

Trust Region Method

Lectures for PHD course on
Non-linear equations and numerical optimization

Enrico Bertolazzi

DIMS – Università di Trento

March 2005

Outline

- 1 The Trust Region method
- 2 The exact solution of trust region step
- 3 The dogleg trust region step

- Newton and quasi-Newton methods search a solution iteratively by choosing at each step a search direction and minimize in this direction.
- An alternative approach is to find a direction and a step-length, then if the step is successful in some sense the step is accepted. Otherwise another direction and step-length is chosen.
- The choice of the step-length and direction is algorithm dependent but a successful approach is the one based on trust region.

- Newton and quasi-Newton at each step (approximately) solve the minimization problem

$$\min m(x_k + s) = f(x_k) + \nabla f(x_k)s + \frac{1}{2}s^T H_k s$$

in the case H_k is symmetric and positive definite (SPD).

- If H_k is SPD the minimum is

$$s = -H_k^{-1}g_k, \quad g_k = \nabla f(x_k)^T$$

and s is the quasi-Newton step.

- If $H_k = \nabla^2 f(x_k)$ and is SPD, then $s = -H_k^{-1}g_k$ is the Newton step.

- If H_k is not positive definite, the search direction $-H_k^{-1}g_k$ may fail to be a descent direction and the previous minimization problem can have no solution.
- The problem is that the model $m(x_k + s)$ is an approximation of $f(x)$

$$m(x_k + s) \approx f(x_k) + \nabla f(x_k)^T s + \frac{1}{2} s^T H_k s$$

and this approximation is valid only in a small neighbors of x_k .

- So that an alternative minimization problem is the following

$$\min m(x_k + s) = f(x_k) + \nabla f(x_k)^T s + \frac{1}{2} s^T H_k s,$$

$$\text{Subject to } \|s\| \leq \delta_k$$

δ_k is the trust region of the model $m(x)$, i.e. the region where we trust the model is valid.



Algorithm (Generic trust region algorithm)

```

x assigned; δ assigned;
g ← ∇f(x)T; H ← ∇2f(x)-1;
while ‖g‖ > ε do
  s ← arg min‖s‖ ≤ δ m(x + s) = f(x) + gTs + 1/2 sTHs;
  pred ← m(x + s) - m(x);
  ared ← f(x + s) - f(x);
  if (ared/pred) < η1 then
    x ← x; δ ← δγ1; — reject step, reduce δ
  else
    x ← x + s; — accept step, update H
    if (ared/pred) > η2 then
      δ ← max{δ, γ2 ‖s‖}; — enlarge δ
    end if
  end if
end while

```



- The previous algorithm is based on two keys ingredients:
 - 1 The ratio $r = (\text{ared}/\text{pred})$ which is the ratio of the **actual reduction** and the **predicted reduction**.
 - 2 Enlarge or reduce the trust region δ .
- If the ratio r is between $0 < \eta_1 < r < \eta_2 < 1$ we have that the model is quite appropriate; we accept the step and do not modify the trust region.
- If the ratio r is small $r \leq \eta_1$ we have that the model is not appropriate; we do not accept the step and we must reduce the trust region by a factor $\gamma_1 < 1$
- If the ratio r is large $r \geq \eta_2$ we have that the model is very appropriate; we do accept the step and we enlarge the trust region factor $\gamma_2 > 1$
- The algorithm is quite insensitive to the constant η_1 and η_2 . Typical values are $\eta_1 = 0.25$, $\eta_2 = 0.75$, $\gamma_1 = 0.5$ and $\gamma_2 = 3$.



Lemma

Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be twice continuously differentiable, $H \in \mathbb{R}^{n \times n}$ symmetric and positive definite. Then the problem

$$\min m(x + s) = f(x) + \nabla f(x)^T s + \frac{1}{2} s^T H s,$$

$$\text{Subject to } \|s\| \leq \delta$$

is solved by

$$s(\mu) \doteq -(H + \mu I)^{-1} g, \quad g = \nabla f(x)^T$$

for the unique $\mu \geq 0$ such that $\|s(\mu)\| = \delta$, unless $\|s(0)\| \leq \delta$, in which case $s(0)$ is the solution. For any $\mu \geq 0$, $s(\mu)$ defines a descent direction for f from x .



Proof.

(1/2).

If $\|s(0)\| \leq \delta$ then $s(0)$ is the global minimum inside the trust region. Otherwise consider the Lagrangian

$$\mathcal{L}(s, \mu) = a + g^T s + \frac{1}{2} s^T H s + \frac{1}{2} \mu (s^T s - \delta^2),$$

where $a = f(x)$ and $g = \nabla f(x)^T$. Then we have

$$\frac{\partial \mathcal{L}}{\partial s}(s, \mu) = H s + \mu s + g = 0 \Rightarrow s = -(H + \mu I)^{-1} g$$

and $s^T s = \delta^2$. Remember that if H is SPD then $H + \mu I$ is SPD for all $\mu \geq 0$. Moreover the inverse of an SPD matrix is SPD. From

$$g^T s = -g^T (H + \mu I)^{-1} g < 0 \quad \text{for all } \mu \geq 0$$

follows that $s(\mu)$ is a descent direction for all $\mu \geq 0$.



Remark

As a consequence of the previous Lemma we have:

- as the ray of the trust region becomes smaller as the scalar μ becomes larger. This means that the search direction become more and more oriented toward the gradient direction.
- as the ray of the trust region becomes larger as the scalar μ becomes smaller. This means that the search direction become more and more oriented toward the Newton direction.

Thus a trust region technique not only change the size of the step-length but also its direction. This results in a more robust numerical technique. The price to pay is that the solution of the minimization is more costly than the inexact line search.



Proof.

(2/2).

To prove the uniqueness consider expand the gradient g with the eigenvectors of H

$$g = \sum_{i=1}^n \alpha_i u_i$$

H is SPD so that u_i can be chosen orthonormal. It follows

$$(H + \mu I)^{-1} g = (H + \mu I)^{-1} \sum_{i=1}^n \alpha_i u_i = \sum_{i=1}^n \frac{\alpha_i}{\lambda_i + \mu} u_i$$

$$\|(H + \mu I)^{-1} g\|^2 = \sum_{i=1}^n \frac{\alpha_i^2}{(\lambda_i + \mu)^2}$$

and $\|(H + \mu I)^{-1} g\|$ is a monotonically decreasing function of μ .



Solving the constrained minimization problem

As for the line-search problem we have many alternative for solving the constrained minimization problem:

- We can solve **accurately** the constrained minimization problem. For example by an iterative method.
- We can **approximate** the solution of the constrained minimization problem.

as for the line search the accurate solution of the constrained minimization problem is not paying while a good cheap approximations is normally better performing.



Outline

- 1 The Trust Region method
- 2 The exact solution of trust region step
- 3 The dogleg trust region step



The Newton approach

(1/5)

- Consider the Lagrangian

$$\mathcal{L}(s, \mu) = a + g^T s + \frac{1}{2} s^T H s + \frac{1}{2} \mu (s^T s - \delta^2),$$

where $a = f(x)$ and $g = \nabla f(x)^T$.

- Then we can try to solve the nonlinear system

$$\frac{\partial \mathcal{L}}{\partial (s, \mu)}(s, \mu) = \begin{pmatrix} Hs + \mu s + g \\ (s^T s - \delta^2)/2 \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ 0 \end{pmatrix}$$

- Using Newton method we have

$$\begin{pmatrix} s_{k+1} \\ \mu_{k+1} \end{pmatrix} = \begin{pmatrix} s_k \\ \mu_k \end{pmatrix} - \begin{pmatrix} H + \mu I & s \\ s^T & 0 \end{pmatrix}^{-1} \begin{pmatrix} Hs_k + \mu_k s_k + g \\ (s_k^T s_k - \delta^2)/2 \end{pmatrix}$$



The Newton approach

(2/5)

- A better approach is given by solving $\Phi(\mu) = 0$ where

$$\Phi(\mu) = \|s(\mu)\| - \delta, \quad \text{and} \quad s(\mu) = -(H + \mu I)^{-1} g$$

- To build Newton method we need to evaluate

$$\Phi(\mu)' = \frac{s(\mu)^T s(\mu)'}{\|s(\mu)\|}, \quad s(\mu)' = (H + \mu I)^{-2} g$$

where to evaluate $s(\mu)'$ we differentiate the relation

$$Hs(\mu) + \mu s(\mu) = g \quad \Rightarrow \quad Hs(\mu)' + \mu s(\mu)' + s(\mu) = \mathbf{0}$$

- Putting all in a Newton step we obtain

$$\mu_{k+1} = \mu_k - \frac{\|s(\mu_k)\|}{s(\mu_k)^T s(\mu_k)'} (\|s(\mu_k)\| - \delta)$$



The Newton approach

(3/5)

- Newton step can be reorganized as follows

$$s_k = -(H + \mu I)^{-1} g$$

$$s_k' = -(H + \mu I)^{-1} s_k$$

$$\beta = \sqrt{s_k^T s_k}$$

$$\mu_{k+1} = \mu_k - \frac{\beta(\beta - \delta)}{s_k^T s_k'}$$

- Thus Newton step require **two** linear system solution per step. However the coefficient matrix is the same so that only **one** LU factorization, thus the cost per step is essentially due to the LU factorization.



The Newton approach

(4/5)

- Evaluating $\Phi(\mu)''$ we have

$$\Phi(\mu)'' = \frac{\|s(\mu)\|^2 + s(\mu)^T s(\mu)''}{\|s(\mu)\|} + \frac{(s(\mu)^T s(\mu)')^2}{\|s(\mu)\|^2}$$

where

$$s(\mu)'' = \mathbf{0}$$

- In fact, from

$$(H + \mu I)s(\mu)' = s(\mu)$$

we have

$$Hs(\mu)'' + \mu s(\mu)'' + s(\mu)' = s(\mu)' \Rightarrow s(\mu)'' = \mathbf{0}.$$

- Then for all $\mu \geq 0$ we have $\Phi''(\mu) > 0$.



- If we develop the vector g with the orthonormal bases given by the eigenvectors of H we have

$$g = \sum_{i=1}^n \alpha_i u_i$$

- Using this expression to evaluate $s(\mu)$ we have

$$s(\mu) = -(H + \mu I)^{-1} g = \sum_{i=1}^n \frac{\alpha_i}{\mu + \lambda_i} u_i$$

$$\|s(\mu)\| = \left(\sum_{i=1}^n \frac{\alpha_i^2}{(\mu + \lambda_i)^2} \right)^{1/2}$$

- This expression suggest to use as a model for $\Phi(\mu)$ the following expression

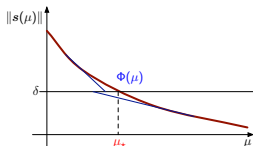
$$m_k(\mu) = \frac{\alpha_k}{\beta_k + \mu} - \delta$$



The Newton approach

(5/5)

- From $\Phi''(\mu) > 0$ we have that Newton step underestimates μ at each step.



- The model consists of two parameter α_k and β_k . To set this parameter we can impose

$$m_k(\mu_k) = \frac{\alpha_k}{\beta_k + \mu_k} - \delta = \Phi(\mu_k)$$

$$m_k(\mu_k)' = -\frac{\alpha_k}{(\beta_k + \mu_k)^2} = \Phi(\mu_k)'$$

- solving for α_k and β_k we have

$$\alpha_k = -\frac{(\Phi(\mu_k) + \delta)^2}{\Phi(\mu_k)'} \quad \beta_k = -\frac{\Phi(\mu_k) + \delta}{\Phi(\mu_k)'} - \mu_k$$

where

$$\Phi(\mu_k) = \|s(\mu_k)\| - \delta \quad \Phi(\mu_k)' = -\frac{s(\mu_k)^T (H + \mu_k I)^{-1} s(\mu_k)}{\|s(\mu_k)\|^2}$$

- Having α_k and β_k it is possible to solve $m_k(\mu) = 0$ obtaining

$$\mu_{k+1} = \frac{\alpha_k}{\delta} - \beta_k$$



- Substituting α_k and β_k the step become

$$\mu_{k+1} = \mu_k - \frac{\Phi(\mu_k)}{\Phi'(\mu_k)} - \frac{\Phi(\mu_k)^2}{\Phi'(\mu_k)\delta} = \mu_k - \frac{\Phi(\mu_k)}{\Phi'(\mu_k)} \left(1 + \frac{\Phi(\mu_k)}{\delta} \right)$$

- Comparing with the Newton step

$$\mu_{k+1} = \mu_k - \frac{\Phi(\mu_k)}{\Phi'(\mu_k)}$$

we see that this method perform larger step by a factor $1 + \Phi(\mu_k)\delta^{-1}$.

- Notice that $1 + \Phi(\mu_k)\delta^{-1}$ converge to 1 as $\mu_k \rightarrow \mu_*$. So that this iteration become the Newton iteration as μ_k becomes near the solution.



Algorithm (Exact trust region algorithm)

```

 $\mu, \mathbf{g}, \mathbf{H}$  assigned;
 $\mathbf{s} \leftarrow (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g}$ ;
while  $\|\mathbf{s}\| - \delta > \epsilon$  do
  — compute the model
   $\mathbf{s}' \leftarrow (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{s}$ ;
   $\Phi \leftarrow \|\mathbf{s}'\| - \delta$ ;
   $\Phi' \leftarrow -(\mathbf{s}'^T \mathbf{s}') / (\mathbf{s}'^T \mathbf{s})$ 
   $\alpha \leftarrow -(\Phi + \delta)^2 / \Phi'$ ;
   $\beta \leftarrow -(\Phi + \delta) / \Phi' - \mu$ ;
  — update  $\mu$  and  $\mathbf{s}$ 
   $\mu \leftarrow \frac{\alpha}{\delta} - \beta$ ;
   $\mathbf{s} \leftarrow (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{g}$ ;
end while

```




Outline

- The Trust Region method
- The exact solution of trust region step
- The dogleg trust region step**



The DogLeg approach

(1/3)

- The computation of the μ such that $\|\mathbf{s}(\mu)\| = \delta$ of the **exact** trust region computation can be very expensive.
- An alternative was proposed by Powell:
 -  **M.J.D. Powell**
A hybrid method for nonlinear equations
in: Numerical Methods for Nonlinear Algebraic Equations
ed. Ph. Rabinowitz, Gordon and Breach, pages 87-114,
1970.
 - where instead of computing exactly the curve $\mathbf{s}(\mu)$ a piecewise linear approximation $\mathbf{s}_{dl}(\mu)$ is used in computation.
- This approximation also permits to solve $\|\mathbf{s}_{dl}(\mu)\| = \delta$ explicitly.



The DogLeg approach

(2/3)

- Form the definition of $s(\mu) = -(H + \mu I)^{-1}g$ it follows

$$s(0) = -H^{-1}g, \quad \lim_{\mu \rightarrow \infty} \frac{s(\mu)'}{\|s(\mu)'\|} = \frac{g}{\|g\|}$$

i.e. the curve start from the Newton step and reduce to zero in the direction of the gradient step.

- The direction $-g$ is a descent direction, so that a first piece of the piecewise approximation should be a straight line from x to the minimum of $m_k(x - \lambda g)$. The minimum λ_* is found at

$$\lambda_* = \frac{\|g\|^2}{g^T H g}$$

- Having reached the minimum if the $-g$ direction we can now go to the point $x + s(0) = x - Hg$ with another straight line.



Lemma

Consider the **dogleg** curve connecting $x + s_n$, $x + s_g$ and x . The curve can be expressed as $x_{dl}(\mu) = x + s_{dl}(\mu)$ where

$$s_{dl}(\mu) = \begin{cases} \mu s_g + (1 - \mu) s_n & \text{for } \mu \in [0, 1] \\ (2 - \mu) s_g & \text{for } \mu \in [1, 2] \end{cases}$$

for this curve if s_g is not parallel to s_n we have that the function

$$d(\mu) = \|x_{dl}(\mu) - x\| = \|s_{dl}(\mu)\|$$

is strictly monotone decreasing, moreover the direction $s(\mu)$ is a descent direction for all $\mu \in [0, 2]$.



The DogLeg approach

(3/3)

- We denote by

$$s_g = -g \frac{\|g\|^2}{g^T H g}, \quad s_n = -H^{-1}g$$

respectively the step due to the unconstrained minimization in the gradient direction and in the Newton direction.

- The piecewise linear curve connecting $x + s_n$, $x + s_g$ and x is the **DogLeg** curve¹ $x_{dl}(\mu) = x + s_{dl}(\mu)$ where

$$s_{dl}(\mu) = \begin{cases} \mu s_g + (1 - \mu) s_n & \text{for } \mu \in [0, 1] \\ (2 - \mu) s_g & \text{for } \mu \in [1, 2] \end{cases}$$

¹notice that $s(\mu)$ is parametrized in the interval $[0, \infty]$ while $s_{dl}(\mu)$ is parametrized in the interval $[0, 2]$



Proof.

(1/5)

In order to have a unique solution to the problem $\|s_{dl}(\mu)\| = \delta$ we must have that $\|s_{dl}(\mu)\|$ is a monotone decreasing function:

$$\|s_{dl}(\mu)\|^2 = \begin{cases} \mu^2 s_g^2 + (1 - \mu)^2 s_n^2 + 2\mu(1 - \mu) s_g^T s_n & \mu \in [0, 1] \\ (2 - \mu)^2 s_g^2 & \mu \in [1, 2] \end{cases}$$

To check monotonicity we take first derivative

$$\begin{aligned} \frac{d}{d\mu} \|s_{dl}(\mu)\|^2 &= \begin{cases} 2\mu s_g^2 - 2(1 - \mu) s_n^2 + (2 - 4\mu) s_g^T s_n & \mu \in [0, 1] \\ (2\mu - 4) s_g^2 & \mu \in [1, 2] \end{cases} \\ &= \begin{cases} 2\mu(s_g^2 + s_n^2 - 2s_g^T s_n) - 2s_n^2 + 2s_g^T s_n & \mu \in [0, 1] \\ (2\mu - 4) s_g^2 & \mu \in [1, 2] \end{cases} \end{aligned}$$



Proof.

(2/5).

Notice that $(2\mu - 4) < 0$ for $\mu \in [1, 2]$ so that we need only to check that

$$2\mu(s_g^2 + s_n^2 - 2s_g^T s_n) - 2s_n^2 + 2s_g^T s_n < 0 \quad \text{for } \mu \in [0, 1]$$

Form the Cauchy-Schwartz inequality we have

$$\begin{aligned} s_g^2 + s_n^2 - 2s_g^T s_n &\geq s_g^2 + s_n^2 - 2\|s_g\| \|s_n\| \\ &= (\|s_g\| - \|s_n\|)^2 \geq 0 \end{aligned}$$

Then it is enough to check the inequality for $\mu = 1$

$$2(s_g^2 + s_n^2 - 2s_g^T s_n) - 2s_n^2 + 2s_g^T s_n = 2s_g^2 - 2s_g^T s_n$$

i.e. we must check $s_g^2 - s_g^T s_n < 0$.



Proof.

(3/5).

From the definition of s_g and s_n we have

$$\begin{aligned} s_g^2 - s_g^T s_n &= \lambda_*^2 \|g\|^2 - \lambda_* g^T H^{-1} g \\ &= \lambda_* \left[\frac{\|g\|^2}{g^T H g} \|g\|^2 - g^T H^{-1} g \right] \\ &= \frac{\lambda_*}{g^T H g} \left[\|g\|^4 - (g^T H g)(g^T H^{-1} g) \right] \end{aligned}$$

So that we must prove that

$$\|g\|^4 < (g^T H g)(g^T H^{-1} g)$$



Proof.

(4/5).

Expanding g by a set of orthonormal eigenvectors of H we have $g = \sum_{i=1}^n \alpha_i u_i$ and the previous inequality becomes

$$\begin{aligned} \|g\|^4 &= \left(\sum_{i=1}^n \alpha_i^2 \right)^2 = \left(\sum_{i=1}^n (\alpha_i \lambda_i^{1/2}) (\alpha_i \lambda_i^{-1/2}) \right)^2 \\ &\leq \left(\sum_{i=1}^n \alpha_i^2 \lambda_i \right) \left(\sum_{i=1}^n \alpha_i^2 \lambda_i^{-1} \right) = (g^T H g) (g^T H^{-1} g) \end{aligned}$$

from the Cauchy-Schwartz inequality the previous inequality is strict unless

$$\alpha_i \lambda_i = c \alpha_i, \quad i = 1, 2, \dots, n$$

this means that $\lambda_i = c$ that for all $\alpha_i \neq 0$. This imply $H^{-1}g = c^{-1}g$, i.e. Newton step and gradient step are parallel. But this is excluded in the lemma hypothesis.



Proof.

(5/5).

To prove that $s_{dl}(\mu)$ is a descent direction it is enough to notice that

- for $\mu \in [0, 1]$ the direction $s_{dl}(\mu)$ is a convex combination of s_g and s_n .
- for $\mu \in [1, 2]$ the direction $s_{dl}(\mu)$ is parallel to s_g .

so that it is enough to verify that s_g and s_n are descent direction. For s_g we have

$$s_g^T g = -\lambda_* g^T g < 0$$

For s_n we have

$$s_n^T g = -g^T H^{-1} g < 0$$



Using the previous Lemma we can prove

Lemma

If $\|s_{dl}(0)\| \geq \delta$ then there is unique point $\mu \in [0, 2]$ such that $\|s_{dl}(\mu)\| = \delta$.

Proof.

It is enough to notice that $s_{dl}(2) = \mathbf{0}$ and that $\|s_{dl}(\mu)\|$ is strictly monotonically descendent. \square

The approximate solution of the constrained minimization can be obtained by this simple algorithm

- 1 if $\delta \leq \|s_g\|$ we set $s_{dl} = -\delta s_g / \|s_g\|$;
- 2 if $\delta \leq \|s_n\|$ we set $s_{dl} = \alpha s_g + (1 - \alpha) s_n$; where α is the root in the interval $[0, 1]$ of:

$$\alpha^2 \|s_g\|^2 + (1 - \alpha)^2 \|s_n\|^2 + 2\alpha(1 - \alpha) s_g^T s_n = \delta^2$$

- 3 if $\delta > \|s_n\|$ we set $s_{dl} = s_n$;

Algorithm (Computing DogLeg step)

```

dogleg( $s_g, s_n, \delta$ );
 $a \leftarrow \|s_g\|^2$ ;
 $b \leftarrow \|s_n\|^2$ ;
 $c \leftarrow \|s_g - s_n\|^2$ ;
 $d \leftarrow (a + b - c)/2$ ;
 $\alpha \leftarrow \frac{b - d + \sqrt{d^2 - ab + \delta^2 c}}{b - \delta^2}$ ;
 $s_{dl} \leftarrow \alpha s_g + (1 - \alpha) s_n$ ;
return  $s_{dl}$ ;

```

Solving

$$\alpha^2 \|s_g\|^2 + (1 - \alpha)^2 \|s_n\|^2 + 2\alpha(1 - \alpha) s_g^T s_n = \delta^2$$

we have that if $\|s_g\| \leq \delta \leq \|s_n\|$ the root in $[0, 1]$ is given by:

$$\Delta = \|s_g\|^2 + \|s_n\|^2 - 2s_g^T s_n = \|s_g - s_n\|^2$$

$$\alpha = \frac{\|s_n\|^2 - s_g^T s_n - \sqrt{(s_g^T s_n)^2 - \|s_g\|^2 \|s_n\|^2 + \delta^2 \Delta}}{\Delta}$$

to avoid cancellation the computation formula is the following

$$\begin{aligned} \alpha &= \frac{1}{\Delta} \frac{\|s_n\|^4 - 2s_g^T s_n \|s_n\|^2 + \|s_g\|^2 \|s_n\|^2 - \delta^2 \Delta}{\|s_n\|^2 - s_g^T s_n + \sqrt{(s_g^T s_n)^2 - \|s_g\|^2 \|s_n\|^2 + \delta^2 \Delta}} \\ &= \frac{\|s_n\|^2 - \delta^2}{\|s_n\|^2 - s_g^T s_n + \sqrt{(s_g^T s_n)^2 - \|s_g\|^2 \|s_n\|^2 + \delta^2 \|s_g - s_n\|^2}} \end{aligned}$$

References

- 1 J. Stoer and R. Bulirsch
Introduction to numerical analysis
Springer-Verlag, Texts in Applied Mathematics, 12, 2002.
- 2 J. E. Dennis, Jr. and Robert B. Schnabel
Numerical Methods for Unconstrained Optimization and Nonlinear Equations
SIAM, Classics in Applied Mathematics, 16, 1996.