# ISTA 421/521 - Final Take-home Assignment

**Due: Monday, December 15, 5pm**

24 pts total undergraduate / 30 pts total graduate

STUDENT NAME

Undergraduate / Graduate

## Instructions

You must work on the problems **INDEPENDENTLY**, not in groups.

**The work on each problem must be your own.**

Include in your final submission the pdf of written answers along with separate files for any python scripts that you write in support of answering the questions (although no scripts are needed for this assignment; if you do write scripts, clearly note in your pdf written answers which script filenames were used). You are required to create either a .zip or tarball (.tar.gz / .tgz) archive of all of the files for your submission and submit the archive to the D2L dropbox by the date/time deadline above.

NOTE: Problem **??** is **required for graduate students only**; undergraduates may complete them for extra credit equal to the point value.

(FCMA refers to the course text: Rogers and Girolami (2012), *A First Course in Machine Learning*.)

1. [4 points] Adapted from **Exercise 5.4** of FCMA p.204:

   Compute the maximum likelihood estimates of $q_{mc}$ for class $c$ of a Bayesian classifier with multinomial class-conditionals and a set of $N_c$, $M$-dimensional objects belonging to class $c$: $\mathbf{x}_1, ..., \mathbf{x}_{N_c}$.

   **Solution.** <Solution goes here>

2. [4 points] Adapted from **Exercise 5.5** of FCMA p.204:

   For a Bayesian classifier with multinomial class-conditionals with $M$-dimensional parameters $\mathbf{q}_c$, compute the posterior Dirichlet for class $c$ when the prior over $\mathbf{q}_c$ is a Dirichlet with constant parameter $\alpha$ and the observations belonging to class $c$ are the $N_c$ observations $\mathbf{x}_1, ..., \mathbf{x}_{N_c}$.

   **Solution.** <Solution goes here>

3. [4 points] Adapted from **Exercise 6.1** of FCMA p.234:

   Derive the EM update for the variance of the $d$th dimension and the $k$th component, $\sigma_{kd}^2$, when the cluster components have a diagonal Gaussian Likelihood:

   $$p(\mathbf{x}_n | z_{nk} = 1, \mu_{k1}, ..., \mu_{KD}, \sigma_{k1}^2, ..., \sigma_{kD}^2) = \prod_{d=1}^{D} \mathcal{N}(\mu_{kd}, \sigma_{kd}^2)$$

   **Solution.** <Solution goes here>

4. [6 points; <span style="color:red">**Required only for Graduates**</span>] Adapted from **Exercise 6.6** of FCMA p.235:

   Derive an EM algorithm for fitting a mixture of Poisson distributions. Assume you observe $N$ integer counts, $x_1, ..., x_N$. The likelihood is:

   $$p(\mathbf{X}|\mathbf{\Delta}) = \prod_{n=1}^{N} \sum_{k=1}^{K} \pi_k \frac{\lambda_k^{x_n} \exp\{-\lambda_k\}}{x_n!}$$

   **Solution.** <Solution goes here>

5. [2 points]

   For a support vector machine, if we remove one of the support vectors from the training set, does the size of the maximum margin decrease, stay the same, or increase for that dataset? Why? Also justify your answer by providing a simple, hand-designed dataset (no more than 2-D) in which you identify the support vectors, draw the location of the maximum margin hyperplane, remove one of the support vectors, and draw the location of the resulting maximum margin hyperplane. You do not have to run any code – this can be done completely by hand and drawn schematically.

   **Solution.** <Solution goes here>

6. [3 points]

   Consider the 2-bit XOR problem for which the entire instance space is as follows: In each row, $x_1$

   | $t$ | $x_1$ | $x_2$ |
   |-----|-------|-------|
   | $-1$ | $-1$ | $-1$ |
   | $+1$ | $-1$ | $1$ |
   | $+1$ | $1$ | $-1$ |
   | $-1$ | $1$ | $1$ |

   and $x_2$ are the coordinates and $t$ is the class for the point. These instances are not linearly separable,

but they are separable with a polynomial kernel. Recall that the polynomial kernel is of the form $\kappa(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i \top \mathbf{x}_j + c)^d$ where $c$ and $d$ are integers. Select (by hand!) values for $c$ and $d$ that yield a space in which the instances above are linearly separable. Write down the mapping $\Phi$ to which this kernel corresponds, write down $\Phi(x)$ for each instance above, and write down the parameters of a hyperplane in the expanded space that perfectly classifies the instances (there are a range of possible hyperplanes, just pick one set of hyperplane parameters that satisfies separating the points – you are not maximizing the margin here, just coming up with one possible separating hyperplane). Again, this can be done without writing code or deriving an analytic solution!

**Solution.** <Solution goes here>

7. [2 points]

Why does the kernel trick allow us to solve SVMs with high dimensional feature spaces without significantly increasing the running time?

**Solution.** <Solution goes here>

8. [5 points] In this course we introduced the Metropolis-Hastings algorithm. Provide the following: (a) describe the problem it is designed to solve, including how it solves this problem by avoiding a potentially intractable problem; (b) describe the basic procedure; (c) describe the role of the proposal distribution.