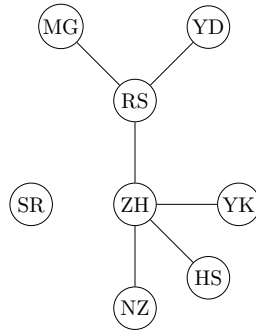


## Modeling Influence in a Network

This problem considers modeling the flow of influence within a social network using an undirected graphical model. We assume there is a new piece of information currently gaining influence in the network. In particular, for this example we consider a viral GIF that is being shared among the CS281 teaching staff. For each person in the network, we will associate a binary random variable indicating whether they are aware of the information, e.g.  $RS=1$  indicates that Rachit Singh has seen the GIF.



We consider the small fixed network of (relatively unpopular) individuals with a set of connections that form a forest.



With each person in the network we associate a unary log-potential e.g.  $\theta_{RS}(0)$  and  $\theta_{RS}(1)$  indicating a “score” for them seeing the GIF on their own.

s	SR	YD	MG	ZH	HS	RS	NZ	YK
$\theta_s(1)$	2	-2	-2	-8	-2	3	-2	1
$\theta_s(0)$	0	0	0	0	0	0	0	0

Additionally we assume a set of attractive binary log-potentials  $\theta_{s-t}(1,1) = 2$  for all connected nodes  $s, t \in E$  with all other values  $\theta_{s-t}(1,0) = 0$ ,  $\theta_{s-t}(0,1) = 0$ ,  $\theta_{s-t}(0,0) = 0$  across the network.

**Problem 1** (30pts)

(Note: for this problem you will be deriving several properties of undirected graphical models and implementing these properties in Python/PyTorch. Before getting started it is worth getting familiar with the `itertools` package in Python and also being sure you understand how PyTorch variables work. In particular try writing a simple scalar function with several vector-valued Variables as inputs and making you understand what the function `.backward` does in practice.)

- (a) Implement a function for computing the global score for an assignment of values to the random variables, i.e. a function that takes in  $\{0, 1\}$  values for each person and returns the *unnormalized* value without the  $A(\theta)$  term. What is the score for an assignment where RS and SR are the *only* people who have seen the GIF?
- (b) Implement a function for computing the log-partition function  $A(\theta)$  using *brute-force*, i.e. take in a vector  $\theta$  and return a scalar value calculated by enumeration. What is its value? Using this value compute the *probability* of the assignment in (a)?

For the problems below we will be interested in computing the marginal probabilities for each random variable, i.e.  $p(\text{SR} = 1)$ ,  $p(\text{RS} = 1)$ , etc.

- (c) First, implement a function to compute  $p(\text{RS} = 1)$  using brute-force marginalization. Use the functions above to enumerate and sum all assignments with  $\text{RS} = 1$ .
- (d) Using what we know about exponential families, derive an expression for computing all marginals directly from the log-partition function  $A(\theta)$ . Simplify the expression.
- (e) Combine parts (b) and (d) to implement a method for computing the marginal probabilities of this model using PyTorch and autograd. Ensure that this gives the same value as part (c).
- (f) Finally compute all the marginals using exact belief propagation with the serial dynamic programming protocol, as described in class and in Murphy section 20.2.1. Use MG as the root of your graph. Verify that your marginals are the same as in the previous two calculations.
- (g) How do the relative values of the marginal probabilities at each node differ from the original log-potentials? Explain which nodes values differ most and why this occurs.
- (h) (Optional) Implement the parallel protocol for exact BP. How does its practical speed compare to serial?

(a) Score = 148.413

(b) Score = 3423.372

$$P(1, 0, 0, 0, 0, 1, 0, 0) = 0.043353 \text{ (ordering as in table on previous page)}$$

(c)  $P(\text{RS} = 1) = 0.9848$

- (d) Let  $\mathbf{X}$  be the random indicator vector for whether or not the individuals in the network saw the GIF. Then  $\mathbf{x}$  is an instance of  $\mathbf{X}$ . Let  $\boldsymbol{\theta}$  be a vector containing  $\theta_s(1) \forall s$  and  $\theta_{s-t}(1, 1) \forall s, t \in E$ . Then we have:

$$p(\mathbf{x}) = \exp \left( \sum_s \theta_s(1) \cdot x_s + \sum_{s,t \in E} \theta_{s-t}(1, 1) \cdot x_s \cdot x_t - A(\boldsymbol{\theta}) \right) \text{ where:}$$

$$\begin{aligned}
A(\boldsymbol{\theta}) &= \log \left( \sum_{\mathbf{x}} \exp \left( \sum_s \theta_s(1) \cdot x_s + \sum_{s,t \in E} \theta_{s-t}(1,1) \cdot x_s \cdot x_t \right) \right) \\
\frac{dA(\boldsymbol{\theta})}{d\theta_c(1)} &= \frac{\sum_{\mathbf{x}} x_c \cdot \exp \left( \sum_s \theta_s(1) \cdot x_s + \sum_{s,t \in E} \theta_{s-t}(1,1) \cdot x_s \cdot x_t \right)}{\sum_{\mathbf{x}} \exp \left( \sum_s \theta_s(1) \cdot x_s + \sum_{s,t \in E} \theta_{s-t}(1,1) \cdot x_s \cdot x_t \right)} \\
&= P(X_c = 1) \quad \forall c
\end{aligned}$$

(e)  $E(\mathbf{X}) = [0.8808, 0.4942, 0.4942, 0.0412, 0.1349, 0.9848, 0.1349, 0.7402]^T$

We see that  $P(X_6 = 1) = P(RS = 1)$  is the same here using autograd as it is in part (c).

- (f) The serial dynamic programming protocol yields the same results as in the previous two calculations. This can be seen in my code.
- (g) There is a positive relationship between the marginal probabilities and the unary log-potentials. Also, holding unary log-potentials constant, those sharing an edge with someone who has a higher unary log-potential have higher marginal probabilities than those sharing an edge with someone who has a lower unary log-potential.

The marginal probabilities of RS=1 and ZH=1 differ the most. This is because the difference in their unary log-potentials is the highest. This difference in the unary log-potentials means that a scenario where RS sees the GIF on her own is much more likely than a scenario where ZH sees the GIF on her own. Being well connected increases the marginal probabilities in general and ZH is well connected, but ZH's unary log-potential is just so low that the difference in the marginals for RS=1 and ZH=1 is still the greatest.

## A Note on Sparse Lookups

For the next two problems it will be beneficial to utilize sparse matrices for computational speed. In particular our input data will consist of sparse indicator vectors that act as lookups into a dense weight matrix. For instance, if we say there are 500 users each associated with a vector  $\mathbb{R}^{100}$  it implies a matrix  $W \in \mathbb{R}^{500 \times 100}$ . If we want a fast sparse lookup of the 1st and 10th user we can run the following code:

```
W = nn.Embedding(500, 100)
x = Variable(torch.LongTensor([ [1], [10] ]))
print(W(x))
```

This same trick can be used to greatly speed up the calculation of bag-of-words features from the last homework. Let's use the same example:

- We like programming. We like food.

Assume that the vocabulary is

```
["We", "like", "programming", "food", "CS281"]
```

In last homework, we converted the above word id sequence to a vector of length of vocab size:

- [2, 2, 1, 1, 0]

Instead we can convert it vector of sentence length (often much shorter):

- [0, 1, 2, 0, 1, 3]

In order to calculate the same dot product between  $\mathbf{x}$  and a weight vector  $\mathbf{w} = [0.5, 0.2, 0.3, 0.1, 0.5]$  we can do sparse lookups and a sum: (verify yourself that it is correct):

```
vocab_size = 4
W = nn.Embedding(vocab_size, 1, padding_idx=text_field.stoi('<pad>'))
W.weight.data = torch.Tensor([ [0.5], [0.2], [0.3], [0.1], [0.5] ])
x = Variable(torch.LongTensor([ [0, 1, 2, 0, 1, 3] ]))
result = torch.sum(W(x), dim=1)).squeeze()
```

This code computes the same dot-product as in HW2 but by doing 6 fast vector lookups into  $w$  and summing them together. (The `padding_idx` part ensures the embedding code works when there are sentences of different length by setting a lookup of its argument to return 0. This lets you use a rectangular matrix even with batches of different sentence lengths.) For more details, see the documentation for this module on the PyTorch website.

**Problem 2** (Modeling users and jokes with a latent linear model, 25pts)

In this problem, we'll use a latent linear model to jointly model the ratings users assign to jokes. The data set we'll use is a modified and preprocessed variant of the Jester data set (version 2) from <http://eigentaste.berkeley.edu/dataset/>. The data we provide you with are user ratings of 150 jokes. There are over 1.7M ratings with values 1, 2, 3, 4 and 5, from about seventy thousand users. **The data we give you is a modified version of the original Jester data set, see the README, please use the files we provide and not the original ones.** The texts of the jokes are also available. Warning: most of the jokes are bad, tasteless, or both. At best, they were funny to late-night TV hosts in 1999-2000. Note also that some of the jokes do not have any ratings and so can be ignored.

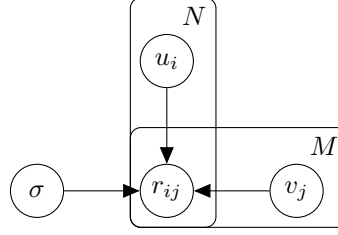
- (a) Let  $r_{i,j} \in \{1, 2, 3, 4, 5\}$  be the rating of user  $i$  on joke  $j$ . A latent linear model introduces a vector  $u_i \in \mathbb{R}^K$  for each user and a vector  $v_j \in \mathbb{R}^K$  for each joke. Then, each rating is modeled as a noisy version of the appropriate inner product. Specifically,

$$r_{i,j} \sim \mathcal{N}(u_i^T v_j, \sigma^2).$$

Draw a directed graphical model with plate diagrams representing this model assuming  $u_i$  and  $v_j$  are random vectors.

- (b) Derive the log-likelihood for this model and the gradients of the log-likelihood with respect to  $u_i$  and  $v_j$ .
- (c) Implement the log-likelihood calculation using PyTorch. Compute the gradients using autograd and confirm that is equal to the manual calculation above. (Hint: Read the documentation for `nn.Embedding` to implement  $u$  and  $v$ ).
- (d) Now set  $K = 2$ ,  $\sigma^2 = 1.0$  and run stochastic gradient descent (`optim.SGD`). We have provided a file `utils.py` to load the data and split it into training, validation and test sets. Note that the maximum likelihood estimate of  $\sigma^2$  is just the mean squared error of your predictions on the training set. Report your MLE of  $\sigma^2$ .
- (e) Evaluate different choices of  $K$  on the provided validation set. Evaluate  $K = 1, \dots, 10$  and produce a plot that shows the root-mean-squared error on both the training set and the validation set for each trial and for each  $K$ . What seems like a good value for  $K$ ?
- (f) We might imagine that some jokes are just better or worse than others. We might also imagine that some users tend to have higher or lower means in their ratings. In this case, we can introduce biases into the model so that  $r_{ij} \approx u_i^T v_j + a_i + b_j + g$ , where  $a_i$ ,  $b_j$  and  $g$  are user, joke and global biases, respectively. Change the model to incorporate these biases and fit it again with  $K = 2$ , learning these new biases as well. Write down the likelihood that you are optimizing. One side-effect is that you should be able to rank the jokes from best to worst. What are the best and worst jokes and their respective biases? What is the value of the global bias?
- (g) Sometimes we have users or jokes that only have a few ratings. We don't want to overfit with these and so we might want to put priors on them. What are reasonable priors for the latent features and the biases? Modify the above directed graphical model that shows all of these variables and their relationships. Note that you are not required to code this new model up, just discuss reasonable priors and write the graphical model.

- a. Let there be  $N$  users and  $M$  jokes.



- b. Let  $R$  be a random vector with the ratings for the user/joke pairs as the elements. Then  $R \in \mathbb{R}^{NM}$  and  $r$  is an instance of  $R$ .

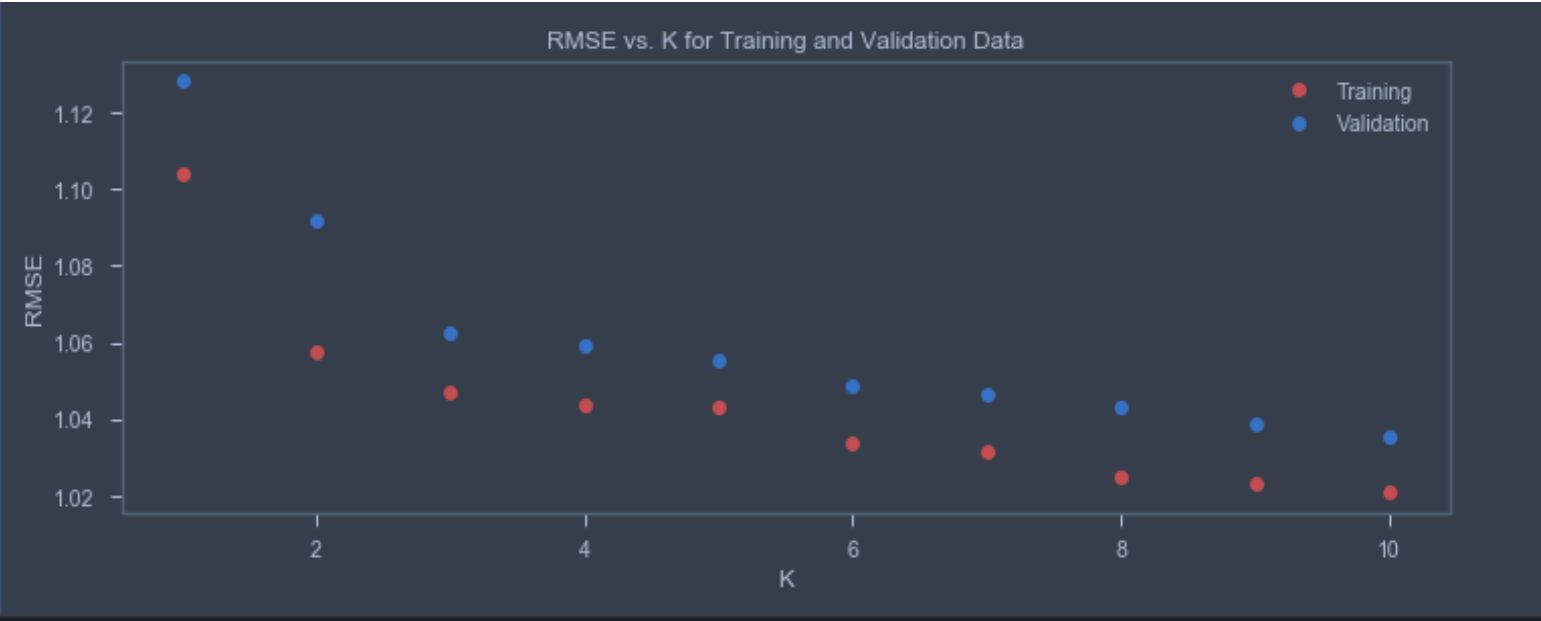
$$p(r) = \prod_{i=1}^N \prod_{j=1}^M \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(r_{ij} - u_i^T v_j)^2\right)$$

$$\Rightarrow \log(p(r)) = -\log(\sqrt{2\pi\sigma^2}) - \frac{1}{2\sigma^2} \sum_{i=1}^N \sum_{j=1}^M (r_{ij} - u_i^T v_j)^2$$

$$\frac{\partial \log(p(r))}{\partial u_i} = \frac{1}{\sigma^2} \sum_{j=1}^M (r_{ij} - u_i^T v_j) v_j$$

$$\frac{\partial \log(p(r))}{\partial v_j} = \frac{1}{\sigma^2} \sum_{i=1}^N (r_{ij} - u_i^T v_j) u_i$$

- c. I verified that autograd and manual computations produce the same results for the values of the gradients using a toy example. This can be seen in my code.
- d. MSE of predictions on training set =  $\hat{\sigma}_{MLE}^2 = 1.1184$
- e.  $K = 10$  is the best choice based on the validation RMSEs (see plot on next page).



f. Now we have  $r_{ij} \sim \mathcal{N}(u_i^T v_j + a_i + b_j + g, \sigma^2)$ .

$$\text{Then we get } \log(p(r)) = -\log(\sqrt{2\pi\sigma^2}) - \frac{1}{2\sigma^2} \sum_{i=1}^N \sum_{j=1}^M (r_{ij} - (u_i^T v_j + a_i + b_j + g))^2$$

Global bias: 3.32949661664

Worst jokes:

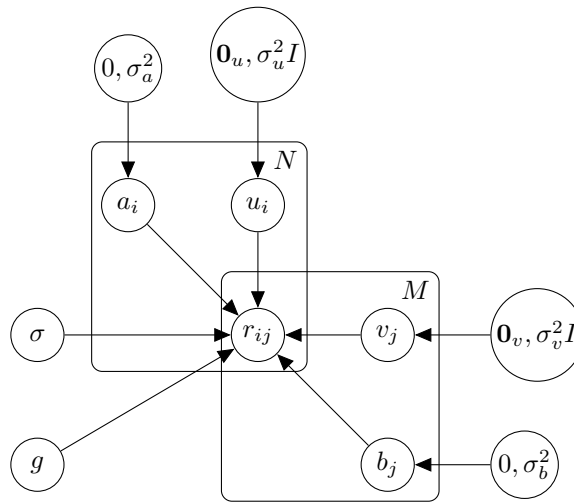
Joke ID	Bias
41	-0.629
16	-0.616
64	-0.534
44	-0.526
50	-0.488

Best jokes:

Joke ID	Bias
89	0.352
53	0.335
116	0.334
91	0.327
29	0.327

- g.  $u_i \sim \mathcal{N}(\mathbf{0}_u, \sigma_u^2 I)$  with  $\mathbf{0}_u \in \mathbb{R}^K$  and  $\sigma_u^2 I \in \mathbb{R}^{K \times K} \forall i \in \{1, \dots, N\}$   
 $v_j \sim \mathcal{N}(\mathbf{0}_v, \sigma_v^2 I)$  with  $\mathbf{0}_v \in \mathbb{R}^K$  and  $\sigma_v^2 I \in \mathbb{R}^{K \times K} \forall j \in \{1, \dots, M\}$   
 $a_i \sim \mathcal{N}(0, \sigma_a^2) \forall i \in \{1, \dots, N\}$   
 $b_j \sim \mathcal{N}(0, \sigma_b^2) \forall j \in \{1, \dots, M\}$

In this model, we can think of  $g$  as being our center for ratings across all users and all jokes. Since we do not have any prior information on the users or jokes, it makes sense to choose zero-centered Gaussians for our priors. This will keep ratings close to  $g$  in the absence of evidence of some alternative and prevent overfitting to the training set. Note that we can choose variances that are relatively small (especially for  $a_i$  and  $b_j$ ) since ratings fall between 1 and 5, inclusive. This will also help prevent overfitting.





## Ordinal Regression

We now address the problem of predicting joke ratings given the text of the joke. The previous models assumed that the ratings were continuous real numbers, while they are actually integers from 1 to 5. To take this into account, we will use an ordinal regression model. Let the rating values be  $r = 1, \dots, R$ . In the ordinal regression model the real line is partitioned into  $R$  contiguous intervals with boundaries

$$b_1 < b_2 < \dots < b_{R+1} = +\infty, \quad (1)$$

such that the interval  $[b_r, b_{r+1})$  corresponds to the  $r$ -th rating value. We will assume that  $b_1 = -4$ ,  $b_2 = -2$ ,  $b_3 = 0$ ,  $b_4 = 2$  and  $b_5 = 4$ . Instead of directly modeling the ratings, we will be modeling them in terms of a hidden variable  $f$ . We have that the rating  $y$  is observed if  $f$  falls in the interval for that rating. The conditional probability of  $y$  given  $f$  is then

$$p(y = r \mid f) = \begin{cases} 1 & \text{if } b_r \leq f < b_{r+1} \\ 0 & \text{otherwise} \end{cases} = \Theta(f - b_r) - \Theta(f - b_{r+1}), \quad (2)$$

where  $\Theta(x)$  is the Heaviside step function, that is,  $\Theta(x) = 1$  if  $x > 0$  and  $\Theta(x) = 0$  otherwise.

Notably there are many possible values of  $f$  that can lead to the same rating. This uncertainty about the exact value of  $f$  can be modeled by adding *additive Gaussian noise* to a noise free prediction. Let  $\sigma^2$  be the variance of this noise. Then  $p(f \mid h) = \mathcal{N}(f \mid h, \sigma^2)$ , where  $h$  is a new latent variable that is the noise free version of  $f$ .

Given some features  $\mathbf{x}$  for a particular joke, we can then combine the previous likelihood with a linear model to make predictions for the possible rating values for the joke. In particular, we can assume that the noise-free rating value  $h$  is a linear function of  $\mathbf{x}$ , that is  $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ , where  $\mathbf{w}$  is a vector of regression coefficients. We will assume that the prior for  $\mathbf{w}$  is  $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ .

*Computational Note for this Problem:* You may need the following numerical trick. If  $a > b$  but  $a - b$  is close to 0, you can compute  $\log(a - b)$  instead as:

$$\log(1 - \exp(\log b - \log a)) + \log a$$

Because  $a$  is always larger than  $b$  we have that  $\exp(\log a - \log b)$  is always smaller than 1 and larger than 0. Therefore,  $\log(1 - \exp(b - a))$  is always well defined.

An alternative way to compute  $\log(a - b)$  is by computing  $\log((1 - b) - (1 - a))$ :

$$\log(1 - \exp(\log(1 - a) - \log(1 - b))) + \log(1 - b)$$

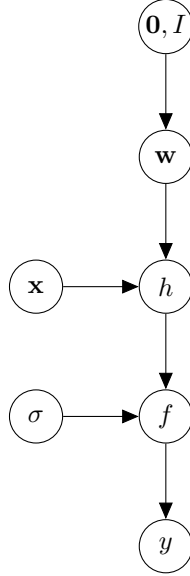
If  $a \rightarrow b$ ,  $b \rightarrow 0$ , then the first method is more numerically stable; if  $a \rightarrow b$ ,  $b \rightarrow 1$ , then the second method is more numerically stable. Therefore, a more numerically stable way of calculating  $\log(a - b)$  is:

$$\text{torch.max}(\log(1 - \exp(\log b - \log a)) + \log a, \log(1 - \exp(\log(1 - a) - \log(1 - b))) + \log(1 - b))$$

**Problem 3** (Ordinal linear regression 25pts)

1. Draw the directed graphical model for applying this model to one data point.
2. Compute the form of  $p(y | h)$  in the ordinal regression model. Explain how this would differ from a model with no additive noise term.
3. Give the equation for the mean of the predictive distribution in the ordinal regression model. How would is this term affected by  $\sigma^2$  (include a diagram).
4. Implement a function to compute the log-posterior distribution of the ordinal regression model up to a normalization constant. Use autograd to compute the gradients of the previous function with respect to  $\mathbf{w}$  and  $\sigma^2$ . Finds the MAP solution for  $\mathbf{w}$  and  $\sigma^2$  in the ordinal regression model given the available training data using SGD. Report the average RMSE on the provided test set. We have provided some helper functions in `utils.py`.
5. Modify the previous model to have a Gaussian likelihood, that is,  $p(y = r | h) = \mathcal{N}(r | h, \sigma^2)$ . Report the resulting average test RMSE of the new model. Does performance increase or decrease? Why?
6. Consider a variant of this model with  $\sigma^2(\mathbf{x}) = \mathbf{w}_\sigma^\top \mathbf{x}$ . How would this model differ? (Optional) Implement this model.
7. How does the performance of the models analyzed in this problem compare to the performance of the model from Problem 2? Which model performs best? Why?

1. DGM for applying ordinal regression model to one data point:



$$\begin{aligned}
 2. \quad p(y|h) &= \prod_{r=1}^R P(Y = r|h)^{\mathbb{1}[y=r]} \\
 &= \prod_{r=1}^R \left( \int_f p(r, f|h) df \right)^{\mathbb{1}[y=r]}
 \end{aligned}$$

$$\begin{aligned}
&= \prod_{r=1}^R \left( \int_f P(Y = r|f, h) p(f|h) df \right)^{\mathbb{1}_{[y=r]}} \\
&= \prod_{r=1}^R \left( \int_{b_r}^{b_{r+1}} P(Y = r|f) p(f|h) df \right)^{\mathbb{1}_{[y=r]}} \\
&= \prod_{r=1}^R \left( \int_{b_r}^{b_{r+1}} p(f|h) df \right)^{\mathbb{1}_{[y=r]}} \\
&= \prod_{r=1}^R \left( \Phi\left(\frac{b_{r+1} - h}{\sigma}\right) - \Phi\left(\frac{b_r - h}{\sigma}\right) \right)^{\mathbb{1}_{[y=r]}} \quad (\Phi \text{ is the standard normal CDF})
\end{aligned}$$

Then we also have  $P(Y = r|h) = \Phi\left(\frac{b_{r+1} - h}{\sigma}\right) - \Phi\left(\frac{b_r - h}{\sigma}\right)$ .

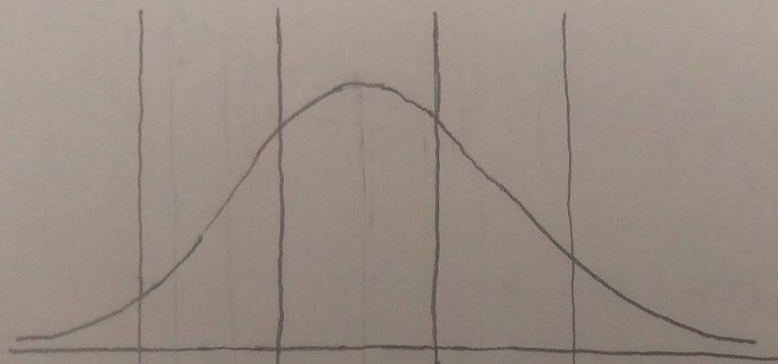
With no additive noise term we have  $P(Y = r|b_r \leq h < b_{r+1}) = 1$  and  $P(Y \neq r|b_r \leq h < b_{r+1}) = 0$ .

3. Using the definition of expectation, we get:

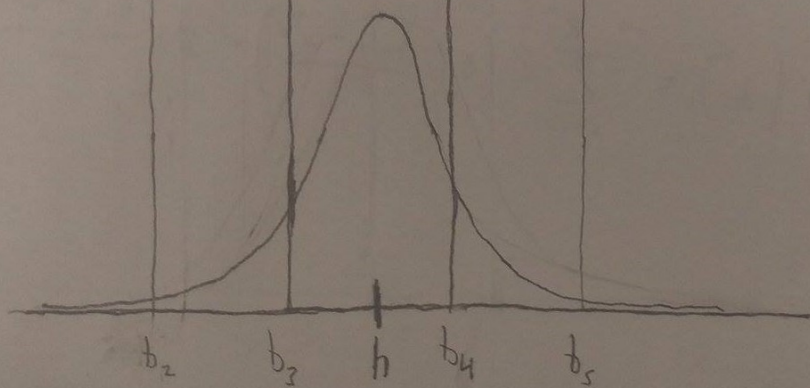
$$E(Y|h) = \sum_{r=1}^R r \left( \Phi\left(\frac{b_{r+1} - h}{\sigma}\right) - \Phi\left(\frac{b_r - h}{\sigma}\right) \right)$$

$P(Y = r|b_r \leq h < b_{r+1})$  will be greater for smaller values  $\sigma$  (see diagram on next page).

Larger  $\sigma$



Smaller  $\sigma$



4. Test set RMSE= 1.2944

$\hat{\sigma}_{MLE}^2 = 6.5794$  (no prior for  $\sigma$  specified, so computing MLE instead of MAP)

$\hat{\mathbf{w}}_{MAP} =$

```
[2.1116e-28 -1.0308e-16 -1.4834e-29 2.1508e-28 -1.3963e-30 -5.6252e-23 -2.8840e-15
-5.4405e-28 2.6699e-28 -6.0757e-29 9.5374e-25 4.2562e-25 3.2805e-27 -1.1288e-28
-1.0367e-28 5.3538e-27 7.5423e-28 6.8906e-27 2.4212e-28 1.6025e-27 1.9136e-27 6.9395e-11
-5.3347e-23 1.2013e-27 8.2377e-30 -4.0415e-27 -1.2174e-27 9.4410e-28 -3.0547e-30 -1.5727e-27
8.3317e-08 -2.2609e-28 8.8179e-26 5.0226e-29 1.1196e-27 -1.3341e-28 -1.5650e-28 -4.4472e-23
-4.3138e-09 1.1640e-26 -7.3984e-29 7.1223e-27 -2.7051e-27 -8.1352e-26 -8.1057e-27 4.8437e-28
-6.2097e-29 -1.0371e-10 6.1987e-28 -3.4184e-18 -1.1628e-27 1.4036e-28 1.1711e-27 -1.7286e-07
-1.8521e-29 2.9590e-28 3.5942e-15 -2.8348e-19 -3.3798e-29 1.6894e-18 9.8680e-21 -7.4688e-31
-5.3513e-28 -7.4957e-28 1.9887e-29 -1.8535e-23 -2.4486e-24 -2.8668e-28 1.8600e-23 -5.2279e-21
2.2171e-24 -2.0490e-29 -1.0312e-24 1.3340e-28 -1.0207e-21 2.7474e-27 2.8038e-27 2.7497e-09
-1.6223e-26 3.5194e-28 -8.9755e-17 3.5600e-29 2.4333e-28 -1.0462e-27 3.7448e-22 -5.3546e-21
8.5839e-28 2.6583e-28 -3.4819e-28 -2.9454e-29 -3.1454e-28 -9.0826e-26 -3.6560e-28 1.1285e-28
-2.3598e-28 -1.4975e-26 -1.1217e-27 -2.4083e-28 -1.7905e-23 2.3561e-25 4.0022e-08 -2.0972e-29
1.7625e-23 -6.6567e-29 -3.1613e-28 -9.2698e-11 1.0426e-10 2.9313e-24 -5.4867e-23 6.7893e-27
-2.0910e-27 -1.5861e-28 -1.9627e-25 -4.8640e-26 2.5109e-20 -2.8539e-27 7.3039e-28 3.3424e-28
-4.4586e-25 -8.8696e-22 1.2661e-28 4.0666e-30 1.2339e-18 -8.9456e-29 -3.7486e-28 -3.6455e-27
-7.7846e-28 1.0433e-24 1.2583e-27 -5.5791e-29 8.2660e-30 -2.5144e-17 -6.1850e-30 6.7540e-27
1.3355e-08 -1.1355e-29 5.9214e-27 3.3451e-27 3.0260e-09 -1.1990e-26 7.5970e-22 -2.2774e-13
2.0166e-24 -5.7185e-28 -1.6015e-28 1.0031e-25 3.8563e-09 3.8764e-26 -4.2684e-18 5.8681e-29
5.4776e-16 4.1298e-27 6.3310e-24 -4.3561e-29 -1.0799e-28 4.7038e-28 -7.5122e-26 3.5160e-28
3.4845e-28 -3.5321e-20 -8.1774e-26 -3.5277e-29]
```

5. Test set RMSE= 1.3170

Performance is slightly worse for the new model relative to the model in part 4. I think that this is because in the new model, we are not taking into account the fact that ratings must take a value in the set  $\{1, 2, 3, 4, 5\}$ .

6. In this variant we allow  $\sigma^2$  to vary based on the joke and model the relation between text and  $\sigma^2$  using a linear model. We then have a new set of parameters  $\mathbf{w}_\sigma$  to learn during training.
7. The models in this problem have worse performances than the models from problem 2. The models in problem 2 learn representations for each user as well as each joke based on the ratings while the models in this problem learn a mapping from text to rating and do not take into account any information about the users. The results indicate that the text is not as useful as the representations learned using the models in problem 2 for predicting joke ratings.