

$$\therefore E(s_1, s_2) = J s_1 s_2 + h_1 s_1 + h_2 s_2 \text{ with } s_1, s_2 \in \{0, 1\}$$

$$\text{Given the Gibbs distribution } \Pr(s_1 = s_1, s_2 = s_2) = \frac{\exp(-\beta E(s_1, s_2))}{Z}$$

$$Z \text{ is partition function } \sum_{s_1} \sum_{s_2} \exp(-\beta E(s_1, s_2))$$

Substitute energy function:

$$Z = \sum_{s_1} \sum_{s_2} \exp(-\beta(-J s_1 s_2 + h_1 s_1 + h_2 s_2))$$

$$= \sum_{s_1} \exp(\beta h_1 s_1) \sum_{s_2} \exp(\beta h_2 s_2) \sum_{s_1} \exp(\beta J s_1 s_2)$$

$$= (1 + \exp(\beta h_1))(1 + \exp(\beta h_2))(1 + 2 \exp(\beta J))$$

$$E(s_1, s_2) = \sum_{s_1} \sum_{s_2} s_1 s_2 \Pr(s_1 = s_1, s_2 = s_2)$$

$$= \sum_{s_1} \sum_{s_2} s_1 s_2 \frac{\exp(\beta J s_1 s_2 - \beta h_1 s_1 - \beta h_2 s_2)}{1 + \exp(\beta h_1)} (1 + \exp(\beta h_2))(1 + 2 \exp(\beta J))$$

$$= \frac{\exp(\beta J) + \exp(2\beta J)}{1 + \exp(\beta h_1)} (1 + \exp(\beta h_2))(1 + 2 \exp(\beta J))$$

2. The KL divergence between p and q is given by

$$D(p \| q) = \sum x p(x) \log\left(p \frac{x}{q}(x)\right)$$

$$L(p, \lambda) = D(p \| q) + \sum i \lambda_i (\sum x p(x) g_i(x) - a_i)$$

$$\log\left(p \frac{x}{q}(x)\right) = -1 - \lambda_1 g_1(x) - \lambda_2 g_2(x) - \dots - \lambda_k g_k(x)$$

$$p(x) = \exp(-1 - \lambda_1 g_1(x) - \lambda_2 g_2(x) - \dots - \lambda_k g_k(x)) q(x)$$

$$\sum x p(x) g_i(x) = a_i, i = 1, 2, \dots$$

$$\sum x \exp(-1 - \lambda_1 g_1(x) - \lambda_2 g_2(x) - \dots - \lambda_k g_k(x)) q(x) g_i(x) = a_i, i = 1, 2, \dots$$

3. We know that $\nabla_{\varphi}(y) = \nabla_{\varphi} \dots$ is the gradient of φ at x .

- We use identity matrix $(\nabla_{\varphi}(x))(y-x)^T(y-x) = \nabla_{\varphi}(x) \cdot \nabla_{\varphi}(x)^T$

$$H = \nabla_{\varphi}(x) \cdot \nabla_{\varphi} \frac{x^T}{(y-x)^T(y-x)}$$

$$u^T \cdot H \cdot u = u^T \cdot \nabla_{\varphi}(x) \cdot \nabla_{\varphi} \frac{x^T}{(y-x)^T(y-x)} \cdot u = (u^T \cdot \nabla_{\varphi}(x)) \cdot \frac{\nabla_{\varphi}(x)^T u}{(y-x)^T(y-x)}$$

$$= \frac{(u^T \cdot \nabla_{\varphi}(x))^2}{(y-x)^T(y-x)} \geq 0$$

$$v^T \cdot H \cdot v = \frac{(v^T \cdot \nabla_{\varphi}(x))^2}{(y-x)^T(y-x)} \geq 0$$

These 2 in conjunction show that the Bregman divergence is convex w.r.t x and y assuming that $\varphi(x)$ is 3 times differentiable.

4. $P(X=k) = p_k$

$$H(X) = -\sum p_k \cdot \log_2(p_k) \leftarrow \text{Entropy}$$

$$p_k = \frac{N_k}{N}$$

Probability of observing the scores of N students: $P(X_1 = x_1, X_2 = x_2, \dots, X_N = x_N)$

→ From this we find the negative log-likelihood of the data:

$$L(\theta) = -\log(P(X_1 = x_1, X_2 = x_2, \dots, X_N = x_N))$$

$$= -(-N \cdot \sum (\frac{N_k}{N}) \cdot \log(N_k) + N \cdot \log(N))$$

$$= N \cdot (-\sum (\frac{N_k}{N}) \cdot \log(N_k)) - N \cdot \log(N)$$

$$= -N \cdot (\sum (\frac{N_k}{N}) \cdot \log(N_k)) + H(X) \cdot N$$

We find that the negative log-likelihood of the data is equal to the negative of entropy multiplied by all students in class. We then add constant factor $N \cdot \log(N)$. They are each related through this factor which is dependent on the total number of students in class. Thus, maximizing the entropy gives the most accurate prediction of the probability distribution of the midterm scores.

$$x_N) = \prod p_{K_K}^N$$