

Optimization techniques for tensor approximation: Newton Trust Region TT \rightarrow Canonical

Vladimir Kazeev

vladimir.kazeev@gmail.com

Institute of Numerical Mathematics
Russian Academy of Sciences

Workshop “Tensor Methods in Multi-Dimensional Boundary-Value
and Spectral Problems”
MPI MiS Leipzig, May 3-5, 2010

Canonical Decomposition

$$\mathbf{S}(i) = \sum_{\beta=1}^R U_1(i_1, \beta) \dots U_d(i_d, \beta)$$

$$\mathbf{S} = \sum_{\beta=1}^R \mathbf{S}_\beta,$$

$$\mathbf{S}_\beta = U_1(:, \beta) \otimes \dots \otimes U_d(:, \beta)$$

dnR parameters
compact representation,
linear complexity in respect to d

Canonical low-rank approximation

$$\mathbf{S} \approx (V_1 \dots V_d) = \mathbf{T}$$

$$\|\mathbf{T} - \mathbf{S}\|^2 = \sum_{\mathbf{i} \in \mathcal{I}}^r (\mathbf{T}(\mathbf{i}) - \mathbf{S}(\mathbf{i}))^2$$

Optimization problem

$$\mathbf{T} = (V_1 \dots V_d) = \underset{\substack{V^{(1)} \in \mathbb{R}^{n \times r} \\ \dots \\ V^{(d)} \in \mathbb{R}^{n \times r}}}{\operatorname{argmin}} \varphi,$$

$$\varphi = \frac{1}{2\|\mathbf{S}\|^2} \|\mathbf{T} - \mathbf{S}\|^2 + \frac{\kappa}{2\|\mathbf{S}\|^2} \sum_{\beta=1}^r \|\mathbf{T}_\beta\|^2$$

Canonical low-rank approximation

The Frobenius norm minimization problem may not have a solution

Let $x_k, y_k \in \mathbb{R}^n$, $k = 1 \dots 3$, be pairs of linearly independent vectors. Consider two tensors:

$$\begin{aligned} \mathbf{S} &= x_1 \otimes x_2 \otimes y_3 + x_1 \otimes y_2 \otimes x_3 + y_1 \otimes x_2 \otimes y_3, \\ \mathbf{T}_\varepsilon &= \frac{1}{\varepsilon} (x_1 + \varepsilon y_1) \otimes (x_2 + \varepsilon y_2) \otimes (x_3 + \varepsilon y_3) - \frac{1}{\varepsilon} x_1 \otimes x_2 \otimes x_3. \end{aligned}$$

Then $\text{rank } \mathbf{S} = 3$ and $\text{rank } \mathbf{T}_\varepsilon = 2$ for all $\varepsilon > 0$ while $\|\mathbf{T}_\varepsilon - \mathbf{S}\| = \mathcal{O}(\varepsilon)$.

(*de Silva and Lim, 2008*)

Canonical low-rank approximation

Solution of the Frobenius norm minimization problem may be essentially ununique

$$\sin\left(\sum_{k=1}^d x_k\right) = \sum_{k=1}^d \sin(x_k) \prod_{\substack{\zeta=1 \\ \zeta \neq k}}^d \frac{\sin(x_\zeta + \alpha_\zeta - \alpha_k)}{\sin(\alpha_\zeta - \alpha_k)}$$

for all $\{\alpha_k\}_{k=1}^d$:

$$\sin(\alpha_k - \alpha_\zeta) \neq 0 \quad \forall k, \zeta = 1 \dots d: \quad \zeta \neq k$$

(Mohlenkamp and Monzón, 2005)

Related works

- (2001) *Zhang, Golub. Rank-one approximation to high order tensors*
- (2005) *Beylkin, Mohlenkamp. Algorithms for numerical analysis in high dimensions*
- (2006) *Oseledets, Savostyanov. Minimization methods for approximating tensors and their comparison*
- (2008) *Espig. Effiziente Bestapproximation mittels Summen von Elementartensoren in hohen Dimensionen*
- (2009) *Acar, Kolda, Dunlavy. An optimization approach for fitting canonical tensor decompositions*
- (2010) *Kazeev, Tyrtysnikov. Structure of the Hessian matrix and an economical implementation of Newton method in the problem of canonical tensor approximation*

Basic ideas of the method

- Newton minimization of quadratic model of objective function
- global convergence due to Trust Region strategy
- PCG as a solver for Trust Region minimization subproblem
- specific strategy of quitting iterations (*Steihaug*, 1983)

Trust Region approach

Given $v_0, \delta_{\max}, \delta_0 \leq \delta_{\max}, \epsilon_g > 0, \kappa, K, \rho, \sigma, \beta$:
 $0 < \rho \leq \sigma < 1, 0 \leq \beta < \frac{1}{4}, 0 < \kappa < 1 < K$.

- ① quit iterations if $\|\mathbf{g}_k\| \leq \epsilon_g$
- ② find (accurate or approximate) solution w_k of a constrained minimization subproblem
- ③ evaluate

$$\text{kred}_k = \frac{\text{ared}_k}{\text{pred}_k} = \frac{\varphi(v_k) - \varphi(v_k + w_k)}{-\langle \mathbf{g}_k, w_k \rangle - \frac{1}{2} \langle w_k, \mathbf{H}_k w_k \rangle}.$$

- ④ let

$$v_{k+1} = \begin{cases} v_k + w_k, & \text{if } \text{kred}_k \geq \beta, \\ v_k & \text{otherwise.} \end{cases}$$

- ⑤ let

$$\delta_{k+1} = \begin{cases} \min \{ \delta_{\max}, K \|w_k\| \}, & \text{if } \text{kred}_k \geq \sigma, \\ \kappa \|w_k\|, & \text{if } \text{kred}_k \leq \rho, \end{cases}$$

Subproblem solution

Preconditioned CG with Trust Region constraints:

- ① if $\langle p_k, \mathbf{H}p_k \rangle \leq 0$, let $w_{k+1} = w_k + \tau p_k$, $\tau > 0$ s. t. $\|w_{k+1}\|_{\mathbf{A}} = \delta$, and break subiterations,
- ② if $\|w_{k+1}\|_{\mathbf{A}} > \delta$, let $w_{k+1} = w_k + \tau p_k$, $\tau > 0$ s. t. $\|w_{k+1}\|_{\mathbf{A}} = \delta$, and break subiterations,

What do we have to do?

- evaluate the objective function

$$\varphi = \frac{1}{2\|\mathbf{S}\|^2} \|\mathbf{T} - \mathbf{S}\|^2 + \frac{\alpha}{2\|\mathbf{S}\|^2} \sum_{\beta=1}^r \|\mathbf{T}_j\|^2$$

- compute its gradient
- compute its Hessian

There is no need to compute the Hessian explicitly, only the matvec operation is required

Convergence theorems

Theorem (Steihaug, 1983)

Let a function f be bounded below and its gradient $g = \nabla f$ be uniformly continuous. If the subproblem of every k -th iteration is solved with relative error ϵ_k such that $\sup_k \epsilon_k < 1$, then

$$\liminf_{k \rightarrow \infty} g_k = 0.$$

Further, if the sequence $\{H_k\}_{k=1}^{\infty}$ is bounded, then

$$\lim_{k \rightarrow \infty} g_k = 0.$$

Convergence theorems

Theorem (Steihaug, 1983)

Let f be bounded below and twice continuously differentiable in the neighbourhood of a solution x^ and $\nabla^2 f(x^*) > 0$. Then there exist such $\varepsilon > 0$ and $\Delta > 0$ that if $\Delta_0 \leq \Delta$ and $\|x_0 - x^*\| \leq \varepsilon$ then $\{x_k\}_{k=1}^{\infty}$ converges to x^* .*

Tensor Train

Oseledets and Tyrtshnikov, 2009:

$$\begin{aligned} \mathbf{S}(\mathbf{i}) = & \sum_{\alpha \in \mathcal{A}} U_1(i_1, \alpha_1) \cdot \\ & \cdot U_2(\alpha_1, i_2, \alpha_2) \dots U_{d-1}(\alpha_{d-2}, i_{d-1}, \alpha_{d-1}) \cdot \\ & \cdot U_d(\alpha_{d-1}, i_d) \end{aligned}$$

for $U_k(\alpha_{k-1}, i_k, \alpha_k) \in \mathbb{R}^{R_{k-1} \times n \times R_k}$

+ robust algorithms for computations and manipulations

TT to Canonical

$$\begin{aligned} \mathbf{S} &= (U_1 \dots U_d) \approx (V_1 \dots V_d) = \mathbf{T} \\ &\quad R_1 \dots R_{d-1} \quad r \\ \|\mathbf{T} - \mathbf{S}\|^2 &= \sum_{i \in \mathcal{I}} (\mathbf{T}(i) - \mathbf{S}(i))^2 \end{aligned}$$

Optimization problem

$$\begin{aligned} \mathbf{T} &= (V_1 \dots V_d) = \operatorname{argmin}_{\substack{V^{(1)} \in \mathbb{R}^{n \times r} \\ \dots \\ V^{(d)} \in \mathbb{R}^{n \times r}}} \varphi, \\ \varphi &= \frac{1}{2\|\mathbf{S}\|^2} \|\mathbf{T} - \mathbf{S}\|^2 + \frac{\kappa}{2\|\mathbf{S}\|^2} \sum_{\beta=1}^r \|\mathbf{T}_\beta\|^2 \end{aligned}$$

Matrices of quadratic model

Define $r \times r$ -matrices, which depend on decomposition of \mathbf{T} :

$$\Psi_k(\beta, \beta') = \sum_{i_k=1}^{n_k} V_k(i_k, \beta) V_k(i_k, \beta'),$$

$$Q(\beta, \beta') = (1 + \kappa \delta(\beta, \beta')) \prod_{k=1}^d \Psi_k(\beta, \beta'),$$

$$Q_\zeta(\beta, \beta') = (1 + \kappa \delta(\beta, \beta')) \prod_{\substack{k=1 \\ k \neq \zeta}}^d \Psi_k(\beta, \beta'),$$

$$Q_{\zeta\zeta'}(\beta, \beta') = (1 + \kappa \delta(\beta, \beta')) \prod_{\substack{k=1 \\ k \neq \zeta, \zeta'}}^d \Psi_k(\beta, \beta').$$

Matrices of quadratic model

Define the following matrices, which depend on decompositions of both S and T

$$\Phi_k(\alpha_{k-1}, \beta, \alpha_k) = \sum_{i_k=1}^{n_k} U_k(\alpha_{k-1}, i_k, \alpha_k) V_k(i_k, \beta)$$

Evaluation of the objective function

In terms of matrices of quadratic model:

$$\varphi = \frac{1}{2} \sum_{\beta=1}^r \sum_{\beta'=1}^r Q(\beta, \beta') - \sum_{\beta=1}^r \sum_{\alpha \in \mathcal{A}} \prod_{k=1}^d \Phi_k(\alpha_{k-1}, \beta, \alpha_k) + \frac{1}{2}$$

Evaluation of the gradient

In terms of matrices of quadratic model:

$$\begin{aligned}
 \mathbf{g}_\zeta(\xi, \eta) &= \frac{\partial \varphi}{\partial V_\zeta(\xi, \eta)} = \\
 &= \sum_{\beta=1}^r V_\zeta(\xi, \beta) Q_\zeta(\beta, \eta) \\
 &- \sum_{\alpha \in \mathcal{A}} U_\zeta(\alpha_{\zeta-1}, \xi, \alpha_\zeta) \prod_{\substack{k=1 \\ k \neq \zeta}}^d \Phi_k(\alpha_{k-1}, \eta, \alpha_k)
 \end{aligned}$$

Evaluation of matvec with the Hessian

In terms of matrices of quadratic model:

$$\begin{aligned}
 \mathbf{H}_{\zeta\zeta'}(\xi, \eta; \xi', \eta') &= \frac{\partial^2 \varphi}{\partial V_{\zeta'}(\xi', \eta') \partial V_{\zeta}(\xi, \eta)}, \\
 \mathbf{H} &= \mathbf{A} + \mathbf{B} + \mathbf{C} - \mathbf{D}, \\
 \mathbf{A}_{\zeta\zeta'}(\xi, \eta; \xi', \eta') &= \boldsymbol{\delta}(\zeta', \zeta) \boldsymbol{\delta}(\xi', \xi) Q_{\zeta}(\eta, \eta'), \\
 \mathbf{B}_{\zeta\zeta'}(\xi, \eta; \xi', \eta') &= \bar{\boldsymbol{\delta}}(\zeta', \zeta) V_{\zeta'}(\xi', \eta) V_{\zeta}(\xi, \eta') Q_{\zeta\zeta'}(\eta, \eta').
 \end{aligned}$$

Evaluation of matvec with the Hessian

In terms of matrices of quadratic model:

$$\begin{aligned}
 \mathbf{C}_{\zeta\zeta'}(\xi, \eta; \xi', \eta') &= \bar{\boldsymbol{\delta}}(\zeta', \zeta) \boldsymbol{\delta}(\eta', \eta) \\
 &\quad \cdot \sum_{\beta=1}^r V_{\zeta}(\xi, \beta) V_{\zeta'}(\xi', \beta) Q_{\zeta}(\beta, \eta), \\
 \mathbf{D}_{\zeta\zeta'}(\xi, \eta; \xi', \eta') &= \bar{\boldsymbol{\delta}}(\zeta', \zeta) \boldsymbol{\delta}(\eta', \eta) \\
 &\quad \cdot \sum_{\alpha \in \mathcal{A}} U_{\zeta}(\alpha_{\zeta-1}, \xi, \alpha_{\zeta}) U_{\zeta'}(\alpha_{\zeta'-1}, \xi', \alpha_{\zeta'}) \\
 &\quad \cdot \prod_{\substack{k=1 \\ k \neq \zeta, \zeta'}}^d \Phi_k(\alpha_{k-1}, \eta, \alpha_k).
 \end{aligned}$$

Computation of the matrices of quadratic model

Computational costs are

$$\mathbf{Work} = \mathcal{O}((d+n)dr^2 + d^2nrR(R+n)),$$

$$\mathbf{Mem} = \mathcal{O}(d^2r^2 + d(dn+R)rR)$$

In case $n = 2$

$$\mathbf{Work} = \mathcal{O}(d^2r^2 + d^2rR^2),$$

$$\mathbf{Mem} = \mathcal{O}(d^2r^2 + d(d+R)rR)$$

Computation costs

Objective function φ

$$\mathbf{Work} = \mathcal{O}(r^2 + rR),$$

$$\mathbf{Mem} = \mathcal{O}(1).$$

Its gradient g

$$\mathbf{Work} = \mathcal{O}(dnr^2 + dnrR),$$

$$\mathbf{Mem} = \mathcal{O}(1).$$

Computation costs

$A > 0$

- preconditioner for the Hessian:

$$\text{Work} = \mathcal{O}(dr^3), \quad \text{Mem} = \mathcal{O}(dr^2)$$

- matvec:

$$\text{Work} = \mathcal{O}(dnr^2), \quad \text{Mem} = \mathcal{O}(1)$$

Computation costs

B, C or $B + C$

- matvec:

$$\text{Work} = \mathcal{O}((d+n)dr^2), \quad \text{Mem} = \mathcal{O}(dr^2)$$

D

- matvec:

$$\text{Work} = \mathcal{O}(d^2n^2r), \quad \text{Mem} = \mathcal{O}(1)$$

The tensor

$$\begin{pmatrix} c_1 & c_3 \\ c_2 & c_4 \end{pmatrix} = \begin{pmatrix} a_1 & a_3 \\ a_2 & a_4 \end{pmatrix} \begin{pmatrix} b_1 & b_3 \\ b_2 & b_4 \end{pmatrix},$$

$$c_k = \sum_{i=1}^4 \sum_{j=1}^4 s_{ijk} a_i b_j,$$

$S = \{s_{ijk}\}$ is a $4 \times 4 \times 4$ -tensor: $d = 3$, $n = 4$, 8 ones and 56 zeros. Once we have a decomposition $S = (U, V, W)$ of rank 7, we may compute

$$c_k = \sum_{s=1}^7 w_{ks} \sum_{i=1}^4 u_{is} a_i \sum_{j=1}^4 v_{js} b_j.$$

Results

2×2 : $d = 3$, $n = 4$, $R = 8$, $r = 7$, $d r n = 84$

	e_0	IT	τ , s.	$\ g\ _2$	e
TT to Can, MATLAB	$1.0 \cdot 10^0$	23	2.4	$2.3 \cdot 10^{-8}$	10^{-8}
Can to Can, C	$1.0 \cdot 10^0$	21	0	$7.5 \cdot 10^{-12}$	$2.6 \cdot 10^{-8}$

Convergence

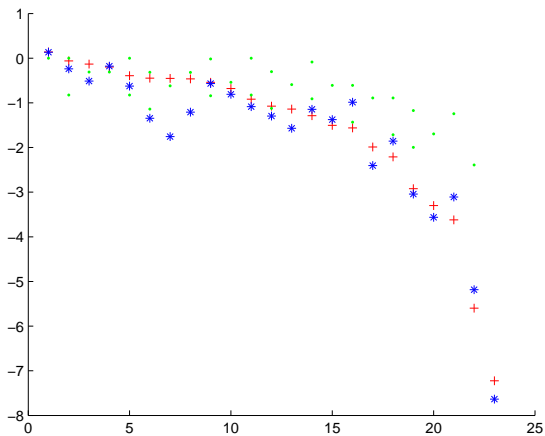


Рис.: e , $\|g\|_2$ and δ vs. iteration no.

Basic idea

- given a function $f : [0, 1] \rightarrow \mathbb{R}$, which is a sum of exponentials, and a mesh $\{x_i\}_{i=1}^N \subset [0, 1]$, $N = 2^d$
- consider a Q-mesh and a Q-vector:

$$\mathbf{S}(i_1 \dots i_d) = f\left(\sum_{k=1}^d 2^{k-1} i_k\right)$$
- convert to TT: $\mathbf{S} = (U_1 \dots U_d)$
- approximate in Canonical with rank r :

$$\mathbf{T} = (V_1 \dots V_d) \approx (U_1 \dots U_d)$$
- hope that $\text{vec}(\mathbf{T}_\beta)$ is an exponential

Two exponentials

$$f(x) = 0.1e^x - 0.003e^{5x}, d = 12$$

e_0	IT	τ , s.	$\ g\ _2$	e
$1.0 \cdot 10^0$	35	9.7	$1.6 \cdot 10^{-7}$	$1.5 \cdot 10^{-8}$

the same f with multiplicative 10% noise, $d = 12$

e_0	IT	τ , s.	$\ g\ _2$	e
$1.0 \cdot 10^0$	42	10.8	$9.7 \cdot 10^{-12}$	$5.8 \cdot 10^{-2}$

Convergence (no noise)

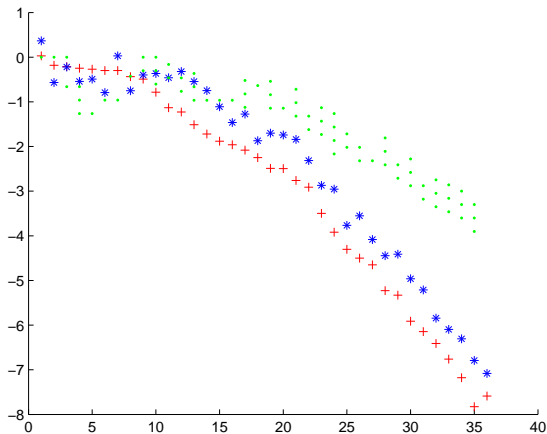
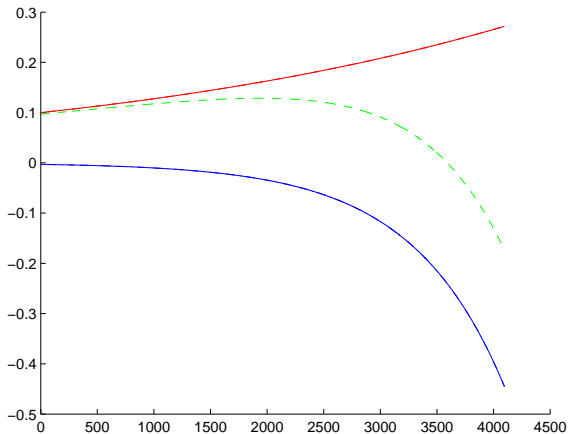
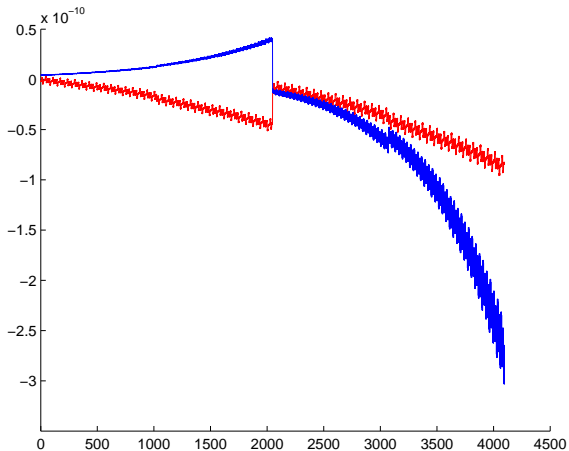


Рис.: e , $\|g\|_2$ and δ vs. iteration no.

Splitting (no noise)



Splitting error (no noise)



Convergence (with noise)

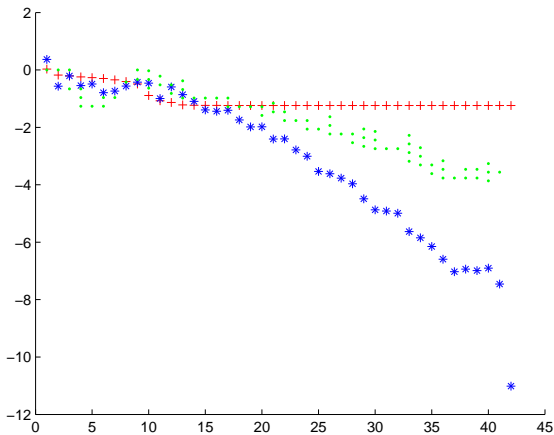
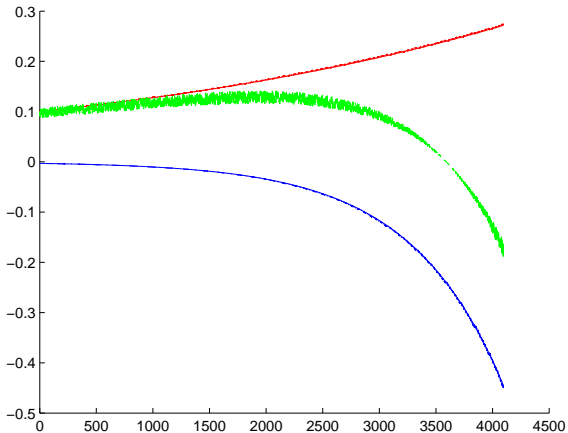
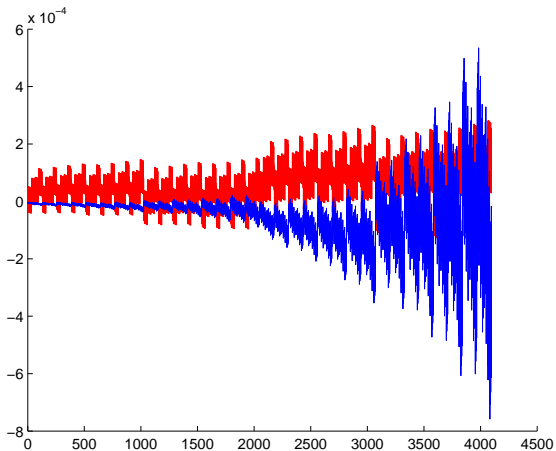


Рис.: e , $\|g\|_2$ and δ vs. iteration no.

Splitting (with noise)



Splitting error (with noise)



Three exponentials

$$f(x) = 0.015e^{3x} - 0.003e^{5x} - 0.0001e^{10x}, \quad d = 12$$

e_0	IT	$\tau, \text{ s.}$	$\ g\ _2$	e
$1.0 \cdot 10^0$	350	115	$3.0 \cdot 10^{-7}$	$6.2 \cdot 10^{-7}$

Convergence (no noise)

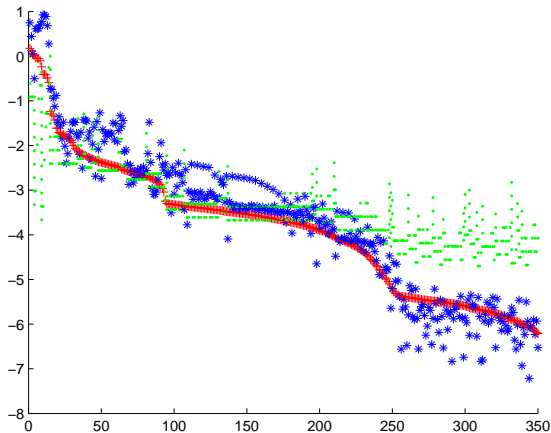
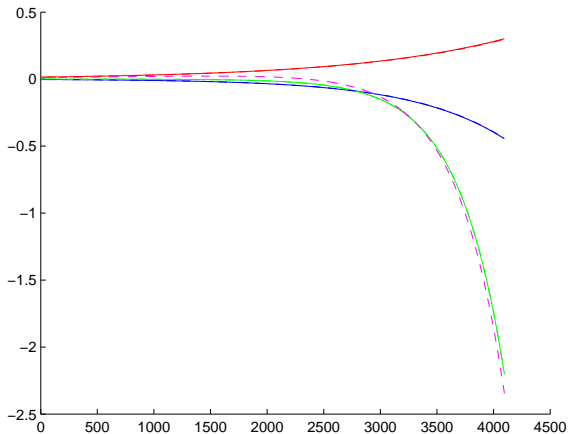
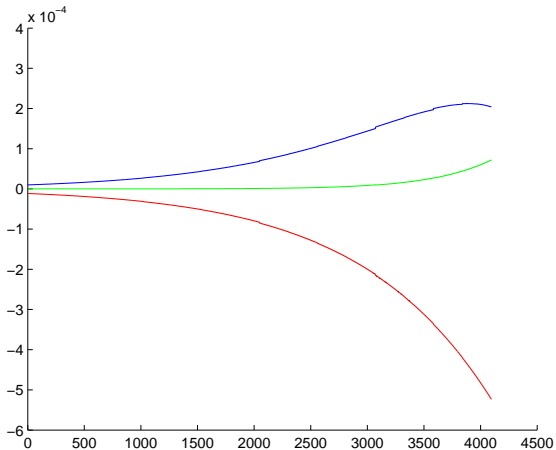


Рис.: e , $\|g\|_2$ and δ vs. iteration no.

Splitting (no noise)



Splitting error (no noise)



Thank you for your attention!