

Probabilistic Programming

Homework 1

- <https://www.cs.ubc.ca/~fwood/CSS32W-539W/homework1.html>

Question 1

Show that the Gamma distribution is conjugate to the Poisson distribution.

The Gamma distribution is over parameter $\lambda \in (0, \infty)$ and the Poisson distribution is over $k \in \{0, 1, 2, \dots\}$ with parameter $\lambda \in (0, \infty)$. Thus it makes sense to treat Gamma as a prior $p(\lambda)$ over $\lambda = x$ and Poisson as the likelihood $p(k|\lambda)$. In this case the questions is asking for us to show that the posterior $p(k|\lambda)$

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

$$p(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} \exp[-\beta\lambda]$$

$$k|\lambda \sim \text{Poisson}[\lambda]$$

$$p(k|\lambda) = \frac{\lambda^k \exp^{-\lambda}}{k!}$$

We can already see a hint of the conjugacy through the λ to some power, and λ in exp.

Let's apply Bayes rule.

$$p(\lambda|k) = \frac{p(k|\lambda)p(\lambda)}{p(k)}$$

Conjugacy is shown by the posterior being Gamma distributed: $\lambda \sim \text{Gamma}[\alpha', \beta']$.

Gamma and Poisson have normalization constants, and the evidence $p(k)$ is unknown (but doesn't depend on λ because it's been marginalized out). However, because the posterior is a distribution over λ , we can drop all multiplicative constants over λ , α, β . If we can show that the posterior has the unnormalized form of the Gamma distribution with some new parameters α' and β' , which are constant in λ , but can contain k, α, β , then we can use the known form for the normalization constant.

$$p(\lambda|k) \propto [\lambda^{\alpha-1} \exp[-\beta\lambda]] [\lambda^k \exp^{-\lambda}] = \lambda^{(\alpha+k)-1} \exp[-(\beta+1)\lambda]$$

Thus we see that the Poisson likelihood of $k|\lambda$ has Gamma as a conjugate prior for $\lambda|\alpha, \beta$, because the Posterior is Gamma distributed with the new parameters $\alpha' = k + \alpha$ and $\beta' = \beta + 1$.

$$\lambda \sim \text{Gamma}[\alpha', \beta']$$

$$p(\lambda) = \frac{\beta^{\alpha'}}{\Gamma(\alpha')} \lambda^{\alpha'-1} \exp[-\beta'\lambda]$$

Question 2

Show that the Gibbs transition operator satisfies the detailed balance equation and as such can be interpreted as an MH transition operator that always accepts.

Let's first clarify what the questions is asking for. In Gibbs there is

- the individual update of a component, based on the conditional for that component with the others frozen at what they are
- the update of all the components after all the components have been updated.

So this questions can be taken in two senses, and I will answer it as follows.

I will first show that each Gibbs update of the one component "satisfies details balance and as such can be interpreted as an MH transition operator that always accepts". Then I will consider what statements we can make on the update of the whole vector (all the components). There's an important subtlety there, but the nuances can get quite academic, and I leave off with some important links that further clarify the issue (which I did read and can proudly say I understand).

First, let's consider updating one component of a vector, as we do in Gibbs, and define some notation (following Detailed balance for Gibbs sampling in <http://www.gatsby.ucl.ac.uk/teaching/courses/mh-2015/lec12-2-handout.pdf>).

- $[x_i, x^*]$ is the transition from vector $x = (x_1, \dots, x_i, \dots, x_n) \rightarrow x^* = (x_1, \dots, x_i^*, \dots, x_n)$, where the i th component of x has been chosen.
- π_i refers to the probability of choosing to update component i . If we randomly choose i components, this would be constant. I'm not sure how to write this down if we cycled through in some fixed order.
- $p(x_i^*|x_{-i})$ is the conditional probability of updating the component x_i to x_i^* given x_{-i} , which is the states of all the components of x , but without the i th component.

$$T[x_i, x^*] = \pi_i p(x_i^*|x_{-i})$$

Detailed balance for a step in Gibbs is defined as, where $p(x)$ is the probability of being in state x .

$$\begin{aligned} T[x_i, x^*]p(x) &= T[x_i, x^*]p(x) \\ &= [\pi_i p(x_i^*|x_{-i})] [p(x_i|x_{-i})p(x_{-i})] \\ &= [\pi_i p(x_i|x_{-i})] [p(x_i|x_{-i})p(x_{-i})] \text{ swapping } x_{-i}^* \rightarrow x_{-i} \\ &= [\pi_i p(x_i|x_{-i})] [p(x_i^*|x_{-i})p(x_{-i})] \text{ rearranging conditional terms} \\ &= T[x_i^*, x]p(x^*) \end{aligned}$$

Note that nothing has been assumed about π_i here, and thus no matter how we choose to update the i th component (randomly, cycle through), details balance will hold for one Gibbs step.

Showing that details balance hold for a whole Gibbs sweep (where we update all the components) is a different matter. Actually it can be shown for the case when the choice of the Gibbs, and the Gibbs sampler for one update from $x_i = (x_1^1, \dots, x_i^1)$ to $x_{i+1} = (x_1^{i+1}, \dots, x_i^{i+1})$ is the average of random updates, instead of their product sum. In the latter case, there can be simple counter examples provided that show detail balance does not hold.

See

- section 12.4 in <https://statweb.stanford.edu/~owen/mc/Ch-MCMC.pdf>
- slide 3 of <http://galton.uchicago.edu/~eichler/stat24600/Handouts/059.pdf>

Now that we've established detailed balance, the next step of the Gibbs updates is to show that the MH acceptance operator is always one

$$A = \min \left(1, \frac{p(x^*)q(x|x^*)}{p(x)q(x^*|x)} \right)$$

p is the stationary distribution. We can show that this transition operator satisfies detailed balance:

$$\begin{aligned} p(x^*)T[x_i^*, x^*]p(x^*) &= p(x)T[x_i, x^*] \text{ (detailed balance)} \\ p(x^*)p(x_i|x^*) &= p(x_i)p(x^*|x_{-i}) \text{ (T is cond. prob in Gibbs)} \\ p(x^*, x_i) &= p(x_i, x^*) \end{aligned}$$

q is the proposal distribution. In Gibbs, the proposal distribution is the stationary distribution, i.e. $p(a|b) = \pi_i q(a|b)$, because the transition operator is the conditional of the stationary distribution (and the probability of choosing to update the i th component).

From detailed balance we have

$$\begin{aligned} \pi_i p(x_i^*|x_{-i})p(x^*) &= \pi_i p(x_i^*|x_{-i})p(x^*) \\ &\implies \frac{\pi_i p(x_i^*|x_{-i})}{p(x_i^*)p(x^*)} = 1 \end{aligned}$$

And since $p = q$ in Gibbs we have

$$A = \min(1, 1) = 1$$

Which means that the Gibbs transition operator in this case always accepts.

Question 3

Write code to compute the probability three ways that it is cloudy given that we observe that the grass is wet using this Bayes net model.

```
In [ ]:
import numpy as np

In [ ]:
##first define the probability distributions as defined in the exercise:

#define 0 as false, 1 as true
def p_C(c):
    #p = np.array([0.5,0.5])
    return p[c]

def p_S_given_C(s,c):
    #p = np.array([[0.5,0.9],[0.5,0.1]])
    return p[s,c]

def p_R_given_C(r,c):
    #p = np.array([[0.6,0.2],[0.2,0.8]])
    return p[r,c]

def p_W_given_S_R(w,s,r):
    #p = np.array([[
    [[0.0,0.1],[0.1,0.01]], #w = False
    [[0.0,0.9],[0.9,0.99]], #w = True
    ]])
    return p[w,s,r]
```

By enumerating all possible world states and conditioning by counting which proportion are cloudy given the observed world characteristic, namely, that the grass is wet.

The joint is the enumeration of all possible world states, and the proportion. It is possible to enumerate all these world states because the we can loop through the small number of discrete states (binary in this example) in this small graphical model (just four nodes in this case).

The complexity of enumerating a graphical model with n_{nodes} with node i having $n_{i,states}$ number of discrete states is $O(\prod_i n_{i,states})$. This is how many multiplications there are. Note that we can precompute these states and don't have to keep calling the functions to compute them.

We just need to

- condition by fixing cloudy as True
- count by taking the ratio of probabilities of the marginals $\frac{p(w=1,c=1)}{p(w=1)} = \frac{\sum_{s,r} p(w=1,s,r,c=1)}{\sum_{s,r} p(s,r,c=1)}$ by Bayes rule.

Once we have the joint, it is easy to get the marginal by just summing over specific axes.

This gives us an "exact" probability, since it uses the joint, which comes from the conditional probability tables. Other sampling procedures (MCMC/Gibbs) that use the conditionals will reach in the limit of many random iterations. This is a simple example, and many times the graphical model has too many nodes, or the nodes have too many states, the joint is not directly available.

```
In [ ]:
##1. enumeration and conditioning:

# compute joint:
p = np.zeros((2,2,2,2)) #c,s,r,w
for c in range(2):
    p_C_precomputed = p_C(c)
    for s in range(2):
        p_S_given_C_precomputed = p_S_given_C(s,c)
        for r in range(2):
            p_R_given_C_precomputed = p_R_given_C(r,c)
            for w in range(2):
                p[c,s,r,w] = p_C_precomputed*p_S_given_C_precomputed*p_R_given_C_precomputed*p_W_given_S_R(w,s,r)

# condition and marginalize:

w=1 # given the grass is wet
c=1 # query 'is it cloudy?'
p_C_given_W = p[c,:, :,w].sum() / p[:, :, :,w].sum() # p(c=1,s,r,w=1)/p(w=1)

print('There is a %.2f%% chance it is cloudy given the grass is wet'.format(p_C_given_W*100))
```

There is a 57.58% chance it is cloudy given the grass is wet

Using ancestral sampling and rejection

We can do ancestral sampling with rejection as follows

- Start at the root of the graph (node with no parents).
- Count the number of cloudy days (prob is proportion of cloudy vs all), while also
- rejecting if grass not wet (if the graph continued on we would stop).

We sample both the the R and S nodes inside the loop with C, because they depend on C, but in parallel with each other.

We sample c by setting it to $\text{Uni}([0,1] > \text{equivalent to } \text{np.random.choice}([0,1],p=[p_C[0],p_C[1]]))$

Here we don't fix wet, so there is a loss of efficiency in so far as the probability of being wet is rare. Quantitatively we reject ~35% of the samples. If it's very rare, then we reject a lot and we have to sample a lot.

We get the probability of cloudy given wet by accumulating the number of wet days and the number of cloudy given wet days. The probability of cloudy given wet, is the number of days that are cloudy given wet (`num_cloudy_given_grass_wet`) divided by the number of wet days (`num_wet`)

```
In [ ]:
num_rv_to_sample = 4 # c,s,r,w
num_total_days = 100000
c_idx, s_idx, r_idx, w_idx = range(num_rv_to_sample)
uni_rv = np.random.uniform(low=0,high=1,size=num_total_days*num_rv_to_sample).reshape(num_total_days,num_rv_to_sample)

num_cloudy_given_grass_wet=0
num_wet = 0
for day in range(num_total_days):
    c = int(uni_rv[day,c_idx] > p_C(0))
    s = int(uni_rv[day,s_idx] > p_S_given_C(0,c))
    r = int(uni_rv[day,r_idx] > p_R_given_C(0,c))
    w = int(uni_rv[day,w_idx] > p_W_given_S_R(0,s,r))

    num_wet += 1
    if c:
        num_cloudy_given_grass_wet += 1
    else:
        rejections += 1
#print('There is a cloudy %.2f%% of the time, given it is wet is %f' % (
print('The chance of it being cloudy given the grass is wet is %.2f%%'.format(100*num_cloudy_given_grass_wet /
print('c %.2f%% of the total samples were rejected'.format(100*rejections/num_total_days))

The chance of it being cloudy given the grass is wet is 57.64%
35.24% of the total samples were rejected
```

Using Gibbs sampling

Since we have the full joint, we can easily get any conditional we want by dividing the joint by the joint marginalized by variable we are conditioning on.

Although the starter code mentions some things about the zero in the conditional probability of $p(w|s=r=0)$, this is not an issue in the conditional $p(c|s,r)$ since we marginalize out w:

$$p(c|s,r) = \frac{p(c,s,r)}{p(s,r)} = \frac{\sum_w p(c,s,r,w)}{\sum_{c,r} p(c,s,r,w)}$$

Doing Gibbs with w=1 means that we always leave it fixed, and thus only spend time doing useful sampling that we won't throw away.

```
In [ ]:
##3. Gibbs
# we can use the joint above to condition on the variables, to create the needed
# conditional distributions:

#we can calculate p(R|C,S,W) and p(S|C,R,W) from the joint, dividing by the right marginal distribution
#findings in (p(R,S,W)) does not depend on C,
# but since p(W|S,R) does not depend on C,
#p(C | R,S) = p(R,S,C)/int(C) (p(R,S,C))

#first create p(R,S,C):
p_C_R = np.zeros((2,2,2)) #c,s,r
for c in range(2):
    for s in range(2):
        for r in range(2):
            p_C_R[c,s,r] = p_C(c)*p_S_given_C(s,c)*p_R_given_C(r,c)

#then create the conditional distribution:
p_C_R_starter_code = p_C_R[:, :,1]/p_C_R[:, :,1].sum(axis=0,keepdims=True)
p_C_given_W = p_C_R[:, :,1], keepdims=True) / p_C_given_W[:, :,1], keepdims=True)
assert np.allclose(p_C_given_W,R_starter_code,p_C_given_W[:, :,1],rtol=0)
```

In []:

Out []:

We run Gibbs for 100'000 samples, and use a burn in of 500, and only take every 10th sample.

At this point I don't know how reasonable these values are, except to compare the final answer with the 57.6% obtained previously.

```
In [ ]:
## Gibbs sampling
num_samples = 1000000
samples = np.zeros(num_samples)
state = np.zeros(4,dtype=int)

num_rv_to_sample=4
uni_rv = np.random.uniform(low=0,high=1,size=num_samples*num_rv_to_sample).reshape(num_samples,num_rv_to_sample)

c,s,r = np.random.choice([0,1],p=[0.5,0.5],size=3)
w=1

for sample in range(num_samples):
    c = int(uni_rv[sample,c_idx] > p_C_given_S_R(0,s,r))
    s = int(uni_rv[sample,s_idx] > p_S_given_C(0,c,w))
    r = int(uni_rv[sample,r_idx] > p_R_given_C(0,c,w))
    words[sample] = c

The chance of it being cloudy given the grass is wet is 58.26%
```

```
In [ ]:
burn_in = 10000
stride = 10
print('The chance of it being cloudy given the grass is wet is %.2f%%'.format(samples[burn_in::stride].mean())

The chance of it being cloudy given the grass is wet is 57.74%
```

Question 4

Consider the Bayesian linear regression model discussed in the lecture on graphical models

I will answer the question "backwards" since the results of part c are used in b and b in a.

c) produce the analytic form of the posterior predictive.

The posterior predictive is the probability of \hat{t} given all the data we know $\{(x_n, t_n), \sigma^2, \alpha, \hat{x}\}$. Notice that w is not in there. That is because \hat{t} is unobserved and we only have a distribution for what it is. Compared with typical least squares regression, in this question we don't use a point estimate for \hat{t} , but the whole probability distribution. This takes some care, because we have a multivariate Gaussian prior on w , and we have an updated posterior on w based on our observed data. Integrating over this will involve a multidimensional integral. Let's see how this works.

In the notation below, we drop α and σ , because they are always given

$$p(\hat{t}|\{x_n, t_n\}, \hat{x}) = \int dw p(\hat{t}|w, \hat{x})p(w|\{x_n, t_n\})$$

The posterior of w is proportional to the posterior of all the $\{t_n\}$ and the prior of w by Bayes rule. The evidence $p(t_n)$ is a constant in w and doesn't affect how we compute the integral over w .

$$p(w|\{x_n, t_n\}) \propto \prod_i [N[t_n|w^T x_n, \sigma^2]] N[w|0, \alpha]$$

We can furthermore show that the posterior of w is a multivariate Gaussian with mean $w_N = \frac{1}{\sigma^2} V_N X^T \hat{t}$ and covariance matrix $V_N = (\frac{X^T X}{\sigma^2} + \frac{1}{\alpha})^{-1}$, where $\hat{t} \in \mathbb{R}^n$ is the vector of all $t_n \in \mathbb{R}$ and $X \in \mathbb{R}^{n \times d}$ is the data matrix of all $x_n \in \mathbb{R}^d$. To show that this is the case, it suffices to show that argument in the exponent of the posterior of w can be written of the form

$$\begin{aligned} &-\frac{1}{2} (w - w_N)^T V_N^{-1} (w - w_N) \\ &= -\frac{1}{2} w^T V_N^{-1} w + w_N^T V_N^{-1} w - \frac{1}{2} w_N^T V_N^{-1} w_N \end{aligned}$$

Neglecting terms that won't affect the integral, and noting that $\sum_n x_n x_n^T = X^T X$ and $\sum_n t_n x_n = X^T \hat{t}$, we have

$$\begin{aligned} \log p(w|\{x_n, t_n\}) &= -\frac{1}{2\sigma^2} \sum_n (t_n - w^T x_n)^2 - \frac{1}{2} w^T \frac{1}{\alpha} w + \dots \\ &= -\frac{1}{2} w^T \left(\frac{X^T X}{\sigma^2} + \frac{1}{\alpha} \right) w + \frac{\hat{t}^T X w}{\sigma^2} + \dots \end{aligned}$$

Completing the square by matching the coefficients in the multidimensional case, we have $V_N^{-1} = \frac{X^T X}{\sigma^2} + \frac{1}{\alpha}$ as desired, and $w_N^T V_N^{-1} = \frac{\hat{t}^T X}{\sigma^2} \implies w_N = \frac{1}{\sigma^2} V_N X^T \hat{t}$ as desired (using the symmetry of V_N and its inverse).

Thus $p(w|\{x_n, t_n\}) = N[w|w_N, V_N]$

To get the posterior in \hat{t} , we have to integrate over w .

$$\int dw N[\hat{t}|w^T x, \sigma^2] N[w|w_N, V_N]$$

This is the form of Eq 7.60 in Murphy, 2012, and the integral is another Gaussian in \hat{t} . Doing this integral by completing the square involves resolving the order of a matrix. Exact details of this are given in Murphy, 2012 (4.3.4.1 inverse of a partitioned matrix using Schur complements; 4.3.4.2 The matrix inversion lemma) By Eq. 7.61 in Murphy, 2012, we have

$$\begin{aligned} p(\hat{t}|w, \hat{x}, \sigma, \{x_n, t_n\}) &= N[\hat{t}|w^T x_N, \sigma_N^2(x) = \sigma^2 + \hat{x}^T V_N \hat{x}] \\ &= \frac{1}{\sqrt{(2\pi\sigma_N^2)}} \exp[-\frac{(\hat{t} - w^T \hat{x})^2}{2\sigma_N^2}] \end{aligned}$$

Note that \hat{t} is centred around the Ridge regression least squares estimate for the weights, but its uncertainty is dependent on its own data point \hat{x} .

b) perform pure Gibbs on both w and t

NB: Instead of using w_i, w_{i+1} (and likewise for \hat{t}), I just use w and w' with the understanding that $\hat{x}, x, t, \sigma, \alpha$ do not change for MCMC steps.

To answer this question, we need to derive conditional probabilities for w and \hat{t} . This will allow us to update them in Gibbs. For the remaining variables in the graphical model, we have observed them and can use these values.

$$p(\hat{t}^+|w, \hat{x}, x, t, \sigma^2, \alpha)$$

$$p(w|\hat{t}^+, \hat{x}, x, t, \sigma^2, \alpha)$$

note the use of the new update for \hat{t} , denoted \hat{t}^+ in the conditional for w

For the conditional of \hat{t} we can drop x, \hat{x} , α since w and σ^2 are given. Then we just have the equation given for $p(t_n|w, x_n, \sigma^2)$ with $t_n \rightarrow \hat{t}$, and $x_n \rightarrow \hat{x}$, and w from the latest MCMC step.

$$p(\hat{t}^+|w, \hat{x}, x, t, \sigma^2, \alpha) = N[w^T \hat{x}, \sigma^2]$$

For the conditional of w , we do supply some information about it, and it will change depending on that MCMC step we are at. As we cycle through Gibbs, we will use our estimate of it at the appropriate time step, and so we want to know how to use it in the conditional of w :

We can use the result derived in part c), with the caveat that the $\{t_n, x_n\}_{N+1}$ data now contains the $\{\hat{t}^+, \hat{x}\}$ pair, indexed $N+1$, making sure to use the most recent samle of \hat{t}^+ from the last MCMC step.

Thus

- $p(w|\{x_n, t_n\}, \hat{t}^+, \hat{x}, \sigma^2, \alpha) = N[w|w_{N+1}, V_{N+1}]$
- with w_{N+1} and V_{N+1} as defined in c), but with X and \hat{t} including $\{\hat{t}^+, \hat{x}\}$, and thus changing every MCMC iteration.

a) perform MH within Gibbs on the blocks w and t

Doing Metropolis-Hastings within Gibbs means that we use MH to sample from the conditionals. In b, we derived the analytical forms, and thus we could use them to sample from. However, if we don't have a nice form, we can use MH. If we further more assume that the proposal q is symmetric, i.e. $q(x|x') = q(x'|x)$, then the MH acceptance probability is even simpler, but we don't have to assume this to do MH in Gibbs.

$$A(x'|x) = \min \left(1, \frac{p(x')q(x|x')}{p(x)q(x'|x)} \right)$$

The probability we are interested in estimating here, $p(x)$ is actually the conditional of \hat{t} or w , so we have two acceptance probabilities with

- $x = \hat{t}^+|w, \hat{x}, x, t, \sigma^2, \alpha$
- and $x' = \hat{t}^-|w, \hat{x}, x, t, \sigma^2, \alpha$

for estimating $p(\hat{t}^+|w, \hat{x}, x, t, \sigma^2, \alpha)$

and

- $x = w^+|\hat{t}^+, \hat{x}, x, t, \sigma^2, \alpha$
- and $x' = w^-|\hat{t}^+, \hat{x}, x, t, \sigma^2, \alpha$

for estimating $p(w^+|\hat{t}^+, \hat{x}, x, t, \sigma^2, \alpha)$

Note that since the conditional is proportional to the joint, up to a normalization constant (the evidence, that cancels out), we can use it in place of the conditional. This is very convenient when getting the evidence is difficult.

$$\begin{aligned} \frac{p(\hat{t}^+)}{p(\hat{t}^-)} &= \frac{p(\hat{t}^+|w, \hat{x}, x, t, \sigma^2, \alpha)}{p(\hat{t}^-|w, \hat{x}, x, t, \sigma^2, \alpha)} = \frac{p(\hat{t}^+|w, \hat{x}, x, t, \sigma^2, \alpha)p(w^+|\hat{t}^+, \hat{x}, x, t, \sigma^2, \alpha)}{p(\hat{t}^-|w, \hat{x}, x, t, \sigma^2, \alpha)p(w^-|\hat{t}^-, \hat{x}, x, t, \sigma^2, \alpha)} \\ &= \frac{p(\hat{t}^+|w, \hat{x}, x, t, \sigma^2, \alpha)}{p(\hat{t}^-|w, \hat{x}, x, t, \sigma^2, \alpha)} \end{aligned}$$

The joints facirize nicely, and everything cancels except for the 1D Gaussian of the new data point - even the Gaussian normalization constant cancels

$$\frac{p(x^*)}{p(x)} = \frac{\exp -\frac{(x^* - w^*)^2}{2\sigma^2}}{\exp -\frac{(x - w^*)^2}{2\sigma^2}}$$

Likewise, for updating w (taking care to use the new \hat{t}^+ we get

$$\begin{aligned} \frac{p(x^*)}{p(x)} &= \frac{p(\hat{t}^+|w^+, \hat{x}, x, t, \sigma^2, \alpha)}{p(\hat{t}^-|w^-, \hat{x}, x, t, \sigma^2, \alpha)} = \frac{p(\hat{t}^+|w^+, \hat{x}, \sigma^2)p(w^+|\alpha)}{p(\hat{t}^-|w^-, \hat{x}, \sigma^2)p(w^-|\alpha)} \\ &= \frac{\exp -\frac{(\hat{t}^+ - w^*)^2}{2\sigma^2}}{\exp -\frac{(\hat{t}^- - w^*)^2}{2\sigma^2}} \frac{w^* \alpha}{2\alpha} \end{aligned}$$

Question 5

Implement a sampler for the LDA model on a corpus consisting of abstracts from NeurIPS in years past. You may wish to start from support code that will be provided on the course website. Use the output to answer questions about the NIPS abstract corpus. The data for this model are in `bagofwords_nips.mat`, words, `nips.mat`, and title, `nips.mat`. When you load these files in the variables DS, WS, WO, and titles will appear in your workspace. The variable WS is the entire corpus vectorized into a long row of words encoded as a row vector of integers. Each integer represents a word, WO is the dictionary; to look up the word corresponding to an integer use WO[i]. The variable DS is a row vector of the same length as WS which indicates from which document each word comes. If the document number is i then the value of that document can be found using `DS[i]`. Train LDA with 20 topics using an MCMC sampler in the collapsed representation (i.e. with all β 's and θ 's integrated out analytically) and use the single sampler with


```
In [69]: def sample_topic_assignment(topic_assignment,
                             topic_counts,
                             doc_counts,
                             topic_N,
                             doc_N,
                             alpha,
                             gamma,
                             document_assignment,
                             sampled_word_idx,
                             topic_integers,
                             n_topics,
                             n_topics,
                             alpha_size):
    """
    Sample the topic assignment for each word in the corpus, one at a time.

    Args:
        topic_assignment: size n array of topic assignments
        topic_counts: n_topics x alphabet_size array of counts per topic of unique words
        doc_counts: n_docs x n_topics array of counts per document of unique topics
        topic_N: array of size n_topics count of total words assigned to each topic
        doc_N: array of size n_docs count of total words in each document, minus 1

        alpha: prior dirichlet parameter on document specific distributions over topics
        gamma: prior dirichlet parameter on topic specific distributions over words.

        words: size n array of words
        document_assignment: size n array of assignments of words to documents

    Returns:
        topic_assignment: updated topic assignment array
        topic_counts: updated topic counts array
        doc_counts: updated doc_counts array
        topic_N: updated count of words assigned to each topic
    """
    # sampled_word_idx = 0
    document_idx = document_assignment[sampled_word_idx]
    word_idx = words[sampled_word_idx]
    old_topic_idx = topic_assignment[sampled_word_idx]

    # remove sampled from counts
    # TODO check how sparse and if sparse encoding helps
    topic_counts[old_topic_idx,word_idx] -= 1
    doc_counts[document_idx,old_topic_idx] -= 1
    topic_N[old_topic_idx] -= 1

    # using notation from
    # https://medium.com/analytics-vizhya/topic-modeling-using-lda-and-gibbs-sampling-explained-49d49b3d10d5
    # n = doc_counts
    # v = topic_counts
    # lam = gamma
    # p = topic_assignment
    # n_alpha = n*alpha
    # v_lam = v*lam

    # unvectorized
    # p = (n_alpha[document_idx,k])*(v_lam[k,w]) / (n_alpha[document_idx].sum() * v_lam[k].sum())
    # vectorize over (k,w)
    # p = n_alpha[document_idx]*v_lam[:,word_idx] / (n_alpha[document_idx].sum()*v_lam.sum())

    # Eq 27.37 in p. 956 of Murphy, 2012
    K = n_topics
    V = alphabet_size
    C_Wk = topic_counts
    C_kk = doc_counts
    c_k = topic_N
    l_w = doc_N

    p = (c_k[document_idx] + alpha)*(c_wk[:,word_idx] + gamma) / ((c_k + V*gamma)*(c_w[document_idx]+K*alpha))

    z_sampled = np.random.choice(topic_integers,size=1,p=p/p.sum())
    new_topic_idx = z_sampled

    # update
    topic_counts[new_topic_idx,word_idx] += 1
    doc_counts[document_idx,new_topic_idx] += 1
    topic_assignment[sampled_word_idx] = new_topic_idx
    #topic_N = topic_counts.sum(axis=1)
    gamma = doc_counts
    topic_N[new_topic_idx] += 1

    return topic_assignment, topic_counts, doc_counts, topic_N
```

```
In [70]: #prior parameters, alpha parameterizes the dirichlet to regularize the
#document specific distributions over topics and gamma parameterizes the
#dirichlet to regularize the topic specific distributions over words.
#These parameters are both scalars and really we use alpha * ones() to
#parameterize each dirichlet distribution. Iters will set the number of
#times your sampler will iterate.
alpha = 0.1
gamma = 0.001
iters = 150

jll = []

ll_max = joint_log_likelihood(doc_counts,topic_counts,alpha,gamma)
jll.append(ll_max)

n_words = words.size
topic_integers = np.arange(20)
iter_ll = []

K = n_topics
M = n_docs
K,V = topic_counts.shape # (20, 13649) (K,V) topic_counts(k,v) is number of times topic k uses word v
K,M = doc_counts.shape # (100, 20) (M,K) doc_counts(m,k) is number of times document m uses topic k
assert K == V

%%timeit
sampled_word_idx=0
sample_topic_assignment(
    topic_assignment,
    topic_counts,
    doc_counts,
    topic_N,
    doc_N,
    alpha,
    gamma,
    words,
    document_assignment,
    sampled_word_idx,
    topic_integers,
    K,
    V)
```

The slowest run took 10.13 times longer than the fastest. This could mean that an intermediate result is being cached.

10000 loops, best of 5: 72.4 µs per loop

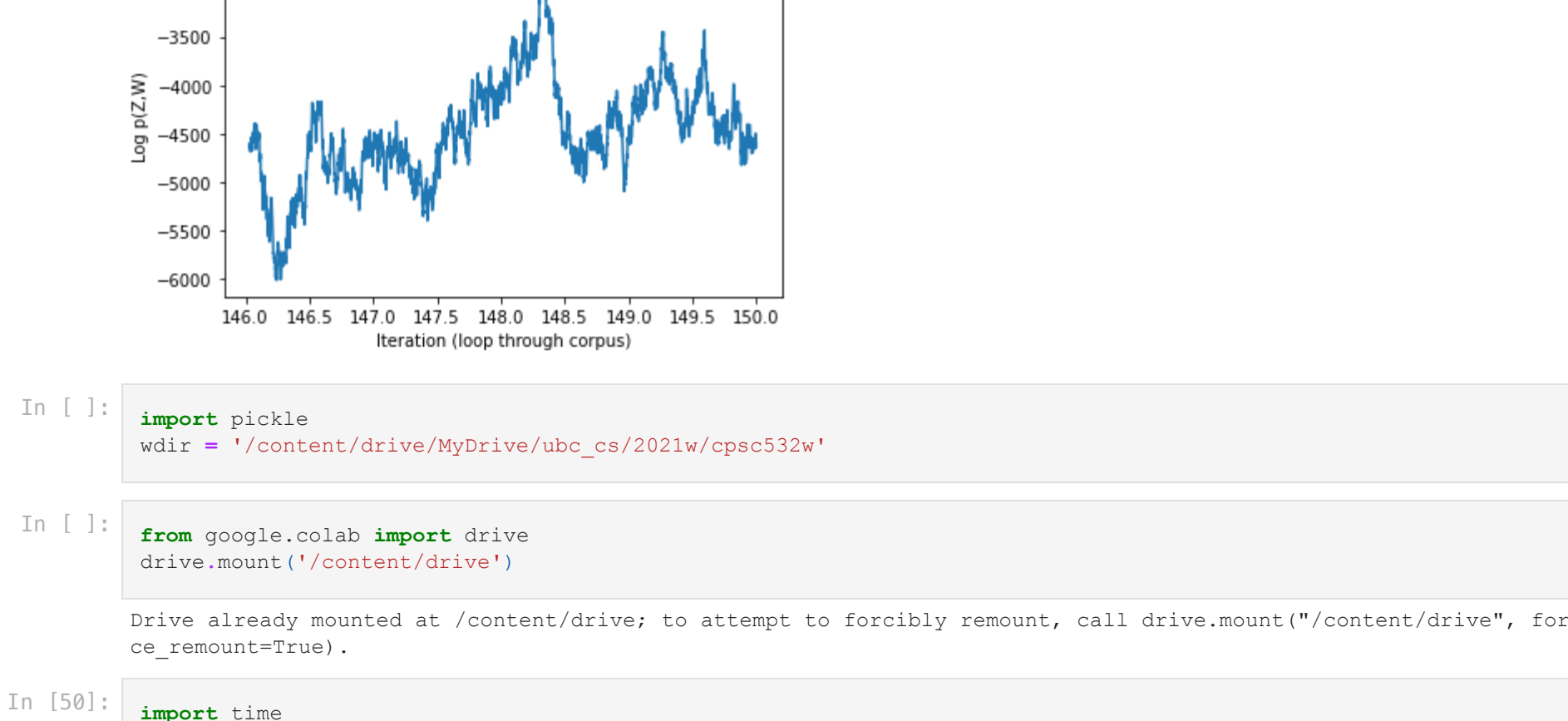
```
In [71]: for iter in range(iters):
    #print('iteration:',iter)
    sampled_word_idx = np.random.permutation(words.size)
    for i,sampled_word_idx in enumerate(sampled_word_idxs):
        if i % (n_words/100) == 0: print(i)
        prm = np.random.permutation(words.shape[0])
        # words = words[prm]
        # document_assignment = document_assignment[prm]
        # topic_assignment = topic_assignment[prm]

        #sampled_word_idx = np.random.choice(n_words)

        topic_assignment, topic_counts, doc_counts, topic_N = sample_topic_assignment(
            topic_assignment,
            topic_counts,
            doc_counts,
            topic_N,
            doc_N,
            alpha,
            gamma,
            words,
            document_assignment,
            sampled_word_idx,
            topic_integers,
            K,
            V)

        if i % (n_words/1000) == 0:
            #print(f'end=')
            ll = joint_log_likelihood(doc_counts,topic_counts,alpha,gamma)
            jll.append(ll)
            iter_ll.append(iter)
            iter_ll.append(iter)
            if ll > ll_max:
                ll_max = ll
                topic_counts_max = topic_counts
            plt.plot(jll)
```

Out[69]: Text(0.5, 1.0, 'Joint distribution')



```
In [69]: sr = pd.Series(jll)
sr.index = np.linspace(0,np.max(iter_ll),len(jll))
sr.iloc[-4000:].plot()
plt.ylabel('Log p(W)')
plt.xlabel('Iteration (loop through corpus)')
plt.title('Joint Distribution (at end of run showing equilibrium)\n')
```

Out[69]: Text(0.5, 1.0, 'Joint distribution (at end of run showing equilibrium)\n')



```
In [70]: import pickle
wdir = '/content/drive/MyDrive/ubc_cs/2021w/cpsc532w'
```

```
In [71]: from google.colab import drive
drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

```
In [50]: import time
fname = ''.join([wdir, '/hw1-q5-',time.strftime('%Y%m%d-%H%M'), '.p'])
```

```
data_d = {'topic_assignment':topic_assignment,
'topic_counts':topic_counts,
'doc_counts':doc_counts,
'topic_N':topic_N,
'alpha':alpha,
'gamma':gamma,
'n_topics':n_topics,
'alphabet_size':alphabet_size,
'n_words':n_words,
'li_max':li_max,
'n_docs':n_docs,
'jll':jll,
'iter_iter_ll,
}

pickle.dump(data_d, open( fname, "wb" ) )
```

```
In [75]: !ls /content/drive/MyDrive/ubc_cs/2021w/cpsc532w
```

```
hw1-q5-20210928-1708.p hw1-q5-20210930-1940.p hw1-q5-20211001-0233.p
hw1-q5-20210929-0530.p hw1-q5-20210930-2348.p hw1-q5-20211001-0525.p
```

```
In [76]: # fname_load = '/content/drive/MyDrive/ubc_cs/2021w/cpsc532w/hw1-q5-20210929-0530.p'
# load_d = pickle.load( open( fname_load, "rb" ) )
```

```
# alpha = load_d['alpha']
# topic_assignment=load_d['topic_assignment']
# topic_counts=load_d['topic_counts']
# doc_counts=load_d['doc_counts']
# topic_N=load_d['topic_N']
# alpha=load_d['alpha']
# gamma=load_d['gamma']
# n_topics=load_d['n_topics']
# alphabet_size=load_d['alphabet_size']
# n_words=load_d['n_words']
# li_max=load_d['li_max']
# n_docs=load_d['n_docs']
# jll=load_d['jll']
```

What are the top ten most probable words for each of the 20 topics?

```
In [70]: ## Find the 10 most probable words of the 20 topics
ind = np.argsort(topic_counts_max, axis=1) # note that this is from best to most probable, and we should look at
```

```
In [71]: n_prob_words = 10
most_popular_words = words_d['W0'][:,ind[:,n_prob_words]].flatten().reshape(n_topics,n_prob_words)
print('most popular words')
```

```
for row in range(most_popular_words.shape[0]):
    print('topic %i'%row,end=" ")
    for col in range(most_popular_words.shape[1]):
        print(most_popular_words[row,col][0],end=" ")
    print()
```

most popular words
topic 0 face pixel objects figure feature features recognition object images image
topic 1 linear distribution mixture space matrix models algorithm gaussian model data
topic 2 fixed model recurrent states neural phase dynamics system state time
topic 3 function functions probability class bound error size examples set number
topic 4 analysis order frequency filter processing signals time noise information signal
topic 5 connectionist graph level representations rule representation set tree rules structure
topic 6 vector classifier error class texts performance classification data set training
topic 7 direction eye cells orientation cortex response stimulus motion model visual
topic 8 approximation matrix gradient case problem functions algorithms vector linear function
topic 9 via input neuron voltage output current figure ship circuit analog
topic 10 number output net nodes layer memory input networks neural network
topic 11 neural log variables prior bayesian parameters models distribution probability model
topic 12 state model input spike firing cells synaptic neuron cell neurons
topic 13 unit weights networks output error hidden input training units network
topic 14 regions target location space region fig position fields map field
topic 15 control time action function optimal states policy reinforcement action state
topic 16 words time hmm training context sequence system word recognition speech
topic 17 search learned number learn problem task algorithms time algorithm learning
topic 18 figure systems controller feedback forward neural motor control system model
topic 19 kernel prediction results regression estimate approach based methods method data

One can consider the distribution over topics as a low dimensional representation of a document. We can use the dot product between topic distributions for two documents as a similarity metric.

What are the ten most similar documents to the first document, "Connectivity Versus Entropy"?

- Normalize counts to distribution (sum to 1). Use cosine similarity to handle confounding length/size.

```
In [72]: query_idx = 0
topical_dist = doc_counts/doc_counts.sum(1).reshape(-1,1)
def norm(arr,axis=None):
    if axis is None:
        return np.sqrt((arr*arr).sum())
    else:
        return np.sqrt((arr*arr).sum(axis))
similarity = topical_dist.dot(topical_dist[query_idx,:]) / (norm(topical_dist,1) * norm(topical_dist[query_idx,:]))
```

```
In [73]: n_sim=10
similar_idx = np.argsort(similarity)[-n_sim:][::-1]
for idx,close_topic in enumerate(similar_idx):
    print(idx,close_topic)
```

```
0 Connectivity Versus Entropy
1 LEARNING STOCHASTIC PERCEPTRONS UNDER K-BLOCKING DISTRIBUTIONS
2 A Method for learning from Bits
3 The Devil and the Network
4 Strong Unimodality and Exact Learning of Constant Depth g-Perceptron Networks
5 On Neural Networks with Minimal Weights
6 Estimating Average-Case Learning Curves Using Bayesian, Statistical Physics and VC Dimension Methods
7 The VC-Dimension versus the Statistical Capacity of Multilayer Networks
8 The Perceptron Algorithm Is Fast for Non-Malicious Distributions
9 Neural Computing with Small Weights
```