# Practical Course: Modeling, Simulation, Optimization

Week 9

**Daniël Veldman**

Chair in Dynamics, Control, and Numerics, Friedrich-Alexander-University Erlangen-Nürnberg

## Contents

# 9.A  Existence and uniqueness of minimizers

# Existence of the infimum

We consider the minimization of a functional $J : U \to \mathbb{R}$ over a normed space $U$.
Note: $U$ can be infinite dimensional.

We assume that $J(u) \geq 0$ for all $u \in U$.

We are also given a subset $U_{\mathrm{ad}} \subseteq U$ of admissible values for $u$.

# Existence of the infimum

We consider the minimization of a functional $J : U \to \mathbb{R}$ over a normed space $U$.
Note: $U$ can be infinite dimensional.

We assume that $J(u) \geq 0$ for all $u \in U$.

We are also given a subset $U_{\mathrm{ad}} \subseteq U$ of admissible values for $u$.

Then $\{J(u) \mid u \in U_{\mathrm{ad}}\}$ is a subset of $\mathbb{R}$ that is bounded from below (by $0$). Therefore,

$$\inf_{u \in U_{\mathrm{ad}}} J(u) = \inf\{J(u) \mid u \in U_{\mathrm{ad}}\},$$

exists.
By definition of the infimum, there thus exists a sequence $u_1, u_2, u_3, \ldots$ in $U_{\mathrm{ad}}$ such that

$$J(u_k) \to \inf_{u \in U_{\mathrm{ad}}} J(u).$$

This sequence is called a *minimizing sequence*.

## Existence of the minimizer (finite dimensional case)

Question: does

$$\min_{u \in U_{\mathrm{ad}}} J(u)$$

exist? In other words, is there a minimizer $u^* \in U_{\mathrm{ad}}$ such that

$$J(u^*) = \inf_{u \in U_{\mathrm{ad}}} J(u)?$$

# Existence of the minimizer (finite dimensional case)

Question: does

$$\min_{u \in U_{\mathrm{ad}}} J(u)$$

exist? In other words, is there a minimizer $u^* \in U_{\mathrm{ad}}$ such that

$$J(u^*) = \inf_{u \in U_{\mathrm{ad}}} J(u)?$$

First consider the case where $U$ is finite dimensional.

Observe, if $U_{\mathrm{ad}}$ is closed and the minimizing sequence $u_1, u_2, u_3, \dots$ is bounded, then it also has a limit in $U_{\mathrm{ad}}$. This limit is a minimizer $u^*$.

Two important cases:
► $U_{\mathrm{ad}}$ is bounded and closed.
  It is immediate that the minimizing sequence is bounded.

## Existence of the minimizer (finite dimensional case)

Question: does

$$\min_{u \in U_{\mathrm{ad}}} J(u)$$

exist? In other words, is there a minimizer $u^* \in U_{\mathrm{ad}}$ such that

$$J(u^*) = \inf_{u \in U_{\mathrm{ad}}} J(u)?$$

First consider the case where $U$ is finite dimensional.

Observe, if $U_{\mathrm{ad}}$ is closed and the minimizing sequence $u_1, u_2, u_3, \ldots$ is bounded, then it also has a limit in $U_{\mathrm{ad}}$. This limit is a minimizer $u^*$.

Two important cases:
- ▶ $U_{\mathrm{ad}}$ is bounded and closed.
  It is immediate that the minimizing sequence is bounded.
- ▶ $J(u)$ is coercive, i.e. $J(u_k) \to \infty$ if $|u_k| \to \infty$. Note: it is sufficient that $J(u) \geq |u|^2$.
  Then we can reason as follows.
  Suppose that the minimizing sequence $u_1, u_2, u_3, \ldots$ is unbounded.
  Then there exists a subsequence $u_{k_1}, u_{k_2}, u_{k_3}, \ldots$ such that $|u_{k_j}| \to \infty$.
  But $J(u_{k_j}) > |u_{k_j}|^2$, so also $J(u_{k_j}) \to \infty$.
  But then $J(u_{k_j})$ is not a minimizing sequence. Contradiction.
  Conclusion: the minimizing sequence must be bounded.

# Existence of the minimizer (infinite dimensional case)

Question: does

$$\min_{u \in U_{\mathrm{ad}}} J(u)$$

exist? In other words, is there a minimizer $u^* \in U_{\mathrm{ad}}$ such that

$$J(u^*) = \inf_{u \in U_{\mathrm{ad}}} J(u)?$$

The infinite dimensional case is much more subtle.

Problem: We can no longer be sure that a bounded sequence has a (strong) limit. In other words, we do no longer have compactness.
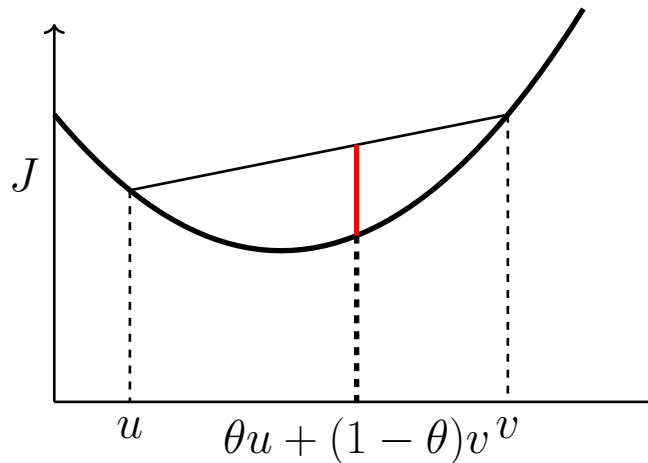
Typical example: consider $U_{\mathrm{ad}} = L^2(0, \pi)$ and consider the sequence $u_k = \sin(kx)$. This sequence converges weakly to zero, but does not have a strong limit.

We will come back to this problem in a few slides.

## Uniqueness of the minimizer (convex analysis)

The functional $J(u)$ is called $\alpha$-convex iff

$$J(\theta u + (1 - \theta)v) \leq \theta J(u) + (1 - \theta)J(v) - \frac{\alpha\theta(1 - \theta)}{2}|u - v|^2, \qquad \theta \in [0, 1].$$



The admissible set $U_{\mathrm{ad}}$ is convex when $u, v \in U_{\mathrm{ad}}$

$$\theta u + (1 - \theta)v \in U_{\mathrm{ad}}, \qquad \theta \in [0, 1].$$

## Uniqueness of the minimizer (convex analysis)

The functional $J(u)$ is called $\alpha$-convex iff

$$J(\theta u + (1 - \theta)v) \le \theta J(u) + (1 - \theta)J(v) - \frac{\alpha\theta(1 - \theta)}{2}|u - v|^2, \qquad \theta \in [0, 1].$$
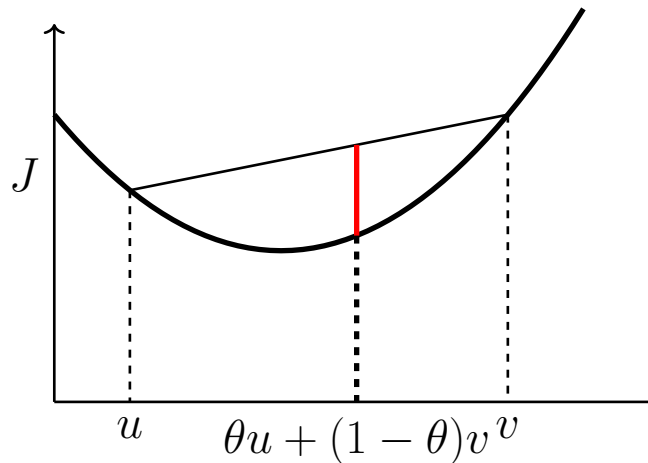


The admissible set $U_{\mathrm{ad}}$ is convex when $u, v \in U_{\mathrm{ad}}$

$$\theta u + (1 - \theta)v \in U_{\mathrm{ad}}, \qquad \theta \in [0, 1].$$

Uniqueness of the minimizer:
Suppose that there are two points $u, v \in U_{\mathrm{ad}}$ such that $J(u) = J(v) = \min_{u \in U_{\mathrm{ad}}} J(u)$.

$$J(\theta u + (1 - \theta)v) \le \min_{u \in U_{\mathrm{ad}}} J(u) - \frac{\alpha\theta(1 - \theta)}{2}|u - v|^2 < \min_{u \in U_{\mathrm{ad}}} J(u),$$

and $\theta u + (1 - \theta)v \in U_{\mathrm{ad}}$. Contradiction.

# Existence of the minimizer (infinite dimensional case, revisited)

Question: does

$$\min_{u \in U_{\mathrm{ad}}} J(u)$$

exist? In other words, is there a minimizer $u^* \in U_{\mathrm{ad}}$ such that

$$J(u^*) = \inf_{u \in U_{\mathrm{ad}}} J(u)?$$

Consider a minimizing sequence $u_1, u_2, u_3, \ldots$.
The minimizing sequence is bounded when $U_{\mathrm{ad}}$ is bounded or when $J$ is coercive.
The bounded minimizing sequence $u_1, u_2, u_3, \ldots$ has a weak limit $v$.

# Existence of the minimizer (infinite dimensional case, revisited)

Question: does

$$\min_{u \in U_{\mathrm{ad}}} J(u)$$

exist? In other words, is there a minimizer $u^* \in U_{\mathrm{ad}}$ such that

$$J(u^*) = \inf_{u \in U_{\mathrm{ad}}} J(u)?$$

Consider a minimizing sequence $u_1, u_2, u_3, \ldots$.
The minimizing sequence is bounded when $U_{\mathrm{ad}}$ is bounded or when $J$ is coercive.
The bounded minimizing sequence $u_1, u_2, u_3, \ldots$ has a weak limit $v$.

Now three problems remain:

▶ Is the weak limit $v \in U_{\mathrm{ad}}$?
If $U_{\mathrm{ad}}$ is strongly closed and convex, it is also weakly closed (Hahn-Banach).

▶ Do we have that $J(v) = \lim_{k \to \infty} J(u_k) = \inf_{u \in U_{\mathrm{ad}}} J(u)$?
This is achieved by assuming that $J$ is weakly lower semi-continuous (by definition).

▶ Does the minimizing sequence $u_1, u_2, u_3, \ldots$ also converge strongly to $v$?
This follows from the previous point and the strong convexity of $J$ (with $\theta = \frac{1}{2}$):

$$J(v) \leq J(\tfrac{u_k + v}{2}) \leq \frac{J(u_k) + J(v)}{2} - \frac{\alpha}{8}|u_k - v|^2, \quad \Rightarrow \quad \frac{\alpha}{8}|u_k - v|^2 \leq \frac{J(u_k) - J(v)}{2} \to 0.$$

# 9.B  A basic gradient descent algorithm

# Gradient descent

Question: How to we compute the minimizer $u^*$ of a (convex) functional $J(u)$.

Basic idea: Start from an initial guess $u_0$.
Compute iterates by updating $u_k$ in the direction of the steepest descent (i.e. $-\nabla J$),

$$u_{k+1} = u_k - \beta_k \nabla J(u_k), \qquad \beta_k > 0,$$

where $\beta$ denotes the step size.

# Gradient descent

Question: How to we compute the minimizer $u^*$ of a (convex) functional $J(u)$.

Basic idea: Start from an initial guess $u_0$.
Compute iterates by updating $u_k$ in the direction of the steepest descent (i.e. $-\nabla J$),

$$u_{k+1} = u_k - \beta_k \nabla J(u_k), \qquad \beta_k > 0,$$

where $\beta$ denotes the step size.

Three problems:
► How to compute $\nabla J$?
► How to choose the stepsize $\beta_k$?
► When do we stop the iterations?

## Computation of the gradient/ sensitivity analysis

By definition of the gradient, we have that

$$\langle \nabla J, \tilde{u} \rangle := \lim_{h \to 0} \frac{J(u + h\tilde{u}) - J(u)}{h} = \frac{\partial J}{\partial u}(u)\tilde{u},$$

for all perturbations $\tilde{u}$.

Note:
- ▶ $\nabla J(u)$ and $\frac{\partial J}{\partial u}$ are not the same:
  $\nabla J(u)$ is a column vector and $\frac{\partial J}{\partial u}$ is a row vector.
- ▶ We can use any innerproduct $\langle \cdot, \cdot \rangle$ at the LHS.
  This will not affect $\frac{\partial J}{\partial u}$ but it will change $\nabla J$!

# Computation of the gradient/ sensitivity analysis

By definition of the gradient, we have that

$$\langle \nabla J, \tilde{u} \rangle := \lim_{h \to 0} \frac{J(u + h\tilde{u}) - J(u)}{h} = \frac{\partial J}{\partial u}(u)\tilde{u},$$

for all perturbations $\tilde{u}$.

Note:

▶ $\nabla J(u)$ and $\frac{\partial J}{\partial u}$ are not the same:
  $\nabla J(u)$ is a column vector and $\frac{\partial J}{\partial u}$ is a row vector.
▶ We can use any innerproduct $\langle \cdot, \cdot \rangle$ at the LHS.
  This will not affect $\frac{\partial J}{\partial u}$ but it will change $\nabla J$!

Two examples:

▶ When $\langle x, y \rangle = x^\top y$, i.e. when we use the standard Euclidean inner product

$$\nabla J = \left(\frac{\partial J}{\partial u}\right)^\top.$$

▶ When we use a weighted inner product $\langle x, y \rangle = x^\top \mathbf{W} y$, for a symmetric and positive definite matrix $\mathbf{W}$, we get that

$$\nabla J = \mathbf{W}^{-1} \left(\frac{\partial J}{\partial u}\right)^\top.$$

# Intermezzo: Why the choice of inner product matters/helps

Suppose that $J(u) = \langle u + b, u \rangle = (u + b)^\top \mathbf{W} u$.
(Any quadratic functional with Hessian $\mathbf{W}$ can be written in this form)

$$\langle \nabla J, \tilde{u} \rangle := \lim_{h \to 0} \frac{J(u + h\tilde{u}) - J(u)}{h} = \lim_{h \to 0} \frac{\langle u + h\tilde{u} + b, u + h\tilde{u} \rangle - \langle u + b, u \rangle}{h},$$

$$= \lim_{h \to 0} \frac{\langle u + b, u \rangle + h\langle u + b, \tilde{u} \rangle + h\langle \tilde{u}, u \rangle + h^2 \langle \tilde{u}, \tilde{u} \rangle - \langle u + b, u \rangle}{h}$$

$$= \langle 2u + b, \tilde{u} \rangle.$$

# Intermezzo: Why the choice of inner product matters/helps

Suppose that $J(u) = \langle u + b, u \rangle = (u + b)^\top \mathbf{W} u$.
(Any quadratic functional with Hessian $\mathbf{W}$ can be written in this form)

$$
\begin{aligned}
\langle \nabla J, \tilde{u} \rangle &:= \lim_{h \to 0} \frac{J(u + h\tilde{u}) - J(u)}{h} = \lim_{h \to 0} \frac{\langle u + h\tilde{u} + b, u + h\tilde{u} \rangle - \langle u + b, u \rangle}{h}, \\
&= \lim_{h \to 0} \frac{\langle u + b, u \rangle + h\langle u + b, \tilde{u} \rangle + h\langle \tilde{u}, u \rangle + h^2 \langle \tilde{u}, \tilde{u} \rangle - \langle u + b, u \rangle}{h} \\
&= \langle 2u + b, \tilde{u} \rangle.
\end{aligned}
$$

We thus see that

$$
\nabla J(u) = 2u + b, \qquad\qquad u^* = -\tfrac{1}{2}b.
$$

Suppose we have an initial guess $u_0$ and take the stepsize $\beta_0 = \frac{1}{2}$. Then

$$
u_1 = u_0 - \tfrac{1}{2}\nabla J(u_0) = u_0 - \tfrac{1}{2}(2u_0 + b) = -\tfrac{1}{2}b = u^*.
$$

**Conclusion:** when we have a quadratic cost functional with Hessian $\mathbf{W}$ and compute the gradient w.r.t. the inner product $\langle \mathbf{u}, \mathbf{v} \rangle = \mathbf{u}^\top \mathbf{W} \mathbf{v}$, the gradient descent algorithm converges in 1 iteration (with $\beta = \frac{1}{2}$).

However, this idea is not directly applicable: often, the Hessian cannot be computed easily and the considered cost functionals are not quadratic.
Even in these situation, choosing $\mathbf{W}$ well can improve the convergence.

# The choice of the step size

We have that

$$J(u_{k+1}) = J(u_k - \beta_k \nabla J(u_k)) = J(u_k) - \beta_k \frac{\partial J}{\partial u_k} \nabla J(u_k) + O(\beta_k^2)$$
$$= J(u_k) - \beta_k \langle \nabla J(u_k), \nabla J(u_k) \rangle + O(\beta_k^2).$$

As long as we are not at a critical point ($\nabla J(u_k) = 0$) $\langle \nabla J(u_k), \nabla J(u_k) \rangle > 0$, so

$$J(u_{k+1}) < J(u_k)$$

for $\beta_k > 0$ small enough.

# The choice of the step size

We have that

$$J(u_{k+1}) = J(u_k - \beta_k \nabla J(u_k)) = J(u_k) - \beta_k \frac{\partial J}{\partial u_k} \nabla J(u_k) + O(\beta_k^2)$$
$$= J(u_k) - \beta_k \langle \nabla J(u_k), \nabla J(u_k) \rangle + O(\beta_k^2).$$

As long as we are not at a critical point ($\nabla J(u_k) = 0$) $\langle \nabla J(u_k), \nabla J(u_k) \rangle > 0$, so

$$J(u_{k+1}) < J(u_k)$$

for $\beta_k > 0$ small enough.

We can thus take the following simple but effective approach (used at every iteration).
- ▶ Choose a step size $\beta_k > 0$.
- ▶ Compute $J(u_k - \beta \nabla J(u_k))$.
- ▶ If $J(u_k - \beta \nabla J(u_k)) < J(u_k)$, we accept this step size.
  If not, we reduce the step size (e.g. by a factor 2) and recompute $J(u_k - \beta \nabla J(u_k))$.

This should always lead to a $\beta_k > 0$ such that $J(u_k - \beta \nabla J(u_k)) < J(u_k)$.
(Provided that $\nabla J(u_k)$ is computed sufficiently accurate)

# Termination/convergence conditions

Typical convergence conditions:

▶ Relative decrease in the cost functional is sufficiently small:

$$J(u_k) - J(u_{k+1}) < \texttt{tol} J(u_k).$$

▶ Relative change in iterates is sufficiently small:

$$|u_{k-1} - u_k| < \texttt{tol}|u_k|.$$

▶ The gradient is sufficiently small:

$$|\nabla J(u_k)| < \texttt{tol}.$$

In the first two conditions, we typically use $\texttt{tol} \in [10^{-6}, 10^{-3}]$.

Often not all three conditions are checked simultaneously, but only one or two are used.

Note: $\texttt{tol}$ in the last condition is an absolute tolerance, while $\texttt{tol}$ in the first two conditions is a relative tolerance.
A reasonable magnitude for the absolute tolerance might be difficult to estimate.

## Pseudo code of the resulting gradient descent algorithm

- ▶ Choose an initial guess $u_0$
- ▶ Choose an initial step size $\beta$
- ▶ Compute $J_0 = J(u_0)$.
- ▶ for $i = 1$: `max_iters`
- ▶      Compute $g_0 = \nabla J(u_0)$.
- ▶      Set $J_1 = \infty$ and $\beta = 4\beta$.
- ▶      while $J_1 > J_0$
- ▶          Set $\beta = \beta/2$.
- ▶          Set $u_1 = u_0 - \beta g_0$.
- ▶          Compute $J_1 = J(u_1)$.
- ▶      if convergence conditions are satisfied
- ▶          Return $u_1$, $J_1$.
- ▶      Set $u_0 = u_1$
- ▶      Set $J_0 = J_1$