

Root Finding

goal: sqrt of $a > 0$

$$y_0 = 1 \quad \boxed{A = a} \\ x_0 = a$$

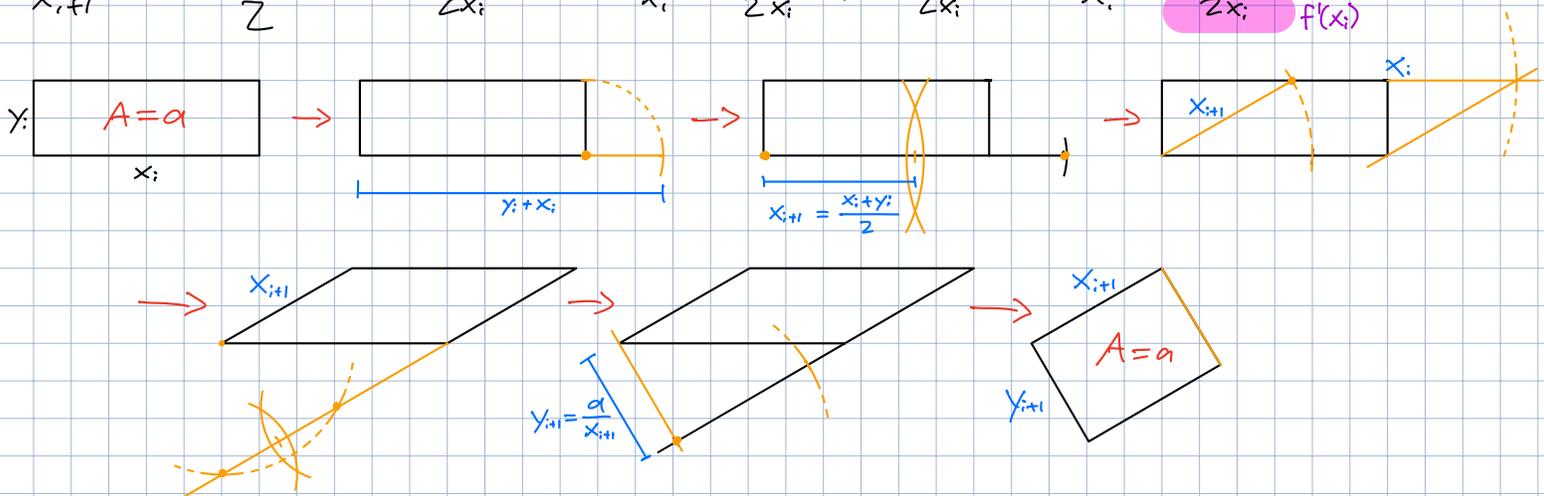
idea: successively transform conserving area until perfect square

straightedge and compass approach:

$$x_{i+1} = \frac{x_i + y_i}{2}$$

since $A = x_i \cdot y_i \stackrel{!}{=} a \Rightarrow y_i = \frac{a}{x_i}$ (can always express like this):

$$x_{i+1} = \frac{x_i + \frac{a}{x_i}}{2} = \frac{x_i^2 + a}{2x_i} = x_i - \frac{2x_i^2}{2x_i} + \frac{x_i^2 + a}{2x_i} = x_i - \frac{x_i^2 - a}{2x_i} \stackrel{f(x_i)}{=} x_i - \frac{f(x_i)}{f'(x_i)}$$



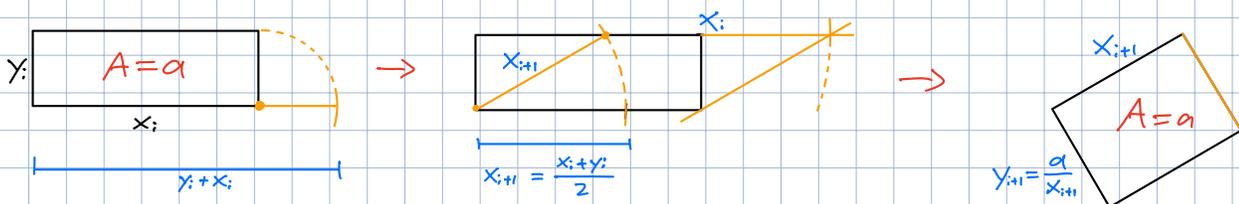
special case of Newton's Method: finding $x \in \mathbb{R}$ st $f(x) = 0$

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$$

$e_i = |x - x_i|$ decreases quadratically in i (show using Taylor Expansion)

BUT careful: can also fail badly if

- x_0 is too far from exact solution x
- $f'(x_i) = 0$ for some x_i
- f' unbounded close to the root x



Linear Systems

LU Decomposition

using Gaussian Elimination, any A can be decomposed into $A = LU$, where L is a lower triangular and U is an upper triangular matrix ^{and "partial pivoting"}

a linear system $A\vec{x} = \vec{b}$ can then be solved:

1. solve $L\vec{z} = \vec{b}$, 2. solve $U\vec{x} = \vec{z}$ "backsubstitution"

$$A = LU$$

decomposition: $O(n^3)$, backsubstitution: $O(n^2)$

sparse system: A has only $O(n)$ nonzero entries

for sparse systems often better to use an iterative method

Iterative Methods

Jacobi Method

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j \neq i} a_{ij} x_j^{(k)} \right), \quad i = 1, \dots, n$$

or in matrix notation:

$$\vec{x}^{(k+1)} = D^{-1} (\vec{b} - L\vec{x}^{(k)} - U\vec{x}^{(k)}), \quad A = L + D + U =$$



can be improved if always the most recent values of x_i are used

Gauss-Seidel Method

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right)$$

or

$$\vec{x}^{(k+1)} = D^{-1} (\vec{b} - L\vec{x}^{(k+1)} - U\vec{x}^{(k)}) \Leftrightarrow \vec{x}^{(k+1)} = (D+L)^{-1} (\vec{b} - U\vec{x}^{(k)})$$

advantages:

- converges faster (often)
- easier to implement, because $x_i^{(k)}$ can simply be overwritten by $x_i^{(k+1)}$

(... further improvement can be obtained by Successive Over-Relaxation (SOR) ...)

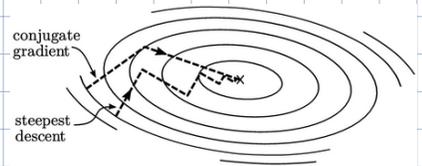
diagonal dominance: $|a_{ii}| > \sum_{j \neq i} |a_{ij}|, \quad i = 1, \dots, n$

both methods are guaranteed to converge to correct solution if A is diagonally dominant
 → Gauss-Seidel also converges if A is spd

cost for each iteration is $O(n^2)$ and only $O(n)$ if A is sparse

Symmetric Positive-Definite Matrices

symmetric: $A^T = A$
 positive def: $\vec{x}^T A \vec{x} > 0$ for $\vec{x} \neq \vec{0}$



Gauss-Siedel also converges if A is spd, but in this case an even better iterative method exists

for such matrices, the solution of $A\vec{x} = \vec{b}$ is also the minimizer of $f(\vec{x}) = \frac{1}{2} \vec{x}^T A \vec{x} - \vec{x}^T \vec{b}$ because $\nabla f(\vec{x}) = \frac{1}{2}(A+A^T)\vec{x} - \vec{b} \stackrel{A=A^T}{=} A\vec{x} - \vec{b}$

Gradient Descent \rightarrow simplest

$$\vec{x}_{k+1} = \vec{x}_k + \alpha_k \vec{r}_k \quad \text{with} \quad \vec{r}_k = -\nabla f(\vec{x}_k) = \vec{b} - A\vec{x}_k$$

α_k is chosen such that $f(\vec{x}_{k+1})$ is minimal wrt α_k :

$$\alpha_k = \frac{\vec{r}_k^T \vec{r}_k}{\vec{r}_k^T A \vec{r}_k}$$

the repeated gradients are not in the least direction \rightarrow not orthogonal

$$\begin{aligned} f(\vec{x}_{k+1}) &= \frac{1}{2} \vec{x}_{k+1}^T A \vec{x}_{k+1} - \vec{x}_{k+1}^T \vec{b} \\ &= \frac{1}{2} (\vec{x}_k + \alpha_k \vec{r}_k)^T A (\vec{x}_k + \alpha_k \vec{r}_k) - (\vec{x}_k + \alpha_k \vec{r}_k)^T \vec{b} \\ \frac{\partial f}{\partial \alpha_k} &= (A(\vec{x}_k + \alpha_k \vec{r}_k))^T \vec{r}_k - \vec{r}_k^T \vec{b} = (A\vec{x}_k + \alpha_k A \vec{r}_k)^T \vec{r}_k - \vec{r}_k^T \vec{b} = \vec{r}_k^T A \vec{x}_k + \alpha_k \vec{r}_k^T A \vec{r}_k - \vec{r}_k^T \vec{b} \\ &\stackrel{!}{=} 0 \Rightarrow \vec{r}_k^T \vec{b} = \vec{r}_k^T A \vec{x}_k + \alpha_k \vec{r}_k^T A \vec{r}_k \\ \rightarrow \alpha_k &= \frac{\vec{r}_k^T \vec{b} - \vec{r}_k^T A \vec{x}_k}{\vec{r}_k^T A \vec{r}_k} = \frac{\vec{r}_k^T (\vec{b} - A\vec{x}_k)}{\vec{r}_k^T A \vec{r}_k} = \frac{\vec{r}_k^T \vec{r}_k}{\vec{r}_k^T A \vec{r}_k} \end{aligned}$$

Conjugate Gradient Descent

idea: take directions d_k which are pairwise conjugate wrt A
 two vectors are A -conjugate if $\vec{v}^T A \vec{w} = 0$ (note: for $A=I$, this is standard orthogonality)

$$\vec{x}_0, \vec{r}_0 = \vec{b} - A\vec{x}_0, \vec{d}_0 = \vec{r}_0$$

then iterate

$$\begin{aligned} \vec{x}_{k+1} &= \vec{x}_k + \alpha_k d_k \\ \alpha_k &= \vec{r}_k^T d_k / d_k^T A d_k = \vec{r}_k^T \vec{r}_k / d_k^T A d_k \\ \vec{r}_{k+1} &= \vec{b} - A\vec{x}_{k+1} = \vec{b} - A(\vec{x}_k + \alpha_k d_k) = \vec{r}_k - \alpha_k A d_k \\ d_{k+1} &= \vec{r}_{k+1} + \beta_k d_k \\ \beta_k &= -\vec{d}_k^T A \vec{r}_{k+1} / \vec{d}_k^T A \vec{d}_k \end{aligned}$$

$$\alpha_k = \vec{d}_k^T \vec{r}_k / \vec{d}_k^T A \vec{d}_k \quad (\text{chosen again such that } f(\vec{x}_{k+1}) \text{ is minimal})$$

$$\beta_k \text{ is chosen such that } \vec{d}_{k+1} \text{ is } A\text{-conjugate to } \vec{d}_k: \vec{d}_{k+1}^T A \vec{d}_k = (\vec{r}_{k+1} + \beta_k d_k)^T A \vec{d}_k = \vec{r}_{k+1}^T A \vec{d}_k + \beta_k \vec{d}_k^T A \vec{d}_k \stackrel{!}{=} 0$$

the following subspaces of \mathbb{R}^n are equal: $\langle \vec{x}_1, \dots, \vec{x}_k \rangle = \langle \vec{r}_0, \dots, \vec{r}_{k-1} \rangle = \langle \vec{d}_0, \dots, \vec{d}_{k-1} \rangle$

residuals \vec{r}_i are pairwise orthogonal: $\vec{r}_i^T \vec{r}_j = 0$ for $j < i$

directions \vec{d}_i are pairwise conjugate: $\vec{d}_i^T A \vec{d}_j = 0$ for $j < i$

residual \vec{r}_{k+1} is orthogonal to \vec{d}_k since $\vec{r}_{k+1}^T \vec{d}_k = -(\vec{x}_{k+1}^T A - \vec{b}^T) \vec{d}_k = -\frac{\partial f}{\partial \alpha_k} \stackrel{\text{by construction}}{=} 0$

\vec{r}_k also orthogonal to previous \vec{d}_i 's: $\vec{r}_k^T \vec{d}_i = 0$ for $j < k$

\vec{d}_k is a linear combination of the residuals $\{\vec{r}_0, \dots, \vec{r}_k\}$: $\vec{d}_k = \sum_{i=0}^k \left(\prod_{j=i+1}^k \beta_j \right) \vec{r}_i$

Arnoldi: $AQ = QH$
 Gram-Schmidt
 Krylov subspace $H = Q^T A Q$
 high accuracy not so reliable but fast values
 similar matrices \rightarrow same eigenvalues

$x_k \in \mathcal{K}_k$ (Krylov space)
 $b \in \mathcal{K}_k$
 $r_k = b - Ax_k \in \mathcal{K}_k$

residuals are orthogonal

① $r_k^T r_i = 0 \quad i < k$ orthogonal in usual sense (since here A already kills the form)
 ② $(r_k - x_k)^T (Ax_k - x_k) = 0 \quad i < k$ orthogonal in A -inner-product

going into cycle with new d : search direction, always α_k is going to go
 α_k is how long along we go in this direction d
 change in r is "step" if I know that things are cyc I know the length in r
 but now we need to look ahead: find direction via local next step number β and search direction
 \rightarrow search direction is in r but with one correction

- 1 $x_k = \frac{r_k^T r_k}{r_k^T A r_k} r_k$
- 2 $x_k = \frac{r_k^T r_k}{r_k^T A r_k} r_k$
- 3 $r_k = r_k - Ax_k$
- 4 $\beta_k = \frac{r_k^T r_{k+1}}{r_k^T A r_k}$
- 5 $d_k = r_k + \beta_k d_{k-1}$ new direction that is A -orthogonal to ALL previous directions (②)

Polynomial Interpolation

most common way of expressing polynomial is to use monomial basis: $\{1, x, x^2, \dots, x^n\}$

$$p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$$

can be evaluated efficiently using Horner's scheme in $O(n)$ operations:

$p := a_n$
for i : from $n-1$ to 0 do
 $p := p x + a_i$
return p

Polynomial Curves $\vec{P}: \mathbb{R} \rightarrow \mathbb{R}^d$

$$\vec{P}(t) = \vec{A}_n t^n + \vec{A}_{n-1} t^{n-1} + \dots + \vec{A}_1 t + \vec{A}_0$$

Horner also works etc...

often interested in interpolating polynomial curve $\vec{P}(t_i) = \vec{P}_i, i = 0, 1, \dots, n$

can be written as linear system

$$\underline{V} \underline{A} = \underline{P}; \quad \underline{V} = \begin{matrix} \text{Vandermonde} \\ \begin{bmatrix} 1 & t_0 & \dots & t_0^n \\ \vdots & \vdots & \dots & \vdots \\ \vdots & \vdots & \dots & \vdots \\ 1 & t_n & \dots & t_n^n \end{bmatrix} \end{matrix} \in \mathbb{R}^{(n+1) \times (n+1)}, \quad \underline{A} = \begin{bmatrix} -\vec{A}_n^T \\ -\vec{A}_{n-1}^T \\ \vdots \\ -\vec{A}_0^T \end{bmatrix} \in \mathbb{R}^{(n+1) \times d}, \quad \underline{P} = \begin{bmatrix} -\vec{P}_0^T \\ -\vec{P}_1^T \\ \vdots \\ -\vec{P}_n^T \end{bmatrix} \in \mathbb{R}^{(n+1) \times d}$$