

# **STATS271/371: Applied Bayesian Statistics**

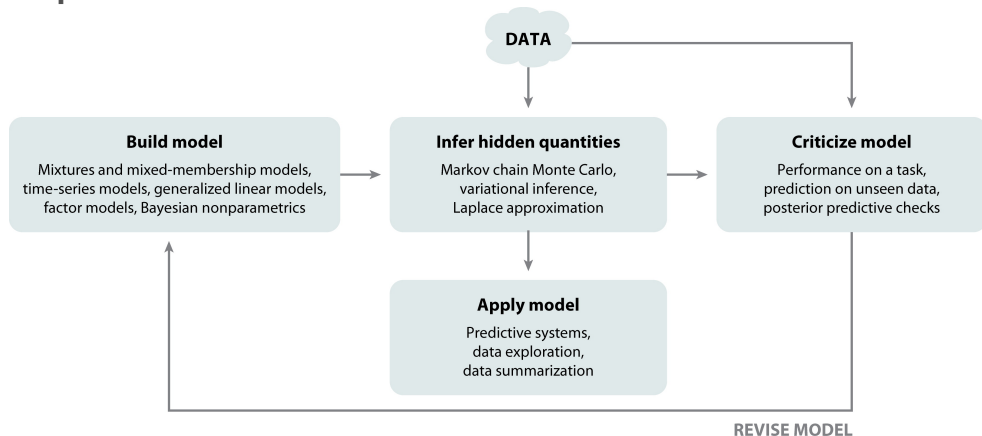
## **Robust regression models and Intro to MCMC**

Scott Linderman

(Some slides are adapted from unpublished notes by David Blei.)

April 19, 2021

# Box's Loop



Blei DM. 2014.

Annu. Rev. Stat. Appl. 1:203–32

Blei, *Ann. Rev. Stat. App.* 2014.

## Lap 3: Robust regress models and Intro to MCMC

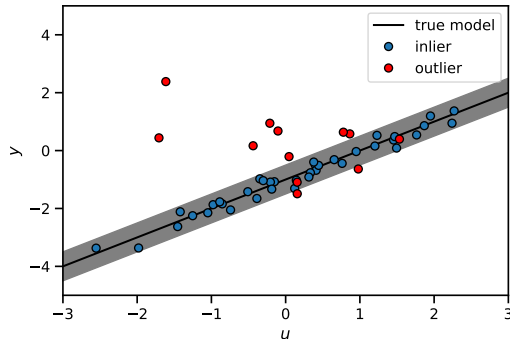
- ▶ **Model:** Robust regression models
- ▶ **Algorithm:** Markov chain Monte Carlo (specifically, HMC)
- ▶ **Criticism:** MCMC Diagnostics

# Robustness to outliers

Consider the following model:

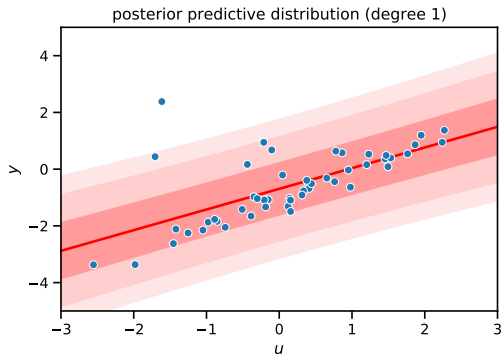
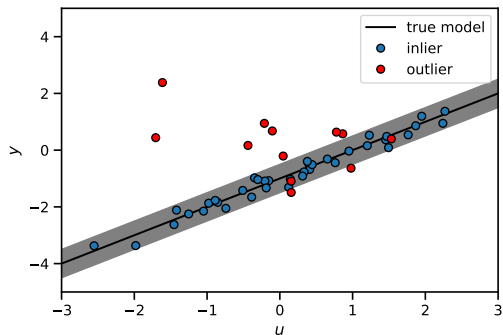
$$y_n \sim \begin{cases} \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma^2) & \text{with prob } 0.85 \\ \mathcal{N}(0, \sigma_{\text{out}}^2) & \text{o.w.} \end{cases}$$

where  $\sigma_{\text{out}}^2 \gg \sigma^2$ .



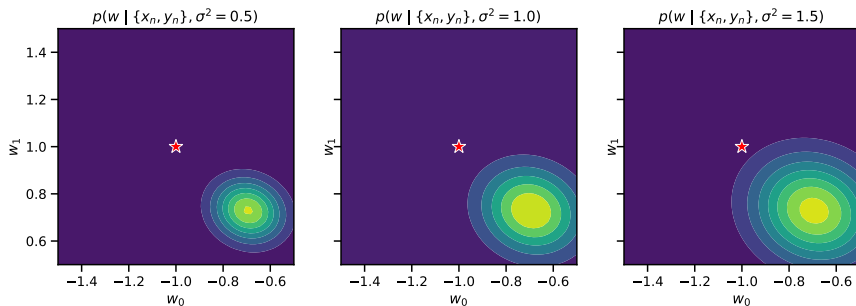
## Robustness to outliers II

A standard Bayesian linear regression (Lap 1) compensates by inferring different weights and a larger variance.



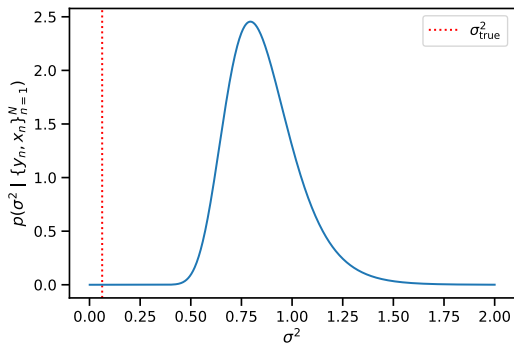
## Robustness to outliers III

We also see this in the posterior distribution of the weights...



## Robustness to outliers IV

...And in the posterior distribution of the variance



## Robust modeling with heavy tailed distributions

One way to allow for outliers is via a discrete mixture distribution, like we used to simulate this data. More on this next week!

Mixtures are one way of creating *heavy tailed* conditional distributions. I.e. they place more probability mass on observations far from the mean or mode.

Heavier tails allow the model to effectively downweight some observations when inferring the weights. In the discrete mixture example, if we knew which data points came from the outlier distribution, we would simply ignore them (set their weights to zero) when calculating the posterior distribution of the regression coefficients.

Of course, we don't *a priori* which datapoints are inliers and which are outliers.

Robust inference largely amounts to *figuring out how to weight different observations* during inference.



## Student's t as a scale-mixture of Gaussians

The Student's t distribution takes the discrete mixture idea to a continuous limit.

We can view the Student's distribution as a continuous *scale-mixture of Gaussians*,

$$t(y \mid \mu, \tau^2, \nu) = \int \mathcal{N}(y \mid \mu, \sigma^2) \text{Inv-}\chi^2(\sigma^2 \mid \nu, \tau^2) d\sigma^2 \quad (1)$$

$$= \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})\sqrt{\pi\nu\tau^2}} \left(1 + \frac{1}{\nu\tau^2} (y - \mu)^2\right)^{-\frac{\nu+1}{2}} \quad (2)$$

Equivalently, we can think of  $y \sim t(\mu, \tau^2, \nu)$  as arising from a two-step sampling procedure,

$$\sigma^2 \sim \text{Inv-}\chi^2(\nu, \tau^2) \quad (3)$$

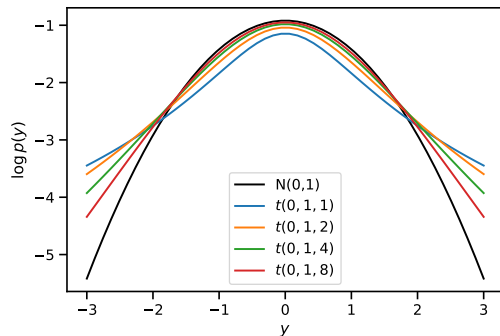
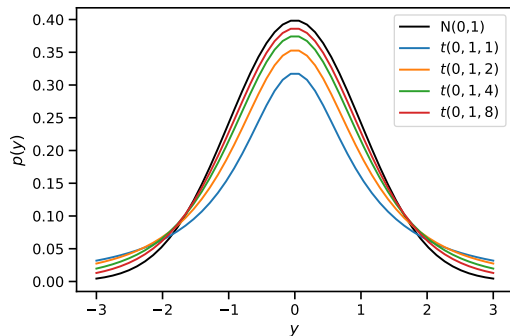
$$y \sim \mathcal{N}(\mu, \sigma^2). \quad (4)$$

First sample a random variance, then sample a Gaussian observation.

Look familiar? Recall the posterior predictive distribution in Bayesian linear regression.

# The Student's $t$ distribution has heavy tails

As the degree of freedom parameter  $\nu$  goes to infinity, the Student's  $t$  converges to a Gaussian.



# Robust regression with Student's t noise model

**Simple idea:** replace the normal distribution in Bayesian linear regression with a Student's t.

$$y_n \sim t(\mathbf{w}^\top \mathbf{x}_n, \tau^2, \nu) \iff \sigma_n^2 \sim \text{Inv-}\chi^2(\nu, \tau^2) \quad (5)$$

$$y_n \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \sigma_n^2). \quad (6)$$

We can think of this as giving each datapoint its own random variance  $\sigma_n^2$ .

How does this relate to “reweighting” various datapoints? Suppose we somehow knew the variances  $\sigma_n^2$  and just needed to infer the parameters  $\mathbf{w}$ . Recall that sufficient statistics for  $\mathbf{w}$  where  $\sum_{n=1}^N \frac{y_n \mathbf{x}_n}{\sigma_n^2}$ ,  $\sum_{n=1}^N \frac{\mathbf{x}_n \mathbf{x}_n^\top}{\sigma_n^2}$ , and  $\sum_{n=1}^N \frac{1}{\sigma_n^2}$ . Effectively, datapoints with larger variance contribute less to the sufficient statistics.

## Probabilistic model

Unfortunately, we don't know the variances (knowing the variances is like knowing which datapoints are outliers) so we have to work with the  $t$  distribution instead.

Our likelihood is,

$$p(\{y_n\}_{n=1}^N \mid \mathbf{w}, \tau^2, \nu, \{\mathbf{x}_n\}_{n=1}^N) = \prod_{n=1}^N t(y_n \mid \mathbf{w}^\top \mathbf{x}_n, \tau^2, \nu) \quad (7)$$

$$= \prod_{n=1}^N \frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2}) \sqrt{\pi \nu \tau^2}} \left( 1 + \frac{1}{\nu \tau^2} (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 \right)^{-\frac{\nu+1}{2}} \quad (8)$$

This doesn't simplify, but it is differentiable and easy to compute any point  $(\mathbf{w}, \tau^2, \nu)$ .

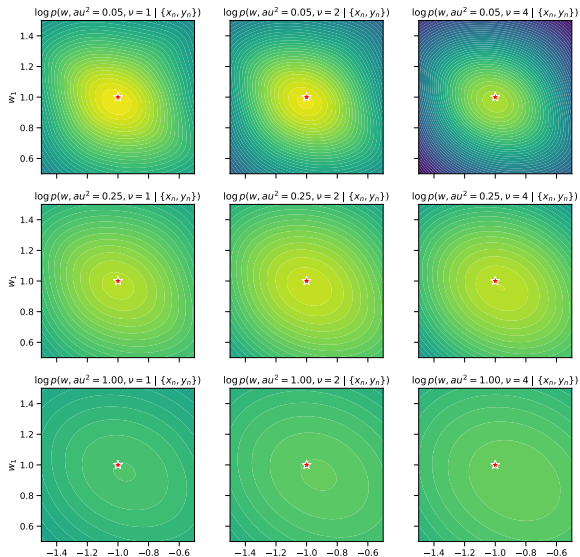
To complete the model, assume a broad Gaussian prior on the *log* of non-negative parameters,

$$p(\log \tau^2, \log \nu) = \mathcal{N}(\log \tau^2 \mid 0, 3) \mathcal{N}(\log \nu \mid 0, 3), \quad (9)$$

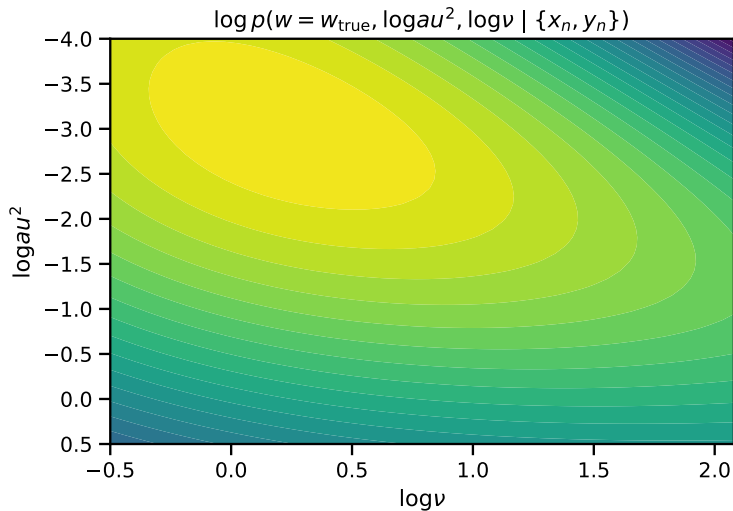
and an improper uniform prior on the weights.

A Gaussian prior on the log of a parameter is equivalent to a *log normal* prior on the parameter.

# Visualizing the log joint probability

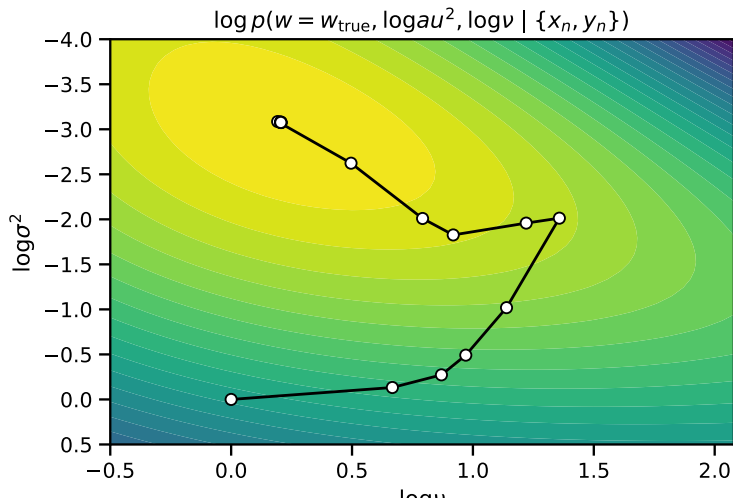


## Visualizing the log joint probability II



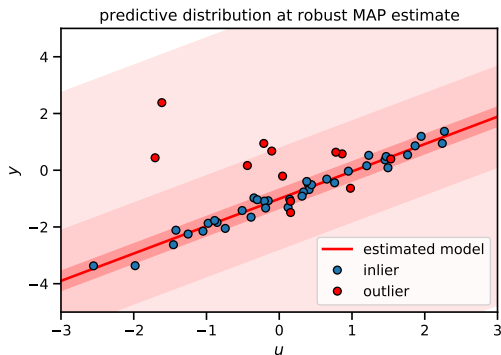
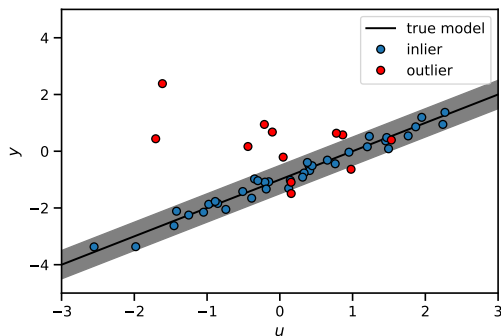
## MAP Estimation

Again, the posterior doesn't have a simple analytical form, but we can still try our tools from last week. E.g. we can use black box optimizers to find the mode.



## Predictive distribution with $w_{\text{MAP}}$

The MAP estimate of the parameters yields a much better estimate of the weights by allowing for a heavy tailed distribution of outliers.





# Notation

- ▶ Let  $\theta \in \Theta$  denote the model parameters.
  - ▶ In robust regression,  $\theta = (\mathbf{w}, \tau^2, \nu) \in \mathbb{R}^p \times \mathbb{R}_+ \times \mathbb{R}_+$ .
- ▶ Let  $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$  denote the observed data.

# Posterior expectations

The central object of Bayesian inference is the posterior distribution,  $p(\boldsymbol{\theta} \mid \mathcal{D})$ .

However, we almost always interact with the posterior distribution through *expectations*.

- ▶  $\mathbb{E}_{p(\boldsymbol{\theta} \mid \mathcal{D})}[\boldsymbol{\theta}]$ , the posterior mean.
- ▶  $\mathbb{E}_{p(\boldsymbol{\theta} \mid \mathcal{D})}[\mathbb{I}[\boldsymbol{\theta} \in \mathcal{A}]]$ , the probability of the parameters being in set  $\mathcal{A}$ .
- ▶  $\mathbb{E}_{p(\boldsymbol{\theta} \mid \mathcal{D})}[p(\mathcal{D}' \mid \boldsymbol{\theta})]$ , the posterior predictive density of new data  $\mathcal{D}'$ .

All of these can be written as  $\mathbb{E}_{p(\boldsymbol{\theta} \mid \mathcal{D})}[f(\boldsymbol{\theta})]$  for some function  $f$ .

(One exception is the posterior mode, which is not so easily expressed as an expectation.)

## Approximating posterior expectations

Generally, we can't analytically compute posterior expectations. (*Why not?*)

In these cases, we need to resort to approximations.

For example, we could use *quadrature methods* like Simpson's rule or the trapezoid rule to numerically approximate the integral over  $\Theta$ .

Roughly,

$$\mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})}[f(\boldsymbol{\theta})] \approx \sum_{m=1}^M p(\boldsymbol{\theta}_m | \mathcal{D}) f(\boldsymbol{\theta}_m) \Delta_m \quad (10)$$

where  $\boldsymbol{\theta}_m \in \Theta$  is a grid of points and  $\Delta_m$  is a volume around that point.

This works for low dimensional problems (say, up to 5 dimensions), but the number of points ( $M$ ) needed to get a good estimate grows exponentially with the parameter dimension.

# Monte Carlo approximations

**Idea:** approximate the expectation via sampling,

$$\mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})}[f(\boldsymbol{\theta})] \approx \frac{1}{S} \sum_{s=1}^S f(\boldsymbol{\theta}_s) \quad \text{where} \quad \boldsymbol{\theta}_s \sim p(\boldsymbol{\theta} | \mathcal{D}). \quad (11)$$

Let  $\hat{f} = \frac{1}{S} \sum_{s=1}^S f(\boldsymbol{\theta}_s)$  denote the Monte Carlo estimate. It is a random variable, since it's a function of random samples  $\boldsymbol{\theta}_s$ .

As such we can reason about its mean and variance. Clearly,

$$\mathbb{E}[\hat{f}] = \frac{1}{S} \sum_{s=1}^S \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})}[f(\boldsymbol{\theta})] = \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})}[f(\boldsymbol{\theta})]. \quad (12)$$

Thus,  $\hat{f}$  is an *unbiased* estimate of the desired expectation.

## Monte Carlo approximations II

What about its variance?

$$\text{Var}[\hat{f}] = \text{Var}\left(\frac{1}{S} \sum_{s=1}^S f(\boldsymbol{\theta}_s)\right) \quad (13)$$

$$= \frac{1}{S^2} \left( \sum_{s=1}^S \text{Var}[f(\boldsymbol{\theta})] + 2 \sum_{1 \leq s < s' \leq S} \text{Cov}[f(\boldsymbol{\theta}_s), f(\boldsymbol{\theta}_{s'})] \right) \quad (14)$$

If the samples are not only identically distributed but also *uncorrelated*, then  $\text{Var}[\hat{f}] = \frac{1}{S} \text{Var}[f(\boldsymbol{\theta})]$ .

In this case, the *root mean squared error* (RMSE) of the estimate is  $\sqrt{\text{Var}[\hat{f}]} = O(S^{-\frac{1}{2}})$ .

Compare this to Simpson's rule, which for smooth 1D problems has error rate  $O(S^{-4})$ . That's roughly 8 times better!

However, for multidimensional problems, Simpson's rule is  $O(S^{-\frac{4}{p}})$ , whereas the **error rate of Monte Carlo does not depend on the dimensionality!**

# The Catch

So far so good: we'll just draw a lot of samples to drive down our Monte Carlo error.

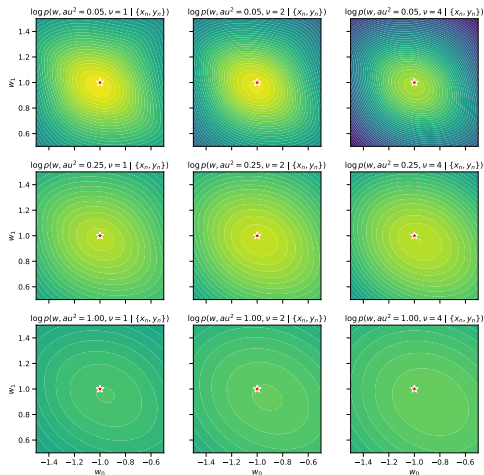
**Here's the catch!** How do you draw samples from the posterior  $p(\boldsymbol{\theta} \mid \mathcal{D})$ ?

We're interested in Monte Carlo for cases where the posterior does not admit a simple closed form!

In general, sampling the posterior is as hard as computing the marginal likelihood.

# Markov chain Monte Carlo (MCMC)

**Idea:** Design a Markov chain whose stationary distribution is the posterior.



# Markov chains

A *Markov chain* is a joint distribution of a sequence of variables,  $\pi(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_S)$ .

(To avoid confusion with the model  $p$ , we denote the densities associated with the Markov chain by  $\pi$ .)

The Markov chain factorizes so that each variable is drawn conditional on the previous variable,

$$\pi(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_S) = \pi_1(\boldsymbol{\theta}_1) \prod_{s=2}^S \pi(\boldsymbol{\theta}_s | \boldsymbol{\theta}_{s-1}). \quad (15)$$

This is called the *Markov property*.

The distribution  $\pi_1(\boldsymbol{\theta}_1)$  is called the *initial distribution*.

The distribution  $\pi(\boldsymbol{\theta}_s | \boldsymbol{\theta}_{s-1})$  is called the *transition distribution*. If the transition distribution is the same for each  $s$ , the Markov chain is *homogenous*.



# Stationary distributions

Let  $\pi_s(\boldsymbol{\theta}_s)$  denote the marginal distribution of sample  $\boldsymbol{\theta}_s$ . It can be obtained recursively as,

$$\pi_s(\boldsymbol{\theta}_s) = \int \pi_{s-1}(\boldsymbol{\theta}_{s-1}) \pi(\boldsymbol{\theta}_s | \boldsymbol{\theta}_{s-1}) d\boldsymbol{\theta}_{s-1}. \quad (16)$$

We are interested in the asymptotic behavior of the marginal distributions as  $s \rightarrow \infty$ .

A distribution  $\pi^*(\boldsymbol{\theta})$  is a **stationary distribution** if,

$$\pi^*(\boldsymbol{\theta}) = \int \pi^*(\boldsymbol{\theta}') \pi(\boldsymbol{\theta} | \boldsymbol{\theta}') d\boldsymbol{\theta}'. \quad (17)$$

That is, suppose the marginal of sample  $\boldsymbol{\theta}'$  is  $\pi^*(\boldsymbol{\theta})$ . Then the marginal of the next time point is also  $\pi^*(\boldsymbol{\theta})$ .

## Detailed balance

How can we relate transition distributions and stationary distributions?

A sufficient (but not necessary) condition for  $\pi^*(\boldsymbol{\theta})$  to be a stationary distribution is that it satisfies *detailed balance*,

$$\pi^*(\boldsymbol{\theta}')\pi(\boldsymbol{\theta} \mid \boldsymbol{\theta}') = \pi^*(\boldsymbol{\theta})\pi(\boldsymbol{\theta}' \mid \boldsymbol{\theta}) \quad (18)$$

In words, the probability of starting at  $\boldsymbol{\theta}'$  and moving to  $\boldsymbol{\theta}$  is the same as that of starting at  $\boldsymbol{\theta}$  and moving to  $\boldsymbol{\theta}'$ , if you draw the starting point from the stationary distribution.

To see that detailed balance is sufficient, integrate both sides to get,

$$\int \pi^*(\boldsymbol{\theta}')\pi(\boldsymbol{\theta} \mid \boldsymbol{\theta}')d\boldsymbol{\theta}' = \int \pi^*(\boldsymbol{\theta})\pi(\boldsymbol{\theta}' \mid \boldsymbol{\theta})d\boldsymbol{\theta}' = \pi^*(\boldsymbol{\theta}) \int \pi(\boldsymbol{\theta}' \mid \boldsymbol{\theta})d\boldsymbol{\theta}' = \pi^*(\boldsymbol{\theta}). \quad (19)$$

Thus,  $\pi^*(\boldsymbol{\theta})$  is a stationary distribution of the Markov chain with transitions  $\pi(\boldsymbol{\theta} \mid \boldsymbol{\theta}')$ .

# Ergodicity

Detailed balance can be used to show that  $\pi^*(\theta)$  is *a* stationary distribution, but not that it is *the* *unique* one.

This is where *ergodicity* comes in. A Markov chain is ergodic if  $\pi_s(\theta_s) \rightarrow \pi^*(\theta)$  regardless of  $\pi_1(\theta_1)$ .

An ergodic chain has only one stationary distribution,  $\pi^*(\theta)$ .

The easiest way to prove ergodicity is to show that it is possible to reach any  $\theta'$  from any other  $\theta$ . E.g. this is trivially so if  $\pi(\theta' | \theta) > 0$ .

Note: a more technical definition is that all pairs of sets *communicate*, in which case the chain is *irreducible*, and that each state is *aperiodic*. The definitions can be a bit overwhelming.

# The Metropolis-Hastings algorithm

Finally we come to our **main objective**: designing a Markov chain for which *the posterior is the unique stationary distribution*.

That is, we want  $\pi^*(\boldsymbol{\theta}) = p(\boldsymbol{\theta} \mid \mathcal{D})$ .

Recall our **constraint**: we can only compute the joint probability (the numerator in Bayes' rule), not the marginal likelihood (the denominator).

Fortunately, that still allows us to compute ratios of posterior densities! We have,

$$\frac{p(\boldsymbol{\theta} \mid \mathcal{D})}{p(\boldsymbol{\theta}' \mid \mathcal{D})} = \frac{p(\boldsymbol{\theta}, \mathcal{D})}{p(\mathcal{D})} \frac{p(\mathcal{D})}{p(\boldsymbol{\theta}', \mathcal{D})} = \frac{p(\boldsymbol{\theta}, \mathcal{D})}{p(\boldsymbol{\theta}', \mathcal{D})}. \quad (20)$$

## The Metropolis-Hastings algorithm II

Now rearrange the detailed balance condition to relate ratios of transition probabilities to ratios of joint probabilities,

$$\frac{\pi(\boldsymbol{\theta} \mid \boldsymbol{\theta}')}{\pi(\boldsymbol{\theta}' \mid \boldsymbol{\theta})} = \frac{\pi^*(\boldsymbol{\theta})}{\pi^*(\boldsymbol{\theta}')} = \frac{p(\boldsymbol{\theta} \mid \mathcal{D})}{p(\boldsymbol{\theta}' \mid \mathcal{D})} = \frac{p(\boldsymbol{\theta}, \mathcal{D})}{p(\boldsymbol{\theta}', \mathcal{D})} \quad (21)$$

To construct such a transition distribution  $\pi(\boldsymbol{\theta} \mid \boldsymbol{\theta}')$ , break it down into two steps.

1. Sample a proposal  $\boldsymbol{\theta}$  from a *proposal distribution*  $q(\boldsymbol{\theta} \mid \boldsymbol{\theta}')$ ,
2. Accept the proposal with *acceptance probability*  $a(\boldsymbol{\theta}' \rightarrow \boldsymbol{\theta})$ . (Otherwise, set  $\boldsymbol{\theta} = \boldsymbol{\theta}'$ .)

Thus,

$$\pi(\boldsymbol{\theta} \mid \boldsymbol{\theta}') = \begin{cases} q(\boldsymbol{\theta} \mid \boldsymbol{\theta}') a(\boldsymbol{\theta}' \rightarrow \boldsymbol{\theta}) & \text{if } \boldsymbol{\theta}' \neq \boldsymbol{\theta} \\ \int q(\boldsymbol{\theta}'' \mid \boldsymbol{\theta}') (1 - a(\boldsymbol{\theta}' \rightarrow \boldsymbol{\theta}'')) d\boldsymbol{\theta}'' & \text{if } \boldsymbol{\theta}' = \boldsymbol{\theta} \end{cases} \quad (22)$$

## The Metropolis-Hastings algorithm III

The constraint in (21) is trivially satisfied when  $\theta = \theta'$ . When  $\theta \neq \theta'$ , we need

$$\frac{\pi(\theta | \theta')}{\pi(\theta' | \theta)} = \frac{q(\theta | \theta') a(\theta' \rightarrow \theta)}{q(\theta' | \theta) a(\theta \rightarrow \theta')} = \frac{p(\theta, \mathcal{D})}{p(\theta', \mathcal{D})} \Rightarrow \frac{a(\theta' \rightarrow \theta)}{a(\theta \rightarrow \theta')} = \underbrace{\frac{p(\theta, \mathcal{D}) q(\theta' | \theta)}{p(\theta', \mathcal{D}) q(\theta | \theta')}}_{\triangleq A(\theta' \rightarrow \theta)} \quad (23)$$

WLOG, assume  $A(\theta' \rightarrow \theta) \leq 1$ . (If it's not, its inverse  $A(\theta \rightarrow \theta')$  must be.)

A simple way to ensure detailed balance is to set  $a(\theta' \rightarrow \theta) = A(\theta' \rightarrow \theta)$  and  $a(\theta \rightarrow \theta') = 1$ .

We can succinctly capture both cases with,

$$a(\theta' \rightarrow \theta) = \min \{1, A(\theta' \rightarrow \theta)\} = \min \left\{ 1, \frac{p(\theta, \mathcal{D}) q(\theta' | \theta)}{p(\theta', \mathcal{D}) q(\theta | \theta')} \right\}. \quad (24)$$

# The Metropolis algorithm

Now consider the special case in which the proposal distribution is symmetric; i.e.  $q(\boldsymbol{\theta} \mid \boldsymbol{\theta}') = q(\boldsymbol{\theta}' \mid \boldsymbol{\theta})$ .

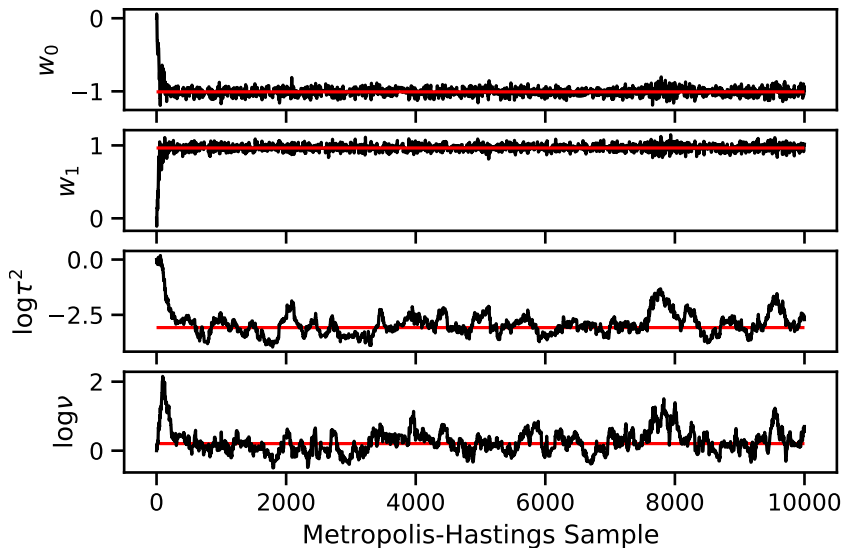
Then the proposal densities cancel in the acceptance probability and,

$$a(\boldsymbol{\theta}' \rightarrow \boldsymbol{\theta}) = \min \left\{ 1, \frac{p(\boldsymbol{\theta}, \mathcal{D})}{p(\boldsymbol{\theta}', \mathcal{D})} \right\}. \quad (25)$$

In other words, you accept any proposal that moves “uphill,” and only accept “downhill” moves with some probability.

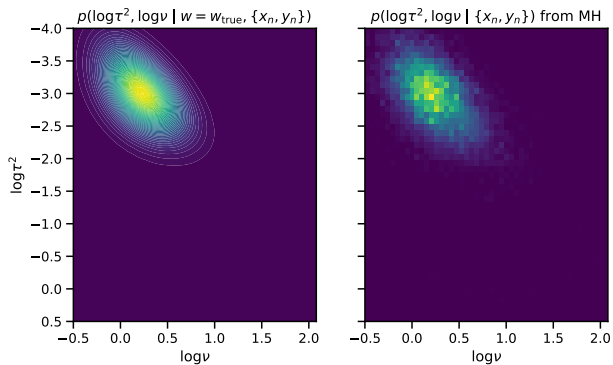
This is called the *Metropolis algorithm* and it has close connections to *simulated annealing*.

# Metropolis-Hastings Trace Plot





# Metropolis-Hastings Density Approximation



## Autocovariance and autocorrelation

Recall that,

$$\text{Var}[\hat{f}] = \frac{1}{S^2} \left( \sum_{s=1}^S \text{Var}[f(\boldsymbol{\theta})] + 2 \sum_{1 \leq s < s' \leq S} \text{Cov}[f(\boldsymbol{\theta}_s), f(\boldsymbol{\theta}_{s'})] \right) \quad (26)$$

$$\approx \frac{1}{S} \left( \text{Var}[f(\boldsymbol{\theta})] + 2 \sum_{\ell=1}^S \text{Cov}[f(\boldsymbol{\theta}_s), f(\boldsymbol{\theta}_{s+\ell})] \right) \quad (27)$$

since the covariance is only a function of the lag  $\ell$  once the chain has reached stationarity.

**Note:** At stationarity, the samples are identically distributed but still correlated!

$\text{Cov}[f(\boldsymbol{\theta}_s), f(\boldsymbol{\theta}_{s+\ell})]$  is called the *autocovariance*. It's a function of the lag  $\ell$  (and the function  $f$ ).

The *autocorrelation function* (ACF) is defined as,  $\text{acf}_f(\ell) = \text{Cov}[f(\boldsymbol{\theta}_s), f(\boldsymbol{\theta}_{s+\ell})] / \text{Var}[f(\boldsymbol{\theta})]$  so

$$\text{Var}[\hat{f}] \approx S^{-1} \text{Var}[f(\boldsymbol{\theta})] \left( 1 + 2 \sum_{\ell=1}^S \text{acf}_f[\ell] \right). \quad (28)$$

## Effective sample size

The *effective sample size* (ESS) approximates the effective number of independent samples you get from an autocorrelated chain, in terms of the variance of the Monte Carlo estimate,

$$S_{\text{eff},f} = S \frac{\text{Var}[f(\theta)]}{\text{Var}[f(\theta)](1 + 2 \sum_{\ell=1}^{\infty} \text{acf}_f[\ell])} = \frac{S}{1 + 2 \sum_{\ell=1}^{\infty} \text{acf}_f[\ell]} \quad (29)$$

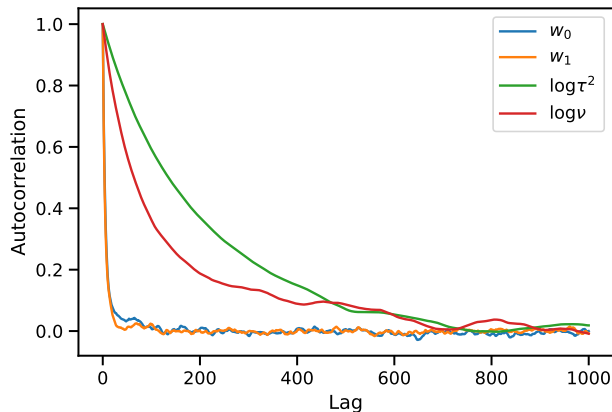
and we let  $S_{\text{eff}}$  denote the ESS of the identity functional  $f(\theta) = \theta$ .

You have to be a bit careful when estimating the ESS—for large values of  $\ell$  the sample correlation is too noisy. Typically, we stop when the sample acf is negative. See Section 11.5 of BDA3 for more details.

In practice, there are already good implementations in Python (c.f. `tfp.stats.effective_sample_size`) and R (c.f. the coda package).

## Autocorrelation of the MH trace above

We used `tfp.stats.auto_correlation` to compute the autocorrelation function for the traces above,



Effective sample sizes (per 100,000 samples): 5629, 7235, 257, and 364.

## Autocorrelation of a Random Walk

[From Geyer, 2011] To get some intuition, consider the following Markov chain,

$$\theta_s = \rho \theta_{s-1} + \epsilon_s \quad (30)$$

where  $\epsilon_s \sim \mathcal{N}(0, \tau^2)$ .

MH isn't quite a mean-reverting random walk, but that's not a terrible model.

In this toy example, we can calculate the autocovariance of the identity functional  $f(\theta) = \theta$ ,

$$\text{Cov}[\theta_s, \theta_{s+\ell}] = \rho \text{Cov}[\theta_s, \theta_{s+\ell-1}] = \rho^{\ell-1} \text{Cov}[\theta_s, \theta_{s+1}] = \rho^\ell \text{Var}[\theta_s] \quad (31)$$

so the autocorrelation function decays geometrically as  $\text{acf}(\ell) = \rho^\ell$ .

At stationarity,

$$\text{Var}[\theta_s] = \text{Var}[\theta_{s+1}] = \rho^2 \text{Var}[\theta_s] + \tau^2 \Rightarrow \text{Var}[\theta_s] = \frac{\tau^2}{1 - \rho^2} \quad (32)$$

For  $\rho^2 < 1$ , the stationary distribution exists and is  $\mathcal{N}(0, \frac{\tau^2}{1 - \rho^2})$ .

## Autocorrelation of a Random Walk II

Letting  $f(\theta) = \theta$ , we have,

$$\text{Var}[\hat{f}] = \frac{1}{S} \left( \text{Var}[\theta] + 2 \sum_{\ell=1}^S \text{Cov}[\theta_s, \theta_{s+\ell}] \right) \quad (33)$$

$$= \frac{1}{S} \cdot \text{Var}[\theta] \left( 1 + 2 \sum_{\ell=1}^S \rho^\ell \right) \quad (34)$$

$$\approx \frac{1}{S} \cdot \text{Var}[\theta] \left( 1 + 2 \frac{\rho}{1-\rho} \right) \quad (\text{for large } S) \quad (35)$$

$$= \frac{1}{S} \cdot \text{Var}[\theta] \cdot \frac{1+\rho}{1-\rho} \quad (36)$$

and

$$S_{\text{eff}} = S \cdot \frac{1-\rho}{1+\rho}. \quad (37)$$

As  $\rho \rightarrow 0$ , we recover ordinary Monte Carlo. As  $\rho \rightarrow 1$ , the autocorrelation causes the variance of the estimator to blow up and the effective sample size to go to zero.

# What about bias?

The mean squared error (MSE) of the estimator is determined by both the bias and the variance.

Ordinary Monte Carlo estimates are unbiased by construction, but MCMC estimates are only *asymptotically unbiased*.

Bias is introduced whenever the initial distribution  $\pi_1(\boldsymbol{\theta})$  differs from the stationary distribution; i.e. in all practical cases!

Fortunately, the bias decays as  $O(S^{-1})$  whereas the variance decays as  $O(S^{-1/2})$ , so asymptotically the MSE is dominated by the variance.

# How can we make smarter proposals?

Metropolis-Hastings with a symmetric Gaussian proposal behaves like a random walk.

Neal (2012) argues that in  $D$  dimensions, random walk MH needs  $O(D^2)$  iterations to get an independent sample.

Can we develop more efficient transition distributions?

Yes, if we have more information about the log probability.

For example, suppose that the log probability  $\log p(\theta)$  is differentiable. We can use the gradient to make proposals that move farther and are more likely to be accepted.



# Metropolis Adjusted Langevin Algorithm (MALA)

The *Metropolis-Adjusted Langevin Algorithm* uses the gradient of the log probability to make asymmetric proposals,

$$q(\boldsymbol{\theta}' | \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} + \tau \nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta}, \mathcal{D}), 2\tau^2 I) \quad (38)$$

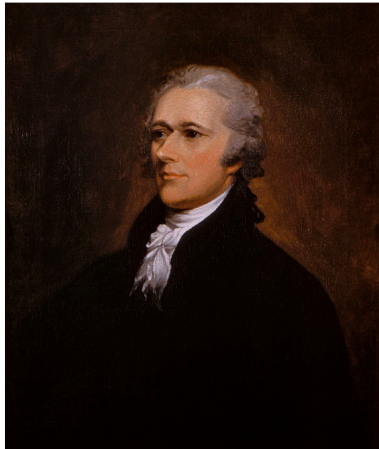
**Note:**  $q(\boldsymbol{\theta}' | \boldsymbol{\theta}) \neq q(\boldsymbol{\theta} | \boldsymbol{\theta}')$ ! To calculate the acceptance probability, you need the gradient at both points.

MALA can be motivated as a discrete-time approximation to the *Langevin* diffusion, a continuous-time stochastic differential equation for modeling molecular dynamics.

In high dimensions, the extra information provided by the gradient can lead to much more efficient chains. Neal argues that MALA needs  $O(D^{4/3})$  computation to produce an independent sample.

But why stop at one gradient step?

# Three Hamiltons



# Hamiltonian Monte Carlo

**Reference:** Neal (2012) *MCMC using Hamiltonian dynamics*.

**Idea:** *Think of negative log probability as an energy landscape. Now imagine a puck sliding around on this bumpy surface. Give it random kicks; it will tend to slide downhill toward points of low potential energy (high probability). Each kick can displace the puck by a large amount. Done properly, the puck will visit points with probability proportional to the posterior probability.*

# Notation

Following Neal (2012), let

- ▶  $\mathbf{q} \in \mathbb{R}^D$  denote the *position*; i.e. the current parameters (previously  $\theta$ )
- ▶  $\mathbf{p} \in \mathbb{R}^D$  denote the *momentum*; auxiliary variables that we don't care about, but which are necessary for HMC.
- ▶  $\mathbf{z} = [\mathbf{q}, \mathbf{p}]^\top \in \mathbb{R}^{2D}$  denote the combined *state of the system*.
- ▶  $\mathbf{M}$  denote the *mass matrix*, another artificial construct. Typically, this will be  $m\mathbf{I}$
- ▶  $U(\mathbf{q})$  denote the *potential energy*
- ▶  $K(\mathbf{p}) = \frac{1}{2}\mathbf{p}^\top \mathbf{M}^{-1}\mathbf{p}$  denote the *kinetic energy*

## Hamiltonian dynamics

The *Hamiltonian* is the sum of the potential  $H(\mathbf{q}, \mathbf{p}) = U(\mathbf{q}) + K(\mathbf{p})$ .

The partial derivatives determine how the state evolves over time,

$$\frac{dq_d}{dt} = \frac{\partial H}{\partial p_d} = [\mathbf{M}^{-1}\mathbf{p}]_d \quad (39)$$

$$\frac{dp_d}{dt} = -\frac{\partial H}{\partial q_d} = -\frac{\partial U}{\partial q_d} \quad (40)$$

for  $d = 1, \dots, D$ .

Compactly,

$$\frac{d\mathbf{z}}{dt} = J\nabla H(\mathbf{z}) \quad (41)$$

where

$$J = \begin{bmatrix} \mathbf{0}, I \\ -I, \mathbf{0} \end{bmatrix} \quad (42)$$

## One dimensional example

Consider the case where  $D = 1$  and  $U(q) = \frac{1}{2}q^2$  and  $K(p) = \frac{1}{2}p^2$ .

The partial derivatives are

$$\frac{\partial H}{\partial p} = p \quad (43)$$

$$-\frac{\partial H}{\partial q} = -q \quad (44)$$

so

$$\frac{d\mathbf{z}}{dt} = J\mathbf{z}. \quad (45)$$

This is a linear dynamical system, and the state at time  $t + \Delta t$  is  $\mathbf{z}(t + \Delta t) = e^{J\Delta t}\mathbf{z}(t)$ .

Since  $J\Delta t$  is skew-symmetric, the matrix exponential  $e^{J\Delta t}$  is orthogonal. More precisely,  $\mathbf{z}(t + \Delta t)$  is a rotation about the origin of  $\mathbf{z}(t)$ .

# Properties of Hamiltonian dynamics

- 1. Reversibility:** The mapping from  $z(t) \rightarrow z(t + \Delta t)$  is one-to-one and invertible. To go from  $z(t + \Delta t)$  to  $z(t)$ , negate  $p(t + \Delta t)$ , apply the the Hamiltonian dynamics for  $\Delta t$  time, and negate the momentum again.
- 2. Conservation of energy:** The Hamiltonian (which is the total energy in a closed system) is conserved,

$$\frac{dH}{dt} = \sum_{d=1}^D \frac{dq_d}{dt} \frac{\partial H}{\partial q_d} + \frac{dp_d}{dt} \frac{\partial H}{\partial p_d} \quad (46)$$

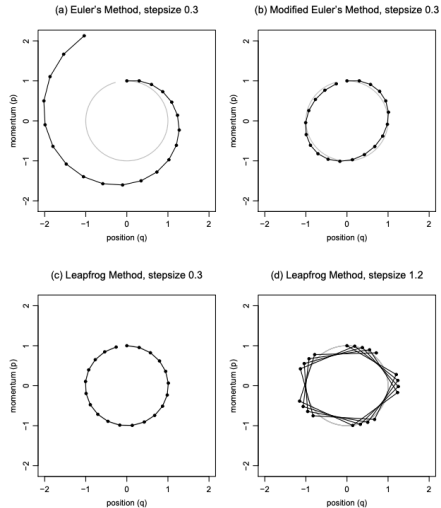
$$= \sum_{d=1}^D \frac{\partial H}{\partial p_d} \frac{\partial H}{\partial q_d} - \frac{\partial p_d}{\partial q_d} \frac{\partial H}{\partial p_d} = 0. \quad (47)$$

- 3. Symplectic/volume preserving** Let  $B$  be the determinant of the transformation from  $z(t) \rightarrow z(t + \Delta t)$ . It turns out that,

$$B^\top J^{-1} B = J^{-1} \quad (48)$$

which implies that  $|B^\top| |J^{-1}| |B| = |J^{-1}|$  and thus  $|B| = 1$ . I.e. the dynamics preserve volume.

# Discretizing Hamilton's equations with the Leapfrog integrator





# Using Hamiltonian dynamics for posterior inference

**Idea:** Let  $\mathbf{q} = \boldsymbol{\theta}$  and  $U(\mathbf{q}) = -\log p(\boldsymbol{\theta}, \mathcal{D})$ .

Then introduce auxiliary momentum variables  $\mathbf{p} \sim \mathcal{N}(0, \mathbf{M})$  so that  $\log p(\boldsymbol{\theta}, \mathcal{D}, \mathbf{p}) = \log p(\boldsymbol{\theta}, \mathcal{D}) + \log p(\mathbf{p})$ .

Use Hamiltonian dynamics to explore the extended state space of  $\mathbf{z} = (\mathbf{q}, \mathbf{p})$ .

Propose new states by sampling  $\mathbf{p}$  from its Gaussian marginal (and stationary) distribution; apply Hamiltonian dynamics for some amount of time; accept or reject based on the change in the Hamiltonian (which should be zero except for numerical integration errors).

The stationary distribution is  $p(\boldsymbol{\theta}, \mathbf{p}) = p(\boldsymbol{\theta} \mid \mathcal{D}) \times \mathcal{N}(\mathbf{p} \mid 0, \mathbf{M})$ .

**More to come next time!**