

EECS 545 Machine Learning Homework 4

Tsung-Hung Yao

March 22, 2016

1) Information Theory

Many algorithms for learning probabilistic models are best understood in terms of *information theory*. Consequently, it is useful to understand and manipulate these quantities in different contexts.

Problem (a). Show that

$$I(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

where $I(X, Y)$ is the mutual information of X and Y , $H(X)$ is the entropy of X , and $H(X|Y)$ is the conditional entropy of X given Y .

Proof. First we write down the following definition from the textbook

$$\begin{aligned} H[X] &= - \int p(x) \ln(p(x)) dx = - \int (- \int p(x, y) dy) \ln(p(x)) dx = - \int \int p(y, x) \ln(p(x)) dy dx \\ H[X|Y] &= - \int \int p(x, y) \ln(p(x|y)) dy dx \\ H[Y|X] &= - \int \int p(y, x) \ln(p(y|x)) dy dx \\ I[x, y] &= - \int \int p(y, x) \ln\left(\frac{p(y)p(x)}{p(x,y)}\right) dy dx \\ \implies H[X] - H[X|Y] &= - \int \int p(y, x) \ln(p(x)) dy dx + \int \int p(x, y) \ln(p(x|y)) dy dx \\ &= - \int \int p(x, y) \ln\left(\frac{p(x)}{p(x|y)}\right) dy dx \\ &= - \int \int p(y, x) \ln\left(\frac{p(y)p(x)}{p(x,y)}\right) dy dx = I[x, y] \implies H[Y] - H[Y|X] = - \int \int p(y, x) \ln(p(y)) dy dx + \\ &\quad \int \int p(x, y) \ln(p(y|x)) dy dx \\ &= - \int \int p(x, y) \ln\left(\frac{p(y)}{p(y|x)}\right) dy dx \\ &= - \int \int p(y, x) \ln\left(\frac{p(y)p(x)}{p(x,y)}\right) dy dx = I[x, y] \end{aligned}$$

□

Problem (b). Prove that if X and Y are related by a bijection f (i.e., $X = f(Y)$ and $Y = f^{-1}(X)$), then $I(X, Y) = H(X) = H(Y)$.

Proof. Consider discrete case

$$\begin{aligned} P(Y = y) &= P(f^{-1}(X) = y) = P(X = f(y)) \\ P(X = x) &= P(f(Y) = x) = P(Y = f^{-1}(x)) \\ P(X = x, Y = y) &= P(X = x, f^{-1}(X) = y) = P(X = x, X = f(y)) = P(X = x) \\ P(X = x, Y = y) &= P(f(Y) = x, Y = y) = P(Y = f^{-1}(x), Y = y) = P(Y = y) \end{aligned}$$

$$\begin{aligned}
&\implies I[X, Y] = -\sum_x \sum_y P(X=x, Y=y) \ln\left(\frac{P(X=x)P(Y=y)}{P(X=x, Y=y)}\right) \\
&= -\sum_x \sum_y P(X=x) \ln\left(\frac{P(X=x)P(X=f(y))}{P(X=x)}\right) \\
&= -\sum_x P(X=x) \ln(P(X=x)) = H[X] \\
&\implies I[X, Y] = -\sum_x \sum_y P(X=x, Y=y) \ln\left(\frac{P(X=x)P(Y=y)}{P(X=x, Y=y)}\right) \\
&= -\sum_y \sum_x P(Y=y) \ln\left(\frac{P(Y=f^{-1}(x))}{P(Y=y)}\right) \\
&= -\sum_y P(Y=y) \ln(P(Y=y)) = H[Y]
\end{aligned}$$

□

Problem (c). Suppose we observe N samples $\mathcal{D} = (x_1, \dots, x_N)$ from some unknown distribution. Define $\hat{p}(x)$ to be the empirical probability density estimate,

$$\hat{p}(x) \triangleq \frac{1}{N} \sum_{i=1}^N \mathbb{I}[x = x_i]$$

Let $q(x|\theta)$ be the probability density corresponding to some known probabilistic model with parameter θ . Show that the minimum Kullback–Leibler divergence

$$\min_{\theta} D_{KL}(\hat{p}||q)$$

is obtained by the maximum likelihood estimate θ_{ML} given the data \mathcal{D} .

Proof. Since the empirical probability density estimate is discrete, we consider KL divergence in discrete form

$$\begin{aligned}
KL(p||q) &= -\sum_x p(x) \ln\left(\frac{q(x)}{p(x)}\right) \\
\implies KL(\hat{p}||q(x|\theta)) &= -\sum_x \hat{p}(x) \ln\left(\frac{q(x|\theta)}{\hat{p}(x)}\right) \\
&= \sum_x \left(\frac{1}{N} \sum_{i=1}^N \mathbb{I}[x = x_i] \right) (-\ln(q(x|\theta)) + \frac{1}{N} \sum_{i=1}^N \mathbb{I}[x = x_i])
\end{aligned}$$

Since there's only one term $-\ln(q(x|\theta))$ related to θ , to $\min_{\theta} KL(p||q)$, we only have to $\max_{\theta} \sum_x -\ln(q(x|\theta))$, which is the exactly the MLE problem. □

Problem (d). Let $p = \mathcal{N}(\mu, \sigma^2)$ be a Gaussian distribution and q be any probability density with mean μ and variance σ^2 . Prove that $H(q) \leq H(p)$, that is, the Gaussian distribution has maximum entropy among all distributions of the same variance. *Hint: Refer to the textbook (PRML by Bishop §1.6) for a proof outline.*

Proof. We can rewrite our problem as

$$\begin{aligned}
&\max_{\mathbf{x}} && H[X] \\
&\text{subject to} && \int_{-\infty}^{\infty} p(x) dx = 1 \\
&&& \int_{-\infty}^{\infty} xp(x) dx = \mu \\
&&& \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx = \sigma^2
\end{aligned} \tag{1}$$

Therefore, we can use Lagrange multiplier to solve this problem.

Let $\tilde{H} = \int_{-\infty}^{\infty} -p(x) \ln(p(x)) dx + \lambda_1(\int_{-\infty}^{\infty} p(x) dx - 1) + \lambda_2(\int_{-\infty}^{\infty} xp(x) dx - \mu) + \lambda_3(\int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx - \sigma^2)$

$$\begin{aligned}
& \mu)^2 p(x) dx - \sigma^2 \\
&= \int_{-\infty}^{\infty} (-p(x) \ln(p(x)) + \lambda_1 p(x) + \lambda_2 x p(x) + \lambda_3 (x - \mu)^2 p(x)) dx - \lambda_1 - \lambda_2 \mu - \lambda_3 \sigma^2 \\
&\text{To find extreme value of the function above, we only have to focus on the first term. Also, by the calculus of variation, we know we can find the extremum of } I = \int_b^a f(y, y', x) dx \text{ when } y(x) \text{ satisfies Euler-Lagrange equation: } \frac{\partial f}{\partial y} - \frac{d}{dx} \left(\frac{\partial f}{\partial y'} \right) = 0 \\
&\implies f(p(x), p(x)', x) = -p(x) \ln(p(x)) + \lambda_1 p(x) + \lambda_2 x p(x) + \lambda_3 (x - \mu)^2 p(x) \\
&\frac{\partial f}{\partial p} - \frac{d}{dx} \left(\frac{\partial f}{\partial p'} \right) \\
&= -\ln(p(x)) - 1 + \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 = 0 \\
&\implies e^{(-1+\lambda_1+\lambda_2 x+\lambda_3(x-\mu)^2)} = p(x) \text{ would maximize the original problem } I = \int_b^a f(y, y', x) dx \\
&= \int_{-\infty}^{\infty} (-p(x) \ln(p(x)) + \lambda_1 p(x) + \lambda_2 x p(x) + \lambda_3 (x - \mu)^2 p(x)) dx \\
&\text{Consider } X \sim N(\mu, \sigma^2) \implies \text{constraints satisfied and } \lambda_1 = 1 + \ln(\frac{1}{\sqrt{2\pi\sigma^2}}), \lambda_2 = 0, \lambda_3 = -\frac{1}{2\sigma^2} \quad \square
\end{aligned}$$

2) Dirichlet Maximum Likelihood

In this problem, you will derive and implement a Newton-Raphson algorithm for maximizing the Dirichlet log-likelihood function. Unlike for the simple distributions we have encountered in the past (Multinomial, Poisson, etc.), no closed-form solution exists for the Maximum Likelihood estimate of Dirichlet parameters.

Recall a Dirichlet-distributed random vector $\mathbf{p} = (p_1, \dots, p_m) \in \Delta^{m-1}$ governed by nonnegative concentration parameters $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_m)$ has following distribution,

$$\text{Dirichlet}(\mathbf{p}|\boldsymbol{\alpha}) = \frac{\Gamma(\sum_{k=1}^m \alpha_k)}{\prod_{k=1}^m \Gamma(\alpha_k)} \prod_{k=1}^m p_k^{\alpha_k - 1}$$

where $\Gamma(t)$ is the *Gamma function* and Δ^{m-1} is the unit simplex in \mathbb{R}^m .

Problem (a). Show that the Dirichlet distribution belongs to the *exponential family*, that is, find a natural parameter $\eta \triangleq \eta(\boldsymbol{\alpha})$, a sufficient statistics function $T(\mathbf{p})$, and a log-partition function $A(\boldsymbol{\alpha})$ such that

$$\text{Dirichlet}(\mathbf{p}|\boldsymbol{\alpha}) = \exp [\eta(\boldsymbol{\alpha})^T T(\mathbf{p}) - A(\boldsymbol{\alpha})]$$

This guarantees that the log-likelihood is convex and Newton's method will converge to a global optimum.

$$\begin{aligned}
&\text{Proof. } \text{Dirichlet}(\mathbf{p}|\boldsymbol{\alpha}) = \frac{\Gamma(\sum_{k=1}^m \alpha_k)}{\prod_{k=1}^m \Gamma(\alpha_k)} \prod_{k=1}^m p_k^{\alpha_k - 1} \\
&= \exp \left(\ln \left(\Gamma \left(\sum_{k=1}^m \alpha_k \right) \right) - \sum_{k=1}^m \ln \left(\Gamma(\alpha_k) \right) + \sum_{k=1}^m (\alpha_k - 1) \log p_k \right) \\
&= \exp [\eta(\boldsymbol{\alpha})^T T(\mathbf{p}) - A(\boldsymbol{\alpha})]
\end{aligned}$$

, where

$$\begin{aligned}
A(\boldsymbol{\alpha}) &= \sum_{k=1}^m \ln(\Gamma(\alpha_k)) - \ln(\Gamma(\sum_{k=1}^m \alpha_k)) \\
\eta(\boldsymbol{\alpha})^T &= [\alpha_1 - 1 \quad \alpha_2 - 1 \quad \alpha_3 - 1 \quad \cdots \quad \alpha_k - 1] \\
T(\mathbf{p})^T &= [\log p_1 \quad \log p_2 \quad \log p_3 \quad \cdots \quad \log p_k] \quad \square
\end{aligned}$$

Problem (b). Given observations $\mathcal{D} = (\mathbf{p}^{(1)}, \dots, \mathbf{p}^{(N)})$, derive an expression for the Dirichlet log-likelihood function $F(\boldsymbol{\alpha}) = \log P(\mathcal{D}|\boldsymbol{\alpha})$ in terms of the *observed sufficient statistics*,

$$\hat{t}_k = \frac{1}{N} \sum_{j=1}^N \log p_k^{(j)}$$

$$\begin{aligned} F(\boldsymbol{\alpha}) &= \log P(\mathcal{D}|\boldsymbol{\alpha}) = \sum_{j=1}^N (\ln(\Gamma(\sum_{k=1}^m \alpha_k)) - \sum_{k=1}^m \ln(\Gamma(\alpha_k)) + \sum_{k=1}^m (\alpha_k - 1) \ln(p_k^{(j)})) \\ &= N \ln(\Gamma(\sum_{k=1}^m \alpha_k)) - N \sum_{k=1}^m \ln(\Gamma(\alpha_k)) + \sum_{k=1}^m [(\alpha_k - 1) \sum_{j=1}^N \ln(p_k^{(j)})] \\ &= N \ln(\Gamma(\sum_{k=1}^m \alpha_k)) - N \sum_{k=1}^m \ln(\Gamma(\alpha_k)) + N \sum_{k=1}^m [(\alpha_k - 1) \hat{t}_k] \end{aligned}$$

Problem (c). Derive an expression for the gradient of the log-likelihood in terms of the observed sufficient statistics and the *digamma function*,

$$\Psi(t) = \frac{d \log \Gamma(t)}{dt} = \frac{\Gamma'(t)}{\Gamma(t)}$$

(Hint: Specify each component $\frac{\partial F}{\partial \alpha_k}$ separately, rather than trying to use matrix operations.)

$$\begin{aligned} \text{Proof. } \frac{\partial F(\boldsymbol{\alpha})}{\partial \alpha_j} &= \frac{\partial}{\partial \alpha_j} N \ln(\Gamma(\sum_{k=1}^m \alpha_k)) - N \sum_{k=1}^m \ln(\Gamma(\alpha_k)) + N \sum_{k=1}^m [(\alpha_k - 1) \hat{t}_k] \\ &= N \frac{\partial}{\partial \alpha_j} \ln(\Gamma(\sum_{k=1}^m \alpha_k)) - N \frac{\partial}{\partial \alpha_j} \sum_{k=1}^m \ln(\Gamma(\alpha_k)) + N \frac{\partial}{\partial \alpha_j} \sum_{k=1}^m [(\alpha_k - 1) \hat{t}_k] \\ &= N \frac{\Gamma'(\sum_{k=1}^m \alpha_k)}{\Gamma(\sum_{k=1}^m \alpha_k)} - N \frac{\Gamma'(\alpha_j)}{\Gamma(\alpha_j)} + N \hat{t}_j \\ &= N \Psi(\sum_{k=1}^m \alpha_k) - N \Psi(\alpha_j) + N \hat{t}_j \end{aligned}$$

Problem (d). Show that the Hessian matrix $H \triangleq \nabla_{\boldsymbol{\alpha}}^2 F(\boldsymbol{\alpha})$ of the log-likelihood can be written as the sum of a diagonal matrix and a matrix whose elements are all the same,

$$H = \nabla_{\boldsymbol{\alpha}}^2 F(\boldsymbol{\alpha}) = Q + c \mathbf{1} \mathbf{1}^T$$

where $c \in \mathbb{R}$ is a constant and $Q \in \mathbb{R}^{m \times m}$ is diagonal with entries q_{11}, \dots, q_{mm} . (Note: It is okay to simply write Ψ' for the derivative of the digamma function.)

$$\text{Proof. } H_{i,j} = (\nabla_{\boldsymbol{\alpha}}^2 F(\boldsymbol{\alpha}))_{i,j} = \frac{\partial^2 F(\boldsymbol{\alpha})}{\partial \alpha_i \partial \alpha_j} = \frac{\partial}{\partial \alpha_i} (N \Psi(\sum_{k=1}^m \alpha_k) - N \Psi(\alpha_j) + N \hat{t}_j)$$

When $i \neq j$, we can have the off-diagonal term:

$$H_{i,j} = \frac{\partial}{\partial \alpha_i} (N \Psi(\sum_{k=1}^m \alpha_k) - N \Psi(\alpha_j) + N \hat{t}_j) = N \Psi'(\sum_{k=1}^m \alpha_k)$$

When $i=j$, we can have the diagonal term:

$$H_{i,i} = \frac{\partial}{\partial \alpha_i} (N \Psi(\sum_{k=1}^m \alpha_k) - N \Psi(\alpha_j) + N \hat{t}_j) = N \Psi'(\sum_{k=1}^m \alpha_k) - N \Psi'(\alpha_j)$$

$\implies H = Q + c \mathbf{1} \mathbf{1}^T$, where

$$c = N \Psi'(\sum_{k=1}^m \alpha_k)$$

$$Q = \begin{bmatrix} -N \Psi'(\alpha_1) & 0 & \cdots & 0 \\ 0 & -N \Psi'(\alpha_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & -N \Psi'(\alpha_M) \end{bmatrix}$$

□

Problem (e). The Newton-Raphson method provides a quadratically converging method for parameter estimation. The general update rule can be written in terms of the Hessian matrix as follows:

$$\boldsymbol{\alpha}^{new} = \boldsymbol{\alpha}^{old} - [H_F(\boldsymbol{\alpha}^{old})]^{-1} \cdot \nabla F(\boldsymbol{\alpha}^{old})$$

Use the Sherman-Morrison matrix inversion lemma to derive a closed-form update for $\boldsymbol{\alpha}$ (any remaining matrix inversions should be diagonal).

Proof. Sherman-Morrison lemma: $(\mathbf{A} + c\mathbf{u}\mathbf{v}^T)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1}\mathbf{u}\mathbf{v}^T\mathbf{A}^{-1}}{1/c + \mathbf{v}^T\mathbf{A}^{-1}\mathbf{u}}$, if $1/c + \mathbf{v}^T\mathbf{A}^{-1}\mathbf{u} \neq 0$ where \mathbf{A} is a invertible squared matrix and \mathbf{u} and \mathbf{v} are column vector.

$$\begin{aligned} \Rightarrow H^{-1} &= (Q + c\mathbf{1}\mathbf{1}^T)^{-1} = Q^{-1} - \frac{Q^{-1}\mathbf{1}\mathbf{1}^TQ^{-1}}{1/c + \mathbf{1}^TQ^{-1}\mathbf{1}} \\ \Rightarrow \mathbf{1}^TQ^{-1}\mathbf{1} &= [1 \ 1 \ 1 \ \dots \ 1] \begin{bmatrix} 1/q_1 & 0 & 0 & \dots & 0 \\ 0 & 1/q_2 & 0 & \dots & 0 \\ 0 & 0 & 1/q_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1/q_m \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \sum_{k=1}^m \frac{1}{q_k} \\ \Rightarrow Q^{-1}\mathbf{1}\mathbf{1}^TQ^{-1} &= \begin{bmatrix} 1/q_1 \\ 1/q_2 \\ 1/q_3 \\ \vdots \\ 1/q_m \end{bmatrix} [1/q_1 \ 1/q_2 \ /q_3 \ \dots \ 1/q_m] = \begin{bmatrix} 1/q_1 q_1 & 1/q_1 q_2 & 1/q_1 q_3 & \dots & 1/q_1 q_m \\ 1/q_2 q_1 & 1/q_2 q_2 & 1/q_2 q_3 & \dots & 1/q_2 q_m \\ 1/q_3 q_1 & 1/q_3 q_2 & 1/q_3 q_3 & \dots & 1/q_3 q_m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1/q_m q_1 & 1/q_m q_2 & 1/q_m q_3 & \dots & 1/q_m q_m \end{bmatrix} \end{aligned}$$

$$\text{Let } \nabla F(\boldsymbol{\alpha}^{old}) = \mathbf{g} \implies [H_F(\boldsymbol{\alpha}^{old})]^{-1} \cdot \nabla F(\boldsymbol{\alpha}^{old}) = Q^{-1}\mathbf{g} - \frac{Q^{-1}\mathbf{1}\mathbf{1}^TQ^{-1}\mathbf{g}}{1/c + \mathbf{1}^TQ^{-1}\mathbf{1}}$$

$$\begin{aligned} &= \begin{bmatrix} 1/q_1 & 0 & 0 & \dots & 0 \\ 0 & 1/q_2 & 0 & \dots & 0 \\ 0 & 0 & 1/q_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1/q_m \end{bmatrix} \begin{bmatrix} g_1 \\ g_2 \\ g_3 \\ \vdots \\ g_m \end{bmatrix} - \frac{1}{\frac{1}{c} + \sum_{k=1}^m \frac{1}{q_k}} \begin{bmatrix} 1/q_1 q_1 & 1/q_1 q_2 & 1/q_1 q_3 & \dots & 1/q_1 q_m \\ 1/q_2 q_1 & 1/q_2 q_2 & 1/q_2 q_3 & \dots & 1/q_2 q_m \\ 1/q_3 q_1 & 1/q_3 q_2 & 1/q_3 q_3 & \dots & 1/q_3 q_m \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1/q_m q_1 & 1/q_m q_2 & 1/q_m q_3 & \dots & 1/q_m q_m \end{bmatrix} \begin{bmatrix} g_1 \\ g_2 \\ g_3 \\ \vdots \\ g_m \end{bmatrix} \\ &= \begin{bmatrix} g_1/q_1 \\ g_2/q_2 \\ g_3/q_3 \\ \vdots \\ g_m/q_m \end{bmatrix} - \frac{1}{\frac{1}{c} + \sum_{k=1}^m \frac{1}{q_k}} \begin{bmatrix} \frac{g_1}{q_1} \sum_{k=1}^m \frac{1}{q_k} \\ \frac{g_2}{q_2} \sum_{k=1}^m \frac{1}{q_k} \\ \frac{g_3}{q_3} \sum_{k=1}^m \frac{1}{q_k} \\ \vdots \\ \frac{g_m}{q_m} \sum_{k=1}^m \frac{1}{q_k} \end{bmatrix}, \text{ where} \end{aligned}$$

$$q_i = -N\Psi'(\alpha_i^{old}) \text{ and } g_i = N\Psi(\sum_{k=1}^m \alpha_k^{old}) - N\Psi(\alpha_i^{old}) + N\hat{t}_i$$

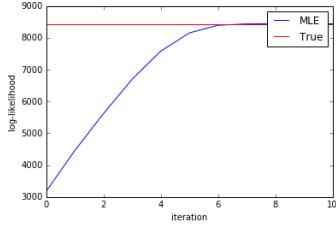
□

Problem (f). Implement this Newton-Raphson update in your language of choice:

- Generate $N = 1000$ samples from the Dirichlet distribution with parameter $\boldsymbol{\alpha} = (10, 5, 15, 20, 50)$ (see `hw4p2.py` on Canvas). Use an initial estimate of $\hat{\boldsymbol{\alpha}} = (1, 1, 1, 1, 1)$.
- The following Python functions may be useful: `from scipy.special import gammaln, polygamma`.

Deliverables:

- A plot of the log-likelihood as a function of iteration number. Also plot the log-likelihood given the true parameters as a constant (horizontal) line. Terminate your algorithm when the log-likelihood increases by less than 10^{-4} .
- The estimated model parameters.
- Please submit your code, as usual.



Proof. The estimated parameter: $(10.14556184 \quad 5.06132558 \quad 15.0496422 \quad 20.35846602 \quad 50.72249073)$

□

3) Graphical Models

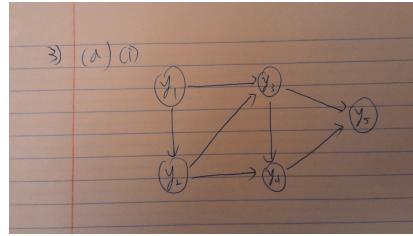
In this problem, you will explore the independence properties of directed graphical models and practice translating them to factored probability distributions and back.

Problem (a). Draw a directed graphical model for each of the following factored distributions. Take advantage of plate notation when convenient, and represent as many independencies with your graph as possible (i.e., don't draw a fully connected graph!).

$$(i) P(y_1, y_2, y_3, y_4, y_5) = P(y_1)P(y_2|y_1)\prod_{k=3}^5 P(y_k|y_{k-1}, y_{k-2})$$

$$\begin{aligned} \text{Proof. } P(y_1, y_2, y_3, y_4, y_5) &= P(y_1)P(y_2|y_1)\prod_{k=3}^5 P(y_k|y_{k-1}, y_{k-2}) \\ &= P(y_1)P(y_2|y_1)P(y_3|y_2, y_1)P(y_4|y_3, y_2)P(y_5|y_4, y_3) \end{aligned}$$

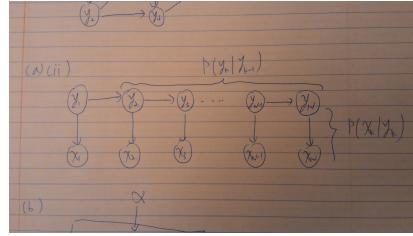
□



$$(ii) P(x_1, \dots, x_N, y_1, \dots, y_N) = P(y_1) \prod_{k=2}^N P(y_k|y_{k-1}) \prod_{k=1}^N P(x_k|y_k)$$

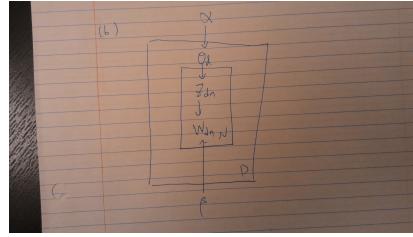
Proof.

□



Problem (b). Draw a directed graphical model for the following model specification, where $\alpha \in \mathbb{R}^K$ and $\beta \in \mathbb{R}^{K \times W}$ are fixed hyperparameters, and β_k denotes the k^{th} row of β :

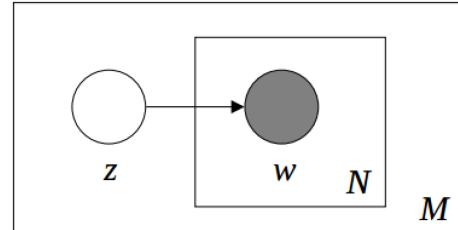
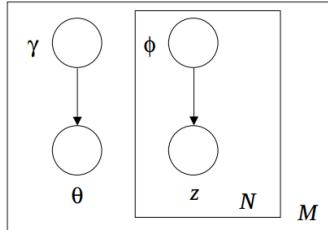
$$\begin{aligned} \theta_1, \dots, \theta_d &\stackrel{iid}{\sim} \text{Dirichlet}(\alpha) \\ z_{d1}, \dots, z_{dN} | \theta_d &\stackrel{iid}{\sim} \text{Categorical}(\theta_d) \quad \forall d \in \{1, \dots, D\} \\ w_{dj} | z_{dj} &\sim \text{Categorical}(\beta_{z_{dj}}) \quad \forall d \in \{1, \dots, D\}, j \in \{1, \dots, N\} \end{aligned}$$



Proof.

□

Problem (c). For each directed model below, write down the factorized joint distribution over all variables.



Proof. joint probability of pic on the left:

$$\prod_{m=1}^M [p(\gamma_m)p(\theta_m|\gamma_m)(\prod_{n=1}^N p(\phi_n|z_n)p(z_n))]$$

joint probability of the pic on the right:

$$\prod_{j=1}^M p(z_j)[\prod_{i=1}^N p(w_{ij}|z_j)]$$

□

4) Clustering

Download the image `mandrill.png` from Canvas. In this problem you will apply the k -means algorithm to image compression. In this context it is also known as the Lloyd-Max algorithm.

Problem (a). First, partition the image into blocks of size $M \times M$ and reshape each block into a vector of length $3M^2$ (see `hw4p4.py` on Canvas). The 3 comes from the fact that this is a color image, and so there are three intensities for each pixel. Assume that M , like the image dimensions, is a power of 2.

Next, write a program that will cluster the vectors from (a) using the k -means algorithm. **You should implement the k -means algorithm yourself.** Please initialize the cluster means to be randomly selected data points, sampled without replacement.

Finally, reconstruct a quantized version of the original image by replacing each block in the original image by the nearest centroid. Test your code using $M = 2$ and $k = 64$.

Deliverables:

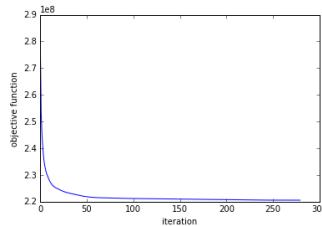
- A plot of the k -means objective function value versus iteration number.
- A description of how the compressed image looks compared to the original. What regions are best preserved, and which are not?
- A picture of the difference of the two images. You should add a neutral gray (128, 128, 128) to the difference before generating the image.
- The compression ratio (use your formula from (b)).
- The relative mean absolute error of the compressed image, defined as

$$\frac{\frac{1}{3N^2} \sum_{i=1}^N \sum_{j=1}^N \sum_{r=1}^3 |\tilde{I}(i, j, r) - I(i, j, r)|}{255}$$

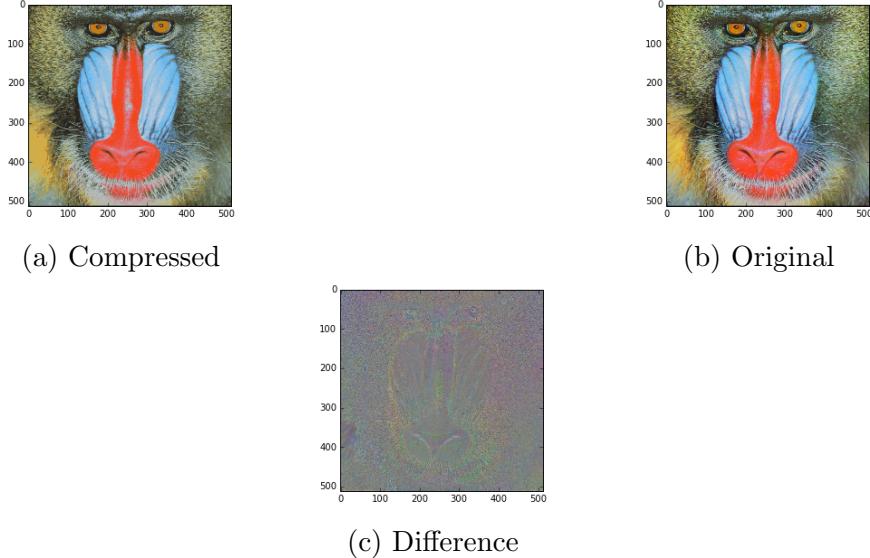
where \tilde{I} and I are the compressed and original images, respectively, viewed as 3-D arrays. This quantity can be viewed as the average error in pixel intensity relative to the range of pixel intensities.

- Please submit your code, as usual.

Proof. 1. Objective Function



2. Compressed image versus Original one: From the picture below, it is really hard to see the difference between the compressed pic and original pic. (Also, it depends on the computer you use.) As for me, I can only see the tiny difference between the color of the nose.
3. Difference of 2 image is shown below



4. Compression Ratio is 0.0634765625 and relative mean absolute error is 0.0500294082598

1

Problem (b). The original uncompressed image uses 24 bits per pixel (bpp), 8 bits for each color. Assuming an image of size $N \times N$, where N is a power of 2, what is the number of bits per pixel, as a function of M , N , and k , needed to store the compressed image? What is the compression ratio, which is defined as the ratio of bpp in the compressed image relative to the original uncompressed image? *Hint: To calculate the bpp of the compressed image, imagine you need to transmit the compressed image to a friend who knows the compression strategy as well as the values of M , N , and k .*

Proof. Uncompressed image uses 24 bits per pixel \implies For the uncompressed image with $N \times N$ pixel, it uses $24N^2$ bits.

Also, assume the size of the compressed image is $M \times M$ and also has value $24M^2$, but the compressed image needs to store some more information to uncompress the image. First, we need to know each pixel belongs to which cluster. Therefore, the need $24M^2 \times k$ to store the compressed image. Last, the compressed image also needs to store the transformation relation between the original and compressed, which takes additional $\frac{N^2}{M^2} \times \log_2 k$

\Rightarrow The total space need to store by compressed image is $\frac{N^2}{M^2} \times \log_2 k + 24M^2 \times k$
 \Rightarrow Compression ratio = $\frac{\text{bits per pixel in the compressed image}}{\text{bits per pixel in the uncompressed image}} = \frac{\log_2 k}{24M^2} + k \frac{M^2}{N^2}$

Problem (c). (Optional, ungraded) Play around with M and k .

5) EM Algorithm for Mixed Linear Regression

Consider regression training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, iid realizations of $(\mathbf{X}, Y) \in \mathbb{R}^d \times \mathbb{R}$, where the conditional distribution of Y given \mathbf{X} is modeled by the pdf

$$f(y|\mathbf{x}; \boldsymbol{\theta}) \sim \sum_{k=1}^K \pi_k \phi(y; \mathbf{w}_k^T \mathbf{x} + b_k, \sigma_k^2),$$

where $\boldsymbol{\theta} = (\pi_1, \dots, \pi_K, \mathbf{w}_1, \dots, \mathbf{w}_K, b_1, \dots, b_K, \sigma_1^2, \dots, \sigma_K^2)$ is a list of the model parameters such that $\pi_k \geq 0$ and $\sum_{k=1}^K \pi_k = 1$, and $\phi(y; \mu, \sigma^2)$ is the pdf of a Gaussian random variable with mean μ and variance σ^2 evaluated at y . Viewing the \mathbf{x}_i as deterministic, derive an EM algorithm for maximizing the likelihood by following these steps.

Problem (a). Denote $\underline{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ and $\underline{y} = (y_1, \dots, y_N)$. Write down the formula for the log-likelihood, $\log f(\underline{y}|\underline{\mathbf{x}}; \boldsymbol{\theta})$, where $f(\underline{y}|\underline{\mathbf{x}}; \boldsymbol{\theta})$ is the model (with parameters $\boldsymbol{\theta}$) for \underline{y} given $\underline{\mathbf{x}}$.

$$\begin{aligned} \text{Proof. } f(y|\mathbf{x}; \boldsymbol{\theta}) &\sim \sum_{k=1}^K \pi_k \phi(y; \mathbf{w}_k^T \mathbf{x} + b_k, \sigma_k^2) = \sum_{k=1}^K \pi_k \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(\frac{-(y - \mathbf{w}_k^T \mathbf{x} + b_k)^2}{2\sigma_k^2}\right) \\ \therefore (\mathbf{x}_i, y_i) \text{ are iid} \\ \implies f(\underline{y}|\underline{\mathbf{x}}; \boldsymbol{\theta}) &= \prod_{i=1}^N (f(y_i|\mathbf{x}_i; \boldsymbol{\theta})) \\ \implies \log f(\underline{y}|\underline{\mathbf{x}}; \boldsymbol{\theta}) &= \sum_{i=1}^N \log(f(y_i|\mathbf{x}_i; \boldsymbol{\theta})) = \sum_{i=1}^N \left(\log\left(\sum_{k=1}^K \pi_k \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(\frac{-(y_i - \mathbf{w}_k^T \mathbf{x}_i + b_k)^2}{2\sigma_k^2}\right)\right) \right) \end{aligned}$$

□

Problem (b). Introduce hidden variable $\underline{z} = (z_1, \dots, z_N)$ which determines the mixture component that y is drawn from, i.e., $f(y|\mathbf{x}, z = k; \boldsymbol{\theta}) = \phi(y; \mathbf{w}_k^T \mathbf{x} + b_k, \sigma_k^2)$. Write down the complete-data log-likelihood, $\log f(\underline{y}, \underline{z}|\underline{\mathbf{x}}; \boldsymbol{\theta})$. Hint: Define a random variable $\Delta_{ik} \triangleq \mathbb{I}[z_i = k]$ so that you can write the log-likelihood as a double sum and so that \underline{z} only appears in the expression through Δ_{ik} .

$$\begin{aligned} \text{Proof. } \log f(y_i, z_i|\mathbf{x}_i, \theta) &= \log \prod_{k=1}^K (f(y_i|z_i = k, \mathbf{x}_i, \theta) f(z_i = k))^{\Delta_{ik}} \\ &= \log \prod_{k=1}^K \pi_k \phi(y_i; \mathbf{w}_k^T \mathbf{x}_i + b_k, \sigma_k^2)^{\Delta_{ik}} \\ &= \sum_{k=1}^K \Delta_{ik} (\log \pi_k \phi(y_i; \mathbf{w}_k^T \mathbf{x}_i + b_k, \sigma_k^2)) \\ \implies \log f(\underline{y}, \underline{z}|\underline{\mathbf{x}}; \boldsymbol{\theta}) &= \log \prod_{i=1}^N f(y_i, z_i|\mathbf{x}_i, \theta) \\ &= \sum_{i=1}^N \log f(y_i, z_i|\mathbf{x}_i, \theta) \\ &= \sum_{i=1}^N \sum_{k=1}^K \Delta_{ik} (\log \pi_k \phi(y_i; \mathbf{w}_k^T \mathbf{x}_i + b_k, \sigma_k^2)) \end{aligned}$$

□

Problem (c). Determine the E-step. Give an explicit formula for

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \mathbb{E}_{\underline{z}} [\log f(\underline{y}, \underline{z}|\underline{\mathbf{x}}; \boldsymbol{\theta}) | \underline{y}, \underline{\mathbf{x}}; \boldsymbol{\theta}^{\text{old}}]$$

in terms of $\boldsymbol{\theta}$, $\boldsymbol{\theta}^{\text{old}}$, and the data. Remember that in this expectation, you should treat only \underline{z} as a random variable, and the distribution of \underline{z} is conditioned on \underline{y} and $\underline{\mathbf{x}}$ and governed by the parameters $\boldsymbol{\theta}^{\text{old}}$. The log-likelihood inside the expectation, on the other hand, is parameterized by $\boldsymbol{\theta}$ and is non-random for a fixed \underline{z} .

$$\begin{aligned}
\text{Proof. } Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) &= \mathbb{E}_{\underline{z}} [\log f(\underline{y}, \underline{z} | \underline{x}; \boldsymbol{\theta}) | \underline{y}, \underline{x}; \boldsymbol{\theta}^{\text{old}}] \\
&= \mathbb{E}_{\underline{z}} \left[\sum_{i=1}^N \sum_{k=1}^K \Delta_{ik} (\log \pi_k \phi(y_i; \mathbf{w}_k^T \mathbf{x}_i + b_k, \sigma_k^2)) | \underline{y}, \underline{x}; \boldsymbol{\theta}^{\text{old}} \right] \\
&= \sum_{i=1}^N \sum_{k=1}^K \mathbb{E}_{\underline{z}} [\mathbb{I}[z_i = k] (\log \pi_k \phi(y_i; \mathbf{w}_k^T \mathbf{x}_i + b_k, \sigma_k^2)) | \underline{y}, \underline{x}; \boldsymbol{\theta}^{\text{old}}] \\
&= \sum_{i=1}^N \sum_{k=1}^K \mathbb{E}_{\underline{z}} [\mathbb{I}[z_i = k | \underline{y}, \underline{x}; \boldsymbol{\theta}^{\text{old}}] \log \pi_k \phi(y_i | \mathbf{w}_k^T \mathbf{x}_i + b_k, \sigma_k^2)] \\
&= \sum_{i=1}^N \sum_{k=1}^K p(z_i = k | y_i, \mathbf{x}_i, \boldsymbol{\theta}^{\text{old}}) \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K p(z_i = k | \underline{y}, \underline{x}, \boldsymbol{\theta}^{\text{old}}) \log \phi(y_i | \mathbf{w}_k^T \mathbf{x}_i + b_k, \sigma_k^2) \\
&= \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \phi(y_i | \mathbf{w}_k^T \mathbf{x}_i + b_k, \sigma_k^2)
\end{aligned}$$

, where

$$\begin{aligned}
r_{ik} &= p(z_i = k | y_i, \mathbf{x}_i, \boldsymbol{\theta}^{\text{old}}) = \frac{p(z_i=k, y_i | \mathbf{x}_i, \boldsymbol{\theta}^{\text{old}})}{p(y_i | \mathbf{x}_i, \boldsymbol{\theta}^{\text{old}})} \\
&= \frac{p(z_i=k | \boldsymbol{\theta}^{\text{old}}) p(y_i | z_i=k, \mathbf{x}_i, \boldsymbol{\theta}^{\text{old}})}{p(y_i | \mathbf{x}_i, \boldsymbol{\theta}^{\text{old}})} \\
&= \frac{\pi_k^{\text{old}} \phi(y_i | \mathbf{w}_k^{T, \text{old}} \mathbf{x}_i + b_k^{\text{old}}, \sigma_k^{2, \text{old}})}{\sum_{k'=1}^K \pi_k^{\text{old}} \phi(y_i | \mathbf{w}_k^{T, \text{old}} \mathbf{x}_i + b_k^{\text{old}}, \sigma_k^{2, \text{old}})}
\end{aligned}$$

□

Problem (d). Determine the M-step. That is, determine the $\boldsymbol{\theta}$ that maximizes $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}})$.
Suggestions: Use Lagrange multiplier theory to optimize the weights π_k . To optimize $(\mathbf{w}_k, b_k, \sigma_k^2)$, first hold σ_k^2 fixed and find the optimal (\mathbf{w}_k, b_k) , then plug that in and find the optimal σ_k^2 . Just treat σ_k^2 as a variable (not the square of a variable).

Proof. We can see this as a max problem:

$$\begin{aligned}
\max_{\pi} \quad & Q[\theta, \theta^{\text{old}}] \\
\text{subject to} \quad & \sum_{k=1}^K \pi_k = 1
\end{aligned} \tag{2}$$

$$\begin{aligned}
\text{Let } \Lambda(\pi, \lambda) &= Q(\theta, \theta^{\text{old}}) + \lambda (\sum_{k=1}^K (\pi_k - 1)) \\
\implies \frac{\partial \Lambda}{\partial \pi_k} &= \sum_{i=1}^N \left(\frac{r_{ik}}{\pi_k} \right) + \lambda = 0 \\
\implies \pi_k &= \frac{-\sum_{i=1}^N r_{ik}}{\lambda} \\
\implies \frac{\partial \Lambda}{\partial \lambda} &= \sum_{k=1}^K \pi_k - 1 = 0 \\
\implies \sum_{k=1}^K \pi_k &= \sum_{k=1}^K \frac{-\sum_{i=1}^N r_{ik}}{\lambda} = 1 \implies \lambda = -\sum_{k=1}^K \sum_{i=1}^N r_{ik} \\
\implies \pi_k &= \frac{\sum_{i=1}^N r_{ik}}{\lambda} = \frac{\sum_{i=1}^N r_{ik}}{\sum_{k=1}^K \sum_{i=1}^N r_{ik}} = \frac{\sum_{i=1}^N r_{ik}}{N}
\end{aligned}$$

$$Q(\theta, \theta^{\text{old}}) = \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \pi_k + \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \phi(y_i | \mathbf{w}_k^T \mathbf{x}_i + b_k, \sigma_k^2)$$

Since the first term of the Q is not related to w and b, we can only focus on the second term.

$$\text{Let } \mathbf{s}_i^T = (1, \mathbf{x}_i^T) \implies \mathbf{S}_{N \times (1+d)} = \begin{pmatrix} \mathbf{s}_1^T \\ \mathbf{s}_2^T \\ \mathbf{s}_3^T \\ \vdots \\ \mathbf{s}_N^T \end{pmatrix}, \beta_{k,(1+d) \times 1} = (b_k, \mathbf{w}_k^T)^T$$

$$\text{and } \mathbf{R}_k = \begin{pmatrix} r_{1k} & 0 & 0 & \cdots & 0 \\ 0 & r_{2k} & 0 & \cdots & 0 \\ 0 & 0 & r_{3k} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & r_{Nk} \end{pmatrix}$$

$$\begin{aligned} \sum_{i=1}^N \sum_{k=1}^K r_{ik} \log \phi(y_i | \mathbf{w}_k^T \mathbf{x}_i + b_k, \sigma_k^2) &= \sum_{k=1}^K \sum_{i=1}^N r_{ik} \log \frac{1}{\sqrt{2\pi\sigma_k^2}} - \sum_{k=1}^K \sum_{i=1}^N \frac{(\sqrt{r_{ik}}(y_i - \mathbf{w}_k^T \mathbf{x}_i - b_k))^2}{2\pi\sigma_k^2} \\ &= \sum_{k=1}^K \sum_{i=1}^N r_{ik} \log \frac{1}{\sqrt{2\pi\sigma_k^2}} - \sum_{k=1}^K \frac{1}{2\pi\sigma_k^2} (\mathbf{R}_k^{1/2} \mathbf{Y} - \mathbf{R}_k^{1/2} \mathbf{S} \boldsymbol{\beta}_k)^T (\mathbf{R}_k^{1/2} \mathbf{Y} - \mathbf{R}_k^{1/2} \mathbf{S} \boldsymbol{\beta}_k) \\ \frac{\partial \Lambda}{\partial \boldsymbol{\beta}_k} &= \frac{\partial Q}{\partial \boldsymbol{\beta}_k} = \frac{\partial}{\partial \boldsymbol{\beta}_k} \left(-\frac{1}{2\pi\sigma_k^2} (\mathbf{R}_k^{1/2} \mathbf{Y} - \mathbf{R}_k^{1/2} \mathbf{S} \boldsymbol{\beta}_k)^T (\mathbf{R}_k^{1/2} \mathbf{Y} - \mathbf{R}_k^{1/2} \mathbf{S} \boldsymbol{\beta}_k) \right) = 0 \\ \implies \frac{\partial}{\partial \boldsymbol{\beta}_k} (\mathbf{R}_k^{1/2} \mathbf{Y} - \mathbf{R}_k^{1/2} \mathbf{S} \boldsymbol{\beta}_k)^T (\mathbf{R}_k^{1/2} \mathbf{Y} - \mathbf{R}_k^{1/2} \mathbf{S} \boldsymbol{\beta}_k) &= 0 \\ \implies \frac{\partial}{\partial \boldsymbol{\beta}_k} (\mathbf{Y}^T \mathbf{R}_k \mathbf{Y} - 2\boldsymbol{\beta}_k^T \mathbf{S}^T \mathbf{R}_k \mathbf{Y} + \boldsymbol{\beta}_k^T \mathbf{S}^T \mathbf{R}_k \mathbf{S} \boldsymbol{\beta}_k) &= 0 \\ \implies -2\mathbf{S}^T \mathbf{R}_k \mathbf{Y} + 2\mathbf{S}^T \mathbf{R}_k \mathbf{S} \boldsymbol{\beta}_k &= 0 \\ \implies \mathbf{S}^T \mathbf{R}_k \mathbf{Y} &= \mathbf{S}^T \mathbf{R}_k \mathbf{S} \boldsymbol{\beta}_k \\ \implies \hat{\boldsymbol{\beta}}_k &= (\mathbf{S}^T \mathbf{R}_k \mathbf{S})^{-1} \mathbf{S}^T \mathbf{R}_k \mathbf{Y} = (\hat{b}_k, \hat{\mathbf{w}}_k^T)^T \end{aligned}$$

Last, to find σ_k^2 , we have to find

$$\begin{aligned} \frac{\partial \Lambda}{\partial \sigma_k^2} &= \frac{\partial Q}{\partial \sigma_k^2} = \sum_{i=1}^N r_{ik} \left(\frac{-1}{2\sigma_k^2} + \frac{(y_i - \hat{\mathbf{w}}_k^T \mathbf{x}_i - \hat{b}_k)^2}{2(\sigma_k^2)^2} \right) = 0 \\ \implies \hat{\sigma}_k^2 &= \frac{\sum_{i=1}^N r_{ik} (y_i - \hat{\mathbf{w}}_k^T \mathbf{x}_i - \hat{b}_k)^2}{\sum_{i=1}^N r_{ik}} \end{aligned}$$

□

Problem (e). Now let's put these ideas into practice. Generate the data as follows (or see `hw4p5.py`):

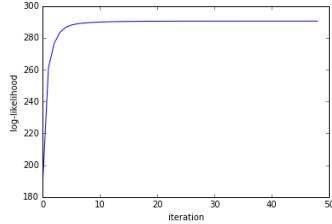
- Use $d = 1$ and $N = 500$. Let \underline{x} be sampled independently and uniformly from the interval $[0, 1]$.
- Use $\pi_1 = 0.7, \pi_2 = 0.3$.
- Use $w_1 = -2, w_2 = 1$.
- Use $b_1 = 0.5, b_2 = -0.5$.
- Use $\sigma_1 = 0.4, \sigma_2 = 0.3$. Note: This is σ , not σ^2 .
- Draw \underline{y} from the distribution using the above parameters and your already sampled \underline{x} .

Implement the EM algorithm using the updates you derived in (d), and estimate the model parameters, initializing your estimates with the following values:

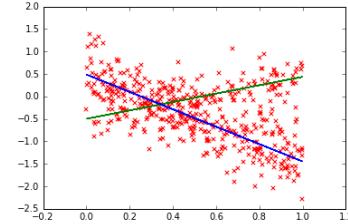
- Use $\hat{\pi}_1 = \hat{\pi}_2 = 0.5$.
- Use $\hat{w}_1 = 1, \hat{w}_2 = -1$.
- Use $\hat{b}_1 = \hat{b}_2 = 0$.
- Use $\hat{\sigma}_1 = \hat{\sigma}_2 = \text{std}(\underline{y})$. This is the standard deviation of your generated \underline{y} .

Deliverables:

- A plot of the log-likelihood as a function of iteration number. Terminate your algorithm when the log-likelihood increases by less than 10^{-4} .
- The estimated model parameters.
- A plot showing the data and estimated lines together.
- Please submit your code, as usual.



(a) Log-Likelihood



(b) Estimated Line

Proof. Estimated Parameter:

$$\begin{aligned}\pi &= [0.29104249836745466, 0.70895750163254534] \\ b &= [-0.49839639559032389, 0.48679655659431775] \\ w &= [0.93434248214276361, -1.9333629149516742] \\ \sigma &= [0.29089784122763501, 0.40147041897116265]\end{aligned}$$

□

6) Code Appendix

Problem 2) Dirichlet Maximum Likelihood (f). Code

```
from __future__ import division
import numpy as np
from scipy.special import polygamma, gammaln
from matplotlib import pyplot as plt
# Generate the data according to the specification in the homework description

N = 1000
m = 5
alpha = np.array([10, 5, 15, 20, 50])
P = np.random.dirichlet(alpha, N)
P_f=P
print np.shape(P)
# TODO: Implement the Newton-Raphson algorithm for estimating the parameters of
# the Dirichlet distribution given observances (rows of P).
LogData=np.log(P_f)
SuffStat=np.mean(LogData, axis=0)

def LogLike(alphaIn):
```

```

    sumAlpha=np.sum(alphaIn)
    output=N*(gammaln(sumAlpha)-np.sum(gammaln(alphaIn))+np.dot((alphaIn-1),Suf
    return output

def Gradient(alphaIn):
    sumAlpha=np.sum(alphaIn)
    gradient=N*(polygamma(0,sumAlpha)-polygamma(0,alphaIn)+SuffStat)
    return gradient

def QMatrix(alphaIn):
    Q_Matrix=np.diag(-N*polygamma(1,alphaIn))
    return Q_Matrix

def RealNumC(alphaIn):
    ConstC=N*polygamma(1,np.sum(alphaIn))
    return ConstC

def Update(alphaIn):
    ConstC=RealNumC(alphaIn)
    oneV=np.ones((m,m))
    QInv=np.diag((1.0/QMatrix(alphaIn)).diagonal())
    GV=Gradient(alphaIn)
    output=np.dot(QInv,GV)-1.0/(1.0/ConstC+np.sum(QInv.diagonal()))*np.dot(np.c
    return output

alphaInit=np.array([1,1,1,1,1])
logLike_lst=[]
True_lst=[]
TrueLogLike=LogLike(alpha)
alphaOld=alphaInit
logLike_lst.append(LogLike(alphaOld))
True_lst.append(TrueLogLike)
alphaNew=alphaOld-Update(alphaOld)
logLike_lst.append(LogLike(alphaNew))
True_lst.append(TrueLogLike)
count=1

while abs(LogLike(alphaNew)-LogLike(alphaOld))>10**(-4):
    alphaOld=alphaNew
    alphaNew=alphaOld-Update(alphaOld)
    print LogLike(alphaNew), LogLike(alpha)
    print alphaNew
    logLike_lst.append(LogLike(alphaNew))
    True_lst.append(TrueLogLike)
    count+=1
else:
    print alphaNew

```

```

TrueLogLike=LogLike(alpha)

count_lst=np.arange(count+1)
plt.plot(count_lst,logLike_lst,'-',label="MLE")
plt.plot(count_lst,True_lst,'r-',label="True")
plt.xlabel("iteration")
plt.ylabel("log-likelihood")
plt.legend()
plt.show()

```

Problem 4) Clustering (a). Code

```

from __future__ import division
from scipy.ndimage import imread
import numpy as np
from matplotlib import pyplot as plt
from numpy import linalg as LA

# Load the mandrill image as an NxNx3 array. Values range from 0.0 to 255.0.
mandrill = imread('mandrill.png', mode='RGB').astype(float)
N = int(mandrill.shape[0])

M = 2
k = 64

# Store each MxM block of the image as a row vector of X
X = np.zeros((N**2//M**2, 3*M**2))
for i in range(N//M):
    for j in range(N//M):
        X[i*N//M+j,:] = mandrill[i*M:(i+1)*M,j*M:(j+1)*M,:].reshape(3*M**2)

# TODO: Implement k-means and cluster the rows of X, then reconstruct the
# compressed image using the cluster center for each block, as specified in
# the homework description.

# To show a color image using matplotlib, you have to restrict the color
# color intensity values to between 0.0 and 1.0. For example,

def AssignLabel(allDataLabel,allCentroid):
    sampleNum=len(allDataLabel)
    labelIndex=allDataLabel.shape[1]-1
    dataVector=allDataLabel[:,0:labelIndex]
    for i in range(sampleNum):
        idx = (LA.norm(dataVector[i]-allCentroid, axis=1)).argmin()
        allDataLabel[i][labelIndex]=idx

```

```

        return allDataLabel

def getCentroid(allDatawithLabel ,numOfGr):
    labelIndex=np.shape(allDatawithLabel)[1]-1
    dataCentroid=[]
    for i in range(numOfGr):
        group_i=allDatawithLabel [allDatawithLabel[:,labelIndex]==i]
        centroid_i=np.mean(group_i[:,0:labelIndex],axis=0)
        dataCentroid.append(centroid_i.tolist())
    dataArrayCentroid=np.array(dataCentroid)
    return dataArrayCentroid

def initStep(data ,numOfGr):
    init_cent_label = np.random.choice(data.shape[0] ,numOfGr ,replace=False)
    centroid_init=data[init_cent_label ,:]
    zero=np.zeros(data.shape[0])
    datawithLabel=np.column_stack((data ,zero))
    return datawithLabel , centroid_init

def checkTerminal(centroidOld ,centroidNew):
    cenOldsort=np.sort(centroidOld)
    cenNewsort=np.sort(centroidNew)
    return np.all(cenOldsort==cenNewsort)

def genClustered_data(dataLabel ,centroid):
    labelIndex=dataLabel.shape[1]-1
    label=dataLabel[:,labelIndex]
    clusterData=[]
    for i in range(len(label)):
        cIdx=int(label[i])
        dataTobeApp=centroid[cIdx].tolist()
        dataTobeApp.append(cIdx)
        clusterData.append(dataTobeApp)
    clusterData=np.array(clusterData)
    return clusterData

def obj_fn(orgD ,cluD):
    value=np.sum((orgD-cluD)**2)
    return value

iteNum=0
obj_list=[]
dataLabel_old ,centroid_old=initStep(X ,k)
dataNewLabel=AssignLabel(dataLabel_old ,centroid_old)
centroid_New=getCentroid(dataNewLabel ,k)
cluD=genClustered_data(dataNewLabel ,centroid_New)
obj_value=obj_fn(dataNewLabel ,cluD)

```

```

obj_list.append(obj_value)

while not checkTerminal(centroid_old, centroid_New):
    iteNum+=1
    centroid_old=centroid_New
    dataLabel_old=dataNewLabel
    dataNewLabel=AssignLabel(dataLabel_old,centroid_old)
    centroid_New=getCentroid(dataNewLabel,k)
    cluD=genClustered_data(dataNewLabel,centroid_New)
    obj_value=obj_fn(dataNewLabel,cluD)
    obj_list.append(obj_value)

print centroid_New

ite=np.arange(iteNum+1)
plt.plot(ite,obj_list,'-')
plt.xlabel("iteration")
plt.ylabel("objective function")
plt.legend()
plt.show()

cluD_mandrill=cluD[:,0:12]
comp_mandrill = np.zeros(mandrill.shape)
for i in range(N//M):
    for j in range(N//M):
        comp_mandrill[i*M:(i+1)*M,j*M:(j+1)*M,:,:] = cluD_mandrill[i*N//M+j,:,:]

plt.imshow(comp_mandrill/255)
plt.show()

plt.imshow(mandrill/255)
plt.show()

plt.imshow((comp_mandrill - mandrill + 128)/255)
plt.show()

CompR = np.log2(k)/(24*M**2)+k*M**2/N**2
print(CompR)

RMAE = np.sum(abs(comp_mandrill-mandrill))*(1/(3*N**2))/255
print(RMAE)

```

Problem 5) EM Algorithm for Mixed Linear Regression (e). Code

```

from __future__ import division
import numpy as np
from matplotlib import pyplot as plt
from scipy.stats import norm

```

```

# Generate the data according to the specification in the homework description

N = 500
x = np.random.rand(N)

pi0 = np.array([0.7, 0.3])
w0 = np.array([-2, 1])
b0 = np.array([0.5, -0.5])
sigma0 = np.array([.4, .3])

y = np.zeros_like(x)
for i in range(N):
    k = 0 if np.random.rand() < pi0[0] else 1
    y[i] = w0[k]*x[i] + b0[k] + np.random.randn()*sigma0[k]
std_y=np.std(y,ddof=1)

# TODO: Implement the EM algorithm for Mixed Linear Regression based on observed
# x and y values.

# Here's the data plotted
plt.scatter(x, y, c='r', marker='x')
plt.show()

def LogLike(y_D,x_D,pi_old,w_old, b_old, sigma_old):
    output=0
    for n in range(N):
        for i in range(len(pi_old)):
            output+=pi_old[i]*norm.pdf(y_D[n],w_old[i]*x_D[n]+b_old[i],sigma_old[i])
    return output

def r_ik(x_i,y_i,pi_old,w_old, b_old, sigma_old):
    numer=pi_old[0]*norm.pdf(y_i,w_old[0]*x_i+b_old[0],sigma_old[0])
    deno=0
    for i in range(len(pi_old)):
        deno+=pi_old[i]*norm.pdf(y_i,w_old[i]*x_i+b_old[i],sigma_old[i])
    r_i0=numer/deno
    r_i1=1-r_i0
    output=np.vstack((r_i0,r_i1))
    return output

def pi_update(r_ikOld):
    piNew=[]
    for k in range(len(pi0)):
        piNew.append(np.sum(r_ikOld[k]/500))
    return piNew

```

```

def beta_update(x_D,y_D,r_ikOld):
    bNew=[]
    wNew=[]
    One=np.ones(N)
    S_D=np.vstack((One,x)).T
    for k in range(2):
        R_k=np.diag(r_ikOld[k])
        beta_k=np.dot(np.linalg.pinv(np.dot(np.dot(S_D.T,R_k),S_D)),np.dot(S_D.T,r_ikOld[k]))
        bNew.append(beta_k[0])
        wNew.append(beta_k[1])
    return bNew, wNew

def sigma_update(x_D,y_D,w_old,b_old,r_ikOld):
    sigmaNew=[]
    for k in range(2):
        deno=np.sum(r_ikOld[k])
        nume=np.sum(r_ikOld[k]*(y_D-w_old[k]*x_D-b_old[k])**2)
        sigmaNew.append((nume/deno)**0.5)
    return sigmaNew

pi_init=np.array([0.5,0.5])
w_init=np.array([1,-1])
b_init=np.array([0,0])
sigma_init=np.array([std_y,std_y])
iteCount=0
LogLike_lst=[]
piOld=pi_init
wOld=w_init
bOld=b_init
sigmaOld=sigma_init
LogLike_lst.append(LogLike(y,x,piOld,wOld, bOld, sigmaOld))
r_ikOld=r_ik(x,y,piOld,wOld,bOld,sigmaOld)

piNew=pi_update(r_ikOld)
bNew, wNew=beta_update(x,y,r_ikOld)
sigmaNew=sigma_update(x,y,wNew,bNew,r_ikOld)
iteCount+=1
LogLike_lst.append(LogLike(y,x,piNew,wNew, bNew, sigmaNew))

while abs(LogLike(y,x,piOld,wOld, bOld, sigmaOld)-LogLike(y,x,piNew,wNew, bNew, sigmaNew))>0.001:
    piOld=piNew
    wOld=wNew
    bOld=bNew
    sigmaOld=sigmaNew
    r_ikOld=r_ik(x,y,piOld,wOld,bOld,sigmaOld)
    piNew=pi_update(r_ikOld)
    bNew, wNew=beta_update(x,y,r_ikOld)

```

```

sigmaNew=sigma_update(x,y,wNew,bNew,r_ikOld)
LogLike_lst.append(LogLike(y,x,piNew,wNew, bNew, sigmaNew))
iteCount+=1

print piNew
print bNew
print wNew
print sigmaNew

count_lst=np.arange(iteCount+1)
plt.plot(count_lst,LogLike_lst,'-')
plt.xlabel("iteration")
plt.ylabel("log-likelihood")
plt.legend()
plt.show()

Line1 = wNew[0]*x+bNew[0]
Line2 = wNew[1]*x+bNew[1]
plt.scatter(x, y, c='r', marker='x')
plt.plot(x,Line1,'g-')
plt.plot(x,Line2,'b-')
plt.show()

```