Using This MRF we can see that all $v_1, ..., v_n$ we (conditionally independent given C. This will allow us to use the known (appearance models to get $P(C | v_1, ..., v_n)$

If $v_i \perp v_j | C$ for all $i, j$ then

$$P(v_1, ..., v_n | C) = \prod_{i=1}^{h} P(v_i | C) = \prod_{i=1}^{h} P(v_i | C)$$

Using Bayes Rule:

$$P(C | v_1, ..., v_n) = \frac{P(C, v_1, ..., v_n)}{P(v_1, ..., v_n)} = \frac{P(v_1, ..., v_n | C) \, P(C)}{P(v_1, ..., v_n)}$$

Now suppose $C = \{C_{-i} = i, C_{-j} = j, C...\}$
and we swap $C_i$ and $C_j$ to get $C' = \{C_{-i} = j, C_{-j} = i, C...\}$

Then $A(C'|C) = \min\left(1, \frac{P(v_1, ..., v_n | C') \cdot P(C')}{P(v_1, ..., v_n)}\right)$
$$\frac{P(v_1, ..., v_n | C) \cdot P(C)}{P(v_1, ..., v_n)}$$

$$= \min\left(1, \frac{P(v_1, ..., v_n | C') \cdot P(C')}{P(v_1, ..., v_n | C) \cdot P(C)}\right) \rightarrow P(C) = P(C') = \frac{1}{N!} \rightarrow \frac{P(C')}{P(C)}$$

$$= \min\left(1, \frac{P(v_1, ..., v_n | C_{-i} = j, C_{-j} = i, C...)}{P(v_1, ..., v_n | C_{-i} = i, C_{-j} = j, C...)}\right)$$

$$\boxed{= \min\left(1, \frac{P(v_i | C_{-i} = j) \cdot P(v_j | C_{-j} = i)}{P(v_i | C_{-i} = i) \cdot P(v_j | C_{-j} = j)}\right)}$$

# 1) Data Association

- $K$ objects $u_1, \ldots, u_k$
- $H$ observations $v_1, \ldots, v_H$, $\text{Val}(v_i) = \{a_1, \ldots, a_k\}$
- $H$ correspondence variables $c_1, \ldots, c_H$, where $\text{Val}(c_i) = \{1, \ldots, H\}$
- known appearance model for each object $u_k$, $P_H(v_i = a_j | c_i = H)$

## a) Compute acceptance probability for each MH step

$$A(x'|x) = \min\left(1, \frac{P(x') \cdot Q(x|x')}{P(x) \cdot Q(x'|x)}\right)$$
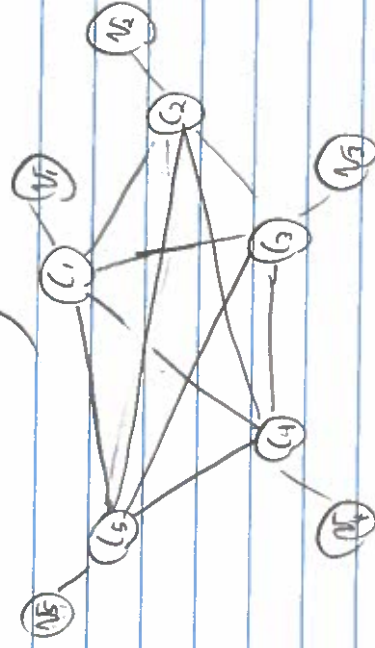
$$= \min\left(1, \frac{P(c'|v_1, \ldots, v_H) \cdot Q(c|c')}{P(c|v_1, \ldots, v_H) \cdot Q(c'|c)}\right)$$

$Q$ models the probability of sampling a new assignment
$c' = \{c_1', c_2', \ldots, c_H'\}$ given $c = \{c_1, c_2, \ldots, c_H\}$
$Q(c'|c)$ is $0$ if $c'$ is not $c$ w/ exactly two values swapped
$Q(c'|c)$ is uniform otherwise

Since $Q(c|c')$ and $Q(c'|c)$ are equally likely,
these cancel to $1$ in $A$.

$$A(c'|c) = \min\left(1, \frac{P(c'|v_1, \ldots, v_H)}{P(c|v_1, \ldots, v_H)}\right)$$

We can use the following MRF to describe the problem:

$c_i$ are fully connected
$v_i$ are only connected to
corresponding $c_i$

b) We run MH samples u long time and collect $M$ samples $(c_1[m], \ldots, c_k[m])$ for $m=1,\ldots,M$ after mixing

$$P(c_i | v_1, \ldots, v_n) = P(c_i = h | v_1, \ldots, v_n) \text{ for } h=1,\ldots,K$$

Let $g(c_i) = \mathbb{1}[c_i = h]$

$$\hat{g}(c_i[1], \ldots, c_i[M]) = \frac{1}{M}\sum_{m=1}^{M} g(c_i[m])$$

c) For Gibbs sampling to work, it must be the that

$$A(c'|c) = 1 = \min(1, 1)$$

$$= \frac{P(v_i | c_{-i} = j) \cdot P(v_j | c_i = i)}{P(v_i | c_{-i}) \cdot P(v_j | c_{-j})}$$

This will not work, because these is no guarantee this quantity $= 1$

b) Derive $\dfrac{\partial}{\partial \theta_i} J(\theta; D)$

$$\dfrac{\partial}{\partial \theta_i} J(\theta; D) = \dfrac{1}{|D|} \sum_{x \in D} f_i(x, y) - (1-a) \cdot \left[ \dfrac{1}{|D|} \sum_{x \in D} E_Q[f_i(x,y)] \right]$$

$$- a \cdot \left[ \dfrac{1}{|D|} \sum_{y \in D} E_Q[f_i(x,y)] \right]$$

If $Q$ is the empirical distribution of our dataset $\hat{P}$, or a conditional distribution of the form $P_\theta(W | Z = z)$, this $E_Q$ will apply.

## M-Step 1:

$$\Theta_c' = \frac{\sum \text{weight}(i)}{\#\text{datapoints}} = \frac{M \cdot |\text{Val}(c)|}{M \cdot |\text{Val}(c)|} = \frac{1}{|\text{Val}(c)|^2}$$

$$\Theta_{x_i c}' = \frac{\#\text{datapoints where } X = x_i \text{ and } C = c \cdot \text{weight}}{\#\text{datapoints where } C = c \cdot \text{weight}} \quad \text{— denote as } t_{ij}$$

$$= \frac{t_{ij} \cdot \frac{1}{|\text{Val}(c)|}}{M \cdot \frac{1}{|\text{Val}(c)|}} = \frac{t_{ij}}{M}, \quad \text{where } t_{ij} \leq M$$

To show EM has converged, we will run one more E-step:

## E-step 2:

$$w = \frac{1}{2} \cdot \frac{1}{|\text{Val}(c)|^2} \cdot \prod_{i=1}^{n} \frac{t_{ij}}{M} = \frac{1}{2} \cdot \frac{1}{|\text{Val}(c)|^2} \cdot \frac{1}{M} \cdot \prod_{i=1}^{n} t_{ij}$$

$$z = \sum_{h=1}^{|\text{Val}(c)|} \frac{1}{|\text{Val}(c)|^2} \cdot \frac{1}{M} \cdot \prod_{i=1}^{n} t_{ij} = \frac{1}{|\text{Val}(c)|} \cdot \frac{1}{M} \cdot \prod_{i=1}^{n} t_{ij}$$

$$w = \frac{\frac{1}{|\text{Val}(c)|^2} \cdot \frac{1}{M} \cdot \prod_{i=1}^{n} t_{ij}}{\frac{1}{|\text{Val}(c)|} \cdot \frac{1}{M} \cdot \prod_{i=1}^{n} t_{ij}} = \frac{|\text{Val}(c)|}{|\text{Val}(c)|^2} = \frac{1}{|\text{Val}(c)|}$$

Since the value for all weights is the same as the previous E-step, EM has converged.

The final parameter values are:

$$\Theta_c = \frac{1}{|\text{Val}(c)|^2}$$

$$\Theta_{x_i c} = \frac{t_{ij}}{M}$$

## 3) Expectation Maximization in a Naive Bayes Model



Class variable $C$, discrete evidence variables $X_1,...,X_n$

CPDs parametrized: $P(C=c) = \theta_c$, $P(X_i=x|C=c) = \theta_{x|c}$
for $i=1,...,n$
and for all assignments $x_i \in Val(X_i)$, classes $c \in Val(C)$

Given dataset $D = \{x[1],...,x[M]\}$ each $x[m]$ assigned to $X_1,...,X_n$
$C$ not observed

E-Step 1: $X[m] = (X_1=x_1, X_2=x_2,...,X_n=x_n)$ $(C=c_1)$
$($ " $)$ $(C=c_2)$
$($ " $)$ $\vdots$
$($ " $)$ $(C=c_{|Val(C)|})$

Expand $X[m]$ into $|Val(C)|$ datapoints, one for each $c \in C$

$$U = \frac{1}{Z} \theta_c^0 \prod_{i=1}^n \theta_{x_i|c} = \frac{1}{Z} \cdot \frac{1}{|Val(C)|} \cdot \prod_{i=1}^n \frac{1}{|Val(X_i)|}$$

$$Z = \sum_{j=1}^{|Val(C)|} \frac{1}{|Val(C)|} \cdot \prod_{i=1}^n \frac{1}{|Val(X_i)|} = \prod_{i=1}^n \frac{1}{|Val(X_i)|}$$

$$\rightarrow U = \frac{\frac{1}{|Val(C)|} \cdot \prod_{i=1}^n \frac{1}{|Val(X_i)|}}{\prod_{i=1}^n \frac{1}{|Val(X_i)|}} = \frac{1}{|Val(C)|}$$

weight $U$ is same for all datapoints

---

M-step:
$$\theta_{c_i} = \frac{\sum weight(c_i)}{\# data} = \frac{M \cdot \frac{1}{|Val(C)|}}{M \cdot |Val(C)|} = \frac{1}{|Val(C)|^2}$$

$$\theta_{x_i|c_i} = \frac{M \cdot \frac{1}{|Val(C)|}}{M \cdot \frac{1}{|Val(C)|}} = \frac{\frac{1}{|Val(C)|}}{\frac{1}{|Val(C)|}}$$

E-step: $U = \frac{1}{Z} \cdot \frac{1}{|Val(C)|^2} \cdot \prod_{i=1}^n t_i = \frac{1}{Z} \cdot \frac{1}{|Val(C)|^3} \cdot \frac{1}{M} \cdot \prod_{i=1}^n t_i$

$$Z = |Val(C)| \cdot \frac{1}{|Val(C)|^3} \cdot \frac{1}{M} \cdot \prod_{i=1}^n t_i$$

$$= \frac{1}{|Val(C)|} \cdot \frac{1}{M} \cdot \prod_{i=1}^n t_i$$

$$\rightarrow U = \frac{\frac{1}{|Val(C)|^2} \cdot \frac{1}{M} \cdot \prod_{i=1}^n t_i}{\frac{1}{|Val(C)|} \cdot \frac{1}{M} \cdot \prod_{i=1}^n t_i} = \frac{\frac{1}{|Val(C)|^2}}{\frac{1}{|Val(C)|}} = \frac{1}{|Val(C)|}$$

a) Derive expression for cond. prob pixel $(i,j)$ black given MB

$$E(y,x) = -\eta \sum_{i,j} y_{i,j} x_{i,j} - \beta \sum_{(i,j),(i',j') \in E} y_{i,j} y_{i',j'}$$

$$P(y,x) = \frac{1}{Z} \exp(-E(y,x)) \quad, \quad \text{let} \quad y_{-\ell} = y_{-y_{i,j}}$$

$$P(y_{i,j} = 1 \mid y_m, (i,j)) = P(y_{i,j} = 1 \mid y_{-\ell}, x)$$

$$= \frac{P(y_{i,j} = 1, y_{-\ell}, x)}{P(y_{-\ell}, x)}$$

$$= \frac{\exp(-E(y_{i,j} = 1, y_{-\ell}, x))}{\exp(-E(y_{i,j} = 1, y_{-\ell}, x)) + \exp(-E(y_{i,j} = -1, y_{-\ell}, x))}$$

$$= \frac{\exp\left(\eta\left(\sum_{i,j} y_{i,j} x_{i,j} + x_\ell\right) + \beta\left(\sum_{(i,j),(i',j') \in E} y_{i,j} y_{i',j'} + \sum_{y_m} y_{i,j}\right)\right)}{\text{---}}$$

$$\frac{}{1 + \exp\left(\eta\left(\sum_{i,j} y_{i,j} x_{i,j} - x_\ell\right) + \beta\left(\sum_{\in E} y_{i,j} y_{i',j'} - \sum_{y_m} y_{i,j}\right)\right)}$$

$$\eta \sum_{i,j} y_{i,j} x_{i,j} + \eta x_\ell + \beta \sum_{\in E} y_{i,j} y_{i',j'} + \beta y_{i,j}$$

$$\rightarrow \exp\left(\eta x_\ell + \beta \sum_{y_m} y_{i,j}\right) \cdot \exp\left(\eta \sum_{i,j} y_{i,j} x_{i,j} + \beta \sum_{\in E} y_{i,j} y_{i',j'}\right)$$

Let $a = \eta x_\ell + \beta \sum_{y_m} y_{i,j}$, $b = \eta \sum_{i,j} y_{i,j} x_{i,j} + \beta \sum_{\in E} y_{i,j} y_{i',j'}$

Then $P(y_{i,j} = 1 \mid y_m, (i,j)) = \dfrac{e^a \cdot e^b}{e^a \cdot e^b + e^{-a} \cdot e^b} = \dfrac{e^a \cdot e^b}{e^b \cdot (e^a + e^{-a})}$

$$= \frac{e^a}{e^a + e^{-a}} \cdot \frac{e^{-a}}{e^{-a}} = \frac{1}{1 + e^{-2a}} = \sigma(2a)$$

$$= \sigma\left(2\left(\eta x_\ell + \beta \sum_{y_m} y_{i,j}\right)\right)$$

# CS228: Probabilistic Graphical Models

## Homework 4

Luke Jaffe

Due: 03/03/2017
Submitted: 03/03/2017

**Problem 4: Programming Assignment**

a) **TODO:** rewrite and scan

b)

i. Outline a Gibbs sampling algorithm (in pseudocode) that iterates over the pixels in the image and samples each $y_{ij}$ given its Markov Blanket.
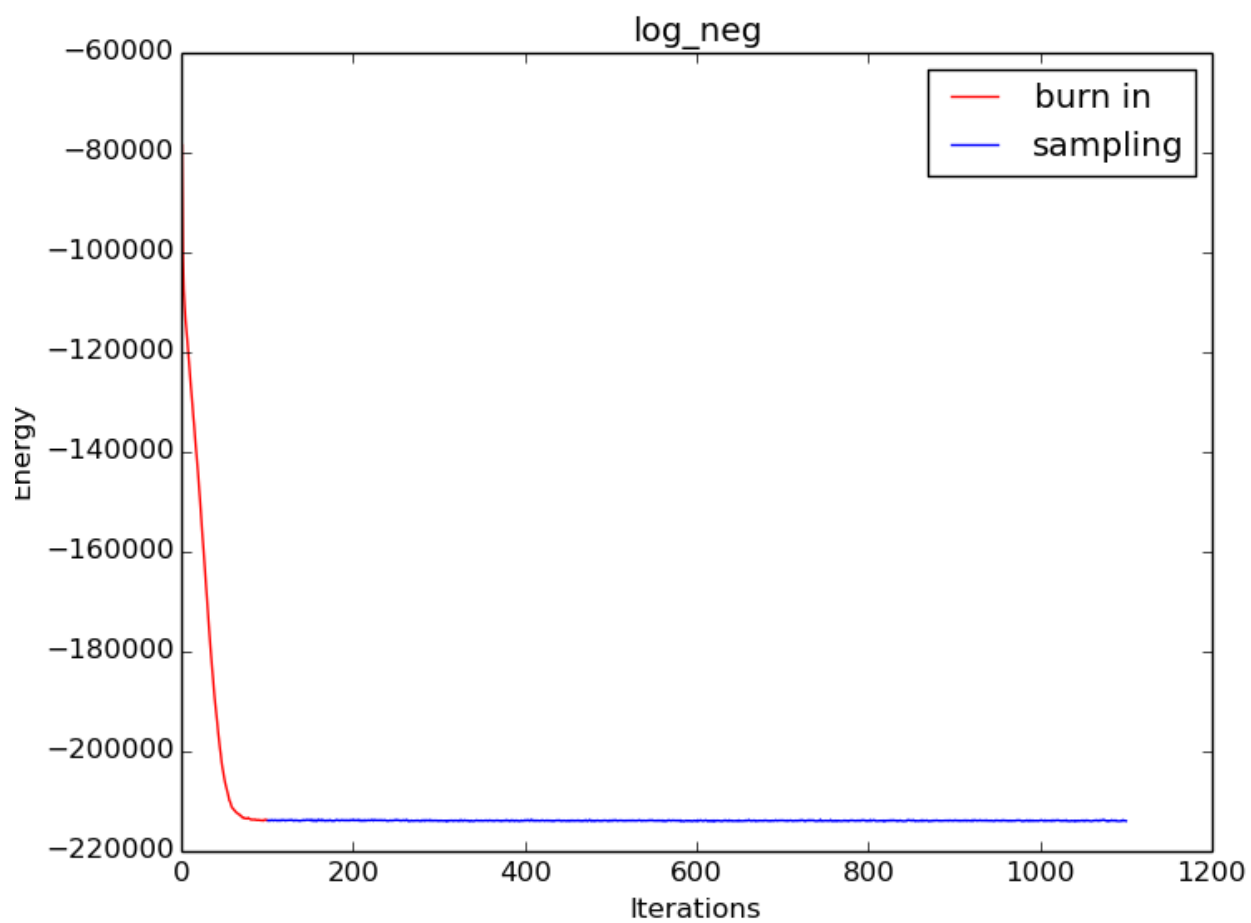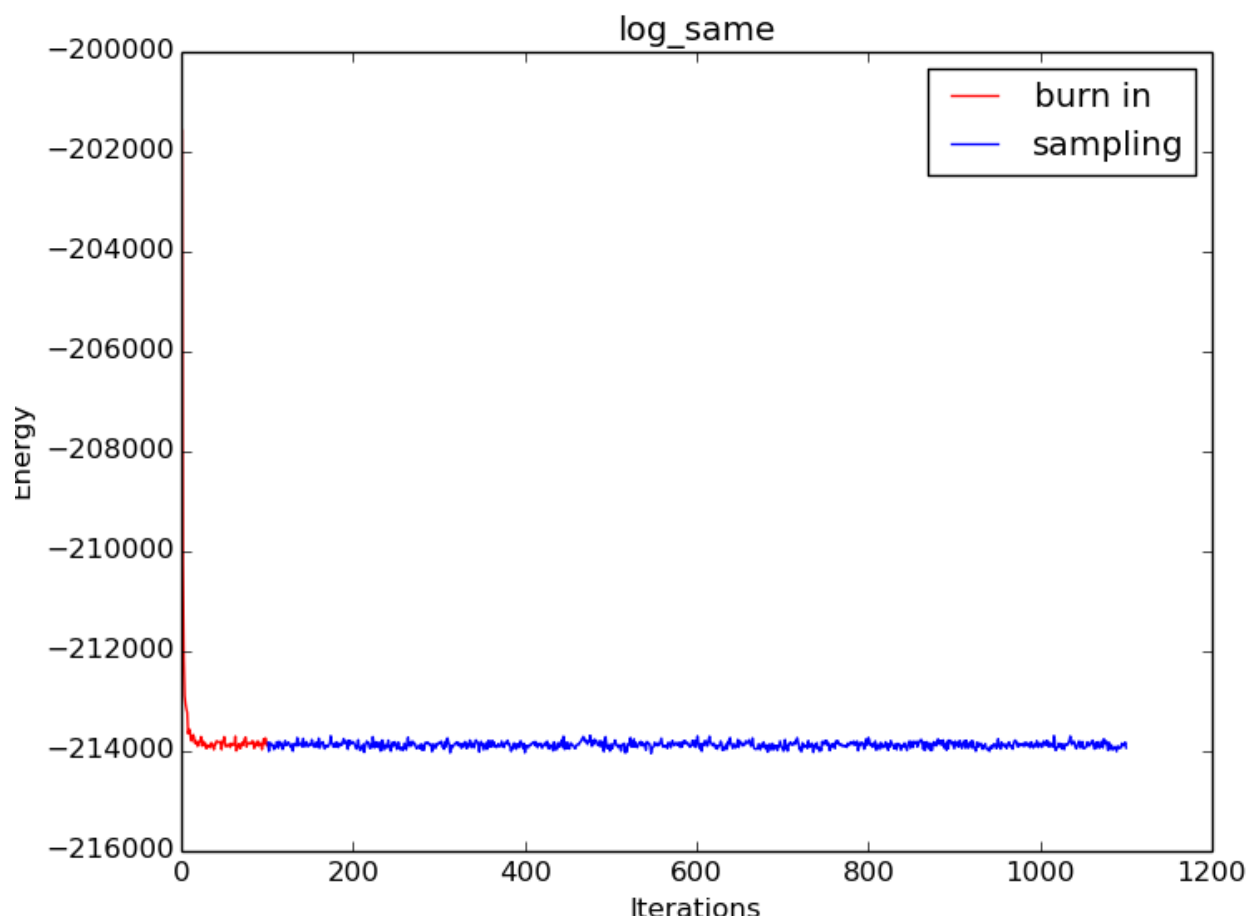
```
given constants eta, beta

function sample(Y, X):
    initialize matrix X to the noisy image
    initialize matrix Y randomly (same size as X)
    for I = 1 to N:
        for j = 1 to M:
            x_term = eta*X[i][j]
            y_term = 0
            for y in markov_blanket(Y, i, j):
                y_term += y
            y_term *= beta
            a = x_term + y_term
            s = sigmoid(2*a)
            Y[i][j] = s

function gibbs(Y, X, B, N):
    for t = 1 to B:
        sample(Y, X)
    for t = 1 to N
        sample(Y, X)
    return Y
```
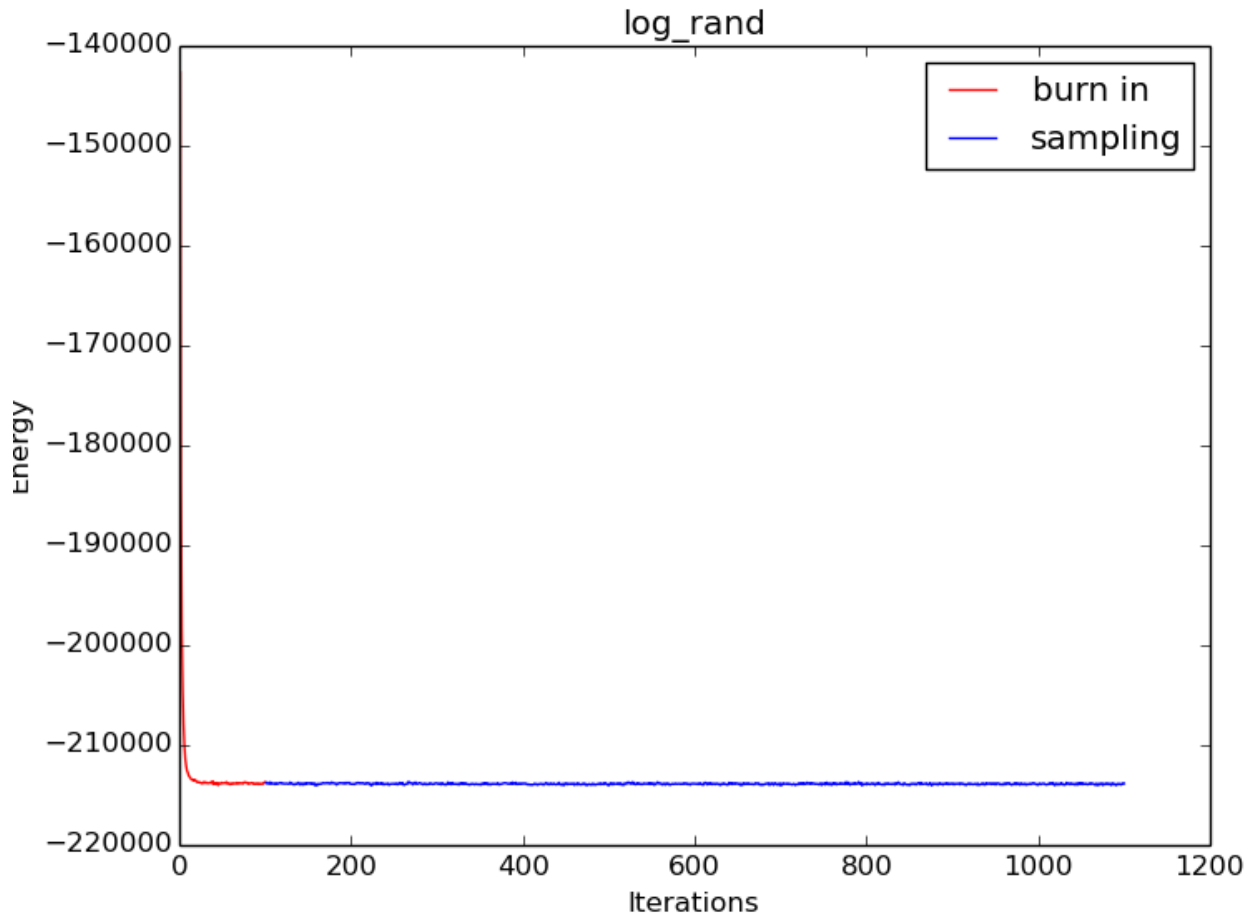
ii. How can we show in our case that the equilibrium distribution is in fact the posterior distribution $p(y|x)$?

We can show this using the Monte Carol equation.

c)

i. Do all three seem to be converging to the same general region of the posterior, or are some obviously suboptimal?

While all three methods do not converge in exactly the same way, they all converge to the same general region of the posterior after some iterations. The negative initialization may be considered suboptimal since it takes longer to converge than the same or random initialization methods.

ii. Does the burn-in seem to be adequate in length?

Yes, the burn-in is more than sufficient in length. The same and random initializations seem to converge in less than 20 iterations, while the negative initialization needs 80 or so.

iii. Is there substantial fluctuation from iteration to iteration, indicating that the chain is mixing well, or does it become stuck at particular energies for several iterations at a time?
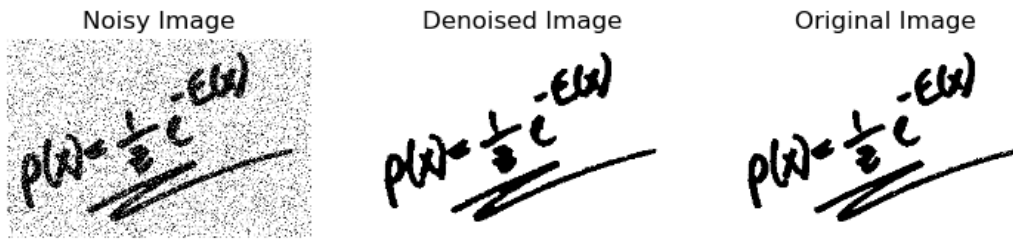
Yes, there is substantial fluctuation between iterations for all methods. In the plots above, the "same" method may appear to fluctuate the most, but this is just because of the y-axis scale.
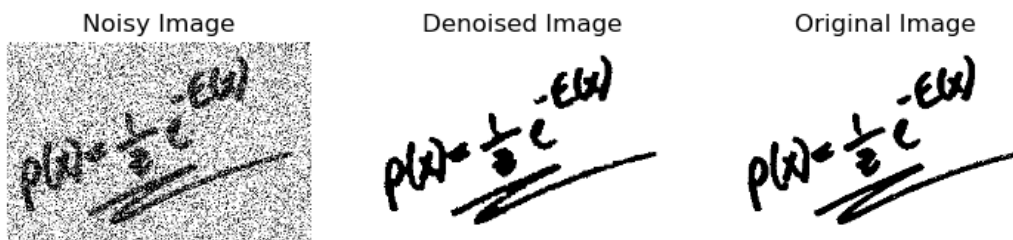
d)

Denoised 10% error: 0.0059
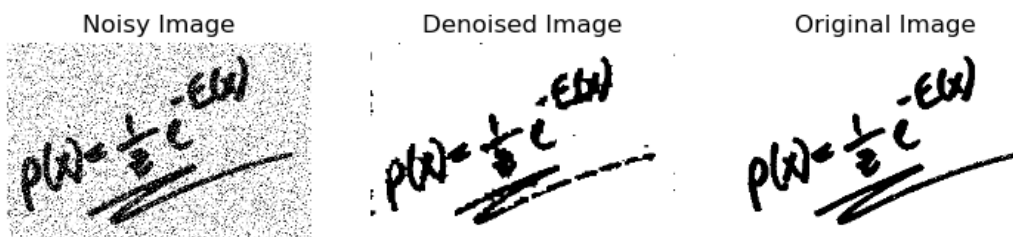Denoised 20% error: 0.0104

## denoised_10% plots
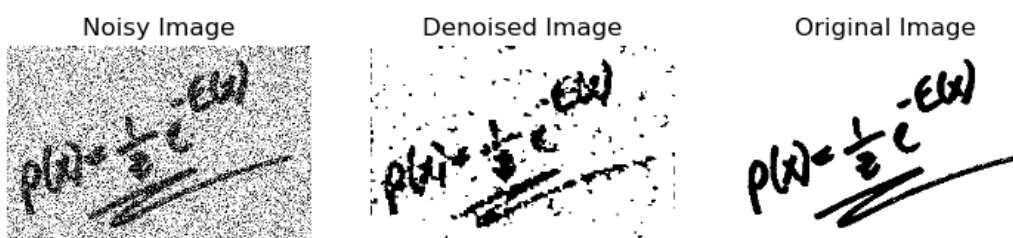


## denoised_20% plots



e)

Denoised dumb 10% error: 0.0213
Denoised dumb 20% error: 0.0588
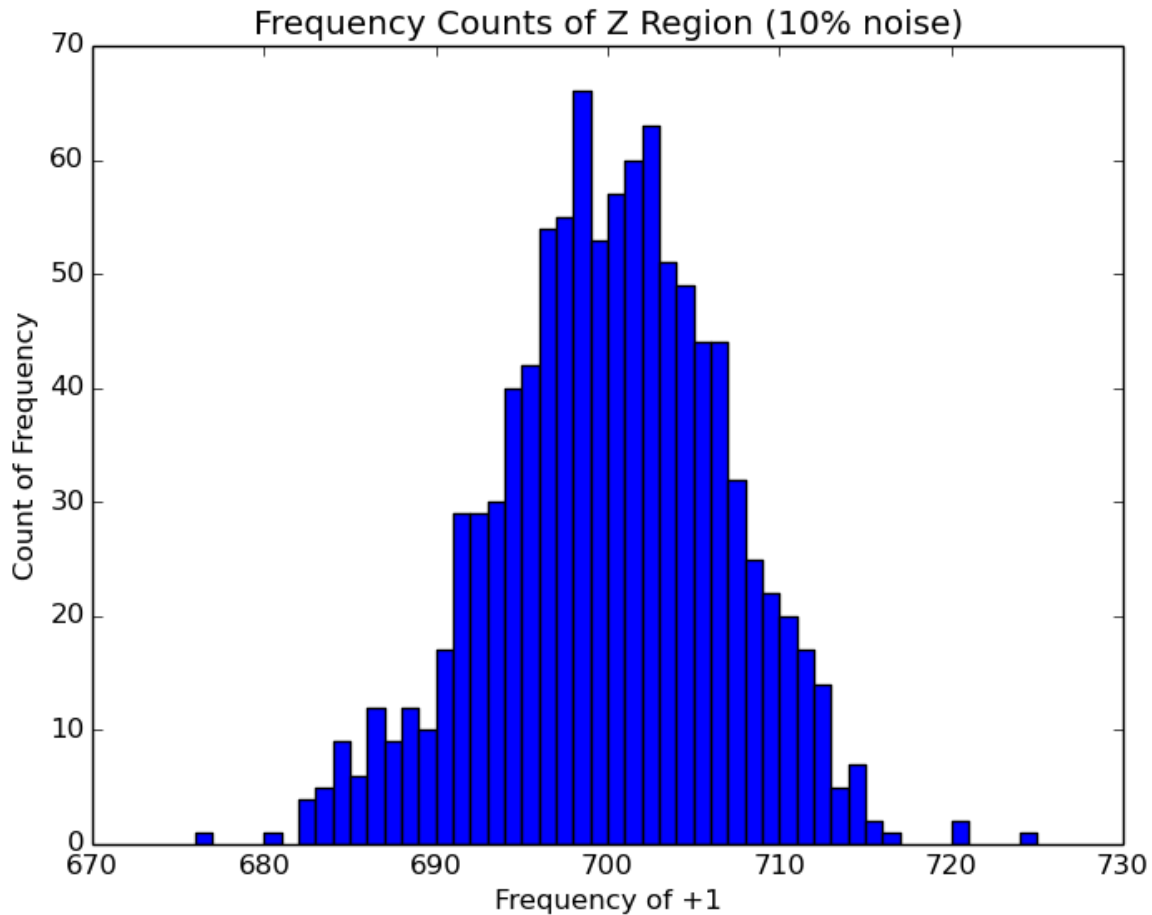
## denoised_dumb_10% plots
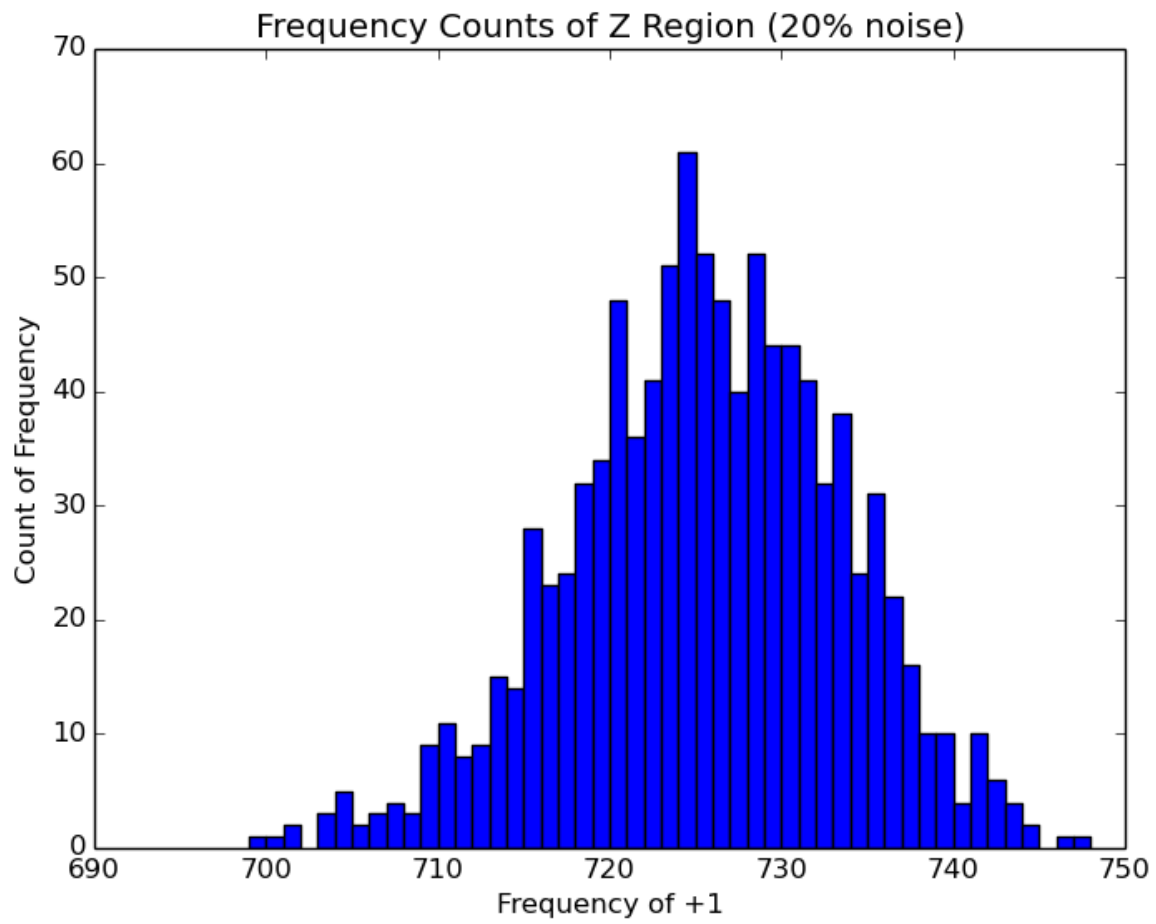


## denoised_dumb_20% plots

i. Does the Gibbs sampler do better than the trivial algorithm? Why or why not?

Yes, the Gibbs sampler does better than the trivial algorithm. The error is lower in the 10% noise case (0.5% v. 2%) and the 20% noise case (1% v. 6%). The images for the Gibbs sampler are also visibly better denoised. This is because the Gibbs sampler is built on a model which to some degree correctly codifies the spatial locality relationships of the actual imagery, whereas the other method is naive.

f)



Frequency Counts of Z Region (10% noise)

Frequency Counts of Z Region (20% noise)

i. The distribution of frequencies for the noisier case is slightly wider (greater variance), and has a greater mean (~725 v. ~700).