

MONTE CARLO INFERENCE

$$x \sim p(x|D)$$

GENERATE UNWEIGHTED SAMPLES FROM THE POSTERIOR \rightarrow USE THEM TO COMPUTE STUFF EG MARGINALS, LIKELIHOOD OR POSTERIOR PREDICTIVE

STANDARD SAMPLING

• **INVERSE PROBABILITY TRANSFORM:** SAMPLE FROM UNIFORM $U \sim U(0,1)$, MAP TO DISTRIBUTION OF INTEREST WITH $F^{-1}(U)$

\rightarrow **EXPONENTIAL:** $F^{-1}(U) = -\ln(U)/\lambda$

\rightarrow **2D GAUSSIAN:** SAMPLE FROM $z_1, z_2 \in (-1,1)$. REJECT $z_1^2 + z_2^2 > 1$. TRANSFORM TO $x_1 = z_1 \left(\frac{-2 \ln r^2}{n^2} \right)^{1/2}$, CHANGE OF VARIABLE FORMULA

FOR MVG: $\Sigma = LL^T$ CHOLSKY DECOMP. L IS LOWER TRIANGULAR. SAMPLE x AS ABOVE $\rightarrow y = Lx + \mu$

$$p(x_1, x_2) = \frac{1}{2\pi} \exp\left(-\frac{1}{2} x^T \Sigma^{-1} x\right)$$

REJECTION SAMPLING

• CREATE PROPOSAL DISTRIBUTION Q WHERE $Mq(x) \geq \tilde{p}(x)$, M CONSTANT, $\tilde{p}(x)$ IS UNNORMALIZED p . SAMPLE $x \sim Q(x)$, x LOCATION. SAMPLE $U \sim U(0,1)$ y LOCATION.

• **ACCEPTANCE RATE** $P(\text{ACCEPT}) = \frac{1}{M} \int \tilde{p}(x) dx \rightarrow$ PICK M AS SMALL AS POSSIBLE

• CAN USE WITH $M = MLE$; LIKELIHOOD TO DRAW SAMPLE FROM POSTERIOR. OR IF PRIOR IS INFORMATIVE

• **PICKING Q, P :** DOUND LOG DENSITY WITH PIECEWISE LINEAR FCN f AT FIXED GRID POINTS. \rightarrow ENVELOPE IS PIECEWISE EXPONENTIAL. SAMPLE REJECTED \rightarrow MAKE GRID TIGHTER.

• **IN HIGH DIM:** $M = \left(\frac{\sigma_u}{\sigma_f} \right)^D$ IS OPTIMUM. ACC. RATE IS $\frac{1}{M}$ CURSE OF DIM! \rightarrow MCMC YO

IMPORTANCE SAMPLING

• APPROXIMATES INTEGRALS OF THE FORM $I = E[f] = \int f(x)p(x) dx$. SAMPLE FROM HIGH PROBABILITY REGIONS $p(x)$, WHERE $|f(x)|$ IS LARGE.

• **IS SUPER-EFFICIENT:** NEEDS LESS SAMPLES THAN EXACT DISTRIBUTION BECAUSE FOCUS ON RELEVANT PARTS OF THE SPACE

• SAMPLES FROM PROPOSAL: THEN $E[f] = \int f(x) \frac{p(x)}{q(x)} q(x) dx = \frac{1}{S} \sum w_s f(x_s) = \hat{I} \rightarrow \frac{p(x_s)}{q(x_s)} = w_s$ **IMPORTANCE WEIGHT**

• **PICK Q TO MINIMIZE VARIANCE OF ESTIMATE \hat{I}** $\rightarrow Q^*(x) = \frac{|f(x)|p(x)}{\int |f(x')|p(x') dx'}$. TOTAL & BORN WHEN WE HAVE NO f IN MIND

• **ANCESTRAL SAMPLING:** FOR DAGS: - **NO EVIDENCE** \rightarrow ROOT, CHILDREN, CHILDREN OF CHILDREN, ... ON BECAUSE WE HAVE NO EVIDENCE.

- **EVIDENCE:** EV. NODES ARE CHANGED TO OBSERVED VALUES. DO ANCESTRAL BUT REJECT ALL SAMPLING IF INCONSISTENT VALUES. **LONG SAMPLING**

\rightarrow **MORE EFFICIENT:** ONLY USE OBSERVED VALUES FOR OBS. VARS, NO SAMPLING. $Q(x) = \prod p(x_t | x_{pa(t)}) \prod \delta_{x_t^*}(x_t)$

$$W(x) = \prod p(x_t | x_{pa(t)})$$

• **SAMPLE IMPORTANCE RESAMPLING:** DRAW ~~WEIGHTS~~ SAMPLES FROM $p(x)$ WITH I.S. $p(x) \approx \sum w_s \delta_{x_s}(x)$. THEN SAMPLE WITH REPLACEMENT FROM p PICKING WITH PROB $w_s \rightarrow p(x) \approx \frac{1}{S} \sum \delta_{x_s}(x)$ • CAN USE FOR BAYESIAN INFERENCE IN LOW-DIM

PARTICLE FILTERING

DOES RECURSIVE BAYESIAN INFERENCE FOR DYNAMIC BAYESIAN NETWORKS, NONLINEARITY, NON-STATIONARITY

APPROXIMATE STATE TRAJECTORY: $p(z_{1:T} | y_{1:T}) \approx \sum W_t^s \delta_{z_{1:T}}(z_{1:T})$. W_t ARE NORMALIZED WEIGHTS. BELIEF STATE UPDATED WITH IMPORTANCE SAMPLING. $W_t^s \propto \frac{p(z_{1:T}^s | y_{1:T})}{Q(z_{1:T}^s | y_{1:T})}$. USUAL MARKOV ASSUMPTIONS $\rightarrow W_t^s \propto W_{t-1}^s \frac{p(y_t | z_t^s) p(z_t^s | z_{t-1}^s)}{Q(z_t^s | z_{t-1}^s, y_t)}$. CAN NOW APPROX POSTERIOR DENSITY $p(z_t | y_{1:T}) \approx \sum W_t^s \delta_{z_t}(z_t)$

DEGENERACY \rightarrow AS IS PF FAILS AFTER A FEW STEPS BECAUSE MOST PARTICLES HAVE NEGLIGIBLE WEIGHT. BECAUSE WE SAMPLE IN INCREASINGLY HIGH DIM SPACE. ESTIMATE USING VARIANCE OF WEIGHTS OVER STEPS. HIGH VAR \rightarrow NO GOOD

FIXES: - **RESAMPLING:** DROP PARTICLES WITH LOW WEIGHT AND REPLACE WITH GOOD ONES ACCORDING TO THEIR WEIGHTS, $O(S)$

BUT NOW WE LOSE DIVERSITY \rightarrow BOOTSTRAP, MCMC, OR SAMPLE NEW PARTS FROM A NODE ESTIMATE FOR SMOOTHING **REGULARIZED PF**

- **PROPOSAL DISTRIBUTION:** SAMPLE FROM PRIOR $\rightarrow W_t^i \propto W_{t-1}^i p(y_t | z_t^i)$ **COMBINATION ALGO.** NOT GOOD IF LIKELIHOOD IS NARROWER THAN PRIOR (SENSOR MORE INFORMATIVE THAN MOTION MODEL) \rightarrow **SAMPLE FROM DATA**

• SAMPLE FROM PROPOSAL, USUALLY FROM PREV TIME STEP

• NORMALIZE SAMPLE WEIGHTS $W_t^i = W_{t-1}^i \cdot \text{ACTUAL WEIGHT}$

$$S = \frac{1}{\text{VAR}(W_t)} \quad W_t^i = \frac{W_{t-1}^i}{\sum W_{t-1}^i}$$

• RESAMPLE THOSE WITH LOW S

OPTIMAL PROPOSAL: $W_t^i \propto W_{t-1}^i \int p(y_t | z_t^i) p(z_t^i | z_{t-1}^i) dz_t^i$ CONDITIONAL ON OUR VALUES WEIGHT VAR IS 0
GENERALLY UNTRACTABLE BUT OK FOR DISCRETE STATE SPACES OR GAUSSIAN.
NOT GAUSSIAN \rightarrow USE GAUSSIAN APPROXIMATION VIA UNSCENTED TRANSFORM

PF APPLICATIONS

ROBOT LOCALIZATION: DISCRETE OCCUPANCY GRID. COULD USE HMM BUT $O(N^2)$ SPACE IS LARGE. PF AS SPARSE APPROXIMATION.

VISUAL TRACKING: OH BUT SENSITIVE TO COLOR \rightarrow MORE PARTICLES OR FRAMESHIP

TIME SENSES: NEURODYNAMIC STUFF

RAO-BLACKWELLIZED PF

TWO TYPES OF PARTICLE. Q_t, Z_t . WHEN NEW Q_t , WE INTEGRATE Z_t OUT. WE SAMPLE Q_t AM $P(Z_t|Q_t)$ IS A DISTRIBUTION **DISTRIBUTIONAL PARTICLES**

REDUCES SAMPLING SPACE \rightarrow REDUCES VARIANCE. \bullet ~~M/E~~ FOR EACH PARTICLE \bullet EXAMPLE: SLDS MODEL, USE WMM FOR PREDICTION \rightarrow MIXTURE OF WMMs

LOOK-AHEAD: NEW WEIGHTS ARE INDEPENDENT FROM NEW Q_t VALUES. COMPUTE THEM, ASSES S FOR PARTICLES TO REPLACE. INITIALIZE RESAMPLES WITH ALREADY COMPUTED WEIGHTS POSTERIOR $\cdot O(N)$

APPLICATIONS: \rightarrow MANEUVERING TARGET TRACKING \cdot UNKNOWN CHAIN CONTINUOUS STATE. IE MISSILES.

\rightarrow FASTSLAM. STAMING WMM IS $O(L^3)$. SAMPLE ROBOT TRAJECTORY Q AM RUN WMM 2D INSIDE EACH PARTICLE, $O(L)$ PER PARTICLE

MARKOV CHAIN MONTE CARLO

IDEA: CONSTRUCT MARKOV CHAIN WHOSE STATIONARY DISTRIBUTION IS TARGET DENSITY $P^*(x)$, PRIOR OR POSTERIOR.
BY PERFORMING RANDOM WALK WE DRAW CORRELATED SAMPLES AND WE CAN INTEGRATE WRT P^*

MCMC | VARIATIONAL METHODS

- EASIER TO IMPLEMENT
- WORKS ON MANY MODELS
- FASTER ON REALLY HUGE MODELS
- FASTER FOR SMALL/MEDIUM PROBLEMS
- DETERMINISTIC
- EASY TO KNOW WHEN TO STOP
- GIVES LOWER BOUND ON LL

GIBBS SAMPLING

- ANALOGUE OF COORDINATE DESCENT
- SAMPLE EACH VARIABLE IN TURN, CONDITIONED ON VALUES OF ALL OTHER VARS ON DISTRIBUTION. $x_1^{s+1} \sim P(x_1 | x_2^s, x_3^s); x_2^{s+1} \sim P(x_2 | x_1^{s+1}, x_3^s); x_3^{s+1} \sim P(x_3 | x_1^{s+1}, x_2^{s+1})$
- $P(x_i | x_{-i})$ IS FULL CONDITIONAL FOR i
- IF $P(x)$ IS GRAPHICAL MODEL, FOR i WE ONLY NEED ITS MARKOV BLANKET
- NEED TO BURN-IN THE CHAIN, RUN AND DISCARD SAMPLES, TO MAKE IT ENTER ITS STATIONARY DISTRIBUTION | DEDICATE BURN-IN

GIBBS ON ISING

- $P(x_T | x_{-T}, \theta) \propto \prod_{\text{NODE}} \psi_T(x_i, x_j)$ FULL CONDITIONAL: $P(x_T = +1 | x_{-T}, \theta) = \frac{\text{EXP}(\eta_T)}{\text{EXP}(\eta_T) + \text{EXP}(-\eta_T)} = \text{SIGM}(2\eta_T)$ COUPLING STRENGTH, $\eta = x_T(\hat{A}_T - D_T)$
- COMBINE W/ LOCAL EVIDENCE: $\psi_T(x_T) = \mathcal{N}(x_T | \mu_T, \sigma^2) \rightarrow P(x_T = +1 | x_{-T}, \psi, \theta) = \text{SIGM}(2\eta_T - \log \frac{\psi_T(+1)}{\psi_T(-1)})$ IS MORE OVERCONFIDENT THAN VANILNE

GIBBS FOR INFERRING GMM PARAMS

- SEMI-CONJUGATE PRIOR: SAME PRIORS FOR EACH MIXTURE COMPONENT
- INDICATORS: $P(z_i = u | x_i, \mu, \Sigma, \pi) \propto N(x_i | \mu_u, \Sigma_u) \times \pi_u$ MIXING WEIGHTS: $P(\pi | z) = \text{Dir}(\{\alpha_u + \sum_{i=1}^n \mathbb{1}(z_i = u)\}_{u=1}^K)$
- MEANS: $P(\mu_u | \Sigma_u, z, x) = N(\mu_u | m_u, V_u)$ COVARIANCES: $P(\Sigma_u | \mu_u, z, x) = \text{IW}(\Sigma_u | S_u, \nu_u)$
- ISSUE: LABEL SWITCHING! MODEL PARAMS θ AND INDICATORS z ARE UNIDENTIFIABLE. WHAT ONE SAMPLE CONSIDERS C1, ANOTHER CONSIDERS C2.
- $O(NKD)$
- CANNOT AVERAGE OVER SAMPLES, NOT AN ISSUE IN EM OR VBEM SINCE THEY LOCK ON A SINGLE MODE
- BEST NOT TO ASK QUESTIONS WITHOUT IDENTIFIABLE ANSWERS, IE CLUSTER MEMBERSHIP, BUT LIKE CLUSTER SAMENESS INSTEAD

COLLAPSED GIBBS

- INTEGRATE SOME UNKNOWNS OUT AND SAMPLE THE REST. WE SAMPLE z WHILE INTEGRATING OUT θ . \rightarrow WE CAN DRAW C.O. SAMPLES $\theta^j \sim P(\theta | z^j, D)$, LOWER VARIANCE THAN JOINT STATE SPACE.
- RAO-BLACKWELL THEOREM: z, θ DEPENDENT RV, $f(z, \theta)$ SCALAR FCN; $\rightarrow \text{VAR}_{z, \theta}[f(z, \theta)] \geq \text{VAR}_z[E_\theta[f(z, \theta) | z]]$ VARIANCE WILL NEVER BE HIGHER THAN FULL SPACE
- FOR A GMM: WE COLLAPSE ON μ_u, Σ_u, π , WE SAMPLE z . z, x_i BECOME INTERDEPENDENT BUT IS OK; AT EACH x REMOVE SUFF STATS, COMPUTE, UPDATE
- IF CONJUGATE PRIORS $\rightarrow P(x_i | D_{-i}, u)$ IN CLOSED FORM. • GENERALLY BETTER THAN VANILLA

GIBBS FOR HIERARCHICAL GLM

- SO TO BORROW STATISTICAL STRENGTH $y_{ij} = x_{ij}^T w_j + \epsilon_{ij}$ • w_j COME FROM COMMON PRIOR $w_j \sim N(\mu_w, \Sigma_w)$. USEFUL COMMON PRIORS
- PRIORS FOR SHARED PARAMS $\mu_w \sim N(\mu_0, V_0); \Sigma_w \sim \text{IW}(\eta_0, S_0^{-1}); \sigma^2 \sim \text{IG}(\nu_0/2, \nu_0 S_0^2/2)$
- FULL CONDITIONALS: $P(w_j | D_j, \theta) \sim N$; $P(\mu_w | w, \Sigma_w) \sim N$; $P(\Sigma_w | \mu_w, w) \sim \text{IW}$; $P(\sigma^2 | D, w) \sim \text{IG}$
- POSTERIOR PREDICTIVE: $E[y_i | x_{ij}] = x_{ij}^T \hat{w}_j$, $\hat{w}_j = E[w_j | D] = \frac{1}{S} \sum w_j^s$ • VERY NICE REGULARIZATION

BUGS/JAGS

GIBBS IS COOL BECAUSE MANY MODELS. OPTIMIZER ONLY NEEDS MODEL SPECIFICATION (A DGM) AND METHODS FOR SAMPLING FROM FULL CONDITIONALS. BUGS/JAGS ARE TOOLS FOR THIS

IMPUTATION POSTERIOR

GIBBS WHERE WE SPLIT z s AND θ s. MCMC VERSION OF EM. $E \rightarrow I$; $M \rightarrow P$ IS DATA AUGMENTATION TECH.

BLOCKING GIBBS

IF VARS ARE SIGNIFICANTLY CORRELATED IT'LL TAKE A WHILE FOR MOVING AWAY FROM SPACE. \rightarrow SAMPLE GROUPS OF VARS AT SAME TIME, FOR WARM MOVES

METROPOLIS / HASTINGS

GIBBS IS LIMITED BECAUSE REQUIRES MARKOV STRUCTURE.

IDEA: WE PROPOSE A MOVE FROM x TO x' WITH $q(x'|x)$ KERNEL/PROPOSAL DISTRIBUTION. COMMON CHOICE: $q(x'|x) = N(x'|x, \Sigma) \rightarrow$ RANDOM WALK METROPOLIS

- IF $q(x'|x) = q(x|x')$ \rightarrow IMPEDENCE SAMPLER

- AFTER MOVE IS PROPOSED, WE ACCEPT IT OR REJECT IT. SYMMETRIC $q(x'|x) = q(x|x')$, ACCEPT $R = \min\left(1, \frac{p^*(x')}{p^*(x)}\right)$

- SAMPLE $u \sim U(0,1)$ AND CHECK AGAINST R

- WE OCCASIONALLY ALLOW MOVES DOWNHILL TO LESS PROBABLE STATES, AND ALWAYS MOVE TO MORE PROBABLE

ASYMMETRIC: $R = \min(1, a)$; $a = \frac{p^*(x')/q(x'|x)}{p^*(x)/q(x|x')}$

HASTINGS CORRECTION

PROPOSAL MIGHT FAVOR STATES INSTEAD OF TARGET

GIBBS IS SPECIAL CASE OF MH WHERE x_i IS FROM FULL CONDITIONAL, x_{-i} UNCHANGED, AND ACCEPTANCE RATE 100%.

- DOES NOT REQUIRE KNOWLEDGE OF z BECAUSE THEY CANCEL OUT IN q

CHOOSING THE PROPOSAL DISTRIBUTE ON IF NONZERO PROB OF MOVING TO NONZERO PROB TARGET SPACES. IN PRACTICE WE HAVE TO TUNE STUFF BECAUSE VAR TOO LITTLE \rightarrow WE DON'T EXPLORE; VAR TOO LARGE \rightarrow WE DON'T GET ENOUGH SAMPLE, CHUNKY HISTOS. PLAY AROUND WITH KERNEL WIDTH. DO PILOT RUNS. 25% - 40% ACCEPTANCE RATE IS GOOD.

- GAUSSIAN PROPOSALS: IMPEDENCE PROPOSALS OR RW PROPOSALS $\rightarrow N(w|w, s^2 I^{-1})$. SET s^2 AT $2.38^2/0$ FOR 23.4% ACC. RATE. CONTINUOUS STATE SPACES, USE HESSIAN TO DEFINE COVARIANCE.

- MIXTURE PROPOSALS: $q(x'|x) = \sum_n w_n q_n(x'|x)$

FEATURES

- DATA-DRIVEN PROPOSALS $q(x'|x, D)$ DEPEND ON DATA AS WELL. SAMPLE (x_i, D) FROM FWD MODEL. TRAIN DISCRIMINATIVE CLASSIF FOR $p(x|f(D))$

CAN DO ON JUST PARTS OF DATA STATE SPACE. $q(x'|x, D) = \pi_0 q_0(x'|x) + \sum_n \pi_n q_n(x'|f_n(D))$

\rightarrow GENERATE AND TEST DE/ PROPOSAL ARE TESTED ON POSTERIOR RATIO.

ADAPTIVE MCMC: ALLOW CHANGES OF PROPOSAL PARAMS AS ALGO RUNS. IE LARGE \rightarrow SMALL COVARIANCE. EXPLORATION \rightarrow EXPLOITATION. WATCH OUT NOT TO VIOLATE MARKOV PROPERTY. NO DEPENDENCY ON WHOLE CHAIN.

INITIALIZATION: IT HAS TO START ~~TR~~ WITH NONZERO PROBABILITY. TRYING TO FIND IN CLOSED FORM. USUALLY FIND A LOCAL MODE W/OPTIMIZER AND INITIALIZE THERE. DISCRETE STATE SPACES: MULTIPLE RESTARTS. CONTINUOUS SPACES: DO LOCAL EXPLORATION TO VISIT ENOUGH PROB MASS OF POSTERIOR

AUXILIARY VARIABLE MCMC

IMPROVES SAMPLING EFFICIENCY BY INTRODUCING DUMMY VARS TO REDUCE CORR. BETWEEN ORIGINALS. $\sum_z P(x, z) = P(x)$, $P(x|z)$ EASIER TO SAMPLE

EXAMPLE: LOGISTIC REGRESSION

$$\begin{cases} z_i = w^T x_i + \epsilon_i \\ \epsilon_i \sim N(0, 1) \\ P(y=1) = \text{SIGM}(w^T x_i) \end{cases}$$

APPROXIMATE LOGISTIC WITH SCALE MIXTURE OF STUDENT'S T

$$G_i \sim T(0, 1, \nu) \rightarrow G_i \sim N(0, \lambda_i^{-1}), \lambda_i \sim \text{Ga}(\nu/2, \nu/2)$$

- SUCH FORMULATION CAN RANGE BETWEEN PROBIT AND LOGIT, ESTIMATING ν
- CAN BE USED FOR SETTING STRENGTH OF REGULARIZER, 'BAYESIAN CV'

Slice Sampling

UNIDIMENSIONAL, MULTIMODAL DISTRIBUTION. ADD AUX VAR TO MAKE UNIFORM MOVES. SAMPLE x^t , SAMPLE U^{t+1} UNIF. ON $[0, f(x^t)]$
SAMPLE x^{t+1} ON SLICE WHERE $f(x) \geq U^{t+1}$. FOR MULTIVARIATE: AUX VAR FOR EACH DIM. DOES NOT NEED FULL CONDITIONALS (GIBBS) OR USER-SPECIFIED PROPOSAL (MH)

SWEDENSEN-WANG

FOR ISING MODELS. SPEEDS UP MIXING. WHEN $J > 0$, FRUSTRATES SYSTEM FOR $J < 0$: EXPONENTIALLY MANY MODES.

ONE PER EDGE AUX. VARIABLES \rightarrow **BOND VARIABLES**. SAMPLE FROM EXTENDED MODELS AND THROW z AWAY, KEEPING x . BONDS ARE CI W/ NO NODES.

EASY FACTORIZATION. IF NODES BETWEEN EDGES ARE IN SAME STATE SET BOND ON WITH $p = 1 - e^{-2J}$. ASSES CONNECTED COMPONENTS.

FORCE NODES IN CC TO ASSUME SAME STATE, RANDOM.

HYBRID/HAMILTONIAN MCMC

CONTINUOUS STATE SPACES. WE CAN GRADIENT OF UNNORMALIZED LOG-POSTERIOR. IE NEURAL NETS.

AUX VARIABLES FOR 'MOMENTUM' OF PARTICLE | PARAMS. UPDATE POSITION/MOMENTUM. SET HOW MANY AND HOW BIG USAPFROG STEPS.

ANNEALING

TEMPERATURE. IN THE BEGINNING DISTRIBUTION IS SMOOTHER.

SIMULATED ANNEALING

FINDS GLOBAL OPTIMUM OF BINCH-BOX FUNCTION f .

BOLTZMANN DISTRIBUTION: $P(x) \propto \frac{\exp(-f(x))}{\int \exp(-f(x)) dx} \exp(-f(x)/T)$ f IS ENERGY. T IS TEMP. $T \rightarrow 0$

- AT HIGH TS FCN IS APPROX. FLAT \rightarrow EASY TO MOVE AROUND. THEN WAVE PEAKS WIDER AND SMALL (LOCAL) ONES DISAPPEAR.
- CONTINUATION METHOD \rightarrow WE CAN TRACK THE PEAK AND FIND GLOBAL OPTIMUM
- SAMPLE ACCORDING TO PROPOSAL IE RW $x' = x_t + \epsilon_u \rightarrow \alpha = \exp((f(x) - f(x'))/T) \rightarrow$ **ACCEPT** WITH PROB $\min(1, \alpha)$
- WE MIGHT STILL ACCEPT A WORSE STATE IF TEMP HIGH ENOUGH
- **COOLING SCHEDULE** IS CRITICAL. USUALLY EXPONENTIAL $T_n = T_0 C^n$, $T_0 = 1 \sim (20, 2)$

ANNEALED IMPORTANCE SAMPLING

SIMULATED ANNEALING + IMPORTANCE SAMPLING. DRAW SAMPLES FROM 'DIFFICULT' DISTRIBUTIONS.

START FROM EASY DISTRIBUTION AND GO OVER STEPS TO THE DIFFICULT $P_j(x) = f_0(x)^{p_j} f_m(x)^{1-p_j}$, $1 = p_0 > p_1 > \dots > p_n = 0$

HAVE MARKOV CHAINS $T_j(x, x')$ LEAVING P_j INVARIANT.

SAMPLE $z_{n-1} \sim P_n \dots z_0 \sim T_1(z_1)$; $x = z_0$ AM GIVE WEIGHT $w = \frac{f_{n-1}(z_{n-1}) \dots f_0(z_0)}{f_n(z_{n-1}) \dots f_1(z_0)}$

PARALLEL TEMPERING

RUN MULTIPLE CHAINS IN PARALLEL AT DIFFERENT TEMPERATURES. ALLOW CHAINS TO SAMPLE FROM NEIGHBORING ONES.

WHY MH WORKS BECAUSE IT DEFINES A TRANSITION FUNCTION WHICH SATISFIES DETAILED BALANCE. $\rightarrow P^*$ IS STATIONARY DISTRIBUTION
CHAIN IS ERGODIC AND IRREDUCIBLE $\rightarrow P^*$ IS UNIQUE

REVERSIBLE JUMP (TRANSDIMENSIONAL MCMC)

WE SAMPLE IN SPACES OF DIFFERENT DIMENSIONALITY. EG MIXTURE MODELS W/ UNKNOWN NUMBER OF MIXTURES.

TROUBLE WHEN COMPUTING MH ACCEPTANCE RATIO \rightarrow AUGMENT LOW-DIM SPACE WITH EXTRA RV SO MEASURE IS COMMON.

BURN-IN

MIXING TIME: TIME TO CONVERGE TO STATIONARY DISTRIBUTION. $\chi_e = \max \chi_e(x_0)$. DETERMINED BY $\gamma = \lambda_1 - \lambda_2$ EIGENVALS OF TRANSITION MATRIX

$\rightarrow \chi_e \leq O\left(\frac{1}{\gamma} \log \frac{n}{\epsilon}\right)$. n IS STATES. HARD FOR HIGH-DIM/CONTINUOUS SPACES.

\rightarrow **CONDUCTANCE** OF CHAIN; MIN PROB OVER ALL SUBSETS OF STATES TO TRANSITION TO ITS COMPLEMENT. $\phi \chi_e \leq O\left(\frac{1}{\phi^2} \log \frac{n}{\epsilon}\right)$
TOO LOW CONDUCTANCE \rightarrow MCMC NOT GOOD

DIAGNOSTICS: TRACE PLOTS. START FROM OVERDISPERSED STARTING POINT AND ASSESS CONVERGENCE TO SAME DISTRIBUTION OF SOME VARS OF INTEREST

ESPR: ESTIMATED POINT SCALE REDUCTION. VARIANCE WITHIN CHAIN VS VARIANCE ACROSS CHAINS. THINK ANOVA. B, W .

$\hat{V} = \frac{S-1}{S} W + \frac{1}{2} S$ UNBIASED OVER STATIONARITY, OVERESTIMATES IF OVERDISPERSION. $R = \sqrt{\frac{\hat{V}}{W}}$ MEASURES HOW MUCH POSTERIOR VARIANCE WOULD DECREASE IF SAMPLING CONTINUES.

ANOVA ANALYSIS BETWEEN/WITHIN CHAINS

$R \approx 1 \checkmark$

ACCURACY

MCMC PRODUCES AUTOCORRELATED SAMPLES **USE AUTOCORRELATION FCN** $VAR_{MCMC} = VAR_{MC}(\bar{f}) + \frac{1}{S^2} \sum E[(f_s - f^*)(f_t - f^*)]$

\rightarrow **USE THINNING** + SAMPLE SHUFFLING. WAS THIS GIBBS IS BETTER

$\rightarrow S_{EFF} = \frac{VAR_{MC}(f)}{VAR_{MCMC}(f)}$

IMPORTANCE OF CORR DISPERSES ON CORR

CHAINS

RUN 'MEDIUM' NUMBER OF CHAINS OF 'MEDIUM' LENGTH. DISCARD FIRST HALF. \rightarrow **TOO N**

MCMC APPROXIMATION (MCMC CHAPTER)

$$P(D|M) = \int P(D|\theta, M) P(\theta|M) d\theta \quad \text{IS OFTEN INTRACTABLE IF UNKNOWN VARIATE PRIOR / HIDDEN VARS,}$$

CANDIDATE METHOD $P(D|M) = \frac{P(D|\theta, M) P(\theta|M)}{P(\theta|D, M)}$ $\forall \theta$. NUMERATOR IS USUALLY GIVEN. APPROX DENOMINATOR WITH MCMC. STRONG ASSUMPTION THAT DREAM. HAS MAXIMIZERS OVER ALL MODES
→ POOR RESULTS

HARMONIC MEAN SAMPLING

$$1/P(D) \approx \frac{1}{S} \sum \frac{1}{P(D|\theta^s)} ; \theta^s = P(\theta|D) \quad \text{HARMONIC MEAN OF DATA LIKELIHOOD UNDER EACH SAMPLE.}$$

THEORETICALLY CORRECT BUT
WORST MONTE-CARLO METHOD EVER

- ONLY DEPENDS ON POSTERIOR SAMPLE
- OFTEN INSENSITIVE TO PRIOR.

ANNEALED IMPORTANCE SAMPLING

TO EVALUATE RATIO OF PARTITION FUNCTIONS $Z_0 = \int f_0(x) dx = \int f_0(z) dz$, $Z_N = \int f_N(x) dx = \int g(z) dz$.

$$\frac{Z_0}{Z_N} = E_g \left[\frac{f_0(z)}{g(z)} \right] \approx \frac{1}{S} \sum_{i=1}^S w_i \quad \bullet \text{ IF } F_N \text{ PRIOR AND } F_0 \text{ POSTERIOR } Z_N = P(D) \text{ GIVEN WE KNOW POST. NORM. CONSTANT } Z_0$$

HAMILTONIAN MCMC

- IDEA: USE HAMILTONIAN DYNAMICS TO PRODUCE MCMC PROPOSALS → FASTER STATE SPACE EXPLORATION THAN WITH RANDOM WALKS
- MOMENTUM VARIABLES - VOLUME PRESERVING NO NEED FOR JACOBIANS

HAMILTONIAN DYNAMICS

- Q POSITION • P MOMENTUM • HAMILTONIAN $H(Q, P)$, $\frac{dQ_i}{dt} = \frac{\partial H}{\partial P_i}$, $\frac{dP_i}{dt} = -\frac{\partial H}{\partial Q_i}$ • $z = (Q, P)$ VECTOR → $\frac{dz}{dt} = \nabla H(z)$
- ENERGY: $H(Q, P) = U(Q) + K(P)$ • $U(Q) = -\log p(Q) + K$, Q DISTRIBUTION POTENTIAL • $K(P) = P^T M^{-1} P / 2$, M MASS; SYMMETRIC POS DEF, DIAGONAL, OFTEN MULTIPLE OF I , $-\log p$ OF 0-MEAN GAUSSIAN WITH COVARIANCE M IS $N(0, I)$ FOR Q
- IN 1D: $H(Q, P) = U(Q) + K(P)$, $U(Q) = Q^2/2$, $K(P) = P^2/2$; $\frac{dQ}{dt} = P$, $\frac{dP}{dt} = -Q$, $Q = R \cos(t+\omega)$, $P = -R \sin(t+\omega)$
- $\frac{dQ}{dt} = M^{-1}P$
- $\frac{dP}{dt} = -\frac{\partial U}{\partial Q_i}$

- PROPERTIES: • REVERSIBILITY: CAN GO BACK FROM T TO T^{-1} , IS INVERTIBLE. LEAVE DISTRIBUTION INVARIANT (MARKOV)
- CONSERVATION H IS PRESERVED, INVARIANT. IN MH ACCEPTANCE PROB = 1 IF INVARIANT. IN PRACTICE APPROXIMATELY INVARIANT
- VOLUME PRESERVATION IN (Q, P) SPACE → GOOD BECAUSE NO NEED TO JACOBIAN. DIVERGENCE = 0 → PRESERVES VOLUME. $\det(\cdot) = 1$ PRESERVES VOLUME
- SYMPLECTICNESS $z = (Q, P)$, $B^T)^{-1} B = J^{-1}$ IMPLES VOL CONSERVATION. SYMPLECTIC CONDITION

HOW TO DISCRETIZE?

- EULER'S METHOD → POOR CONVERGENCE
 - MODIFIED EULER
 - LEAPFROG METHOD → EVEN BETTER
- $$\begin{cases} P_i(t+\epsilon) = P_i(t) - \epsilon \frac{\partial U}{\partial Q_i}(Q(t)) \\ Q_i(t+\epsilon) = Q_i(t) + \epsilon \frac{P_i(t+\epsilon)}{m_i} \end{cases}$$
- CAN ALSO DO Q BEFORE P
 - BETTER BECAUSE EXACT VOLUME PRESERVATION
 - USE UPDATED P VALUE → A LA GIBBS
- $$\begin{cases} P(t+\epsilon/2) = P(t) - (\epsilon/2) \frac{\partial U}{\partial Q_i}(Q(t)) \\ Q(t+\epsilon) = Q(t) + \epsilon \frac{P(t+\epsilon/2)}{m_i} \\ P(t+\epsilon) = P(t+\epsilon/2) - (\epsilon/2) \frac{\partial U}{\partial Q_i}(Q(t+\epsilon)) \end{cases}$$
- DO HALFWAY STEP INCREMENT FOR P
 - VOL PRESERVING, STABLE TRAJECTORIES

LET'S DO SOME MCMC NOW

- DENSITY → POTENTIAL ENERGY TO SHAPE
- MOMENTUM FOR ORIGINAL VARIABLES
- $U(Q) = -\log[\pi(Q)L(Q|D)]$ POSTERIOR FROM MODEL $\pi(Q)$ PRIOR π ; LIKELIHOOD L
- Q IS DISTRIBUTION OF INTEREST, WE CHOOSE DISTRIBUTION OF P VIA FORM OF H → USUALLY A QUADRATIC FOR P P IS 0 MEAN MVN, m_i IS VARIABLE
- CANONICAL DISTRIBUTION $P(x) = \frac{1}{Z} \exp[-E(x)/T] \rightarrow P(Q, P) = \frac{1}{Z} \exp[-H(Q, P)/T]$ E ENERGY FOR STATE x → $E(x) = -\log P(x) - \log(Z) \rightarrow P(Q, P) = \frac{1}{Z} \exp[-U(Q)/T] \exp[-K(P)/T]$ JOINTS ARE FACTORIZED!
- P, Q INDEPENDENT
- MOVES ON HYPERSURFACE OF CONSTANT PROBABILITY DENSITY
- E VARIABLES, P NEEDS TO MAKE IT WORK A MOMENTUM PER DIMENSION
- STEP 1: DRAW P_i FROM THEIR MVNS, INDEPENDENTLY OF Q . → JOINT NOT CHANGED, P_i INDEPENDENT
- STEP 2: DO A METROPOLIS UPDATE USING $H(Q, P)$ AS PROPOSAL. RUN HAMILTONIAN TRAJECTORY FOR L STEPS, STEPSIZE ϵ HYPERPARAMS LEAPFROG METHOD.
 - NEGATE MOMENTUM VARIABLES AT END: L . ACCEPT WITH PROB $\min[1, \exp(-H(Q^*, P^*) + H(Q, P))]$ → MAKES PROPOSAL SYMMETRICAL. NOT DONE BECAUSE WOULD NOT RUNNING THE CHAIN FOR > 1 STEPS. NEXT ITERATION WILL REDRAW P_i 'S
- $H(Q, P)$ ALMOST UNCHANGED
- REVERSIBILITY → SATISFIES DETAILED BALANCE CONDITIONS → CANONICAL DISTRIBUTION UNCHANGED/INVARIANT
- BEHAVIOR IS VERY DIFFERENT FROM RANDOM WALK BECAUSE GRADIENTS/MOMENTUMS, TURN BY MAGNITUDE OF STDDEV. UNLESS PROBABLY IN SOLUTIONS → RANDOM WALK
- IS GOOD BECAUSE TENDS TO MOVE IN CONSISTENT DIRECTIONS. BETTER ACCEPTANCE PROBABILITY. BETTER EXPLORATION. RW → INCREASE OF POSITION GROWS WITH N OF ITERATIONS
- IMPROVEMENT WRT RW IS RATIO OF MAX/MIN VARIANCES OF P_i
- CAN USE DIFFERENT K FUNCTIONS | TRANSFORMATION TO MAKE IT EFFICIENT | TRANSFORMATION INVARIANT
- SETTING L AND ϵ IS TRICKY, TRAIL & ERROR. CAN HAVE DIFFERENT ϵ FOR DIFFERENT P, Q PAIRS/DIMENSIONS
- FOR NEWELL NBS; CAN ALTERNATE $H(M)$ STEPS FOR P ONLY, AND STEPS FOR $H(Q)$ ONLY
- DIFFERENT 'RULE OF THUMB' VALUES THAN VANILLA MCMC
- SCALES BETTER IN INCREASING DIM SPACES

EXTENSIONS

- DISCRETIZATIONS OTHER THAN LEAPFROG METHOD $\rightarrow H$ SPLIT IN OTHER WAYS
- LANGEVIN MC: EQUIVALENT TO HMC WITH 1 LEAPFROG STEP, CAN NOT HAVE EXACT P_i' , IMPLY THEM BY MAKING Q_i 110 MD SAMPLE FROM GAUSSIAN.
- REFRESH MOMENTUM: $P' = \alpha P + (1 - \alpha)^{1/2} N$; N IS MVN VECTOR. REPLACE MOMENTUM AT EACH ITERATION: BUT SOFT, GRADUAL MOVE
- ACCEPT / REJECT WINDOWS: HAVE ACCEPTANCE DEPEND ON 'WINDOWS' OF STATES TOO CLOSE/OVR
- USE APPROXIMATIONS FOR $U(x)$
- DO SLOPE TO ϵ ADAPTIVELY
- TEMPERING: EITHER USE EXPLICIT T , CONTINUATION METHOD, OR MULTIPLY/DIVIDE MOMENTUM VARIABLES IN FIRST/SECOND HALF OF LEAPFROG TRAJECTORY. PRESERVES VOLUME
STEP $(t + \epsilon/2)$
KEEP IT AROUND 1