

# Unit 11: Numerical linear algebra

November 5, 2014

References:

- Gentle: Numerical Linear Algebra for Applications in Statistics (my notes here are based primarily on this source) [Gentle-NLA]
  - Unfortunately, this is not in the UCB library system - I have a copy that you can take a look at.
- Gentle: Computational Statistics [Gentle-CS]
- Lange: Numerical Analysis for Statisticians
- Monahan: Numerical Methods of Statistics

In working through how to compute something or understanding an algorithm, it can be very helpful to depict the matrices and vectors graphically. We'll see this on the board in class.

## 1 Preliminaries

### 1.1 Goals

Here's what I'd like you to get out of this unit:

1. How to think about the computational order (number of computations involved) of a problem
2. How to choose a computational approach to a given linear algebra calculation you need to do.
3. An understanding of how issues with computer numbers (Unit 7) play out in terms of linear algebra.

## 1.2 Key principle

**The form of a mathematical expression and how it should be evaluated on a computer may be very different.** Better computational approaches can increase speed and improve the numerical properties of the calculation.

Example 1: We do not compute  $(X^T X)^{-1} X^T Y$  by computing  $X^T X$  and finding its inverse. In fact, perhaps more surprisingly, we may never actually form  $X^T X$  in some implementations.

Example 2: Suppose I have a matrix  $A$ , and I want to permute (switch) two rows. I can do this with a permutation matrix,  $P$ , which is mostly zeroes. On a computer, in general I wouldn't need to even change the values of  $A$  in memory in some cases (e.g., if I were to calculate  $PAB$ ). Why not?

## 1.3 Computational complexity

We can assess the computational complexity of a linear algebra calculation by counting the number multiplies/divides and the number of adds/subtracts. Sidenote: addition is a bit faster than multiplication, so some algorithms attempt to trade multiplication for addition.

In general we do not try to count the actual number of calculations, but just their order, though in some cases in this unit we'll actually get a more exact count. In general, we denote this as  $O(f(n))$  which means that the number of calculations approaches  $cf(n)$  as  $n \rightarrow \infty$  (i.e., we know the calculation is approximately proportional to  $f(n)$ ). Consider matrix multiplication,  $AB$ , with matrices of size  $a \times b$  and  $b \times c$ . Each column of the second matrix is multiplied by all the rows of the first. For any given inner product of a row by a column, we have  $b$  multiplies. We repeat these operations for each column and then for each row, so we have  $abc$  multiplies so  $O(abc)$  operations. We could count the additions as well, but there's usually an addition for each multiply, so we can usually just count the multiplies and then say there are such and such {multiply and add}s. This is Monahan's approach, but you may see other counting approaches where one counts the multiplies and the adds separately.

For two symmetric,  $n \times n$  matrices, this is  $O(n^3)$ . Similarly, matrix factorization (e.g., the Cholesky decomposition) is  $O(n^3)$  unless the matrix has special structure, such as being sparse. As matrices get large, the speed of calculations decreases drastically because of the scaling as  $n^3$  and memory use increases drastically. In terms of memory use, to hold the result of the multiply indicated above, we need to hold  $ab + bc + ac$  total elements, which for symmetric matrices sums to  $3n^2$ . So for a matrix with  $n = 10000$ , we have  $3 \cdot 10000^2 \cdot 8/1e9 = 2.4\text{Gb}$ .

When we have  $O(n^q)$  this is known as polynomial time. Much worse is  $O(b^n)$  (exponential time), while much better is  $O(\log n)$  (log time). Computer scientists talk about NP-complete problems; these are essentially problems for which there is not a polynomial time algorithm - it turns

out all such problems can be rewritten such that they are equivalent to one another.

In real calculations, it's possible to have the actual time ordering of two approaches differ from what the order approximations tell us. For example, something that involves  $n^2$  operations may be faster than one that involves  $1000(n \log n + n)$  even though the former is  $O(n^2)$  and the latter  $O(n \log n)$ . The problem is that the constant,  $c = 1000$ , can matter (depending on how big  $n$  is), as can the extra calculations from the lower order term(s), in this case  $1000n$ .

A note on terminology: *flops* stands for both floating point operations (the number of operations required) and floating point operations per second, the speed of calculation.

## 1.4 Notation and dimensions

I'll try to use capital letters for matrices,  $A$ , and lower-case for vectors,  $x$ . Then  $x_i$  is the  $i$ th element of  $x$ ,  $A_{ij}$  is the  $i$ th row,  $j$ th column element, and  $A_{.j}$  is the  $j$ th column and  $A_{i.}$  the  $i$ th row. By default, we'll consider a vector,  $x$ , to be a one-column matrix, and  $x^\top$  to be a one-row matrix. Some of the textbook resources also use  $a_{ij}$  for  $A_{ij}$  and  $a_j$  for the  $j$ th column.

Throughout, we'll need to be careful that the matrices involved in an operation are conformable: for  $A + B$  both matrices need to be of the same dimension, while for  $AB$  the number of columns of  $A$  must match the number of rows of  $B$ . Note that this allows for  $B$  to be a column vector, with only one column,  $Ab$ . Just checking dimensions is a good way to catch many errors. Example: is  $\text{Cov}(Ax) = A\text{Cov}(x)A^\top$  or  $\text{Cov}(Ax) = A^\top\text{Cov}(x)A$ ? Well, if  $A$  is  $m \times n$ , it must be the former, as the latter is not conformable.

The inner product of two vectors is  $\sum_i x_i y_i = x^\top y \equiv \langle x, y \rangle \equiv x \cdot y$ . The outer product is  $xy^\top$ , which comes from all pairwise products of the elements.

When the indices of summation should be obvious, I'll sometimes leave them implicit. Ask me if it's not clear.

## 1.5 Norms

$\|x\|_p = (\sum_i |x_i|^p)^{1/p}$  and the standard (Euclidean) norm is  $\|x\|_2 = \sqrt{\sum x_i^2} = \sqrt{x^\top x}$ , just the length of the vector in Euclidean space, which we'll refer to as  $\|x\|$ , unless noted otherwise. The standard norm for a matrix is the Frobenius norm,  $\|A\|_F = (\sum_{i,j} a_{ij}^2)^{1/2}$ . There is also the induced matrix norm, corresponding to any chosen vector norm,

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|}$$

So we have

$$\|A\|_2 = \sup_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \sup_{\|x\|_2=1} \|Ax\|_2$$

A property of any legitimate matrix norm (including the induced norm) is that  $\|AB\| \leq \|A\|\|B\|$ . Recall that norms must obey the triangle inequality,  $\|A + B\| \leq \|A\| + \|B\|$ .

A normalized vector is one with “length”, i.e., Euclidean norm, of one. We can easily normalize a vector:  $\tilde{x} = x/\|x\|$

The angle between two vectors is

$$\theta = \cos^{-1} \left( \frac{\langle x, y \rangle}{\sqrt{\langle x, x \rangle \langle y, y \rangle}} \right)$$

## 1.6 Orthogonality

Two vectors are orthogonal if  $x^\top y = 0$ , in which case we say  $x \perp y$ . An orthogonal matrix is a matrix in which all of the columns are orthogonal to each other and normalized. Orthogonal matrices can be shown to have full rank. Furthermore if  $A$  is orthogonal,  $A^\top A = I$ , so  $A^{-1} = A^\top$ . Given all this, the determinant of orthogonal  $A$  is either 1 or -1. Finally the product of two orthogonal matrices,  $A$  and  $B$ , is also orthogonal since  $(AB)^\top AB = B^\top A^\top AB = B^\top B = I$ .

**Permutations** Sometimes we make use of matrices that permute two rows (or two columns) of another matrix when multiplied. Such a matrix is known as an elementary permutation matrix and is an orthogonal matrix with a determinant of -1. You can multiply such matrices to get more general permutation matrices that are also orthogonal. If you premultiply by  $P$ , you permute rows, and if you postmultiply by  $P$  you permute columns. Note that on a computer, you wouldn’t need to actually do the multiply (and if you did, you should use a sparse matrix routine), but rather one can often just rework index values that indicate where relevant pieces of the matrix are stored (more in the next section).

## 1.7 Some vector and matrix properties

$AB \neq BA$  but  $A + B = B + A$  and  $A(BC) = (AB)C$ .

In R, recall the syntax is

```
A + B
A %*% B
```

You don’t need the spaces, but they’re nice for code readability.

## 1.8 Trace and determinant of square matrices

The trace of a matrix is the sum of the diagonal elements. For square matrices,  $\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B)$ ,  $\text{tr}(A) = \text{tr}(A^\top)$ .

We also have  $\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$  - basically you can move a matrix from the beginning to the end or end to beginning, provided they are conformable for this operation. This is helpful for a couple reasons:

1. We can find the ordering that reduces computation the most if the individual matrices are not square.
2.  $x^\top Ax = \text{tr}(x^\top Ax)$  since the quadratic form,  $x^\top Ax$ , is a scalar, and this is equal to  $\text{tr}(xx^\top A)$  where  $xx^\top A$  is a matrix. It can be helpful to be able to go back and forth between a scalar and a trace in some statistical calculations.

For square matrices, the determinant exists and we have  $|AB| = |A||B|$  and therefore,  $|A^{-1}| = 1/|A|$  since  $|I| = |AA^{-1}| = 1$ . Also  $|A| = |A^\top|$ .

**Other matrix multiplications** The Hadamard or direct product is simply multiplication of the corresponding elements of two matrices by each other. In R this is simply `A * B`.

**Challenge:** How can I find  $\text{tr}(AB)$  without using `A %*% B` ?

The Kronecker product is the product of each element of one matrix with the entire other matrix”

$$A \otimes B = \begin{pmatrix} A_{11}B & \cdots & A_{1m}B \\ \vdots & \ddots & \vdots \\ A_{n1}B & \cdots & A_{nm}B \end{pmatrix}$$

The inverse of a Kronecker product is the Kronecker product of the inverses,

$$B^{-1} \otimes A^{-1}$$

which is obviously quite a bit faster because the inverse (i.e., solving a system of equations) in this special case is  $O(n^3 + m^3)$  rather than the naive approach being  $O((nm)^3)$ .

## 1.9 Linear independence, rank, and basis vectors

A set of vectors,  $v_1, \dots, v_n$ , is linearly independent (LIN) when none of the vectors can be represented as a linear combination,  $\sum c_i v_i$ , of the others for scalars,  $c_1, \dots, c_n$ . If we have vectors of

length  $n$ , we can have at most  $n$  linearly independent vectors. The rank of a matrix is the number of linearly independent rows (or columns - it's the same), and is at most the minimum of the number of rows and number of columns. We'll generally think about it in terms of the dimension of the column space - so we can just think about the number of linearly independent columns.

Any set of linearly independent vectors (say  $v_1, \dots, v_n$ ) span a space made up of all linear combinations of those vectors ( $\sum_{i=1}^n c_i v_i$ ). The spanning vectors are known as basis vectors. We can express a vector  $x$  that is in the space with respect to (as a linear combination of) basis vectors as  $x = \sum_i c_i v_i$ , where if the basis vectors are normalized and orthogonal, we can find the weights as  $c_i = \langle x, v_i \rangle$ .

Consider a regression context. We have  $p$  covariates ( $p$  columns in the design matrix,  $X$ ), of which  $q$  are linearly independent covariates. This means that  $p - q$  of the vectors can be written as linear combos of the  $q$  vectors. The space spanned by the covariate vectors is of dimension  $q$ , rather than  $p$ , and  $X^\top X$  has  $p - q$  eigenvalues that are zero. The  $q$  LIN vectors are basis vectors for the space - we can represent any point in the space as a linear combination of the basis vectors. You can think of the basis vectors as being like the axes of the space, except that the basis vectors are not orthogonal. So it's like denoting a point in  $\mathbb{R}^q$  as a set of  $q$  numbers telling us where on each of the axes we are - this is the same as a linear combination of axis-oriented vectors. When we have  $n \leq q$ , a vector of  $n$  observations can be represented exactly as a linear combination of the  $q$  basis vectors, so there is no residual. If  $n = q$ , then we have a single unique solution, while if  $n < q$  we have multiple possible solutions and the system is ill-determined (under-determined). Of course we usually have  $n > q$ , so the system is overdetermined - there is no exact solution, but regression is all about finding solutions that minimize some criterion about the differences between the observations and linear combinations of the columns of the  $X$  matrix (such as least squares or penalized least squares). In standard regression, we project the observation vector onto the space spanned by the columns of the  $X$  matrix, so we find the point in the space closest to the observation vector.

## 1.10 Invertibility, singularity, rank, and positive definiteness

For square matrices, let's consider how invertibility, singularity, rank and positive (or non-negative) definiteness relate.

Square matrices that are "regular" have an eigendecomposition,  $A = \Gamma \Lambda \Gamma^{-1}$  where  $\Gamma$  is a matrix with the eigenvectors as the columns and  $\Lambda$  is a diagonal matrix of eigenvalues,  $\Lambda_{ii} = \lambda_i$ . Symmetric matrices and matrices with unique eigenvalues are regular, as are some other matrices. The number of non-zero eigenvalues is the same as the rank of the matrix. Square matrices that have an inverse are also called nonsingular, and this is equivalent to having full rank. If the matrix is

symmetric, the eigenvectors and eigenvalues are real and  $\Gamma$  is orthogonal, so we have  $A = \Gamma\Lambda\Gamma^\top$ . The determinant of the matrix is the product of the eigenvalues (why?), which is zero if it is less than full rank. Note that if none of the eigenvalues are zero then  $A^{-1} = \Gamma\Lambda^{-1}\Gamma^\top$ .

Let's focus on symmetric matrices. The symmetric matrices that tend to arise in statistics are either positive definite (p.d.) or non-negative definite (n.n.d.). If a matrix is positive definite, then by definition  $x^\top Ax > 0$  for any  $x$ . Note that if  $\text{Cov}(y) = A$  then  $x^\top Ax = x^\top \text{Cov}(y)x = \text{Cov}(x^\top y) = \text{Var}(x^\top y)$  if so positive definiteness amounts to having linear combinations of random variables having positive variance. So we must have that positive definite matrices are equivalent to variance-covariance matrices (I'll just refer to this as a variance matrix or as a covariance matrix). If  $A$  is p.d. then it has all positive eigenvalues and it must have an inverse, though as we'll see, from a numerical perspective, we may not be able to compute it if some of the eigenvalues are very close to zero. In R, `eigen(A)$vectors` is  $\Gamma$ , with each column a vector, and `eigen(A)$values` contains the ordered eigenvalues.

Let's interpret the eigendecomposition in a generative context as a way of generating random vectors. We can generate  $y$  s.t.  $\text{Cov}(y) = A$  if we generate  $y = \Gamma\Lambda^{1/2}z$  where  $\text{Cov}(z) = I$  and  $\Lambda^{1/2}$  is formed by taking the square roots of the eigenvalues. So  $\sqrt{\lambda_i}$  is the standard deviation associated with the basis vector  $\Gamma_{\cdot i}$ . That is, the  $z$ 's provide the weights on the basis vectors, with scaling based on the eigenvalues. So  $y$  is produced as a linear combination of eigenvectors as basis vectors, with the variance attributable to the basis vectors determined by the eigenvalues.

If  $x^\top Ax \geq 0$  then  $A$  is nonnegative definite (also called positive semi-definite). In this case one or more eigenvalues can be zero. Let's interpret this a bit more in the context of generating random vectors based on non-negative definite matrices,  $y = \Gamma\Lambda^{1/2}z$  where  $\text{Cov}(z) = I$ . Questions:

1. What does it mean when one or more eigenvalue (i.e.,  $\lambda_i = \Lambda_{ii}$ ) is zero?
2. Suppose I have an eigenvalue that is very small and I set it to zero? What will be the impact upon  $y$  and  $\text{Cov}(y)$ ?
3. Now let's consider the inverse of a covariance matrix, known as the precision matrix,  $A^{-1} = \Gamma\Lambda^{-1}\Gamma^\top$ . What does it mean if a  $(\Lambda^{-1})_{ii}$  is very large? What if  $(\Lambda^{-1})_{ii}$  is very small?

Consider an arbitrary  $n \times p$  matrix,  $X$ . Any crossproduct or sum of squares matrix, such as  $X^\top X$  is positive definite (non-negative definite if  $p > n$ ). This makes sense as it's just a scaling of an empirical covariance matrix.

## 1.11 Generalized inverses

Suppose I want to find  $x$  such that  $Ax = b$ . Mathematically the answer (provided  $A$  is invertible, i.e. of full rank) is  $x = A^{-1}b$ .

Generalized inverses arise in solving equations when  $A$  is not full rank. A generalized inverse is a matrix,  $A^-$  s.t.  $AA^-A = A$ . The Moore-Penrose inverse (the pseudo-inverse),  $A^+$ , is a (unique) generalized inverse that also satisfies some additional properties.  $x = A^+b$  is the solution to the linear system,  $Ax = b$ , that has the shortest length for  $x$ .

We can find the pseudo-inverse based on an eigendecomposition (or an SVD) as  $\Gamma\Lambda^+\Gamma^\top$ . We obtain  $\Lambda^+$  from  $\Lambda$  as follows. For values  $\lambda_i > 0$ , compute  $1/\lambda_i$ . All other values are set to 0. Let's interpret this statistically. Suppose we have a precision matrix with one or more zero eigenvalues and we want to find the covariance matrix. A zero eigenvalue means we have no precision, or infinite variance, for some linear combination (i.e., for some basis vector). We take the pseudo-inverse and assign that linear combination zero variance.

Let's consider a specific example. Autoregressive models are often used for smoothing (in time, in space, and in covariates). A first order autoregressive model for  $y_1, y_2, \dots, y_T$  has  $E(y_i|y_{-i}) = \frac{1}{2}(y_{i-1} + y_{i+1})$ . Another way of writing the model is in time-order:  $y_i = y_{i-1} + \epsilon_i$ . A second order autoregressive model has  $E(y_i|y_{-i}) = \frac{1}{6}(4y_{i-1} + 4y_{i+1} - y_{i-2} - y_{i+2})$ . These constructions basically state that each value should be a smoothed version of its neighbors. One can figure out that the **precision** matrix for  $y$  in the first order model is

$$\begin{pmatrix} \ddots & & & & \\ & -1 & 2 & -1 & 0 \\ & \cdots & -1 & 2 & -1 & \cdots \\ & & 0 & -1 & 2 & -1 \\ & & & \vdots & & \ddots \end{pmatrix}$$

and in the second order model is

$$\begin{pmatrix} \ddots & & & & & \\ & 1 & -4 & 6 & -4 & 1 \\ & \cdots & 1 & -4 & 6 & -4 & 1 & \cdots \\ & & & 1 & -4 & 6 & -4 & 1 \\ & & & & \vdots & & & \ddots \end{pmatrix}.$$

If we look at the eigendecomposition of such matrices, we see that in the first order case, the eigenvalue corresponding to the constant eigenvector is zero.

```
precMat <- matrix(c(1,-1,0,0,0,-1,2,-1,0,0,0,-1,2,-1,
  0,0,0,-1,2,-1,0,0,0,-1,1), 5)
e <- eigen(precMat)
```



```
e$values
## [1] 3.618034 2.618034 1.381966 0.381966 0.000000

e$vectors[, 5]
## [1] 0.4472136 0.4472136 0.4472136 0.4472136 0.4472136
```

This means we have no information about the overall level of  $y$ . So how would we generate sample  $y$  vectors? We can't put infinite variance on the constant basis vector and still generate samples. Instead we use the pseudo-inverse and assign ZERO variance to the constant basis vector. This corresponds to generating realizations under the constraint that  $\sum y_i$  has no variation, i.e.,  $\sum y_i = \bar{y} = 0$  - you can see this by seeing that  $\text{Var}(\Gamma_i^\top y) = 0$  when  $\lambda_i = 0$ .

```
# generate a realization
e$values[1:4] <- 1 / e$values[1:4]
y <- e$vec %*% (e$values * rnorm(5))
sum(y)

## [1] -2.553513e-15
```

In the second order case, we have two non-identifiabilities: for the sum and for the linear component of the variation in  $y$  (linear in the indices of  $y$ ).

I could parameterize a statistical model as  $\mu + y$  where  $y$  has covariance that is the generalized inverse discussed above. Then I allow for both a non-zero mean and for smooth variation governed by the autoregressive structure. In the second-order case, I would need to add a linear component as well, given the second non-identifiability.

## 1.12 Matrices arising in regression

In regression, we work with  $X^\top X$ . Some properties of this matrix are that it is symmetric and non-negative definite (hence our use of  $(X^\top X)^{-1}$  in the OLS estimator). When is it not positive definite?

Fitted values are  $X\hat{\beta} = X(X^\top X)^{-1}X^\top Y = HY$ . The “hat” matrix,  $H$ , projects  $Y$  into the column space of  $X$ .  $H$  is idempotent:  $HH = H$ , which makes sense - once you've projected into the space, any subsequent projection just gives you the same thing back.  $H$  is singular. Why? Also, under what special circumstance would it not be singular?

## 2 Computational issues

### 2.1 Storing matrices

We've discussed column-major and row-major storage of matrices. First, retrieval of matrix elements from memory is quickest when multiple elements are contiguous in memory. So in a column-major language (e.g., R, Fortran), it is best to work with values in a common column (or entire columns) while in a row-major language (e.g., C) for values in a common row.

In some cases, one can save space (and potentially speed) by overwriting the output from a matrix calculation into the space occupied by an input. This occurs in some clever implementations of matrix factorizations.

### 2.2 Algorithms

Good algorithms can change the efficiency of an algorithm by one or more orders of magnitude, and many of the improvements in computational speed over recent decades have been in algorithms rather than in computer speed.

Most matrix algebra calculations can be done in multiple ways. For example, we could compute  $b = Ax$  in either of the following ways, denoted here in pseudocode.

1. Stack the inner products of the rows of  $A$  with  $x$ .

```
for ( i = 1 : n ) {  
    b_i = 0  
    for ( j = 1 : m ) {  
        b_i = b_i + a_{ i j } x_j  
    }  
}
```

2. Take the linear combination (based on  $x$ ) of the columns of  $A$

```
for ( i = 1 : n ) {  
    b_i = 0  
}  
for ( j = 1 : m ) {  
    for ( i = 1 : n ) {  
        b_i = b_i + a_{ i j } x_j  
    }  
}
```

```

    }
  }
}

```

In this case the two approaches involve the same number of operations but the first might be better for row-major matrices (so might be how we would implement in C) and the second for column-major (so might be how we would implement in Fortran). **Challenge:** check whether the second approach is faster in R. (Write the code just doing the outer loop and doing the inner loop using vectorized calculation.)

**General computational issues** The same caveats we discussed in terms of computer arithmetic hold naturally for linear algebra, since this involves arithmetic with many elements. Good implementations of algorithms are aware of the danger of catastrophic cancellation and of the possibility of dividing by zero or by values that are near zero.

## 2.3 Ill-conditioned problems

**Basics** A problem is ill-conditioned if small changes to values in the computation result in large changes in the result. This is quantified by something called the *condition number* of a calculation. For different operations there are different condition numbers.

Ill-conditionedness arises most often in terms of matrix inversion, so the standard condition number is the “condition number with respect to inversion”, which when using the  $L_2$  norm is the ratio of the absolute values of the largest to smallest eigenvalue. Here’s an example:

$$A = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix}.$$

The solution of  $Ax = b$  for  $b = (32, 23, 33, 31)$  is  $x = (1, 1, 1, 1)$ , while the solution for  $b + \delta b = (32.1, 22.9, 33.1, 30.9)$  is  $x + \delta x = (9.2, -12.6, 4.5, -1.1)$ , where  $\delta$  is notation for a perturbation to the vector or matrix. What’s going on?

```

norm2 <- function(x) sqrt(sum(x^2))
A <- matrix(c(10, 7, 8, 7, 7, 5, 6, 5, 8, 6, 10, 9, 7, 5, 9, 10), 4)
e <- eigen(A)
b <- c(32, 23, 33, 31)
bPerturb <- c(32.1, 22.9, 33.1, 30.9)

```

```

x <- solve(A, b)
xPerturb <- solve(A, bPerturb)
norm2(x - xPerturb)

## [1] 16.39695

norm2(b - bPerturb)

## [1] 0.2

norm2(x - xPerturb)/norm2(x)

## [1] 8.198475

(e$val[1]/e$val[4])*norm2(b - bPerturb)/norm2(b)

## [1] 9.942834

```

Some manipulations with inequalities involving the induced matrix norm (for any chosen vector norm, but we might as well just think about the Euclidean norm) (see Gentle-CS Sec. 5.1) give

$$\frac{\|\delta x\|}{\|x\|} \leq \|A\| \|A^{-1}\| \frac{\|\delta b\|}{\|b\|}$$

where we define the condition number w.r.t. inversion as  $\text{cond}(A) \equiv \|A\| \|A^{-1}\|$ . We'll generally work with the  $L_2$  norm, and for a nonsingular square matrix the result is that the condition number is the ratio of the absolute values of the largest and smallest magnitude eigenvalues. This makes sense since  $\|A\|_2$  is the absolute value of the largest magnitude eigenvalue of  $A$  and  $\|A^{-1}\|_2$  that of the inverse of the absolute value of the smallest magnitude eigenvalue of  $A$ . We see in the code above that the large disparity in eigenvalues of  $A$  leads to an effect predictable from our inequality above, with the condition number helping us find an upper bound.

The main use of these ideas for our purposes is in thinking about the numerical accuracy of a linear system solution (Gentle-NLA Sec 3.4). On a computer we have the system

$$(A + \delta A)(x + \delta x) = b + \delta b$$

where the 'perturbation' is from the inaccuracy of computer numbers. Our exploration of computer numbers tells us that

$$\frac{\|\delta b\|}{\|b\|} \approx 10^{-p}; \quad \frac{\|\delta A\|}{\|A\|} \approx 10^{-p}$$

where  $p = 16$  for standard double precision floating points. Following Gentle, one gets the approximation

$$\frac{\|\delta x\|}{\|x\|} \approx \text{cond}(A)10^{-p},$$

so if  $\text{cond}(A) \approx 10^t$ , we have accuracy of order  $10^{t-p}$  instead of  $10^{-p}$ . (Gentle cautions that this holds only if  $10^{t-p} \ll 1$ ). So we can think of the condition number as giving us the number of digits of accuracy lost during a computation relative to the precision of numbers on the computer. E.g., a condition number of  $10^8$  means we lose 8 digits of accuracy relative to our original 16 on standard systems. One issue is that estimating the condition number is itself subject to numerical error and requires computation of  $A^{-1}$  (albeit not in the case of  $L_2$  norm with square, nonsingular  $A$ ) but see Golub and van Loan (1996; p. 76-78) for an algorithm.

**Improving conditioning** Ill-conditioned problems in statistics often arise from collinearity of regressors. Often the best solution is not a numerical one, but re-thinking the modeling approach, as this generally indicates statistical issues beyond just the numerical difficulties.

A general comment on improving conditioning is that we want to avoid large differences in the magnitudes of numbers involved in a calculation. In some contexts such as regression, we can center and scale the columns to avoid such differences - this will improve the condition of the problem. E.g., in simple quadratic regression with  $x = \{1990, \dots, 2010\}$  (e.g., regressing on calendar years), we see that centering and scaling the matrix columns makes a huge difference on the condition number

```
x1 <- 1990:2010
x2 <- x1 - 2000 # centered
x3 <- x2/10 # centered and scaled
X1 <- cbind(rep(1, 21), x1, x1^2)
X2 <- cbind(rep(1, 21), x2, x2^2)
X3 <- cbind(rep(1, 21), x3, x3^2)
e1 <- eigen(crossprod(X1))
e1$values

## [1] 3.360186e+14 7.699100e+02 -3.833498e-08

e2 <- eigen(crossprod(X2))
e2$values

## [1] 50677.704275 770.000000 9.295725
```

```
e3 <- eigen(crossprod(X3))
e3$values

## [1] 24.112935  7.700000  1.953665
```

The basic story is that simple strategies often solve the problem, and that you should be cognizant of the absolute and relative magnitudes involved in your calculations.

One rule of thumb is to try to work with numbers whose magnitude is around 1. We can often scale the values in our problem in order to do this. I.e., change the units of your variables. Instead of personal income in dollars, use personal income in thousands or hundreds of thousands of dollars.

### 3 Matrix factorizations (decompositions) and solving systems of linear equations

Suppose we want to solve the following linear system:

$$\begin{aligned} Ax &= b \\ x &= A^{-1}b \end{aligned}$$

Numerically, this is never done by finding the inverse and multiplying. Rather we solve the system using a matrix decomposition (or equivalent set of steps). One approach uses Gaussian elimination (equivalent to the LU decomposition), while another uses the Cholesky decomposition. There are also iterative methods that generate a sequence of approximations to the solution but reduce computation (provided they are stopped before the exact solution is found).

Gentle-CS has a nice table overviewing the various factorizations (Table 5.1, page 219).

#### 3.1 Triangular systems

As a preface, let's figure out how to solve  $Ax = b$  if  $A$  is upper triangular. The basic algorithm proceeds from the bottom up (and therefore is called a 'backsolve'. We solve for  $x_n$  trivially, and then move upwards plugging in the known values of  $x$  and solving for the remaining unknown in each row (each equation).

1.  $x_n = b_n / A_{nn}$

2. Now for  $k < n$ , use the already computed  $\{x_n, x_{n-1}, \dots, x_{k+1}\}$  to calculate  $x_k = \frac{b_k - \sum_{j=k+1}^n x_j A_{kj}}{A_{kk}}$ .

3. Repeat for all rows.

How many multiplies and adds are done? Solving lower triangular systems is very similar and involves the same number of calculations.

In R, *backsolve()* solves upper triangular systems and *forwardsolve()* solves lower triangular systems:

```
n <- 20
X <- crossprod(matrix(rnorm(n^2), n))
b <- rnorm(n)
U <- chol(crossprod(X)) # U is upper-triangular
L <- t(U) # L is lower-triangular
out1 <- backsolve(U, b)
out2 <- forwardsolve(L, b)
all.equal(out1, c(solve(U) %*% b))

## [1] TRUE

all.equal(out2, c(solve(L) %*% b))

## [1] TRUE
```

We can also solve  $(U^T)^{-1}b$  and  $(L^T)^{-1}b$  as

```
backsolve(U, b, transpose = TRUE)
forwardsolve(L, b, transpose = TRUE)
```

**To reiterate the distinction between matrix inversion and solving a system of equations, when we write  $U^{-1}b$ , what we mean on a computer is to carry out the above algorithm, not to find the inverse and then multiply.**

## 3.2 Gaussian elimination (LU decomposition)

Gaussian elimination is a standard way of directly computing a solution for  $Ax = b$ . It is equivalent to the LU decomposition. LU is primarily done with square matrices, but not always. Also LU decompositions do exist for some singular matrices.

The idea of Gaussian elimination is to convert the problem to a triangular system. In class, we'll walk through Gaussian elimination in detail and see how it relates to the LU decomposition. I'll describe it more briefly here. Following what we learned in algebra when we have multiple

equations, we preserve the solution,  $x$ , when we add multiples of rows (i.e., add multiples of equations) together. This amounts to doing  $L_1Ax = L_1b$  for a lower-triangular matrix  $L_1$  that produces all zeroes in the first column of  $L_1A$  except for the first row. We proceed to zero out values below the diagonal for the other columns of  $A$ . The result is  $L_{n-1} \cdots L_1A \equiv U = L_{n-1} \cdots L_1b \equiv b^*$  where  $U$  is upper triangular. This is the forward reduction step of Gaussian elimination. Then the backward elimination step solves  $Ux = b^*$ .

If we're just looking for the solution of the system, we don't need the lower-triangular factor  $L = (L_{n-1} \cdots L_1)^{-1}$  in  $A = LU$ , but it turns out to have a simple form that is computed as we go along, it is unit lower triangular and the values below the diagonal are the negative of the values below the diagonals in  $L_1, \dots, L_{n-1}$  (note that each  $L_j$  has non-zeroes below the diagonal only in the  $j$ th column). As a side note related to storage, it turns out that as we proceed, we can store the elements of  $L$  and  $U$  in the original  $A$  matrix, except for the implicit 1s on the diagonal of  $L$ .

In class, we'll work out the computational complexity of the LU and see that it is  $O(n^3)$ .

If we look at `solve.default()` in R, we see that it uses `dgesv`. A Google search indicates that this is a Lapack routine that does the LU decomposition with partial pivoting and row interchanges (see below on what these are), so R is using the algorithm we've just discussed.

One additional complexity is that we want to avoid dividing by very small values to avoid introducing numerical inaccuracy (we would get large values that might overwhelm whatever they are being added to, and small errors in the divisor will have large effects on the result). This can be done on the fly by interchanging equations to use the equation (row) that produces the largest value to divide by. For example in the first step, we would switch the first equation (first row) for whichever of the remaining equations has the largest value in the first column. This is called partial pivoting. The divisors are called pivots. Complete pivoting also considers interchanging columns, and while theoretically better, partial pivoting is generally sufficient and requires fewer computations. Note that partial pivoting can be expressed as multiplying along the way by permutation matrices,  $P_1, \dots, P_{n-1}$  that switch rows. Based on pivoting, we have  $PA = LU$ , where  $P = P_{n-1} \cdots P_1$ . In the demo code, we'll see a toy example of the impact of pivoting.

Finally  $|PA| = |P||A| = |L||U| = |U|$  (why?) so  $|A| = |U|/|P|$  and since the determinant of each permutation matrix,  $P_j$  is -1 (except when  $P_j = I$  because we don't need to switch rows), we just need to multiply by minus one if there is an odd number of permutations. Or if we know the matrix is non-negative definite, we just take the absolute value of  $|U|$ . So Gaussian elimination provides a fast stable way to find the determinant.



### 3.3 Cholesky decomposition

When  $A$  is p.d., we can use the Cholesky decomposition to solve a system of equations. Positive definite matrices can be decomposed as  $U^T U = A$  where  $U$  is upper triangular.  $U$  is called a square root matrix and is unique (apart from the sign, which we fix by requiring the diagonals to be positive). One algorithm for computing  $U$  is:

1.  $U_{11} = \sqrt{A_{11}}$
2. For  $j = 2, \dots, n$ ,  $U_{1j} = A_{1j}/U_{11}$
3. For  $i = 2, \dots, n$ ,
  - $U_{ii} = \sqrt{A_{ii} - \sum_{k=1}^{i-1} U_{ki}^2}$
  - for  $j = i + 1, \dots, n$ :  $U_{ij} = (A_{ij} - \sum_{k=1}^{i-1} U_{ki} U_{kj})/U_{ii}$

We can then solve a system of equations as:  $U^{-1}(U^T)^{-1}b$ , which in R can be done in either of the following ways:

```
backsolve(U, backsolve(U, b, transpose = TRUE))  
backsolve(U, forwardsolve(t(U), b)) # equivalent but less efficient
```

The Cholesky has some nice advantages over the LU: (1) while both are  $O(n^3)$ , the Cholesky involves only half as many computations,  $n^3/6 + O(n^2)$  and (2) the Cholesky factorization has only  $(n^2 + n)/2$  unique values compared to  $n^2 + n$  for the LU. Of course the LU is more broadly applicable. The Cholesky does require computation of square roots, but it turns out this is not too intensive. There is also a method for finding the Cholesky without square roots.

**Uses of the Cholesky** The standard algorithm for generating  $y \sim \mathcal{N}(0, A)$  is:

```
U <- chol(A)  
y <- crossprod(U, rnorm(n)) # i.e., t(U) %*% rnorm(n), but much faster
```

**Question:** where will most of the time in this two-step calculation be spent?

If a regression design matrix,  $X$ , is full rank, then  $X^T X$  is positive definite, so we could find  $\hat{\beta} = (X^T X)^{-1} X^T Y$  using either the Cholesky or Gaussian elimination. **Challenge:** write efficient R code to carry out the OLS solution using either LU or Cholesky factorization.

However, it turns out that the standard approach is to work with  $X$  using the QR decomposition rather than working with  $X^T X$ ; working with  $X$  is more numerically stable, though in most situations without extreme collinearity, either of the approaches will be fine.

**Numerical issues with eigendecompositions and Cholesky decompositions for positive definite matrices** Monahan comments that in general Gaussian elimination and the Cholesky decomposition are very stable. However, in the Cholesky case, if the matrix is very ill-conditioned we can get  $A_{ii} - \sum_k U_{ki}^2$  being negative and then the algorithm stops when we try to take the square root. In this case, the Cholesky decomposition does not exist numerically although it exists mathematically. It's not all that hard to produce such a matrix, particularly when working with high-dimensional covariance matrices with large correlations.

```
# require(fields)
locs <- runif(100)
rho <- .1
C <- exp(-rdist(locs)^2/rho^2)
e <- eigen(C)
e$values[96:100]

## [1] -1.905153e-16 -2.071019e-16 -2.445624e-16 -2.502047e-16 -5.247692e-16

U <- chol(C)

## Error in chol.default(C): the leading minor of order 31 is not positive
definite

vals <- abs(e$values)
max(vals)/min(vals)

## [1] 4.753109e+18

U <- chol(C, pivot = TRUE)

## Warning in chol.default(C, pivot = TRUE): the matrix is either rank-defi
or indefinite
```

We can think about the accuracy here as follows. Suppose we have a matrix whose diagonal elements (i.e., the variances) are order of magnitude 1 and that the true value of a  $U_{ii}$  is less than  $1 \times 10^{-16}$ . From the given  $A_{ii}$  we are subtracting  $\sum_k U_{ki}^2$  and trying to calculate this very small number but we know that we can only represent the values  $A_{ii}$  and  $\sum_k U_{ki}^2$  accurately to 16 places, so the difference is garbage starting in the 17th position and could well be negative. Now realize that  $\sum_k U_{ki}^2$  is the result of a potentially large set of arithmetic operations, and is likely represented accurately to fewer than 16 places. Now if the true value of  $U_{ii}$  is smaller than the accuracy to

which  $\sum_k U_{ki}^2$  is represented, we can get a difference that is negative.

Note that when the Cholesky fails, we can still compute an eigendecomposition, but we have negative numeric eigenvalues. Even if all the eigenvalues are numerically positive (or equivalently, we're able to get the Cholesky), errors in small eigenvalues near machine precision could have large effects when we work with the inverse of the matrix. This is what happens when we have columns of the  $X$  matrix nearly collinear. We cannot statistically distinguish the effect of two (or more) covariates, and this plays out numerically in terms of unstable results.

A strategy when working with mathematically but not numerically positive definite  $A$  is to set eigenvalues or singular values to zero when they get very small, which amounts to using a pseudo-inverse and setting to zero any linear combinations with very small variance. We can also use pivoting with the Cholesky and accumulate zeroes in the last  $n - q$  rows (for cases where we try to take the square root of a negative number), corresponding to the columns of  $A$  that are numerically linearly dependent. See the *pivot* argument to R's *chol()*.

## 3.4 QR decomposition

### 3.4.1 Introduction

The QR decomposition is available for any matrix,  $X = QR$ , with  $Q$  orthogonal and  $R$  upper triangular. If  $X$  is non-square,  $n \times p$  with  $n > p$  then the leading  $p$  rows of  $R$  provide an upper triangular matrix ( $R_1$ ) and the remaining rows are 0. (I'm using  $p$  because the QR is generally applied to design matrices in regression). In this case we really only need the first  $p$  columns of  $Q$ , and we have  $X = Q_1 R_1$ , the 'skinny' QR (this is what R's QR provides). For uniqueness, we can require the diagonals of  $R$  to be nonnegative, and then  $R$  will be the same as the upper-triangular Cholesky factor of  $X^\top X$ :

$$\begin{aligned} X^\top X &= R^\top Q^\top Q R \\ &= R^\top R \end{aligned}$$

There are three standard approaches for computing the QR, using (1) reflections (Householder transformations), (2) rotations (Givens transformations), or (3) Gram-Schmidt orthogonalization (see below for details).

For  $n \times n$   $X$ , the QR (for the Householder approach) requires  $2n^3/3$  flops, so QR is less efficient than LU or Cholesky.

We can also obtain the pseudo-inverse of  $X$  from the QR:  $X^+ = [R_1^{-1} \ 0]Q^\top$ . In the case that  $X$  is not full-rank, there is a version of the QR that will work (involving pivoting) and we end up with some additional zeroes on the diagonal of  $R_1$ .

### 3.4.2 Regression and the QR

Often QR is used to fit linear models, including in R. Consider the linear model in the form  $Y = X\beta + \epsilon$ , finding  $\hat{\beta} = (X^\top X)^{-1}X^\top Y$ . Let's consider the skinny QR and note that  $R^\top$  is invertible. Therefore, we can express the normal equations as

$$\begin{aligned}X^\top X\beta &= X^\top Y \\R^\top Q^\top QR\beta &= R^\top Q^\top Y \\R\beta &= Q^\top Y\end{aligned}$$

and solving for  $\beta$  is just a backsolve since  $R$  is upper-triangular. Furthermore the standard regression quantities, such as the hat matrix, the SSE, the residuals, etc. can be easily expressed in terms of  $Q$  and  $R$ .

Why use the QR instead of the Cholesky on  $X^\top X$ ? The condition number of  $X$  is the square root of that of  $X^\top X$ , and the  $QR$  factorizes  $X$ . Monahan has a discussion of the condition of the regression problem, but from a larger perspective, the situations where numerical accuracy is a concern are generally cases where the OLS estimators are not particularly helpful anyway (e.g., highly collinear predictors).

What about computational order of the different approaches to least squares? The Cholesky is  $np^2 + \frac{1}{3}p^3$ , an algorithm called sweeping is  $np^2 + p^3$ , the Householder method for QR is  $2np^2 - \frac{2}{3}p^3$ , and the modified Gram-Schmidt approach for QR is  $2np^2$ . So if  $n \gg p$  then Cholesky (and sweeping) are faster than the QR approaches. According to Monahan, modified Gram-Schmidt is most numerically stable and sweeping least. In general, regression is pretty quick unless  $p$  is large since it is linear in  $n$ , so it may not be worth worrying too much about computational differences of the sort noted here.

### 3.4.3 Regression and the QR in R

Regression in R uses the QR decomposition via `qr()`, which calls a Fortran function. `qr()` (and the Fortran functions that are called) is specifically designed to output quantities useful in fitting linear models. Note that by default you get the skinny QR, namely only the first  $p$  rows of  $R$  and the first  $p$  columns of  $Q$ , where the latter form an orthonormal basis for the column space of  $X$ . The remaining columns form an orthonormal basis for the null space of  $X$  (the space orthogonal to the column space of  $X$ ). The analogy in regression is that we get the basis vectors for the regression, while adding the remaining columns gives us the full  $n$ -dimensional space of the observations.

`qr()` returns the result as a list meant for use by other tools. R stores the  $R$  matrix in the upper triangle of `$qr`, while the lower triangle of `$qr` and `$aux` store the information for constructing  $Q$

(this relates to the Householder-related vectors  $u$  below). One can multiply by  $Q$  using `qr.qy()` and by  $Q^\top$  using `qr.qty()`. If you want to extract  $R$  and  $Q$ , the following will work:

```
X.qr = qr(X)
Q = qr.Q(X.qr)
R = qr.R(X.qr)
```

As a side note, there are QR-based functions that provide regression-related quantities, such as `qr.resid()`, `qr.fitted()` and `qr.coef()`. These functions (and their Fortran counterparts) exist because one can work through the various regression quantities of interest and find their expressions in terms of  $Q$  and  $R$ , with nice properties resulting from  $Q$  being orthogonal and  $R$  triangular.

### 3.4.4 Computing the QR decomposition

We'll work through some of the details of the different approaches to the QR, in part because they involve some concepts that may be useful in other contexts.

One approach involves reflections of vectors and a second rotations of vectors. Reflections and rotations are transformations that are performed by orthogonal matrices. The determinant of a reflection matrix is -1 and the determinant of a rotation matrix is 1. We'll see some of the details in the demo code.

**Reflections** If  $u$  and  $v$  are orthonormal vectors and  $x$  is in the space spanned by  $u$  and  $v$ ,  $x = c_1u + c_2v$ , then  $\tilde{x} = -c_1u + c_2v$  is a reflection (a *Householder* reflection) along the  $u$  dimension (since we are using the negative of that basis vector). We can think of this as reflecting across the plane perpendicular to  $u$ . This extends simply to higher dimensions with orthonormal vectors,  $u, v_1, v_2, \dots$

Suppose we want to formulate the reflection in terms of a “Householder” matrix,  $Q$ . It turns out that

$$Qx = \tilde{x}$$

if  $Q = I - 2uu^\top$ .  $Q$  has the following properties: (1)  $Qu = -u$ , (2)  $Qv = v$  for  $u^\top v = 0$ , (3)  $Q$  is orthogonal and symmetric.

One way to create the QR decomposition is by a series of Householder transformations that create an upper triangular  $R$  from  $X$ :

$$\begin{aligned} R &= Q_p \cdots Q_1 X \\ Q &= (Q_p \cdots Q_1)^\top \end{aligned}$$

where we make use of the symmetry in defining  $Q$ .

Basically  $Q_1$  reflects the first column of  $X$  with respect to a carefully chosen  $u$ , so that the result is all zeroes except for the first element. We want  $Q_1 x = \tilde{x} = (\|x\|, 0, \dots, 0)$ . This can be achieved with  $u = \frac{x - \tilde{x}}{\|x - \tilde{x}\|}$ . Then  $Q_2$  makes the last  $n - 2$  rows of the second column equal to zero. We'll work through this a bit in class.

In the regression context, as we work through the individual transformations,  $Q_j = I - 2u_j u_j^\top$ , we apply them to  $X$  and  $Y$  to create  $R$  (note this would not involve doing the full matrix multiplication - think about what calculations are actually needed) and  $QY = Q^\top Y$ , and then solve  $R\beta = Q^\top Y$ . To find  $\text{Cov}(\hat{\beta}) \propto (X^\top X)^{-1} = (R^\top R)^{-1} = R^{-1} R^{-\top}$  we do need to invert  $R$ , but it's upper-triangular and of dimension  $p \times p$ . It turns out that  $Q^\top Y$  can be partitioned into the first  $p$  and the last  $n - p$  elements,  $z^{(1)}$  and  $z^{(2)}$ . The SSR is  $\|z^{(1)}\|^2$  and SSE is  $\|z^{(2)}\|^2$ .

**Rotations** A Givens rotation matrix rotates a vector in a two-dimensional subspace to be axis oriented with respect to one of the two dimensions by changing the value of the other dimension. E.g. we can create  $\tilde{x} = (x_1, \dots, \tilde{x}_p, \dots, 0, \dots, x_n)$  from  $x = (x_1, \dots, x_p, \dots, x_q, \dots, x_n)$  using a matrix multiplication:  $\tilde{x} = Qx$ .  $Q$  is orthogonal but not symmetric.

We can use a series of Givens rotations to do the QR but unless it is done carefully, more computations are needed than with Householder reflections. The basic story is that we apply a series of Givens rotations to  $X$  such that we zero out the lower triangular elements.

$$\begin{aligned} R &= Q_{pn} \cdots Q_{23} Q_{1n} \cdots Q_{13} Q_{12} X \\ Q &= (Q_{pn} \cdots Q_{12})^\top \end{aligned}$$

Note that we create the  $n - p$  zero rows in  $R$  (because the calculations affect the upper triangle of  $R$ ), but we can then ignore those rows and the corresponding columns of  $Q$ .

**Gram-Schmidt Orthogonalization** Gram-Schmidt involves finding a set of orthonormal vectors to span the same space as a set of LIN vectors,  $x_1, \dots, x_p$ . If we take the LIN vectors to be the columns of  $X$ , so that we are discussing the column space of  $X$ , then G-S yields the QR decomposition. Here's the algorithm:

1.  $\tilde{x}_1 = \frac{x_1}{\|x_1\|}$  (normalize the first vector)
2. Orthogonalize the remaining vectors with respect to  $\tilde{x}_1$ :
  - (a)  $\tilde{x}_2 = \frac{x_2 - \tilde{x}_1^\top x_2 \tilde{x}_1}{\|x_2 - \tilde{x}_1^\top x_2 \tilde{x}_1\|}$ , which orthogonalizes with respect to  $\tilde{x}_1$  and normalizes. Note that  $\tilde{x}_1^\top x_2 \tilde{x}_1 = \langle \tilde{x}_1, x_2 \rangle \tilde{x}_1$ . So we are finding a scaling,  $c\tilde{x}_1$ , where  $c$  is based on the inner product, to remove the variation in the  $x_1$  direction from  $x_2$ .

- (b) For  $k > 2$ , find interim vectors,  $x_k^{(2)}$ , by orthogonalizing with respect to  $\tilde{x}_1$
3. Proceed for  $k = 3, \dots$ , in turn orthogonalizing and normalizing the first of the remaining vectors w.r.t.  $\tilde{x}_{k-1}$  and orthogonalizing the remaining vectors w.r.t.  $\tilde{x}_{k-1}$  to get new interim vectors

Mathematically, we could instead orthogonalize  $x_2$  w.r.t.  $\tilde{x}_1$ , then orthogonalize  $x_3$  w.r.t.  $\{\tilde{x}_1, \tilde{x}_2\}$ , etc. The algorithm above is the *modified* G-S, and is known to be more numerically stable if the columns of  $X$  are close to collinear, giving vectors that are closer to orthogonal. The resulting  $\tilde{x}$  vectors are the columns of  $Q$ . The elements of  $R$  are obtained as we proceed: the diagonal values are the the normalization values in the denominators, while the off-diagonals are the inner products with the already-computed columns of  $Q$  that are computed as part of the numerators.

Another way to think about this is that  $R = Q^\top X$ , which is the same as regressing the columns of  $X$  on  $Q$ , since  $(Q^\top Q)^{-1} Q^\top X = Q^\top X$ . By construction, the first column of  $X$  is a scaling of the first column of  $Q$ , the second column of  $X$  is a linear combination of the first two columns of  $Q$ , etc., so  $R$  being upper triangular makes sense.

### 3.4.5 The “tall-skinny” QR

Suppose you have a very large regression problem, with  $n$  very large, and  $n \gg p$ . There is a variant of the QR, called the tall-skinny QR (see <http://arxiv.org/pdf/0808.2664v1.pdf> for details) that allows us to find the decomposition in a parallel fashion. The basic idea is to do a nested set of QR decompositions on blocks of rows of  $X$ :

$$X = \begin{pmatrix} X_0 \\ X_1 \\ X_2 \\ X_3 \end{pmatrix} = \begin{pmatrix} Q_0 R_0 \\ Q_1 R_1 \\ Q_2 R_2 \\ Q_3 R_3 \end{pmatrix},$$

followed by ‘reduction’ steps (this can be done in a map-reduce context) that do the  $QR$  of pairs of the  $R$  factors:

$$\begin{pmatrix} R_0 \\ R_1 \\ R_2 \\ R_3 \end{pmatrix} = \begin{pmatrix} \begin{pmatrix} R_0 \\ R_1 \end{pmatrix} \\ \begin{pmatrix} R_2 \\ R_3 \end{pmatrix} \end{pmatrix} = \begin{pmatrix} Q_{01} R_{01} \\ Q_{23} R_{23} \end{pmatrix}$$

and

$$\begin{pmatrix} R_{01} \\ R_{23} \end{pmatrix} = Q_{0123} R_{0123}.$$

The full decomposition is then

$$X = \begin{pmatrix} Q_0 & 0 & 0 & 0 \\ 0 & Q_1 & 0 & 0 \\ 0 & 0 & Q_2 & 0 \\ 0 & 0 & 0 & Q_3 \end{pmatrix} \begin{pmatrix} Q_{01} & 0 \\ 0 & Q_{23} \end{pmatrix} Q_{0123} R_{0123} = QR.$$

The computation can be done in parallel (in particular it can be done with map-reduce) and the  $Q$  matrix for big problems would generally not be computed explicitly but would be stored in its constituent pieces.

Alternatively, there is a variant on the algorithm that processes the row-blocks of  $X$  serially, allowing you to do QR on a large tall-skinny matrix that you can't fit in memory (or possibly even on disk). First you do  $QR$  on  $X_0$  to get  $Q_0 R_0$ . Then you stack  $R_0$  on top of  $X_1$  and do QR to get  $R_{01}$ . Then stack  $R_{01}$  on top of  $X_2$  to get  $R_{012}$ , etc.

### 3.5 Determinants

The absolute value of the determinant of a square matrix can be found from the product of the diagonals of the triangular matrix in any factorization that gives a triangular (including diagonal) matrix times an orthogonal matrix (or matrices) since the determinant of an orthogonal matrix is either one or minus one.

$$\begin{aligned} |A| &= |QR| = |Q||R| = \pm |R| \\ |A^\top A| &= |(QR)^\top QR| = |R^\top R| = |R_1^\top R_1| = |R_1|^2 \end{aligned}$$

In R, the following will do it (on the log scale), since  $R$  is stored in the upper triangle of the  $qr$  element.

```
myqr = qr(A)
magn = sum(log(abs(diag(myqr$qr))))
```

An alternative is the product of the diagonal elements of  $D$  (the singular values) in the SVD factorization,  $A = UDV^\top$ .

For non-negative definite matrices, we know the determinant is non-negative, so the uncertainty about the sign is not an issue. For positive definite matrices, a good approach is to use the product of the diagonal elements of the Cholesky decomposition.

One can also use the product of the eigenvalues:  $|A| = |\Gamma \Lambda \Gamma^{-1}| = |\Gamma| |\Gamma^{-1}| |\Lambda| = |\Lambda|$

**Computation** Computing from any of these diagonal or triangular matrices as the product of the diagonals is prone to overflow and underflow, so we **always** work on the log scale as the sum of the



log of the values. When some of these may be negative, we can always keep track of the number of negative values and take the log of the absolute values.

Often we will have the factorization as a result of other parts of the computation, so we get the determinant for free.

R's *determinant()* uses the LU decomposition. Supposedly *det()* just wraps *determinant()*, but I can't seem to pass the *logarithm* argument into *det()*, so *determinant()* seems more useful.

## 4 Eigendecomposition and SVD

### 4.1 Eigendecomposition

The eigendecomposition (spectral decomposition) is useful in considering convergence of algorithms and of course for statistical decompositions such as PCA. We think of decomposing the components of variation into orthogonal patterns (the eigenvectors) with variances (eigenvalues) associated with each pattern.

Square symmetric matrices have real eigenvectors and eigenvalues, with the factorization into orthogonal  $\Gamma$  and diagonal  $\Lambda$ ,  $A = \Gamma\Lambda\Gamma^\top$ , where the eigenvalues on the diagonal of  $\Lambda$  are ordered in decreasing value. Of course this is equivalent to the definition of an eigenvalue/eigenvector pair as a pair such that  $Ax = \lambda x$  where  $x$  is the eigenvector and  $\lambda$  is a scalar, the eigenvalue. The inverse of the eigendecomposition is simply  $\Gamma\Lambda^{-1}\Gamma^\top$ . On a similar note, we can create a square root matrix,  $\Gamma\Lambda^{1/2}$ , by taking the square roots of the eigenvalues.

The spectral radius of  $A$ , denoted  $\rho(A)$ , is the maximum of the absolute values of the eigenvalues. As we saw when talking about ill-conditionedness, for symmetric matrices, this maximum is the induced norm, so we have  $\rho(A) = \|A\|_2$ . It turns out that  $\rho(A) \leq \|A\|$  for any induced matrix norm. The spectral radius comes up in determining the rate of convergence of some iterative algorithms.

**Computation** There are several methods for eigenvalues; a common one for doing the full eigendecomposition is the *QR algorithm*. The first step is to reduce  $A$  to upper Hessenberg form, which is an upper triangular matrix except that the first subdiagonal in the lower triangular part can be non-zero. For symmetric matrices, the result is actually tridiagonal. We can do the reduction using Householder reflections or Givens rotations. At this point the QR decomposition (using Givens rotations) is applied iteratively (to a version of the matrix in which the diagonals are shifted), and the result converges to a diagonal matrix, which provides the eigenvalues. It's more work to get the eigenvectors, but they are obtained as a product of Householder matrices (required for the initial reduction) multiplied by the product of the  $Q$  matrices from the successive QR decompositions.

We won't go into the algorithm in detail, but note that it involves manipulations and ideas we've seen already.

If only the largest (or the first few largest) eigenvalues and their eigenvectors are needed, which can come up in time series and Markov chain contexts, the problem is easier and can be solved by the *power method*. E.g., in a Markov chain context, steady state is reached through  $x_t = A^t x_0$ . One can find the largest eigenvector by multiplying by  $A$  many times, normalizing at each step.  $v^{(k)} = Az^{(k-1)}$  and  $z^{(k)} = v^{(k)} / \|v^{(k)}\|$ . There is an extension to find the  $p$  largest eigenvalues and their vectors. See the demo code for an implementation.

## 4.2 Singular value decomposition

Let's consider an  $n \times m$  matrix,  $A$ , with  $n \geq m$  (if  $m > n$ , we can always work with  $A^\top$ ). This often is a matrix representing  $m$  features of  $n$  observations. We could have  $n$  documents and  $m$  words, or  $n$  gene expression levels and  $m$  experimental conditions, etc.  $A$  can always be decomposed as

$$A = UDV^\top$$

where  $U$  and  $V$  are matrices with orthonormal columns (left and right eigenvectors) and  $D$  is diagonal with non-negative values (which correspond to eigenvalues in the case of square  $A$  and to squared eigenvalues of  $A^\top A$ ).

The SVD can be represented in more than one way. One representation is

$$A_{n \times m} = U_{n \times k} D_{k \times k} V_{k \times m}^\top = \sum_{j=1}^k D_{jj} u_j v_j^\top$$

where  $u_j$  and  $v_j$  are the columns of  $U$  and  $V$  and where  $k$  is the rank of  $A$  (which is at most the minimum of  $n$  and  $m$  of course). The diagonal elements of  $D$  are the singular values.

If  $A$  is positive semi-definite, the eigendecomposition is an SVD. Furthermore,  $A^\top A = VD^2V^\top$  and  $AA^\top = UD^2U^\top$ , so we can find the eigendecomposition of such matrices using the SVD of  $A$ . Note that the squares of the singular values of  $A$  are the eigenvalues of  $A^\top A$  and  $AA^\top$ .

We can also fill out the matrices to get

$$A = U_{n \times n} D_{n \times m} V_{m \times m}^\top$$

where the added rows and columns of  $D$  are zero with the upper left block the  $D_{k \times k}$  from above.

**Uses** The SVD is an excellent way to determine a matrix rank and to construct a pseudo-inverse ( $A^+ = VD^+U^\top$ ).

We can use the SVD to approximate  $A$  by taking  $A \approx \tilde{A} = \sum_{j=1}^p D_{jj} u_j v_j^\top$  for  $p < m$ . This approximation holds in terms of the Frobenius norm for  $A - \tilde{A}$ . As an example if we have a large image of dimension  $n \times m$ , we could hold a compressed version by a rank- $p$  approximation using the SVD. The SVD is used a lot in clustering problems. For example, the Netflix prize was won based on a variant of SVD (in fact all of the top methods used variants on SVD, I believe).

**Computation** The basic algorithm (Golub-Reinsch) is similar to the QR method for the eigen-decomposition. We use a series of Householder transformations on the left and right to reduce  $A$  to an upper bidiagonal matrix,  $A^{(0)}$ . The post-multiplications (the transformations on the right) generate the zeros in the upper triangle. (An upper bidiagonal matrix is one with non-zeroes only on the diagonal and first subdiagonal above the diagonal). Then the algorithm produces a series of upper bidiagonal matrices,  $A^{(0)}$ ,  $A^{(1)}$ , etc. that converge to a diagonal matrix,  $D$ . Each step is carried out by a sequence of Givens transformations:

$$\begin{aligned} A^{(j+1)} &= R_{m-2}^\top R_{m-3}^\top \cdots R_0^\top A^{(j)} T_0 T_1 \cdots T_{m-2} \\ &= R A^{(j)} T \end{aligned}$$

This eventually gives  $A^{(\infty)} = D$  and by construction,  $U$  (the product of the pre-multiplied Householder matrices and the  $R$  matrices) and  $V$  (the product of the post-multiplied Householder matrices and the  $T$  matrices) are orthogonal. The result is then transformed by a diagonal matrix to make the elements of  $D$  non-negative and by permutation matrices to order the elements of  $D$  in nonincreasing order.

**Computation for large tall-skinny matrices** The SVD can also be generated from a QR decomposition. Take  $X = QR$  and then do an SVD on the  $R$  matrix to get  $X = QUDV^\top = U^*DV^\top$ . This is particularly helpful for the case when  $X$  is tall and skinny (suppose  $X$  is  $n \times p$  with  $n \gg p$ ), because we can do the tall-skinny QR, and the resulting SVD on  $R$  is easy computationally if  $p$  is manageable.

## 5 Computation

### 5.1 Linear algebra in R

Speedups and storage savings can be obtained by working with matrices stored in special formats when the matrices have special structure. E.g., we might store a symmetric matrix as a full matrix but only use the upper or lower triangle. Banded matrices and block diagonal matrices are

other common formats. Banded matrices are all zero except for  $A_{i,i+c_k}$  for some small number of integers,  $c_k$ . Viewed as an image, these have bands. The bands are known as co-diagonals.

Note that for many matrix decompositions, you can change whether all of the aspects of the decomposition are returned, or just some, which may speed calculations.

Some useful packages in R for matrices are *Matrix*, *spam*, and *bdsmatrix*. *Matrix* can represent a variety of rectangular matrices, including triangular, orthogonal, diagonal, etc. and provides methods for various matrix calculations that are specific to the matrix type. *spam* handles general sparse matrices with fast matrix calculations, in particular a fast Cholesky decomposition. *bdsmatrix* focuses on block-diagonal matrices, which arise frequently in contexts where there is clustering that induces within-cluster correlation and cross-cluster independence.

In general, matrix operations in R go to compiled C or Fortran code without much intermediate R code, so they can actually be pretty efficient and are based on the best algorithms developed by numerical experts. The core libraries that are used are LAPACK and BLAS (the Linear Algebra PACKage and the Basic Linear Algebra Subroutines). As we've discussed in the parallelization unit, one way to speed up code that relies heavily on linear algebra is to make sure you have a BLAS library tuned to your machine. These include OpenBLAS (free; formerly called GotoBLAS), Intel's MKL, AMD's ACML, and Apple's vecLib. These can be installed and R can be linked to the shared object library file (.so file or .dylib on a Mac) for the fast BLAS. These BLAS libraries are also available in threaded versions that farm out the calculations across multiple cores or processors that share memory.

BLAS routines do vector operations (level 1), matrix-vector operations (level 2), and dense matrix-matrix operations (level 3). Often the name of the routine has as its first letter “d”, “s”, “c” to indicate the routine is double precision, single precision, or complex. LAPACK builds on BLAS to implement standard linear algebra routines such as eigendecomposition, solutions of linear systems, a variety of factorizations, etc.

## 5.2 Sparse matrices

As an example of exploiting sparsity, here's how the *spam* package in R stores a sparse matrix. Consider the matrix to be row-major and store the non-zero elements in order in a vector called *value*. Then create a vector called *rowptr* that stores the position of the first element of each row. Finally, have a vector, *colindices* that tells the column identity of each element. Here's an example in the *spam* package in R:

```
require(spam)
mat = matrix(c(0,0,1,0,10,0,0,0,100,0,rep(0,5),1000,rep(0,4)), nrow = 4, byrow = FALSE)
mat = as.spam(mat)
```

```

mat@entries
## [1]      1    10   100  1000

mat@rowpointers
## [1] 1 3 4 4 5

mat@colindices
## [1] 3 5 4 1

```

We can do a fast matrix multiply,  $Ab$ , as follows in pseudo-code:

```

for(i in 1:nrows(A)){
    x[i] = 0
    # should also check that row is not empty...
    for(j in (rowptr[i]:(rowptr[i+1]-1)) {
        x[i] = x[i] + value[j] * b[colindices[j]]
    }
}

```

How many computations have we done? Only  $k$  multiplies and  $O(k)$  additions where  $k$  is the number of non-zero elements of  $A$ . Compare this to the usual  $O(n^2)$  for dense multiplication.

Note that for the Cholesky of a sparse matrix, if the sparsity pattern is fixed, but the entries change, one can precompute an optimal re-ordering that retains as much sparsity in  $U$  as possible. Then multiple Cholesky decompositions can be done more quickly as the entries change.

**Banded matrices** Suppose we have a banded matrix  $A$  where the lower bandwidth is  $p$ , namely  $A_{ij} = 0$  for  $i > j + p$  and the upper bandwidth is  $q$  ( $A_{ij} = 0$  for  $j > i + q$ ). An alternative to reducing to  $Ux = b^*$  is to compute  $A = LU$  and then do two solutions,  $U^{-1}(L^{-1}b)$ . One can show that the computational complexity of the LU factorization is  $O(npq)$  for banded matrices, while solving the two triangular systems is  $O(np + nq)$ , so for small  $p$  and  $q$ , the speedup can be dramatic.

Banded matrices come up in time series analysis. E.g., MA models produce banded covariance structures because the covariance is zero after a certain number of lags.

### 5.3 Low rank updates

A transformation of the form  $A - uv^\top$  is a rank-one update because  $uv^\top$  is of rank one.

More generally a low rank update of  $A$  is  $\tilde{A} = A - UV^\top$  where  $U$  and  $V$  are  $n \times m$  with  $n \geq m$ . The Sherman-Morrison-Woodbury formula tells us that

$$\tilde{A}^{-1} = A^{-1} + A^{-1}U(I_m - V^\top A^{-1}U)^{-1}V^\top A^{-1}$$

so if we know  $x_0 = A^{-1}b$ , then the solution to  $\tilde{A}x = b$  is  $x + A^{-1}U(I_m - V^\top A^{-1}U)^{-1}V^\top x$ . Provided  $m$  is not too large, and particularly if we already have a factorization of  $A$ , then  $A^{-1}U$  is not too bad computationally, and  $I_m - V^\top A^{-1}U$  is  $m \times m$ . As a result  $A^{-1}(U(\cdots)^{-1}V^\top x)$  isn't too bad.

This also comes up in working with precision matrices in Bayesian problems where we may have  $A^{-1}$  but not  $A$  (we often add precision matrices to find conditional normal distributions). An alternative expression for the formula is  $\tilde{A} = A + UCV^\top$ , and the identity tells us

$$\tilde{A}^{-1} = A^{-1} - A^{-1}U(C^{-1} + V^\top A^{-1}U)^{-1}V^\top A^{-1}$$

Basically Sherman-Morrison-Woodbury gives us matrix identities that we can use in combination with our knowledge of smart ways of solving systems of equations.

## 6 Iterative solutions of linear systems

**Gauss-Seidel** Suppose we want to iteratively solve  $Ax = b$ . Here's the algorithm, which sequentially updates each element of  $x$  in turn.

- Start with an initial approximation,  $x^{(0)}$ .
- Hold all but  $x_1^{(0)}$  constant and solve to find  $x_1^{(1)} = \frac{1}{a_{11}}(b_1 - \sum_{j=2}^n a_{1j}x_j^{(0)})$ .
- Repeat for the other rows of  $A$  (i.e., the other elements of  $x$ ), finding  $x^{(1)}$ .
- Now iterate to get  $x^{(2)}$ ,  $x^{(3)}$ , etc. until a convergence criterion is achieved, such as  $\|x^{(k)} - x^{(k-1)}\| \leq \epsilon$  or  $\|r^{(k)} - r^{(k-1)}\| \leq \epsilon$  for  $r^{(k)} = b - Ax^{(k)}$ .

Let's consider how many operations are involved in a single update:  $O(n)$  for each element, so  $O(n^2)$  for each update. Thus if we can stop well before  $n$  iterations, we've saved computation relative to exact methods.

If we decompose  $A = L + D + U$  where  $L$  is strictly lower triangular,  $U$  is strictly upper triangular, then Gauss-Seidel is equivalent to solving

$$(L + D)x^{(k+1)} = b - Ux^{(k)}$$

and we know that solving the lower triangular system is  $O(n^2)$ .

It turns out that the rate of convergence depends on the spectral radius of  $(L + D)^{-1}U$ .

Gauss-Seidel amounts to optimizing by moving in axis-oriented directions, so it can be slow in some cases.

**Conjugate gradient** For positive definite  $A$ , conjugate gradient (CG) reexpresses the solution to  $Ax = b$  as an optimization problem, minimizing

$$f(x) = \frac{1}{2}x^\top Ax - x^\top b,$$

since the derivative of  $f(x)$  is  $Ax - b$  and at the minimum this gives  $Ax - b = 0$ .

Instead of finding the minimum by following the gradient at each step (so-called steepest descent, which can give slow convergence - we'll see a demonstration of this in the optimization unit), CG chooses directions that are mutually conjugate w.r.t.  $A$ ,  $d_i^\top A d_j = 0$  for  $i \neq j$ . The method successively chooses vectors giving the direction,  $d_k$ , in which to move down towards the minimum and a scaling of how much to move,  $\alpha_k$ . If we start at  $x_{(0)}$ , the  $k$ th point we move to is  $x_{(k)} = x_{(k-1)} + \alpha_k d_k$  so we have

$$x_{(k)} = x_{(0)} + \sum_{j \leq k} \alpha_j d_j$$

and we use a convergence criterion such as given above for Gauss-Seidel. The directions are chosen to be the residuals,  $b - Ax_{(k)}$ . Here's the basic algorithm:

- Choose  $x_{(0)}$  and define the residual,  $r_{(0)} = b - Ax_{(0)}$  (the error on the scale of  $b$ ) and the direction,  $d_0 = r_{(0)}$  and set  $k = 0$ .
- Then iterate
  - $\alpha_k = \frac{r_{(k)}^\top r_{(k)}}{d_k^\top A d_k}$  (choose step size so next error will be orthogonal to current direction - which we can express in terms of the residual, which is easily computable)
  - $x_{(k+1)} = x_{(k)} + \alpha_k d_k$  (update current value)
  - $r_{(k+1)} = r_{(k)} - \alpha_k A d_k$  (update current residual)
  - $d_{k+1} = r_{(k+1)} + \frac{r_{(k+1)}^\top r_{(k)}}{r_{(k)}^\top r_{(k)}} d_k$  (choose next direction by conjugate Gram-Schmidt, starting with  $r_{(k+1)}$  and removing components that are not  $A$ -orthogonal to previous directions, but it turns out that  $r_{(k+1)}$  is already  $A$ -orthogonal to all but  $d_k$ ).
- Stop when  $\|r_{(k+1)}\|$  is sufficiently small.

The convergence of the algorithm depends in a complicated way on the eigenvalues, but in general convergence is faster when the condition number is smaller (the eigenvalues are not too spread out). CG will in principle give the exact answer in  $n$  steps (where  $A$  is  $n \times n$ ). However, computationally we lose accuracy and interest in the algorithm is really as an iterative approximation where we stop before  $n$  steps. The approach basically amounts to moving in axis-oriented directions in a space stretched by  $A$ .

In general, CG is used for large sparse systems.

See the [extensive description from Shewchuk](#) for more details and for the figures shown in class, as well as the use of CG when  $A$  is not positive definite.

**Updating a solution** Sometimes we have solved a system,  $Ax = b$  and then need to solve  $Ax = c$ . If we have solved the initial system using a factorization, we can reuse that factorization and solve the new system in  $O(n^2)$ . Iterative approaches can do a nice job if  $c = b + \delta b$ . Start with the solution  $x$  for  $Ax = b$  as  $x^{(0)}$  and use one of the methods above.