

FN 5012  
5th. class

①

### Minimization Algorithms.

When the matrix  $A$  is symmetric and positive definite, the problem of obtaining the solution of the linear equation

$$Ax = b,$$

is equivalent to finding the point at which the quadratic functional

$$J(y) = \frac{1}{2} y^T A y - y^T b$$

achieves its minimum. We can thus use a minimization algorithm for  $J$  to obtain the solution of the linear equation  $x$ .

We shall study the gradient method and the conjugate gradient method. We shall see that, even for sparse matrices, the performance of the gradient method is not necessarily better than that of the Choleski decomposition. We shall also see that for the sparse matrices of finite element approximations, the conjugate gradient method is the method of choice.

The gradient method and the Conjugate gradient method have both algorithms of the following form

$$x^0 = x_0, \quad r^0 = Ax^0 - b.$$

```

For k=0, ..., M do
  If  $r^k \approx 0$  then set  $x = x^k$  and stop.
  else
    Compute a search direction  $d^k$ ,
    compute the step length  $\alpha^k$ ,
    and compute  $x^{k+1}$  as follows:
      
$$x^{k+1} = x^k + \alpha^k d^k.$$

  endIf
endFor

```

Notice that once the search direction is fixed we can only play with the value of  $\alpha^k$  to set  $x^{k+1}$ . The step length  $\alpha^k$  is called optimal when

$$J(x^k + \alpha^k d^k) \leq J(x^k + \alpha d^k) \quad \forall \alpha \in \mathbb{R}.$$

Since

$$J(x^k + \alpha d^k) = J(x^k) + \alpha d^{kT} (Ax^k - b) + \frac{1}{2} \alpha^2 d^{kT} A d^k$$

we have that the optimal step length is

$$\alpha^k = - \frac{d^{kT} r^k}{d^{kT} A d^k}$$

- ② the Gradient method. This method picks the search direction  $d^k$  to be the direction of fastest (local) descent of the functional  $J$ , that is,

$$\begin{aligned} d^k &= - \nabla J(x^k) \\ &= - (Ax^k - b) \\ &= - r^k \end{aligned}$$

the algorithm for this method is thus the following:

$$x^0 = x_0, \quad r^0 = Ax^0 - b.$$

For  $k=0, \dots, M$  do

If  $r^k \approx 0$  then set  $x = x^k$  and stop.

else

$$\alpha^k = \frac{r^{kT} r^k}{r^{kT} A r^k},$$

$$x^{k+1} = x^k - \alpha^k r^k$$

end If

end For

The gradient method generates a minimizing sequence  $\{x^k\}_{k \geq 0}$  since

$$\begin{aligned} J(x^{k+1}) &= J(x^k) - \alpha^k r^{kT} r^k + \frac{1}{2} \alpha^{2k} r^{kT} A r^k \\ &= J(x^k) - \frac{1}{2} \frac{(r^{kT} r^k)^2}{r^{kT} A r^k} \\ &< J(x^k) \end{aligned}$$

if  $r^k \neq 0$ . Moreover, the error in the energy norm decreases exponentially with the number of iterations. To see this, first notice that

$$\begin{aligned} e^{k+1} &= x^{k+1} - x \\ &= (x^k - \alpha^k (Ax^k - Ax)) - x \\ &= e^k - \alpha^k A e^k, \end{aligned}$$

and so

$$\|e^{k+1}\|_A^2 = e^{kT} A e^k - 2\alpha^k e^{kT} A^2 e^k + \alpha_k^2 e^{kT} A^3 e^k.$$

Since

$$\alpha^k = \frac{e^{kT} A^2 e^k}{e^{kT} A^3 e^k},$$

we get

$$\begin{aligned}\|e^{k+1}\|_A^2 &= e^{kT} A e^k - \frac{(e^{kT} A^2 e^k)^2}{e^{kT} A^3 e^k} \\ &= \left[ 1 - \frac{e^{kT} A^2 e^k}{e^{kT} A^3 e^k} \right] \frac{e^{kT} A^2 e^k}{e^{kT} A e^k} \|e^k\|_A^2\end{aligned}$$

and by the so-called Kantorovich inequality,

$$\|e^{k+1}\|_A^2 \leq \left( \frac{\kappa - 1}{\kappa + 1} \right)^2 \|e^k\|_A^2,$$

where  $\kappa = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$ .

thus, we have

$$\|e^k\|_A \leq \left( \frac{\kappa - 1}{\kappa + 1} \right)^k \|e^0\|_A.$$

this inequality tells us that if the matrix  $A$  has a very big condition number, which is usually the case for matrices associated with elliptic problems, then the gradient method converges extremely slowly.

For example, if we want to make sure that

$$\|e^M\|_A \leq \epsilon,$$

Then we must take

$$M \geq 2 \frac{\log(\epsilon / \|e^0\|_A)}{\log(1 - 1/k)}$$

$$\approx \left[ -2 \log\left(\frac{\epsilon}{\|e^0\|_A}\right) \right] \cdot k.$$

$$= c_0 k.$$

this estimate is sharp. To see this, let us consider the following simple case:

$$A = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \quad \text{and} \quad b = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

If  $k = \frac{\lambda_1}{\lambda_2} = 1$ , then it is easy to verify that the gradient method converges in one iteration, as the error estimate predicts.

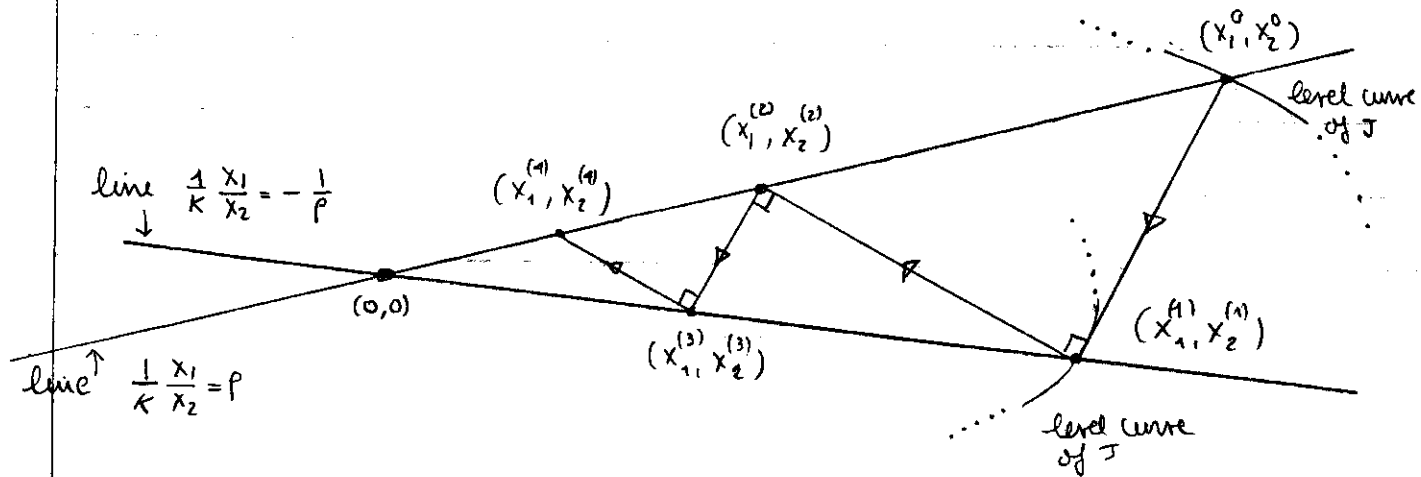
If  $k > 1$ , the gradient method still converges in a single iteration if  $x_1^0 = 0$  or if  $x_2^0 = 0$ . However, in the other cases, we have that

$$\frac{1}{k} \frac{x_1^{(2n)}}{x_2^{(2n)}} = \rho, \quad \frac{1}{k} \frac{x_1^{(2n+1)}}{x_2^{(2n+1)}} = -\frac{1}{\rho},$$

where

$$\rho = \frac{1}{k} \frac{x_1^{(0)}}{x_2^{(0)}}.$$

this implies that we have the following situation:



It can be shown that

$$x_1^{(k+2)} = \left( \frac{1 - 2/k + 1/k^2}{1 + (p^2 + \bar{p}^2)/k + 1/k^2} \right) x_1^{(k)},$$

and so,

$$|x_1^{(k+2)}| \leq \left( \frac{1 - 1/k}{1 + 1/k} \right)^2 |x_1^{(k)}|.$$

this shows that the number of iterations for convergence is proportional to  $k$  when  $k \gg 1$ , as expected.

Let us conclude our study of the gradient method by obtaining the number of operations it needs. This number is clearly  $M$  times the number of operations needed to compute  $x^{k+1}$  from  $x^k$ .

We can rewrite the gradient method's algorithm as follows, to minimize the number of operations it requires:

$$x^0 = x_0; \quad r^0 = Ax^0 - b.$$

For  $k=0, \dots, M$  do

$$p^k = A r^k$$

$$\alpha^k = \frac{r^{kT} r^k}{r^{kT} p^k}$$

$$x^{k+1} = x^k - \alpha^k r^k$$

$$r^{k+1} = r^k - \alpha^k p^k$$

endfor

thus, the number of operations is  $M$  times the number of operations a multiplication by  $A$  requires!



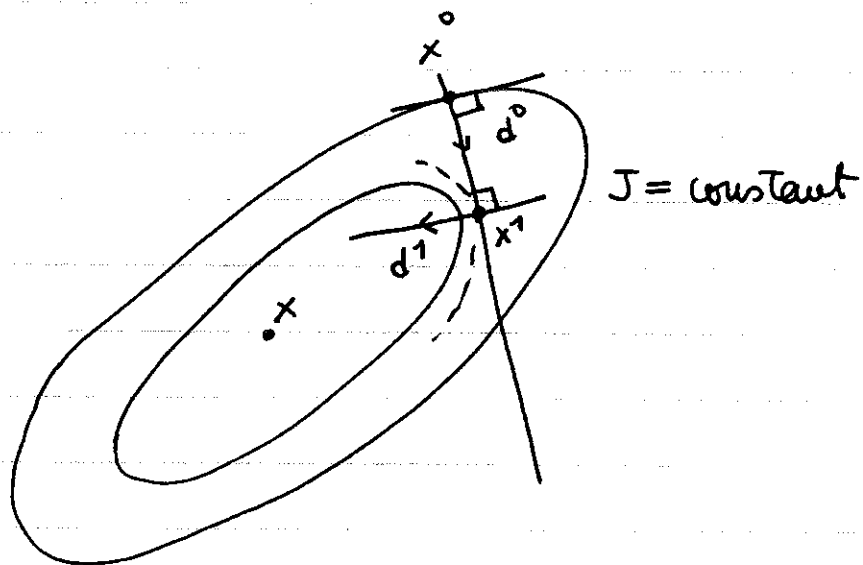
How can we improve the steepest descent method?  
 To see how, let us begin by noting that

$$\begin{aligned} J(y) &= \frac{1}{2} y^T A y - b^T y \\ &= \frac{1}{2} y^T A y - x^T A y \\ &= \frac{1}{2} (y-x)^T A (y-x) - \frac{1}{2} x^T A x \end{aligned}$$

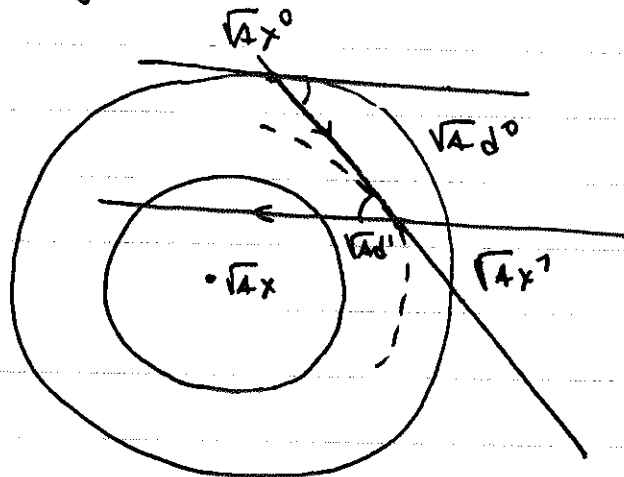
since  $A$  is symmetric (verify this computation!).  
 Hence

$$J(y) = J(x) + \frac{1}{2} (y-x)^T A (y-x) \quad \forall y \in \mathbb{R}^n.$$

Since  $x$  is fixed, the level curves of  $J$  are those of the mapping  $y \mapsto \frac{1}{2} (y-x)^T A (y-x)$  which are ellipses centered at the point  $x$ :



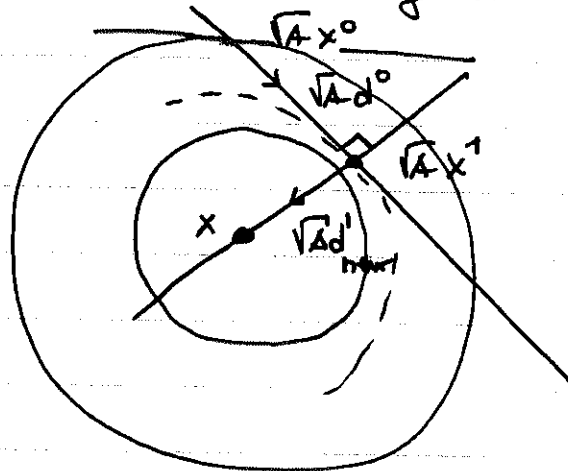
If we carry out the linear transformation  $z = \sqrt{A} y$ , the previous figure becomes



We now see that we could improve the steepest descent method by taking a different search direction  $d^1$ . Note that if the new search direction,  $d_{\text{new}}^1$  is such that

$$0 = (\sqrt{A}d^0)^T (\sqrt{A}d_{\text{new}}^1) = d^{0T} A d_{\text{new}}^1,$$

the method would converge in only two iterations:



The conjugate gradient is based in this idea and is indeed superior to the steepest descent method as we see next.

- ② the Conjugate gradient method. This method picks the search direction  $d^k$  as follows:

$$d^k = -r^k + \beta^k d^{k-1}, \quad (d^0 = -r^0),$$

where the factor  $\beta^k$  is obtained by enforcing the A-orthogonality of  $d^k$  and  $d^{k-1}$ , that is:

$$\beta^k = \frac{r^{kT} A d^{k-1}}{d^{kT} A d^{k-1}}.$$

With this choice of search direction, the following remarkable property holds:

If  $r^{m+1} \neq 0$  then, for  $i < j \leq m+1$ ,

$$\boxed{\begin{array}{l} d^{i^T} A d^j = 0, \\ r^{i^T} r^j = 0, \\ d^{i^T} r^j = 0. \end{array}}$$

the most important consequence of this property is the following result:

Theorem. The conjugate gradient method converges in at most  $n$  iterations.

Proof. Suppose the result not true. then  $r^n \neq 0$ . But  $\{r^k\}_{k=0}^n$  is a set of orthogonal vectors, by our orthogonality property. Since this is not possible, we must have  $r^n = 0$ . This completes the proof.

Another consequence is that we can rewrite  $\alpha^k$  and  $\beta^k$  as follows:

$$\begin{aligned}\alpha^k &= -\frac{r^{kT} d^k}{d^{kT} A d^k} \\ &= -\frac{r^{kT} (-r^k + \beta^k d^{k-1})}{d^{kT} A d^k} \quad (\text{definition of } d^k) \\ &= \frac{r^{kT} r^k}{d^{kT} A d^k},\end{aligned}$$

$$\begin{aligned}\beta^k &= \frac{r^{kT} A d^{k-1}}{d^{k-1T} A d^{k-1}} \\ &= \frac{r^{kT} \{r^k - r^{k-1}\}}{d^{k-1T} A d^{k-1}} \cdot \frac{1}{\alpha^{k-1}} \quad (\text{def. of } \alpha^{k-1}, r^k) \\ &= \frac{r^{kT} r^k}{d^{k-1T} A d^{k-1}} \cdot \frac{d^{k-1T} A d^{k-1}}{r^{k-1T} r^{k-1}} \\ &= \frac{r^{kT} r^k}{r^{k-1T} r^{k-1}}.\end{aligned}$$

We can thus write the algorithm as follows:

$$x^0 = x_0, \quad r^0 = Ax^0 - b.$$

For  $k=0, \dots, n-1$  do

If  $r^k \approx 0$  then set  $x = x^k$  and stop.

else

$$\beta^k = \frac{r^{kT} r^k}{r^{kT} r^{k-1}} \quad (\beta^0 = 0)$$

$$d^k = -r^k + \beta^k d^{k-1} \quad (d^0 = -r^0)$$

$$p^k = Ad^k$$

$$\alpha^k = \frac{r^{kT} r^k}{d^{kT} p^k}$$

$$x^{k+1} = x^k + \alpha^k d^k$$

$$r^{k+1} = r^k + \alpha^k p^k$$

end If

end For

Thus, we can see that at each iteration, only one multiplication by the matrix  $A$  is required.

Next, we show that we do not have to carry out  $n-1$  iterations and stop the algorithm a considerable amount of iterations before. To do this, we use the fact that

$$\begin{aligned} \text{span} \{ d^0, \dots, d^{m-1} \} &= \text{span} \{ r^0, r^1, \dots, r^{m-1} \} \\ &= \text{span} \{ r^0, Ar^0, \dots, A^{m-1} r^0 \} =: W_m, \end{aligned}$$

To estimate the error  $\|x^k - x\|_A$ .

Since

$$W_k = \text{span} \{r^0, r^1, \dots, r^{k-1}\},$$

we have the following error equation:

$$e^k{}^T A w = r^k{}^T w = 0 \quad \forall w \in W_k.$$

thus:

$$(x^k - x^0)^T A w = (x - x^0)^T A w \quad \forall w \in W_k,$$

and since

$$x^k - x^0 = \sum_{j=0}^{k-1} \alpha^j d^j \in W_k,$$

we see that

$x^k - x^0$  is the  $A$ -orthogonal projection of  $x - x^0$  into  $W_k$ .

As a consequence, we have

$$\|e^k\|_A \leq \inf_{w \in W_k} \|x - x^0 - w\|_A.$$

Since

$$\begin{aligned} W_k &= \text{span} \{ r^0, A r^0, \dots, A^{k-1} r^0 \} \\ &= \text{span} \{ A(x^0 - x), A^2(x^0 - x), \dots, A^k(x^0 - x) \}, \end{aligned}$$

We can write

$$\|x^k - x\|_A \leq \inf_{p \in \mathcal{P}_k} \|p(A)(x - x^0)\|_A,$$

where  $\mathcal{P}_k$  is the set of polynomials of degree at most  $k$  such that  $p(0) = 1$ . It can be proven that

$$\|x^k - x\|_A \leq 2 \left[ \frac{\sqrt{k} - 1}{\sqrt{k} + 1} \right]^k \|x^0 - x\|_A.$$

#### ④ Proof of the orthogonality properties of the CG method.

First, we prove that the following properties are equivalent

- (1)  $r^{m+1} \neq 0$
- (2)  $\beta^{m+1} \neq 0$
- (3)  $d^{m+1} \neq 0$
- (4)  $\alpha^{m+1} \neq 0$ .

Let us prove that (1)  $\Rightarrow$  (2). Suppose that  $\beta^{m+1} = 0$ . By the orthogonality relations,

$$\beta^{m+1} = \frac{Y^{m+1T} Y^{m+1}}{Y^{m+1T} Y^m},$$

and so  $Y^{m+1} = 0$ , which is not possible. Hence, (1)  $\Rightarrow$  (2).

Let us prove that (2)  $\Rightarrow$  (3). Suppose that  $d^{m+1} = 0$ . Then  $Y^{m+1} = \beta^{m+1} d^m$  and hence,

$$\begin{aligned} Y^{m+1T} Y^{m+1} &= \beta^{m+1} Y^{m+1T} d^m \\ &= 0, \end{aligned}$$

by the orthogonality relations, this implies  $Y^{m+1} = 0$  and hence  $\beta^{m+1} = 0$ , which is not possible. Hence (2)  $\Rightarrow$  (3).

Let us prove that (3)  $\Rightarrow$  (4). Assume that  $\alpha^{m+1} = 0$ . By the orthogonality properties

$$\begin{aligned} \alpha^{m+1} &= \frac{Y^{m+1T} Y^{m+1}}{d^{m+1T} \Delta d^{m+1}} \\ &= 0, \end{aligned}$$

which implies  $Y^{m+1} = 0$ , which in turn implies  $d^{m+1} = 0$ .

Hence (3)  $\Rightarrow$  (4).



Finally, let us prove that (4)  $\Rightarrow$  (1). Assume that  $r^{m+1} = 0$ .  
 then we clearly have  $\alpha^{m+1} = 0$ , which is not possible  
 this completes the proof.

this shows that if  $r^{m+1} \neq 0$ , the corresponding new  
 step of the method is well-defined.

Now, we prove that we have

$$\begin{aligned} & \text{span} \{ d^0, d^1, \dots, d^{m+1} \} \\ &= \text{span} \{ r^0, r^1, \dots, r^{m+1} \} \\ &= \text{span} \{ r^0, Ar^0, \dots, A^m r^0 \}, \end{aligned}$$

provided that  $r^{m+1} \neq 0$ . We proceed by induction on  
 $m$ . For  $m = -1$ , the result is trivially true since  
 $d^0 = -r^0$ . Now assume the result true for  $m = n$  and  
 let us prove it is true for  $m = n+1$ . But

$$\begin{aligned} \text{span} \{ d^0, d^1, \dots, d^{n+1} \} &= \text{span} \{ d^0, d^1, \dots, -r^{n+1} + \beta_{n+1} d^n \} \\ &= \text{span} \{ d^0, d^1, \dots, d^n, -r^{n+1} \} \\ &= \text{span} \{ r^0, r^1, \dots, r^n, r^{n+1} \} \\ &= \text{span} \{ r^0, r^1, \dots, r^n, r^n + \alpha_n A d^n \} \\ & \text{(since } \alpha_n \neq 0!) \\ &= \text{span} \{ r^0, r^1, \dots, r^n, A d^n \} \\ &= \text{span} \{ r^0, Ar^0, \dots, A^n r^0, A d^n \} \\ & \text{(since } \beta^n \neq 0!) \\ &= \text{span} \{ r^0, Ar^0, \dots, A^n r^0, A^{n+1} r^0 \}. \end{aligned}$$

this proves the result.

Now, we prove the orthogonality property by induction on  $m$ . For  $m=0$ , we have

$$d^{0T} A d^1 = 0,$$

since this is enforced to compute  $\beta^1$ . We also have, that

$$\begin{aligned} r^{0T} r^1 &= r^{0T} (r^0 + \alpha^0 A d^0) \\ &= r^{0T} r^0 + \alpha^0 r^{0T} A d^0 \\ &= r^{0T} r^0 + \frac{r^{0T} r^0}{r^{0T} A r^0} \cdot (r^{0T} A (-r^0)) \\ &= 0. \end{aligned}$$

Finally

$$d^{0T} r^1 = -r^{0T} r^1 = 0.$$

Suppose that the orthogonality properties hold for  $m=n$ , and let us prove they also hold for  $m=n+1$ . We have to prove that the property holds for  $j=n+1$  and  $i \leq n$ . We have that

$$d^{nT} A d^{n+1} = 0,$$

by definition of  $\beta^{n+1}$ . Also, for  $i < n$ ,

$$\begin{aligned}
d^{iT} A d^{n+1} &= d^{iT} A (-r^{n+1} + \beta^n d^n) \\
&= -d^{iT} A r^{n+1} \\
&= \frac{1}{\alpha^i} (r^i - r^{i+1})^T r^{n+1} \\
&= 0
\end{aligned}$$

provided that  $r^{jT} r^{n+1} = 0 \quad \forall j \leq n$ . But

$$\begin{aligned}
r^{jT} r^{n+1} &= r^{jT} (r^n + \alpha^n A d^n) \\
&= r^{jT} r^n + \alpha^n r^{jT} A d^n \\
&= r^{jT} r^n + \alpha^n (-d^j + \beta^j d^{j-1}) A d^n \\
&= 0
\end{aligned}$$

for  $j < n$ , by the inductive hypothesis. For  $j = n$ , we have

$$\begin{aligned}
r^{nT} r^{n+1} &= r^{nT} r^n - \alpha^n d^n A d^n \\
&= r^{nT} r^n + r^{nT} d^n \\
&= r^{nT} (r^n + d^n) \\
&= r^{nT} d^{n-1} \cdot \beta_n \\
&= 0,
\end{aligned}$$

since  $d^{n-1} \in \text{span} \{r^0, \dots, r^{n-1}\}$  and by the inductive hypothesis.

Finally, the last property follows from the second and the "span" property.