

### Choose your path wisely

#### Gradient descent in a Bregman distance framework

Martin Benning<sup>1</sup> Marta M. Betcke<sup>2</sup> Matthias J. Ehrhardt<sup>1</sup> Carola-Bibiane Schönlieb<sup>1</sup>

<sup>1</sup>University of Cambridge, UK <sup>2</sup>University College London, UK

INI-SLB Workshop on large scale optimisation, Schlumberger Cambridge Research, 28.09.2017



Properties of the linearised Bregman iteration

Global convergence of the linearised Bregman iteration

Examples

Conclusions & outlook



# Outline

### Introduction

Properties of the linearised Bregman iteration

Global convergence of the linearised Bregman iteration

#### Examples

Conclusions & outlook



### Gradient descent

- We are interested in finding stationary points of smooth but not necessarily convex functionals E : U → R over Banach spaces U.
- A standard numerical approach for finding these in finite dimensions or Hilbert spaces is gradient descent, i.e.

$$u^{k+1} = u^k - \tau^k \nabla E(u^k),$$

for all  $k \in \mathbb{N}$  and some initial value  $u^0$ .

How can this be extended to more general Banach spaces? And why would we be interested in such an extension?



#### Gradient descent

- We are interested in finding stationary points of smooth but not necessarily convex functionals E : U → R over Banach spaces U.
- A standard numerical approach for finding these in finite dimensions or Hilbert spaces is gradient descent, i.e.

$$u^{k+1} = u^k - \tau^k \nabla E(u^k),$$

for all  $k \in \mathbb{N}$  and some initial value  $u^0$ .

How can this be extended to more general Banach spaces? And why would we be interested in such an extension?



### Gradient descent

- We are interested in finding stationary points of smooth but not necessarily convex functionals E : U → R over Banach spaces U.
- A standard numerical approach for finding these in finite dimensions or Hilbert spaces is gradient descent, i.e.

$$u^{k+1} = u^k - \tau^k \nabla E(u^k),$$

for all  $k \in \mathbb{N}$  and some initial value  $u^0$ .

How can this be extended to more general Banach spaces? And why would we be interested in such an extension?



Imagine non-convex problems such as blind deconvolution, i.e.

$$(\hat{u}, \hat{h}) \in \underset{\substack{u \in \mathbb{R}^n \\ h \in \mathbb{R}^r}}{\arg\min} \left\{ \frac{1}{2} \| u * h - f \|_2^2 + \alpha \mathsf{TV}(u) + \chi_{\mathcal{S}}(h) \right\}.$$
(1)

Here \* denotes the discrete two-dimensional convolution,  $\mathsf{TV}(u) := \| \| \nabla u \|_2 \|_1$  the total variation, and  $\chi_S(h) = \begin{cases} 0 & h \in S \\ \infty & \text{else} \end{cases}$  the

characteristic function over the simplex constraint

$$S = \left\{ h \in \mathbb{R}^r \; \left| \; h_j \geq 0, \; \sum_j h_j = 1, \; orall j 
ight\} 
ight.$$



Imagine non-convex problems such as blind deconvolution, i.e.

$$(\hat{u}, \hat{h}) \in \underset{\substack{u \in \mathbb{R}^n \\ h \in \mathbb{R}^r}}{\arg\min} \left\{ \frac{1}{2} \| u * h - f \|_2^2 + \alpha \mathsf{TV}(u) + \chi_{\mathcal{S}}(h) \right\}.$$
(1)

Here \* denotes the discrete two-dimensional convolution,  $\mathsf{TV}(u) := \| \| \nabla u \|_2 \|_1$  the total variation, and  $\chi_S(h) = \begin{cases} 0 & h \in S \\ \infty & \text{else} \end{cases}$  the characteristic function over the simplex constraint

$$S = \left\{ h \in \mathbb{R}^r \; \left| \; h_j \geq 0, \; \sum_j h_j = 1, \; orall j 
ight\} 
ight.$$



### We discover several issues here:

- Problem (1) is non-smooth  $\Rightarrow$  gradient descent not applicable
- Remedy: use forward backward splitting [Lions, Mercier, 1979]
- But: real issue is that total variation of blurred images is small



TV = 22.46 TV = 6.03 TV = 12.04▶ How do we prevent an optimisation procedure from picking a blurry solution?



### We discover several issues here:

- Problem (1) is non-smooth  $\Rightarrow$  gradient descent not applicable
- Remedy: use forward backward splitting [Lions, Mercier, 1979]
- But: real issue is that total variation of blurred images is small



TV = 22.46 TV = 6.03 TV = 12.04▶ How do we prevent an optimisation procedure from picking a blurry solution?



#### We discover several issues here:

- Problem (1) is non-smooth  $\Rightarrow$  gradient descent not applicable
- ▶ Remedy: use forward backward splitting [Lions, Mercier, 1979]
- But: real issue is that total variation of blurred images is small



TV = 22.46 TV = 6.03 TV = 12.04
▶ How do we prevent an optimisation procedure from picking a blurry solution?



#### We discover several issues here:

- Problem (1) is non-smooth  $\Rightarrow$  gradient descent not applicable
- ▶ Remedy: use forward backward splitting [Lions, Mercier, 1979]
- But: real issue is that total variation of blurred images is small



TV = 22.46

TV = 6.03

TV = 12.04

How do we prevent an optimisation procedure from picking a blurry solution?



#### We discover several issues here:

- Problem (1) is non-smooth  $\Rightarrow$  gradient descent not applicable
- ▶ Remedy: use forward backward splitting [Lions, Mercier, 1979]
- But: real issue is that total variation of blurred images is small



TV = 22.46 TV = 10.79 TV = 2.41
 ▶ How do we prevent an optimisation procedure from picking a blurry solution?



#### We discover several issues here:

- Problem (1) is non-smooth  $\Rightarrow$  gradient descent not applicable
- ▶ Remedy: use forward backward splitting [Lions, Mercier, 1979]
- But: real issue is that total variation of blurred images is small



TV = 22.46 TV = 10.79 TV = 2.41
▶ How do we prevent an optimisation procedure from picking a blurry solution?



Gradient descent cast as an optimisation problem

Gradient descent can be cast as

$$u^{k+1} = \arg\min_{u} \left\{ \tau^k \langle \nabla E(u^k), u - u^k \rangle + \frac{1}{2} \|u - u^k\|^2 \right\},$$

where  $\|\cdot\| = \sqrt{\langle \cdot, \cdot \rangle}$  is the norm of some Hilbert space.



### Linearised Bregman iteration for non-convex optimisation

Gradient descent can be cast as

$$u^{k+1} = \arg\min_{u} \left\{ \tau^{k} \langle \nabla E(u^{k}), u - u^{k} \rangle + D_{J}^{p^{k}}(u, u^{k}) \right\},$$
  
$$p^{k+1} = p^{k} - \tau^{k} \nabla E(u^{k}),$$

where  $D_J^{p^k}(u^{k+1}, u^k) := J(u^{k+1}) - J(u^k) - \langle p^k, u^{k+1} - u^k \rangle$  is the (generalised) Bregman distance for  $p^k \in \partial J(u^k)$ .



### Linearised Bregman iteration for non-convex optimisation

Gradient descent can be cast as

$$u^{k+1} = \arg\min_{u} \left\{ \tau^{k} \langle \nabla E(u^{k}), u - u^{k} \rangle + D_{J}^{p^{k}}(u, u^{k}) \right\},$$
  
$$p^{k+1} = p^{k} - \tau^{k} \nabla E(u^{k}),$$

where  $D_J^{p^k}(u^{k+1}, u^k) := J(u^{k+1}) - J(u^k) - \langle p^k, u^{k+1} - u^k \rangle$  is the (generalised) Bregman distance for  $p^k \in \partial J(u^k)$ .

▶ Here  $J : U \to \mathbb{R} \cup \{\infty\}$  is proper, lower semi-continuous (l.s.c) and convex, and

$$\partial J(u) := \{ p \in \mathcal{U}^* \mid D^p_J(v, u) \ge 0, \ \forall v \in \mathcal{U} \} \subset \mathcal{U}^* \ .$$



We use the linearised Bregman iteration to (approximately) minimise

$$E(u,h) = \frac{1}{2} ||u * h - f||_2^2,$$

where \* denotes the two-dimensional, discrete convolution. We make the following choice for J:

$$J(u,h) := \frac{1}{2} \|u\|_2^2 + \alpha \mathsf{TV}(u) + \sum_{j=1}^r h_j \log(h_j) - h_j + \chi_S(h)$$
  
with  $\chi_S(h) = \begin{cases} 0 & h \in S \\ \infty & \text{else} \end{cases}$  and the simplex constraint set  
 $S = \begin{cases} h \in \mathbb{R}^r & h_j \ge 0, \ \sum_j h_j = 1, \ \forall j \end{cases}$ .



We use the linearised Bregman iteration to (approximately) minimise

$$E(u,h) = \frac{1}{2} ||u * h - f||_2^2,$$

where \* denotes the two-dimensional, discrete convolution. We make the following choice for *J*:

$$J(u,h) := \frac{1}{2} \|u\|_2^2 + \alpha \mathsf{TV}(u) + \sum_{j=1}^r h_j \log(h_j) - h_j + \chi_S(h)$$
  
with  $\chi_S(h) = \begin{cases} 0 & h \in S \\ \infty & \text{else} \end{cases}$  and the simplex constraint set  
 $S = \begin{cases} h \in \mathbb{R}^r & h_j \ge 0, \ \sum_j h_j = 1, \ \forall j \end{cases}$ .



Note that the choice of J results in the following algorithm:

$$u^{k+1} = \arg\min_{u} \left\{ \frac{1}{2} \left\| u - (u^{k} - (u^{k} * h^{k} - f) \bar{*} h^{k} + \alpha q^{k}) \right\|_{2}^{2} + \alpha \mathsf{TV}(u) \right\},$$
  

$$h_{j}^{k+1} = \frac{h_{j}^{k} \exp(-\tau^{k} (u^{k} \bar{*} (u^{k} * h^{k} - f))_{j})}{\sum_{j=1}^{r} h_{j}^{k} \exp(-\tau^{k} (u^{k} \bar{*} (u^{k} * h^{k} - f))_{j})} \quad \forall j \in \{1, \dots, r\}, \ (\dagger)$$
  

$$q^{k+1} = q^{k} - \frac{1}{\alpha} \left( u^{k+1} - u^{k} + (u^{k} * h^{k} - f) \bar{*} h^{k} \right), \quad q^{k} \in \partial \mathsf{TV}(u^{k});$$

here  $\overline{*}$  denotes the transpose convolution operation.

(†) is basically entropic descent, see Nemirovsky, Yudin, 1982 & Beck, Teboulle, Operations Research Letters, 2003



Note that the choice of J results in the following algorithm:

$$u^{k+1} = \arg\min_{u} \left\{ \frac{1}{2} \left\| u - (u^{k} - (u^{k} * h^{k} - f) \bar{*} h^{k} + \alpha q^{k}) \right\|_{2}^{2} + \alpha \mathsf{TV}(u) \right\},$$
  

$$h_{j}^{k+1} = \frac{h_{j}^{k} \exp(-\tau^{k} (u^{k} \bar{*} (u^{k} * h^{k} - f))_{j})}{\sum_{j=1}^{r} h_{j}^{k} \exp(-\tau^{k} (u^{k} \bar{*} (u^{k} * h^{k} - f))_{j})} \quad \forall j \in \{1, \dots, r\}, \ (\dagger)$$
  

$$q^{k+1} = q^{k} - \frac{1}{\alpha} \left( u^{k+1} - u^{k} + (u^{k} * h^{k} - f) \bar{*} h^{k} \right), \quad q^{k} \in \partial \mathsf{TV}(u^{k});$$

here  $\overline{*}$  denotes the transpose convolution operation.

(†) is basically entropic descent, see Nemirovsky, Yudin, 1982 & Beck, Teboulle, Operations Research Letters, 2003



### We consider the same (noisy) motion blur example as before:





Reconstruction via linearised Bregman iteration



#### A non-exhaustive list of works on Bregman iteration:

- ► Censor, Zenios. Journal of Opt. Theory and App. 1992
- ► Teboulle. Mathematics of Operations Research 1992
- Eckstein. Mathematics of Operations Research 1993
- ► Kiwiel. SIAM Journal on Control and Optimization 1997
- ▶ Beck, Teboulle. Operations Research Letter 2003
- ▶ Osher, Burger, Goldfarb, Xu, Yin. Multiscale Mod. & Sim. 2005
- ► Yin, Osher, Goldfarb, Darbon. SIAM Journal on Imaging Sci. 2008
- ► Cai, Osher, Shen. Mathmatics of Computation 2009
- Bachmayr, Burger. Inverse Problems 2009
- Yin. SIAM Journal on Imaging Sciences 2010





### Properties of the linearised Bregman iteration

Global convergence of the linearised Bregman iteration

#### Examples

Conclusions & outlook



### Linearised Bregman iteration

$$u^{k+1} = \arg\min_{u \in \mathcal{U}} \left\{ \tau^k \langle \nabla E(u^k), u - u^k \rangle + D_J^{p^k}(u, u^k) \right\},$$
  
$$p^{k+1} = p^k - \tau^k \nabla E(u^k),$$

for  $k \in \mathbb{N}$ , and some  $u^0 \in \mathcal{U}$  and  $p^0 \in \partial J(u^0)$ .

Note that this Algorithm can be simplified to

$$u^{k+1} = \operatorname*{arg\,min}_{u \in \mathcal{U}} \left\{ \tau^k \left\langle \nabla E(u^k) - \frac{1}{\tau^k} p^k, u - u^k \right\rangle + J(u) \right\},$$
$$p^{k+1} = p^k - \tau^k \nabla E(u^k).$$



### Linearised Bregman iteration

$$u^{k+1} = \arg\min_{u \in \mathcal{U}} \left\{ \tau^k \langle \nabla E(u^k), u - u^k \rangle + D_J^{p^k}(u, u^k) \right\},$$
  
$$p^{k+1} = p^k - \tau^k \nabla E(u^k),$$

for  $k \in \mathbb{N}$ , and some  $u^0 \in \mathcal{U}$  and  $p^0 \in \partial J(u^0)$ .

Note that this Algorithm can be simplified to

$$u^{k+1} = \operatorname*{arg\,min}_{u \in \mathcal{U}} \left\{ \tau^k \left\langle \nabla E(u^k) - \frac{1}{\tau^k} p^k, u - u^k \right\rangle + J(u) \right\},$$
  
$$p^{k+1} = p^k - \tau^k \nabla E(u^k).$$



Specialised linearised Bregman iteration (I)

Specific choice  $J^k(u) = \frac{1}{2} ||u||_2^2 + \tau^k \alpha R(u)$ :

$$\begin{split} u^{k+1} &= \arg\min_{u \in \mathcal{U}} \left\{ \frac{1}{2} \| u - (u^k - \tau^k \nabla E(u^k)) \|_2^2 + \tau^k \alpha D_R^{q^k}(u, u^k) \right\} ,\\ &= \arg\min_{u \in \mathcal{U}} \left\{ \frac{1}{2} \| u - (u^k + \tau^k (\alpha q^k - \nabla E(u^k))) \|_2^2 + \tau^k \alpha R(u) \right\} ,\\ &= (I + \alpha \tau^k \partial R)^{-1} (u^k + \tau^k (\alpha q^k - \nabla E(u^k)))\\ q^{k+1} &= q^k - \frac{1}{\tau^k \alpha} \left( u^{k+1} - u^k + \tau^k \nabla E(u^k) \right) , \qquad q^k \in \partial R(u^k) \end{split}$$



### Specialised linearised Bregman iteration (II)

More general specific choice  $J^k(u) = H(u) + \tau^k \alpha R(u)$ :

$$u^{k+1} = \arg\min_{u \in \mathcal{U}} \left\{ H(u) - \langle u, \nabla H(u^k) - \tau^k \nabla E(u^k) \rangle + \tau^k \alpha D_R^{q^k}(u, u^k) \right\},$$
  

$$= \arg\min_{u \in \mathcal{U}} \left\{ H(u) - \langle u, \nabla H(u^k) + \tau^k (\alpha q^k - \nabla E(u^k)) \rangle + \tau^k \alpha R(u) \right\},$$
  

$$= (\nabla H + \alpha \tau^k \partial R)^{-1} (\nabla H^* (\nabla H(u^k) + \tau^k (\alpha q^k - \nabla E(u^k))))$$
  

$$(2a)$$
  

$$q^{k+1} = q^k - \frac{1}{\tau^k \alpha} \left( \nabla H(u^{k+1}) - \nabla H(u^k) + \tau^k \nabla E(u^k) \right), q^k \in \partial R(u^k)$$
  

$$(2b)$$



•  $E: \mathcal{U} \to \mathbb{R}$  is a smooth functional that is bounded from below and has *L*-Lipschitz-continuous gradient, i.e.

 $\|\nabla E(u) - \nabla E(v)\|_{\mathcal{U}^*} \leq L \|u - v\|_{\mathcal{U}}, \qquad \forall u, v \in \mathcal{U}.$ 

•  $H: \mathcal{U} \to \mathbb{R}$  is a  $\gamma$ -strongly convex functional, i.e.

 $\gamma \|u - v\|_{\mathcal{U}}^2 \leq \boldsymbol{D}_{\boldsymbol{H}}^{\mathsf{symm}}(u, v) := D_{\boldsymbol{H}}(u, v) + D_{\boldsymbol{H}}(v, u), \quad \forall u, v \in \mathcal{U},$ 

and has  $\delta$ -Lipschitz gradient  $\nabla H$ .

▶ The conditions on *E* and *H* imply the useful estimate

$$E(u) \leq E(v) + \langle \nabla E(v), u - v \rangle + \frac{L}{\gamma} D_H(u, v).$$
 (3)



•  $E: \mathcal{U} \to \mathbb{R}$  is a smooth functional that is bounded from below and has *L*-Lipschitz-continuous gradient, i.e.

 $\|\nabla E(u) - \nabla E(v)\|_{\mathcal{U}^*} \leq L \|u - v\|_{\mathcal{U}}, \qquad \forall u, v \in \mathcal{U}.$ 

•  $H: \mathcal{U} \to \mathbb{R}$  is a  $\gamma$ -strongly convex functional, i.e.

 $\gamma \| u - v \|_{\mathcal{U}}^2 \leq D_H^{\text{symm}}(u, v) := D_H(u, v) + D_H(v, u), \quad \forall u, v \in \mathcal{U},$ 

and has  $\delta$ -Lipschitz gradient  $\nabla H$ .

▶ The conditions on *E* and *H* imply the useful estimate

$$E(u) \le E(v) + \langle \nabla E(v), u - v \rangle + \frac{L}{\gamma} D_H(u, v).$$
 (3)



 E: U → ℝ is a smooth functional that is bounded from below and has L-Lipschitz-continuous gradient, i.e.

$$\|
abla E(u) - 
abla E(v)\|_{\mathcal{U}^*} \leq L \|u - v\|_{\mathcal{U}}, \qquad \forall u, v \in \mathcal{U} \;.$$

•  $H: \mathcal{U} \to \mathbb{R}$  is a  $\gamma$ -strongly convex functional, i.e.

$$\gamma \| u - v \|_{\mathcal{U}}^2 \leq \mathcal{D}_H^{\mathsf{symm}}(u, v) := \mathcal{D}_H(u, v) + \mathcal{D}_H(v, u), \quad \forall u, v \in \mathcal{U},$$

and has  $\delta$ -Lipschitz gradient  $\nabla H$ .

▶ The conditions on *E* and *H* imply the useful estimate

$$E(u) \leq E(v) + \langle \nabla E(v), u - v \rangle + \frac{L}{\gamma} D_H(u, v).$$
 (3)



•  $E: \mathcal{U} \to \mathbb{R}$  is a smooth functional that is bounded from below and has *L*-Lipschitz-continuous gradient, i.e.

 $\|\nabla E(u) - \nabla E(v)\|_{\mathcal{U}^*} \leq L \|u - v\|_{\mathcal{U}}, \qquad \forall u, v \in \mathcal{U}.$ 

•  $H: \mathcal{U} \to \mathbb{R}$  is a  $\gamma$ -strongly convex functional, i.e.

 $\gamma \| u - v \|_{\mathcal{U}}^2 \leq D_H^{\text{symm}}(u, v) := D_H(u, v) + D_H(v, u), \quad \forall u, v \in \mathcal{U},$ 

and has  $\delta$ -Lipschitz gradient  $\nabla H$ .

▶ The conditions on *E* and *H* imply the useful estimate

$$E(u) \leq E(v) + \langle \nabla E(v), u - v \rangle + \frac{L}{\gamma} D_H(u, v).$$
 (3)



### Surrogate objective

For the iterate  $u^{k-1} \in \mathcal{U}$  of (2) and corresponding subgradient  $q^{k-1} \in \partial R(u^{k-1})$  we define  $E^k : \mathcal{U} \to \mathbb{R} \cup \{\infty\}$  as

$$E^{k}(u) := E(u) + \alpha D_{R}^{q^{k-1}}(u, u^{k-1}).$$
(4)

#### Stepsize constraint

Choose  $0 < \tau^k$  such that

$$D_H(u^{k+1}, u^k) \le \frac{\gamma(1 - \tau^k \rho)}{\tau^k L} D_H^{\mathsf{symm}}(u^{k+1}, u^k) \tag{5}$$

holds true for all  $k \in \mathbb{N}$  and  $\mathsf{0} < \delta < 1/$  max $_k au^k$  .



### Surrogate objective

For the iterate  $u^{k-1} \in \mathcal{U}$  of (2) and corresponding subgradient  $q^{k-1} \in \partial R(u^{k-1})$  we define  $E^k : \mathcal{U} \to \mathbb{R} \cup \{\infty\}$  as

$$E^{k}(u) := E(u) + \alpha D_{R}^{q^{k-1}}(u, u^{k-1}).$$
(4)

#### Stepsize constraint

Choose  $0 < \tau^k$  such that

$$D_H(u^{k+1}, u^k) \le \frac{\gamma(1 - \tau^k \rho)}{\tau^k L} D_H^{\text{symm}}(u^{k+1}, u^k)$$
(5)

holds true for all  $k \in \mathbb{N}$  and  $0 < \delta < 1/\max_k \tau^k$ .



### Surrogate objective

For the iterate  $u^{k-1} \in \mathcal{U}$  of (2) and corresponding subgradient  $q^{k-1} \in \partial R(u^{k-1})$  we define  $E^k : \mathcal{U} \to \mathbb{R} \cup \{\infty\}$  as

$$E^{k}(u) := E(u) + \alpha D_{R}^{q^{k-1}}(u, u^{k-1}).$$
(4)

### Stepsize constraint

Choose  $0 < \tau^k$  such that

$$\frac{1}{2} \|u^{k+1} - u^k\|_2^2 \le \left(\frac{1 - \tau^k \rho}{\tau^k L}\right) \|u^{k+1} - u^k\|_2^2 \tag{5}$$

holds true for all  $k \in \mathbb{N}$  and  $0 < \delta < 1/\max_k \tau^k$ .



### Surrogate objective

For the iterate  $u^{k-1} \in \mathcal{U}$  of (2) and corresponding subgradient  $q^{k-1} \in \partial R(u^{k-1})$  we define  $E^k : \mathcal{U} \to \mathbb{R} \cup \{\infty\}$  as

$$E^{k}(u) := E(u) + \alpha D_{R}^{q^{k-1}}(u, u^{k-1}).$$
(4)

#### Stepsize constraint

Choose  $0 < \tau^k$  such that

$$D_H(u^{k+1}, u^k) \le \frac{\gamma(1 - \tau^k \rho)}{\tau^k L} D_H^{\text{symm}}(u^{k+1}, u^k)$$
(5)

holds true for all  $k \in \mathbb{N}$  and  $0 < \delta < 1/\max_k \tau^k$ .


### Sufficient decrease property

Let  $0 < \tau^k < 1/\delta$  satisfy (5). Then the iterates of the linearised Bregman iteration (2) satisfy

$$E^{k+1}(u^{k+1}) + \rho D_{H}^{\text{symm}}(u^{k+1}, u^{k}) + \alpha \left( D_{R}^{q^{k+1}}(u^{k}, u^{k+1}) + D_{R}^{q^{k-1}}(u^{k}, u^{k-1}) \right)$$
(6)  
$$\leq E^{k}(u^{k}).$$

In addition, we observe

$$\lim_{k\to\infty} D_H^{\mathsf{symm}}(u^{k+1},u^k) = 0 \quad \text{and} \quad \lim_{k\to\infty} D_R^{\mathsf{symm}}(u^{k+1},u^k) = 0 \,.$$



### Sufficient decrease property

Let  $0 < \tau^k < 1/\delta$  satisfy (5). Then the iterates of the linearised Bregman iteration (2) satisfy

$$E^{k+1}(u^{k+1}) + \rho \| u^{k+1} - u^{k} \|_{2}^{2} + \alpha \left( D_{R}^{q^{k+1}}(u^{k}, u^{k+1}) + D_{R}^{q^{k-1}}(u^{k}, u^{k-1}) \right)$$
(6)  
$$\leq E^{k}(u^{k}).$$

In addition, we observe

$$\lim_{k\to\infty} \|\boldsymbol{u}^{k+1} - \boldsymbol{u}^k\|_2^2 = 0 \quad \text{and} \quad \lim_{k\to\infty} D_R^{\text{symm}}(\boldsymbol{u}^{k+1}, \boldsymbol{u}^k) = 0.$$



### Sufficient decrease property

Let  $0 < \tau^k < 1/\delta$  satisfy (5). Then the iterates of the linearised Bregman iteration (2) satisfy

$$E^{k+1}(u^{k+1}) + \rho D_{H}^{\text{symm}}(u^{k+1}, u^{k}) + \alpha \left( D_{R}^{q^{k+1}}(u^{k}, u^{k+1}) + D_{R}^{q^{k-1}}(u^{k}, u^{k-1}) \right)$$
(6)  
$$\leq E^{k}(u^{k}).$$

In addition, we observe

$$\lim_{k\to\infty} D_H^{\mathsf{symm}}(u^{k+1},u^k) = 0 \quad \text{and} \quad \lim_{k\to\infty} D_R^{\mathsf{symm}}(u^{k+1},u^k) = 0 \,.$$



#### A subgradient lower bound of the iterates gap

Let the same assumptions hold true as in the previous theorem. Then the iterates (2) satisfy

$$\|\nabla E(u^{k+1}) + \alpha(q^{k+1} - q^k)\|_{\mathcal{U}^*} \le \rho_2 \|u^{k+1} - u^k\|_{\mathcal{U}},$$
(7)

for  $q^k \in \partial R(u^k)$ ,  $q^{k+1} \in \partial R(u^{k+1})$ ,  $\rho_2 := (L + \delta/(\min_k \tau^k))$  and  $k \in \mathbb{N}$ .

Proof: Use Lipschitz-continuity of  $\nabla E$ , Lipschitz-continuity of  $\nabla H$  and the triangle inequality for norms.



#### A subgradient lower bound of the iterates gap

Let the same assumptions hold true as in the previous theorem. Then the iterates (2) satisfy

$$\|\nabla E(u^{k+1}) + \alpha(q^{k+1} - q^k)\|_{\mathcal{U}^*} \le \rho_2 \|u^{k+1} - u^k\|_{\mathcal{U}},$$
(7)

for  $q^k \in \partial R(u^k)$ ,  $q^{k+1} \in \partial R(u^{k+1})$ ,  $\rho_2 := (L + \delta/(\min_k \tau^k))$  and  $k \in \mathbb{N}$ .

Proof: Use Lipschitz-continuity of  $\nabla E$ , Lipschitz-continuity of  $\nabla H$  and the triangle inequality for norms.





Introduction

Properties of the linearised Bregman iteration

#### Global convergence of the linearised Bregman iteration

Examples

Conclusions & outlook



## Convergence of the surrogate functional

## Definition (Set of limiting points)

From now on we set  $\mathcal{U} = \mathbb{R}^n$ . The set of limiting points is defined as  $\omega(u^0) := \left\{ \overline{u} \in \mathbb{R}^n \ \Big| \ \text{there exists an increasing sequence of integers} \\ \{k_j\}_{j \in \mathbb{N}} \ \text{such that } \lim_{j \to \infty} u^{k_j} = \overline{u} \right\}.$ 

Lemma (Convergence of the surrogate functional

Let  $\overline{u} \in \omega(u^0)$ . Then we obtain

$$\lim_{k \to \infty} E^k(u^k) = E(\overline{u}).$$
(8)



## Convergence of the surrogate functional

## Definition (Set of limiting points)

From now on we set  $\mathcal{U} = \mathbb{R}^n$ . The set of limiting points is defined as  $\omega(u^0) := \left\{ \overline{u} \in \mathbb{R}^n \mid \text{there exists an increasing sequence of integers} \\ \{k_j\}_{j \in \mathbb{N}} \text{ such that } \lim_{j \to \infty} u^{k_j} = \overline{u} \right\}.$ 

Lemma (Convergence of the surrogate functional)

Let  $\overline{u} \in \omega(u^0)$ . Then we obtain

$$\lim_{k \to \infty} E^k(u^k) = E(\overline{u}).$$
(8)



### Theorem (Finite length property)

Let  $\mathcal{U} = \mathbb{R}^n$  be finite dimensional, and let the sequences  $\{u^k\}_{k \in \mathbb{N}}$  and  $\{q^k\}_{k \in \mathbb{N}}$  generated by Algorithm (2) be bounded. Then,  $\{u^k\}_{k \in \mathbb{N}}$ 

$$\sum_{k=0}^{\infty} \|u^{k+1} - u^k\|_{\mathcal{U}} < \infty ..$$
(9)

2. converges to a critical point  $\hat{u}$  of E, i.e.  $\nabla E(\hat{u}) = 0$ , with  $\hat{\rho} \in \partial J(\hat{u})$ .



### Theorem (Finite length property)

Let  $\mathcal{U} = \mathbb{R}^n$  be finite dimensional, and let the sequences  $\{u^k\}_{k \in \mathbb{N}}$  and  $\{q^k\}_{k \in \mathbb{N}}$  generated by Algorithm (2) be bounded. Then,  $\{u^k\}_{k \in \mathbb{N}}$ 1. has finite length, i.e.

$$\sum_{k=0}^{\infty} \|u^{k+1} - u^k\|_{\mathcal{U}} < \infty ..$$
 (9)

2. converges to a critical point  $\hat{u}$  of E, i.e.  $\nabla E(\hat{u}) = 0$ , with  $\hat{\rho} \in \partial J(\hat{u})$ .



### Theorem (Finite length property)

Let  $\mathcal{U} = \mathbb{R}^n$  be finite dimensional, and let the sequences  $\{u^k\}_{k \in \mathbb{N}}$  and  $\{q^k\}_{k \in \mathbb{N}}$  generated by Algorithm (2) be bounded. Then,  $\{u^k\}_{k \in \mathbb{N}}$ 

1. has finite length, i.e.

$$\sum_{k=0}^{\infty} \|u^{k+1} - u^k\|_{\mathcal{U}} < \infty ..$$
 (9)

2. converges to a critical point  $\hat{u}$  of E, i.e.  $\nabla E(\hat{u}) = 0$ , with  $\hat{p} \in \partial J(\hat{u})$ .



### Theorem (Finite length property)

Let  $\mathcal{U} = \mathbb{R}^n$  be finite dimensional, and let the sequences  $\{u^k\}_{k \in \mathbb{N}}$  and  $\{q^k\}_{k \in \mathbb{N}}$  generated by Algorithm (2) be bounded. Then,  $\{u^k\}_{k \in \mathbb{N}}$ 

1. has finite length, i.e.

$$\sum_{k=0}^{\infty} \|u^{k+1} - u^k\|_{\mathcal{U}} < \infty ..$$
 (9)

2. converges to a critical point  $\hat{u}$  of E, i.e.  $\nabla E(\hat{u}) = 0$ , with  $\hat{p} \in \partial J(\hat{u})$ .



If  $\{q^k\}_{k\in\mathbb{N}}$  is unbounded we can find the following counter example for the previous theorem:



If  $\{q^k\}_{k\in\mathbb{N}}$  is unbounded we can find the following counter example for the previous theorem:



If  $\{q^k\}_{k\in\mathbb{N}}$  is unbounded we can find the following counter example for the previous theorem:



If  $\{q^k\}_{k\in\mathbb{N}}$  is unbounded we can find the following counter example for the previous theorem:

► 
$$E(u) := \frac{1}{2}(u+1)^2$$
►  $H(u) := \frac{1}{2}u^2$ 
►  $R(u) := \chi_{\geq 0}(x) := \begin{cases} 0 & x \geq 0 \\ \infty & x < 0 \end{cases}$ 
►  $u^0 > 0, \ q^0 = 0$ 
⇒  $\lim_{k \to \infty} u^k = 0 \neq -1 = \hat{u}$  with  $E'(\hat{u}) = 0$ 
 $q^k \to -\infty$ 



If  $\{q^k\}_{k\in\mathbb{N}}$  is unbounded we can find the following counter example for the previous theorem:

► 
$$E(u) := \frac{1}{2}(u+1)^2$$
►  $H(u) := \frac{1}{2}u^2$ 
►  $R(u) := \chi_{\geq 0}(x) := \begin{cases} 0 & x \geq 0 \\ \infty & x < 0 \end{cases}$ 
►  $u^0 > 0, q^0 = 0$ 
⇒  $\lim_{k \to \infty} u^k = 0 \neq -1 = \hat{u}$  with  $E'(\hat{u}) = 0$ 
 $q^k \to -\infty$ 





Introduction

Properties of the linearised Bregman iteration

Global convergence of the linearised Bregman iteration

#### Examples

Conclusions & outlook



In parallel MRI one aims to minimise

$$E(u, b_1, \ldots, b_n) = \frac{1}{2} \sum_{j=1}^n \| (S(\mathcal{F}(K(u, b_1, \ldots, b_n))))_j - f_j \|_2^2.$$

- $\mathcal{F}$  is the (discrete) Fourier transform.
- ► *S* is a sub-sampling operator.
- K is the operator  $K(u, b_1, \ldots, b_2) = (ub_1, ub_2, \ldots, ub_n)^T$ , with
  - the spin-proton density u,
  - the coil sensitivities  $b_1, b_2, \ldots, b_n$ .
- sub-sampled k-space data  $f_1, \ldots, f_n$ .



In parallel MRI one aims to minimise

$$E(u, b_1, \ldots, b_n) = \frac{1}{2} \sum_{j=1}^n \| (S(\mathcal{F}(K(u, b_1, \ldots, b_n))))_j - f_j \|_2^2.$$

- $\mathcal{F}$  is the (discrete) Fourier transform.
- ► *S* is a sub-sampling operator.
- K is the operator  $K(u, b_1, \ldots, b_2) = (ub_1, ub_2, \ldots, ub_n)^T$ , with
  - the spin-proton density u,
  - the coil sensitivities  $b_1, b_2, \ldots, b_n$ .
- sub-sampled k-space data  $f_1, \ldots, f_n$ .



In parallel MRI one aims to minimise

$$E(u, b_1, \ldots, b_n) = \frac{1}{2} \sum_{j=1}^n \| (S(\mathcal{F}(K(u, b_1, \ldots, b_n))))_j - f_j \|_2^2.$$

- $\mathcal{F}$  is the (discrete) Fourier transform.
- ► *S* is a sub-sampling operator.
- K is the operator  $K(u, b_1, \ldots, b_2) = (ub_1, ub_2, \ldots, ub_n)^T$ , with
  - the spin-proton density u,
  - the coil sensitivities  $b_1, b_2, \ldots, b_n$ .
- sub-sampled k-space data  $f_1, \ldots, f_n$ .



In parallel MRI one aims to minimise

$$E(u, b_1, \ldots, b_n) = \frac{1}{2} \sum_{j=1}^n \|(S(\mathcal{F}(K(u, b_1, \ldots, b_n))))_j - f_j\|_2^2.$$

- $\mathcal{F}$  is the (discrete) Fourier transform.
- ► *S* is a sub-sampling operator.
- K is the operator  $K(u, b_1, \ldots, b_2) = (ub_1, ub_2, \ldots, ub_n)^T$ , with
  - ► the spin-proton density *u*,
  - the coil sensitivities  $b_1, b_2, \ldots, b_n$ .
- sub-sampled k-space data  $f_1, \ldots, f_n$ .



In parallel MRI one aims to minimise

$$E(u, b_1, \ldots, b_n) = \frac{1}{2} \sum_{j=1}^n \|(S(\mathcal{F}(K(u, b_1, \ldots, b_n))))_j - f_j\|_2^2.$$

- $\mathcal{F}$  is the (discrete) Fourier transform.
- S is a sub-sampling operator.
- K is the operator  $K(u, b_1, \ldots, b_2) = (ub_1, ub_2, \ldots, ub_n)^T$ , with
  - ► the spin-proton density *u*,
  - the coil sensitivities  $b_1, b_2, \ldots, b_n$ .
- sub-sampled k-space data  $f_1, \ldots, f_n$ .



## Parallel MRI



## Parallel MRI

Potential choice for H and R:

$$H(u, b_1, \dots, b_n) = H_0(u) + \sum_{k=1}^n H_k(b_k),$$

$$R(u, b_1, \dots, b_n) = R_0(u) + \sum_{k=1}^n R_k(b_k)$$
with
$$H_k(v) = \frac{1}{2} ||v||_2^2 \quad \forall k \in \{0, \dots, n\},$$
and
$$R_0(u) = \mathsf{TV}(u),$$

$$R_k(b_k) = ||Cb_k||_1 \quad \forall k \in \{1, \dots, n\}.$$

Here C denotes the discrete two-dimensional cosine transform.



W



### Classical gradient descent

Linearised Bregman iteration



We want to train the two-layer neural network

## $K(U_1, U_2) := \min(\max(U_1 \max(U_2 D, 0), 0), 1)$

where  $\rho_1$  and  $\rho_2$  are smooth approximations of min(max( $\cdot$ , 0), 1) and max( $\cdot$ , 0), respectively. Matrices  $U_1 \in \mathbb{R}^{10 \times 784}$  and  $U_2 \in \mathbb{R}^{784 \times 784}$  are the unknowns;  $D \in \mathbb{R}^{784 \times 50000}$  is data matrix that contains 50000 training images of MNIST (LeCun et al. 2010). We want to solve the inverse problem

 $K(U_1, U_2) = F,$ 

where  $F \in \{0,1\}^{10 imes 50000}$  is the labelling matrix.

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0</t



We want to train the two-layer neural network

$$K(U_1, U_2) := \rho_1(U_1 \rho_2(U_2 D))$$

where  $\rho_1$  and  $\rho_2$  are smooth approximations of min(max( $\cdot, 0$ ), 1) and max( $\cdot, 0$ ), respectively. Matrices  $U_1 \in \mathbb{R}^{10 \times 784}$  and  $U_2 \in \mathbb{R}^{784 \times 784}$  are the unknowns;  $D \in \mathbb{R}^{784 \times 50000}$  is data matrix that contains 50000 training images of MNIST (LeCun et al. 2010). We want to solve the inverse problem

 $K(U_1, U_2) = F,$ 

where  $F \in \{0,1\}^{10 imes 50000}$  is the labelling matrix.

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0



We want to train the two-layer neural network

$$K(U_1, U_2) := \rho_1(U_1 \rho_2(U_2 D))$$

where  $\rho_1$  and  $\rho_2$  are smooth approximations of min(max( $\cdot$ , 0), 1) and max( $\cdot$ , 0), respectively. Matrices  $U_1 \in \mathbb{R}^{10 \times 784}$  and  $U_2 \in \mathbb{R}^{784 \times 784}$  are the unknowns;  $D \in \mathbb{R}^{784 \times 50000}$  is data matrix that contains 50000 training images of MNIST (LeCun et al. 2010). We want to solve the inverse

problem

 $K(U_1, U_2) = F,$ 

where  $F \in \{0,1\}^{10 imes 50000}$  is the labelling matrix.

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0



We want to train the two-layer neural network

$$K(U_1, U_2) := \rho_1(U_1 \rho_2(U_2 D))$$

where  $\rho_1$  and  $\rho_2$  are smooth approximations of min(max( $\cdot$ , 0), 1) and max( $\cdot$ , 0), respectively. Matrices  $U_1 \in \mathbb{R}^{10 \times 784}$  and  $U_2 \in \mathbb{R}^{784 \times 784}$  are the unknowns;  $D \in \mathbb{R}^{784 \times 50000}$  is data matrix that contains 50000 training images of MNIST (LeCun et al. 2010). We want to solve the inverse problem

 $K(U_1, U_2) = F,$ 

where  $F \in \{0,1\}^{10 \times 50000}$  is the labelling matrix.

 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0



We compute  $U_1$  and  $U_2$  via (2) with  $H(U_1, U_2) = \sum_{i=1}^2 \|U_i\|_{\mathsf{Fro}}^2$  and

$$R(U_1, U_2) := \sum_{i=1}^2 \|U_i\|_* = \sum_{i=1}^{10} \sigma_{1,i} + \sum_{j=1}^{784} \sigma_{2,j}.$$

Here  $\{\sigma_{1,i}\}_{i \in \{1,...,10\}}$  and  $\{\sigma_{2,j}\}_{j \in \{1,...,784\}}$  denote the singular values of  $U_1$  and  $U_2$ , respectively.







Introduction

Properties of the linearised Bregman iteration

Global convergence of the linearised Bregman iteration

Examples

Conclusions & outlook



## Conclusions

- We have generalised gradient descent to Banach spaces
- We have proven global convergence for this generalisation.
- We shown different solution paths for same initialisation.

#### Outlook

- ▶ Convergence analysis for unbounded {q<sup>k</sup>}<sub>k∈N</sub>
- Variable splitting, BFGS-style updates, Newton-type methods
- Incremental or even stochastic versions



## Conclusions

- We have generalised gradient descent to Banach spaces
- ► We have proven global convergence for this generalisation.
- We shown different solution paths for same initialisation.

#### Outlook

- ▶ Convergence analysis for unbounded {q<sup>k</sup>}<sub>k∈N</sub>
- Variable splitting, BFGS-style updates, Newton-type methods
- Incremental or even stochastic versions



## Conclusions

- We have generalised gradient descent to Banach spaces
- ► We have proven global convergence for this generalisation.
- ▶ We shown different solution paths for same initialisation.

#### Outlook

- ▶ Convergence analysis for unbounded {q<sup>k</sup>}<sub>k∈N</sub>
- Variable splitting, BFGS-style updates, Newton-type methods
- Incremental or even stochastic versions



## Conclusions

- We have generalised gradient descent to Banach spaces
- ► We have proven global convergence for this generalisation.
- ▶ We shown different solution paths for same initialisation.

## Outlook

- ► Convergence analysis for unbounded {q<sup>k</sup>}<sub>k∈ℕ</sub>
- Variable splitting, BFGS-style updates, Newton-type methods
- Incremental or even stochastic versions


### Conclusions & outlook

#### Conclusions

- We have generalised gradient descent to Banach spaces
- ► We have proven global convergence for this generalisation.
- ▶ We shown different solution paths for same initialisation.

#### Outlook

- Convergence analysis for unbounded  $\{q^k\}_{k\in\mathbb{N}}$
- Variable splitting, BFGS-style updates, Newton-type methods
- Incremental or even stochastic versions

Preprint: http://www.imi.kyushu-u.ac.jp/eng/files/ imipublishattachment/file/math\_58ec341a238fe.pdf



### Conclusions & outlook

#### Conclusions

- We have generalised gradient descent to Banach spaces
- ► We have proven global convergence for this generalisation.
- ▶ We shown different solution paths for same initialisation.

#### Outlook

- Convergence analysis for unbounded  $\{q^k\}_{k\in\mathbb{N}}$
- Variable splitting, BFGS-style updates, Newton-type methods
- Incremental or even stochastic versions

Preprint: http://www.imi.kyushu-u.ac.jp/eng/files/ imipublishattachment/file/math\_58ec341a238fe.pdf



### Thank you for your attention





Isaac<sup>®</sup> Newton

Trust







## The Kurdyka-Łojasiewicz property

### Definition (Kurdyka-Łojasiewicz (KL) property)

We assume for  $\eta > 0$  that  $\varphi : [0, \eta[ \rightarrow \mathbb{R}_{\geq 0}]$  is a concave function that is continuous at zero and satisfies  $\varphi(0) = 0$ ,  $\varphi \in C^1(]0, \eta[)$ .

1. *E* fulfils the KL property at a point  $\overline{u} \in \mathbb{R}^n$  if there exists  $\eta \in ]0, \infty[$ , a neighbourhood *U* of  $\overline{u}$  such that for all

$$u \in U \cap \{u \mid E(\overline{u}) < E(u) < E(\overline{u}) + \eta\}$$

we observe

$$\varphi'(E(u) - E(\overline{u})) \|\nabla E(u)\|_2 \ge 1.$$
(10)

 If E satisfies the KL property for all arguments in ℝ<sup>n</sup>, E is called a KL functional.



## The Kurdyka-Łojasiewicz property

### Definition (Kurdyka-Łojasiewicz (KL) property)

We assume for  $\eta > 0$  that  $\varphi : [0, \eta[ \rightarrow \mathbb{R}_{\geq 0}]$  is a concave function that is continuous at zero and satisfies  $\varphi(0) = 0$ ,  $\varphi \in C^1(]0, \eta[)$ .

1. *E* fulfils the KL property at a point  $\overline{u} \in \mathbb{R}^n$  if there exists  $\eta \in ]0, \infty[$ , a neighbourhood *U* of  $\overline{u}$  such that for all

$$u \in U \cap \{u \mid E(\overline{u}) < E(u) < E(\overline{u}) + \eta\}$$

we observe

$$\varphi'(E(u) - E(\overline{u})) \|\nabla E(u)\|_2 \ge 1.$$
(10)

 If E satisfies the KL property for all arguments in ℝ<sup>n</sup>, E is called a KL functional.



# The Kurdyka-Łojasiewicz property

#### Example - Łojasiewicz 1965

Choose  $\varphi(x) = \frac{1}{1-\theta} |x|^{1-\theta}$ , for  $\theta \in ]0,1[$ . Then (10) reads as

 $|E(u) - E(\overline{u})|^{-\theta} \|\nabla E(u)\|_2 \ge 1,$ 

respectively

$$|E(u) - E(\overline{u})|^{\theta} \leq \|\nabla E(u)\|_2.$$

Thus, the norm of the gradient is a bound for the functional values.

