

$$\begin{aligned}
 \text{la, } L(\theta) &= \theta_1^{\sum_{i=1}^N x_{1i}^{(1)}} \theta_2^{\sum_{i=1}^N x_{2i}^{(1)}} \dots \theta_k^{\sum_{i=1}^N x_{ki}^{(1)}} \\
 &= \theta_1^{N_1} \cdot \theta_2^{N_2} \dots \theta_k^{N_k} \\
 &= \prod_{k=1}^{k-1} \theta_k^{N_k} \cdot \left(1 - \sum_{k=1}^{k-1} \theta_k\right)^{N_k}
 \end{aligned}$$

$$\begin{aligned}
 \text{Then } l(\theta) &= \sum_{k=1}^{k-1} \log \theta_k^{N_k} + \log \left(1 - \sum_{i=1}^{k-1} \theta_i\right)^{N_k} \\
 &= \sum_{k=1}^{k-1} N_k \log \theta_k + N_k \log \left(1 - \sum_{i=1}^{k-1} \theta_i\right) \\
 &= \sum_{k=1}^{k-1} N_k \log \theta_k + \left(N - \sum_{k=1}^{k-1} N_k\right) \log \left(1 - \sum_{i=1}^{k-1} \theta_i\right)
 \end{aligned}$$

$$\text{Let } \frac{\partial}{\partial \theta_k} l(\theta) = \frac{N_k}{\theta_k} - \left(N - \sum_{k=1}^{k-1} N_k\right) \frac{1}{1 - \sum_{i=1}^{k-1} \theta_i} = 0$$

$$\begin{aligned}
 \text{Then } \theta_k &= \frac{N_k \left(1 - \sum_{i=1}^{k-1} \theta_i\right)}{\left(N - \sum_{k=1}^{k-1} N_k\right)} \\
 &= \frac{N_k}{N} \cdot \frac{\left(1 - \sum_{i=1}^{k-1} \theta_i\right)}{\left(1 - \sum_{i=1}^{k-1} N_i/N\right)} \\
 &= \frac{N_k}{N} \cdot \frac{\theta_k}{N_k/N} \quad \text{for all } k \in \{1, \dots, k-1\} \\
 &= \frac{N_k \theta_k}{N_k}
 \end{aligned}$$

$$\text{Since } \sum_{i=1}^k \theta_i = 1$$

$$\text{Then } \frac{\theta_k}{N_k} \sum_{i=1}^{k-1} N_i + \theta_k = 1$$

$$\theta_k \left(\frac{\sum_{i=1}^{k-1} N_i}{N_k} + 1 \right) = 1$$

$$\theta_k \frac{N_k + \sum_{i=1}^{k-1} N_i}{N_k} = 1$$

$$\theta_k \frac{N}{N_k} = 1$$

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

$$\text{So for } k \in \{1, \dots, k-1\}, \theta_k = \frac{N_k \frac{N_k}{N}}{N_k} = \frac{N_k}{N}$$

$$\text{So } \forall k \in \{1, \dots, k\}, \theta_k = \frac{N_k}{N}$$

$$\begin{aligned}
 1b) \quad p(\theta | D) &\propto p(\theta) p(D | \theta) \\
 &\propto \prod_{i=1}^k \theta_i^{a_i-1} \prod_{i=1}^k \theta_i^{\sum_{j=1}^N x_i^{(j)}} \\
 &\propto \prod_{i=1}^k \theta_i^{a_i-1} \prod_{i=1}^k \theta_i^{N_i} \\
 &= \prod_{i=1}^k \theta_i^{a_i + N_i - 1}
 \end{aligned}$$

$$\hookrightarrow p(\theta | D) \sim \text{Dirichlet}(a_1 + N_1, \dots, a_k + N_k)$$

1c) For $i \in \{1, \dots, k\}$

$$\begin{aligned}
 \hat{\theta}_{i, \text{map}} &= \underset{\theta_i}{\text{argmax}} p(\theta | D) \\
 &= \underset{\theta_i}{\text{argmax}} p(\theta) + p(D | \theta)
 \end{aligned}$$

$$\text{Since } p(\theta) \propto \theta_1^{a_1-1} \dots \theta_k^{a_k-1}$$

$$p(D | \theta) = \prod_{k=1}^k \theta_k^{N_k}$$

$$\begin{aligned}
 \hookrightarrow \hat{\theta}_{i, \text{map}} &= \underset{\theta_i}{\text{argmax}} \prod_{k=1}^k \theta_k^{a_k + N_k - 1} \\
 &= \underset{\theta_i}{\text{argmax}} \sum_{i=1}^k (a_i + N_i - 1) \log \theta_i
 \end{aligned}$$

$$\text{Set } \theta_k = 1 - \sum_{z=1}^{k-1} \theta_z$$

$$\hat{\theta}_{i, \text{map}} = \underset{\theta_i}{\text{argmax}} \sum_{j=1}^{k-1} (a_j + N_j - 1) \log \theta_j + (a_k + N_k - 1) \log \left(1 - \sum_{z=1}^{k-1} \theta_z \right)$$

$$\text{Set } \frac{d}{d\theta_i} \sum_{j=1}^{k-1} (a_j + N_j - 1) \log \theta_j + (a_k + N_k - 1) \log \left(1 - \sum_{z=1}^{k-1} \theta_z \right) = 0$$

$$a_i + N_i - 1 / \theta_i - a_k + N_k - 1 / \theta_k = 0$$

$$\Rightarrow \theta_i = \frac{a_i + N_i - 1}{a_k + N_k - 1} \cdot \theta_k$$

$$\text{Since } \sum_{j=1}^{k-1} \theta_j + \theta_k = 1$$

$$\begin{aligned}
 \hookrightarrow \frac{\sum_{j=1}^{k-1} a_j + \sum_{j=1}^{k-1} N_j - k + 1}{a_k + N_k - 1} \theta_k + \theta_k &= 1 \\
 \theta_k &= \frac{a_k + N_k - 1}{\sum_{j=1}^k a_j + N - k}
 \end{aligned}$$

$$\text{Substitute } \theta_k \text{ into } \theta_i, \quad \theta_i = \frac{a_i + N_i - 1}{\sum_{j=1}^k a_j + N - k}$$

$$\hookrightarrow \forall i \in \{1, \dots, k\}$$

$$\hat{\theta}_{i, \text{map}} = \frac{N_i + a_i}{N - k + \sum_{j=1}^k a_j}$$

$$1 d) \theta_{k|red} = P(x_k^{N+1} = 1 | D)$$

$$= \int P(\theta | D) P(x_k^{N+1} = 1 | \theta) d\theta$$

$$= \int \text{Dirichlet}(a_1 + N_1, \dots, a_k + N_k) \prod_{i=1}^k \theta_i^{x_i^{N+1}} d\theta$$

$$= \int \text{Dirichlet}(a_1 + N_1, \dots, a_k + N_k) \theta_k d\theta$$

$$= E(\theta_k) = \frac{\frac{a_k + N_k}{\sum_{j=1}^k a_j + N_j}}{\sum_{j=1}^k \frac{a_j + N_j}{\sum_{j=1}^k a_j + N_j}} = \frac{a_k + N_k}{\sum_{j=1}^k a_j + N}$$

$$\text{so } E[\theta_k | D] = \frac{N_k + a_k}{N + \sum_{j=1}^k a_j}$$

$$\begin{aligned}
 2a) \quad p(x) &= \sum_{i=1}^k p(x|T=i) p(T=i) \\
 &= \sum_{i=1}^k \left(\prod_{d=1}^D 2\pi\sigma_d^2 \right)^{-1/2} \exp \left\{ -\sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d - \mu_{id})^2 \right\} \cdot \alpha_i \\
 &= \left(\prod_{d=1}^D 2\pi\sigma_d^2 \right)^{-1/2} \sum_{i=1}^k \exp \left\{ -\sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d - \mu_{id})^2 \right\} \cdot \alpha_i
 \end{aligned}$$

$$\begin{aligned}
 p(T=k|x) &= p(x|T=k) \cdot p(T=k) / p(x) \\
 &= \frac{\exp \left\{ -\sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d - \mu_{kd})^2 \right\} \cdot \alpha_k}{\sum_{i=1}^k \exp \left\{ -\sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d - \mu_{id})^2 \right\} \cdot \alpha_i}
 \end{aligned}$$

$$\begin{aligned}
 2b) \quad l(\theta) &= \sum_{i=1}^N \log p(T^{(i)}, x^{(i)}) \\
 &= \sum_{i=1}^N \log p(T^{(i)}) p(x^{(i)}|T^{(i)}) \\
 &= \sum_{i=1}^N \log p(T^{(i)}) + \log p(x^{(i)}|T^{(i)}) \\
 &= \sum_{i=1}^N \sum_{k=1}^k \log p(T^{(i)}=k) \mathbb{1}(T^{(i)}=k) + \log p(x^{(i)}|T^{(i)}=k) \mathbb{1}(T^{(i)}=k) \\
 &= \sum_{i=1}^N \sum_{k=1}^k \left(\log \alpha_k + \log \left(\prod_{d=1}^D 2\pi\sigma_d^2 \right)^{-1/2} \exp \left\{ -\sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{kd})^2 \right\} \right) \mathbb{1}(T^{(i)}=k) \\
 &= \sum_{i=1}^N \sum_{k=1}^k \left(\log \alpha_k - \frac{1}{2} \sum_{d=1}^D \log 2\pi\sigma_d^2 - \sum_{d=1}^D \frac{1}{2\sigma_d^2} (x_d^{(i)} - \mu_{kd})^2 \right) \mathbb{1}(T^{(i)}=k)
 \end{aligned}$$

$$2C) \quad \frac{\partial}{\partial \sigma_d^2} l(\theta) = \sum_{i=1}^N \sum_{k=1}^K \left(-\frac{1}{2} \cdot \frac{2\sigma_d^2}{2\sigma_d^4} + \frac{1}{2\sigma_d^4} (x_d^{(i)} - \mu_{kd})^2 \right) \mathbb{1}(T^{(i)} = k)$$

$$\text{Let } \frac{\partial}{\partial \sigma_d^2} l(\theta) = 0$$

$$\text{Then } \frac{1}{\sigma_d^4} \sum_{i=1}^N \sum_{k=1}^K \frac{1}{2} \mathbb{1}(T^{(i)} = k) = \frac{1}{2\sigma_d^4} \sum_{i=1}^N \sum_{k=1}^K (x_d^{(i)} - \mu_{kd})^2 \mathbb{1}(T^{(i)} = k)$$

$$\sigma_d^2 = \frac{\sum_{i=1}^N \sum_{k=1}^K (x_d^{(i)} - \mu_{kd})^2 \mathbb{1}(T^{(i)} = k)}{\sum_{i=1}^N \sum_{k=1}^K \mathbb{1}(T^{(i)} = k)}$$

$$\sigma_d^2 = \frac{\sum_{i=1}^N \sum_{k=1}^K (x_d^{(i)} - \mu_{kd})^2 \mathbb{1}(T^{(i)} = k)}{N}$$

$$\frac{\partial}{\partial \mu_{kd}} l(\theta) = \sum_{i=1}^N \frac{1}{\sigma_d^4} (x_d^{(i)} - \mu_{kd}) \mathbb{1}(T^{(i)} = k)$$

$$\text{Let } \frac{\partial}{\partial \mu_{kd}} l(\theta) = 0$$

$$\sum_{i=1}^N \frac{1}{\sigma_d^4} (x_d^{(i)} - \mu_{kd}) \mathbb{1}(T^{(i)} = k) = 0$$

$$\sum_{i=1}^N \left(\frac{1}{\sigma_d^4} \right) x_d^{(i)} \cdot \mathbb{1}(T^{(i)} = k) = \mu_{kd} \sum_{i=1}^N \left(\frac{1}{\sigma_d^4} \right) \mathbb{1}(T^{(i)} = k)$$

$$\sum_{i=1}^N x_d^{(i)} \mathbb{1}(T^{(i)} = k) = \mu_{kd} N_k$$

$$\mu_{kd} = \frac{\sum_{i=1}^N x_d^{(i)} \mathbb{1}(T^{(i)} = k)}{N_k}$$

$$\text{So } \mu_{kd} = \frac{\sum_{i=1}^N x_d^{(i)} \mathbb{1}(T^{(i)} = k)}{N_k}$$

$$\sigma_d^2 = \frac{\sum_{i=1}^N \sum_{k=1}^K (x_d^{(i)} - \mu_{kd})^2 \mathbb{1}(T^{(i)} = k)}{N}$$

3 a)

```
def compute_mean_mles(train_data, train_labels):  
    """  
    Compute the mean estimate for each digit class. You may iterate over  
    the possible digits (0 to 9), but otherwise make sure that your code  
    is vectorized.  
  
    Arguments  
        train_data: size N x 64 numpy array with the images  
        train_labels: size N numpy array with corresponding labels  
  
    Returns  
        means: size 10 x 64 numpy array with the ith row corresponding  
            to the mean estimate for digit class i  
    """  
    # Initialize array to store means  
    means = np.zeros((10, 64))  
    for t in range(10):  
        k = t # initialize the label  
        td = train_data.copy()  
        tl = train_labels.copy()  
        tl[train_labels != k] = 0  
        tl[train_labels == k] = 1  
        tl = np.matrix(tl)  
        nk = tl.sum() # find the total number of label = k in training set  
        m = (tl * td) / nk  
        means[k, :] = m  
    return means
```

3a)

```
def compute_sigma_mles(train_data, train_labels):  
    """  
    Compute the covariance estimate for each digit class. You may iterate over  
    the possible digits (0 to 9), but otherwise make sure that your code  
    is vectorized.  
  
    Arguments  
        train_data: size N x 64 numpy array with the images  
        train_labels: size N numpy array with corresponding labels  
  
    Returns  
        covariances: size 10 x 64 x 64 numpy array with the ith row corresponding  
            to the covariance matrix estimate for digit class i  
    """  
    # Initialize array to store covariances  
    covariances = np.zeros((10, 64, 64))  
    # == YOUR CODE GOES HERE ==  
    for t in range(10):  
        k = t  
        mean_k = compute_mean_mles(train_data, train_labels)[k, :]  
        td = train_data.copy()  
        tl = train_labels.copy()  
        tl[train_labels != k] = 0  
        tl[train_labels == k] = 1  
        ind = np.where(tl == 1)[0] # Find the index of training label k  
        nk = ind.shape[0] # Find Nk  
        td = td[ind] # Find the data whose training label is k  
        c = np.dot((td - mean_k).T, (td - mean_k))/nk \  
            + 0.01 * np.identity(64)  
        covariances[k, :, :] = c  
    # ===  
    return covariances
```

3 a)

```
def generative_likelihood(digits, means, covariances):
    """
    Compute the generative log-likelihood  $\log p(x|t)$ . You may iterate over
    the possible digits (0 to 9), but otherwise make sure that your code
    is vectorized.

    Arguments
        digits: size N x 64 numpy array with the images
        means: size 10 x 64 numpy array with the 10 class means
        covariances: size 10 x 64 x 64 numpy array with the 10 class covariances

    Returns
        likelihoods: size N x 10 numpy array with the ith row corresponding
        to  $\log p(x^{(i)} | t)$  for  $t$  in {0, ..., 9}
    """
    N = digits.shape[0]
    likelihoods = np.zeros((N, 10))
    d = 64
    # == YOUR CODE GOES HERE ==
    for t in range(10):
        m = means[t, :]
        c = covariances[t, :, :]
        x = digits
        cinv = np.linalg.inv(c) # compute the inverse of covariance matrix
        p = -(d/2)*np.log(2*math.pi) - (1/2)*np.log(np.linalg.det(c)) \
            - 0.5 * np.diagonal(np.dot(np.dot(x-m, cinv), (x-m).T))
        likelihoods[:, t] = p
    # ===
    return likelihoods
```


3 a)

```
def conditional_likelihood(digits, means, covariances):  
    """  
    Compute the generative log-likelihood  $\log p(t|x)$ . Make sure that your code  
    is vectorized.  
  
    Arguments  
        digits: size N x 64 numpy array with the images  
        means: size 10 x 64 numpy array with the 10 class means  
        covariances: size 10 x 64 x 64 numpy array with the 10 class covariances  
  
    Returns  
        likelihoods: size N x 10 numpy array with the ith row corresponding  
        to  $\log p(t | x^{(i)})$  for  $t$  in  $\{0, \dots, 9\}$   
    """  
    # == YOUR CODE GOES HERE ==  
    N = digits.shape[0]  
    likelihoods = np.zeros((N, 10))  
    px = np.sum(np.exp(generative_likelihood(digits, means, covariances))  
                * (1/10), axis=1) # Find px using the law of total probability  
    pk = 1/10  
    px_t = generative_likelihood(digits, means, covariances)  
    for t in range(10):  
        p = px_t[:, t] - np.log(px) + np.log(pk)  
        likelihoods[:, t] = p  
    return likelihoods  
    # ==
```

```
def classify_data(digits, means, covariances):  
    """  
    Classify new points by taking the most likely posterior class.  
    Make sure that your code is vectorized.  
  
    Arguments  
        digits: size N x 64 numpy array with the images  
        means: size 10 x 64 numpy array with the 10 class means  
        covariances: size 10 x 64 x 64 numpy array with the 10 class covariances  
  
    Returns  
        pred: size N numpy array with the ith element corresponding  
        to  $\operatorname{argmax}_t \log p(t | x^{(i)})$   
    """  
    # Compute and return the most likely class  
    log = conditional_likelihood(digits, means, covariances)  
    c = np.argmax(log, axis=1) # Find the label that maximizes the likelihood  
    return c  
    # == YOUR CODE GOES HERE ==  
    # ==
```

3b)

```
Train average conditional log-likelihood: -0.12462443666863028
Test average conditional log-likelihood: -0.19667320325525578
Train posterior accuracy: 0.9814285714285714
Test posterior accuracy: 0.97275
```