

Lecture 17: Stability of linear equations (class)

Admin:

Recall: Singular-value decomposition

$$A = \sum_i \sigma_i \vec{v}_i \vec{u}_i^T$$

condition number
 $K(A) = \frac{\text{largest } \sigma}{\text{smallest } \sigma}$

SVD applications ...

* Solving linear equations

$$\vec{A}\vec{x} = \vec{b}$$

What is the **sensitivity**,
e.g., to numerical errors?

Find the **shortest solution**

When there is no solution,
find \vec{x} to minimize $\|\vec{A}\vec{x} - \vec{b}\|$

Least-squares regression analysis

* Rank minimization

Principal Component Analysis (PCA)

Data mining, clustering, recommendation systems,...

NUMERICAL STABILITY OF $\vec{A}\vec{x} = \vec{b}$

Systems of linear equations

Example:

$$\text{Let } A = \begin{pmatrix} 10^{-4} & 1 \\ 1 & 1 \end{pmatrix}, \vec{b} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

Reading:
 Meyer 1.6
 5.12.1-2
 Strang 7.2

How sensitive is the solution $\vec{x} = A^{-1}\vec{b}$

to small errors (e.g., numerical errors) in \vec{b} ?

$$\text{Let } B = \begin{pmatrix} 1+10^{-4} & 1 \\ 1 & 1 \end{pmatrix} \quad \begin{array}{l} 1.0001x + y = b_1 \\ x + y = b_2 \end{array}$$

How sensitive is $B^{-1}\vec{b}$ to small errors in \vec{b} ?

Why not experiment?

$$B\vec{x} = \vec{b}$$

```
>> A = [10^-4 1; 1 1];
B = [1+10^-4 1; 1 1];
b = [5; 5];
berror = 10^(-5) * randn(2,1);
>> (A\b) - (A\b+berror))
```

ans =

$$\begin{pmatrix} 1.0304e-05 \\ -7.2699e-06 \end{pmatrix}$$

← pretty small
 solution not too sensitive — good!

```
>> (B\b) - (B\b+berror))
```

ans =

$$\begin{pmatrix} -1.0303e-01 \\ 1.0304e-01 \end{pmatrix}$$

← huge error!

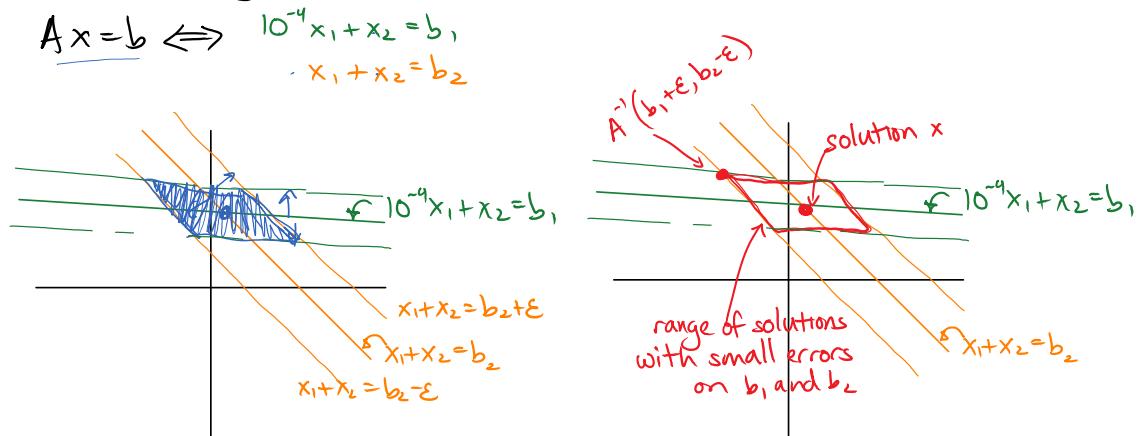
$$\begin{aligned} \vec{b}' &= \vec{b} + \vec{\delta} \\ \vec{A}^{-1}\vec{b}' &= \vec{A}^{-1}\vec{b} + \vec{A}^{-1}\vec{\delta} \\ \|\vec{x}' - \vec{x}\| &= \|\vec{A}^{-1}\vec{b} - \vec{A}^{-1}\vec{b}'\| \\ &= \|\vec{A}^{-1}(\vec{b} - \vec{b}')\| \\ &= \|\vec{A}^{-1}\vec{\delta}\| \end{aligned}$$

$$\leq \|\vec{A}^{-1}\| \cdot \|\vec{\delta}\| \quad \text{ie } \|\vec{\delta}\|_{V_2} \quad \|\vec{\delta}\|_{V_2} = \|\vec{\delta}\|$$

/ with equality if

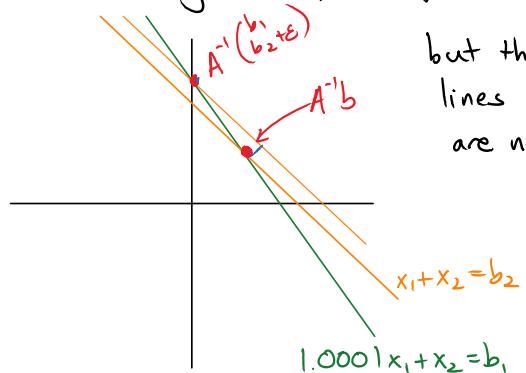
$$\begin{aligned} \|\vec{A}\| &= \sqrt{\sigma_1^2 + \sigma_2^2} \\ \vec{A}^{-1} &= \frac{1}{\sigma_1} \vec{v}_1 \vec{u}_1^T \\ &\quad + \frac{1}{\sigma_2} \vec{v}_2 \vec{u}_2^T \end{aligned}$$

Why?
 ① Geometrically



small changes to b_1, b_2 shift the lines by small amounts

but the intersection of the lines can shift a lot, if they are nearly parallel



② Algebraic explanation for stability/instability

$$(*) \quad \|A^{-1}b - A^{-1}(b + \delta)\| = \|A^{-1}\delta\| \\ \leq \|A^{-1}\| \cdot \|\delta\|$$

\Rightarrow If $\|A^{-1}\|$ is large, then a small error in b can be amplified to a large error in $x = A^{-1}b$

Continuing the previous example:

octave:7> norm(A^-1)

ans = 1.6182

octave:8> norm(B^-1)

ans = 2.0001e+04

In the worst case, the upper bound in (*) is reached!

$$A = \lambda_1 u_1 v_1^T + \lambda_2 u_2 v_2^T + \dots + \lambda_n u_n v_n^T$$

$$A^{-1} = \frac{1}{\lambda_1} v_1 u_1^T + \frac{1}{\lambda_2} v_2 u_2^T + \dots + \frac{1}{\lambda_n} v_n u_n^T$$

$$\Rightarrow \|A^{-1}\| = \lambda_n^{-1}$$

and if error vector $\vec{\delta} = \varepsilon \vec{u}_n$,

then $A^{-1}\vec{f} = \frac{1}{\|A\|} \vec{v}_n$
 $\Rightarrow \|A^{-1}\vec{f}\| = \|A^{-1}\| \|\vec{f}\|$

Easy solution?

Instead of solving $Ax = b$,

solve $(1000A)y = b \Rightarrow y = \frac{1}{1000} A^{-1}b$

and then set $x = 1000y$

$\|(1000A)^{-1}\| = \frac{1}{1000} \|A^{-1}\| \Rightarrow$ errors in b now
lead to much smaller
errors in y

→ But this is cheating; the error in x will be the same

⇒ Definition: The condition number of a matrix A is

$$K(A) = \|A\| \cdot \|A^{-1}\|$$

equivalently, $= \frac{\text{largest singular value of } A}{\text{smallest singular value}} \geq 1$

Observe: This is invariant under scaling:

moral: $\|(1000A)\| \cdot \|(1000A)^{-1}\| = \|A\| \cdot \|A^{-1}\|.$

HIGH CONDITION NUMBER
 $\frac{\max}{\min} \text{sing. value}$

$\leftrightarrow Ax = b$
UNSTABLE

Example: $K\left(\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}\right) = 1$

$$K\left(\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}\right) = 1. \quad \text{unitary} \leftrightarrow \text{singular values all} = 1$$

$$K\left(\begin{pmatrix} 10^3 & 0 \\ 0 & 1 \end{pmatrix}\right) = 10^3 = K\left(\begin{pmatrix} 1 & 0 \\ 0 & 10^{-3} \end{pmatrix}\right)$$

① consider sensitivity to errors in \vec{b}

Formally: We have

$$x = A^{-1}\vec{b} \quad x' = A^{-1}(b + \delta)$$

$$\|x - x'\| = \|A^{-1}\delta\| \leq \|A^{-1}\| \|\delta\|$$

$$\text{and } \|b\| = \|A \cdot A^{-1}b\| \leq \|A\| \cdot \|A^{-1}b\|$$

relative error in \vec{x} → $\frac{\|A^{-1}\delta\|}{\|A^{-1}b\|} \leq \underbrace{\|A\| \cdot \|A^{-1}\|}_{K(A)}$ $\frac{\|\delta\|}{\|b\|}$ } relative error in \vec{b}
correct answer → $\frac{\|A^{-1}\delta\|}{\|A^{-1}b\|} \leq \underbrace{\|A\| \cdot \|A^{-1}\|}_{K(A)} \cdot \frac{\|\delta\|}{\|b\|}$ condition #

$$\Rightarrow \frac{\|A^{-1}\delta\|}{\|A^{-1}b\|} \leq \|A\| \cdot \|A^{-1}\| \cdot \frac{\|f\|}{\|b\|}$$

relative error
 in the solution
 $(\|A^{-1}\delta\|/\|x\|)$

condition
 # of A

relative
 error in b

\Rightarrow condition number bounds the relative errors

Three sources of error in solving $Ax=b$

① Errors in the data $\rightarrow \frac{\|A^{-1}\delta\|}{\|A^{-1}x\|} \leq \kappa \cdot \frac{\|s\|}{\|b\|}$

② Errors in the linear solve routine

Even if b is known exactly, errors such as floating point errors might arise and accumulate while solving the equations. $\kappa(B) \approx 2 \cdot 10^4$

Example:

$$B = \begin{pmatrix} 1+10^{-4} & 1 \\ 1 & 1 \end{pmatrix}, \quad b = \begin{pmatrix} 5 \\ 5 \end{pmatrix} \quad \text{as above}$$

Say you solve these and got $x = \begin{pmatrix} 5 \\ 4.5 \end{pmatrix}$

$\gg B * [.5; 4.5]$

ans =

5.0000e+00
5.0000e+00

\leftarrow Looks good!

$$Bx = \begin{pmatrix} 5.00005 \\ 5 \end{pmatrix}$$

But is it?

The right answer is $\begin{pmatrix} 0 \\ 5 \end{pmatrix}$!

$\gg \text{format long e}$
 $B * [.5; 4.5]$

ans =

5.000050000000000e+00
5.000000000000000e+00

when you check the answer,
it looks great
but it is very far from
the correct answer!

This is actually the same problem we already considered.
Even though the error in b is small, the corresponding error in x is huge. This is only a good check for well-conditioned systems.

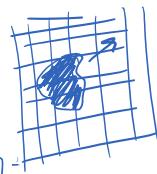
Moral: With an ill-conditioned system of linear equations

- use high-precision arithmetic
- find another way to check the answer?
- be careful! possibly try to reformulate the problem
 → use a preconditioner

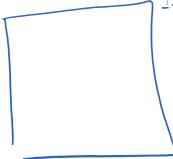
Solve $P^{-1}Ax = P^{-1}b$, where P is such that $P^{-1}A$ is better conditioned, but P^{-1} can be applied efficiently. This is an important, but more advanced, topic, that we might get to later.

$$A \vec{x} = \vec{b}$$

$$\underbrace{PAx}_{\text{hopefully}} = Pb \quad \underline{\kappa(PA) \ll \kappa(A)}$$



$$\begin{pmatrix} \frac{1}{1000} & 0 \\ 0 & \frac{1}{5} \end{pmatrix} \begin{pmatrix} 1000 & 10 \\ 5 & 1 \end{pmatrix}$$



③ Errors in the matrix A

Comparing

$$Ax = b \quad \text{and} \quad (A + \underbrace{\delta A}_{\text{some error in the matrix}})(x + \delta x) = b$$



$$\begin{aligned} A \cdot \delta x &= -\delta A(x + \delta x) \\ \delta x &= -A^{-1} \cdot \delta A \cdot (x + \delta x) \\ \frac{\|\delta x\|}{\|x + \delta x\|} &\leq \underbrace{\|A^{-1}\| \cdot \|A\|}_{\kappa} \cdot \frac{\|\delta A\|}{\|A\|} \end{aligned}$$

$$\frac{\|\delta x\|}{\|x + \delta x\|} \leq \underbrace{\|A^{-1}\| \cdot \|A\|}_{\text{relative error in the solution}} \cdot \underbrace{\frac{\|\delta A\|}{\|A\|}}_{\substack{\text{condition} \\ \# \text{ of } A}} \cdot \underbrace{\frac{\|\delta A\|}{\|A\|}}_{\text{relative error in } A}$$

Example:

```
octave:11> Aerror = 10^(-5) * randn(2,2);
octave:12> (A\b) - ((A+Aerror)\b)
ans =
```

```
-1.1004e-04
7.6036e-05
```

```
octave:13> (B\b) - ((B+Berror)\b)
ans =
```

```
-210.22
210.26
```

Moral: In all three cases, the (relative) effect of the error is bounded by the condition number

$$= \|A\| \cdot \|A^{-1}\| = \frac{\text{largest singular value}}{\text{smallest singular value}}$$

1 1 1 1 1 1

smallest singular value

Note: In practice, computing the condition number can be a lot of work.

It is enough usually to approximate it, or give upper bounds.

Certain classes of matrices reliably have large condition numbers; with experience you learn when to be careful.

Example: What is the condition number of

$$A = \begin{pmatrix} -2 & 1 & & & 1 \\ 1 & -2 & 1 & & \\ & 1 & -2 & 1 & \\ & & 1 & -2 & 1 \\ & & & \ddots & \\ & & & & 1 & -2 & 1 \\ & & & & & 1 & -2 & 1 \end{pmatrix} \quad ? \quad A \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \vec{0}$$

$$\sigma_1 = \|A\| = 4$$

$$\sigma_n = 0$$

A^{-1} doesn't exist

$$\kappa(A) = \frac{\sigma_1}{\sigma_n} = \infty$$

$$B = \begin{pmatrix} 4 & 0 \\ 0 & 0 \end{pmatrix} \quad \kappa = \infty$$

$$\sigma_1 = 4$$

$$\sigma_2 = 0$$

Example: What is the condition number of

$$C = \begin{pmatrix} -2 & 1 & & & 0 \\ 1 & -2 & 1 & & \\ & 1 & -2 & 1 & \\ & & 1 & -2 & 1 \\ & & & \ddots & \\ & & & & 1 & -2 & 1 \\ & & & & & 1 & -2 & 1 \end{pmatrix}_{n \times n} \quad N(C) = \{\vec{0}\}$$

$$C = -2I + U + U^T \Rightarrow \|C\| \leq 2 + 2\|U\| = 4$$

$$\begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

$$\|C^{-1}\| = \frac{1}{\min_{x: \|x\|=1} \|Cx\|}$$

$$C = \sigma_1 \vec{v}_1 \vec{u}_1^T + \dots + \sigma_n \vec{v}_n \vec{u}_n^T$$

$$\sigma_+ = \min_{x: \|x\|=1} \|Cx\| = \sigma_n$$

$$\|C \frac{1}{\sqrt{n}}\| = \left\| \begin{pmatrix} -1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & & & & -1 \end{pmatrix} \right\| = \sqrt{\frac{2}{n}} \Rightarrow \sigma_n \leq \sqrt{\frac{2}{n}}$$

for $x = \underline{v}_n$

$$\Rightarrow \frac{1}{\sigma_n} \geq \sqrt{\frac{n}{2}} \Rightarrow \frac{\sigma_1}{\sigma_n} \geq 2\sqrt{n}$$

Answer: Largest s.v. = norm ≈ 4 still.

Smallest singular value?

Observe:

$x = \frac{1}{\sqrt{n}}(1, 1, \dots, 1)$ has norm 1

$Ax = (-\frac{1}{\sqrt{n}}, 0, 0, 0, \dots, 0, \frac{1}{\sqrt{n}})$

has norm $\sqrt{\frac{2}{n}}$

\Rightarrow smallest singular value $\leq \sqrt{\frac{2}{n}}$

\Rightarrow condition # $\gtrsim \sqrt{n}$

So be careful!

In fact, the condition # is $\sim \frac{4}{\pi^2}(n+1)^2 = \Theta(n^2)$

A better test vector is $x = \frac{1}{\sqrt{n+1}}(\underbrace{\sin \frac{\pi}{n+1}}, \underbrace{\sin \frac{2\pi}{n+1}}, \dots, \underbrace{\sin \frac{n\pi}{n+1}})$

ANOTHER REASON THE CONDITION NUMBER IS IMPORTANT:

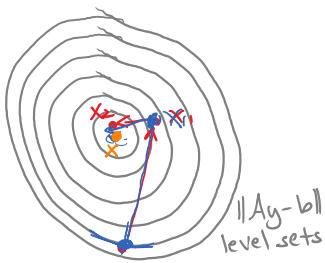
Poorly Conditioned Systems of Equations

TAKE LONGER TO SOLVE

Rough intuition: Iterative solvers for $Ax=b$ start with a guess, and try to improve it.

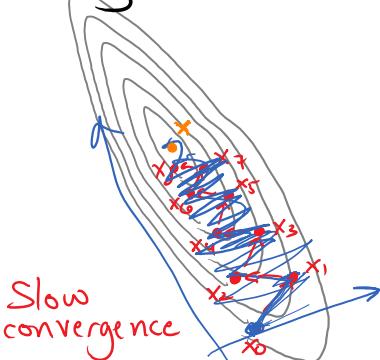
$$x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow x_3 \rightarrow \dots$$

Well-conditioned A



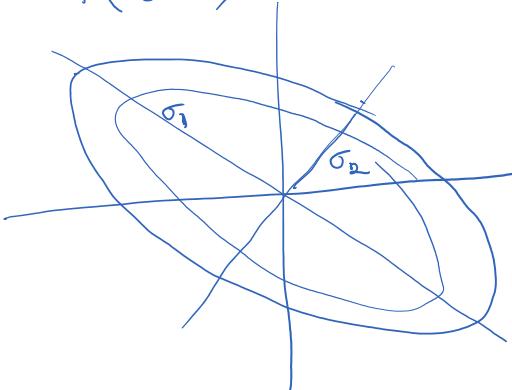
Fast convergence

Poorly conditioned A



Slow convergence

$$A \begin{pmatrix} \cos \theta & \\ & \sin \theta \end{pmatrix}$$



$$\|Cx - b\|^2$$

Theorem:

Conjugate-gradient method on a positive semi-definite matrix A , with condition number K , takes

$$\mathcal{O}(\sqrt{K} \cdot \log \frac{1}{\epsilon})$$

steps to get within ϵ of the true solution x .
(Or, $O(K \log \frac{1}{\epsilon})$ steps for a general matrix A .)

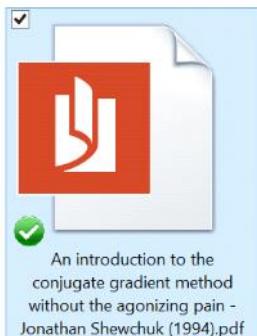
$$\text{or } O(K \log \frac{1}{\epsilon})$$

for a general A

Each step uses a matrix-vector multiplication
 $= \mathcal{O}(\# \text{ of nonzero entries in } A)$ steps.

Time permitting, we'll discuss this later.

Or, see ↗



Example: GRADIENT DESCENT

Goal: Solve $A\vec{x} = \vec{b}$.

Define the cost function

$$\begin{aligned} C(\vec{x}) &= \|A\vec{x} - \vec{b}\|^2 \\ &= \|A(\vec{x} - \vec{x}^*)\|^2 \quad \text{where } A\vec{x}^* = \vec{b} \\ &= (\vec{x} - \vec{x}^*)^T \underbrace{A^T A}_{\text{positive semi-definite}} (\vec{x} - \vec{x}^*) \end{aligned}$$

If A has SVD

$$A = \sum_i \lambda_i u_i v_i^T$$

$$A^T A = \sum_i \lambda_i^2 v_i v_i^T$$

$$\Rightarrow C(\vec{x}) = \sum_i \lambda_i^2 |v_i \cdot (\vec{x} - \vec{x}^*)|^2$$

\Rightarrow level set $\{\vec{x} \mid C(\vec{x}) = b\}$ is an ellipse!
 (like $\lambda_1^2 x_1^2 + \lambda_2^2 x_2^2 + \dots + \lambda_n^2 x_n^2 = b$)

Algorithm: Repeat until convergence:

$x^t = x^{t-1} - \alpha \left(\frac{\partial C}{\partial x_1}(x^{t-1}), \dots, \frac{\partial C}{\partial x_n}(x^{t-1}) \right)$
 i.e., step in the direction opposite the largest increase in C .

$$\begin{aligned} \frac{\partial C}{\partial x_j} &= \frac{\partial}{\partial x_j} \left[(Ax - b)^T (Ax - b) \right] \\ &= \frac{\partial}{\partial x_j} \sum_i (Ax - b)_i^2 \\ &= 2 \sum_i (Ax - b)_i \cdot \frac{\partial}{\partial x_j} [(Ax - b)_i] \\ &= 2 \sum_i (Ax - b)_i \cdot a_{ij} \\ \Rightarrow \left(\frac{\partial C}{\partial x_1}, \dots, \frac{\partial C}{\partial x_n} \right) &= 2 \sum_i \underbrace{(Ax - b)_i \cdot (a_{1i}, \dots, a_{ni})}_{(Ax - b)^T e_i} \cdot e_i^T A \\ &= 2(Ax - b)^T \left(\sum_i e_i e_i^T \right) A \\ &= 2(Ax - b)^T A \end{aligned}$$

(As you might guess!)

Example:

$$A = \begin{pmatrix} 2 & -1.5 \\ -1.5 & 4 \end{pmatrix};$$

```


$$A = \begin{pmatrix} 2 & -1.5 \\ -1.5 & 4 \end{pmatrix};$$

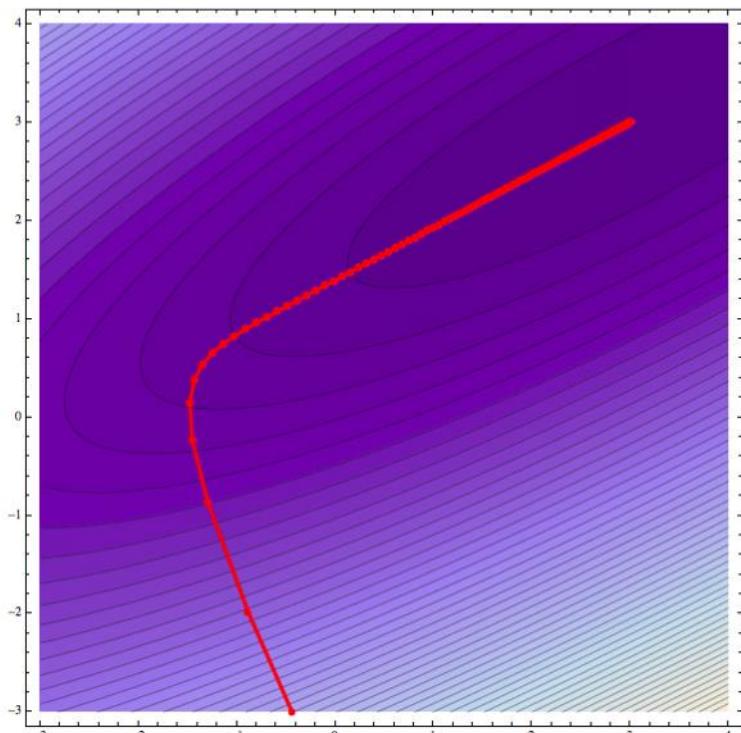
b = A.{3, 3};
Eigenvalues[A]

\alpha = .01;
x = {0, -4};
list = {x};
For[k = 1, k \leq 1000, k++,
  x = x - 2 \alpha (A.x - b).A;
  AppendTo[list, x];
];
"The solution is:"
x
(4.80278, 1.19722)

The solution is:
(3., 3.)

"Display the results:";
iter = ListPlot[list, PlotRange \rightarrow 3 {{-1, 1}, {-1, 1}}, Joined \rightarrow True,
  PlotStyle \rightarrow {Red, Thickness[.005]}, PlotMarkers \rightarrow Automatic];
contour = ContourPlot[Norm[A.{x, y} - b]^2, {x, -3, 4}, {y, -3, 4}, Contours \rightarrow 50];
Show[contour, iter]

```



For small enough α , this has to converge, since there is a unique local and global minimum.

Observe: Smaller condition number
 \Rightarrow Rounder ellipse \Rightarrow Faster convergence

Remark: Although not generally the fastest way of solving a set of linear equations, the gradient descent approach is very robust, and many variants are used in many applications.