



ISTA 421/521

Introduction to Machine Learning

Lecture 19: Estimation: Sampling, Metropolis-Hastings

Clay Morrison

clayton@sista.arizona.edu

Gould-Simpson 819

Phone 621-6609

28 October 2014



$$g(\mathbf{w}; \mathbf{X}, \mathbf{t}, \sigma^2) = p(\mathbf{t}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\sigma^2)$$
$$p(\mathbf{w}|\mathbf{X}, \mathbf{t}, \sigma^2) = Z^{-1}g(\mathbf{w}; \mathbf{X}, \mathbf{t}, \sigma^2)$$
$$Z^{-1} = p(\mathbf{t}|\mathbf{X}, \sigma^2) = \int p(\mathbf{t}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\sigma^2) d\mathbf{w}$$

Our Options (when cannot compute *posterior* directly)

1. Find the single value of \mathbf{w} that corresponds to the highest value of the posterior. As $g(\mathbf{w}; \mathbf{X}, \mathbf{t}, \sigma^2)$ is proportional to the posterior, a maximum of g will also correspond to a maximum of the posterior. Z^{-1} is not a function of \mathbf{w} . MAP
2. Approximate $p(\mathbf{w}; \mathbf{X}, \mathbf{t}, \sigma^2)$ with some other density that we can compute analytically.
3. Sample directly from the posterior $p(\mathbf{w}; \mathbf{X}, \mathbf{t}, \sigma^2)$, knowing only $g(\mathbf{w}; \mathbf{X}, \mathbf{t}, \sigma^2)$



Method 1: MAP point estimate

- While we cannot derive a direct analytic posterior density that we can compute, we can compute something proportional to it:

$$g(\mathbf{w}; \mathbf{X}, \mathbf{t}, \sigma^2) = p(\mathbf{t}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\sigma^2)$$

- We will find the value of \mathbf{w} that maximizes g
- This will correspond to the value at the maximum of the posterior.
- This will be the most likely value $\hat{\mathbf{w}}$ under the posterior.
- In cases like logistic regression, need to estimate $\hat{\mathbf{w}}$ through an approximation method; we introduced gradient methods (Woodrow-Hoff, Newton-Raphson)



Using Newton-Raphson for MAP

$$\log g(\mathbf{w}; \mathbf{X}, \mathbf{t}) = \log p(\mathbf{t}|\mathbf{X}, \mathbf{w}) + \log p(\mathbf{w}|\sigma^2)$$

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$$

$$\mathbf{w}' = \mathbf{w} - \left(\frac{\partial^2 \log g(\mathbf{w}; \mathbf{X}, \mathbf{t})}{\partial \mathbf{w} \partial \mathbf{w}^\top} \right)^{-1} \frac{\partial \log g(\mathbf{w}; \mathbf{X}, \mathbf{t})}{\partial \mathbf{w}}$$

Point is guaranteed maximum if Hessian is negative definite
(as we showed for max likelihood)



Method 2: The Laplace* Approximation

- **The Idea:** *approximate* the density of interest with a Gaussian.
- (Recall that the Gaussian is used quite often in statistics to approximate other distributions!)
- However, **keep in mind:** our predictions will only be as good as our approximation – if the true posterior is not very Gaussian, then our predictions will be easy to compute but not very useful.

*Following the note in the book: the Machine Learning community has come to refer to the method this way, but this is elsewhere referred to as **saddle-point approx.**, and in statistics, the Laplace approx. is something different.



$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}) = \prod_{n=1}^N \left(\frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x}_n)} \right)^{t_n} \left(\frac{\exp(-\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(-\mathbf{w}^\top \mathbf{x}_n)} \right)^{1-t_n} \quad p(\mathbf{w}|\sigma^2) = \mathcal{N}(0, \sigma^2 \mathbf{I})$$

Approximating g using the Taylor Expansion

$$g(\mathbf{w}; \mathbf{X}, \mathbf{t}, \sigma^2) = p(\mathbf{t}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\sigma^2)$$

$$\log g(\mathbf{w}; \mathbf{X}, \mathbf{t}) = \log p(\mathbf{t}|\mathbf{X}, \mathbf{w}) + \log p(\mathbf{w}|\sigma^2)$$

$$\sum_{n=0}^{\infty} \frac{(\mathbf{w} - \hat{\mathbf{w}})^n}{n!} \frac{\partial^n f(\mathbf{w})}{\partial \mathbf{w}^n} \bigg|_{\hat{\mathbf{w}}}$$

Recall, the univariate Gaussian:

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

The log of the univ. G
(K is the normalizing constant):

$$\log(K) - \frac{1}{2\sigma^2}(w - \mu)^2$$

similar form!

Univariate version:

$$\mu = \hat{w} \quad \sigma^2 = 1/v$$

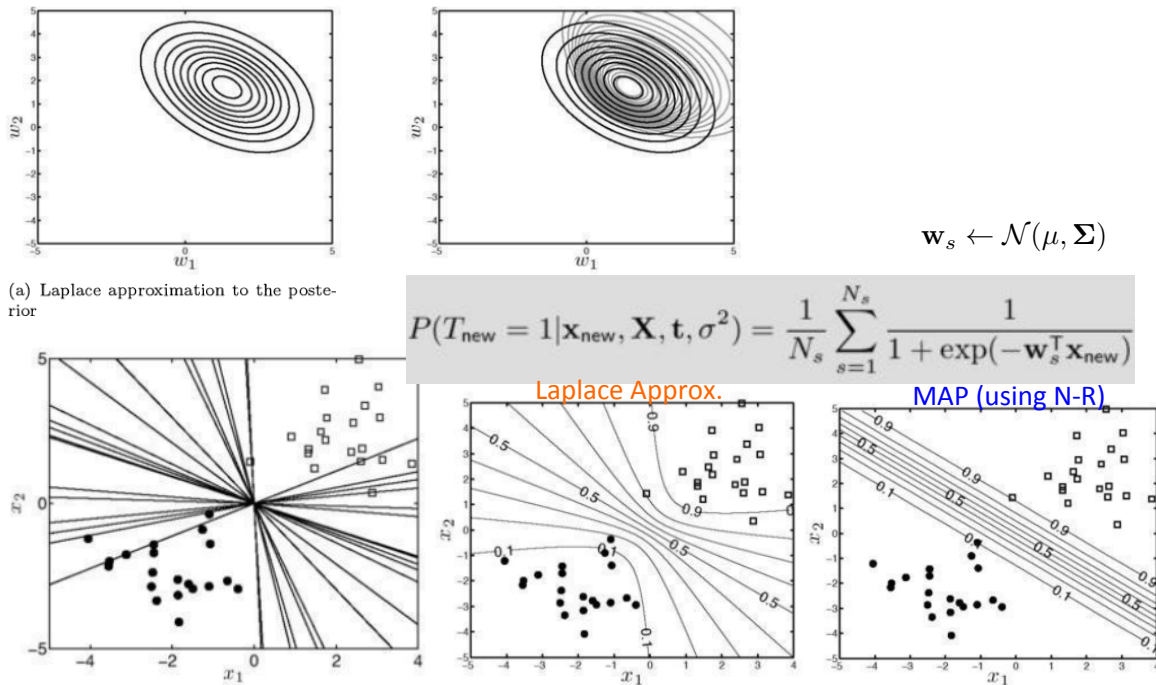
This is the Laplace approximation!
We approximate the posterior with a Gaussian that has its
mean at the posterior **mode** ($\hat{\mathbf{w}}$),
variance inversely proportional to the curvature of the posterior (g'') at its mode.

Multivariate version:

$$\mu = \hat{\mathbf{w}}, \quad \Sigma^{-1} = - \left(\frac{\partial^2 \log g(\mathbf{w}; \mathbf{X}, \mathbf{t})}{\partial \mathbf{w} \partial \mathbf{w}^\top} \right) \bigg|_{\hat{\mathbf{w}}}$$

$$\log g(\mathbf{w}; \mathbf{X}, \mathbf{t}, \sigma^2) \approx \log g(\hat{\mathbf{w}}; \mathbf{X}, \mathbf{t}, \sigma^2) - \frac{v}{2}(\mathbf{w} - \hat{\mathbf{w}})^2 \quad v = - \frac{\partial^2 \log g(\mathbf{w}; \mathbf{X}, \mathbf{t}, \sigma^2)}{\partial w^2} \bigg|_{\hat{\mathbf{w}}}$$

Laplace Approximation Example



Method 3: Sampling from Posterior

- Interest in Posterior density is to allow us to take *all* the uncertainty in \mathbf{w} into account when making predictions.

Posterior density over the parameters

$$P(T_{\text{new}} = 1 | \mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t}, \sigma^2) = \mathbf{E}_{p(\mathbf{w} | \mathbf{X}, \mathbf{t}, \sigma^2)} \{P(T_{\text{new}} = 1 | \mathbf{x}_{\text{new}}, \mathbf{w})\}$$

- Laplace method uses a *similar* density to provide an approximation of the posterior; still had to sample from it to estimate the integral.
- Now we'll look at **sampling directly** from the posterior.

The Intuition



\mathbf{y} : position of the dart

Δ : intended target

$p(\mathbf{y}|\Delta)$ Can be hard to model

$T = f(\mathbf{y})$ A new random variable:
 $T=1$: within 20
 $T=0$: outside of 20

$P(T = 1|\Delta)$

$$P(T = 1|\Delta) = \mathbb{E}_{p(\mathbf{y}|\Delta)}\{f(\mathbf{y})\} = \int f(\mathbf{y})p(\mathbf{y}|\Delta)d\mathbf{y}$$

OR: have your friend try to hit the 20, and find average hits!

$\mathbf{y}_s \leftarrow p(\mathbf{y}|\Delta)$

$$P(T = 1|\Delta) \simeq \frac{1}{N_s} \sum_{s=1}^{N_s} f(\mathbf{y}_s)$$



Expectation

$$\mathbb{E}_{p(z)}\{f(z)\} = \int f(z)p(z)dz$$

Normalization

$$\int p(z)dz = \mathbb{E}_{p(z)}\{\mathbf{1}_{\mathbb{R}}(z)\}$$

$$\mathbf{1}_S(x) = \begin{cases} 1 & \text{if } x \in S \\ 0 & \text{otherwise} \end{cases}$$

Probability

$$\begin{aligned} P(Z \leq z_0) &= \int_{-\infty}^{z_0} p(z)dz = \int_{-\infty}^{\infty} \mathbf{1}_{(-\infty, z_0]}(z)dz \\ &= \mathbb{E}_{p(z)}\{\mathbf{1}_{(-\infty, z_0]}(z)\} \end{aligned}$$

Marginalization

$$p(v) = \int p(v, z)dz = \int p(v|z)p(z)dz = \mathbb{E}_{p(z)}\{p(v|z)\}$$

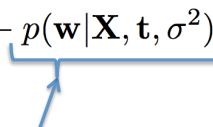
The **predictive distribution** is a marginalization:

$$\begin{aligned} P(T_{new} = 1 | \mathbf{x}_{new}, \mathbf{X}, \mathbf{t}, \sigma^2) &= \int P(T_{new} = 1 | \mathbf{x}_{new}, \mathbf{w}) p(\mathbf{w} | \mathbf{X}, \mathbf{t}, \sigma^2) d\mathbf{w} \\ &= \mathbb{E}_{p(\mathbf{w} | \mathbf{X}, \mathbf{t}, \sigma^2)} \{P(T_{new} = 1 | \mathbf{x}_{new}, \mathbf{w})\} \end{aligned}$$

... put another way: it's the **expectation** of the **likelihood** function under the **posterior** distribution

$$P(T_{new} = 1 | \mathbf{x}_{new}, \mathbf{X}, \mathbf{t}, \sigma^2) \simeq \frac{1}{N_s} \sum_{s=1}^{N_s} P(T_{new} = 1 | \mathbf{x}_{new}, \mathbf{w}_s)$$

$\mathbf{w}_s \leftarrow p(\mathbf{w} | \mathbf{X}, \mathbf{t}, \sigma^2)$



But how do we estimate this?

(A touch of theory on) Markov Chains

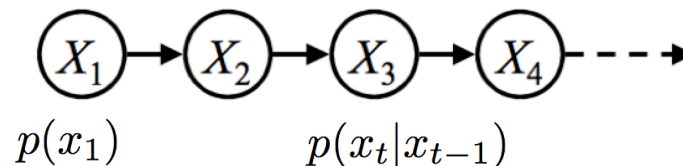
- A **Stochastic Process** is a collection of random variables indexed by a set T ; i.e., $\{X_t \mid t \in T\}$
 - Here, we only care about discrete-time stochastic processes, i.e., when T is a countable set.
 - For example, $T = \mathbb{Z}_+$ or $T = \{0, 1, 2, 3, 4, 5\}$
- **Example**
 - The results of 100 coin tosses is a stochastic process, with random variables C_1, \dots, C_{100}
 - The daily temperature is a stochastic process, represented by random variables T_1, T_2, \dots
- Stochastic processes are like any other sets of variables; we can talk about distributions:
 - $p_t(x_t)$, for any $t \in T$
 - $p(x_{t_1}, \dots, x_{t_n})$, for any $\{t_1, \dots, t_n\} \subset T$

(A touch of theory on) Markov Chains

- A (first-order) **Markov chain** is a discrete-time stochastic process with the **Markov Property**:

$$p(x_t | x_{t-1}, x_{t-2}, \dots, x_1) = p(x_t | x_{t-1})$$

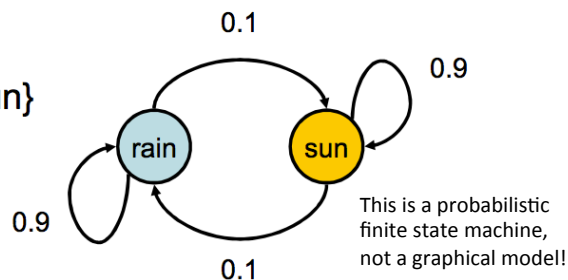
- That is, the variable at time t only depends on the variable at time $t-1$.
- Here, the conditional density $p(x_t | x_{t-1})$ (also called the transition kernel) is the same for all $t \in T$



Example Markov Chain

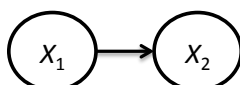
Weather:

- States: $X = \{\text{rain}, \text{sun}\}$
- Transitions:



- Initial distribution: 1.0 sun
- What's the probability distribution after one step?

$$P(X_2 = \text{sun}) = P(X_2 = \text{sun} | X_1 = \text{sun})P(X_1 = \text{sun}) + P(X_2 = \text{sun} | X_1 = \text{rain})P(X_1 = \text{rain})$$

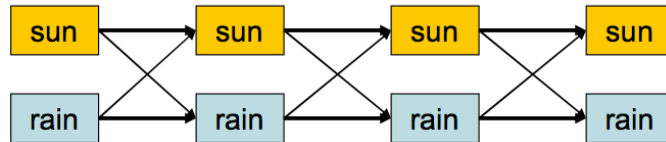


$$0.9 \cdot 1.0 + 0.1 \cdot 0.0 = 0.9$$

Join (product) of X_1 and X_2 , followed by sum (marginalization) of X_1

Mini “Forward” Algorithm

- Question: What's $P(X)$ on some day t ?
 - An instance of variable elimination!



$$P(x_t) = \sum_{x_{t-1}} P(x_t | x_{t-1}) P(x_{t-1})$$

$$P(x_1) = \text{known}$$

Forward simulation

Example Markov Chains

- From initial observation of sun

$$\begin{array}{cccc} \left\langle \begin{array}{c} 1.0 \\ 0.0 \end{array} \right\rangle & \left\langle \begin{array}{c} 0.9 \\ 0.1 \end{array} \right\rangle & \left\langle \begin{array}{c} 0.82 \\ 0.18 \end{array} \right\rangle & \longrightarrow \left\langle \begin{array}{c} 0.5 \\ 0.5 \end{array} \right\rangle \\ P(X_1) & P(X_2) & P(X_3) & P(X_\infty) \end{array}$$

- From initial observation of rain

$$\begin{array}{cccc} \left\langle \begin{array}{c} 0.0 \\ 1.0 \end{array} \right\rangle & \left\langle \begin{array}{c} 0.1 \\ 0.9 \end{array} \right\rangle & \left\langle \begin{array}{c} 0.18 \\ 0.82 \end{array} \right\rangle & \longrightarrow \left\langle \begin{array}{c} 0.5 \\ 0.5 \end{array} \right\rangle \\ P(X_1) & P(X_2) & P(X_3) & P(X_\infty) \end{array}$$

Example Markov Chains

- What if we had a different transition probability ?

$$\begin{array}{c}
 \text{sunny} \\
 \text{rainy}
 \end{array}
 \begin{array}{c}
 t \\
 t+1
 \end{array}
 \begin{array}{cc}
 \text{sunny} & \text{rainy} \\
 \begin{bmatrix} 0.9 & 0.1 \\ 0.5 & 0.5 \end{bmatrix}
 \end{array}$$

- If we set $\pi_0 = [1, 0]^\top$, i.e., today is sunny:
 - $\pi_2 = [0.844, 0.156]^\top$
 - $\pi_5 = [0.834, 0.166]^\top$
 - $\pi_{20} = [0.83333, 0.16667]^\top$
 - $\pi_{50} = [0.83333, 0.16667]^\top$
 - See a pattern?

Stationary Distribution

- A distribution π is **stationary** for a Markov chain if the transition kernel for the chain preserves π ; i.e., if for all $x_t \in R^d$

$$\int p(x_t | x_{t-1}) \pi(x_{t-1}) dx_{t-1} = \pi(x_t)$$

- Implication: if at any time t , $p_t(x_t) = \pi(x_t)$, then the marginals from that point on will be $\pi(x_t)$, since

$$\begin{aligned}
 \int p_{t+1}(x_{t+1}) &= \int p(x_{t+1} | x_t) p_t(x_t) dx_t \\
 &= \int p(x_{t+1} | x_t) \pi(x_t) dx_t \\
 &= \pi(x_{t+1})
 \end{aligned}$$

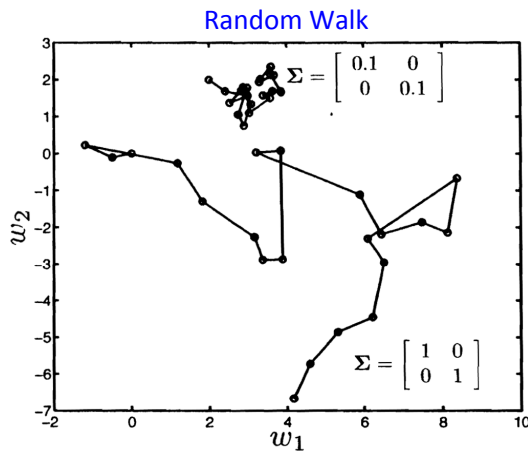
- Ergodicity**: guarantees a stationary distribution **exists** and is **unique**
 - Aperiodic**: can always transition back to state a having just transition from a to b (a state x has a period k if, starting from that state, it is only possible to return to it in multiples of k ; in that case, the x is said to be periodic.)
 - Irreducible**: It is possible to get to any state from any state (i.e., does not end up in a sink)
- Theorem**: If a Markov chain is irreducible and aperiodic, then it will have a unique stationary distribution
- So, how do we construct one for our posterior distribution of interest?

Named in the list of Top 10 Algorithms of the 20th Century (*SIAM News*, Vol 33, No 4, 2000)

Nicholas
physicist

Metropolis-Hastings Algorithm

W. Keith
statistician 1953 1970



Two desirable criteria for proposal distribution:

- (a) Easy to sample from
- (b) Symmetric:

$$p(\tilde{\mathbf{w}}_s | \mathbf{w}_{s-1}, \Sigma) = p(\mathbf{w}_{s-1} | \tilde{\mathbf{w}}_s, \Sigma)$$

$\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{s-1}, \mathbf{w}_s, \dots, \mathbf{w}_{N_s}$

- (0) Getting started: choosing \mathbf{w}_1
Turns out it doesn't matter: in theory, sample long enough and guaranteed to converge

Generating \mathbf{w}_s takes 2 steps:

- (1) Propose a new sample (based on previous)

$$\tilde{\mathbf{w}}_s \leftarrow p(\tilde{\mathbf{w}}_s | \mathbf{w}_{s-1}) \text{ proposal distribution}$$

Does **not** have to be related to target distribution!

Popular to use Gaussian centered on \mathbf{w}_{s-1}

$$p(\tilde{\mathbf{w}}_s | \mathbf{w}_{s-1}, \Sigma) = \mathcal{N}(\mathbf{w}_{s-1}, \Sigma)$$

- (2) Test whether to accept or reject $\tilde{\mathbf{w}}_s$

$$r = \frac{\overset{\text{posterior at proposed sample}}{p(\tilde{\mathbf{w}}_s | \mathbf{X}, \mathbf{t}, \sigma^2)} \overset{\text{ratio of Proposal densities}}{p(\mathbf{w}_{s-1} | \tilde{\mathbf{w}}_s, \Sigma)}}{\underset{\text{posterior at old sample}}{p(\mathbf{w}_{s-1} | \mathbf{X}, \mathbf{t}, \sigma^2)} p(\tilde{\mathbf{w}}_s | \mathbf{w}_{s-1}, \Sigma)}$$

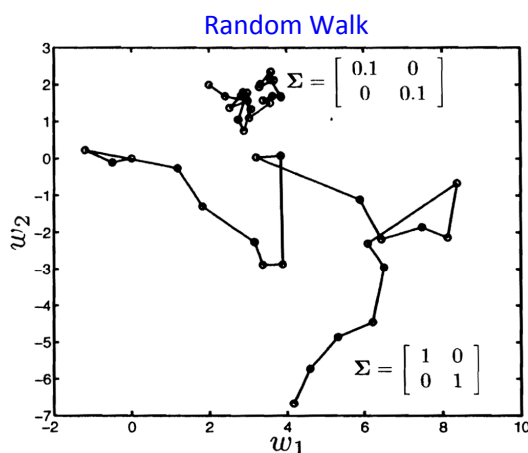
Problem?: Cannot directly compute posteriors!

Named in the list of Top 10 Algorithms of the 20th Century (*SIAM News*, Vol 33, No 4, 2000)

Nicholas
physicist

Metropolis-Hastings Algorithm

W. Keith
statistician 1953 1970



Two desirable criteria for proposal distribution:

- (a) Easy to sample from
- (b) Symmetric:

$$p(\tilde{\mathbf{w}}_s | \mathbf{w}_{s-1}, \Sigma) = p(\mathbf{w}_{s-1} | \tilde{\mathbf{w}}_s, \Sigma)$$

$\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{s-1}, \mathbf{w}_s, \dots, \mathbf{w}_{N_s}$

- (0) Getting started: choosing \mathbf{w}_1
Turns out it doesn't matter: in theory, sample long enough and guaranteed to converge

Generating \mathbf{w}_s takes 2 steps:

- (1) Propose a new sample (based on previous)

$$\tilde{\mathbf{w}}_s \leftarrow p(\tilde{\mathbf{w}}_s | \mathbf{w}_{s-1}) \text{ proposal distribution}$$

Does **not** have to be related to target distribution!

Popular to use Gaussian centered on \mathbf{w}_{s-1}

$$p(\tilde{\mathbf{w}}_s | \mathbf{w}_{s-1}, \Sigma) = \mathcal{N}(\mathbf{w}_{s-1}, \Sigma)$$

- (2) Test whether to accept or reject $\tilde{\mathbf{w}}_s$

$$r = \frac{g(\tilde{\mathbf{w}}_s; \mathbf{X}, \mathbf{t}, \sigma^2)}{g(\mathbf{w}_{s-1}; \mathbf{X}, \mathbf{t}, \sigma^2)} = \frac{p(\tilde{\mathbf{w}}_s | \sigma^2)}{p(\mathbf{w}_{s-1} | \sigma^2)} \frac{p(\mathbf{t} | \tilde{\mathbf{w}}_s, \mathbf{X})}{p(\mathbf{t} | \mathbf{w}_{s-1}, \mathbf{X})}$$

Don't need to calculate posteriors directly because the Marginal Likelihoods cancel in the ratio!

$r \geq 1$? If **yes**: always choose best: $\mathbf{w}_s = \tilde{\mathbf{w}}_s$

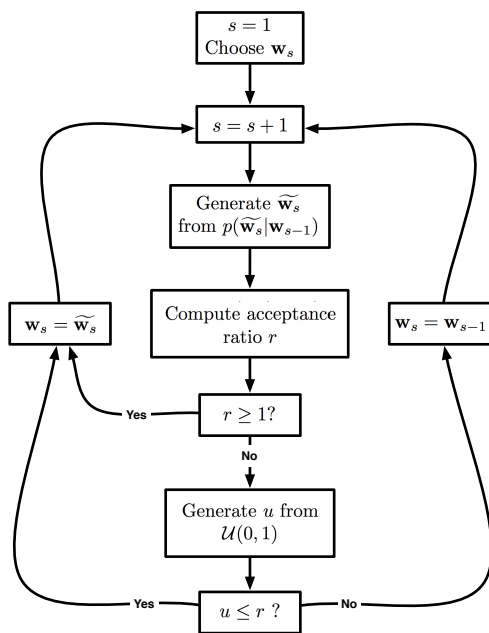
If **no**: possibly accept anyway

$u \leftarrow \mathcal{U}(0, 1), u \leq r$? **yes** $\rightarrow \mathbf{w}_s = \tilde{\mathbf{w}}_s$ **no** $\rightarrow \mathbf{w}_s = \mathbf{w}_{s-1}$

Nicholas
physicist

Metropolis-Hastings Algorithm

W. Keith
1953 statistician 1970



$\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{s-1}, \mathbf{w}_s, \dots, \mathbf{w}_{N_s}$

(0) Getting started: choosing \mathbf{w}_1

Turns out it doesn't matter: in theory, sample long enough and guaranteed to converge

Generating \mathbf{w}_s takes 2 steps:

(1) Propose a new sample (based on previous)

$$\tilde{\mathbf{w}}_s \leftarrow p(\tilde{\mathbf{w}}_s | \mathbf{w}_{s-1}) \text{ proposal distribution}$$

Does **not** have to be related to target distribution!

Popular to use Gaussian centered on \mathbf{w}_{s-1}

$$p(\tilde{\mathbf{w}}_s | \mathbf{w}_{s-1}, \Sigma) = \mathcal{N}(\mathbf{w}_{s-1}, \Sigma)$$

(2) Test whether to accept or reject $\tilde{\mathbf{w}}_s$

$$r = \frac{g(\tilde{\mathbf{w}}_s; \mathbf{X}, \mathbf{t}, \sigma^2)}{g(\mathbf{w}_{s-1}; \mathbf{X}, \mathbf{t}, \sigma^2)} = \frac{p(\tilde{\mathbf{w}}_s | \sigma^2)}{p(\mathbf{w}_{s-1} | \sigma^2)} \frac{p(\mathbf{t} | \tilde{\mathbf{w}}_s, \mathbf{X})}{p(\mathbf{t} | \mathbf{w}_{s-1}, \mathbf{X})}$$

Don't need to calculate posteriors directly because the Marginal Likelihoods cancel in the ratio!

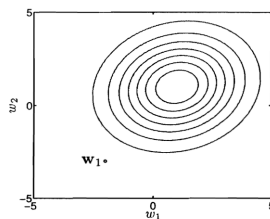
$r \geq 1$? If **yes**: always choose best: $\mathbf{w}_s = \tilde{\mathbf{w}}_s$

If **no**: possibly accept anyway

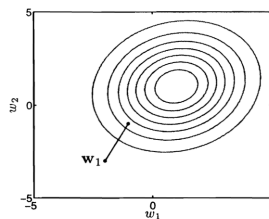
$u \leftarrow \mathcal{U}(0, 1), u \leq r$?

yes $\rightarrow \mathbf{w}_s = \tilde{\mathbf{w}}_s$
no $\rightarrow \mathbf{w}_s = \mathbf{w}_{s-1}$

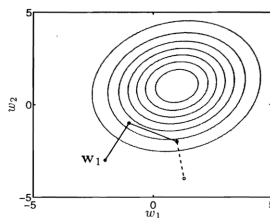
Metropolis-Hastings Example 1



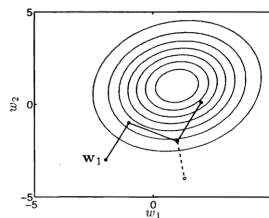
(a) Starting point



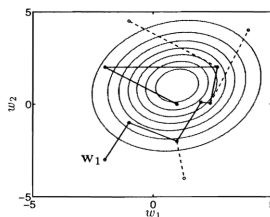
(b) After one sample



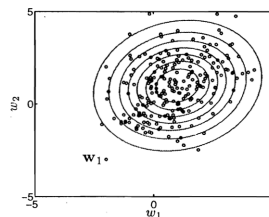
(c) After three samples. $\tilde{\mathbf{w}}_3$ was accepted, $\tilde{\mathbf{w}}_4$ rejected (dashed line)



(d) After four samples



(e) After 10 samples



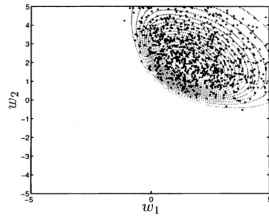
(f) The first 300 samples

$$\boldsymbol{\mu} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \mathbf{S} = \begin{bmatrix} 3 & 0.4 \\ 0.4 & 3 \end{bmatrix}$$

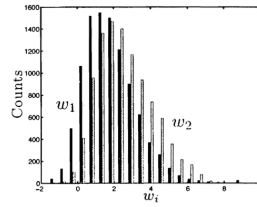
$$\boldsymbol{\mu}' = \frac{1}{N_s} \sum_{s=1}^{N_s} \mathbf{w}_s, \mathbf{S}' = \frac{1}{N_s} \sum_{s=1}^{N_s} (\mathbf{w}_s - \boldsymbol{\mu}')(\mathbf{w}_s - \boldsymbol{\mu}')^T$$

$$\boldsymbol{\mu}' = \begin{bmatrix} 0.9770 \\ 1.0928 \end{bmatrix}, \mathbf{S}' = \begin{bmatrix} 3.0777 & 0.4405 \\ 0.4405 & 2.8983 \end{bmatrix}$$

Metropolis-Hastings Example 2



(a) 1000 of the MH samples along with the posterior contours



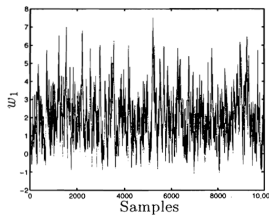
(b) Histograms of the samples for both w_1 (black) and w_2 (grey)

Burn In

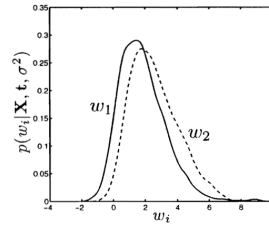
Convergence

Estimating predictions from samples \mathbf{w}_s

$$P(T_{\text{new}} = 1 | \mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t}, \sigma^2) = \frac{1}{N_s} \sum_{s=1}^{N_s} \frac{1}{1 + \exp(-\mathbf{w}_s^T \mathbf{x}_{\text{new}})}$$

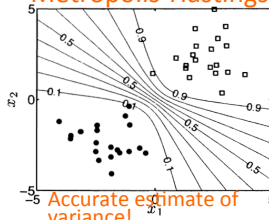


(c) All of the w_1 samples plotted against iteration, s



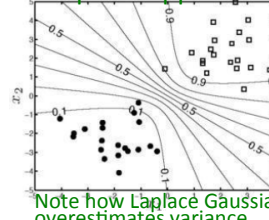
(d) Continuous densities fitted to the w_1 and w_2 samples

Metropolis-Hastings



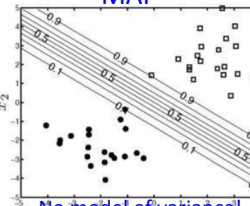
Accurate estimate of variance!

Laplace Approx.



Note how Laplace Gaussian overestimates variance

MAP



No model of variance!