# APTS Statistical Computing:
# Practical Lab 2 (Thursday)

1. (a) For this exercise, we will use the `FLtools` package from:

    https://bitbucket.org/finnlindgren/FLtools/

    developed by Finn Lindgren (the previous lecturer for this course). You can install it with:

    ```
    remotes::install_bitbucket("finnlindgren/FLtools")
    ```

    If you don't have the `remotes` package, first install it with:

    ```
    install.packages("remotes")
    ```

    Once you have installed the `FLtools` package, you should be able to load it with

    ```
    library(FLtools)
    ```

    Make sure you have this package installed before proceeding to the next step.

    (b) Start the `optimisation` shiny app:

    ```
    FLtools::optimisation()
    ```

    This should start a Shiny web application. It will also attempt to start up a tab in your browser connected to the session. If this doesn't work, just connect your browser to the URL of the Shiny app. Make sure the Shiny app is running in a browser window before proceeding to the next step.

    (c) For the "Simple (1D)" and "Simple (2D)" functions, familiarise yourself with the "Step", "Converge", and "Reset" buttons.

    (d) Choose different optimisation starting points by clicking in the figure.

    (e) Explore the different optimisation methods and what they display in the figure for each optimisation step[1][2][3]. Also observe the diagnostic output box and how the number of function, gradient, and Hessian evaluations differ between the methods.

    (f) For the "Rosenbrock (2D)" function, observe the differences in convergence behaviour for the four different optimisation methods.

    (g) For the "Multimodal" functions, explore how the optimisation methods behave for different starting points.

    (h) How far out can the optimisation start for the "Spiral" function? E.g., try the "Newton" method, starting in the top right corner of the figure.

2. Write your own code to optimise Rosenbrock's function by Newton's method. Ensure that you have implemented it correctly by comparing your output (and implementation) with that of the Shiny app from the first exercise. For this question you will want to make use of the preliminary material for the course (and the solutions).

---

[1]LS stands for "line search".

[2]The simplex/triangle shapes are shown for each "Simplex" method step in blue. The "best" points for each simplex are connected (magenta).

[3]The Newton methods display the true quadratic Taylor approximations (red) as well as the approximations used to find the proposed steps (blue).

3. Consider the linear mixed model for a response vector $\mathbf{y}$:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\epsilon}, \quad \mathbf{b} \sim N(\mathbf{0}, \mathbf{I}\sigma_b^2), \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{I}\sigma^2)$$

$\mathbf{X}$ and $\mathbf{Z}$ are (fixed) model matrices, $\boldsymbol{\beta}$, $\sigma_b^2$ and $\sigma^2$ are parameters, and $\mathbf{b}$ and $\boldsymbol{\epsilon}$ are independent.

(a) First simulate some data from a model of this sort, taking care to relate the code back to the mathematical statement of the model. . .

```
set.seed(10)
n <- 100;n.b <- 10;n.beta <- 5
## X and Z are fixed in the model, not random. Random numbers
## used only to generate arbitrary examples, here....
X <- cbind(1,matrix(runif(n*n.beta-n),n,n.beta-1))
Z <- matrix(runif(n*n.b),n,n.b)
beta <- rep(1,n.beta)
b <- rnorm(n.b)
y <- X%*%beta + Z%*%b + rnorm(n)
```

You'll use the data, $\mathbf{y}$, simulated here, along with the corresponding X and Z, to experiment with *fitting* linear mixed models (so from now on pretend that you don't know what values $\boldsymbol{\beta}$, $\sigma_b$ and $\sigma$ had).

(b) With pencil and paper, find the (marginal) expectation, $\boldsymbol{\mu}$, and covariance matrix, $\mathbf{V}$, of $\mathbf{y}$. State the (marginal) distribution of $\mathbf{y}$. Here we marginalise over the random effects, $\mathbf{b}$, and the errors, but still condition on everything else.

(c) The following R function evaluates the log likelihood of $\boldsymbol{\theta}^\mathsf{T} = (\boldsymbol{\beta}^\mathsf{T}, \sigma_b^2, \sigma^2)$ given data $\mathbf{y}$. Note that $\boldsymbol{\theta}$ is the first argument of the function.

```
logLik <- function(theta,y,X,Z) {
## somewhat plodding linear mixed model log
## likelihood with theta partitioned
## [beta,sig2.b,sig2]
   n <- length(y)
   beta <- theta[1:ncol(X)]
   theta <- theta[-(1:ncol(X))]
   V <- diag(n)*theta[2] + Z %*% t(Z)*theta[1]
   R <- chol(V)
   z <- forwardsolve(t(R), y-X %*% beta)
   ll <- -n*log(2*pi)/2 - sum(log(diag(R))) - sum(z*z)/2
   ll
}
```

To maximise the log likelihood of the model using unconstrained methods, it is better to use a parameterization that guarantees positive variances.[4] Modify the function to accept a parameter vector $\boldsymbol{\theta}^\mathsf{T} = (\boldsymbol{\beta}^\mathsf{T}, \rho_b, \rho)$ where $\rho = \log(\sigma)$ and $\rho_b = \log(\sigma_b)$.

(d) Use `optim` to maximise your likelihood (note that `optim` *minimizes* by default; see the documentation for how to do maximisation, and for how to choose optimisation method).

---

[4]Such reparameterisation can often have the added benefit of leading posterior distributions closer to Gaussian, enabling accurate and precise Bayesian approximations not relying on Monte Carlo simulations.

(e) In fact, using general purpose optimisation methods to find the optimising $\beta$ is a bit wasteful. Given the variance parameters, closed form expressions for the $\beta$ maximising the likelihood are available, and might as well be used. Then it is only necessary to use general methods for the variance parameters. The likelihood considered only as a function of the variance parameters, with the corresponding MLEs of $\beta$ 'plugged in' is termed a 'profile likelihood'. Show that, given the variance parameters, the log-likelihood is maximised by the $\beta$ minimising

$$(\mathbf{y} - \mathbf{X}\beta)^\mathsf{T}\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\beta) = \|\mathbf{R}^{-\mathsf{T}}(\mathbf{y} - \mathbf{X}\beta)\|^2$$

where $\mathbf{R}^\mathsf{T}\mathbf{R} = \mathbf{V}$. ($\|\mathbf{x}\|^2 = \mathbf{x}^\mathsf{T}\mathbf{x}$ here.) Hence, produce a 'profile log likelihood' function equivalent to your previous log likelihood function. Your function should accept a vector of variance parameters as its first argument, and should return the corresponding profile log likelihood value. You might want to return the corresponding $\beta$ values as an attribute of the return value, e.g.

```
  .

  .
  attr(ll,"beta") <- beta
  ll
}
```

(f) Use `optim` to maximise your profiled log likelihood function (or copy the code from the online solutions!), and confirm that you get near identical parameter estimates to those from part (d).

4. Note: `base::chol()` always returns a matrix in dense storage format, so use `Matrix::chol()` instead, to obtain sparse storage output for sparse storage input.

(a) Run the following code that measures the time it takes to compute dense Cholesky factorisations for matrices of varying size (from $10$ to $1000$):

```r
An <- c(10, 20, 50, 100, 200, 500, 1000)
Atime <- c()
for (n in An) {
  ## Use several repeated runs for small matrices:
  loop.max <- max(1, 10000/n)
  ## Construct a random symmetric positive definite matrix:
  A <- matrix(rnorm(n^2), n, n); A <- t(A) %*% A
  ## Compute and time the Cholesky calculations.
  ## Use B <- A*1 to make sure R doesn't use any hidden
  ## precomputations
  Atime <- rbind(Atime,
                 system.time({
                   for (loop in 1:loop.max) {
                     B <- A*1
                     Matrix::chol(B)
                 }}) / loop.max)
}
Atime
plot(An, Atime[,1], log="")
plot(An, Atime[,1], type="l", log="xy")
```

(b) Adapt the code to measure the time it takes to compute dense Cholesky factorisations of covariance matrices of an AR$(1)$ process, for size $n = 10$ to $1000$. Note that you can create such a matrix (for standard deviation `sd` and auto-regressive parameter `a` $\in [0, 1)$) with something like:

```r
S <- as.matrix(dist(0:(n - 1)))
S <- sd^2 * a^S
```

Choosing `sd=10` and `a=0.9` should be fine, but feel free to explore alternatives.

(c) Adapt the code to measure the time it takes to compute sparse Cholesky factorisations of precision matrices of an AR$(1)$ process. Let $n$ vary between $10$ and $10^6$. Note that you can create such a matrix with something like:

```r
Q <- Matrix::sparseMatrix(i = c(1:n, 2:n, 1:(n - 1)),
          j = c(1:n, 1:(n - 1), 2:n), x = rep(c(1, 1 + a^2,
            1, -a), c(1, n - 2, 1, 2 * (n - 1)))/(1 - a^2)/sd^2,
          dims = c(n, n))
```

(d) Graphically compare the computational costs of dense and sparse Cholesky factorisations.