

Gradient Methods (§8.1)

- Given $\mathbf{x}^{(k)}$.
- The vector $-\nabla f(\mathbf{x}^{(k)})$ points in the direction of maximum rate of decrease.
- Makes sense to choose $\mathbf{d}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$.
- Gradient algorithm:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)}).$$

- Step size α_k can be chosen in many different ways.
- For sufficiently small step size, the gradient algorithm has descent property.
- Prop.: Suppose $\nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$. There exists $\bar{\alpha} > 0$ such that for all $\alpha_k \in (0, \bar{\alpha})$, we have

$$f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)}).$$

- Remark: if $\nabla f(\mathbf{x}^{(k)}) = \mathbf{0}$, the FONC holds. Can use as basis for stopping.

Proof of prop.:

- Proof: Consider $\phi(\alpha) = f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}))$.
- By chain rule, we have

$$\phi'(0) = -\|\nabla f(\mathbf{x}^{(k)})\|^2 < 0.$$

- Hence, there exists $\bar{\alpha} > 0$ such that for all $\alpha_k \in (0, \bar{\alpha})$, we have

$$\phi(\alpha_k) < \phi(0).$$

- Rewriting, we obtain

$$f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)}).$$

- Several possible choices for α_k .
- If α_k too small, we need to iterate many times to get to the solution.
- If α_k too big, algorithm may zig-zag around the solution (overshoot).
- We can either fix $\alpha_k = \alpha$ for all k , or let α_k vary from iteration to iteration.
- Greedy scheme:

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)})).$$

Name: *Steepest descent algorithm*

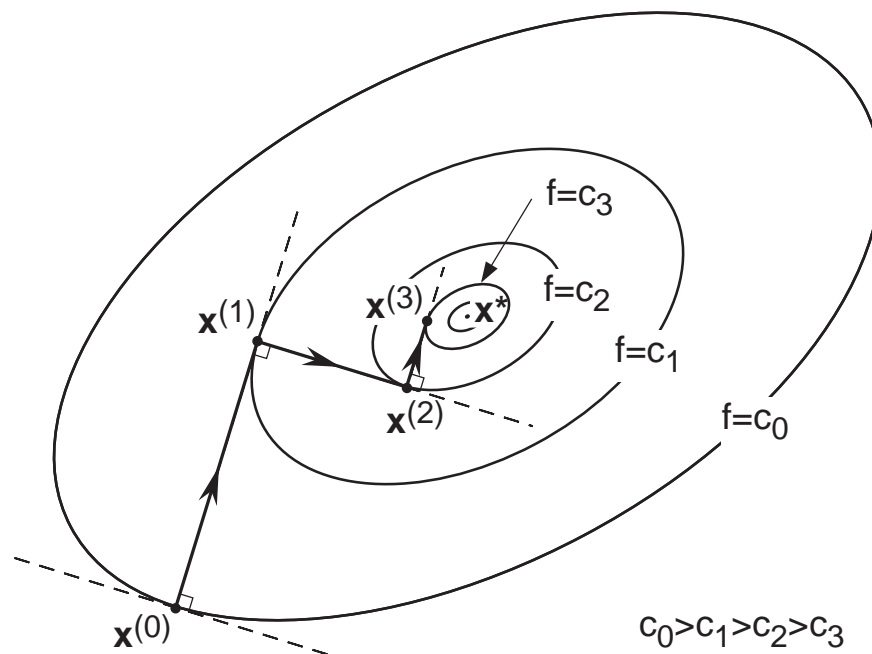
Steepest descent algorithm (§8.2)

- See Example 8.1.
- The steepest descent algorithm has the descent property. Why?
- Prop. (8.1): In the steepest descent algorithm,

$$\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$$

is orthogonal to

$$\mathbf{x}^{(k+2)} - \mathbf{x}^{(k+1)}.$$



Proof of prop.:

- We have

$$\begin{aligned} \langle \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}, \mathbf{x}^{(k+2)} - \mathbf{x}^{(k+1)} \rangle \\ = \alpha_k \alpha_{k+1} \langle \nabla f(\mathbf{x}^{(k)}), \nabla f(\mathbf{x}^{(k+1)}) \rangle. \end{aligned}$$

- To complete the proof it is enough to show that

$$\langle \nabla f(\mathbf{x}^{(k)}), \nabla f(\mathbf{x}^{(k+1)}) \rangle = 0.$$

- Let $\phi(\alpha) = f(\mathbf{x}^{(k)} - \alpha \nabla f(\mathbf{x}^{(k)}))$.
- By FONC, $\phi'(\alpha_k) = 0$.

- By chain rule, $\nabla f(\mathbf{x}^{(k)})^T \nabla f(\mathbf{x}^{(k+1)}) = 0$.

- Typical stopping criteria:

$$\|\nabla f(\mathbf{x}^{(k)})\| \leq \varepsilon,$$

or

$$\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| \leq \varepsilon,$$

or

$$\frac{\|\nabla f(\mathbf{x}^{(k)})\|}{\|\nabla f(\mathbf{x}^{(0)})\|} \leq \varepsilon,$$

or

$$\frac{\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|}{\|\mathbf{x}^{(k)}\|} \leq \varepsilon,$$

where ε is prespecified.

- The latter two are usually preferable, because they are “scale-free.”

Analysis of optimization algorithms

- Rely heavily on mathematical tools.
- “Do we really have to go through this?”
- Analysis provides insight into:
 - Range of applicability of an algorithm.
 - Appropriate choice of algorithm for a given problem.
 - Qualitative behavior of an algorithm.
- We must be able to answer:
 - Does the method work?
 - When does it work?
 - How well does it work?
- Not good enough to superficially use commercial optimization software package.

Several characterizations of performance:

- *Globally convergent*: starting from any initial point, the algorithm converges to a “solution.”
- Usually, by “solution” we mean a point satisfying the FONC.
- *Locally convergent*: starting from an initial point that is close enough to a solution, the algorithm converges to the solution.
- *Rate of convergence*: how fast an algorithm converges.

Analysis of gradient methods (§8.3)

- We analyze gradient algorithms applied to quadratics only:

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

where $\mathbf{Q} > 0$.

- We restrict our attention to quadratics because:
 - Simplifies analysis.
 - Local behavior near solution. (Global convergence for quadratics tells us something about local convergence in more general functions. How?)

Steepest descent method applied to quadratics

- Consider

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{b}^T \mathbf{x},$$

where $\mathbf{Q} > 0$.

- We have $\nabla f(\mathbf{x}) = \mathbf{Q} \mathbf{x} - \mathbf{b}$ and $\mathbf{F}(\mathbf{x}) = \mathbf{Q}$.
- For simplicity, write $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$.
- We can find an explicit formula for α_k .
- Let $\phi(\alpha) = f(\mathbf{x}^{(k)} - \alpha \mathbf{g}^{(k)})$.
- $\phi(\alpha)$ is a quadratic:

$$\phi(\alpha) = \left(\frac{1}{2} \mathbf{g}^{(k)T} \mathbf{Q} \mathbf{g}^{(k)} \right) \alpha^2 - (\mathbf{g}^{(k)T} \mathbf{g}^{(k)}) \alpha + \text{constant}.$$

- Hence, we get

$$\alpha_k = \frac{\mathbf{g}^{(k)T} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)T} \mathbf{Q} \mathbf{g}^{(k)}}.$$

- Algorithm applied to quadratic:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left(\frac{\mathbf{g}^{(k)T} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)T} \mathbf{Q} \mathbf{g}^{(k)}} \right) \mathbf{g}^{(k)}.$$

- For convenience, instead of working with f , we work with

$$\begin{aligned} V(\mathbf{x}) &= f(\mathbf{x}) + \frac{1}{2} \mathbf{x}^{*T} \mathbf{Q} \mathbf{x}^* \\ &= \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^T \mathbf{Q} (\mathbf{x} - \mathbf{x}^*), \end{aligned}$$

where $\mathbf{x}^* = \mathbf{Q}^{-1} \mathbf{b}$.

- The constant we add to f does not change solution. Why?

- Lemma (8.1): We have

$$V(\mathbf{x}^{(k+1)}) = (1 - \gamma_k) V(\mathbf{x}^{(k)}),$$

where γ_k is defined as

$$\gamma_k = \alpha_k \frac{\mathbf{g}^{(k)T} \mathbf{Q} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)T} \mathbf{Q}^{-1} \mathbf{g}^{(k)}} \left(2 \frac{\mathbf{g}^{(k)T} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)T} \mathbf{Q} \mathbf{g}^{(k)}} - \alpha_k \right)$$

if $\mathbf{g}^{(k)} \neq \mathbf{0}$, and $\gamma_k = 1$ if $\mathbf{g}^{(k)} = \mathbf{0}$.

- Proof: By substitution and algebraic manipulations.

Remarks:

- γ_k is simply a (complicated) function of α_k .
- Note that $\gamma_k \leq 1$ always.
- $\gamma_k = 1$ implies that $V(\mathbf{x}^{k+1}) = 0$, which means $\mathbf{x}^{k+1} = \mathbf{x}^*$.
- The previous lemma has the following strong consequence.
- Theorem (8.1): Suppose $\gamma_k > 0$ for all k . Then, $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ for any initial condition $\mathbf{x}^{(0)}$ if and only if

$$\sum_{k=0}^{\infty} \gamma_k = \infty.$$

- To apply the theorem, we just check if our step size sequence $\{\alpha_k\}$ satisfies the above.

Proof of theorem:

- Note that $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ if and only if $V(\mathbf{x}^{(k)}) \rightarrow 0$.
- By previous lemma,

$$V(\mathbf{x}^{(k)}) = \left(\prod_{i=0}^{k-1} (1 - \gamma_i) \right) V(\mathbf{x}^{(0)}).$$

- Assume $\gamma_k < 1$ (otherwise, the result holds trivially).
- Hence,

$$\begin{aligned} \mathbf{x}^{(k)} \rightarrow \mathbf{x}^* \text{ for all } \mathbf{x}^{(0)} &\Leftrightarrow \prod_{i=0}^{\infty} (1 - \gamma_i) = 0 \\ &\Leftrightarrow \sum_{i=0}^{\infty} \gamma_i = \infty, \end{aligned}$$

where the last line is obtained by taking logs.

Application of convergence theorem

- We can apply the previous theorem to answer the following questions (for quadratics):
 - Is the steepest descent algorithm globally convergent?
 - In a fixed step size gradient algorithm (i.e., $\alpha_k = \alpha$ for all k), for what values of the step size α is the algorithm globally convergent?

Convergence of steepest descent algorithm

- We now apply the previous theorem to show convergence of the steepest descent algorithm.
- Recall that in this case α_k is given by

$$\alpha_k = \frac{\mathbf{g}^{(k)T} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)T} \mathbf{Q} \mathbf{g}^{(k)}}.$$

- Substituting into the formula for γ_k yields

$$\gamma_k = \frac{(\mathbf{g}^{(k)T} \mathbf{g}^{(k)})^2}{(\mathbf{g}^{(k)T} \mathbf{Q} \mathbf{g}^{(k)}) (\mathbf{g}^{(k)T} \mathbf{Q}^{-1} \mathbf{g}^{(k)})}.$$

Rayleigh's inequality

- Given $\mathbf{Q} > 0$.
- Rayleigh's inequality:

$$\lambda_{\min}(\mathbf{Q}) \|\mathbf{x}\|^2 \leq \mathbf{x}^T \mathbf{Q} \mathbf{x} \leq \lambda_{\max}(\mathbf{Q}) \|\mathbf{x}\|^2,$$

where

$\lambda_{\min}(\mathbf{Q})$ is the smallest eigenvalue of \mathbf{Q}

$\lambda_{\max}(\mathbf{Q})$ is the largest eigenvalue of \mathbf{Q}

- See p. 34.
- Note that $\lambda_{\max}(\mathbf{Q}^{-1}) = 1/\lambda_{\min}(\mathbf{Q})$.
- Applying Rayleigh's inequality to \mathbf{Q} and \mathbf{Q}^{-1} , we obtain

$$\gamma_k \geq \frac{\lambda_{\min}(\mathbf{Q})}{\lambda_{\max}(\mathbf{Q})} > 0.$$

- Hence, $\gamma_k > 0$ for all k , and also

$$\sum_{k=0}^{\infty} \gamma_k = \infty.$$

- By previous theorem, $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ for all $\mathbf{x}^{(0)}$.
[Theorem 8.2].

Convergence of fixed step size gradient algorithm

- Consider the case where we fix $\alpha_k = \alpha$ for all k .
- Theorem (8.3): $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$ for any $\mathbf{x}^{(0)}$ if and only if

$$0 < \alpha < \frac{2}{\lambda_{\max}(\mathbf{Q})}.$$

- The theorem gives some idea of how large the step size is allowed to be.

Proof of theorem:

- We have

$$\gamma_k = \alpha \frac{\mathbf{g}^{(k)T} \mathbf{Q} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)T} \mathbf{Q}^{-1} \mathbf{g}^{(k)}} \left(2 \frac{\mathbf{g}^{(k)T} \mathbf{g}^{(k)}}{\mathbf{g}^{(k)T} \mathbf{Q} \mathbf{g}^{(k)}} - \alpha \right).$$

- To prove \Leftrightarrow , we have two directions to prove:
 \Rightarrow and \Leftarrow .

To prove \Leftarrow :

- Apply Rayleigh's inequality to get

$$\gamma_k \geq \alpha (\lambda_{\min}(\mathbf{Q}))^2 \left(\frac{2}{\lambda_{\max}(\mathbf{Q})} - \alpha \right) > 0.$$

- Hence, $\gamma_k > 0$ for all k , and also

$$\sum_{k=0}^{\infty} \gamma_k = \infty.$$

To prove \Rightarrow :

- Use contraposition.
- Suppose either $\alpha \leq 0$ or $\alpha \geq 2/\lambda_{\max}(\mathbf{Q})$.
- It suffices to find a single $\mathbf{x}^{(0)}$ that makes the algorithm not converge.
- So, choose $\mathbf{x}^{(0)}$ such that $\mathbf{x}^{(0)} - \mathbf{x}^*$ is an eigenvector of \mathbf{Q} corresponding to the eigenvalue $\lambda_{\max} = \lambda_{\max}(\mathbf{Q})$.
- Note that for all k , $\mathbf{x}^{(k)} - \mathbf{x}^*$ is an eigenvector of \mathbf{Q} corresponding to λ_{\max} .
- Moreover, for all k , $\mathbf{g}^{(k)}$ is an eigenvector of \mathbf{Q} corresponding to λ_{\max} .
- Hence, using the formula for γ_k , we get

$$\gamma_k = \alpha \lambda_{\max}^2 \left(\frac{2}{\lambda_{\max}} - \alpha \right) \leq 0.$$

- By previous lemma, $V(\mathbf{x}^{(k+1)}) \geq V(\mathbf{x}^{(k)})$, which implies that $\mathbf{x}^{(k)} \not\rightarrow \mathbf{x}^*$.

Example (8.3):

- Consider

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \begin{bmatrix} 8 & 2\sqrt{2} \\ 2\sqrt{2} & 10 \end{bmatrix} \mathbf{x} + \mathbf{x}^T \begin{bmatrix} 3 \\ 6 \end{bmatrix} + 24.$$

- Eigenvalues of \mathbf{Q} : 6 and 12.
- Fixed step size gradient algorithm converges if and only if $0 < \alpha < 2/12 = 1/6$.

Other insights into convergence

- Theorem (8.4): In the steepest descent algorithm,

$$V(\mathbf{x}^{(k+1)}) \leq \left(1 - \frac{1}{r} \right) V(\mathbf{x}^{(k)}),$$

where

$$r = \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})}$$

(called the *condition number* of \mathbf{Q}).

- If r is small (close to 1), convergence is fast.
- If r is large, convergence is slow.
- See book for proof.

Order of convergence

- Given: a sequence $\{\mathbf{x}^{(k)}\}$ converging to \mathbf{x}^* .
- The *order of convergence* is p (where $1 \leq p < \infty$) if

$$0 < \lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} < \infty.$$

We say that the order of convergence is ∞ if for all $p \geq 1$,

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} = 0.$$

- Order of convergence is one measure of “speed” of convergence.

Example:

- Given: $x^{(k)} = 1/k$.
- Thus, $x^{(k)} \rightarrow 0$. Then,

$$\frac{|x^{(k+1)}|}{|x^{(k)}|^p} = \frac{1/(k+1)}{1/k^p} = \frac{k^p}{k+1}.$$

- If $p > 1$, it grows to ∞ .
- If $p = 1$, the sequence converges to 1.
- Hence, the order of convergence is 1.

Example:

- Given: $x^{(k)} = \gamma^k$, where $0 < \gamma < 1$.
- Thus, $x^{(k)} \rightarrow 0$. Then,

$$\frac{|x^{(k+1)}|}{|x^{(k)}|^p} = \frac{\gamma^{k+1}}{(\gamma^k)^p} = \gamma^{k+1-kp} = \gamma^{k(1-p)+1}.$$

- If $p > 1$, it grows to ∞ .

- If $p = 1$, the sequence converges to γ (actually, it remains constant).
- Hence, the order of convergence is 1.

Example:

- Given: $x^{(k)} = \gamma^{(q^k)}$, where $q > 1$ and $0 < \gamma < 1$.
- Thus, $x^{(k)} \rightarrow 0$. Then,

$$\frac{|x^{(k+1)}|}{|x^{(k)}|^p} = \frac{\gamma^{(q^{k+1})}}{(\gamma^{(q^k)})^p} = \gamma^{(q^{k+1} - pq^k)} = \gamma^{(q-p)q^k}.$$

- If $p < q$, the above sequence converges to 0.
- If $p > q$, it grows to ∞ .
- If $p = q$, the sequence converges to 1 (actually, it remains constant).
- Hence, the order of convergence is q .

Bounding the order of convergence

- $g(h) = O(h)$ means that there exists a constant c such that $|g(h)| \leq c|h|$ for sufficiently small h .
- Given $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$. If

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = O(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p),$$

then the order of convergence is at least p .

- Example: Order of convergence is at least 2 if

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = O(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2).$$

- $g(h) = \Omega(h)$ means there exists a constant $c > 0$ such that $|g(h)| \geq c|h|$ for sufficiently small h .
- Given $\mathbf{x}^{(k)} \rightarrow \mathbf{x}^*$. If

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = \Omega(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p),$$

then the order of convergence is at most p .

Order of convergence of steepest descent

Theorem: The steepest descent algorithm has order of convergence of 1 in the worst case.

Proof:

- Consider only quadratic case. Assume $\lambda_{\max}(\mathbf{Q}) > \lambda_{\min}(\mathbf{Q})$.
- Suffices to show that there exists $\mathbf{x}^{(0)}$ such that $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = \Omega(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|)$; i.e.,

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \geq c \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$$

for some $c > 0$.

- By Rayleigh's inequality,

$$\begin{aligned} V(\mathbf{x}^{(k+1)}) &= \frac{1}{2}(\mathbf{x}^{(k+1)} - \mathbf{x}^*)^T \mathbf{Q}(\mathbf{x}^{(k+1)} - \mathbf{x}^*) \\ &\leq \frac{\lambda_{\max}(\mathbf{Q})}{2} \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|^2. \end{aligned}$$

- Similarly,

$$V(\mathbf{x}^{(k)}) \geq \frac{\lambda_{\min}(\mathbf{Q})}{2} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2.$$

- Therefore, by Lemma 8.1, it suffices to show that $\gamma_k \leq d$ for some $d < 1$.
- Recall that for the steepest descent algorithm, γ_k depends on $\mathbf{g}^{(k)}$:

$$\gamma_k = \frac{(\mathbf{g}^{(k)T} \mathbf{g}^{(k)})^2}{(\mathbf{g}^{(k)T} \mathbf{Q} \mathbf{g}^{(k)})(\mathbf{g}^{(k)T} \mathbf{Q}^{-1} \mathbf{g}^{(k)})}.$$

- First consider the case where $n = 2$.
- Suppose $\mathbf{x}^{(0)}$ is chosen such that $\mathbf{g}^{(0)}$ is not an eigenvector of \mathbf{Q} . By Prop. 8.1, $\mathbf{g}^{(k)}$ is also not an eigenvector of \mathbf{Q} for all k (because any two eigenvectors corresponding to $\lambda_{\max}(\mathbf{Q})$ and $\lambda_{\min}(\mathbf{Q})$ are mutually orthogonal).
- Also, $\mathbf{g}^{(k)}$ lies in one of 2 mutually orthogonal directions. Therefore, the value of γ_k is one of 2 numbers, both of which are < 1 . This proves the $n = 2$ case.
- For the general n case, let \mathbf{v}_1 and \mathbf{v}_2 be mutually orthogonal eigenvectors corresponding to $\lambda_{\max}(\mathbf{Q})$ and $\lambda_{\min}(\mathbf{Q})$. Choose $\mathbf{x}^{(0)}$ such that $\mathbf{g}^{(0)}$ lies in the span of \mathbf{v}_1 and \mathbf{v}_2 but is not equal to either. Then, $\mathbf{g}^{(k)}$ lies in the span of \mathbf{v}_1 and \mathbf{v}_2 for all k . To see this, note that $\mathbf{g}^{(k+1)} = (\mathbf{I} - \alpha_k \mathbf{Q})\mathbf{g}^{(k)}$, and the span of \mathbf{v}_1 and \mathbf{v}_2 is invariant under operations of $\mathbf{I} - \alpha_k \mathbf{Q}$. We can now proceed as in the $n = 2$ case.