tau^z)$$

$$= \text{IG}(a_2+1, b_2 + \frac{1}{2}(y - x_\beta - w)^T(y - x_\beta - w))$$

$$p(\sigma^2|w) \propto p(w|\alpha^z, \phi) p(\sigma^2)$$

$$\text{IG}(a_1+1, b_1 + \frac{1}{2}w^T w)$$

$$p(\phi|w) \propto p(w|\sigma^2, \phi) p(\phi)$$

$$N(0, \sigma^2 H(\phi)) \phi^{e-1} e^{-\phi/f}$$

$$\frac{1}{\sigma^2 H(\phi)} \exp\left(-\frac{1}{2\sigma^2} W^T H(\phi)^{-1} W\right)$$

$$H(\phi) = U(\phi) D(\phi) U^T$$

$$D(\phi) = \begin{pmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{pmatrix}$$

need Metropolis

step'

~~$\frac{1}{\prod_{i=1}^n \sigma^2 \lambda_i} \prod_{i=1}^n \lambda_i^{-1/2}$~~

$$\frac{1}{|H(\phi)|} = \frac{1}{\prod_{i=1}^n \lambda_i} = \prod_{i=1}^n \lambda_i$$

$$\left| \frac{1}{\prod_{i=1}^n \sigma^2 \lambda_i} \right|^{-1/2} \exp\left(-\frac{1}{2\sigma^2} \underbrace{W^T U(\phi) D(\phi)^{-1} U(\phi)^T W}_{} \right)$$

$$p(w|m) \propto p(y|\theta, w) p(w|\theta)$$

$$\begin{pmatrix} y_{\lambda_1} & \dots & y_{\lambda_n} \end{pmatrix}$$

$$p(\theta|y) \propto p(y|\theta, w) p(w|\theta) p(\theta)$$

$$\left(\begin{array}{l} \tau^* \sim \text{IG}(\alpha_1, b_1) \\ \tau^* \sim \text{IG}(\alpha_2, b_2) \quad \phi \sim \text{Gamma}(\alpha_3, b_3) \\ \beta \sim N(0, \Sigma) \quad (\text{Vague}) \\ \phi \sim \cancel{\text{IG}(\alpha_4, b_4)} \end{array} \right) \quad \text{informative}$$

Bayesian way

$$y|x \sim N(x, \tau^2)$$

$$x \sim N(0, \sigma^2)$$

$$x|y \sim N(y, (I - B)^T \Sigma)$$

$$\hat{\beta}|\beta \sim N(\beta, \Sigma)$$

$$\beta \sim N(0, \Sigma)$$

$$\beta|\hat{\beta} \sim N(\hat{\beta}, (I - C)^T \Sigma)$$

~~q(t)~~ Compute MLEs

$$\text{CDL} = p(y|\theta, w) p(w|\theta)$$

$$N(X\beta + w, \tau^2 I) \quad N(0, \sigma^2 H(\phi))$$

Need

$$p(w|\theta, y) \propto p(y|\theta, w) p(w|\theta)$$

$$\begin{aligned} \text{CDLL} = & -\frac{1}{2} \log |\tau^2 I| - \frac{1}{2} (Y - X\beta - W)^T (Y - X\beta - W) / \tau^2 \\ & - \frac{1}{2} \log |\sigma^2 H(\phi)| - \frac{1}{2} W^T H(\phi)^{-1} W / \sigma^2 \end{aligned}$$

$$p(w|\theta, y) = N(B y, (I - B)^T \Sigma H(\phi))$$

$$B = \Sigma H(\phi)^{-1} [\tau^2 I + \Sigma H(\phi)]^{-1}$$

Can iterate EM algorithm directly by filling in w and max Q. Estimating β, τ^2 okay but σ^2 and ϕ are harder (ϕ nonlinear)

Prediction at new location.

Suppose we have $x(s_0)$ and we want $y(s_0)$

↳ Reg (Bayesian) Kriging, Interpolation

We want predictive distribution;

$$p(y(s_0) | Y, X, x(s_0))$$

\downarrow
obs y \nwarrow
obs X

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \\ y_0 \end{pmatrix} \sim N \begin{pmatrix} x_1 - \\ x_2 - \\ \vdots \\ x_n - \\ x_0 - \end{pmatrix}$$

$$= \int p(y_0 | Y, X, \theta, x_0) \underbrace{p(\theta | Y, X)}_{\text{posterior}} d\theta$$

$$= \int \left[\int p(y_0 | Y, X, \theta, x_0, w) p(w | \theta) dw \right] p(\theta | Y, X) d\theta$$

If both normal, just

$$use p(y_0 | Y, X, \theta, x_0)$$

$$= N(X\beta, \tau^2 I + \sigma^2 H(\phi))$$

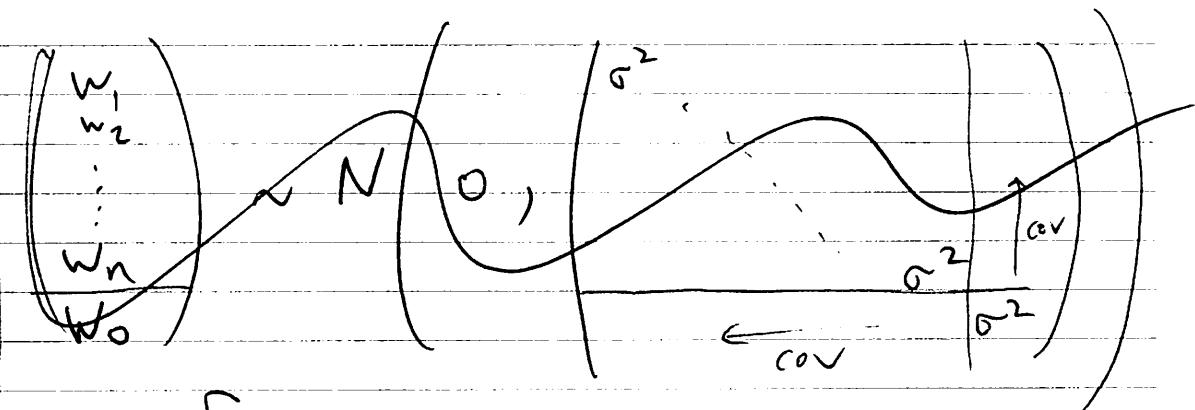
But if $p(y | \theta, w)$ is not normal

(maybe Poisson, Bernoulli) Then we are stuck,

need to sample.

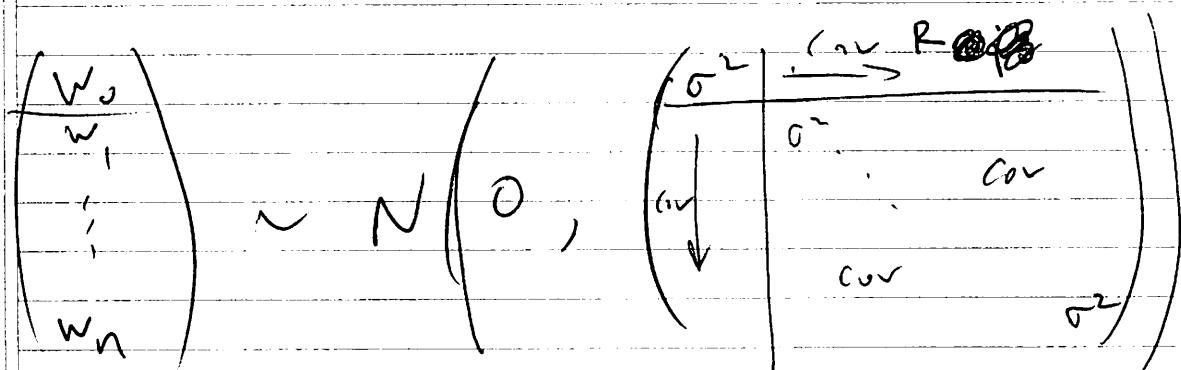
$$p(w_0, w | \theta)$$

$$\sim p(w_0 | w, \theta) p(w | \theta)$$

y_1 

$$W_0 | W_1, \dots, W_n \sim \text{multivariate normal}$$

$$N(0 + B w, (I - B)^{-1} \sigma^2 I_{n+1})$$



$$W_0 | W_1, \dots, W_n \sim N(R^T H(\phi)^{-1} W, \sigma^2 - R^T H(\phi)^{-1} R)$$

Forward Stagewise Additive Modeling

For the regression problem

$$y_i = f(x_i) + \varepsilon_i$$

y_i = continuous
 $\varepsilon_i \sim \mathcal{E}(0, \sigma^2)$

We want to minimize ~~$\sum L(y_i, f(x_i))$~~

x_i = p-dim vector

$$\sum_{i=1}^N L(y_i, f(x_i))$$

~~over all~~ f in some class. For model f , often useful to use an additive model

$$f(x) = \sum_{m=1}^M \beta_m b(x | \gamma_m)$$

Regression: $b(x | \gamma_m) = x_m$

~~Neural Net.~~: $b(x | \gamma_m) = r(\gamma_{0m} + \gamma_{1m}' x)$

Proj Pursuit: $b(x | \gamma_m) = g_m(\gamma' x)$ (g_m smooth)

Trees: ~~parametrizes~~ $b(x | \gamma_m) = \mathbb{I}\{x \in R_m(\gamma_m)\}$

Therefore:

$$\sum_{i=1}^N L\left(y_i, \sum_{m=1}^M \beta_m b(x_i; \gamma_m)\right) \quad \textcircled{*}$$

Loss functions

$$L(y, f(x))$$

least squares :

$$(y - f(x))^2$$

Robust regression:

$$|y - f(x)|$$

Deviance :

$$\log(1 + \exp(-y f(x)))$$

Exponential:

$$\exp(-y f(x))$$

Support vector:

$$(1 - y f(x))_+$$

\Rightarrow Solving $\textcircled{*}$ Globally can be hard

\Rightarrow What if Solving $\sum_{i=1}^N L(y_i, \beta b(x_i; \gamma_{**}))$
is easy?

Forward stagewise Additive Model

① Initialize $f_0(x) = 0$

② For $m = 1, \dots, M$

③ $(\beta_m, \gamma_m) =$

$$\arg \min_{\beta, \gamma} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + \beta b(x_i | r))$$

④ Set $f_m(x) = f_{m-1}(x) + \beta_m b(x | \gamma_m)$

\Rightarrow An approximation (greedy) to the gl. ls problem.

\Rightarrow Not backfitting b/c previously add terms not modified

Linear regression w/ least squares

$$L(y_i, f_{m-1}(x_i) + \beta b(x_i | \gamma))$$

$$= (\underbrace{y_i - f_{m-1}(x_i) - \beta b(x_i | \gamma)}_{r_i})^2$$

~~x_i~~ ~~β~~

$$(\underbrace{y_i - \beta'_{m-1} X_{im}}_{r_i} - \beta X_{im})^2$$

$$\hat{\beta} = \frac{\text{Cov}(X_m, r)}{\text{Var}(X_m)}$$

Boosting (Ada Boost)

$$L(y, f(x)) = \exp\left(-\underbrace{y f(x)}_{\text{margin}}\right)$$

$$f(x) = \sum_{m=1}^M \underbrace{\alpha_m G_m(x)}_{\text{Individual trees/ classifiers}}$$

\Rightarrow We must solve at each $m = 1, \dots, M$

$$(\beta_m, G_m) = \underset{\beta, G}{\operatorname{arg\,min}} \sum_{i=1}^N \exp\left[Y_i (f_{m-1}(x_i) + \beta G(x_i))\right]$$

$$= \underset{\beta, G}{\operatorname{arg\,min}} \sum_{i=1}^N w_i^{(m)} \exp\left(-Y_i \beta G(x_i)\right)$$

G_m minimizes $\sum_{i=1}^N w_i^{(m)} \mathbb{1}\{Y_i \neq G(x_i)\}$ for $\beta > 0$

$$\beta_m = \frac{1}{2} \log \frac{1 - \text{err}_m}{\text{err}_m}, \quad \text{err}_m = \frac{\sum w_i^{(m)} \mathbb{1}\{Y_i \neq G_m(x_i)\}}{\sum w_i^{(m)}}$$

$$f_m(x) = f_{m-1}(x) + \beta_m G_m(x)$$

Gradient Boosting

$$\textcircled{1} \quad f_0(x) = \arg \min_f \sum_{i=1}^N L(y_i, f)$$

\textcircled{2} For $m = 1, \dots, M$

$$\textcircled{a} \quad \text{Compute } r_{im} = -\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \Big|_{f=f_{m-1}}$$

\textcircled{b} regress r_{im} on predictors to produce

J_m regions of leaf nodes

$$\Rightarrow \hat{f}_m(x) \quad \forall x$$

$$J_{im} = \arg \min_j \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + j)$$

$$\textcircled{d} \quad f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} r_{jm} \mathbf{1}\{x \in R_{jm}\}$$

$$\textcircled{e} \quad \hat{f}(x) = f_M(x)$$

$$\begin{aligned} & \sqrt{\sum_{i=1}^n (y_i - \hat{f}(x_i))^2} \\ & \approx \sqrt{\sum_{i=1}^n (y_i - f_M(x_i))^2} \end{aligned}$$

Bayes

- ① For $b = 1 \rightarrow B$
 - a) Draw z^* w/ replace
 - b) Fit tree to z^* to get T_b

- ② Output $T_b, 1 \rightarrow B$

$$\hat{f}(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$$

not iid \Rightarrow id

$$\text{Var}\left(\frac{1}{B}(X_1 \rightarrow X_B)\right) = \rho \sigma^2 + \frac{1-\rho}{B} \sigma^2$$

If iid $\text{Var}\left(\frac{1}{B}(X_1 \rightarrow X_B)\right) = \frac{1}{B} \sigma^2$

Random Forest

Given $(Y_i, X_i) \quad i=1 \rightarrow N$

① For $b = 1 \rightarrow B$

④ Draw bootstrap sample Z^* of size N
from train data

⑤ RF tree T_b

① Select m vars from p vars

② Pick best split variable

③ Split node into 2 nodes

④ etc.

② Output T_b for $i \rightarrow B$

$$\hat{f}(x) = \frac{1}{B} \sum T_b(x) \quad \text{Regression}$$

$$\hat{f}(x) = \text{Majority vote } \left\{ T_b(x) \right\}_1^B \quad \begin{matrix} \text{Sample} \\ \text{n° explain!} \end{matrix}$$

Use 50B samples to estimate error rate