**Problem 1:** Consider a generative classification model for $C$ classes defined by class probabilities $p(y = c) = \pi_c$ and general class-conditional densities $p(x \mid y = c, \theta_c)$ where $x \in \mathbb{R}^D$ is the input feature vector and $\theta = \{\theta_c\}_{c=1}^C$ are further model parameters. Suppose we are given a training set $\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$ where $y^{(n)}$ is a binary target vector of length $C$ that uses the 1-of-$C$ (one-hot) encoding scheme, so that it has components $y_c^{(n)} = \delta_{ck}$ if pattern $n$ is from class $y = k$. Assuming that the data points are i.i.d., show that the maximum-likelihood solution for the class probabilities $\pi$ is given by

$$\pi_c = \frac{N_c}{N}$$

where $N_c$ is the number of data points assigned to class $c$.

---

The data likelihood given the parameters $\{\pi_c, \theta_c\}_{c=1}^C$ is

$$p(\mathcal{D} \mid \{\pi_c, \theta_c\}_{c=1}^C) = \prod_{n=1}^N \prod_{c=1}^C (p(x^{(n)} \mid \theta_c)\pi_c)^{y_c^{(n)}}$$

and so the data log-likelihood is given by

$$\log p(\mathcal{D} \mid \{\pi_c, \theta_c\}_{c=1}^C) = \boxed{\sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \pi_c} + \text{const w.r.t. } \pi_c$$

In order to maximize the log likelihood with respect to $\pi_c$ we need to preserve the constraint $\sum_c \pi_c = 1$. For this we use the method of Lagrange multipliers where we introduce $\lambda$ as an unconstrained additional parameter and find a local extremum of the unconstrained function

$$\sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \pi_c - \lambda \left( \sum_{c=1}^C \pi_c - 1 \right).$$

instead. See [wikipedia article on Lagrange multipliers](#) for an intuition of why this works. This function is a sum of concave terms in $\pi_c$ as well as $\lambda$ and is therefore itself concave in these variables. We can find the extremum by finding the root of the derivatives. Setting the derivative with respect to $\pi_c$ equal to zero, we obtain

$$\pi_c = \frac{1}{\lambda} \sum_{n=1}^N y_c^{(n)} = \frac{N_c}{\lambda}.$$

Setting the derivative with respect to $\lambda$ equal to zero, we obtain the original constraint

$$\sum_{c=1}^C \pi_c = 1$$

where we can now plug in the previous result $\pi_c = \frac{N_c}{\lambda}$ and obtain $\lambda = \sum_c N_c = N$. Plugging this in turn into the expression for $\pi_c$ we obtain

$$\pi_c = \frac{N_c}{N}$$

which we wanted to show.

---

**Problem 2:** Using the same classification model as in the previous question, now suppose that the class-conditional densities are given by Gaussian distributions with a *shared* covariance matrix, so that

$$p(x \mid y = c, \theta) = p(x \mid \theta_c) = \mathcal{N}(x \mid \mu_c, \Sigma).$$

Show that the maximum likelihood estimate for the mean of the Gaussian distribution for class $c$ is given by

$$\mu_c = \frac{1}{N_c} \sum_{\substack{n=1 \\ y^{(n)}=c}}^N x^{(n)}$$

which represents the mean of the observations assigned to class $c$.

Similarly, show that the maximum likelihood estimate for the shared covariance matrix is given by

$$\Sigma = \sum_{c=1}^C \frac{N_c}{N} S_c \quad \text{where} \quad S_c = \frac{1}{N_c} \sum_{\substack{n=1 \\ y^{(n)}=c}}^N (x^{(n)} - \mu_c)(x^{(n)} - \mu_c)^{\mathrm{T}}.$$

Thus $\Sigma$ is given by a weighted average of the sample covariances of the data associated with each class, in which the weighting coefficients $N_c/N$ are the prior probabilities of the classes.

---

We begin by writing out the data log-likelihood.

$$\log p(\mathcal{D} \mid \{\pi_c, \theta_c\}_{c=1}^C)$$

$$= \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \pi_c \cdot p(x^{(n)} \mid \mu_c, \Sigma)$$

Then we plug in the definition of the multivariate Gaussian

$$= \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \left( (2\pi)^{-\frac{D}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp\left( -\frac{1}{2}(x^{(n)} - \mu_c)^{\mathrm{T}}\Sigma^{-1}(x^{(n)} - \mu_c) \right) \right) + y^{(n)} \log \pi_c$$

and simplify.

$$= -\frac{1}{2} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \left( D \log 2\pi + \log \det(\Sigma) + (x^{(n)} - \mu_c)^{\mathrm{T}}\Sigma^{-1}(x^{(n)} - \mu_c) - 2\log \pi_c \right)$$

1-hit

$C = 5 \rightarrow$ $1 \rightarrow 0\,0\,0\,0\,1$
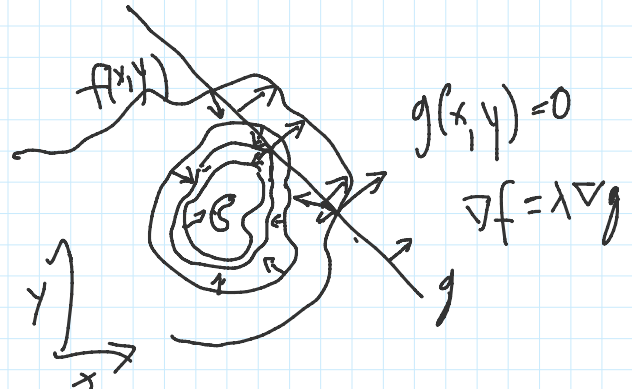$2 \rightarrow 0\,0\,0\,1\,0$
$3 \rightarrow 0\,0\,1\,0\,0$
$\vdots$

$\delta_{ij} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$   $x^0 = 1$

$p(D \mid \pi, \theta) = \prod_n p(x^n, y^n)$

$= \prod_n p(x^n \mid y_n, \theta) \cdot p(y_n \mid \pi)$

$= \prod_n \prod_c \left[ p(x_n \mid y=c, \theta) \cdot p(y=c \mid \pi) \right]^{y_c^{(n)}}$

$g(x, y) = 0$

$\nabla f = \lambda \nabla g$

$-\log \tilde{x} = \log x^{-1}$

This expression is concave in $\mu_c$, so we can obtain the maximizer by finding the root of the derivative. With the help of the matrix cookbook, we identify the derivative with respect to $\mu_c$ as

$$\frac{1}{2}\sum_{n=1}^{N} y_c^{(n)}\Sigma^{-1}(x^{(n)} - \mu_c)$$

which we can set to 0 and solve for $\mu_c$ to obtain

$$\mu_c = \frac{1}{\sum_{n=1}^{N} y_c^{(n)}}\sum_{n=1}^{N} y_c^{(n)}x^{(n)} = \frac{1}{N_c}\sum_{\substack{n=1 \\ y^{(n)}=c}}^{N} x^{(n)}.$$

To find the optimal $\Sigma$, we need the trace trick

$$a = \mathrm{Tr}(a) \text{ for all } a \in \mathbb{R} \quad \text{and} \quad \mathrm{Tr}(ABC) = \mathrm{Tr}(BCA).$$

With this we can rewrite

$$(x^{(n)} - \mu_c)^{\mathrm{T}}\Sigma^{-1}(x^{(n)} - \mu_c) = \mathrm{Tr}\left(\overbrace{\Sigma^{-1}}^{A}\overbrace{(x^{(n)} - \mu_c)(x^{(n)} - \mu_c)^{\mathrm{T}}}^{B}\right)$$

and use the matrix-trace derivative rule $\frac{\partial}{\partial A}\mathrm{Tr}(AB) = B^{\mathrm{T}}$ to find the derivative of the data log-likelihood with respect to $\Sigma$. Because the log-likelihood contains both $\Sigma$ and $\Sigma^{-1}$, we convert one into the other with $\log\det A = -\log\det A^{-1}$ to obtain

$$-\frac{1}{2}\sum_{n=1}^{N}\sum_{c=1}^{C} y_c^{(n)}\left(-\log\det\Sigma^{-1} + \mathrm{Tr}\left(\Sigma^{-1}(x^{(n)} - \mu_c)(x^{(n)} - \mu_c)^{\mathrm{T}}\right)\right) + \text{const w.r.t. } \Sigma.$$

Finally, we use rule (57) from the matrix cookbook $\frac{\partial\log|\det X|}{\partial X} = (X^{-1})^{\mathrm{T}}$ and compute the derivative of the log-likelihood with respect to $\Sigma^{-1}$ as

$$-\frac{1}{2}\sum_{n=1}^{N}\sum_{c=1}^{C} y_c^{(n)}\left(-\Sigma^{\mathrm{T}} + (x^{(n)} - \mu_c)(x^{(n)} - \mu_c)^{\mathrm{T}}\right).$$

We find the root with respect to $\Sigma$ and find

$$\Sigma = \frac{1}{\sum_{n=1}^{N}\sum_{c=1}^{C} y_c^{(n)}}\left(\sum_{n=1}^{N}\sum_{c=1}^{C} y_c^{(n)}(x^{(n)} - \mu_c)(x^{(n)} - \mu_c)^{\mathrm{T}}\right)^{\mathrm{T}} = \frac{1}{N}\sum_{c=1}^{C}\sum_{\substack{n=1 \\ y^{(n)}=c}}^{N} (x^{(n)} - \mu_c)(x^{(n)} - \mu_c)^{\mathrm{T}}$$

which we can immediately break apart into the representation in the instructions.