

# Numerical Optimal Control

## Lecture 3: Newton method & SQP

Sébastien Gros

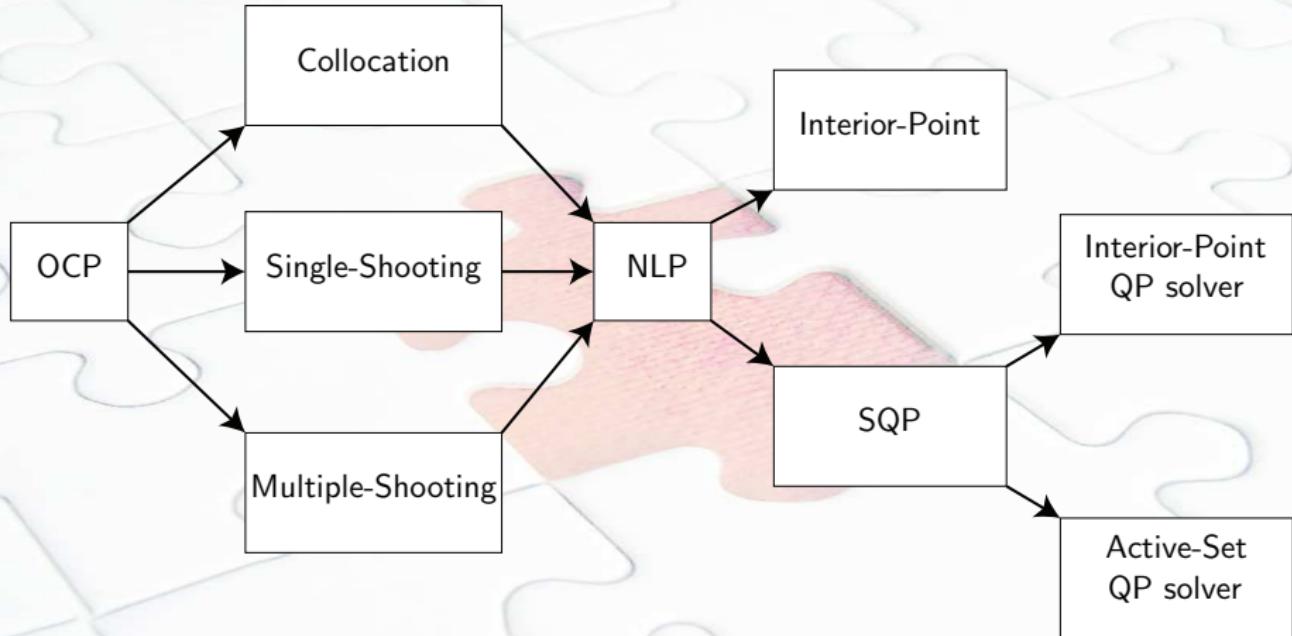
ITK NTNU

NTNU PhD course

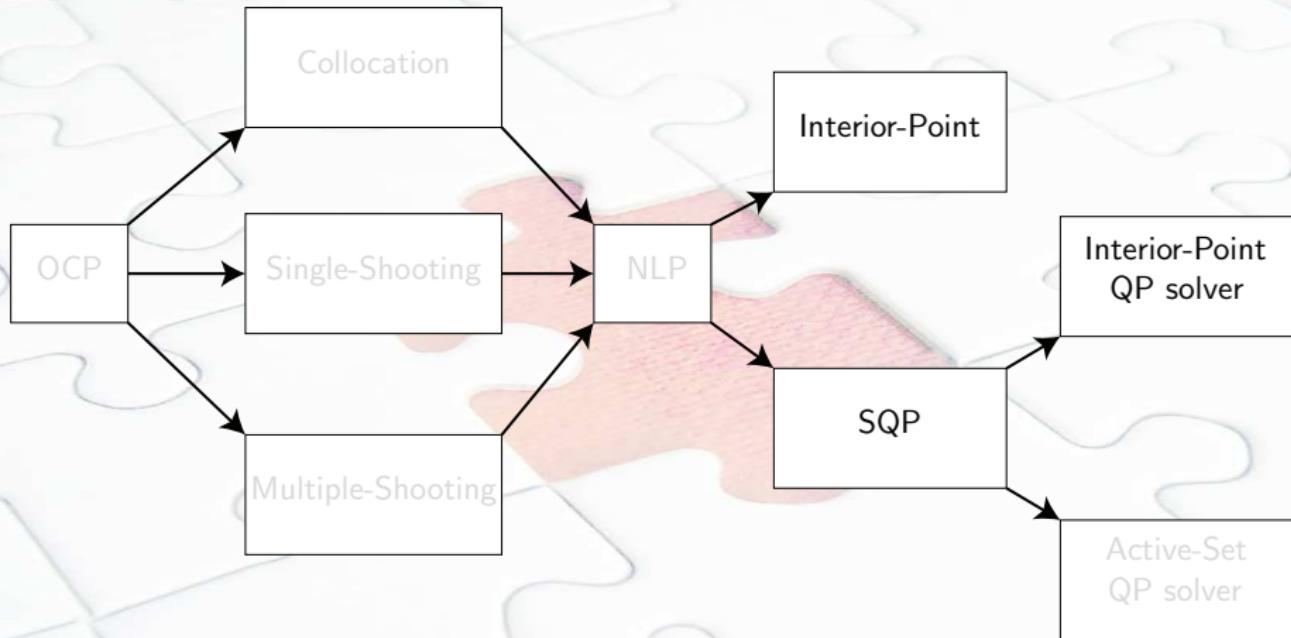
# Outline

- 1 The Newton method
- 2 Newton on the KKT conditions
- 3 The reduced Newton step (unconstrained problems)
- 4 The merit function - Line-search for constrained problems
- 5 Newton-type methods
- 6 Sequential Quadratic Programming

# Survival map of Direct Optimal Control

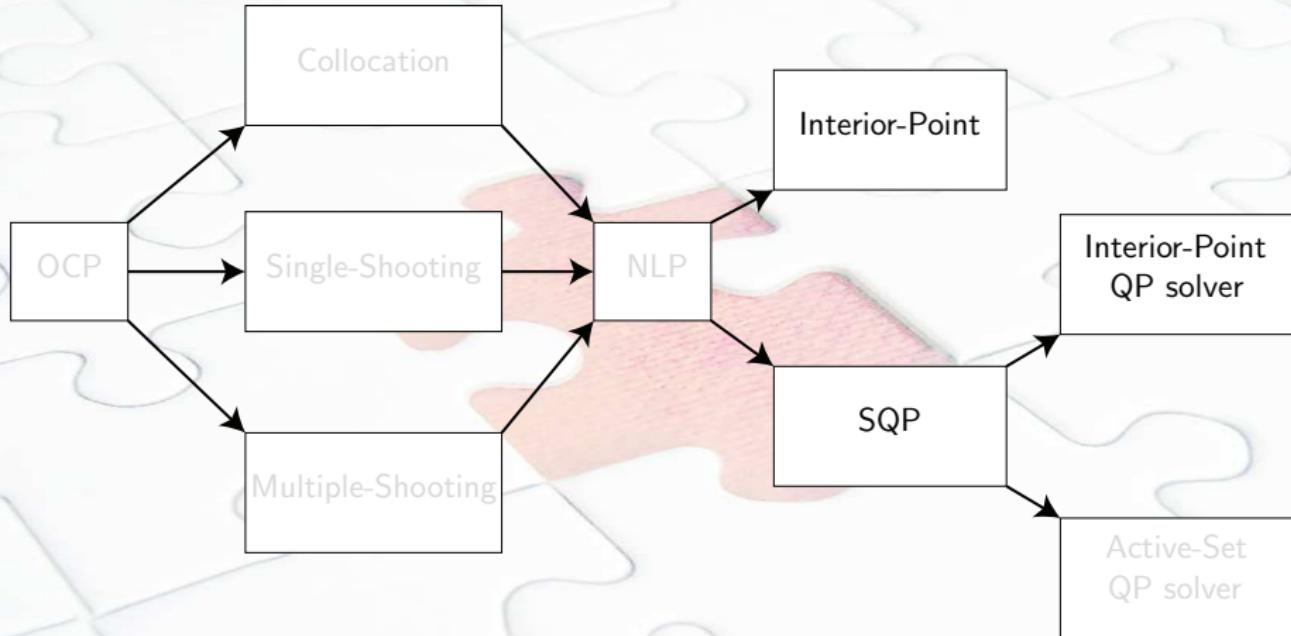


# Survival map of Direct Optimal Control



**Newton - a general-purpose sledgehammer for algebraic equations...**

# Survival map of Direct Optimal Control



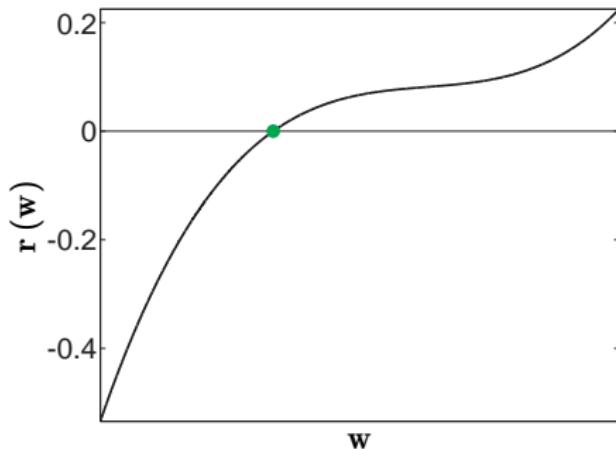
**Newton - a general-purpose sledgehammer for algebraic equations...  
... will be used to solve the KKT conditions !!**

# Outline

- 1 The Newton method
- 2 Newton on the KKT conditions
- 3 The reduced Newton step (unconstrained problems)
- 4 The merit function - Line-search for constrained problems
- 5 Newton-type methods
- 6 Sequential Quadratic Programming

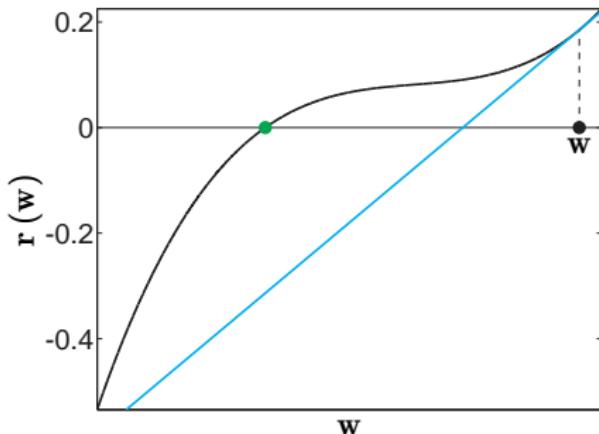
## Core idea

**Goal:** solve  $r(\mathbf{w}) = 0 \dots \text{how } ?!?$



## Core idea

**Goal:** solve  $r(w) = 0 \dots$  how ?!?

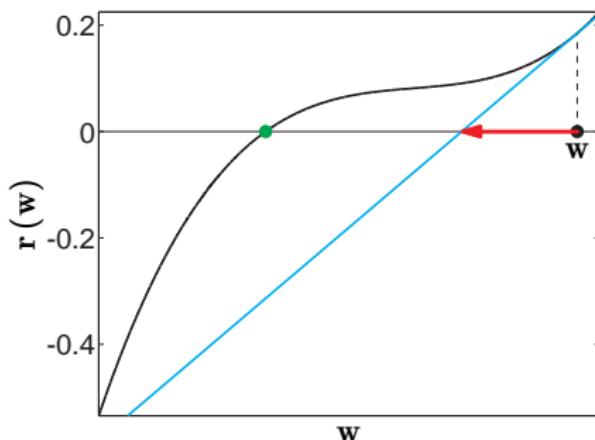


**Key idea:** guess  $w$ , iterate the [linear model](#):

$$r(w + \Delta w) \approx r(w) + \nabla r(w)^\top \Delta w = 0$$

## Core idea

Goal: solve  $r(\mathbf{w}) = 0 \dots$  how ?!?

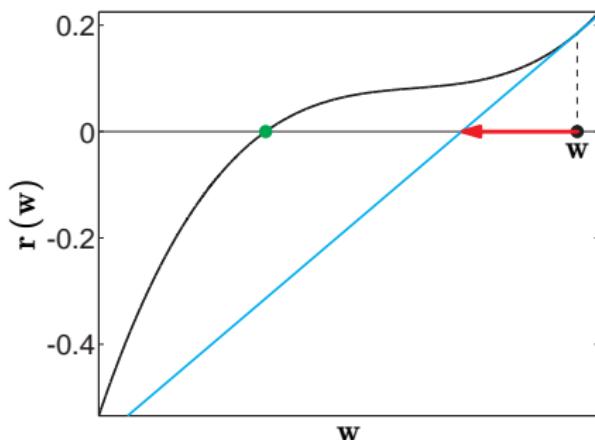


Key idea: guess  $\mathbf{w}$ , iterate the linear model:

$$r(\mathbf{w} + \Delta\mathbf{w}) \approx r(\mathbf{w}) + \nabla r(\mathbf{w})^\top \Delta\mathbf{w} = 0$$

## Core idea

Goal: solve  $r(w) = 0$ ... how ?!?



Key idea: guess  $w$ , iterate the linear model:

$$r(w + \Delta w) \approx r(w) + \nabla r(w)^T \Delta w = 0$$

---

### Algorithm: Newton method

---

**Input:**  $w$ , Tol

**while**  $\|r(w)\| \geq \text{tol}$  **do**

    Compute

$r(w)$  and  $\nabla r(w)$

    Compute the **Newton direction**

$$\nabla r(w)^T \Delta w = -r(w)$$

    Newton step

$$w \leftarrow w + \Delta w$$

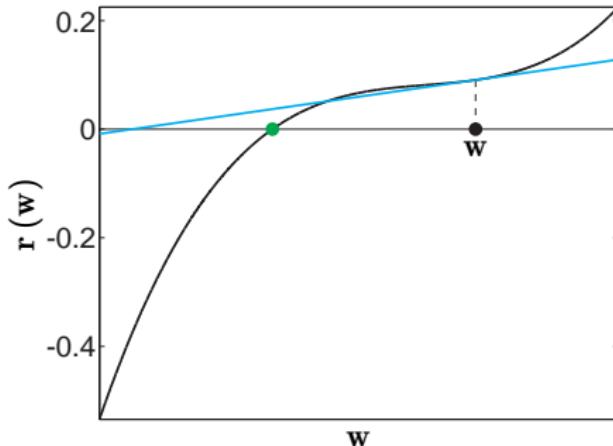
---

**return**  $w$

---

## Core idea

Goal: solve  $\mathbf{r}(\mathbf{w}) = 0 \dots$  how ?!



Key idea: guess  $\mathbf{w}$ , iterate the [linear model](#):

$$\mathbf{r}(\mathbf{w} + \Delta\mathbf{w}) \approx \mathbf{r}(\mathbf{w}) + \nabla\mathbf{r}(\mathbf{w})^\top \Delta\mathbf{w} = 0$$

---

### Algorithm: Newton method

---

**Input:**  $\mathbf{w}$ , Tol

**while**  $\|\mathbf{r}(\mathbf{w})\| \geq \text{tol}$  **do**

    Compute

$$\mathbf{r}(\mathbf{w}) \quad \text{and} \quad \nabla\mathbf{r}(\mathbf{w})$$

    Compute the **Newton direction**

$$\nabla\mathbf{r}(\mathbf{w})^\top \Delta\mathbf{w} = -\mathbf{r}(\mathbf{w})$$

    Newton step

$$\mathbf{w} \leftarrow \mathbf{w} + \Delta\mathbf{w}$$

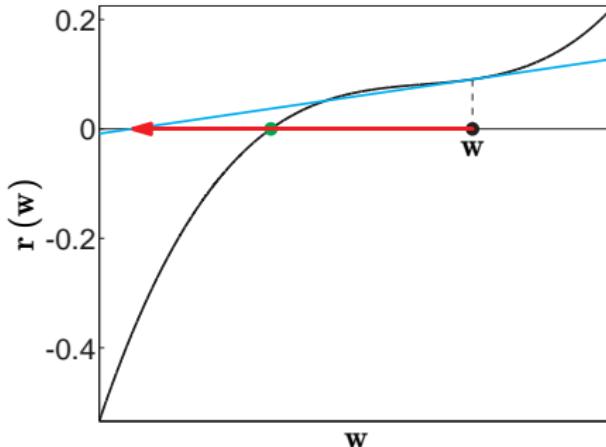
---

**return**  $\mathbf{w}$

---

## Core idea

Goal: solve  $r(w) = 0$ ... how ?!



Key idea: guess  $w$ , iterate the linear model:

$$r(w + \Delta w) \approx r(w) + \nabla r(w)^T \Delta w = 0$$

---

### Algorithm: Newton method

---

**Input:**  $w$ , Tol

**while**  $\|r(w)\| \geq \text{tol}$  **do**

    Compute

$$r(w) \quad \text{and} \quad \nabla r(w)$$

    Compute the **Newton direction**

$$\nabla r(w)^T \Delta w = -r(w)$$

    Newton step

$$w \leftarrow w + \Delta w$$

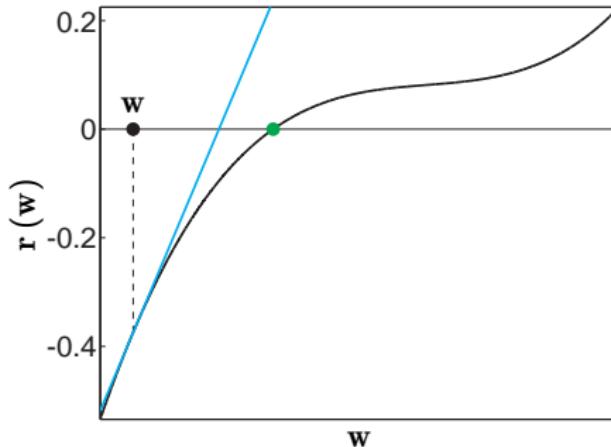
---

**return**  $w$

---

## Core idea

Goal: solve  $\mathbf{r}(\mathbf{w}) = 0 \dots$  how ?!



Key idea: guess  $\mathbf{w}$ , iterate the linear model:

$$\mathbf{r}(\mathbf{w} + \Delta\mathbf{w}) \approx \mathbf{r}(\mathbf{w}) + \nabla\mathbf{r}(\mathbf{w})^\top \Delta\mathbf{w} = 0$$

---

### Algorithm: Newton method

---

**Input:**  $\mathbf{w}$ , Tol

**while**  $\|\mathbf{r}(\mathbf{w})\| \geq \text{tol}$  **do**

    Compute

$$\mathbf{r}(\mathbf{w}) \quad \text{and} \quad \nabla\mathbf{r}(\mathbf{w})$$

    Compute the **Newton direction**

$$\nabla\mathbf{r}(\mathbf{w})^\top \Delta\mathbf{w} = -\mathbf{r}(\mathbf{w})$$

    Newton step

$$\mathbf{w} \leftarrow \mathbf{w} + \Delta\mathbf{w}$$

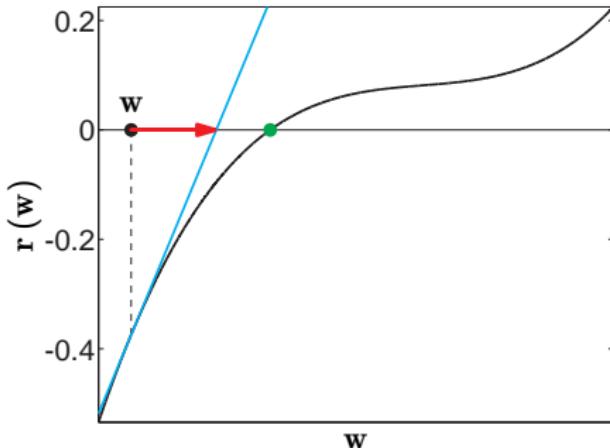
---

**return**  $\mathbf{w}$

---

## Core idea

Goal: solve  $r(w) = 0 \dots$  how ?!



Key idea: guess  $w$ , iterate the linear model:

$$r(w + \Delta w) \approx r(w) + \nabla r(w)^T \Delta w = 0$$

---

### Algorithm: Newton method

---

**Input:**  $w$ , Tol

**while**  $\|r(w)\| \geq \text{tol}$  **do**

    Compute

$$r(w) \quad \text{and} \quad \nabla r(w)$$

    Compute the **Newton direction**

$$\nabla r(w)^T \Delta w = -r(w)$$

    Newton step

$$w \leftarrow w + \Delta w$$

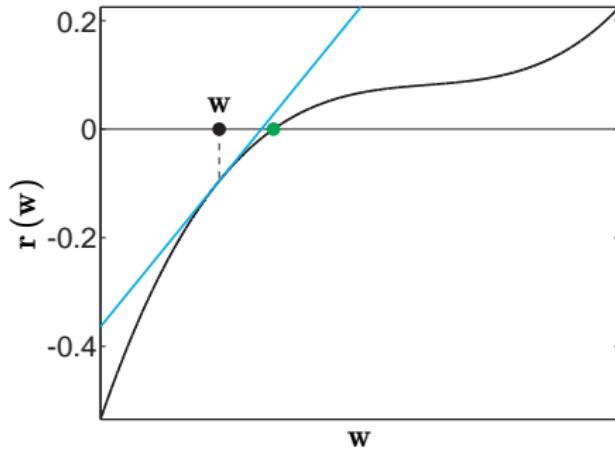
---

**return**  $w$

---

## Core idea

Goal: solve  $r(\mathbf{w}) = 0 \dots$  how ?!



Key idea: guess  $\mathbf{w}$ , iterate the linear model:

$$r(\mathbf{w} + \Delta\mathbf{w}) \approx r(\mathbf{w}) + \nabla r(\mathbf{w})^\top \Delta\mathbf{w} = 0$$

---

### Algorithm: Newton method

---

**Input:**  $\mathbf{w}$ , Tol

**while**  $\|r(\mathbf{w})\| \geq \text{tol}$  **do**

    Compute

$$r(\mathbf{w}) \quad \text{and} \quad \nabla r(\mathbf{w})$$

    Compute the **Newton direction**

$$\nabla r(\mathbf{w})^\top \Delta\mathbf{w} = -r(\mathbf{w})$$

    Newton step

$$\mathbf{w} \leftarrow \mathbf{w} + \Delta\mathbf{w}$$

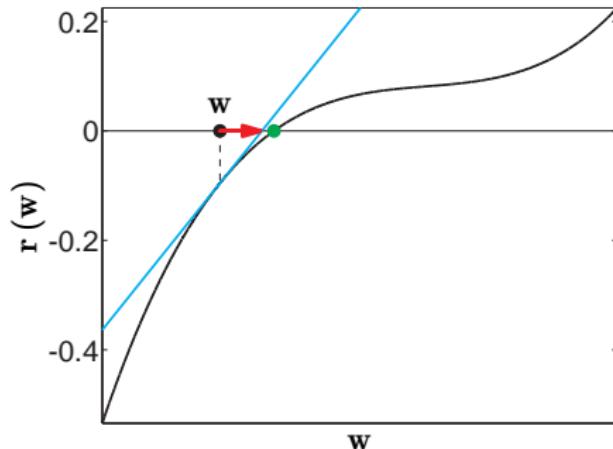
---

**return**  $\mathbf{w}$

---

## Core idea

Goal: solve  $r(w) = 0 \dots$  how ?!



Key idea: guess  $w$ , iterate the linear model:

$$r(w + \Delta w) \approx r(w) + \nabla r(w)^T \Delta w = 0$$

---

### Algorithm: Newton method

---

**Input:**  $w$ , Tol

**while**  $\|r(w)\| \geq \text{tol}$  **do**

    Compute

$$r(w) \quad \text{and} \quad \nabla r(w)$$

    Compute the **Newton direction**

$$\nabla r(w)^T \Delta w = -r(w)$$

    Newton step

$$w \leftarrow w + \Delta w$$

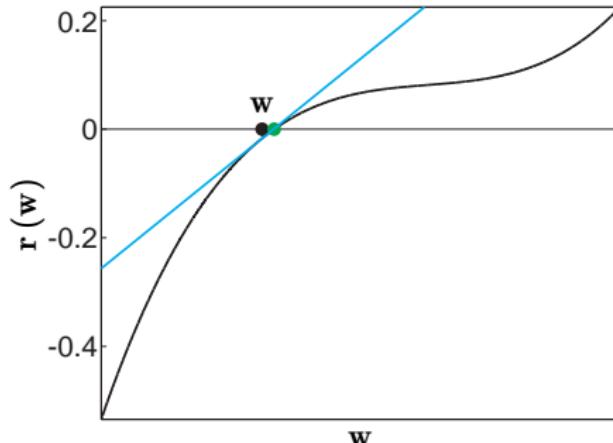
---

**return**  $w$

---

## Core idea

Goal: solve  $\mathbf{r}(\mathbf{w}) = 0 \dots$  how ?!?



Key idea: guess  $\mathbf{w}$ , iterate the [linear model](#):

$$\mathbf{r}(\mathbf{w} + \Delta\mathbf{w}) \approx \mathbf{r}(\mathbf{w}) + \nabla\mathbf{r}(\mathbf{w})^\top \Delta\mathbf{w} = 0$$

---

### Algorithm: Newton method

---

**Input:**  $\mathbf{w}$ , Tol

**while**  $\|\mathbf{r}(\mathbf{w})\| \geq \text{tol}$  **do**

    Compute

$$\mathbf{r}(\mathbf{w}) \quad \text{and} \quad \nabla\mathbf{r}(\mathbf{w})$$

    Compute the **Newton direction**

$$\nabla\mathbf{r}(\mathbf{w})^\top \Delta\mathbf{w} = -\mathbf{r}(\mathbf{w})$$

    Newton step

$$\mathbf{w} \leftarrow \mathbf{w} + \Delta\mathbf{w}$$

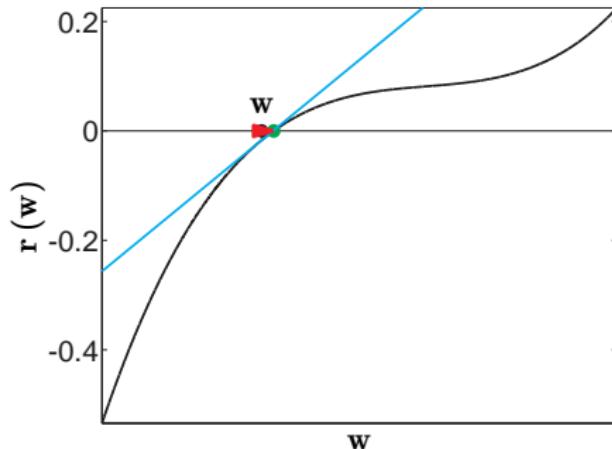
---

**return**  $\mathbf{w}$

---

## Core idea

Goal: solve  $r(w) = 0$ ... how ?!?



Key idea: guess  $w$ , iterate the linear model:

$$r(w + \Delta w) \approx r(w) + \nabla r(w)^T \Delta w = 0$$

---

### Algorithm: Newton method

---

**Input:**  $w$ , Tol

**while**  $\|r(w)\| \geq \text{tol}$  **do**

    Compute

$$r(w) \quad \text{and} \quad \nabla r(w)$$

    Compute the **Newton direction**

$$\nabla r(w)^T \Delta w = -r(w)$$

    Newton step

$$w \leftarrow w + \Delta w$$

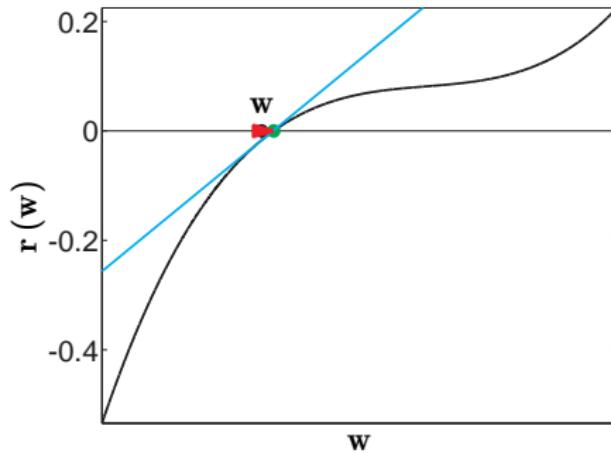
---

**return**  $w$

---

## Core idea

Goal: solve  $r(w) = 0 \dots$  how ?!?



Key idea: guess  $w$ , iterate the linear model:

$$r(w + \Delta w) \approx r(w) + \nabla r(w)^T \Delta w = 0$$

• This is a full-step Newton iteration

---

### Algorithm: Newton method

---

**Input:**  $w$ , Tol

**while**  $\|r(w)\| \geq \text{tol}$  **do**

    Compute

$$r(w) \quad \text{and} \quad \nabla r(w)$$

    Compute the **Newton direction**

$$\nabla r(w)^T \Delta w = -r(w)$$

    Newton step

$$w \leftarrow w + \Delta w$$

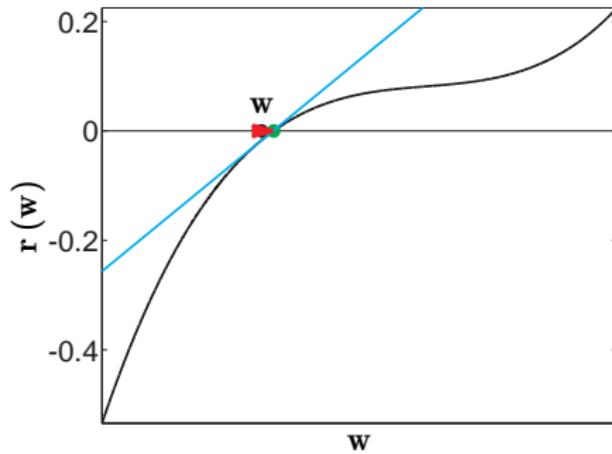
---

**return**  $w$

---

## Core idea

Goal: solve  $r(w) = 0 \dots$  how ?!?



Key idea: guess  $w$ , iterate the linear model:

$$r(w + \Delta w) \approx r(w) + \nabla r(w)^T \Delta w = 0$$

---

### Algorithm: Newton method

---

**Input:**  $w$ , Tol

**while**  $\|r(w)\| \geq \text{tol}$  **do**

    Compute

$$r(w) \quad \text{and} \quad \nabla r(w)$$

    Compute the **Newton direction**

$$\nabla r(w)^T \Delta w = -r(w)$$

    Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w$$

---

**return**  $w$

---

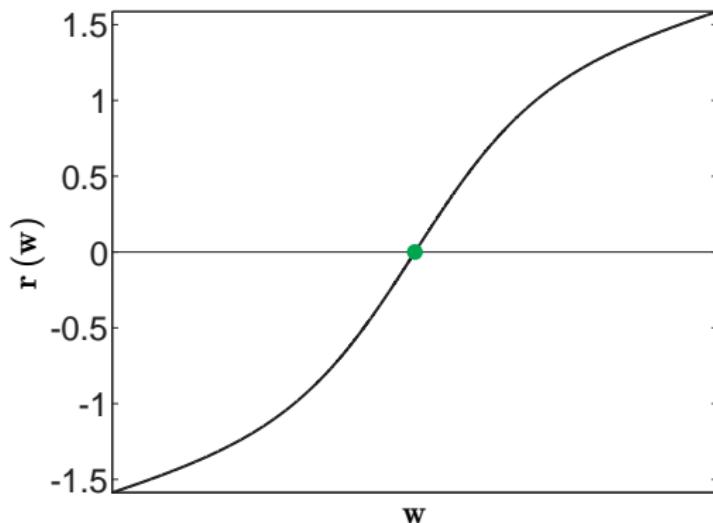
- This is a **full-step** Newton iteration
- Reduced steps are often needed

## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$

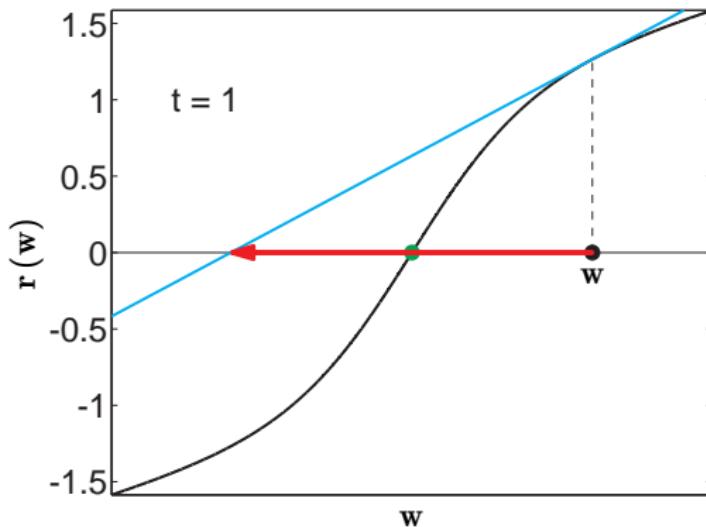


## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$

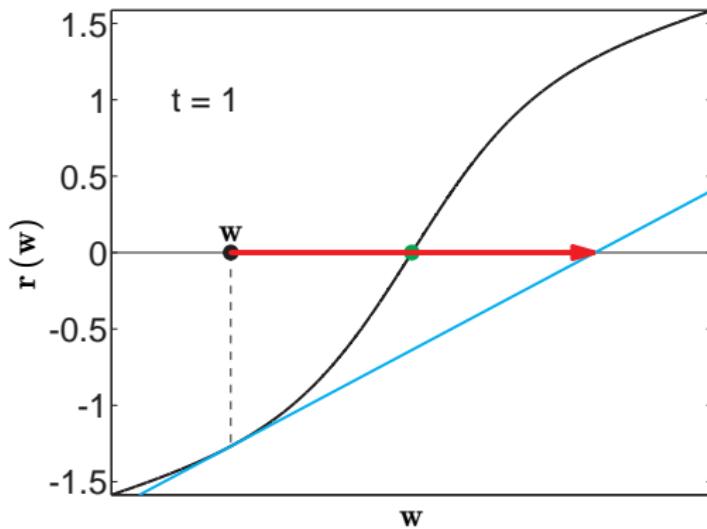


## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$

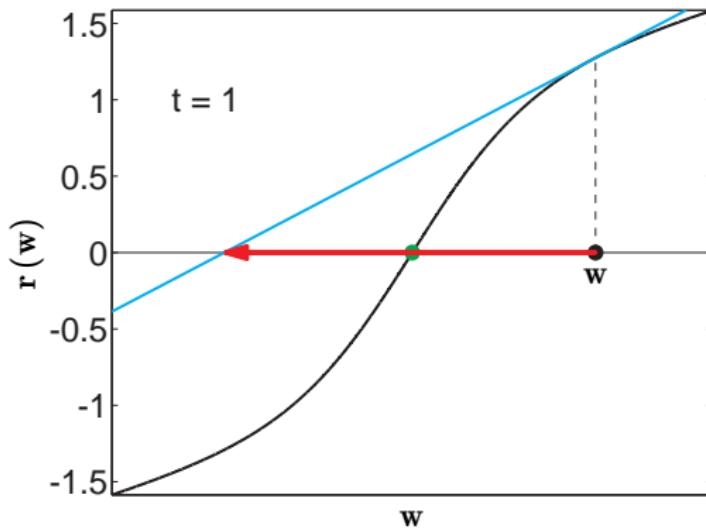


## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$

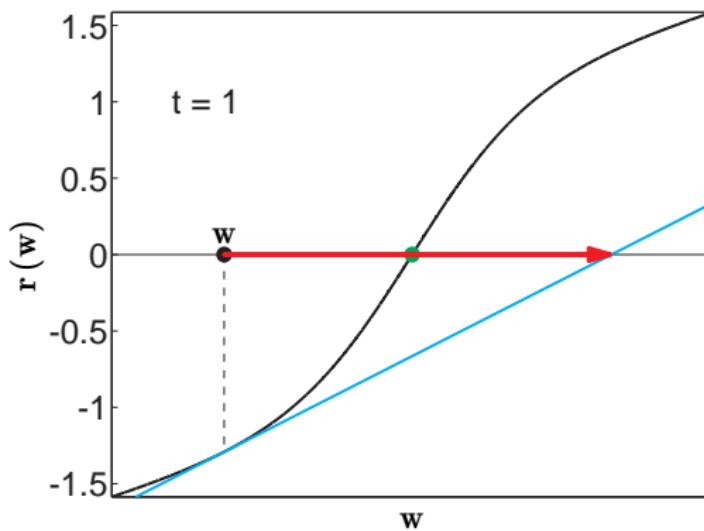


## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$

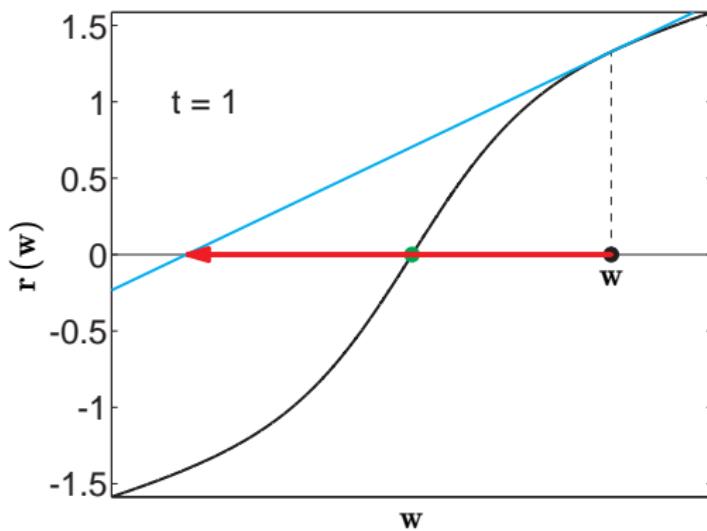


## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$

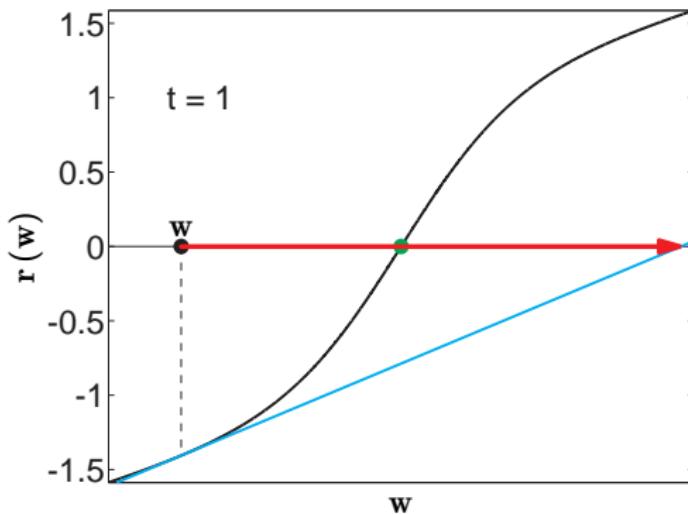


## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$

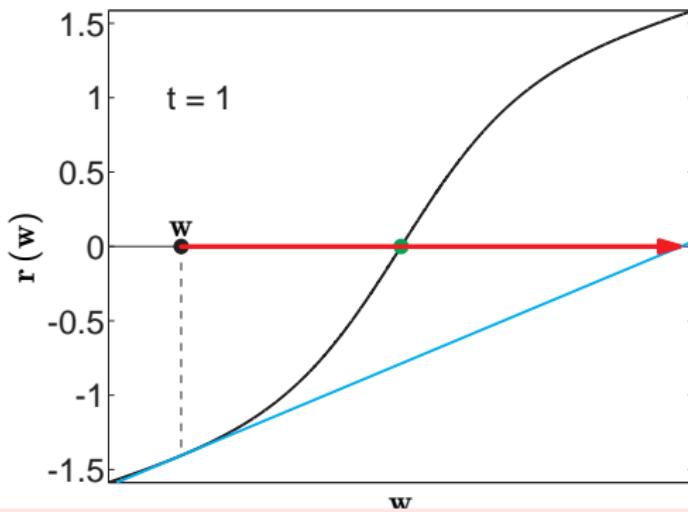


## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$



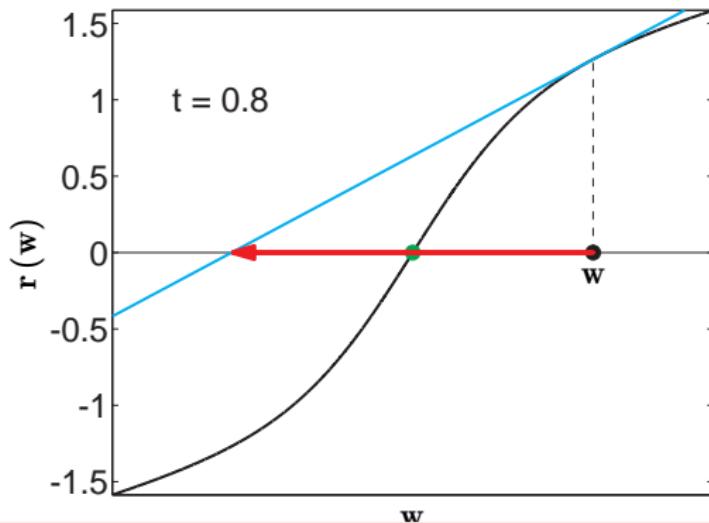
The full-step Newton iteration can be unstable !!

## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$



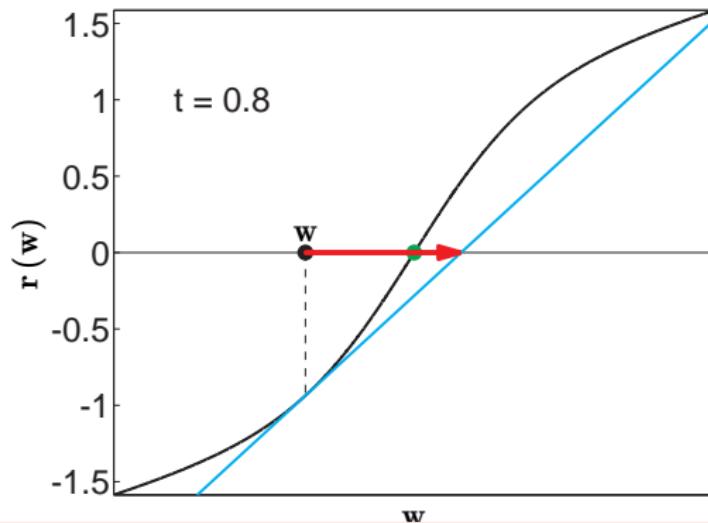
The full-step Newton iteration can be unstable !!

## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$



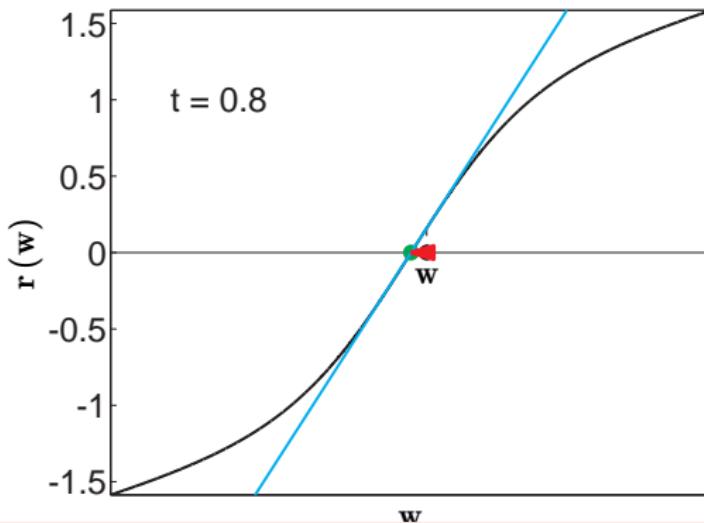
The full-step Newton iteration can be unstable !!

## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$



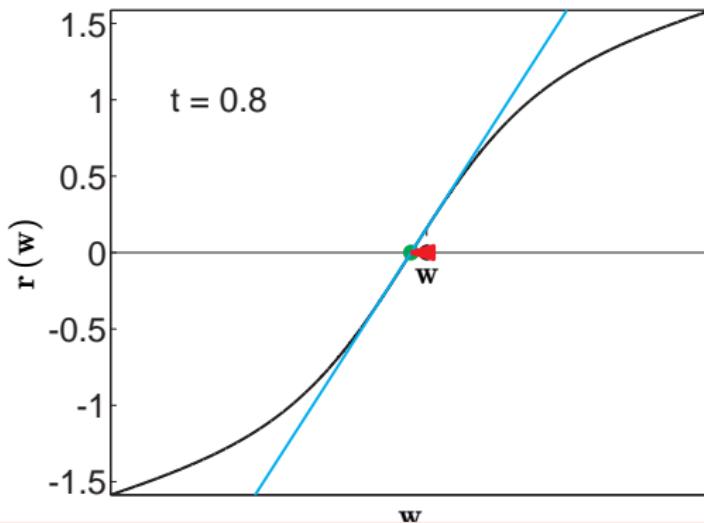
The full-step Newton iteration can be unstable !!

## Why reduced steps ?

Newton step with  $t \in ]0, 1]$ :

$$\nabla r(w) \Delta w = -r(w)$$

$$w \leftarrow w + t \Delta w$$



The full-step Newton iteration can be unstable !! Newton iteration with adequately reduced steps converges (does it ?)

## Does Newton always work ?

Is there always a reduced Newton step  $t\Delta w$  "improving"  $r(w)$  ?

## Does Newton always work ?

Is there always a reduced Newton step  $t\Delta w$  "improving"  $r(w)$  ?  
I.e. is there always a  $t > 0$  s.t.  $\|r(w + t\Delta w)\| < \|r(w)\|$  ?

## Does Newton always work ?

Is there always a reduced Newton step  $t\Delta w$  "improving"  $r(w)$  ?  
I.e. is there always a  $t > 0$  s.t.  $\|r(w + t\Delta w)\| < \|r(w)\|$  ? Yes... sort of

## Does Newton always work ?

Is there always a reduced Newton step  $t\Delta w$  "improving"  $r(w)$  ?  
I.e. is there always a  $t > 0$  s.t.  $\|r(w + t\Delta w)\| < \|r(w)\|$  ? Yes... sort of

**Proof:**  $\|r(w + t\Delta w)\| < \|r(w)\|$  holds for some  $t > 0$  if

$$\frac{d}{dt} \|r(w + t\Delta w)\|_2^2 \Big|_{t=0} < 0$$

with  $\|r(w)\|_2^2$  differentiable.

## Does Newton always work ?

Is there always a reduced Newton step  $t\Delta w$  "improving"  $r(w)$  ?  
I.e. is there always a  $t > 0$  s.t.  $\|r(w + t\Delta w)\| < \|r(w)\|$  ? Yes... sort of

**Proof:**  $\|r(w + t\Delta w)\| < \|r(w)\|$  holds for some  $t > 0$  if

$$\frac{d}{dt} \|r(w + t\Delta w)\|_2^2 \Big|_{t=0} < 0$$

with  $\|r(w)\|_2^2$  differentiable. I.e.

$$2r(w)^\top \frac{d}{dt} r(w + t\Delta w)_{t=0} < 0$$

## Does Newton always work ?

Is there always a reduced Newton step  $t\Delta w$  "improving"  $r(w)$ ?  
I.e. is there always a  $t > 0$  s.t.  $\|r(w + t\Delta w)\| < \|r(w)\|$ ? Yes... sort of

**Proof:**  $\|r(w + t\Delta w)\| < \|r(w)\|$  holds for some  $t > 0$  if

$$\frac{d}{dt} \|r(w + t\Delta w)\|_2^2 \Big|_{t=0} < 0$$

with  $\|r(w)\|_2^2$  differentiable. I.e.

$$2r(w)^\top \frac{d}{dt} r(w + t\Delta w)_{t=0} < 0$$

We have

$$\frac{d}{dt} r(w + t\Delta w)_{t=0} = \nabla r(w)^\top \Delta w = -\nabla r(w)^\top \nabla r(w)^\top r(w) = -r(w)$$

## Does Newton always work ?

Is there always a reduced Newton step  $t\Delta w$  "improving"  $r(w)$ ?  
I.e. is there always a  $t > 0$  s.t.  $\|r(w + t\Delta w)\| < \|r(w)\|$ ? Yes... sort of

**Proof:**  $\|r(w + t\Delta w)\| < \|r(w)\|$  holds for some  $t > 0$  if

$$\frac{d}{dt} \|r(w + t\Delta w)\|_2^2 \Big|_{t=0} < 0$$

with  $\|r(w)\|_2^2$  differentiable. I.e.

$$2r(w)^\top \frac{d}{dt} r(w + t\Delta w)_{t=0} < 0$$

We have

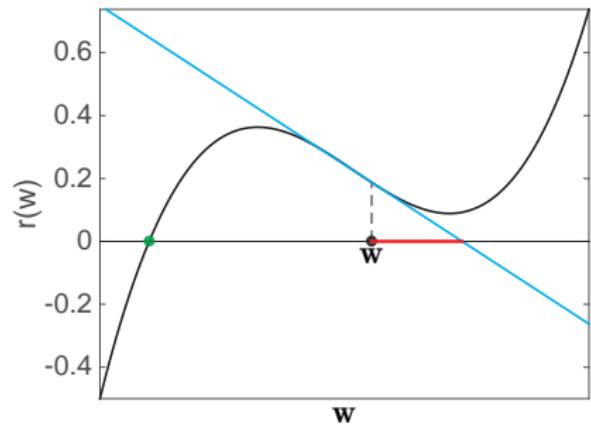
$$\frac{d}{dt} r(w + t\Delta w)_{t=0} = \nabla r(w)^\top \Delta w = -\nabla r(w)^\top \nabla r(w)^\top r(w) = -r(w)$$

Then

$$\frac{d}{dt} \|r(w + t\Delta w)\|_2^2 \Big|_{t=0} = -2\|r(w)\|_2^2 < 0$$

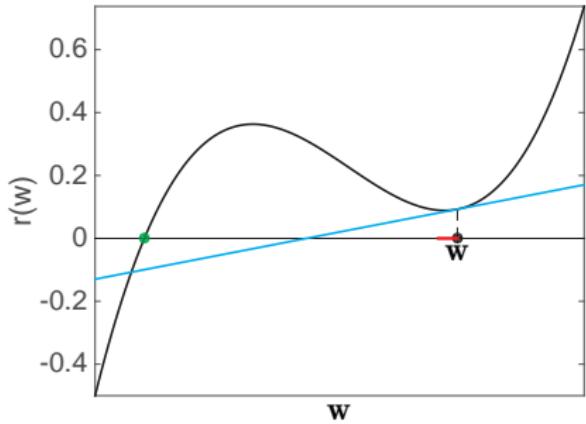
But still, Newton can fail...

Solve  $\mathbf{r}(\mathbf{w}) = 0$



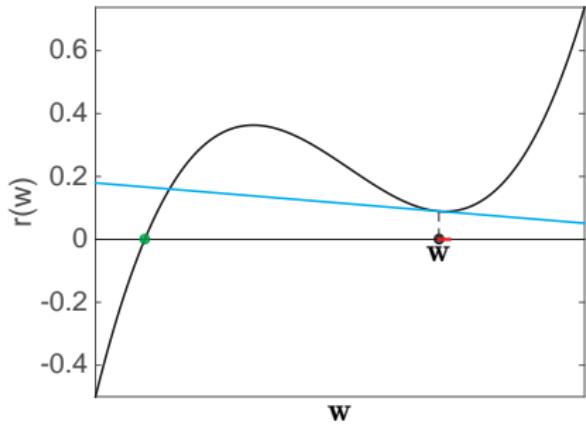
But still, Newton can fail...

Solve  $\mathbf{r}(\mathbf{w}) = 0$



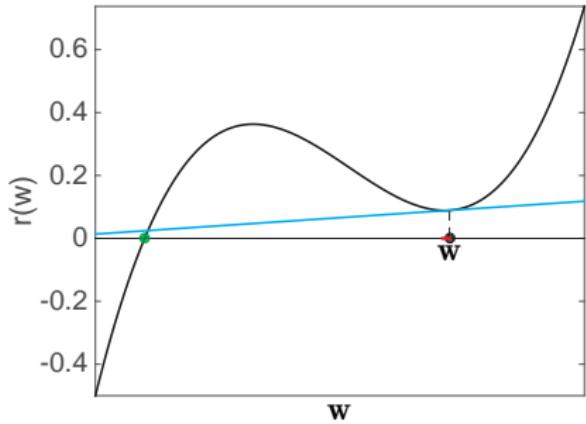
But still, Newton can fail...

Solve  $\mathbf{r}(\mathbf{w}) = 0$



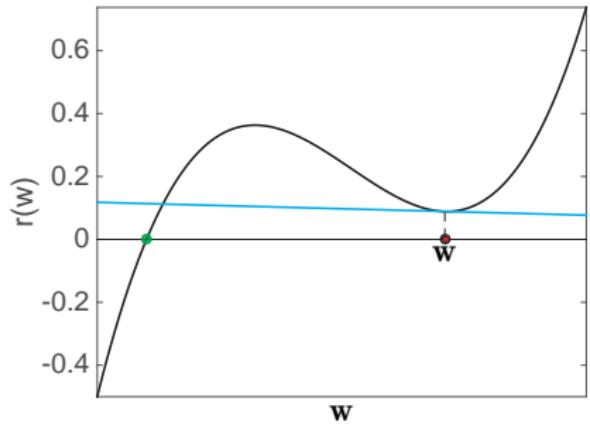
But still, Newton can fail...

Solve  $\mathbf{r}(\mathbf{w}) = 0$



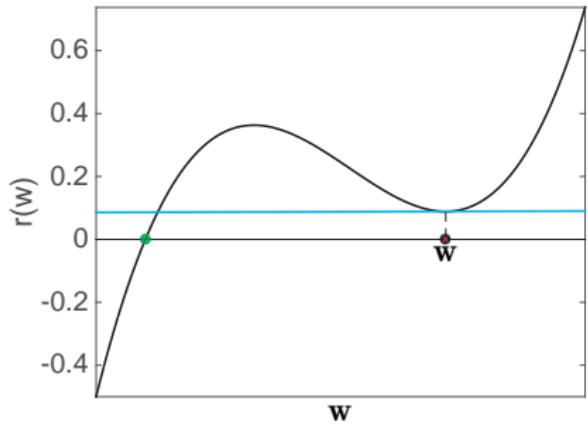
But still, Newton can fail...

Solve  $\mathbf{r}(\mathbf{w}) = 0$



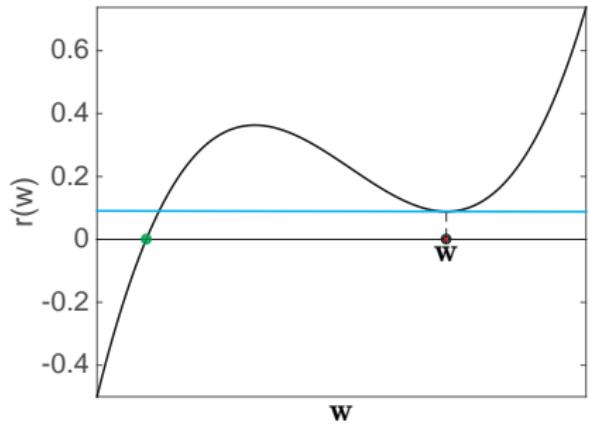
But still, Newton can fail...

Solve  $\mathbf{r}(\mathbf{w}) = 0$



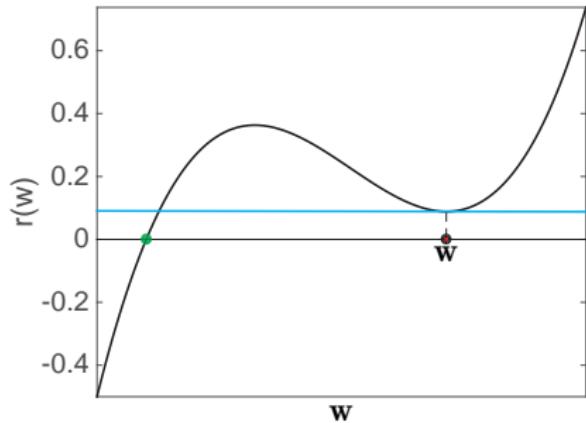
But still, Newton can fail...

Solve  $\mathbf{r}(\mathbf{w}) = 0$



But still, Newton can fail...

Solve  $\mathbf{r}(\mathbf{w}) = 0$



Newton stops with

$\mathbf{r}(\mathbf{w}) \neq 0$  and  $\nabla\mathbf{r}(\mathbf{w})$  singular

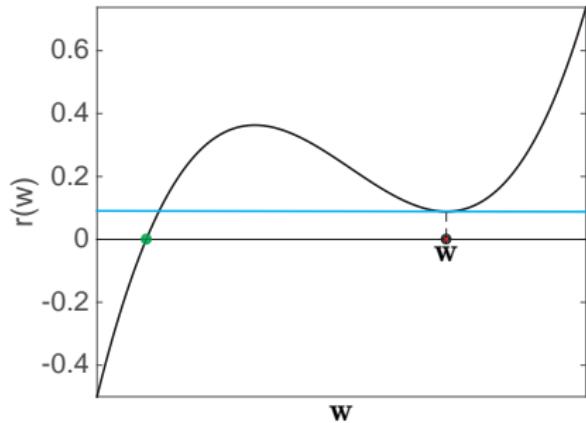
i.e. the Newton direction  $\Delta\mathbf{w}$  given by

$$\nabla\mathbf{r}(\mathbf{w}) \Delta\mathbf{w} = -\mathbf{r}(\mathbf{w})$$

is undefined

But still, Newton can fail...

Solve  $\mathbf{r}(\mathbf{w}) = 0$



Newton stops with

$\mathbf{r}(\mathbf{w}) \neq 0$  and  $\nabla\mathbf{r}(\mathbf{w})$  singular

i.e. the Newton direction  $\Delta\mathbf{w}$  given by

$$\nabla\mathbf{r}(\mathbf{w}) \Delta\mathbf{w} = -\mathbf{r}(\mathbf{w})$$

is undefined

Reduced-step Newton iteration converges to a solution of  $\mathbf{r}(\mathbf{w}) = 0$  if  $\nabla\mathbf{r}$  is full rank everywhere

## Convergence of Newton methods & Inexact Jacobian

Consider  $\mathbf{r}(\mathbf{w}^*) = 0$ , and the Newton iteration  $\mathbf{w}^+ = \mathbf{w} - \mathbf{M}^{-1}\mathbf{r}(\mathbf{w})$ .

## Convergence of Newton methods & Inexact Jacobian

Consider  $\mathbf{r}(\mathbf{w}^*) = 0$ , and the Newton iteration  $\mathbf{w}^+ = \mathbf{w} - \mathbf{M}^{-1}\mathbf{r}(\mathbf{w})$ . We then have:

$$\mathbf{w}^+ - \mathbf{w}^* = \mathbf{w} - \mathbf{w}^* - \mathbf{M}^{-1}(\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*))$$

## Convergence of Newton methods & Inexact Jacobian

Consider  $\mathbf{r}(\mathbf{w}^*) = 0$ , and the Newton iteration  $\mathbf{w}^+ = \mathbf{w} - \mathbf{M}^{-1}\mathbf{r}(\mathbf{w})$ . We then have:

$$\mathbf{w}^+ - \mathbf{w}^* = \mathbf{w} - \mathbf{w}^* - \mathbf{M}^{-1}(\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*))$$

We use the following result from analysis:

$$\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*) = \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top \cdot d\mathbf{t} \right) (\mathbf{w} - \mathbf{w}^*)$$

## Convergence of Newton methods & Inexact Jacobian

Consider  $\mathbf{r}(\mathbf{w}^*) = 0$ , and the Newton iteration  $\mathbf{w}^+ = \mathbf{w} - \mathbf{M}^{-1}\mathbf{r}(\mathbf{w})$ . We then have:

$$\mathbf{w}^+ - \mathbf{w}^* = \mathbf{w} - \mathbf{w}^* - \mathbf{M}^{-1}(\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*))$$

We use the following result from analysis:

$$\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*) = \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top \cdot d\mathbf{t} \right) (\mathbf{w} - \mathbf{w}^*)$$

$$\text{To obtain: } \mathbf{w}^+ - \mathbf{w}^* = \left( I - \mathbf{M}^{-1} \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top dt \right) (\mathbf{w} - \mathbf{w}^*)$$

## Convergence of Newton methods & Inexact Jacobian

Consider  $\mathbf{r}(\mathbf{w}^*) = 0$ , and the Newton iteration  $\mathbf{w}^+ = \mathbf{w} - \mathbf{M}^{-1}\mathbf{r}(\mathbf{w})$ . We then have:

$$\mathbf{w}^+ - \mathbf{w}^* = \mathbf{w} - \mathbf{w}^* - \mathbf{M}^{-1}(\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*))$$

We use the following result from analysis:

$$\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*) = \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top \cdot d\mathbf{t} \right) (\mathbf{w} - \mathbf{w}^*)$$

$$\text{To obtain: } \mathbf{w}^+ - \mathbf{w}^* = \left( I - \mathbf{M}^{-1} \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top dt \right) (\mathbf{w} - \mathbf{w}^*)$$

And equivalently:

$$\mathbf{w}^+ - \mathbf{w}^* = \mathbf{M}^{-1} \left( \mathbf{M} - \nabla \mathbf{r}(\mathbf{w})^\top - \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top - \nabla \mathbf{r}(\mathbf{w})^\top dt \right) (\mathbf{w} - \mathbf{w}^*)$$

## Convergence of Newton methods & Inexact Jacobian

Consider  $\mathbf{r}(\mathbf{w}^*) = 0$ , and the Newton iteration  $\mathbf{w}^+ = \mathbf{w} - \mathbf{M}^{-1}\mathbf{r}(\mathbf{w})$ . We then have:

$$\mathbf{w}^+ - \mathbf{w}^* = \mathbf{w} - \mathbf{w}^* - \mathbf{M}^{-1}(\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*))$$

We use the following result from analysis:

$$\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*) = \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top \cdot d\mathbf{t} \right) (\mathbf{w} - \mathbf{w}^*)$$

$$\text{To obtain: } \mathbf{w}^+ - \mathbf{w}^* = \left( I - \mathbf{M}^{-1} \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top dt \right) (\mathbf{w} - \mathbf{w}^*)$$

And equivalently:

$$\mathbf{w}^+ - \mathbf{w}^* = \mathbf{M}^{-1} \left( \mathbf{M} - \nabla \mathbf{r}(\mathbf{w})^\top - \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top - \nabla \mathbf{r}(\mathbf{w})^\top dt \right) (\mathbf{w} - \mathbf{w}^*)$$

So that:

$$\begin{aligned} \|\mathbf{w}^+ - \mathbf{w}^*\| &\leq \left\| \mathbf{M}^{-1} \left( \mathbf{M} - \nabla \mathbf{r}(\mathbf{w})^\top \right) \right\| \cdot \|(\mathbf{w} - \mathbf{w}^*)\| \\ &\quad + \left\| \mathbf{M}^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top - \nabla \mathbf{r}(\mathbf{w})^\top dt \right) \right\| \cdot \|(\mathbf{w} - \mathbf{w}^*)\| \end{aligned}$$

## Convergence of Newton methods & Inexact Jacobian

Consider  $\mathbf{r}(\mathbf{w}^*) = 0$ , and the Newton iteration  $\mathbf{w}^+ = \mathbf{w} - \mathbf{M}^{-1}\mathbf{r}(\mathbf{w})$ . We then have:

$$\mathbf{w}^+ - \mathbf{w}^* = \mathbf{w} - \mathbf{w}^* - \mathbf{M}^{-1}(\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*))$$

We use the following result from analysis:

$$\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*) = \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top \cdot d\mathbf{t} \right) (\mathbf{w} - \mathbf{w}^*)$$

$$\text{To obtain: } \mathbf{w}^+ - \mathbf{w}^* = \left( I - \mathbf{M}^{-1} \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top dt \right) (\mathbf{w} - \mathbf{w}^*)$$

And equivalently:

$$\mathbf{w}^+ - \mathbf{w}^* = \mathbf{M}^{-1} \left( \mathbf{M} - \nabla \mathbf{r}(\mathbf{w})^\top - \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top - \nabla \mathbf{r}(\mathbf{w})^\top dt \right) (\mathbf{w} - \mathbf{w}^*)$$

So that:

$$\begin{aligned} \|\mathbf{w}^+ - \mathbf{w}^*\| &\leq \underbrace{\left\| \mathbf{M}^{-1} \left( \mathbf{M} - \nabla \mathbf{r}(\mathbf{w})^\top \right) \right\|}_{\text{small if } \mathbf{M} \approx \nabla \mathbf{r}(\mathbf{w})^\top} \cdot \|(\mathbf{w} - \mathbf{w}^*)\| \\ &\quad + \left\| \mathbf{M}^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top - \nabla \mathbf{r}(\mathbf{w})^\top dt \right) \right\| \cdot \|(\mathbf{w} - \mathbf{w}^*)\| \end{aligned}$$

## Convergence of Newton methods & Inexact Jacobian

Consider  $\mathbf{r}(\mathbf{w}^*) = 0$ , and the Newton iteration  $\mathbf{w}^+ = \mathbf{w} - M^{-1}\mathbf{r}(\mathbf{w})$ . We then have:

$$\mathbf{w}^+ - \mathbf{w}^* = \mathbf{w} - \mathbf{w}^* - M^{-1}(\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*))$$

We use the following result from analysis:

$$\mathbf{r}(\mathbf{w}) - \mathbf{r}(\mathbf{w}^*) = \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top \cdot d\mathbf{t} \right) (\mathbf{w} - \mathbf{w}^*)$$

$$\text{To obtain: } \mathbf{w}^+ - \mathbf{w}^* = \left( I - M^{-1} \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top dt \right) (\mathbf{w} - \mathbf{w}^*)$$

And equivalently:

$$\mathbf{w}^+ - \mathbf{w}^* = M^{-1} \left( M - \nabla \mathbf{r}(\mathbf{w})^\top - \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top - \nabla \mathbf{r}(\mathbf{w})^\top dt \right) (\mathbf{w} - \mathbf{w}^*)$$

So that:

$$\begin{aligned} \|\mathbf{w}^+ - \mathbf{w}^*\| &\leq \underbrace{\left\| M^{-1} \left( M - \nabla \mathbf{r}(\mathbf{w})^\top \right) \right\|}_{\text{small if } M - \nabla \mathbf{r}(\mathbf{w})^\top} \cdot \|\mathbf{w} - \mathbf{w}^*\| \\ &\quad + \underbrace{\left\| M^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w} + t(\mathbf{w}^* - \mathbf{w}))^\top - \nabla \mathbf{r}(\mathbf{w})^\top dt \right) \right\|}_{\text{small if } \nabla \mathbf{r} \text{ is not changing much between } \mathbf{w} \text{ and } \mathbf{w}^*} \cdot \|\mathbf{w} - \mathbf{w}^*\| \end{aligned}$$

## Convergence of Newton methods

**Theorem:** consider the Newton iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1}\mathbf{r}(\mathbf{w}_k)$  and assume

- Lipschitz condition:  $\|M_k^{-1}(\nabla\mathbf{r}(\mathbf{w})^\top - \nabla\mathbf{r}(\mathbf{w}^*))^\top\| \leq \omega \|\mathbf{w} - \mathbf{w}^*\|$

## Convergence of Newton methods

**Theorem:** consider the Newton iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1}\mathbf{r}(\mathbf{w}_k)$  and assume

- Lipschitz condition:  $\|M_k^{-1}(\nabla\mathbf{r}(\mathbf{w})^\top - \nabla\mathbf{r}(\mathbf{w}^*))^\top\| \leq \omega \|\mathbf{w} - \mathbf{w}^*\|$
- Bound on the Jacobian approximation error  $\|M_k^{-1}(\nabla\mathbf{r}(\mathbf{w}_k)^\top - M_k)\| \leq \kappa_k < \kappa < 1$

## Convergence of Newton methods

**Theorem:** consider the Newton iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1}\mathbf{r}(\mathbf{w}_k)$  and assume

- Lipschitz condition:  $\|M_k^{-1}(\nabla\mathbf{r}(\mathbf{w})^\top - \nabla\mathbf{r}(\mathbf{w}^*))^\top\| \leq \omega \|\mathbf{w} - \mathbf{w}^*\|$
- Bound on the Jacobian approximation error  $\|M_k^{-1}(\nabla\mathbf{r}(\mathbf{w}_k)^\top - M_k)\| \leq \kappa_k < \kappa < 1$
- Good initial guess  $\|\mathbf{w}_0 - \mathbf{w}^*\| \leq \frac{2}{\omega}(1 - \kappa)$

Then  $\mathbf{w}_k \rightarrow \mathbf{w}^*$  with the following linear-quadratic contraction rate:

$$\|\mathbf{w}_{k+1} - \mathbf{w}^*\| \leq \left( \kappa_k + \frac{\omega}{2} \|\mathbf{w}_k - \mathbf{w}^*\| \right) \|\mathbf{w}_k - \mathbf{w}^*\|.$$

## Convergence of Newton methods

**Theorem:** consider the Newton iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1} \mathbf{r}(\mathbf{w}_k)$  and assume

- Lipschitz condition:  $\|M_k^{-1}(\nabla \mathbf{r}(\mathbf{w})^\top - \nabla \mathbf{r}(\mathbf{w}^*))^\top\| \leq \omega \|\mathbf{w} - \mathbf{w}^*\|$
- Bound on the Jacobian approximation error  $\|M_k^{-1}(\nabla \mathbf{r}(\mathbf{w}_k)^\top - M_k)\| \leq \kappa_k < \kappa < 1$
- Good initial guess  $\|\mathbf{w}_0 - \mathbf{w}^*\| \leq \frac{2}{\omega}(1 - \kappa)$

Then  $\mathbf{w}_k \rightarrow \mathbf{w}^*$  with the following linear-quadratic contraction rate:

$$\|\mathbf{w}_{k+1} - \mathbf{w}^*\| \leq \left( \kappa_k + \frac{\omega}{2} \|\mathbf{w}_k - \mathbf{w}^*\| \right) \|\mathbf{w}_k - \mathbf{w}^*\|.$$

**Proof**

$$\begin{aligned} \|\mathbf{w}_{k+1} - \mathbf{w}^*\| &\leq \overbrace{\left\| M_k^{-1} \left( M_k - \nabla \mathbf{r}(\mathbf{w}_k)^\top \right) \right\|}^{\leq \kappa_k} \cdot \|\mathbf{w}_k - \mathbf{w}^*\| \\ &\quad + \left\| M_k^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top dt \right) \right\| \cdot \|\mathbf{w}_k - \mathbf{w}^*\| \end{aligned}$$

## Convergence of Newton methods

**Theorem:** consider the Newton iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1} \mathbf{r}(\mathbf{w}_k)$  and assume

- Lipschitz condition:  $\|M_k^{-1}(\nabla \mathbf{r}(\mathbf{w})^\top - \nabla \mathbf{r}(\mathbf{w}^*))^\top\| \leq \omega \|\mathbf{w} - \mathbf{w}^*\|$
- Bound on the Jacobian approximation error  $\|M_k^{-1}(\nabla \mathbf{r}(\mathbf{w}_k)^\top - M_k)\| \leq \kappa_k < \kappa < 1$
- Good initial guess  $\|\mathbf{w}_0 - \mathbf{w}^*\| \leq \frac{2}{\omega}(1 - \kappa)$

Then  $\mathbf{w}_k \rightarrow \mathbf{w}^*$  with the following linear-quadratic contraction rate:

$$\|\mathbf{w}_{k+1} - \mathbf{w}^*\| \leq \left( \kappa_k + \frac{\omega}{2} \|\mathbf{w}_k - \mathbf{w}^*\| \right) \|\mathbf{w}_k - \mathbf{w}^*\|.$$

**Proof**

$$\begin{aligned} \|\mathbf{w}_{k+1} - \mathbf{w}^*\| &\leq \overbrace{\left\| M_k^{-1} \left( M_k - \nabla \mathbf{r}(\mathbf{w}_k)^\top \right) \right\|}^{\leq \kappa_k} \cdot \|\mathbf{w}_k - \mathbf{w}^*\| \\ &\quad + \left\| M_k^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top dt \right) \right\| \cdot \|\mathbf{w}_k - \mathbf{w}^*\| \\ \text{and } &\left\| M_k^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top dt \right) \right\| \leq \\ &\int_0^1 \left\| M_k^{-1} \left( \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top \right) \right\| dt \end{aligned}$$

## Convergence of Newton methods

**Theorem:** consider the Newton iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1} \mathbf{r}(\mathbf{w}_k)$  and assume

- Lipschitz condition:  $\|M_k^{-1}(\nabla \mathbf{r}(\mathbf{w})^\top - \nabla \mathbf{r}(\mathbf{w}^*))^\top\| \leq \omega \|\mathbf{w} - \mathbf{w}^*\|$
- Bound on the Jacobian approximation error  $\|M_k^{-1}(\nabla \mathbf{r}(\mathbf{w}_k)^\top - M_k)\| \leq \kappa_k < \kappa < 1$
- Good initial guess  $\|\mathbf{w}_0 - \mathbf{w}^*\| \leq \frac{2}{\omega}(1 - \kappa)$

Then  $\mathbf{w}_k \rightarrow \mathbf{w}^*$  with the following linear-quadratic contraction rate:

$$\|\mathbf{w}_{k+1} - \mathbf{w}^*\| \leq \left( \kappa_k + \frac{\omega}{2} \|\mathbf{w}_k - \mathbf{w}^*\| \right) \|\mathbf{w}_k - \mathbf{w}^*\|.$$

**Proof**

$$\begin{aligned} \|\mathbf{w}_{k+1} - \mathbf{w}^*\| &\leq \overbrace{\left\| M_k^{-1} \left( M_k - \nabla \mathbf{r}(\mathbf{w}_k)^\top \right) \right\|}^{\leq \kappa_k} \cdot \|\mathbf{w}_k - \mathbf{w}^*\| \\ &\quad + \left\| M_k^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top dt \right) \right\| \cdot \|\mathbf{w}_k - \mathbf{w}^*\| \end{aligned}$$

$$\text{and } \left\| M_k^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top dt \right) \right\| \leq \int_0^1 \left\| M_k^{-1} \left( \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top \right) \right\| dt$$

$$\leq \int_0^1 \omega \|\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k) - \mathbf{w}_k\| dt$$

## Convergence of Newton methods

**Theorem:** consider the Newton iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1} \mathbf{r}(\mathbf{w}_k)$  and assume

- Lipschitz condition:  $\|M_k^{-1}(\nabla \mathbf{r}(\mathbf{w})^\top - \nabla \mathbf{r}(\mathbf{w}^*))^\top\| \leq \omega \|\mathbf{w} - \mathbf{w}^*\|$
- Bound on the Jacobian approximation error  $\|M_k^{-1}(\nabla \mathbf{r}(\mathbf{w}_k)^\top - M_k)\| \leq \kappa_k < \kappa < 1$
- Good initial guess  $\|\mathbf{w}_0 - \mathbf{w}^*\| \leq \frac{2}{\omega}(1 - \kappa)$

Then  $\mathbf{w}_k \rightarrow \mathbf{w}^*$  with the following linear-quadratic contraction rate:

$$\|\mathbf{w}_{k+1} - \mathbf{w}^*\| \leq \left( \kappa_k + \frac{\omega}{2} \|\mathbf{w}_k - \mathbf{w}^*\| \right) \|\mathbf{w}_k - \mathbf{w}^*\|.$$

**Proof**

$$\begin{aligned} \|\mathbf{w}_{k+1} - \mathbf{w}^*\| &\leq \overbrace{\left\| M_k^{-1} \left( M_k - \nabla \mathbf{r}(\mathbf{w}_k)^\top \right) \right\|}^{\leq \kappa_k} \cdot \|\mathbf{w}_k - \mathbf{w}^*\| \\ &\quad + \left\| M_k^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top dt \right) \right\| \cdot \|\mathbf{w}_k - \mathbf{w}^*\| \\ \text{and } &\left\| M_k^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top dt \right) \right\| \leq \\ &\int_0^1 \left\| M_k^{-1} \left( \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top \right) \right\| dt \\ &\leq \int_0^1 \omega \|\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k) - \mathbf{w}_k\| dt \leq \int_0^1 t\omega \|\mathbf{w}^* - \mathbf{w}_k\| dt \end{aligned}$$

## Convergence of Newton methods

**Theorem:** consider the Newton iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1} \mathbf{r}(\mathbf{w}_k)$  and assume

- Lipschitz condition:  $\|M_k^{-1}(\nabla \mathbf{r}(\mathbf{w})^\top - \nabla \mathbf{r}(\mathbf{w}^*))^\top\| \leq \omega \|\mathbf{w} - \mathbf{w}^*\|$
- Bound on the Jacobian approximation error  $\|M_k^{-1}(\nabla \mathbf{r}(\mathbf{w}_k)^\top - M_k)\| \leq \kappa_k < \kappa < 1$
- Good initial guess  $\|\mathbf{w}_0 - \mathbf{w}^*\| \leq \frac{2}{\omega}(1 - \kappa)$

Then  $\mathbf{w}_k \rightarrow \mathbf{w}^*$  with the following linear-quadratic contraction rate:

$$\|\mathbf{w}_{k+1} - \mathbf{w}^*\| \leq \left( \kappa_k + \frac{\omega}{2} \|\mathbf{w}_k - \mathbf{w}^*\| \right) \|\mathbf{w}_k - \mathbf{w}^*\|.$$

**Proof**

$$\begin{aligned} \|\mathbf{w}_{k+1} - \mathbf{w}^*\| &\leq \overbrace{\left\| M_k^{-1} \left( M_k - \nabla \mathbf{r}(\mathbf{w}_k)^\top \right) \right\|}^{\leq \kappa_k} \cdot \|\mathbf{w}_k - \mathbf{w}^*\| \\ &\quad + \left\| M_k^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top dt \right) \right\| \cdot \|\mathbf{w}_k - \mathbf{w}^*\| \end{aligned}$$

$$\text{and } \left\| M_k^{-1} \left( \int_0^1 \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top dt \right) \right\| \leq \int_0^1 \left\| M_k^{-1} \left( \nabla \mathbf{r}(\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k))^\top - \nabla \mathbf{r}(\mathbf{w}_k)^\top \right) \right\| dt$$

$$\leq \int_0^1 \omega \|\mathbf{w}_k + t(\mathbf{w}^* - \mathbf{w}_k) - \mathbf{w}_k\| dt \leq \int_0^1 t\omega \|\mathbf{w}^* - \mathbf{w}_k\| dt = \frac{\omega}{2} \|\mathbf{w}^* - \mathbf{w}_k\|$$

## Affine invariance of the (exact) Newton method

### Affine change of coordinates

Consider:  $\mathbf{w} = A\mathbf{v} + \mathbf{a}$  with  $A \in \mathbb{R}^{n \times n}$  non-singular and  $\mathbf{a} \in \mathbb{R}^n$ .

## Affine invariance of the (exact) Newton method

### Affine change of coordinates

Consider:  $\mathbf{w} = A\mathbf{v} + \mathbf{a}$  with  $A \in \mathbb{R}^{n \times n}$  non-singular and  $\mathbf{a} \in \mathbb{R}^n$ .

Define  $\tilde{\mathbf{r}}(\mathbf{v}) = \mathbf{r}(A\mathbf{v} + \mathbf{a})$  then:  $\nabla_{\mathbf{v}} \tilde{\mathbf{r}}(\mathbf{v}) = A^\top \nabla_{\mathbf{w}} \mathbf{r}(\mathbf{w})$

## Affine invariance of the (exact) Newton method

### Affine change of coordinates

Consider:  $\mathbf{w} = A\mathbf{v} + \mathbf{a}$  with  $A \in \mathbb{R}^{n \times n}$  non-singular and  $\mathbf{a} \in \mathbb{R}^n$ .

Define  $\tilde{\mathbf{r}}(\mathbf{v}) = \mathbf{r}(A\mathbf{v} + \mathbf{a})$  then:  $\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v}) = A^\top \nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})$

The Newton step in  $\mathbf{v}$  given by:

$$\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v})^\top \Delta\mathbf{v} = -\tilde{\mathbf{r}}(\mathbf{v})$$

## Affine invariance of the (exact) Newton method

### Affine change of coordinates

Consider:  $\mathbf{w} = A\mathbf{v} + \mathbf{a}$  with  $A \in \mathbb{R}^{n \times n}$  non-singular and  $\mathbf{a} \in \mathbb{R}^n$ .

Define  $\tilde{\mathbf{r}}(\mathbf{v}) = \mathbf{r}(A\mathbf{v} + \mathbf{a})$  then:  $\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v}) = A^\top \nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})$

The Newton step in  $\mathbf{v}$  given by:

$$\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v})^\top \Delta\mathbf{v} = -\tilde{\mathbf{r}}(\mathbf{v})$$

also reads as

$$\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})^\top A\Delta\mathbf{v} = -\mathbf{r}(A\mathbf{v} + \mathbf{a}) = -\mathbf{r}(\mathbf{w})$$

## Affine invariance of the (exact) Newton method

### Affine change of coordinates

Consider:  $\mathbf{w} = A\mathbf{v} + \mathbf{a}$  with  $A \in \mathbb{R}^{n \times n}$  non-singular and  $\mathbf{a} \in \mathbb{R}^n$ .

Define  $\tilde{\mathbf{r}}(\mathbf{v}) = \mathbf{r}(A\mathbf{v} + \mathbf{a})$  then:  $\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v}) = A^\top \nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})$

The Newton step in  $\mathbf{v}$  given by:

$$\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v})^\top \Delta\mathbf{v} = -\tilde{\mathbf{r}}(\mathbf{v})$$

also reads as

$$\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})^\top A\Delta\mathbf{v} = -\mathbf{r}(A\mathbf{v} + \mathbf{a}) = -\mathbf{r}(\mathbf{w})$$

It follows that  $\Delta\mathbf{w} = A\Delta\mathbf{v}$  is the Newton step on  $\mathbf{r}(\mathbf{w})$ :

$$\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})^\top \Delta\mathbf{w} = -\mathbf{r}(\mathbf{w})$$

## Affine invariance of the (exact) Newton method

### Affine change of coordinates

Consider:  $\mathbf{w} = A\mathbf{v} + \mathbf{a}$  with  $A \in \mathbb{R}^{n \times n}$  non-singular and  $\mathbf{a} \in \mathbb{R}^n$ .

Define  $\tilde{\mathbf{r}}(\mathbf{v}) = \mathbf{r}(A\mathbf{v} + \mathbf{a})$  then:  $\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v}) = A^\top \nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})$

The Newton step in  $\mathbf{v}$  given by:

$$\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v})^\top \Delta\mathbf{v} = -\tilde{\mathbf{r}}(\mathbf{v})$$

also reads as

$$\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})^\top A\Delta\mathbf{v} = -\mathbf{r}(A\mathbf{v} + \mathbf{a}) = -\mathbf{r}(\mathbf{w})$$

It follows that  $\Delta\mathbf{w} = A\Delta\mathbf{v}$  is the Newton step on  $\mathbf{r}(\mathbf{w})$ :

$$\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})^\top \Delta\mathbf{w} = -\mathbf{r}(\mathbf{w})$$

Original Newton step  $\Delta\mathbf{w}$  and "transformed" step  $\Delta\mathbf{v}$  are linearly related !!

## Affine invariance of the (exact) Newton method

### Affine change of coordinates

Consider:  $\mathbf{w} = A\mathbf{v} + \mathbf{a}$  with  $A \in \mathbb{R}^{n \times n}$  non-singular and  $\mathbf{a} \in \mathbb{R}^n$ .

Define  $\tilde{\mathbf{r}}(\mathbf{v}) = \mathbf{r}(A\mathbf{v} + \mathbf{a})$  then:  $\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v}) = A^\top \nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})$

The Newton step in  $\mathbf{v}$  given by:

$$\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v})^\top \Delta\mathbf{v} = -\tilde{\mathbf{r}}(\mathbf{v})$$

also reads as

$$\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})^\top A\Delta\mathbf{v} = -\mathbf{r}(A\mathbf{v} + \mathbf{a}) = -\mathbf{r}(\mathbf{w})$$

It follows that  $\Delta\mathbf{w} = A\Delta\mathbf{v}$  is the Newton step on  $\mathbf{r}(\mathbf{w})$ :

$$\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})^\top \Delta\mathbf{w} = -\mathbf{r}(\mathbf{w})$$

Original Newton step  $\Delta\mathbf{w}$  and "transformed" step  $\Delta\mathbf{v}$  are linearly related !!

- Scaling does not affect the behaviour of the exact Newton method (iterations linearly related). Inexact Newton method can be affected !

## Affine invariance of the (exact) Newton method

### Affine change of coordinates

Consider:  $\mathbf{w} = A\mathbf{v} + \mathbf{a}$  with  $A \in \mathbb{R}^{n \times n}$  non-singular and  $\mathbf{a} \in \mathbb{R}^n$ .

Define  $\tilde{\mathbf{r}}(\mathbf{v}) = \mathbf{r}(A\mathbf{v} + \mathbf{a})$  then:  $\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v}) = A^\top \nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})$

The Newton step in  $\mathbf{v}$  given by:

$$\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v})^\top \Delta\mathbf{v} = -\tilde{\mathbf{r}}(\mathbf{v})$$

also reads as

$$\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})^\top A\Delta\mathbf{v} = -\mathbf{r}(A\mathbf{v} + \mathbf{a}) = -\mathbf{r}(\mathbf{w})$$

It follows that  $\Delta\mathbf{w} = A\Delta\mathbf{v}$  is the Newton step on  $\mathbf{r}(\mathbf{w})$ :

$$\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})^\top \Delta\mathbf{w} = -\mathbf{r}(\mathbf{w})$$

Original Newton step  $\Delta\mathbf{w}$  and "transformed" step  $\Delta\mathbf{v}$  are linearly related !!

- Scaling does not affect the behaviour of the exact Newton method (iterations linearly related). Inexact Newton method can be affected !
- Convergence proof is not scale-invariant ! Check out self-concordance theory to address that question !

## Affine invariance of the (exact) Newton method

### Affine change of coordinates

Consider:  $\mathbf{w} = A\mathbf{v} + \mathbf{a}$  with  $A \in \mathbb{R}^{n \times n}$  non-singular and  $\mathbf{a} \in \mathbb{R}^n$ .

Define  $\tilde{\mathbf{r}}(\mathbf{v}) = \mathbf{r}(A\mathbf{v} + \mathbf{a})$  then:  $\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v}) = A^\top \nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})$

The Newton step in  $\mathbf{v}$  given by:

$$\nabla_{\mathbf{v}}\tilde{\mathbf{r}}(\mathbf{v})^\top \Delta\mathbf{v} = -\tilde{\mathbf{r}}(\mathbf{v})$$

also reads as

$$\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})^\top A\Delta\mathbf{v} = -\mathbf{r}(A\mathbf{v} + \mathbf{a}) = -\mathbf{r}(\mathbf{w})$$

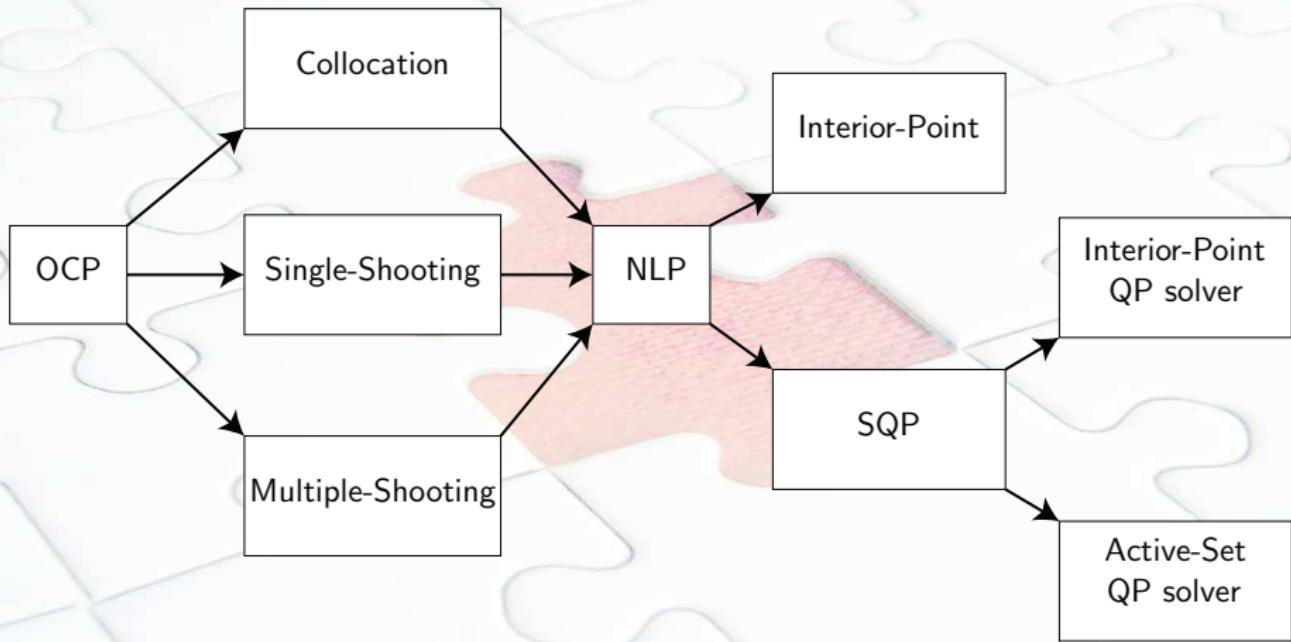
It follows that  $\Delta\mathbf{w} = A\Delta\mathbf{v}$  is the Newton step on  $\mathbf{r}(\mathbf{w})$ :

$$\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})^\top \Delta\mathbf{w} = -\mathbf{r}(\mathbf{w})$$

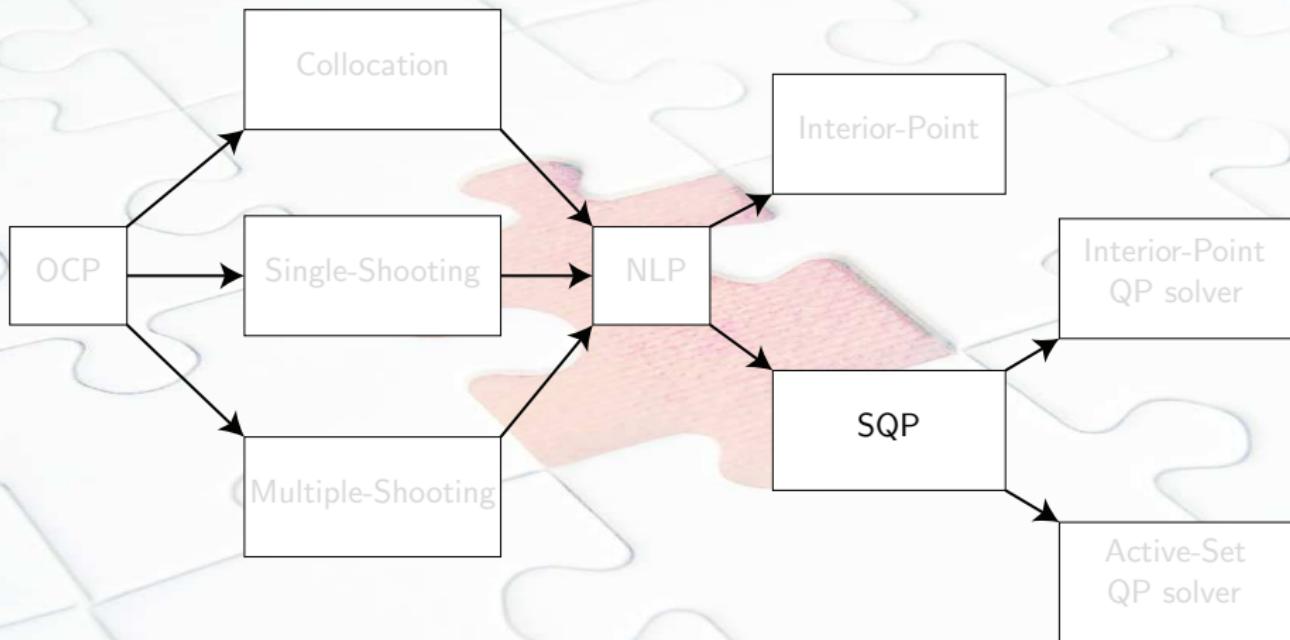
Original Newton step  $\Delta\mathbf{w}$  and "transformed" step  $\Delta\mathbf{v}$  are linearly related !!

- Scaling does not affect the behaviour of the exact Newton method (iterations linearly related). Inexact Newton method can be affected !
- Convergence proof is not scale-invariant ! Check out self-concordance theory to address that question !
- Rescaling** can still be very useful to **improve conditioning** of  $\nabla_{\mathbf{w}}\mathbf{r}(\mathbf{w})$

# Survival map of Direct Optimal Control



# Survival map of Direct Optimal Control



Let's approach the problem of solving the KKT conditions

# Outline

- 1 The Newton method
- 2 Newton on the KKT conditions
- 3 The reduced Newton step (unconstrained problems)
- 4 The merit function - Line-search for constrained problems
- 5 Newton-type methods
- 6 Sequential Quadratic Programming

## Core idea

A vast majority of solvers try to find a point  $w, \mu, \lambda$  satisfying the KKT conditions:

**Primal Feasibility:**  $g(w) = 0, h(w) \leq 0,$

**Dual Feasibility:**  $\nabla_w \mathcal{L}(w, \mu, \lambda) = 0, \mu \geq 0,$

**Complementarity Slackness:**  $\mu_i h_i(w) = 0, i = 1, \dots$

where  $\mathcal{L} = \Phi(w) + \lambda^\top g(w) + \mu^\top h(w)$

## Core idea

A vast majority of solvers try to find a point  $\mathbf{w}, \boldsymbol{\mu}, \boldsymbol{\lambda}$  satisfying the KKT conditions:

$$\text{Primal Feasibility: } g(\mathbf{w}) = 0, \quad h(\mathbf{w}) \leq 0,$$

$$\text{Dual Feasibility: } \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = 0, \quad \boldsymbol{\mu} \geq 0,$$

$$\text{Complementarity Slackness: } \boldsymbol{\mu}_i h_i(\mathbf{w}) = 0, \quad i = 1, \dots$$

$$\text{where } \mathcal{L} = \Phi(\mathbf{w}) + \boldsymbol{\lambda}^\top g(\mathbf{w}) + \boldsymbol{\mu}^\top h(\mathbf{w})$$

Let's consider for now equality constrained problems, i.e. find  $\mathbf{w}, \boldsymbol{\lambda}$  s.t.:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}) = 0$$

$$g(\mathbf{w}) = 0$$

## Core idea

A vast majority of solvers try to find a point  $w, \mu, \lambda$  satisfying the KKT conditions:

$$\text{Primal Feasibility: } g(w) = 0, \quad h(w) \leq 0,$$

$$\text{Dual Feasibility: } \nabla_w \mathcal{L}(w, \mu, \lambda) = 0, \quad \mu \geq 0,$$

$$\text{Complementarity Slackness: } \mu_i h_i(w) = 0, \quad i = 1, \dots$$

$$\text{where } \mathcal{L} = \Phi(w) + \lambda^\top g(w) + \mu^\top h(w)$$

Let's consider for now equality constrained problems, i.e. find  $w, \lambda$  s.t.:

$$\nabla_w \mathcal{L}(w, \lambda) = 0$$

$$g(w) = 0$$

Idea: apply the Newton method on the KKT conditions, i.e.

Solve...

$$r(w, \lambda) = \begin{bmatrix} \nabla_w \mathcal{L}(w, \lambda) \\ g(w) \end{bmatrix} = 0$$

## Core idea

A vast majority of solvers try to find a point  $w, \mu, \lambda$  satisfying the KKT conditions:

$$\text{Primal Feasibility: } g(w) = 0, \quad h(w) \leq 0,$$

$$\text{Dual Feasibility: } \nabla_w \mathcal{L}(w, \mu, \lambda) = 0, \quad \mu \geq 0,$$

$$\text{Complementarity Slackness: } \mu_i h_i(w) = 0, \quad i = 1, \dots$$

$$\text{where } \mathcal{L} = \Phi(w) + \lambda^\top g(w) + \mu^\top h(w)$$

Let's consider for now equality constrained problems, i.e. find  $w, \lambda$  s.t.:

$$\nabla_w \mathcal{L}(w, \lambda) = 0$$

$$g(w) = 0$$

Idea: apply the Newton method on the KKT conditions, i.e.

Solve...

... by iterating

$$r(w, \lambda) = \begin{bmatrix} \nabla_w \mathcal{L}(w, \lambda) \\ g(w) \end{bmatrix} = 0 \quad \nabla r(w, \lambda)^\top \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

## Newton method on the KKT conditions

**KKT conditions:**

$$\begin{aligned}\nabla_w \mathcal{L}(w, \lambda) &= 0 \\ g(w) &= 0\end{aligned}$$

**Newton iteration:**

$$\nabla r(w, \lambda)^\top \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

**Newton direction:**

## Newton method on the KKT conditions

KKT conditions:

$$\begin{aligned}\nabla_w \mathcal{L}(w, \lambda) &= 0 \\ g(w) &= 0\end{aligned}$$

Newton iteration:

$$\nabla r(w, \lambda)^T \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

Newton direction:

$$\begin{aligned}\nabla_w^2 \mathcal{L}(w, \lambda) \Delta w + \nabla_{w,\lambda} \mathcal{L}(w, \lambda) \Delta \lambda &= -\nabla_w \mathcal{L}(w, \lambda) \\ \nabla g(w)^T \Delta w &= -g(w)\end{aligned}$$

## Newton method on the KKT conditions

KKT conditions:

$$\begin{aligned}\nabla_w \mathcal{L}(w, \lambda) &= 0 \\ g(w) &= 0\end{aligned}$$

Newton iteration:

$$\nabla r(w, \lambda)^T \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

Newton direction: using  $\nabla_w \mathcal{L}(w, \lambda) = \nabla \Phi(w) + \nabla g(w)\lambda$

$$\begin{array}{lclclcl}\nabla_w^2 \mathcal{L}(w, \lambda) \Delta w &+& \nabla_{w, \lambda} \mathcal{L}(w, \lambda) \Delta \lambda &=& -\nabla_w \mathcal{L}(w, \lambda) \\ \nabla g(w)^T \Delta w &&&=& -g(w)\end{array}$$

## Newton method on the KKT conditions

KKT conditions:

$$\begin{aligned}\nabla_w \mathcal{L}(w, \lambda) &= 0 \\ g(w) &= 0\end{aligned}$$

Newton iteration:

$$\nabla r(w, \lambda)^T \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

Newton direction: using  $\nabla_w \mathcal{L}(w, \lambda) = \nabla \Phi(w) + \nabla g(w)\lambda$

$$\begin{array}{lclclcl}\nabla_w^2 \mathcal{L}(w, \lambda) \Delta w &+& \nabla g(w) \Delta \lambda &=& -\nabla_w \mathcal{L}(w, \lambda) \\ \nabla g(w)^T \Delta w & & & & = & -g(w)\end{array}$$

## Newton method on the KKT conditions

KKT conditions:

$$\begin{aligned}\nabla_w \mathcal{L}(w, \lambda) &= 0 \\ g(w) &= 0\end{aligned}$$

Newton iteration:

$$\nabla r(w, \lambda)^\top \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

Newton direction: using  $\nabla_w \mathcal{L}(w, \lambda) = \nabla \Phi(w) + \nabla g(w)\lambda$

$$\begin{array}{lcl}\nabla_w^2 \mathcal{L}(w, \lambda) \Delta w &+& \nabla g(w) \Delta \lambda \\ \nabla g(w)^\top \Delta w &=& -\nabla \Phi(w) - \nabla g(w) \lambda \\ &=& -g(w)\end{array}$$

## Newton method on the KKT conditions

KKT conditions:

$$\begin{aligned}\nabla_w \mathcal{L}(w, \lambda) &= 0 \\ g(w) &= 0\end{aligned}$$

Newton iteration:

$$\nabla r(w, \lambda)^\top \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

Newton direction: using  $\nabla_w \mathcal{L}(w, \lambda) = \nabla \Phi(w) + \nabla g(w)\lambda$

$$\begin{array}{lcl} \nabla_w^2 \mathcal{L}(w, \lambda) \Delta w &+& \nabla g(w)(\lambda + \Delta \lambda) \\ \nabla g(w)^\top \Delta w && \\ &=& -\nabla \Phi(w) \\ &=& -g(w) \end{array}$$

## Newton method on the KKT conditions

KKT conditions:

$$\begin{aligned}\nabla_w \mathcal{L}(w, \lambda) &= 0 \\ g(w) &= 0\end{aligned}$$

Newton iteration:

$$\nabla r(w, \lambda)^\top \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

Newton direction: using  $\nabla_w \mathcal{L}(w, \lambda) = \nabla \Phi(w) + \nabla g(w)\lambda$

$$\begin{array}{lcl}\nabla_w^2 \mathcal{L}(w, \lambda) \Delta w &+& \nabla g(w)(\lambda + \Delta \lambda) \\ \nabla g(w)^\top \Delta w &&\end{array} = \begin{array}{l}-\nabla \Phi(w) \\ -g(w)\end{array}$$

### The Newton direction on the KKT conditions

$$\begin{bmatrix} \nabla_w^2 \mathcal{L}(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda + \Delta \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

## Newton method on the KKT conditions

KKT conditions:

$$\begin{aligned}\nabla_w \mathcal{L}(w, \lambda) &= 0 \\ g(w) &= 0\end{aligned}$$

Newton iteration:

$$\nabla r(w, \lambda)^T \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

Newton direction: using  $\nabla_w \mathcal{L}(w, \lambda) = \nabla \Phi(w) + \nabla g(w)\lambda$

$$\begin{array}{lcl}\nabla_w^2 \mathcal{L}(w, \lambda) \Delta w &+& \nabla g(w)(\lambda + \Delta \lambda) \\ \nabla g(w)^T \Delta w &&\end{array} = \begin{array}{l}-\nabla \Phi(w) \\ -g(w)\end{array}$$

### The Newton direction on the KKT conditions

$$\underbrace{\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix}}_{\text{KKT matrix (symmetric indefinite)}} \begin{bmatrix} \Delta w \\ \lambda + \Delta \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

where  $H(w, \lambda) = \nabla_w^2 \mathcal{L}(w, \lambda)$  is the Hessian of the problem.

## Newton method on the KKT conditions

KKT conditions:

$$\begin{aligned}\nabla_w \mathcal{L}(w, \lambda) &= 0 \\ g(w) &= 0\end{aligned}$$

Newton iteration:

$$\nabla r(w, \lambda)^\top \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

Newton direction: using  $\nabla_w \mathcal{L}(w, \lambda) = \nabla \Phi(w) + \nabla g(w)\lambda$

$$\begin{array}{lcl} \nabla_w^2 \mathcal{L}(w, \lambda) \Delta w &+& \nabla g(w)(\lambda + \Delta \lambda) \\ \nabla g(w)^\top \Delta w && = -\nabla \Phi(w) \\ && = -g(w) \end{array}$$

### The Newton direction on the KKT conditions

$$\underbrace{\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix}}_{\text{KKT matrix (symmetric indefinite)}} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

where  $H(w, \lambda) = \nabla_w^2 \mathcal{L}(w, \lambda)$  is the Hessian of the problem.

- In the following, we label  $\lambda^+ = \lambda + \Delta \lambda$  (full dual Newton step)

## Newton method on the KKT conditions

KKT conditions:

$$\begin{aligned}\nabla_w \mathcal{L}(w, \lambda) &= 0 \\ g(w) &= 0\end{aligned}$$

Newton iteration:

$$\nabla r(w, \lambda)^\top \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

Newton direction: using  $\nabla_w \mathcal{L}(w, \lambda) = \nabla \Phi(w) + \nabla g(w)\lambda$

$$\begin{array}{lcl} \nabla_w^2 \mathcal{L}(w, \lambda) \Delta w &+& \nabla g(w)(\lambda + \Delta \lambda) \\ \nabla g(w)^\top \Delta w && = -\nabla \Phi(w) \\ && = -g(w) \end{array}$$

### The Newton direction on the KKT conditions

$$\underbrace{\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix}}_{\text{KKT matrix (symmetric indefinite)}} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

where  $H(w, \lambda) = \nabla_w^2 \mathcal{L}(w, \lambda)$  is the Hessian of the problem.

- In the following, we label  $\lambda^+ = \lambda + \Delta \lambda$  (full dual Newton step)
- $\mathcal{L}(w, \lambda)$  enters only in the Hessian !

## Newton method on the KKT conditions

KKT conditions:

$$\begin{aligned}\nabla_w \mathcal{L}(w, \lambda) &= 0 \\ g(w) &= 0\end{aligned}$$

Newton iteration:

$$\nabla r(w, \lambda)^\top \begin{bmatrix} \Delta w \\ \Delta \lambda \end{bmatrix} = -r(w, \lambda)$$

Newton direction: using  $\nabla_w \mathcal{L}(w, \lambda) = \nabla \Phi(w) + \nabla g(w)\lambda$

$$\begin{array}{lcl}\nabla_w^2 \mathcal{L}(w, \lambda) \Delta w &+& \nabla g(w)(\lambda + \Delta \lambda) \\ \nabla g(w)^\top \Delta w &&\end{array} = \begin{array}{l}-\nabla \Phi(w) \\ -g(w)\end{array}$$

### The Newton direction on the KKT conditions

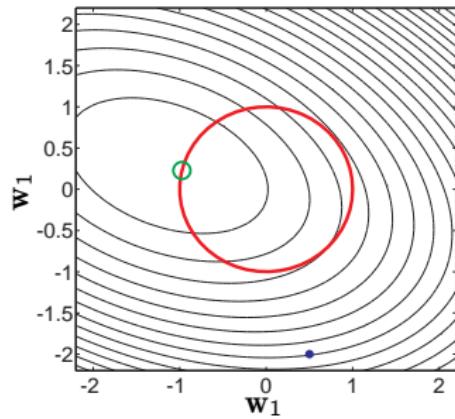
$$\underbrace{\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix}}_{\text{KKT matrix (symmetric indefinite)}} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

where  $H(w, \lambda) = \nabla_w^2 \mathcal{L}(w, \lambda)$  is the Hessian of the problem.

- In the following, we label  $\lambda^+ = \lambda + \Delta \lambda$  (full dual Newton step)
- $\mathcal{L}(w, \lambda)$  enters only in the Hessian !
- The updated multiplier  $\lambda^+$  is readily provided...

## Newton Iteration for Optimization - Example

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \mathbf{w}^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} \mathbf{w} + \mathbf{w}^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(\mathbf{w}) = \mathbf{w}^T \mathbf{w} - 1 = 0 \end{aligned}$$

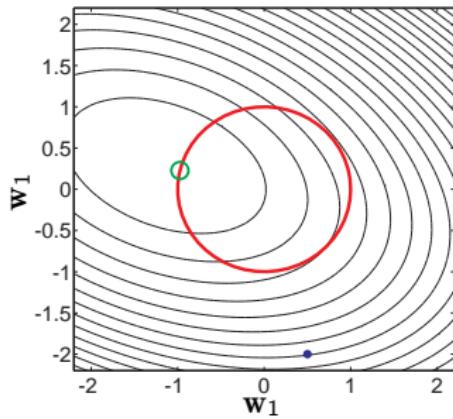


## Newton Iteration for Optimization - Example

Iterate:

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$



## Newton Iteration for Optimization - Example

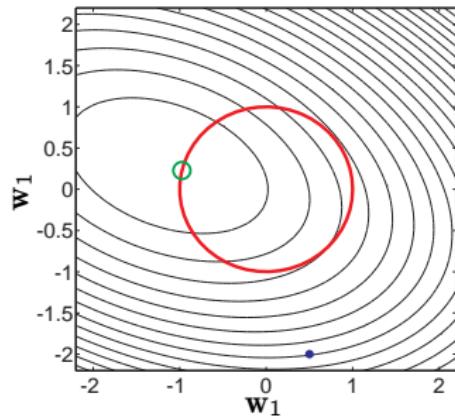
Iterate:

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

with:

$$\nabla g(w) = 2w = \begin{bmatrix} 2w_1 \\ 2w_2 \end{bmatrix}$$

$$\begin{aligned} & \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ & \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$



## Newton Iteration for Optimization - Example

Iterate:

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

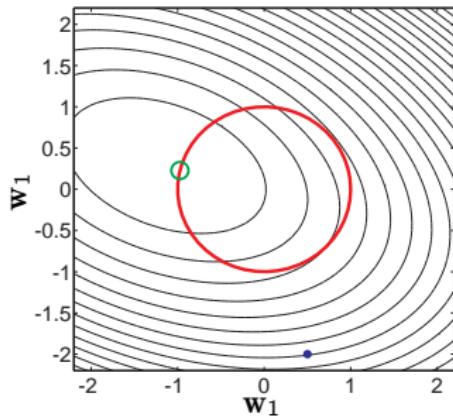
with:

$$\nabla g(w) = 2w = \begin{bmatrix} 2w_1 \\ 2w_2 \end{bmatrix}$$

$$\mathcal{L}(w, \lambda) = \Phi(w) + \lambda g(w)$$

$$\nabla_w \mathcal{L}(w, \lambda) = \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 2\lambda w$$

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$



## Newton Iteration for Optimization - Example

Iterate:

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

with:

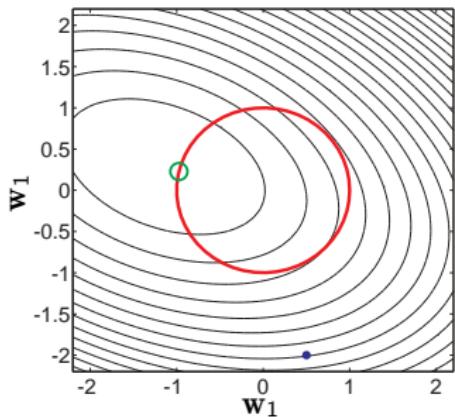
$$\nabla g(w) = 2w = \begin{bmatrix} 2w_1 \\ 2w_2 \end{bmatrix}$$

$$\mathcal{L}(w, \lambda) = \Phi(w) + \lambda g(w)$$

$$\nabla_w \mathcal{L}(w, \lambda) = \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 2\lambda w$$

$$H(w, \lambda) = \begin{bmatrix} 2 + 2\lambda & 1 \\ 1 & 4 + 2\lambda \end{bmatrix}$$

$$\begin{aligned} & \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ & \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$



## Newton Iteration for Optimization - Example

Iterate:

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

with:

$$\nabla g(w) = 2w = \begin{bmatrix} 2w_1 \\ 2w_2 \end{bmatrix}$$

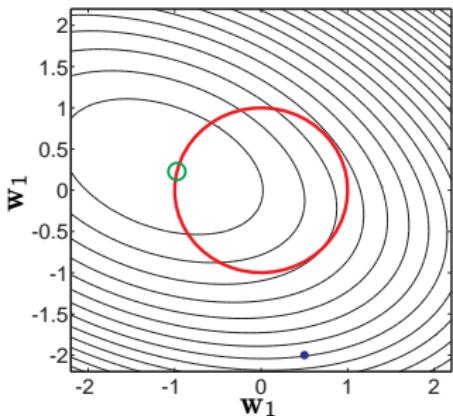
$$\mathcal{L}(w, \lambda) = \Phi(w) + \lambda g(w)$$

$$\nabla_w \mathcal{L}(w, \lambda) = \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 2\lambda w$$

$$H(w, \lambda) = \begin{bmatrix} 2 + 2\lambda & 1 \\ 1 & 4 + 2\lambda \end{bmatrix}$$

$$\nabla \Phi(w) = \begin{bmatrix} 2w_1 + w_2 + 1 \\ w_1 + 4w_2 \end{bmatrix}$$

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

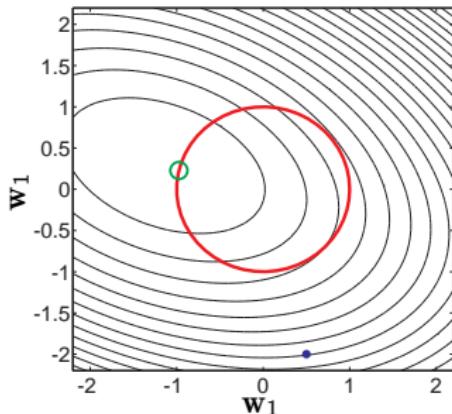
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

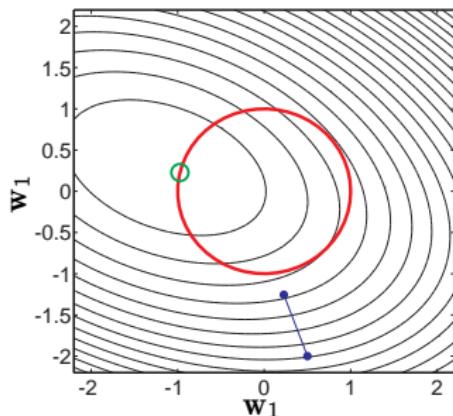
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

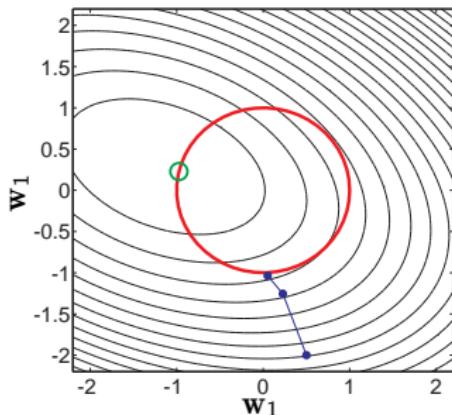
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

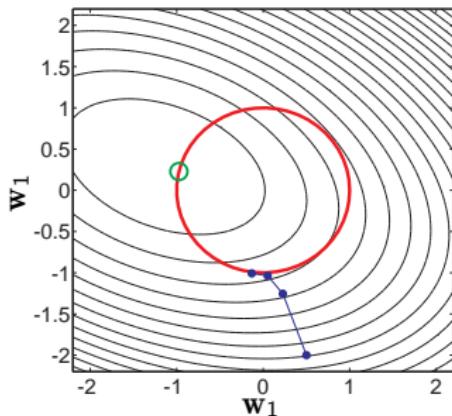
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

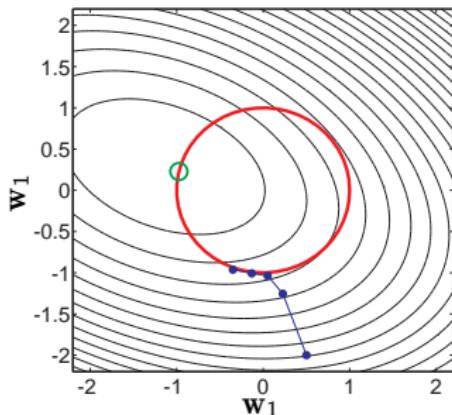
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

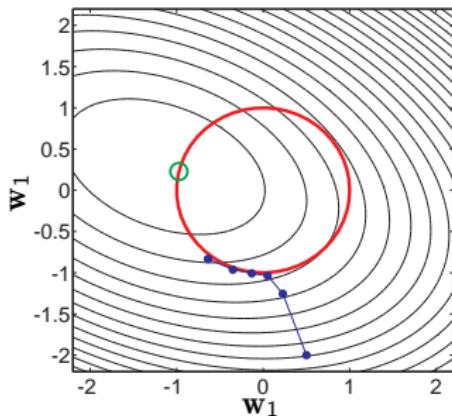
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

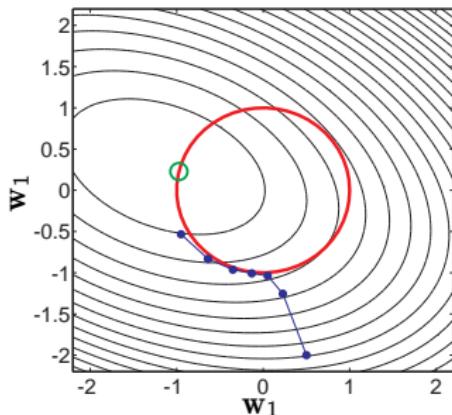
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

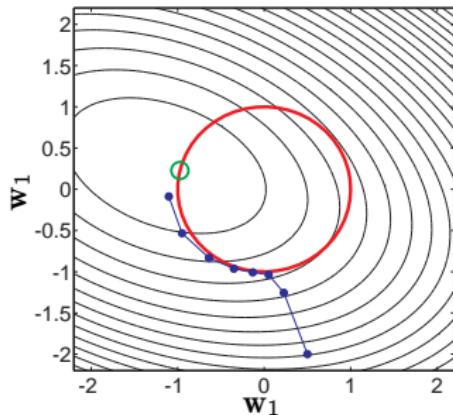
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

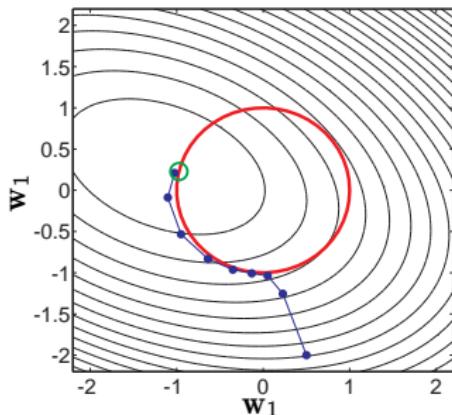
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

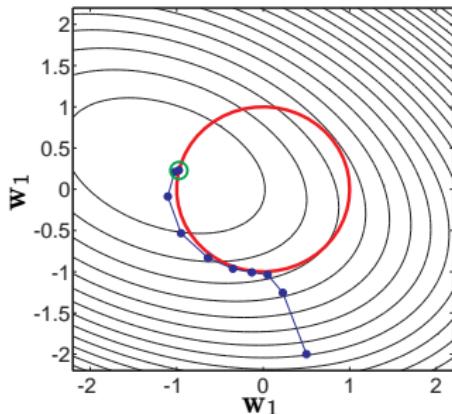
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

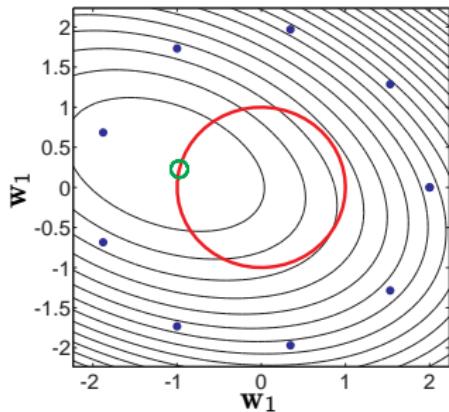
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

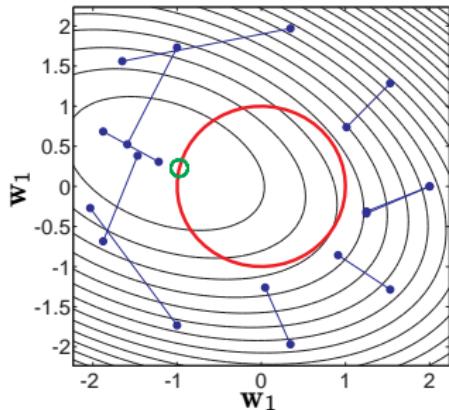
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

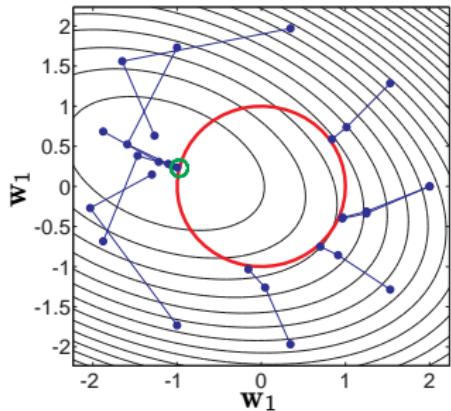
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

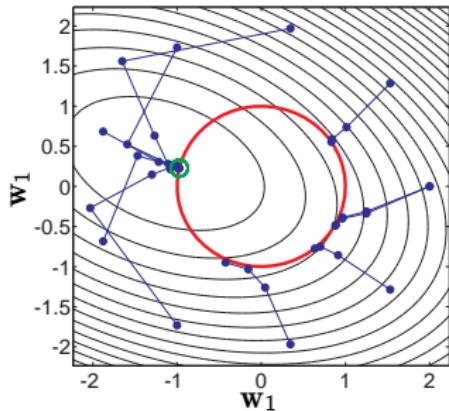
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

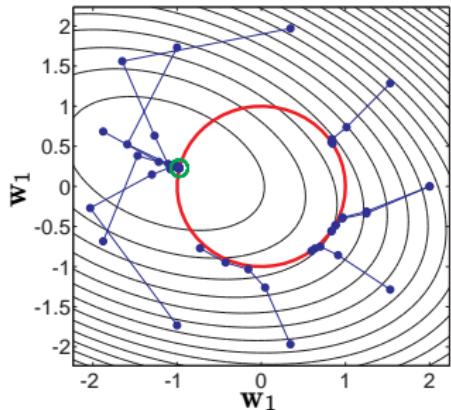
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

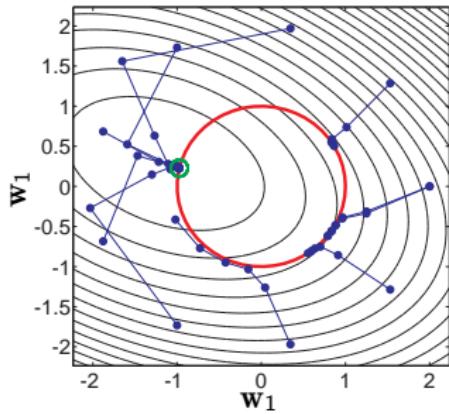
$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

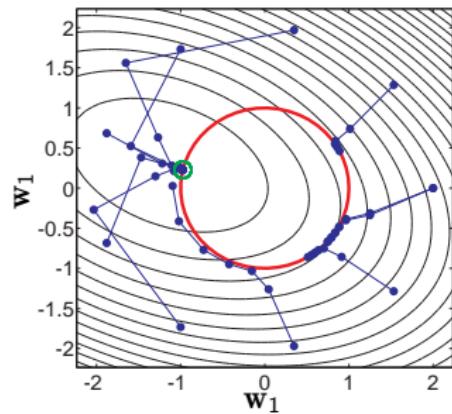
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

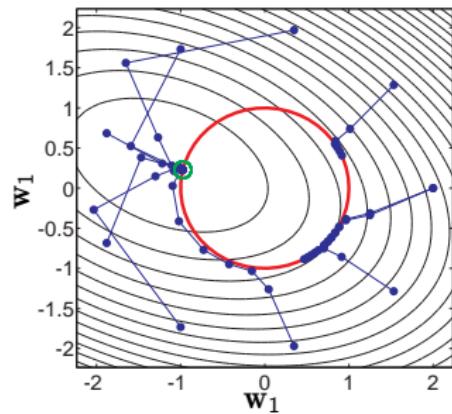
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

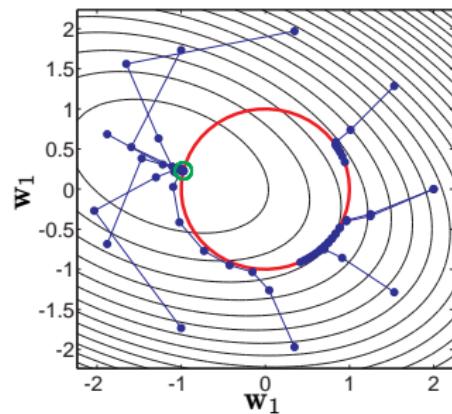
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

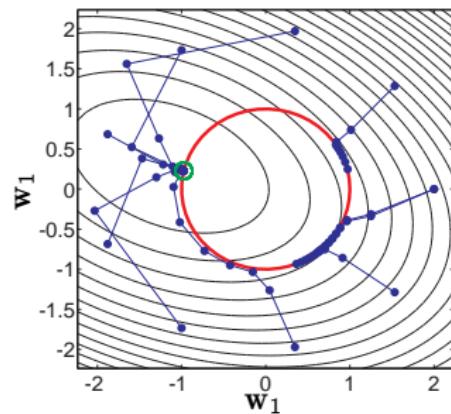
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

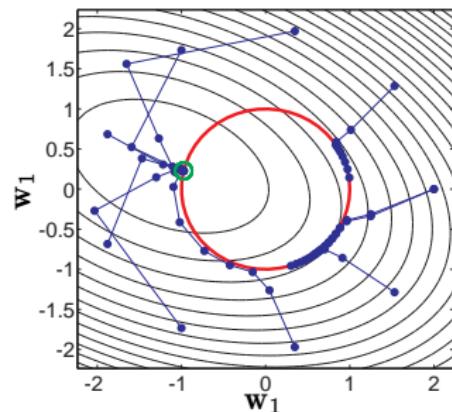
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

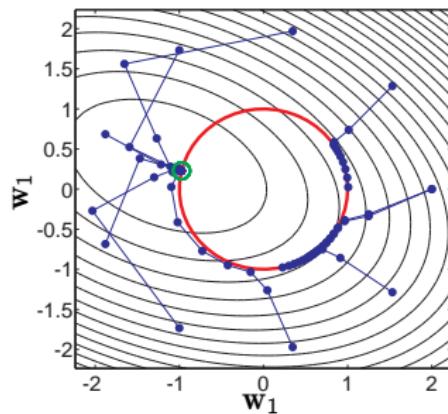
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

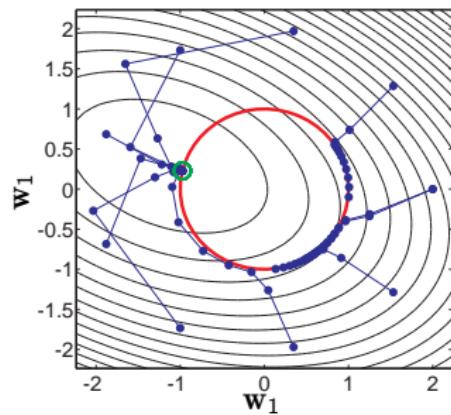
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

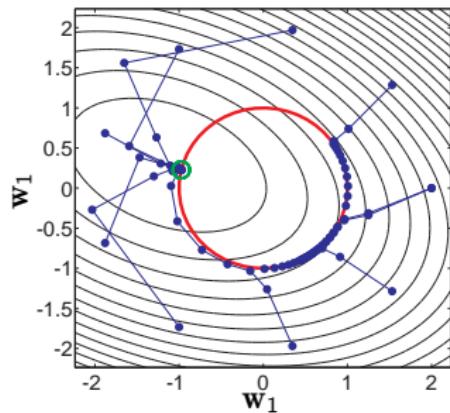
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

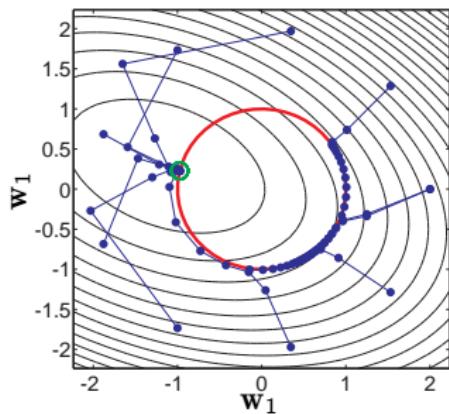
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

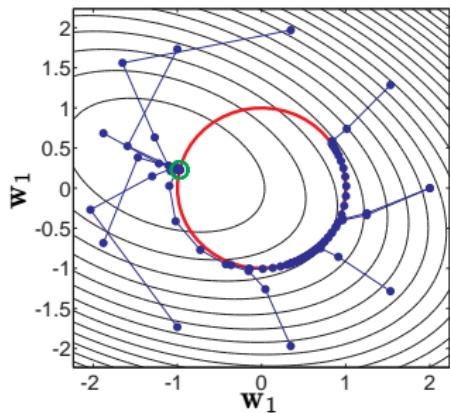
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

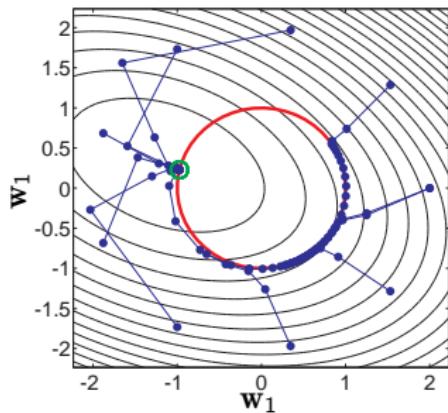
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

---

**Algorithm:** Newton method

---

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

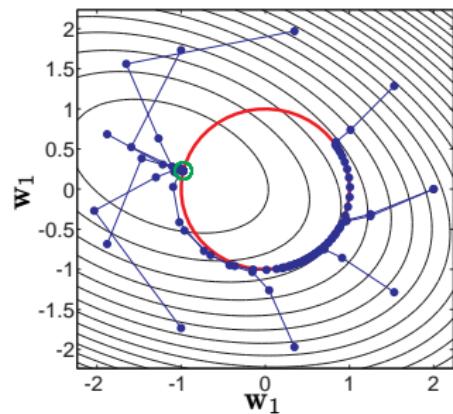
---

**return**  $w, \lambda$

---

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



# Newton Iteration for Optimization - Example

**Algorithm:** Newton method

**Input:** guess  $w$ ,  $\lambda$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|g\| \geq \text{tol}$  **do**

    Compute

$$H(w, \lambda), \nabla g(w), \nabla \Phi(w), g(w)$$

    Compute **Newton direction**

$$\begin{bmatrix} H & \nabla g \\ \nabla g^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi \\ g \end{bmatrix}$$

$$\Delta \lambda = \lambda^+ - \lambda$$

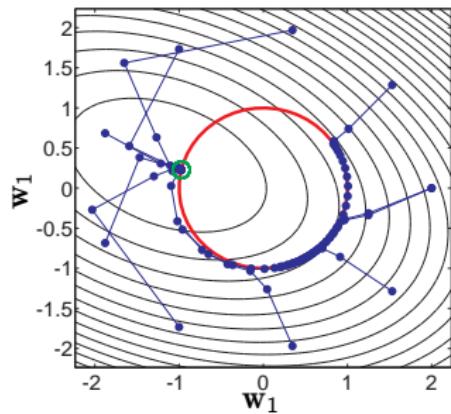
    Compute Newton step,  $t \in ]0, 1]$

$$w \leftarrow w + t \Delta w, \quad \lambda \leftarrow \lambda + t \Delta \lambda$$

**return**  $w, \lambda$

$$\begin{aligned} \min_w \quad & \frac{1}{2} w^T \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} w + w^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{s.t. } & g(w) = w^T w - 1 = 0 \end{aligned}$$

Guess  $\lambda = 0$ , step  $t = 1$



Your initial guess matters !!

# Invertibility of the KKT matrix

## The Newton direction of the KKT conditions

$$\underbrace{\begin{bmatrix} H(\mathbf{w}, \boldsymbol{\lambda}) & \nabla g(\mathbf{w}) \\ \nabla g(\mathbf{w})^\top & 0 \end{bmatrix}}_{\text{KKT matrix (symmetric indefinite)}} \begin{bmatrix} \Delta \mathbf{w} \\ \boldsymbol{\lambda}^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(\mathbf{w}) \\ g(\mathbf{w}) \end{bmatrix}$$

where  $H(\mathbf{w}, \boldsymbol{\lambda}) = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda})$

# Invertibility of the KKT matrix

## The Newton direction of the KKT conditions

$$\underbrace{\begin{bmatrix} H(\mathbf{w}, \boldsymbol{\lambda}) & \nabla g(\mathbf{w}) \\ \nabla g(\mathbf{w})^\top & 0 \end{bmatrix}}_{\text{KKT matrix (symmetric indefinite)}} \begin{bmatrix} \Delta \mathbf{w} \\ \boldsymbol{\lambda}^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(\mathbf{w}) \\ g(\mathbf{w}) \end{bmatrix}$$

where  $H(\mathbf{w}, \boldsymbol{\lambda}) = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda})$

**The KKT matrix is invertible if (sufficient, not necessary)**

- $\nabla g(\mathbf{w})$  is full row rank (LICQ)
- $\forall \mathbf{d} \neq 0$ , such that  $\nabla g(\mathbf{w})^\top \mathbf{d} = 0$

$$\mathbf{d}^\top H(\mathbf{w}, \boldsymbol{\lambda}) \mathbf{d} > 0, \quad (\text{SOSC})$$

# Invertibility of the KKT matrix

## The Newton direction of the KKT conditions

$$\underbrace{\begin{bmatrix} H(\mathbf{w}, \boldsymbol{\lambda}) & \nabla g(\mathbf{w}) \\ \nabla g(\mathbf{w})^\top & 0 \end{bmatrix}}_{\text{KKT matrix (symmetric indefinite)}} \begin{bmatrix} \Delta \mathbf{w} \\ \boldsymbol{\lambda}^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(\mathbf{w}) \\ g(\mathbf{w}) \end{bmatrix}$$

where  $H(\mathbf{w}, \boldsymbol{\lambda}) = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda})$

**The KKT matrix is invertible if (sufficient, not necessary)**

- $\nabla g(\mathbf{w})$  is full row rank (LICQ)
- $\forall \mathbf{d} \neq 0$ , such that  $\nabla g(\mathbf{w})^\top \mathbf{d} = 0$

$$\mathbf{d}^\top H(\mathbf{w}, \boldsymbol{\lambda}) \mathbf{d} > 0, \quad (\text{SOSC})$$

If  $(\mathbf{w}, \boldsymbol{\lambda})$  is LICQ & SOSC, then the KKT matrix  
is invertible in a neighborhood of  $(\mathbf{w}, \boldsymbol{\lambda})$

# Invertibility of the KKT matrix

## The Newton direction of the KKT conditions

$$\underbrace{\begin{bmatrix} H(\mathbf{w}, \boldsymbol{\lambda}) & \nabla g(\mathbf{w}) \\ \nabla g(\mathbf{w})^\top & 0 \end{bmatrix}}_{\text{KKT matrix (symmetric indefinite)}} \begin{bmatrix} \Delta \mathbf{w} \\ \boldsymbol{\lambda}^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(\mathbf{w}) \\ g(\mathbf{w}) \end{bmatrix}$$

where  $H(\mathbf{w}, \boldsymbol{\lambda}) = \nabla_{\mathbf{w}}^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda})$

**The KKT matrix is invertible if (sufficient, not necessary)**

- $\nabla g(\mathbf{w})$  is full row rank (LICQ)
- $\forall \mathbf{d} \neq 0$ , such that  $\nabla g(\mathbf{w})^\top \mathbf{d} = 0$

$$\mathbf{d}^\top H(\mathbf{w}, \boldsymbol{\lambda}) \mathbf{d} > 0, \quad (\text{SOSC})$$

If  $(\mathbf{w}, \boldsymbol{\lambda})$  is LICQ & SOSC, then the KKT matrix  
is invertible in a neighborhood of  $(\mathbf{w}, \boldsymbol{\lambda})$

**Punchline: if you don't have LICQ/SOSC at  $\mathbf{w}^*$ , then the Newton iteration becomes (in general) ill-defined as you approach  $\mathbf{w}^*$**

## Inertia of the KKT matrix

**Matrix Inertia:** # of positive, zero, negative eigenvalue ( $n_+, n_0, n_-$ )

## Inertia of the KKT matrix

**Matrix Inertia:** # of positive, zero, negative eigenvalue ( $n_+, n_0, n_-$ )

**NLP:** with  $w \in \mathbb{R}^n$  and  $g \in \mathbb{R}^m$

$$\begin{aligned} & \underset{w}{\min} \quad \Phi(w) \\ & \text{s.t.} \quad g(w) = 0 \end{aligned}$$

## Inertia of the KKT matrix

**Matrix Inertia:** # of positive, zero, negative eigenvalue ( $n_+, n_0, n_-$ )

**NLP:** with  $w \in \mathbb{R}^n$  and  $g \in \mathbb{R}^m$

$$\begin{aligned} & \underset{w}{\min} \quad \Phi(w) \\ & \text{s.t.} \quad g(w) = 0 \end{aligned}$$

yields **KKT matrix**:

$$M = \begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix}$$

## Inertia of the KKT matrix

**Matrix Inertia:** # of positive, zero, negative eigenvalue ( $n_+, n_0, n_-$ )

**NLP:** with  $w \in \mathbb{R}^n$  and  $g \in \mathbb{R}^m$

$$\begin{array}{ll}\min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0\end{array}$$

yields **KKT matrix**:

$$M = \begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix}$$

Note that:

- $\mathcal{L}(w, \lambda) = \Phi(w) + \lambda^T g(w)$

## Inertia of the KKT matrix

**Matrix Inertia:** # of positive, zero, negative eigenvalue ( $n_+, n_0, n_-$ )

**NLP:** with  $w \in \mathbb{R}^n$  and  $g \in \mathbb{R}^m$

$$\begin{array}{ll}\min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0\end{array}$$

yields **KKT matrix**:

$$M = \begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix}$$

Note that:

- $\mathcal{L}(w, \lambda) = \Phi(w) + \lambda^T g(w)$
- $\nabla_{w\lambda} \mathcal{L}(w, \lambda) = \nabla g(w)$

## Inertia of the KKT matrix

**Matrix Inertia:** # of positive, zero, negative eigenvalue ( $n_+, n_0, n_-$ )

**NLP:** with  $w \in \mathbb{R}^n$  and  $g \in \mathbb{R}^m$

$$\begin{aligned} & \underset{w}{\min} \quad \Phi(w) \\ & \text{s.t.} \quad g(w) = 0 \end{aligned}$$

yields **KKT matrix**:

$$M = \begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix}$$

Note that:

- $\mathcal{L}(w, \lambda) = \Phi(w) + \lambda^T g(w)$
- $\nabla_{w\lambda} \mathcal{L}(w, \lambda) = \nabla g(w)$
- Hence

$$\nabla \mathcal{L}(w, \lambda) = \begin{bmatrix} \nabla_{ww}^2 \mathcal{L} & \nabla_{w\lambda} \mathcal{L} \\ \nabla_{\lambda w} \mathcal{L} & \nabla_{\lambda\lambda} \mathcal{L} \end{bmatrix} = M$$

## Inertia of the KKT matrix

**Matrix Inertia:** # of positive, zero, negative eigenvalue ( $n_+, n_0, n_-$ )

**NLP:** with  $w \in \mathbb{R}^n$  and  $g \in \mathbb{R}^m$

$$\begin{aligned} \min_w \quad & \Phi(w) \\ \text{s.t.} \quad & g(w) = 0 \end{aligned}$$

**Theorem:** if SOSC and LICQ hold, then the inertia of  $M$  is  $(n, 0, m)$

yields **KKT matrix**:

$$M = \begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix}$$

Note that:

- $\mathcal{L}(w, \lambda) = \Phi(w) + \lambda^\top g(w)$
- $\nabla_{w\lambda}\mathcal{L}(w, \lambda) = \nabla g(w)$
- Hence

$$\nabla\mathcal{L}(w, \lambda) = \begin{bmatrix} \nabla_{ww}^2\mathcal{L} & \nabla_{w\lambda}\mathcal{L} \\ \nabla_{\lambda w}\mathcal{L} & \nabla_{\lambda\lambda}\mathcal{L} \end{bmatrix} = M$$

## Inertia of the KKT matrix

**Matrix Inertia:** # of positive, zero, negative eigenvalue ( $n_+, n_0, n_-$ )

**NLP:** with  $w \in \mathbb{R}^n$  and  $g \in \mathbb{R}^m$

$$\begin{aligned} \min_w \quad & \Phi(w) \\ \text{s.t.} \quad & g(w) = 0 \end{aligned}$$

yields **KKT matrix**:

$$M = \begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix}$$

Note that:

- $\mathcal{L}(w, \lambda) = \Phi(w) + \lambda^T g(w)$
- $\nabla_{w\lambda} \mathcal{L}(w, \lambda) = \nabla g(w)$
- Hence

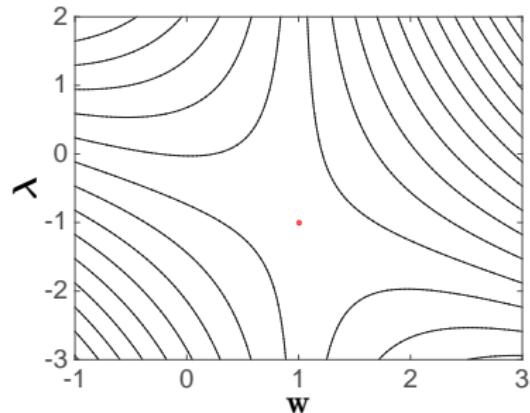
$$\nabla \mathcal{L}(w, \lambda) = \begin{bmatrix} \nabla_{ww}^2 \mathcal{L} & \nabla_{w\lambda} \mathcal{L} \\ \nabla_{\lambda w} \mathcal{L} & \nabla_{\lambda\lambda} \mathcal{L} \end{bmatrix} = M$$

**Theorem:** if SOSC and LICQ hold, then the inertia of  $M$  is  $(n, 0, m)$

### Why ?

Solution is a saddle-point of  $\mathcal{L}(w, \lambda)$ .

- Pos. eig. match  $n$  primal directions
- Neg. eig. match  $m$  dual directions



## Quadratic model interpretation

Problem:

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

The **Newton direction** is given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

## Quadratic model interpretation

Problem:

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

The **Newton direction** is given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

The **Newton direction** is also given by the Quadratic Program (QP):

$$\begin{array}{ll} \min_{\Delta w} & \frac{1}{2} \Delta w^T H(w, \lambda) \Delta w + \nabla \Phi(w)^T \Delta w \\ \text{s.t.} & g(w) + \nabla g(w)^T \Delta w = 0 \end{array}$$

## Quadratic model interpretation

Problem:

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

The **Newton direction** is given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

The **Newton direction** is also given by the Quadratic Program (QP):

$$\begin{array}{ll} \min_{\Delta w} & \frac{1}{2} \Delta w^T H(w, \lambda) \Delta w + \nabla \Phi(w)^T \Delta w \\ \text{s.t.} & g(w) + \nabla g(w)^T \Delta w = 0 \end{array}$$

Dual variables  $\lambda^+$  given by the dual variables of the QP, i.e.  $\lambda^+ = \lambda_{QP}$

## Quadratic model interpretation

Problem:

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

The **Newton direction** is given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

The **Newton direction** is also given by the Quadratic Program (QP):

$$\begin{array}{ll} \min_{\Delta w} & \frac{1}{2} \Delta w^T H(w, \lambda) \Delta w + \nabla \Phi(w)^T \Delta w \\ \text{s.t.} & g(w) + \nabla g(w)^T \Delta w = 0 \end{array}$$

Dual variables  $\lambda^+$  given by the dual variables of the QP, i.e.  $\lambda^+ = \lambda_{QP}$

*Proof: KKT of the QP are equivalent to the KKT system.*

## Quadratic model interpretation

Problem:

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

The **Newton direction** is given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

The **Newton direction** is also given by the Quadratic Program (QP):

$$\begin{array}{ll} \min_{\Delta w} & \frac{1}{2} \Delta w^T H(w, \lambda) \Delta w + \nabla \Phi(w)^T \Delta w \\ \text{s.t.} & g(w) + \nabla g(w)^T \Delta w = 0 \end{array}$$

Dual variables  $\lambda^+$  given by the dual variables of the QP, i.e.  $\lambda^+ = \lambda_{QP}$

*Proof: KKT of the QP are equivalent to the KKT system.*

The Newton direction is given by solving a quadratic models of the original problem !!

# Outline

- 1 The Newton method
- 2 Newton on the KKT conditions
- 3 The reduced Newton step (unconstrained problems)
- 4 The merit function - Line-search for constrained problems
- 5 Newton-type methods
- 6 Sequential Quadratic Programming

## The Newton direction is a descent direction

Consider the unconstrained problem

$$\min_w \Phi(w)$$

The Newton direction  $\Delta w$  minimizes the quadratic model (unconstrained):

$$Q(w, \Delta w) = \Phi(w) + \nabla \Phi(w)^T \Delta w + \frac{1}{2} \Delta w^T H(w) \Delta w$$

## The Newton direction is a descent direction

Consider the unconstrained problem

$$\min_w \Phi(w)$$

The Newton direction  $\Delta w$  minimizes the quadratic model (unconstrained):

$$Q(w, \Delta w) = \Phi(w) + \nabla \Phi(w)^T \Delta w + \frac{1}{2} \Delta w^T H(w) \Delta w$$

where  $H(w) = \nabla^2 \Phi(w)$ . Step  $\Delta w$  reads as:

$$\Delta w = -H(w)^{-1} \nabla \Phi(w)$$

## The Newton direction is a descent direction

Consider the unconstrained problem

$$\min_w \Phi(w)$$

The Newton direction  $\Delta w$  minimizes the quadratic model (unconstrained):

$$Q(w, \Delta w) = \Phi(w) + \nabla \Phi(w)^T \Delta w + \frac{1}{2} \Delta w^T H(w) \Delta w$$

where  $H(w) = \nabla_w^2 \Phi(w)$ . Step  $\Delta w$  reads as:

$$\Delta w = -H(w)^{-1} \nabla \Phi(w)$$

If the Hessian  $H(w) \succ 0$  then for  $t \rightarrow 0_+$ :

$$\Phi(w + t\Delta w) - \Phi(w) = t \nabla \Phi(w)^T \Delta w = -t \nabla \Phi(w)^T H(w)^{-1} \nabla \Phi(w) < 0$$

## The Newton direction is a descent direction

Consider the unconstrained problem

$$\min_w \Phi(w)$$

The Newton direction  $\Delta w$  minimizes the quadratic model (unconstrained):

$$Q(w, \Delta w) = \Phi(w) + \nabla \Phi(w)^T \Delta w + \frac{1}{2} \Delta w^T H(w) \Delta w$$

where  $H(w) = \nabla_w^2 \Phi(w)$ . Step  $\Delta w$  reads as:

$$\Delta w = -H(w)^{-1} \nabla \Phi(w)$$

If the Hessian  $H(w) \succ 0$  then for  $t \rightarrow 0_+$ :

$$\Phi(w + t\Delta w) - \Phi(w) = t \nabla \Phi(w)^T \Delta w = -t \nabla \Phi(w)^T H(w)^{-1} \nabla \Phi(w) < 0$$

i.e.

$$\Phi(w + t\Delta w) < \Phi(w) \text{ for } t \in \mathbb{R}_+ \text{ sufficiently small.}$$

## The Newton direction is a descent direction

Consider the unconstrained problem

$$\min_w \Phi(w)$$

The Newton direction  $\Delta w$  minimizes the quadratic model (unconstrained):

$$Q(w, \Delta w) = \Phi(w) + \nabla \Phi(w)^T \Delta w + \frac{1}{2} \Delta w^T H(w) \Delta w$$

where  $H(w) = \nabla_w^2 \Phi(w)$ . Step  $\Delta w$  reads as:

$$\Delta w = -H(w)^{-1} \nabla \Phi(w)$$

If the Hessian  $H(w) \succ 0$  then for  $t \rightarrow 0_+$ :

$$\Phi(w + t\Delta w) - \Phi(w) = t \nabla \Phi(w)^T \Delta w = -t \nabla \Phi(w)^T H(w)^{-1} \nabla \Phi(w) < 0$$

i.e.

$$\boxed{\Phi(w + t\Delta w) < \Phi(w) \text{ for } t \in \mathbb{R}_+ \text{ sufficiently small.}}$$

However, there's no guarantee that  $\Phi(w + \Delta w) < \Phi(w) !!$

## Failure of the full Newton step

Newton direction  $\Delta w$  minimizes the quadratic model :

$$Q(w, \Delta w) = \Phi(w) + \nabla \Phi(w)^T \Delta w + \frac{1}{2} \Delta w^T H(w) \Delta w$$

Given by:

$$\Delta w = -H(w)^{-1} \nabla \Phi(w)$$

Full Newton step:

$$w \leftarrow w + \Delta w$$

## Failure of the full Newton step

Newton direction  $\Delta w$  minimizes the quadratic model :

$$Q(w, \Delta w) = \Phi(w) + \nabla \Phi(w)^T \Delta w + \frac{1}{2} \Delta w^T H(w) \Delta w$$

Given by:

$$\Delta w = -H(w)^{-1} \nabla \Phi(w)$$

Full Newton step:

$$w \leftarrow w + \Delta w$$

What if the quadratic model is "not good" ?

## Failure of the full Newton step

Newton direction  $\Delta w$  minimizes the quadratic model :

$$Q(w, \Delta w) = \Phi(w) + \nabla \Phi(w)^T \Delta w + \frac{1}{2} \Delta w^T H(w) \Delta w$$

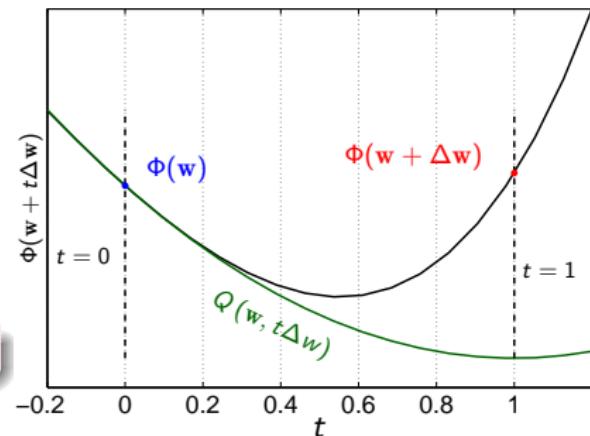
Given by:

$$\Delta w = -H(w)^{-1} \nabla \Phi(w)$$

Full Newton step:

$$w \leftarrow w + \Delta w$$

What if the quadratic model is "not good" ?



## Failure of the full Newton step

Newton direction  $\Delta w$  minimizes the quadratic model :

$$Q(w, \Delta w) = \Phi(w) + \nabla \Phi(w)^T \Delta w + \frac{1}{2} \Delta w^T H(w) \Delta w$$

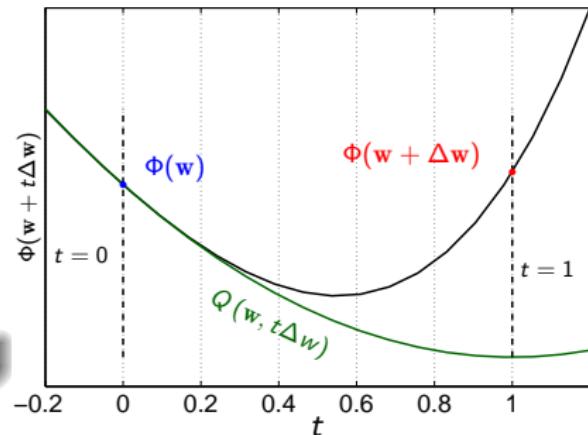
Given by:

$$\Delta w = -H(w)^{-1} \nabla \Phi(w)$$

Full Newton step:

$$w \leftarrow w + \Delta w$$

What if the quadratic model is "not good" ?



A situation with  $\Phi(w + \Delta w) > \Phi(w)$  can easily occur...

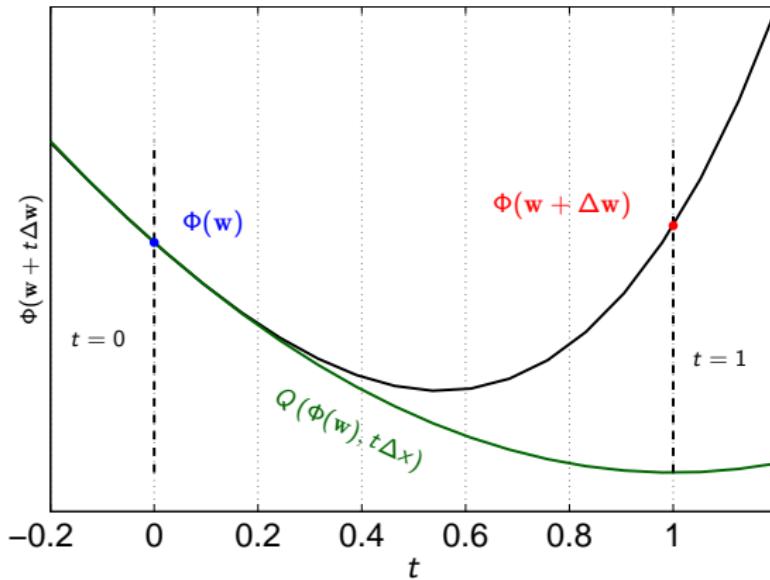
- Strong variation of  $H$
- Nonlinear constraints (if the problem is constrained)

# Globalization - Line search strategies

## Exact line search (for unconstrained optimization)

Find the best step length:

$$t = \arg \min_{t \in [0,1]} \Phi(\mathbf{w} + t\Delta\mathbf{w})$$



## Recap of last lecture

- Newton is a tool to solve nonlinear equations, will be used to solve the KKTs associated to an NLP
- Reduced-step Newton converges to a solution if  $\nabla r$  is full rank
- Convergence of full steps Newton is quadratic (fast), better when equations are “closer to linear”
- Rescaling the decision variables  $w$  has no effect on Newton, may help the linear algebra though
- Newton on KKT conditions has a specific structure, step well posed if LICQ & SOSC
- Newton step is equivalent to solving a QP forming a “linear-quadratic” approximation of the NLP
- In unconstrained optimization problems, this observation is equivalent to “minimizing a quadratic model of the cost”, reduced steps are needed when that model validity is “locally limited”

# Globalization - Line search strategies

## "Armijo's" backtracking line search (for unconstrained optimization)

Given a primal direction  $\Delta w$ , using  $0 < \alpha \leq \frac{1}{2}$  and  $0 < \beta < 1$ , do  $t = 1$ :

---

**Algorithm:** "Armijo's" backtracking

---

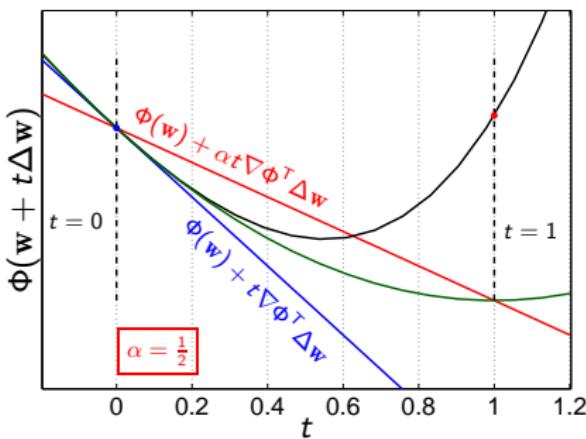
$t \leftarrow 1$

**while**  $\Phi(w + t\Delta w) \geq \Phi(w) + \alpha t \nabla \Phi(w)^T \Delta w$  **do**

  └  $t \leftarrow \beta t$

**return**  $t$

---



# Globalization - Line search strategies

## "Armijo's" backtracking line search (for unconstrained optimization)

Given a primal direction  $\Delta w$ , using  $0 < \alpha \leq \frac{1}{2}$  and  $0 < \beta < 1$ , do  $t = 1$ :

---

### Algorithm: "Armijo's" backtracking

---

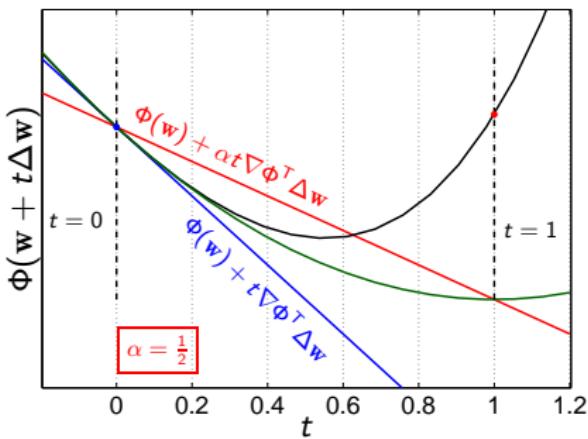
$t \leftarrow 1$

**while**  $\Phi(w + t\Delta w) \geq \Phi(w) + \alpha t \nabla \Phi(w)^T \Delta w$  **do**

  └  $t \leftarrow \beta t$

**return**  $t$

---



- If  $\alpha$  too small we may accept steps yielding only mediocre improvement.
- If  $\Phi$  is quadratic, we want full step, i.e.  $\alpha \leq \frac{1}{2}$
- In practice  $\alpha < \frac{1}{2}$  is usually preferred
- Additional conditions are sometimes used to yield a sufficiently large the step (c.f. Wolfe or Goldstein conditions)

## Convergence of the Newton method with Line-search (unconstrained)

### Theorem

Assume:

- Hessian is bounded  $m \cdot I \leq \nabla^2 \Phi(\mathbf{w}) \leq M \cdot I$ , with  $m, M > 0$
- Hessian is Lipschitz, i.e.  $\|\nabla^2 \Phi(\mathbf{w}) - \nabla^2 \Phi(\mathbf{y})\| \leq L \|\mathbf{w} - \mathbf{y}\|$ ,  $\forall \mathbf{w}, \mathbf{y}$

## Convergence of the Newton method with Line-search (unconstrained)

### Theorem

Assume:

- Hessian is bounded  $m \cdot I \leq \nabla^2 \Phi(\mathbf{w}) \leq M \cdot I$ , with  $m, M > 0$
- Hessian is Lipschitz, i.e.  $\|\nabla^2 \Phi(\mathbf{w}) - \nabla^2 \Phi(\mathbf{y})\| \leq L \|\mathbf{w} - \mathbf{y}\|, \quad \forall \mathbf{w}, \mathbf{y}$

then  $\exists \eta, \gamma > 0$  with  $\eta < \frac{m^2}{L}$  such that  $\forall \mathbf{w}$ :

Damped phase:

$$\Phi(\mathbf{w}^+) - \Phi(\mathbf{w}) \leq -\gamma \quad \text{if } \|\nabla \Phi(\mathbf{w})\| \geq \eta$$

Quadratic phase:

$$\|\nabla \Phi(\mathbf{w}^+)\| \leq \frac{L}{2m^2} \|\nabla \Phi(\mathbf{w})\|^2 \quad \text{if } \|\nabla \Phi(\mathbf{w})\| < \eta$$

where  $\mathbf{w}^+ = \mathbf{w} + t\Delta\mathbf{w}$  is the (possibly) reduced Newton step

## Convergence of the Newton method with Line-search (unconstrained)

### Theorem

Assume:

- Hessian is bounded  $m \cdot I \leq \nabla^2 \Phi(\mathbf{w}) \leq M \cdot I$ , with  $m, M > 0$
- Hessian is Lipschitz, i.e.  $\|\nabla^2 \Phi(\mathbf{w}) - \nabla^2 \Phi(\mathbf{y})\| \leq L \|\mathbf{w} - \mathbf{y}\|, \quad \forall \mathbf{w}, \mathbf{y}$

then  $\exists \eta, \gamma > 0$  with  $\eta < \frac{m^2}{L}$  such that  $\forall \mathbf{w}$ :

Damped phase:

$$\Phi(\mathbf{w}^+) - \Phi(\mathbf{w}) \leq -\gamma \quad \text{if } \|\nabla \Phi(\mathbf{w})\| \geq \eta$$

Quadratic phase:

$$\|\nabla \Phi(\mathbf{w}^+)\| \leq \frac{L}{2m^2} \|\nabla \Phi(\mathbf{w})\|^2 \quad \text{if } \|\nabla \Phi(\mathbf{w})\| < \eta$$

where  $\mathbf{w}^+ = \mathbf{w} + t\Delta\mathbf{w}$  is the (possibly) reduced Newton step

**Practical consequence:** two-phase convergence

- If  $\mathbf{w}$  is *far* from  $\mathbf{w}^*$   $\Rightarrow$  Damped convergence (reduced steps)
- If  $\mathbf{w}$  is *close* to  $\mathbf{w}^*$   $\Rightarrow$  Quadratic convergence (full steps)
- Once Newton has entered the quadratic phase, it stays quadratic !!

# Convergence of the Newton method with Line-search (unconstrained)

## Theorem

Assume:

- Hessian is bounded  $m \cdot I \leq \nabla^2 \Phi(\mathbf{w}) \leq M \cdot I$ , with  $m, M > 0$
- Hessian is Lipschitz, i.e.  $\|\nabla^2 \Phi(\mathbf{w}) - \nabla^2 \Phi(\mathbf{y})\| \leq L \|\mathbf{w} - \mathbf{y}\|, \quad \forall \mathbf{w}, \mathbf{y}$

then  $\exists \eta, \gamma > 0$  with  $\eta < \frac{m^2}{L}$  such that  $\forall \mathbf{w}$ :

Damped phase:

$$\Phi(\mathbf{w}^+) - \Phi(\mathbf{w}) \leq -\gamma \quad \text{if } \|\nabla \Phi(\mathbf{w})\| \geq \eta$$

Quadratic phase:

$$\|\nabla \Phi(\mathbf{w}^+)\| \leq \frac{L}{2m^2} \|\nabla \Phi(\mathbf{w})\|^2 \quad \text{if } \|\nabla \Phi(\mathbf{w})\| < \eta$$

where  $\mathbf{w}^+ = \mathbf{w} + t\Delta\mathbf{w}$  is the (possibly) reduced Newton step

**Practical consequence:** region of quadratic convergence

- If  $L$  is large, then
  - ① the quadratic region ( $\|\nabla \Phi(\mathbf{w})\| < \eta$ ) is small
  - ② the quadratic contraction rate is small

## Convergence of the Newton method with Line-search (unconstrained)

### Theorem

Assume:

- Hessian is bounded  $m \cdot I \leq \nabla^2 \Phi(\mathbf{w}) \leq M \cdot I$ , with  $m, M > 0$
- Hessian is Lipschitz, i.e.  $\|\nabla^2 \Phi(\mathbf{w}) - \nabla^2 \Phi(\mathbf{y})\| \leq L \|\mathbf{w} - \mathbf{y}\|, \quad \forall \mathbf{w}, \mathbf{y}$

then  $\exists \eta, \gamma > 0$  with  $\eta < \frac{m^2}{L}$  such that  $\forall \mathbf{w}$ :

Damped phase:

$$\Phi(\mathbf{w}^+) - \Phi(\mathbf{w}) \leq -\gamma \quad \text{if } \|\nabla \Phi(\mathbf{w})\| \geq \eta$$

Quadratic phase:

$$\|\nabla \Phi(\mathbf{w}^+)\| \leq \frac{L}{2m^2} \|\nabla \Phi(\mathbf{w})\|^2 \quad \text{if } \|\nabla \Phi(\mathbf{w})\| < \eta$$

where  $\mathbf{w}^+ = \mathbf{w} + t\Delta\mathbf{w}$  is the (possibly) reduced Newton step

**What if we have constraints ?** Similar results, with  $\nabla_{\mathbf{w}}^2 \mathcal{L}$ ,  $\nabla_{\mathbf{w}} \mathbf{g}$ . However, Newton can fail on an infeasible point (no direction improving the constraints, c.f. "But still, Newton can fail...")

# Outline

- 1 The Newton method
- 2 Newton on the KKT conditions
- 3 The reduced Newton step (unconstrained problems)
- 4 The merit function - Line-search for constrained problems
- 5 Newton-type methods
- 6 Sequential Quadratic Programming

## Assessing convergence with constraints

Solve the NLP...

$$\min_w \Phi(w)$$

$$\text{s.t. } g(w) = 0$$

## Assessing convergence with constraints

Solve the NLP...                    ... by satisfying the KKT conditions...

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

$$r(w, \lambda) = \begin{bmatrix} \nabla_w \mathcal{L}(w, \lambda) \\ g(w) \end{bmatrix} = 0$$

## Assessing convergence with constraints

Solve the NLP...      ... by satisfying the KKT conditions...

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

$$r(w, \lambda) = \begin{bmatrix} \nabla_w \mathcal{L}(w, \lambda) \\ g(w) \end{bmatrix} = 0$$

... via taking Newton steps along the Newton directions.

## Assessing convergence with constraints

Solve the NLP...      ... by satisfying the KKT conditions...

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

$$r(w, \lambda) = \begin{bmatrix} \nabla_w \mathcal{L}(w, \lambda) \\ g(w) \end{bmatrix} = 0$$

... via taking Newton steps along the Newton directions.

The Newton directions are given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

and  $\Delta \lambda = \lambda^+ - \lambda$

## Assessing convergence with constraints

Solve the NLP...      ... by satisfying the KKT conditions...

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

$$r(w, \lambda) = \begin{bmatrix} \nabla_w \mathcal{L}(w, \lambda) \\ g(w) \end{bmatrix} = 0$$

... via taking Newton steps along the Newton directions.

The **Newton directions** are given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

and  $\Delta \lambda = \lambda^+ - \lambda$

- We know that there exists a Newton step that improves  $\|r\|$ , i.e.

$$\|r(w + t\Delta w, \lambda + t\Delta \lambda)\| < \|r(w, \lambda)\|$$

for some  $t \in ]0, 1]$  (as long as  $\nabla r$  is non-singular). How to select  $t$ ?

## Assessing convergence with constraints

Solve the NLP...      ... by satisfying the KKT conditions...

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

$$r(w, \lambda) = \begin{bmatrix} \nabla_w \mathcal{L}(w, \lambda) \\ g(w) \end{bmatrix} = 0$$

... via taking Newton steps along the Newton directions.

The Newton directions are given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

and  $\Delta \lambda = \lambda^+ - \lambda$

- We know that there exists a Newton step that improves  $\|r\|$ , i.e.

$$\|r(w + t\Delta w, \lambda + t\Delta \lambda)\| < \|r(w, \lambda)\|$$

for some  $t \in ]0, 1]$  (as long as  $\nabla r$  is non-singular). How to select  $t$ ?

- Could ensure the progress of  $\|r(w + t\Delta w, \lambda + t\Delta \lambda)\|$

## Assessing convergence with constraints

Solve the NLP...      ... by satisfying the KKT conditions...

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

$$r(w, \lambda) = \begin{bmatrix} \nabla_w \mathcal{L}(w, \lambda) \\ g(w) \end{bmatrix} = 0$$

... via taking Newton steps along the Newton directions.

The Newton directions are given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

and  $\Delta \lambda = \lambda^+ - \lambda$

- We know that there exists a Newton step that improves  $\|r\|$ , i.e.

$$\|r(w + t\Delta w, \lambda + t\Delta \lambda)\| < \|r(w, \lambda)\|$$

for some  $t \in ]0, 1]$  (as long as  $\nabla r$  is non-singular). How to select  $t$ ?

- Could ensure the progress of  $\|r(w + t\Delta w, \lambda + t\Delta \lambda)\|$
- Some caveats... e.g. need to evaluate  $\nabla g(w + t\Delta w)$  each time we backtrack on  $t$   
!! This can be computationally expensive...

## Line-search with Equality Constraints

$T_1$  merit function...

... accounts for progress on  $\Phi$  and  $g$

$$T_1(w) = \Phi(w) + \nu \|g(w)\|_1$$

We would like  $\min_w T_1(w)$  to coincide with the solution of the NLP !!

# Line-search with Equality Constraints

## $T_1$ merit function...

... accounts for progress on  $\Phi$  and  $g$

$$T_1(w) = \Phi(w) + \nu \|g(w)\|_1$$

We would like  $\min_w T_1(w)$  to coincide with the solution of the NLP !!

### Properties of the $T_1$ merit function:

- If  $w^*, \lambda^*$  is a regular KKT point with SOSC, then  $w^*$  is a (local) minimum of  $T_1$  for  $\nu > \|\lambda^*\|_\infty$

# Line-search with Equality Constraints

## $T_1$ merit function...

... accounts for progress on  $\Phi$  and  $g$

$$T_1(w) = \Phi(w) + \nu \|g(w)\|_1$$

We would like  $\min_w T_1(w)$  to coincide with the solution of the NLP !!

### Properties of the $T_1$ merit function:

- If  $w^*, \lambda^*$  is a regular KKT point with SOSC, then  $w^*$  is a (local) minimum of  $T_1$  for  $\nu > \|\lambda^*\|_\infty$
- If  $w$  is a minimum of  $T_1$  for  $\nu > \|\lambda^*\|_\infty$  with  $g(w) = 0$ , then  $w$  is a KKT point of the original problem. If  $g(w) \neq 0$ , then we have collapsed to an infeasible point.

# Line-search with Equality Constraints

## $T_1$ merit function...

... accounts for progress on  $\Phi$  and  $g$

$$T_1(w) = \Phi(w) + \nu \|g(w)\|_1$$

We would like  $\min_w T_1(w)$  to coincide with the solution of the NLP !!

### Properties of the $T_1$ merit function:

- If  $w^*, \lambda^*$  is a regular KKT point with SOSC, then  $w^*$  is a (local) minimum of  $T_1$  for  $\nu > \|\lambda^*\|_\infty$
- If  $w$  is a minimum of  $T_1$  for  $\nu > \|\lambda^*\|_\infty$  with  $g(w) = 0$ , then  $w$  is a KKT point of the original problem. If  $g(w) \neq 0$ , then we have collapsed to an infeasible point.

**Merit function  $T_1$  cannot be improved from  $w^*$ . Minimising  $T_1$  yields a KKT point.**

## Line-search with Equality Constraints

$T_1$  merit function...

... accounts for progress on  $\Phi$  and  $g$

$$T_1(w) = \Phi(w) + \nu \|g(w)\|_1$$

We would like  $\min_w T_1(w)$  to coincide with the solution of the NLP !!

**Properties of the  $T_1$  merit function:**

- If  $w^*, \lambda^*$  is a regular KKT point with SOSC, then  $w^*$  is a (local) minimum of  $T_1$  for  $\nu > \|\lambda^*\|_\infty$
- If  $w$  is a minimum of  $T_1$  for  $\nu > \|\lambda^*\|_\infty$  with  $g(w) = 0$ , then  $w$  is a KKT point of the original problem. If  $g(w) \neq 0$ , then we have collapsed to an infeasible point.

**Merit function  $T_1$  cannot be improved from  $w^*$ . Minimising  $T_1$  yields a KKT point.**

What about convergence of Newton with line-search on  $T_1$  ?

Similar results as in the unconstrained case if LICQ is respected !!

# Line-search with Equality Constraints

## $T_1$ merit function...

... accounts for progress on  $\Phi$  and  $g$

$$T_1(w) = \Phi(w) + \nu \|g(w)\|_1$$

We would like  $\min_w T_1(w)$  to coincide with the solution of the NLP !!

### Properties of the $T_1$ merit function:

- If  $w^*, \lambda^*$  is a regular KKT point with SOSC, then  $w^*$  is a (local) minimum of  $T_1$  for  $\nu > \|\lambda^*\|_\infty$
- If  $w$  is a minimum of  $T_1$  for  $\nu > \|\lambda^*\|_\infty$  with  $g(w) = 0$ , then  $w$  is a KKT point of the original problem. If  $g(w) \neq 0$ , then we have collapsed to an infeasible point.

**Merit function  $T_1$  cannot be improved from  $w^*$ . Minimising  $T_1$  yields a KKT point.**

### What about convergence of Newton with line-search on $T_1$ ?

Similar results as in the unconstrained case if LICQ is respected !!

Difficulty:  $\lambda^*$  is not known a priori. Parameter  $\nu$  must be chosen adequately, possibly adjusted over the iterations !

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

For  $g_i(w) \neq 0 \forall i$ , using  $\sigma = \text{sign}(g(w))$  (i.e.  $\sigma^\top g(w) = \|g(w)\|_1$ )

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

For  $g_i(w) \neq 0 \forall i$ , using  $\sigma = \text{sign}(g(w))$  (i.e.  $\sigma^\top g(w) = \|g(w)\|_1$ )

$$\frac{d}{dt} \|g(w + t\Delta w)\|_1 \Big|_{t=0} = \sigma^\top \nabla g(w)^\top \Delta w$$

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

For  $g_i(w) \neq 0 \forall i$ , using  $\sigma = \text{sign}(g(w))$  (i.e.  $\sigma^\top g(w) = \|g(w)\|_1$ )

$$\frac{d}{dt} \|g(w + t\Delta w)\|_1 \Big|_{t=0} = \sigma^\top \nabla g(w)^\top \Delta w$$

and since  $\nabla g(w)^\top \Delta w = -g(w)$ , we have:

$$\frac{d}{dt} \|g(w + t\Delta w)\|_1 \Big|_{t=0} = -\sigma^\top g(w) = -\|g(w)\|_1$$

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

We have:

$$\frac{d}{dt} \|g(w + t\Delta w)\|_{t=0} = -\sigma^\top g(w) = -\|g(w)\|_1$$

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

We have:

$$\frac{d}{dt} \|g(w + t\Delta w)\|_{t=0} = -\sigma^\top g(w) = -\|g(w)\|_1$$

Moreover,  $H\Delta w + \nabla g\lambda = -\nabla\Phi$  yields  $\Delta w^\top H\Delta w + \Delta w^\top \nabla g\lambda = -\Delta w^\top \nabla\Phi$

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

We have:

$$\frac{d}{dt} \|g(w + t\Delta w)\|_{t=0} = -\sigma^\top g(w) = -\|g(w)\|_1$$

Moreover,  $H\Delta w + \nabla g\lambda = -\nabla\Phi$  yields  $\Delta w^\top H\Delta w + \Delta w^\top \nabla g\lambda = -\Delta w^\top \nabla\Phi$ , and:

$$\frac{d}{dt} \Phi(w + t\Delta w) \Big|_{t=0} = \nabla\Phi^\top \Delta w$$

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

We have:

$$\frac{d}{dt} \|g(w + t\Delta w)\|_{t=0} = -\sigma^T g(w) = -\|g(w)\|_1$$

Moreover,  $H\Delta w + \nabla g\lambda = -\nabla\Phi$  yields  $\Delta w^T H\Delta w + \Delta w^T \nabla g\lambda = -\Delta w^T \nabla\Phi$ , and:

$$\frac{d}{dt} \Phi(w + t\Delta w) \Big|_{t=0} = \nabla\Phi^T \Delta w = -\Delta w^T H\Delta w - \lambda \nabla g^T \Delta w$$

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

We have:

$$\frac{d}{dt} \|g(w + t\Delta w)\|_{t=0} = -\sigma^\top g(w) = -\|g(w)\|_1$$

Moreover,  $H\Delta w + \nabla g\lambda = -\nabla\Phi$  yields  $\Delta w^\top H\Delta w + \Delta w^\top \nabla g\lambda = -\Delta w^\top \nabla\Phi$ , and:

$$\frac{d}{dt} \Phi(w + t\Delta w) \Big|_{t=0} = \nabla\Phi^\top \Delta w = -\Delta w^\top H\Delta w - \lambda \nabla g^\top \Delta w$$

Such that:

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} = -\Delta w^\top H\Delta w - \lambda \nabla g^\top \Delta w - \nu \|g(w)\|_1$$

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

We have:

$$\frac{d}{dt} \|g(w + t\Delta w)\|_{t=0} = -\sigma^\top g(w) = -\|g(w)\|_1$$

Moreover,  $H\Delta w + \nabla g\lambda = -\nabla\Phi$  yields  $\Delta w^\top H\Delta w + \Delta w^\top \nabla g\lambda = -\Delta w^\top \nabla\Phi$ , and:

$$\frac{d}{dt} \Phi(w + t\Delta w) \Big|_{t=0} = \nabla\Phi^\top \Delta w = -\Delta w^\top H\Delta w - \lambda \nabla g^\top \Delta w$$

Such that:

$$\begin{aligned} \frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} &= -\Delta w^\top H\Delta w - \lambda \nabla g^\top \Delta w - \nu \|g(w)\|_1 = \\ &\quad -\Delta w^\top H\Delta w + \lambda g - \nu \|g(w)\|_1 \end{aligned}$$

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

We have:

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} = -\Delta w^\top H \Delta w + \lambda g - \nu \|g(w)\|_1$$

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

We have:

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} = -\Delta w^\top H \Delta w + \lambda g - \nu \|g(w)\|_1$$

Then, since  $|\lambda g| \leq \|\lambda\|_\infty \|g\|_1$ , we have:

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

We have:

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} = -\Delta w^\top H \Delta w + \lambda g - \nu \|g(w)\|_1$$

Then, since  $|\lambda g| \leq \|\lambda\|_\infty \|g\|_1$ , we have:

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} \leq -\Delta w^\top H \Delta w + \|\lambda\|_\infty \|g\|_1 - \nu \|g(w)\|_1$$

Newton direction is a descent direction for  $T_1$

**Theorem:** for  $g_i(w) \neq 0 \forall i$ , then the Newton direction  $(w, \lambda)$

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

is a descent direction for the  $T_1$  merit function  $T_1 = \Phi(w) + \nu \|g(w)\|_1$ .

**Proof:**  $\Delta w$  is a descent direction if

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} < 0$$

We have:

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} = -\Delta w^T H \Delta w + \lambda g - \nu \|g(w)\|_1$$

Then, since  $|\lambda g| \leq \|\lambda\|_\infty \|g\|_1$ , we have:

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} \leq -\Delta w^T H \Delta w + \|\lambda\|_\infty \|g\|_1 - \nu \|g(w)\|_1$$

i.e.

$$\frac{d}{dt} T_1(w + t\Delta w) \Big|_{t=0} \leq -\Delta w^T H \Delta w + (\|\lambda\|_\infty - \nu) \|g(w)\|_1 < 0$$

holds if SOSC (reduced Hessian  $Z^T H Z \geq 0$ ) and  $\|\lambda\|_\infty > \nu$ .

# Outline

- 1 The Newton method
- 2 Newton on the KKT conditions
- 3 The reduced Newton step (unconstrained problems)
- 4 The merit function - Line-search for constrained problems
- 5 Newton-type methods
- 6 Sequential Quadratic Programming

## Convergence of Newton-type methods

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an [approximation](#) instead !!

**What happens to the convergence of the Newton method ?**

## Convergence of Newton-type methods

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an **approximation** instead !!

**What happens to the convergence of the Newton method ?**

Root finding problem  $\mathbf{r}(\mathbf{w}) = 0$  with iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - \left( \frac{\partial \mathbf{r}}{\partial \mathbf{w}} \right)^{-1} \mathbf{r}(\mathbf{w})$

## Convergence of Newton-type methods

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an **approximation** instead !!

**What happens to the convergence of the Newton method ?**

Root finding problem  $\mathbf{r}(\mathbf{w}) = 0$  with iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - \left( \frac{\partial \mathbf{r}}{\partial \mathbf{w}} \right)^{-1} \mathbf{r}(\mathbf{w})$

The **Newton-type** iteration reads as

$$\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1} \mathbf{r}(\mathbf{w}_k), \quad \text{with} \quad M_k \text{ invertible, approximates } J = \frac{\partial \mathbf{r}(\mathbf{w})}{\partial \mathbf{w}}$$

## Convergence of Newton-type methods

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an **approximation** instead !!

**What happens to the convergence of the Newton method ?**

Root finding problem  $\mathbf{r}(\mathbf{w}) = 0$  with iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - \left( \frac{\partial \mathbf{r}}{\partial \mathbf{w}} \right)^{-1} \mathbf{r}(\mathbf{w})$

The **Newton-type** iteration reads as

$$\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1} \mathbf{r}(\mathbf{w}_k), \quad \text{with} \quad M_k \text{ invertible, approximates } J = \frac{\partial \mathbf{r}(\mathbf{w})}{\partial \mathbf{w}}$$

We assume:

## Convergence of Newton-type methods

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an **approximation** instead !!

### What happens to the convergence of the Newton method ?

Root finding problem  $\mathbf{r}(\mathbf{w}) = 0$  with iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - \left( \frac{\partial \mathbf{r}}{\partial \mathbf{w}} \right)^{-1} \mathbf{r}(\mathbf{w})$

The **Newton-type** iteration reads as

$$\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1} \mathbf{r}(\mathbf{w}_k), \quad \text{with} \quad M_k \text{ invertible, approximates } J = \frac{\partial \mathbf{r}(\mathbf{w})}{\partial \mathbf{w}}$$

We assume:

- Lipschitz condition:  $\|M_k^{-1}(J(\mathbf{w}_k) - J(\mathbf{w}^*))\| \leq \omega \|\mathbf{w}_k - \mathbf{w}^*\|$

## Convergence of Newton-type methods

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an **approximation** instead !!

### What happens to the convergence of the Newton method ?

Root finding problem  $\mathbf{r}(\mathbf{w}) = 0$  with iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - \left( \frac{\partial \mathbf{r}}{\partial \mathbf{w}} \right)^{-1} \mathbf{r}(\mathbf{w})$

The **Newton-type** iteration reads as

$$\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1} \mathbf{r}(\mathbf{w}_k), \quad \text{with} \quad M_k \text{ invertible, approximates } J = \frac{\partial \mathbf{r}(\mathbf{w})}{\partial \mathbf{w}}$$

We assume:

- Lipschitz condition:  $\|M_k^{-1}(J(\mathbf{w}_k) - J(\mathbf{w}^*))\| \leq \omega \|\mathbf{w}_k - \mathbf{w}^*\|$
- Bound on the Jacobian approximation error  $\|M_k^{-1}(J(\mathbf{w}_k) - M_k)\| \leq \kappa_k < \kappa < 1$

## Convergence of Newton-type methods

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an **approximation** instead !!

### What happens to the convergence of the Newton method ?

Root finding problem  $\mathbf{r}(\mathbf{w}) = 0$  with iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - \left( \frac{\partial \mathbf{r}}{\partial \mathbf{w}} \right)^{-1} \mathbf{r}(\mathbf{w})$

The **Newton-type** iteration reads as

$$\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1} \mathbf{r}(\mathbf{w}_k), \quad \text{with} \quad M_k \text{ invertible, approximates } J = \frac{\partial \mathbf{r}(\mathbf{w})}{\partial \mathbf{w}}$$

We assume:

- Lipschitz condition:  $\|M_k^{-1}(J(\mathbf{w}_k) - J(\mathbf{w}^*))\| \leq \omega \|\mathbf{w}_k - \mathbf{w}^*\|$
- Bound on the Jacobian approximation error  $\|M_k^{-1}(J(\mathbf{w}_k) - M_k)\| \leq \kappa_k < \kappa < 1$
- Good initial guess  $\|\mathbf{w}_0 - \mathbf{w}^*\| \leq \frac{2}{\omega}(1 - \kappa)$

## Convergence of Newton-type methods

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an **approximation** instead !!

### What happens to the convergence of the Newton method ?

Root finding problem  $\mathbf{r}(\mathbf{w}) = 0$  with iteration  $\mathbf{w}_{k+1} = \mathbf{w}_k - \left( \frac{\partial \mathbf{r}}{\partial \mathbf{w}} \right)^{-1} \mathbf{r}(\mathbf{w})$

The **Newton-type** iteration reads as

$$\mathbf{w}_{k+1} = \mathbf{w}_k - M_k^{-1} \mathbf{r}(\mathbf{w}_k), \quad \text{with} \quad M_k \text{ invertible, approximates } J = \frac{\partial \mathbf{r}(\mathbf{w})}{\partial \mathbf{w}}$$

We assume:

- Lipschitz condition:  $\|M_k^{-1}(J(\mathbf{w}_k) - J(\mathbf{w}^*))\| \leq \omega \|\mathbf{w}_k - \mathbf{w}^*\|$
- Bound on the Jacobian approximation error  $\|M_k^{-1}(J(\mathbf{w}_k) - M_k)\| \leq \kappa_k < \kappa < 1$
- Good initial guess  $\|\mathbf{w}_0 - \mathbf{w}^*\| \leq \frac{2}{\omega}(1 - \kappa)$

Then  $\mathbf{w}_k \rightarrow \mathbf{w}^*$  with the following linear contraction in each iteration:

$$\|\mathbf{w}_{k+1} - \mathbf{w}^*\| \leq \left( \kappa_k + \frac{\omega}{2} \|\mathbf{w}_k - \mathbf{w}^*\| \right) \|\mathbf{w}_k - \mathbf{w}^*\|.$$

## Newton-type Techniques (unconstrained)

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an [approximation](#) instead !!

## Newton-type Techniques (unconstrained)

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an [approximation](#) instead !!

### Descent direction

If  $M_k \succ 0$  then  $\Delta \mathbf{w} = -M_k^{-1} \nabla \Phi(\mathbf{w})$  is a descent direction.

## Newton-type Techniques (unconstrained)

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an [approximation](#) instead !!

### Descent direction

If  $M_k \succ 0$  then  $\Delta\mathbf{w} = -M_k^{-1}\nabla\Phi(\mathbf{w})$  is a descent direction.

### Local convergence for all Newton Methods

Assume

- $\mathbf{w}^*$  is **SOSC**
- **Lipschitz Hessian:**  $\|M_k^{-1}(\nabla^2\Phi(\mathbf{w}_k) - \nabla^2\Phi(\mathbf{y}))\| \leq \omega \|\mathbf{w}_k - \mathbf{y}\|$  holds on the iterations,  $\forall \mathbf{y}$ .
- **Compatibility:**  $\|M_k^{-1}(\nabla^2\Phi(\mathbf{w}_k) - M_k)\| \leq \kappa_k$  with  $\kappa_k \leq \kappa < 1$

## Newton-type Techniques (unconstrained)

Computing  $H = \nabla_w^2 \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$  can be expensive, use an [approximation](#) instead !!

### Descent direction

If  $M_k \succ 0$  then  $\Delta \mathbf{w} = -M_k^{-1} \nabla \Phi(\mathbf{w})$  is a descent direction.

### Local convergence for all Newton Methods

Assume

- $\mathbf{w}^*$  is **SOSC**
- **Lipschitz Hessian:**  $\|M_k^{-1} (\nabla^2 \Phi(\mathbf{w}_k) - \nabla^2 \Phi(\mathbf{y}))\| \leq \omega \|\mathbf{w}_k - \mathbf{y}\|$  holds on the iterations,  $\forall \mathbf{y}$ .
- **Compatibility:**  $\|M_k^{-1} (\nabla^2 \Phi(\mathbf{w}_k) - M_k)\| \leq \kappa_k$  with  $\kappa_k \leq \kappa < 1$

Then if  $\mathbf{w}_k$  is close to  $\mathbf{w}^*$ ,  $\mathbf{w}_k \rightarrow \mathbf{w}^*$  and convergence is

- **Quadratic** for  $\kappa = 0$ :  $\|\mathbf{w}_{k+1} - \mathbf{w}^*\| \leq C \|\mathbf{w}_k - \mathbf{w}^*\|^2$  with  $C = \omega/2$
- **Superlinear** for  $\kappa_k \rightarrow 0$ :  $\|\mathbf{w}_{k+1} - \mathbf{w}^*\| \leq C_k \|\mathbf{w}_k - \mathbf{w}^*\|$  with  $C_k \rightarrow 0$
- **Linear** for  $\kappa_k > \rho$ :  $\|\mathbf{w}_{k+1} - \mathbf{w}^*\| \leq C \|\mathbf{w}_k - \mathbf{w}^*\|$  with  $C < 1$

## Steepest descent

Use  $M_k = \alpha_k^{-1} I$ , then:

$$\Delta \mathbf{w}_k = -M_k^{-1} \nabla \Phi(\mathbf{w}_k) = -\alpha_k \nabla \Phi(\mathbf{w}_k)$$

Step size  $\alpha_k$  is chosen sufficiently small by the line-search.

## Convergence

- Compatibility:  $\|\alpha_k \nabla^2 \Phi(\mathbf{w}_k) - I\| \leq \kappa_k$  with  $\kappa_k \leq \kappa < 1$
- Constant  $\kappa_k$  does not converge to 0, i.e.  $\kappa_k > \rho$ ,  $\forall k$
- Linear convergence when  $\mathbf{w}_k$  is close to  $\mathbf{w}^*$

## Newton-type Methods (unconstrained)

Cost function of the type  $\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{R}(\mathbf{w})\|^2$ , with  $\mathbf{R}(\mathbf{w}) \in \mathbb{R}^m$

### Gauss-Newton Hessian approximation

Observe that

$$H(\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} (\nabla \mathbf{R}(\mathbf{w}) \mathbf{R}(\mathbf{w})) = \nabla \mathbf{R}(\mathbf{w}) \nabla \mathbf{R}(\mathbf{w})^\top + \sum_{i=1}^m \nabla^2 \mathbf{R}_i(\mathbf{w}) \mathbf{R}_i(\mathbf{w})$$

## Newton-type Methods (unconstrained)

Cost function of the type  $\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{R}(\mathbf{w})\|^2$ , with  $\mathbf{R}(\mathbf{w}) \in \mathbb{R}^m$

### Gauss-Newton Hessian approximation

Observe that

$$H(\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} (\nabla \mathbf{R}(\mathbf{w}) \mathbf{R}(\mathbf{w})) = \nabla \mathbf{R}(\mathbf{w}) \nabla \mathbf{R}(\mathbf{w})^\top + \sum_{i=1}^m \nabla^2 \mathbf{R}_i(\mathbf{w}) \mathbf{R}_i(\mathbf{w})$$

Gauss-Newton method proposes to use:

$$B_k = \nabla \mathbf{R}(\mathbf{w}_k) \nabla \mathbf{R}(\mathbf{w}_k)^\top + \alpha_k I$$

## Newton-type Methods (unconstrained)

Cost function of the type  $\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{R}(\mathbf{w})\|^2$ , with  $\mathbf{R}(\mathbf{w}) \in \mathbb{R}^m$

### Gauss-Newton Hessian approximation

Observe that

$$H(\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} (\nabla \mathbf{R}(\mathbf{w}) \mathbf{R}(\mathbf{w})) = \nabla \mathbf{R}(\mathbf{w}) \nabla \mathbf{R}(\mathbf{w})^\top + \sum_{i=1}^m \nabla^2 \mathbf{R}_i(\mathbf{w}) \mathbf{R}_i(\mathbf{w})$$

Gauss-Newton method proposes to use:

$$\boxed{B_k = \nabla \mathbf{R}(\mathbf{w}_k) \nabla \mathbf{R}(\mathbf{w}_k)^\top + \alpha_k I}$$

$B_k$  is a good approximation if:

- all  $\nabla^2 \mathbf{R}_i(\mathbf{w})$  are small ( $\mathbf{R}$  close to linear), or
- all  $\mathbf{R}_i(\mathbf{w})$  are small

## Newton-type Methods (unconstrained)

Cost function of the type  $\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{R}(\mathbf{w})\|^2$ , with  $\mathbf{R}(\mathbf{w}) \in \mathbb{R}^m$

### Gauss-Newton Hessian approximation

Observe that

$$H(\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} (\nabla \mathbf{R}(\mathbf{w}) \mathbf{R}(\mathbf{w})) = \nabla \mathbf{R}(\mathbf{w}) \nabla \mathbf{R}(\mathbf{w})^\top + \sum_{i=1}^m \nabla^2 \mathbf{R}_i(\mathbf{w}) \mathbf{R}_i(\mathbf{w})$$

Gauss-Newton method proposes to use:

$$\boxed{B_k = \nabla \mathbf{R}(\mathbf{w}_k) \nabla \mathbf{R}(\mathbf{w}_k)^\top + \alpha_k I}$$

$B_k$  is a good approximation if:

- all  $\nabla^2 \mathbf{R}_i(\mathbf{w})$  are small ( $\mathbf{R}$  close to linear), or
- all  $\mathbf{R}_i(\mathbf{w})$  are small

Typical application to **fitting problems** or **control problems**:  $\mathbf{R}(\mathbf{w}) = \mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}}$

## Newton-type Methods (unconstrained)

Cost function of the type  $\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{R}(\mathbf{w})\|^2$ , with  $\mathbf{R}(\mathbf{w}) \in \mathbb{R}^m$

### Gauss-Newton Hessian approximation

Observe that

$$H(\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} (\nabla \mathbf{R}(\mathbf{w}) \mathbf{R}(\mathbf{w})) = \nabla \mathbf{R}(\mathbf{w}) \nabla \mathbf{R}(\mathbf{w})^\top + \sum_{i=1}^m \nabla^2 \mathbf{R}_i(\mathbf{w}) \mathbf{R}_i(\mathbf{w})$$

Gauss-Newton method proposes to use:

$$\boxed{B_k = \nabla \mathbf{R}(\mathbf{w}_k) \nabla \mathbf{R}(\mathbf{w}_k)^\top + \alpha_k I}$$

$B_k$  is a good approximation if:

- all  $\nabla^2 \mathbf{R}_i(\mathbf{w})$  are small ( $\mathbf{R}$  close to linear), or
- all  $\mathbf{R}_i(\mathbf{w})$  are small

Typical application to **fitting problems** or **control problems**:  $\mathbf{R}(\mathbf{w}) = \mathbf{y}(\mathbf{w}) - \bar{\mathbf{y}}$

### Convergence

- If  $\sum_{i=1}^m \nabla^2 \mathbf{R}_i(\mathbf{w}) \mathbf{R}_i(\mathbf{w}) \rightarrow 0$  then  $\kappa_k \rightarrow 0$
- Superlinear convergence in some cases...

## Newton-type Methods (unconstrained)

Compute  $H(\mathbf{w})$  numerically in an efficient (iterative) way

### BFGS

Define

$$\mathbf{s}_k = \mathbf{w}_{k+1} - \mathbf{w}_k$$

$$\mathbf{y}_k = \nabla\Phi(\mathbf{w}_{k+1}) - \nabla\Phi(\mathbf{w}_k)$$

Idea: Update  $M_k \rightarrow M_{k+1}$  such that  $M_{k+1}\mathbf{s}_k = \mathbf{y}_k$  (secant condition)

## Newton-type Methods (unconstrained)

Compute  $H(\mathbf{w})$  numerically in an efficient (iterative) way

### BFGS

Define

$$\mathbf{s}_k = \mathbf{w}_{k+1} - \mathbf{w}_k$$

$$\mathbf{y}_k = \nabla\Phi(\mathbf{w}_{k+1}) - \nabla\Phi(\mathbf{w}_k)$$

Idea: Update  $M_k \rightarrow M_{k+1}$  such that  $M_{k+1}\mathbf{s}_k = \mathbf{y}_k$  (secant condition)

BFGS formula:

$$M_{k+1} = M_k - \frac{M_k \mathbf{s}_k \mathbf{s}_k^\top M_k}{\mathbf{s}_k^\top M_k \mathbf{s}_k} + \frac{\mathbf{y}_k \mathbf{y}_k^\top}{\mathbf{s}_k^\top \mathbf{y}_k}, \quad M_0 = I$$

## Newton-type Methods (unconstrained)

Compute  $H(\mathbf{w})$  numerically in an efficient (iterative) way

### BFGS

Define

$$\mathbf{s}_k = \mathbf{w}_{k+1} - \mathbf{w}_k$$

$$\mathbf{y}_k = \nabla\Phi(\mathbf{w}_{k+1}) - \nabla\Phi(\mathbf{w}_k)$$

Idea: Update  $M_k \rightarrow M_{k+1}$  such that  $M_{k+1}\mathbf{s}_k = \mathbf{y}_k$  (secant condition)

BFGS formula:

$$M_{k+1} = M_k - \frac{M_k \mathbf{s}_k \mathbf{s}_k^\top M_k}{\mathbf{s}_k^\top M_k \mathbf{s}_k} + \frac{\mathbf{y}_k \mathbf{y}_k^\top}{\mathbf{s}_k^\top \mathbf{y}_k}, \quad M_0 = I$$

See "Powell's trick" ensure  $M_{k+1} > 0$

## Newton-type Methods (unconstrained)

Compute  $H(\mathbf{w})$  numerically in an efficient (iterative) way

### BFGS

Define

$$\mathbf{s}_k = \mathbf{w}_{k+1} - \mathbf{w}_k$$

$$\mathbf{y}_k = \nabla\Phi(\mathbf{w}_{k+1}) - \nabla\Phi(\mathbf{w}_k)$$

Idea: Update  $M_k \rightarrow M_{k+1}$  such that  $M_{k+1}\mathbf{s}_k = \mathbf{y}_k$  (secant condition)

BFGS formula: 
$$M_{k+1} = M_k - \frac{M_k \mathbf{s}_k \mathbf{s}_k^\top M_k}{\mathbf{s}_k^\top M_k \mathbf{s}_k} + \frac{\mathbf{y}_k \mathbf{y}_k^\top}{\mathbf{s}_k^\top \mathbf{y}_k}, \quad M_0 = I$$

See "Powell's trick" ensure  $M_{k+1} > 0$

### Convergence

- It can be shown that  $M_k \rightarrow H(\mathbf{w})$ , then  $\kappa_k \rightarrow 0$
- Superlinear convergence after enough (and “sufficiently rich”) iterations

## Newton-type Methods (unconstrained)

Compute  $H(\mathbf{w})$  numerically in an efficient (iterative) way

### BFGS

Define

$$\mathbf{s}_k = \mathbf{w}_{k+1} - \mathbf{w}_k$$

$$\mathbf{y}_k = \nabla\Phi(\mathbf{w}_{k+1}) - \nabla\Phi(\mathbf{w}_k)$$

Idea: Update  $M_k \rightarrow M_{k+1}$  such that  $M_{k+1}\mathbf{s}_k = \mathbf{y}_k$  (secant condition)

BFGS formula: 
$$M_{k+1} = M_k - \frac{M_k \mathbf{s}_k \mathbf{s}_k^\top M_k}{\mathbf{s}_k^\top M_k \mathbf{s}_k} + \frac{\mathbf{y}_k \mathbf{y}_k^\top}{\mathbf{s}_k^\top \mathbf{y}_k}, \quad M_0 = I$$

See "Powell's trick" ensure  $M_{k+1} > 0$

### Convergence

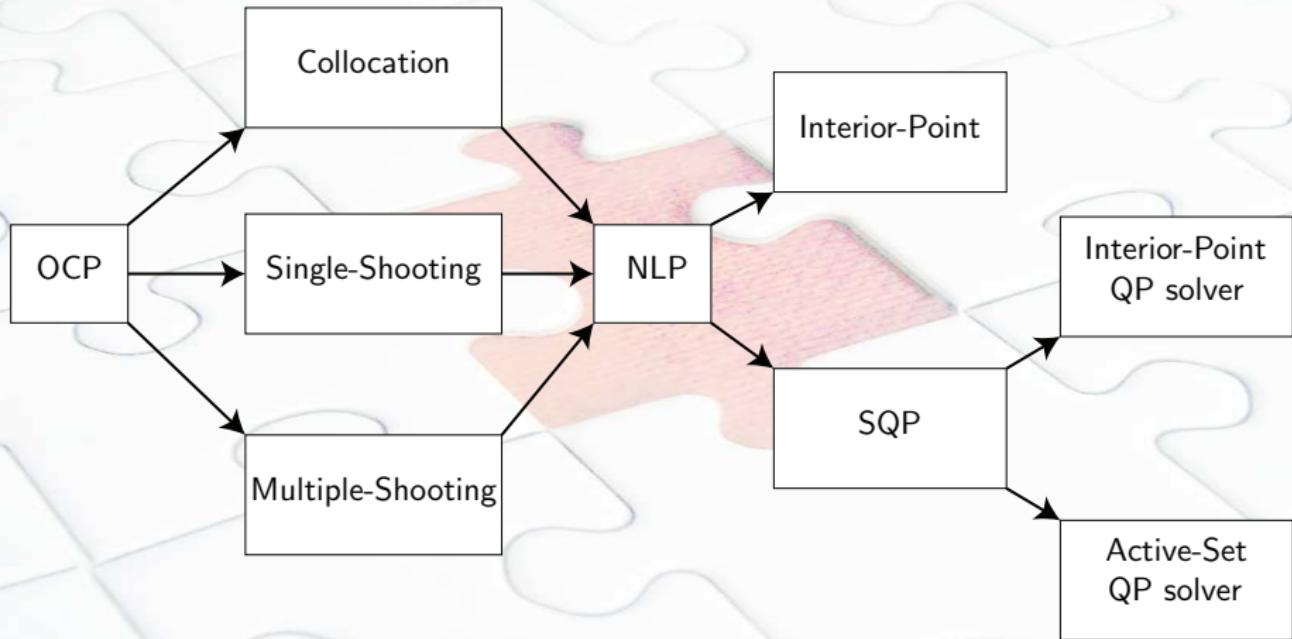
- It can be shown that  $M_k \rightarrow H(\mathbf{w})$ , then  $\kappa_k \rightarrow 0$
- Superlinear convergence after enough (and “sufficiently rich”) iterations

Obs: BFGS typically yields **dense Hessian** approximations, whereas the exact Hessian may be sparse. Factorization can become *unnecessarily expensive* !!

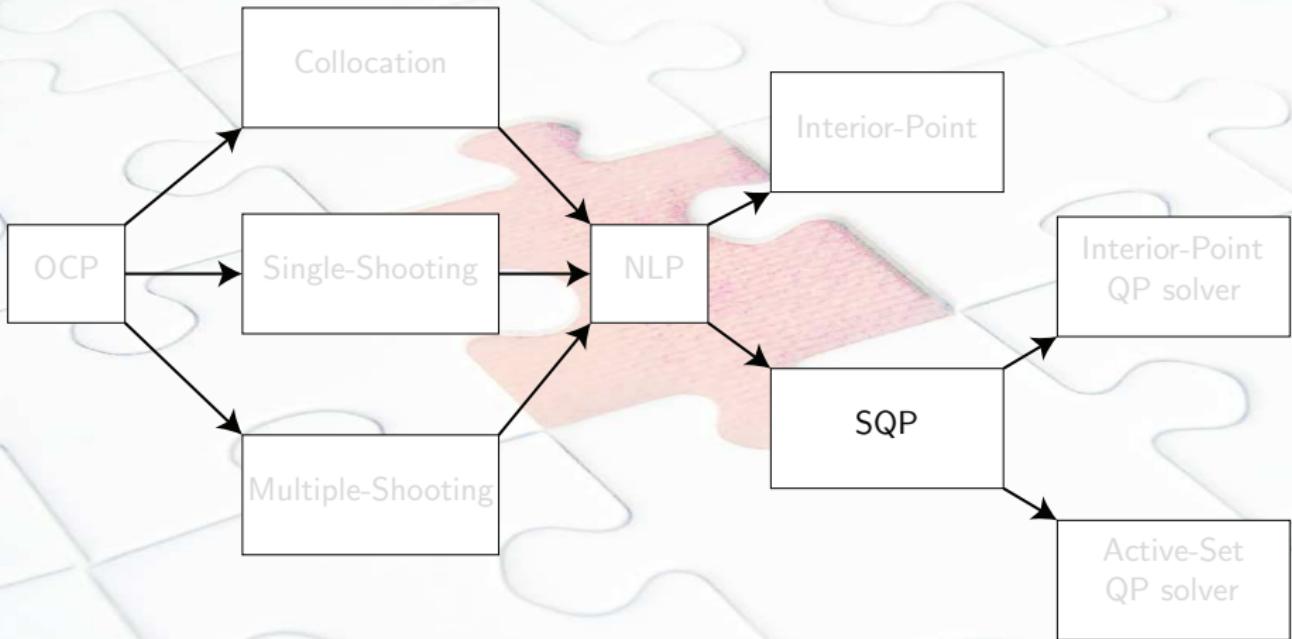
# Outline

- 1 The Newton method
- 2 Newton on the KKT conditions
- 3 The reduced Newton step (unconstrained problems)
- 4 The merit function - Line-search for constrained problems
- 5 Newton-type methods
- 6 Sequential Quadratic Programming

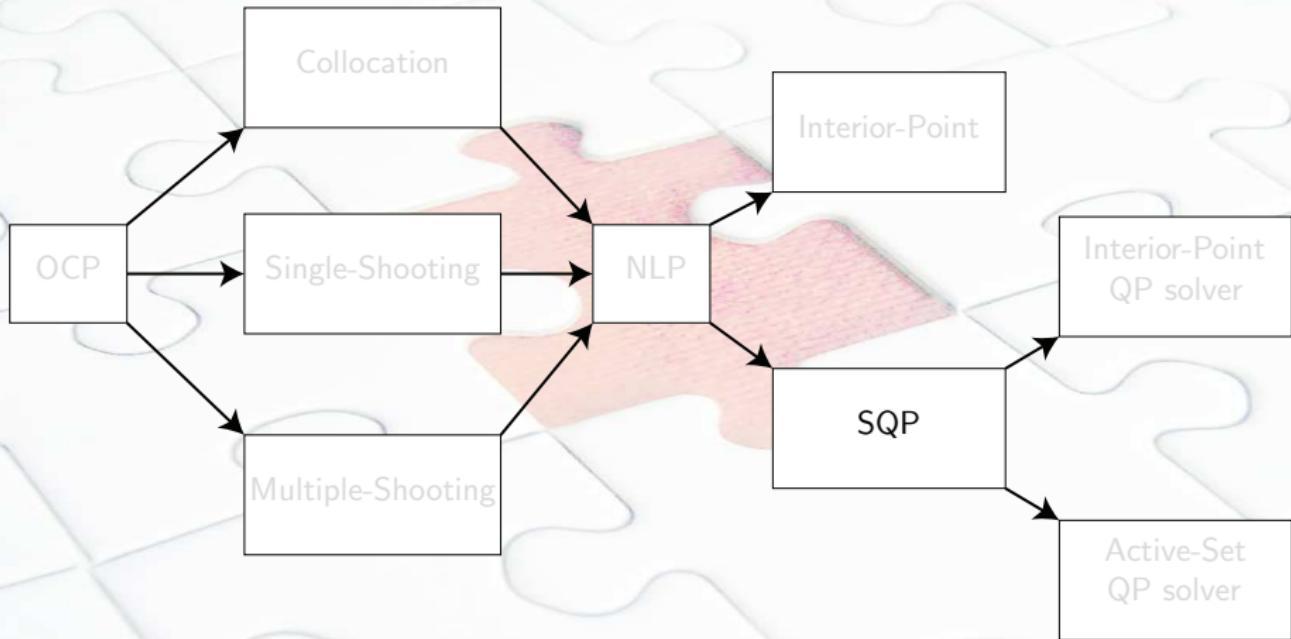
# Survival map of Direct Optimal Control



# Survival map of Direct Optimal Control



# Survival map of Direct Optimal Control



## What about inequality constraints ?

Find the "primal-dual" variables  $x, \mu, \lambda$  such that:

**Primal Feasibility:**  $g(w) = 0, h(w) \leq 0,$

**Dual Feasibility:**  $\nabla_w \mathcal{L}(w, \mu, \lambda) = 0, \mu \geq 0,$

**Complementarity Slackness:**  $\mu_i h_i(w) = 0, i = 1, \dots$

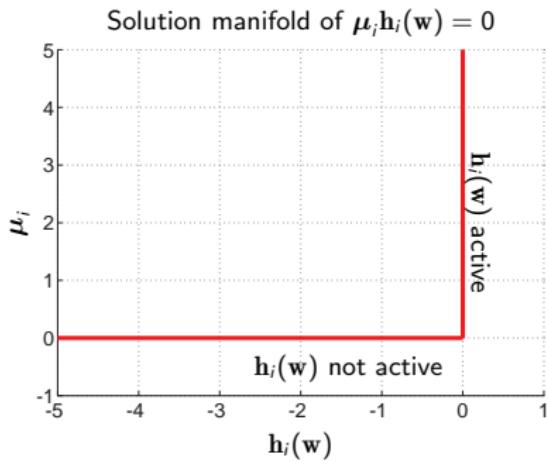
## What about inequality constraints ?

Find the "primal-dual" variables  $x, \mu, \lambda$  such that:

**Primal Feasibility:**  $g(w) = 0, h(w) \leq 0,$

**Dual Feasibility:**  $\nabla_w \mathcal{L}(w, \mu, \lambda) = 0, \mu \geq 0,$

**Complementarity Slackness:**  $\mu_i h_i(w) = 0, i = 1, \dots$



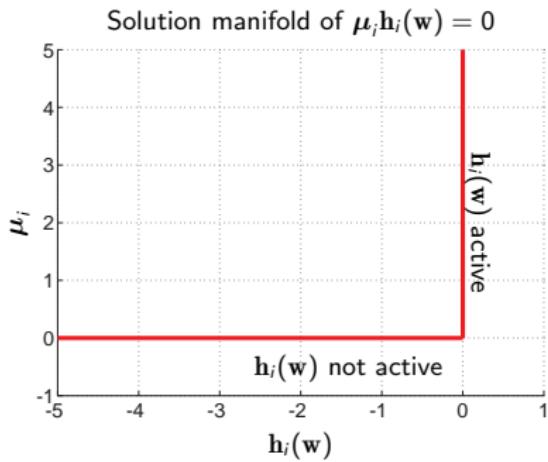
## What about inequality constraints ?

Find the "primal-dual" variables  $x, \mu, \lambda$  such that:

**Primal Feasibility:**  $g(w) = 0, h(w) \leq 0,$

**Dual Feasibility:**  $\nabla_w \mathcal{L}(w, \mu, \lambda) = 0, \mu \geq 0,$

**Complementarity Slackness:**  $\mu_i h_i(w) = 0, i = 1, \dots$



Manifold generated by the Complementarity Slackness condition is **not-smooth**, Newton can not be used !!

## Quadratic model interpretation

Problem:

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

The **Newton direction** is given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

with  $H(w, \lambda) = \nabla_w^2 \mathcal{L}(w, \lambda)$

## Quadratic model interpretation

Problem:

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

The **Newton direction** is given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

with  $H(w, \lambda) = \nabla_w^2 \mathcal{L}(w, \lambda)$

The **Newton direction** is given by the Quadratic Program (QP):

$$\begin{array}{ll} \min_{\Delta w} & \frac{1}{2} \Delta w^T H(w, \lambda) \Delta w + \nabla \Phi(w)^T \Delta w \\ \text{s.t.} & g(w) + \nabla g(w)^T \Delta w = 0 \end{array}$$

## Quadratic model interpretation

Problem:

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

The **Newton direction** is given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

with  $H(w, \lambda) = \nabla_w^2 \mathcal{L}(w, \lambda)$

The **Newton direction** is given by the Quadratic Program (QP):

$$\begin{array}{ll} \min_{\Delta w} & \frac{1}{2} \Delta w^T H(w, \lambda) \Delta w + \nabla \Phi(w)^T \Delta w \\ \text{s.t.} & g(w) + \nabla g(w)^T \Delta w = 0 \end{array}$$

Dual variables  $\lambda^+$  given by the dual variables of the QP, i.e.  $\lambda^+ = \lambda_{QP}$

## Quadratic model interpretation

Problem:

$$\begin{array}{ll} \min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \end{array}$$

The **Newton direction** is given by

$$\begin{bmatrix} H(w, \lambda) & \nabla g(w) \\ \nabla g(w)^T & 0 \end{bmatrix} \begin{bmatrix} \Delta w \\ \lambda^+ \end{bmatrix} = - \begin{bmatrix} \nabla \Phi(w) \\ g(w) \end{bmatrix}$$

with  $H(w, \lambda) = \nabla_w^2 \mathcal{L}(w, \lambda)$

The **Newton direction** is given by the Quadratic Program (QP):

$$\begin{array}{ll} \min_{\Delta w} & \frac{1}{2} \Delta w^T H(w, \lambda) \Delta w + \nabla \Phi(w)^T \Delta w \\ \text{s.t.} & g(w) + \nabla g(w)^T \Delta w = 0 \end{array}$$

Dual variables  $\lambda^+$  given by the dual variables of the QP, i.e.  $\lambda^+ = \lambda_{QP}$

*Proof: KKT of the QP are equivalent to the linearised KKT system of the original problem.*

## Quadratic interpretation for inequality constraints

Problem:

$$\min_w \Phi(w)$$

$$\text{s.t. } g(w) = 0$$

$$\text{s.t. } h(w) \leq 0$$

## Quadratic interpretation for inequality constraints

Problem:

$$\begin{array}{ll}\min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \\ \text{s.t.} & h(w) \leq 0\end{array}$$

The **Newton direction** is given by the Quadratic Program (QP):

$$\begin{array}{ll}\min_{\Delta w} & \frac{1}{2} \Delta w^T H(w, \lambda, \mu) \Delta w + \nabla \Phi(w)^T \Delta w \\ \text{s.t.} & g(w) + \nabla g(w)^T \Delta w = 0 \\ & h(w) + \nabla h(w)^T \Delta w \leq 0\end{array}$$

with  $H(w, \lambda) = \nabla_w^2 \mathcal{L}(w, \lambda)$

## Quadratic interpretation for inequality constraints

Problem:

$$\begin{array}{ll}\min_w & \Phi(w) \\ \text{s.t.} & g(w) = 0 \\ \text{s.t.} & h(w) \leq 0\end{array}$$

The **Newton direction** is given by the Quadratic Program (QP):

$$\begin{array}{ll}\min_{\Delta w} & \frac{1}{2} \Delta w^T H(w, \lambda, \mu) \Delta w + \nabla \Phi(w)^T \Delta w \\ \text{s.t.} & g(w) + \nabla g(w)^T \Delta w = 0 \\ & h(w) + \nabla h(w)^T \Delta w \leq 0\end{array}$$

with  $H(w, \lambda) = \nabla_w^2 \mathcal{L}(w, \lambda)$

Dual variables  $\lambda^+$  and  $\mu^+$  given by the dual variables of the QP, i.e.

$$\lambda^+ = \lambda_{QP}, \quad \mu^+ = \mu_{QP}$$

## SQP - Monitoring progress with the $T_1$ merit function:

$$T_1(\mathbf{w}) = \Phi(\mathbf{w}) + \nu \|\mathbf{g}(\mathbf{w})\|_1 + \nu \sum_{i=1}^m \max(0, h_i(\mathbf{w}))$$

**Algorithm:** SQP with line-search

**Input:** guess  $\mathbf{w}$ ,  $\lambda$ ,  $\mu$

**while**  $\|\nabla \mathcal{L}\|$  or  $\|\mathbf{g}\|$  or  $\max(0, \mathbf{h}) \geq \text{tol}$  **do**

    Compute  $\mathbf{g}$ ,  $\mathbf{h}$ ,  $\nabla \Phi(\mathbf{w})$ ,  $\nabla \mathbf{g}(\mathbf{w})$ ,  $\nabla \mathbf{h}(\mathbf{w})$ ,  $H(\mathbf{w}, \mu, \lambda)$

    Compute **Newton direction** by solving the QP

$$\min_{\Delta \mathbf{w}} \quad \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \lambda, \mu) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w}$$

$$\text{s.t.} \quad \mathbf{g}(\mathbf{w}) + \nabla \mathbf{g}(\mathbf{w})^\top \Delta \mathbf{w} = 0$$

$$\mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0$$

    Perform line-search on  $T_1(\mathbf{w} + \alpha \Delta \mathbf{w})$ , get step length  $\alpha$

    Take primal step:  $\mathbf{w} \leftarrow \mathbf{w} + \alpha \Delta \mathbf{w}$

    Take dual step:  $\lambda \leftarrow (1 - \alpha)\lambda + \alpha \lambda_{QP}$ ,  $\mu \leftarrow (1 - \alpha)\mu + \alpha \mu_{QP}$

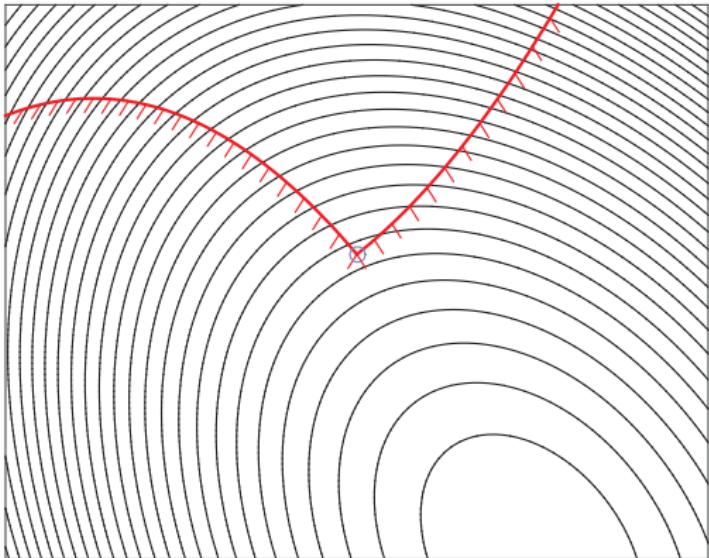
**return**  $\mathbf{w}$ ,  $\lambda$ ,  $\mu$

If  $H(\mathbf{x}_k, \mu_k, \lambda_k) > 0$  and  $\nu > \max\{\|\mu_{k+1}\|_\infty, \|\lambda_{k+1}\|_\infty\}$  then  $\Delta \mathbf{w}$  is a descent direction for  $T_1(\mathbf{w})$

## SQP - Illustration

**NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$



**QP:**

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

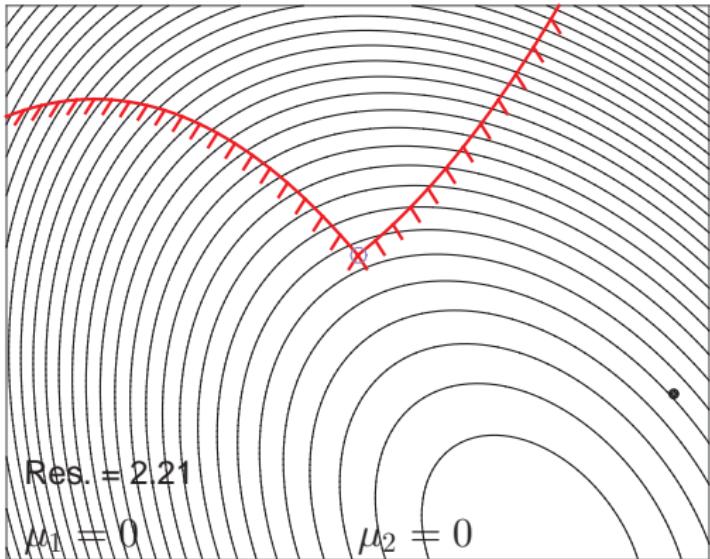
**Hessian:**

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## SQP - Illustration

NLP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$



QP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

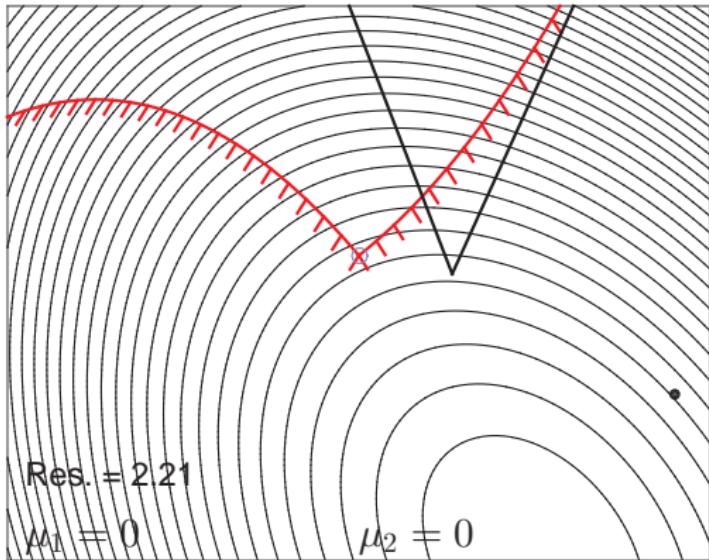
Hessian:

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## Linearized constraints

**NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

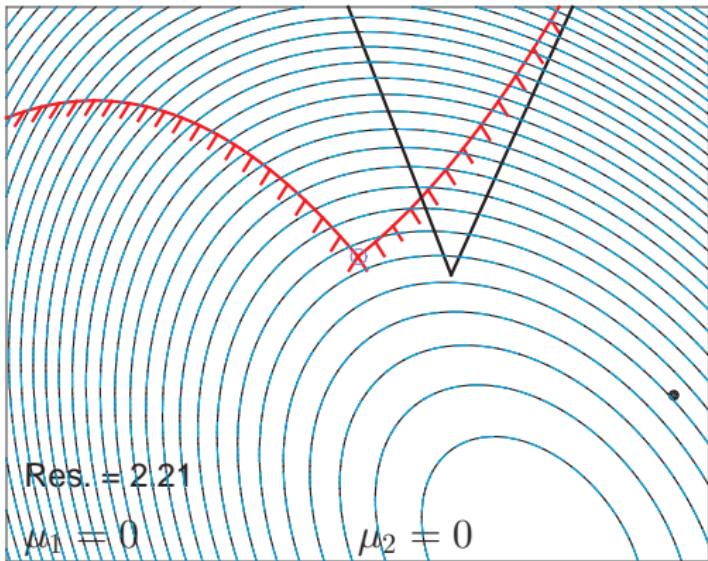
**Hessian:**

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## Contours of QP cost

**NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

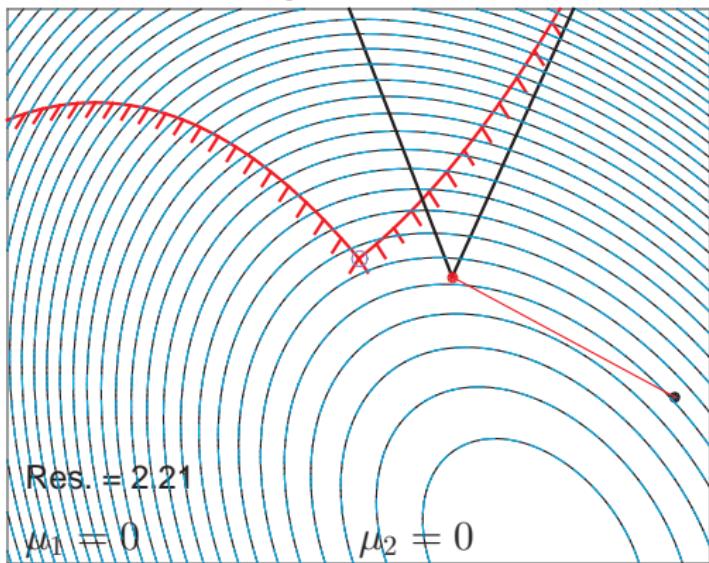
$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

**Hessian:**

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

Step with  $t = 1$ **NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

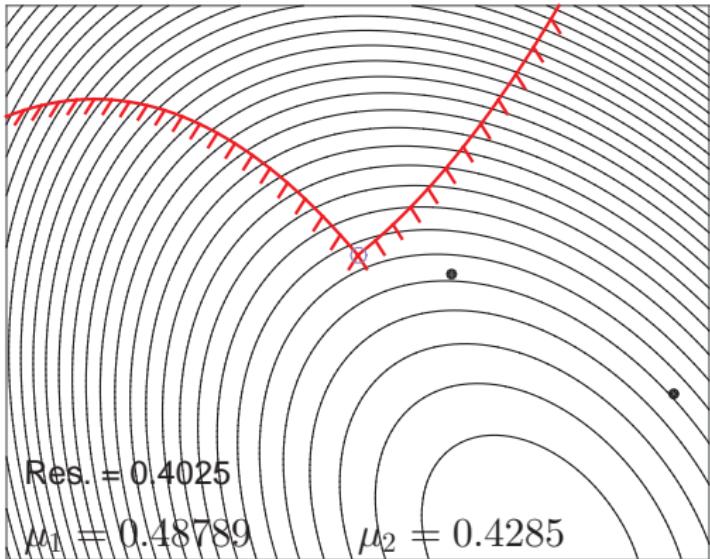
**Hessian:**

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## SQP - Illustration

NLP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$



QP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

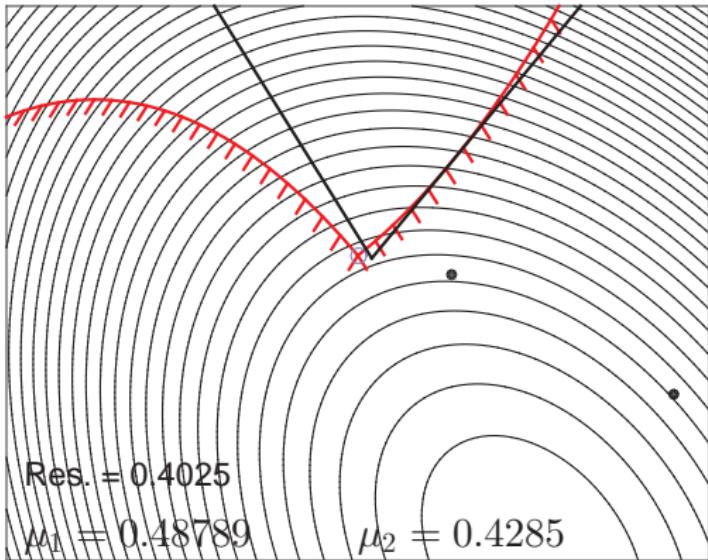
Hessian:

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## Linearized constraints

**NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

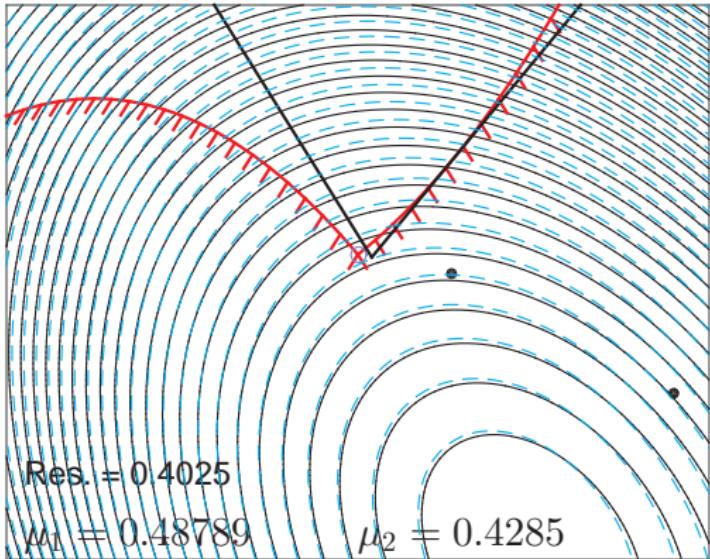
**Hessian:**

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## Contours of QP cost

**NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

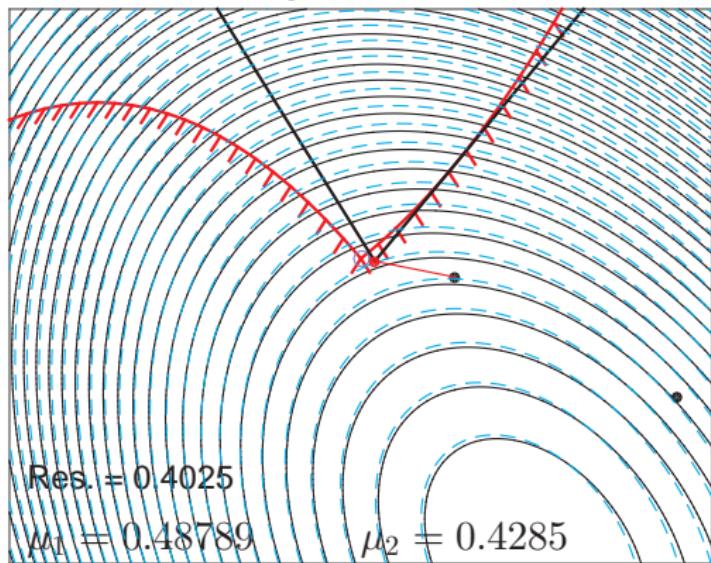
$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \mu) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

**Hessian:**

$$H(\mathbf{w}, \mu) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\mu^\top \mathbf{h}(\mathbf{w}))$$

**Step with  $t = 1$** **NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \mu) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

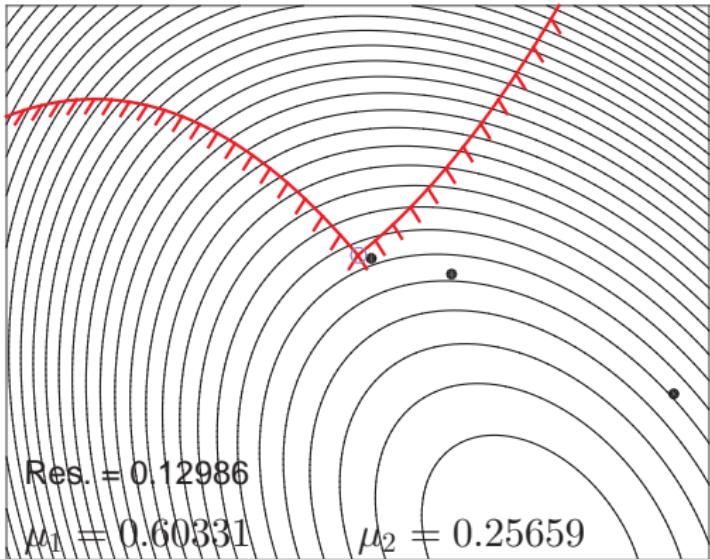
**Hessian:**

$$H(\mathbf{w}, \mu) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\mu^\top \mathbf{h}(\mathbf{w}))$$

## SQP - Illustration

NLP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$



QP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

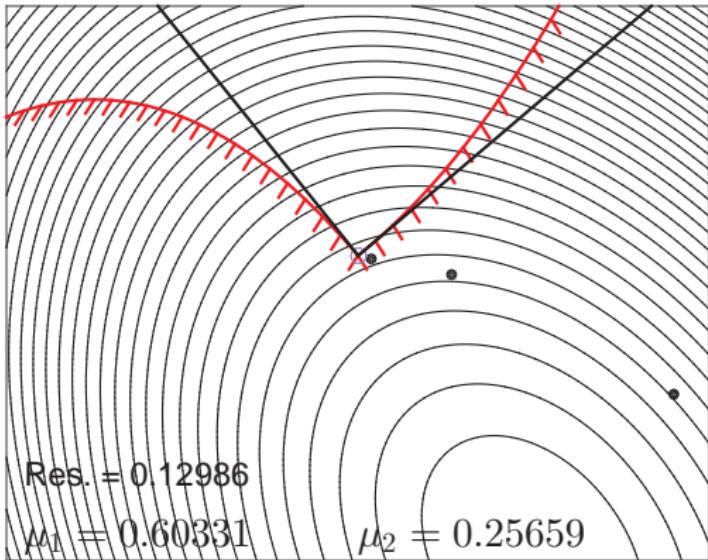
Hessian:

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## Linearized constraints

**NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

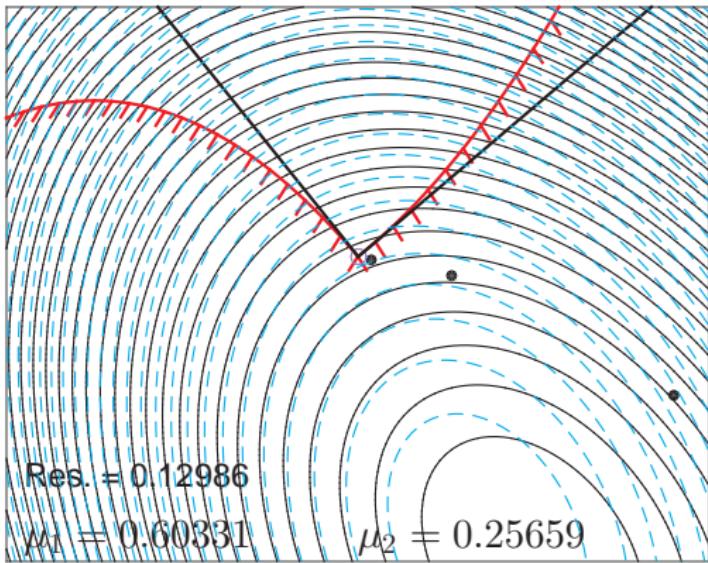
$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

**Hessian:**

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

**NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

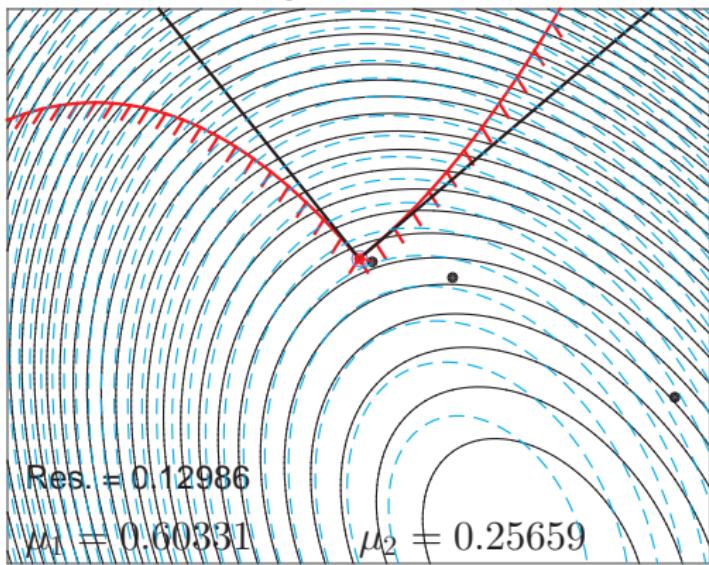
$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

**Hessian:**

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

Step with  $t = 1$ **NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

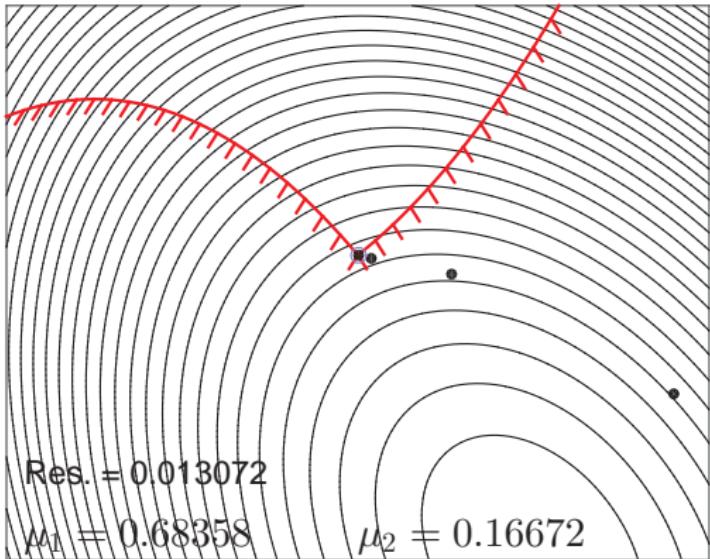
**Hessian:**

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## SQP - Illustration

NLP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$



QP:

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

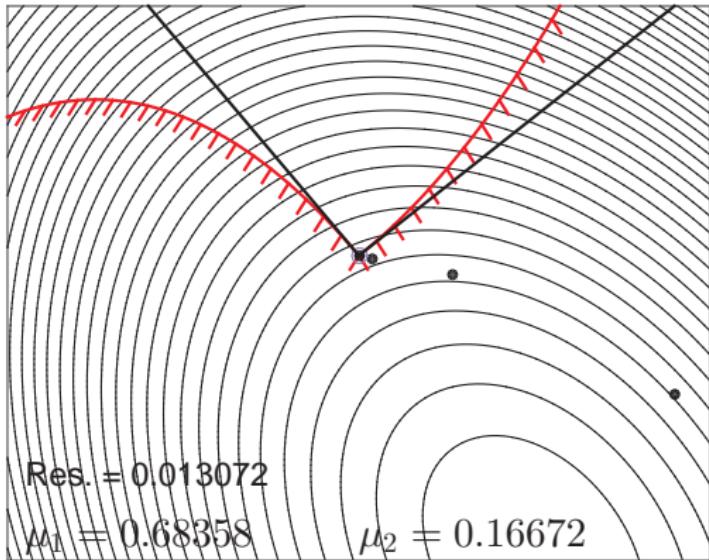
Hessian:

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## Linearized constraints

**NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

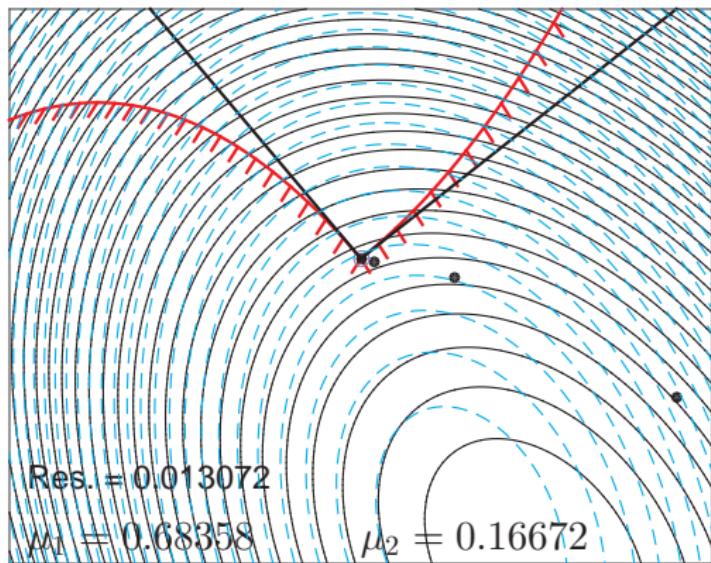
**Hessian:**

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## Contours of QP cost

**NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

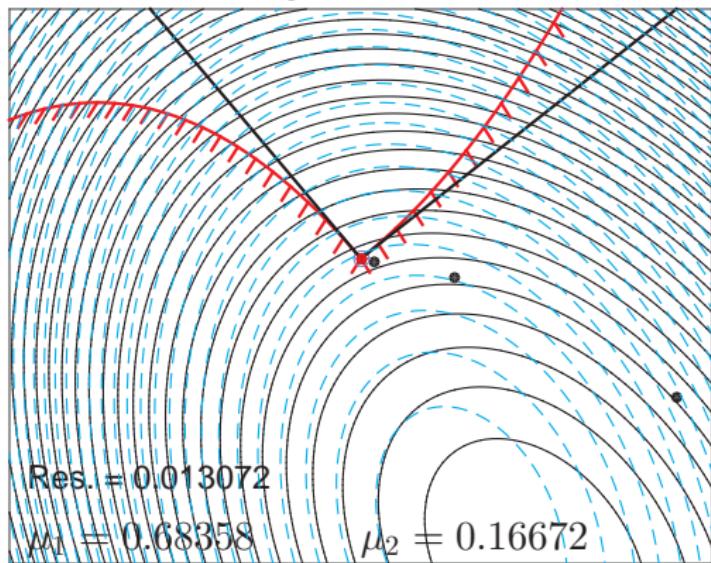
$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \mu) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

**Hessian:**

$$H(\mathbf{w}, \mu) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\mu^\top \mathbf{h}(\mathbf{w}))$$

Step with  $t = 1$ **NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

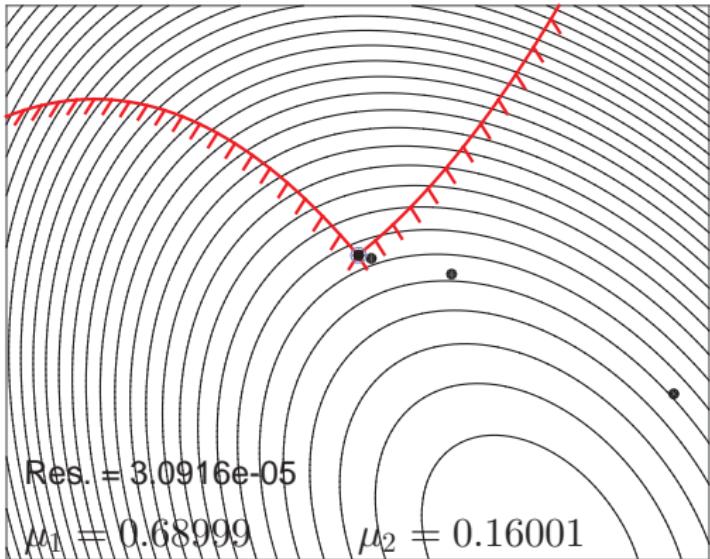
**Hessian:**

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## SQP - Illustration

NLP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$



QP:

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

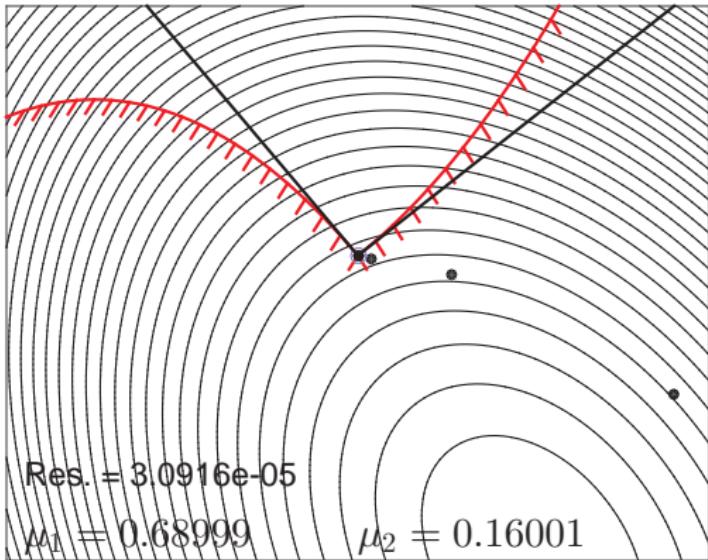
Hessian:

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## Linearized constraints

**NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

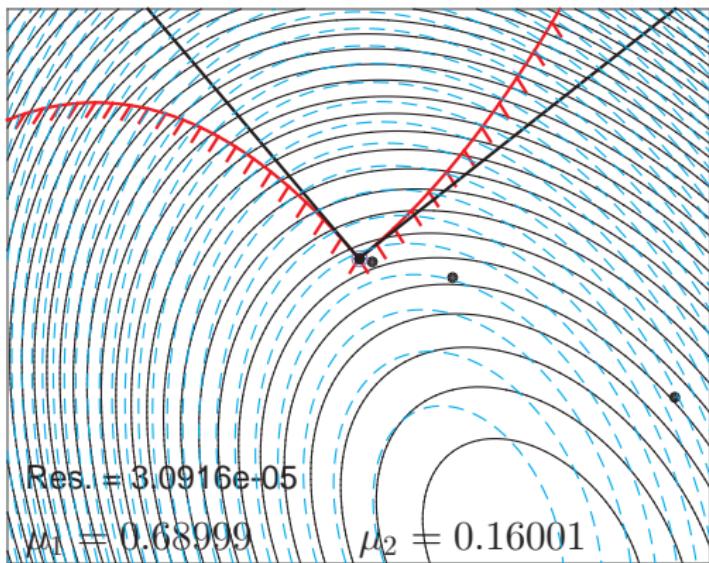
**Hessian:**

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## Contours of QP cost

**NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

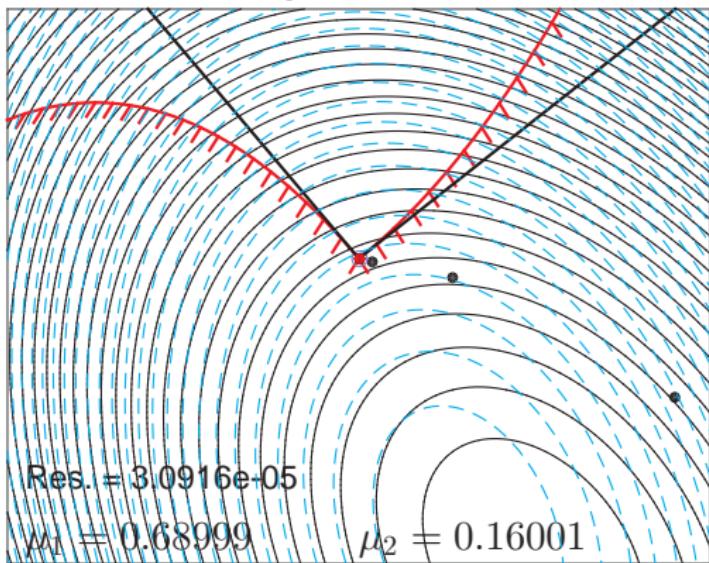
$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \mu) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

**Hessian:**

$$H(\mathbf{w}, \mu) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\mu^\top \mathbf{h}(\mathbf{w}))$$

Step with  $t = 1$ **NLP:**

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$

**QP:**

$$\begin{aligned} \min_{\Delta \mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \mu) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

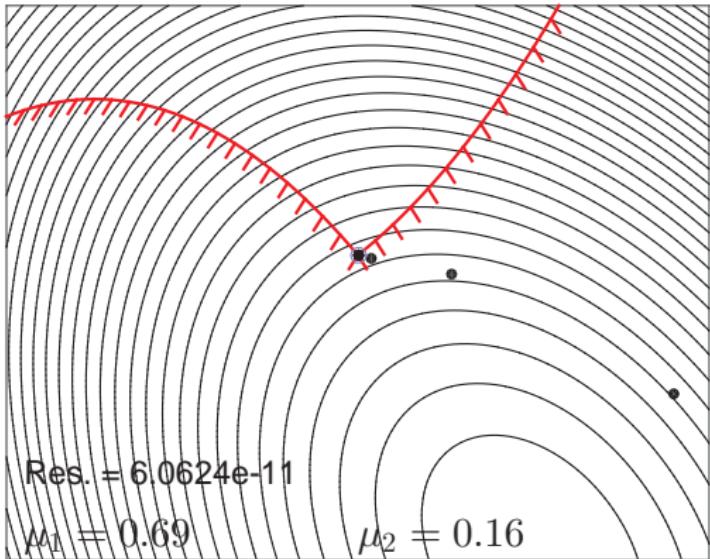
**Hessian:**

$$H(\mathbf{w}, \mu) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\mu^\top \mathbf{h}(\mathbf{w}))$$

## SQP - Illustration

NLP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$



QP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

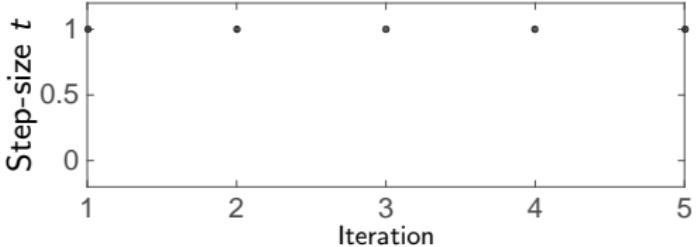
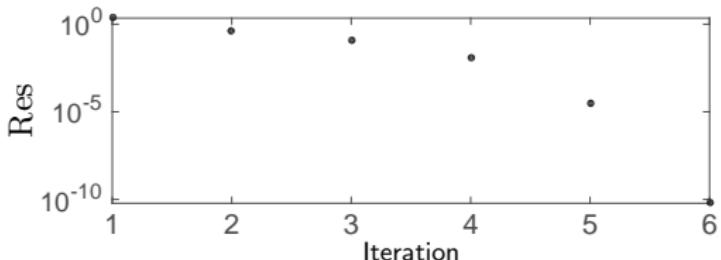
Hessian:

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## SQP - Illustration

NLP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w} - \mathbf{w}_0\|_Q^2 \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) \leq 0 \end{aligned}$$



QP:

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \Delta \mathbf{w}^\top H(\mathbf{w}, \boldsymbol{\mu}) \Delta \mathbf{w} + \nabla \Phi(\mathbf{w})^\top \Delta \mathbf{w} \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{w}) + \nabla \mathbf{h}(\mathbf{w})^\top \Delta \mathbf{w} \leq 0 \end{aligned}$$

Hessian:

$$H(\mathbf{w}, \boldsymbol{\mu}) = \nabla^2 \Phi(\mathbf{w}) + \nabla^2 (\boldsymbol{\mu}^\top \mathbf{h}(\mathbf{w}))$$

## Maratos effect

Consider the NLP :

$$\min_{u,v} \Phi = 3v^2 - 2u$$

$$\text{s.t. } g = u - v^2 = 0$$

Optimum:  $\mathbf{w}^* = [ \begin{array}{cc} 0 & 0 \end{array} ]$ .

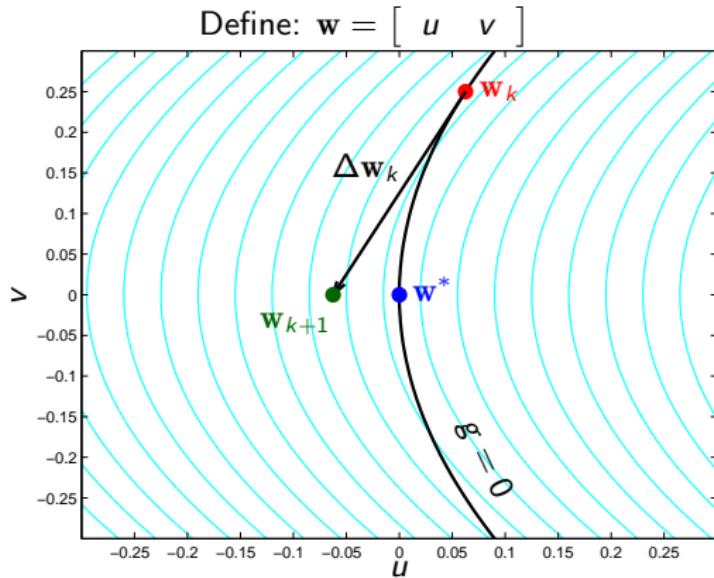
Consider the iterate:

$$\mathbf{w}_k = [ \begin{array}{cc} a^2 & a \end{array} ]$$

The Newton step is:

$$\Delta \mathbf{w}_k = - [ \begin{array}{cc} 2a^2 & a \end{array} ]$$

for  $\lambda = 2 \dots$



## Maratos effect

Consider the NLP :

$$\min_{u,v} \Phi = 3v^2 - 2u$$

$$\text{s.t. } g = u - v^2 = 0$$

Optimum:  $\mathbf{w}^* = [ \begin{array}{cc} 0 & 0 \end{array} ]$ .

Consider the iterate:

$$\mathbf{w}_k = [ \begin{array}{cc} a^2 & a \end{array} ]$$

The Newton step is:

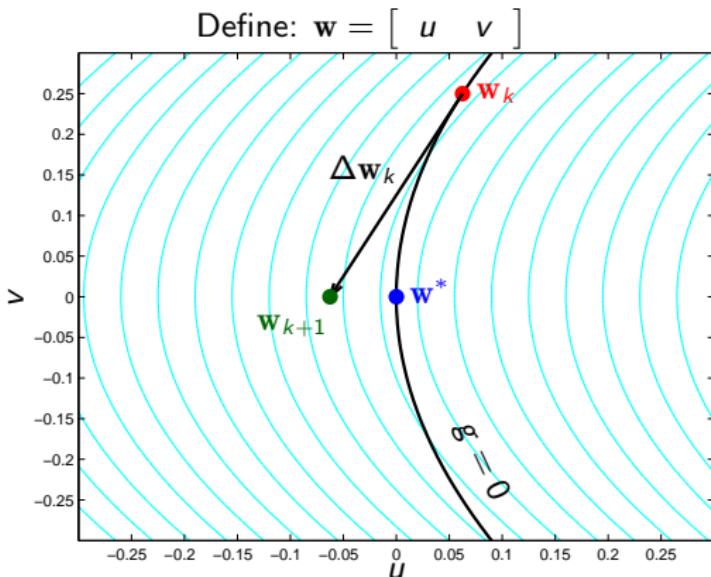
$$\Delta \mathbf{w}_k = - [ \begin{array}{cc} 2a^2 & a \end{array} ]$$

for  $\lambda = 2$ ...

The full Newton step:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \Delta \mathbf{w}_k$$

is much closer to  $\mathbf{w}^*$  than  $\mathbf{w}_k$ .



## Maratos effect

Consider the NLP :

$$\min_{u,v} \Phi = 3v^2 - 2u$$

$$\text{s.t. } g = u - v^2 = 0$$

Optimum:  $\mathbf{w}^* = [ \begin{array}{cc} 0 & 0 \end{array} ]$ .

Consider the iterate:

$$\mathbf{w}_k = [ \begin{array}{cc} a^2 & a \end{array} ]$$

The Newton step is:

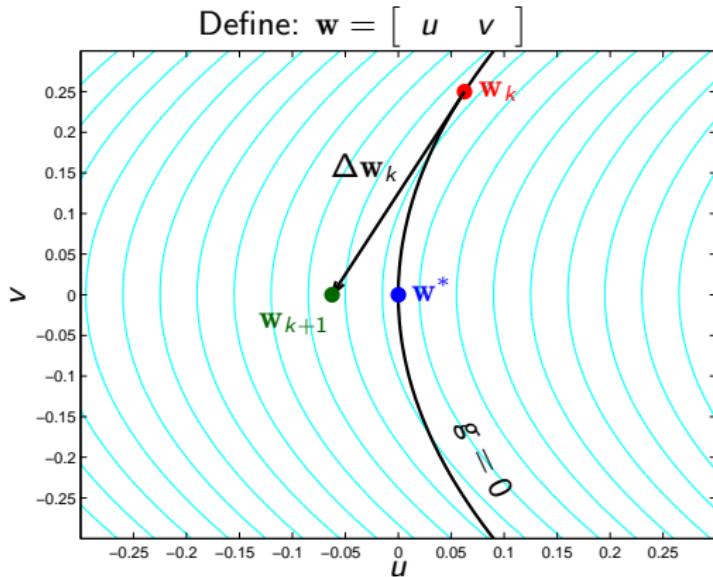
$$\Delta \mathbf{w}_k = - [ \begin{array}{cc} 2a^2 & a \end{array} ]$$

for  $\lambda = 2$ ...

The full Newton step:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \Delta \mathbf{w}_k$$

is much closer to  $\mathbf{w}^*$  than  $\mathbf{w}_k$ .



But:

$$\Phi(\mathbf{w}_{k+1}) > \Phi(\mathbf{w}_k)$$

$$|g(\mathbf{w}_{k+1})| > |g(\mathbf{w}_k)|$$

No penalty function can accept  $\Delta \mathbf{w}_k$  !!

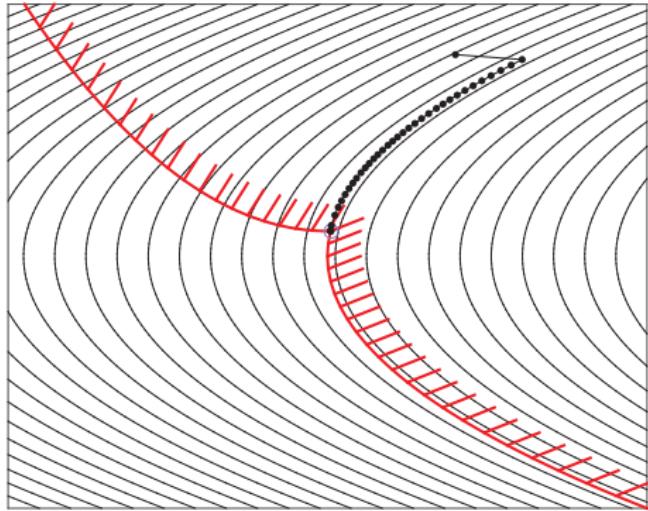
## Maratos effect - Illustration

Some NLPs can yield  
"creeping" convergence

NLP:

$$\min_w \Phi(w)$$

$$\text{s.t. } h(w) \leq 0$$



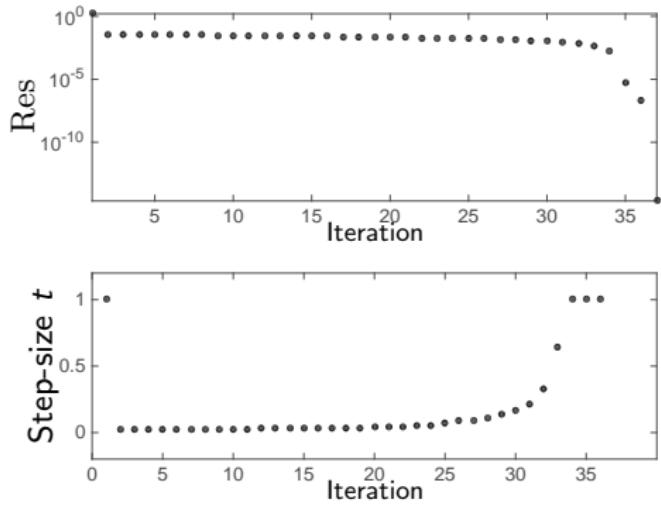
## Maratos effect - Illustration

Some NLPs can yield  
"creeping" convergence

NLP:

$$\min_w \Phi(w)$$

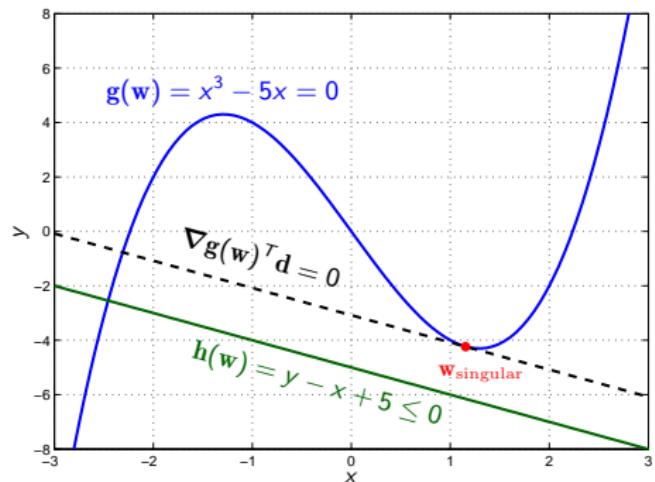
$$\text{s.t. } h(w) \leq 0$$



$$\text{where Res} = \left\| \frac{\nabla \mathcal{L}}{\max(h, 0)} \right\|_{\infty}$$

## Failure of the Newton methods - Convergence to infeasible points

Problem:



$$\begin{array}{ll} \min_w & x^2 + y^2 \\ \text{s.t.} & x^3 - 5x = 0 \\ & y - x + 5 \leq 0 \end{array}$$

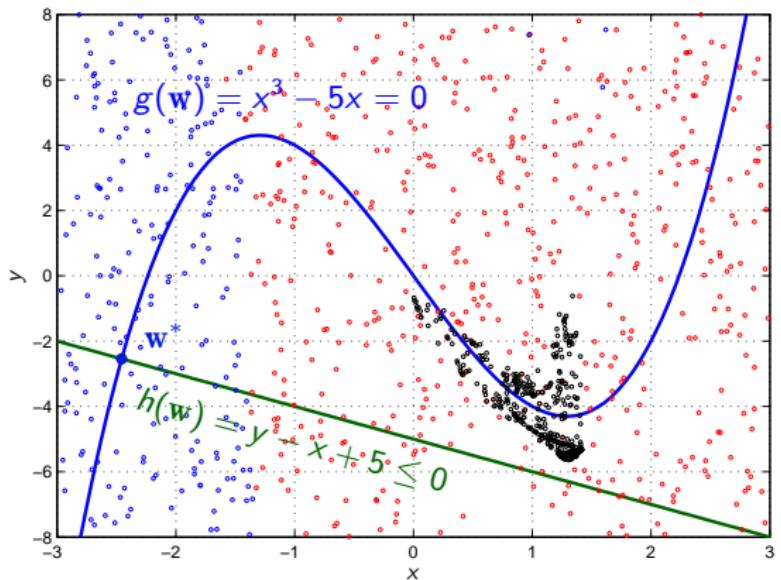
Vectors  $\nabla g(w)$  and  $\nabla h(w)$  are not Linearly Independent for some  $w$ , i.e.

$\nexists \Delta w$  such that:

$$\begin{bmatrix} \nabla g(w)^\top \\ \nabla h(w)^\top \end{bmatrix} \Delta w + \begin{bmatrix} g(w) \\ h(w) \end{bmatrix} = 0$$

There is no feasible direction  $\Delta w$

## Failure of the Newton methods - Convergence to infeasible points



Problem:

$$\begin{array}{ll} \min_{x,y} & x^2 + y^2 \\ \text{s.t.} & x^3 - 5x = 0 \\ & y - x + 5 \leq 0 \end{array}$$

Red dots: failed starting points  $\rightarrow$  black dots

Blue dots: successfull starting points  $\rightarrow$   $(-2.46, -2.54)$