

# AN EASY FINITE ELEMENT IMPLEMENTATION OF THE OBSTACLE PROBLEM FOR POISSON'S EQUATION

ED BUELER

Suppose an elastic membrane is attached to a flat wire frame which encloses a region  $\Omega$  in the plane. Suppose this membrane is subject to a distributed load  $f(x, y)$ . The equilibrium position  $z = u(x, y)$  of the membrane (assuming small displacements, etc.) solves the Poisson problem

$$(1) \quad -\Delta u = f \text{ on } \Omega, \quad u|_{\partial\Omega} = 0,$$

where  $z = 0$  is the plane of the wire frame.

Now suppose that an obstacle is placed underneath the membrane. Specifically, suppose the obstacle has a continuous and differentiable surface  $z = \psi(x, y)$  and that  $\psi|_{\partial\Omega} \leq 0$ . The problem now becomes to find the region  $R$  where  $u$  coincides with  $\psi$  and to solve  $-\Delta u = f$  in the complementary region  $\Omega \setminus R$ . That is, the problem is to find the minimal energy configuration of the membrane “stretched over” the obstacle. Figure 1 shows an example where  $\Omega$  is the disc of radius two centered at the origin, the obstacle is the sphere of radius one centered at the origin (in  $\mathbf{R}^3$ , that is) and  $f \equiv 0$ .

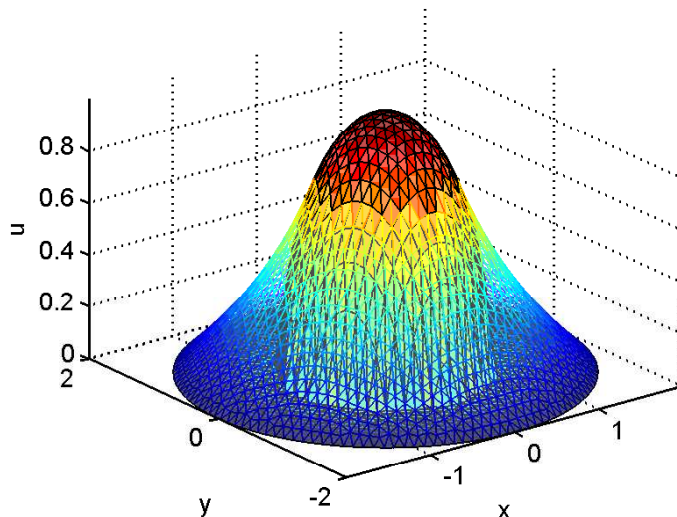


FIGURE 1. A membrane stretched over an obstacle. Figure produced by `obstacle.m`.

In particular the problem is now of *free boundary* type, that is, we seek to impose  $u = \psi$  and tangency ( $\nabla u = \nabla \psi$ ) on the boundary of  $R$ , but determining the location of that boundary is part of the problem.

The “free boundary” description of the problem suggests the difficulty one has in constructing a finite difference or other local approximation to this PDE problem. Our immediate goal, therefore, is to formulate this obstacle problem “globally” as a calculus-of-variations (minimization) problem in a closed, convex subspace of a function space.

**Definition.** Let

$$K_\psi = \left\{ v \in H_0^1(\Omega) \mid v \geq \psi \right\}.$$

Recall that  $H_0^1(\Omega)$  is the space of (weakly) differentiable functions  $v$  on  $\Omega$  for which  $\int_\Omega |v|^2 < \infty$ ,  $\int_\Omega |\nabla v|^2 < \infty$  and  $v|_{\partial\Omega} = 0$  (in a trace sense).

**Lemma 1.**  $K_\psi$  is closed and convex.

*Proof.* Suppose  $v_j \rightarrow v$  in  $H_0^1(\Omega)$  for  $v_j \in K_\psi$ , but suppose  $v \notin K_\psi$ . That is, suppose  $v < \psi$  on a set of positive measure. By standard Lebesgue measure methods there is  $\epsilon > 0$  and a positive measure set  $A \subset \Omega$  so that  $v \leq \psi - \epsilon$  on  $A$ . But then

$$\int_\Omega |v_j - v|^2 \geq \int_A |v_j - v|^2 \geq \int_A |\psi - v|^2 \geq \epsilon^2 m(A) > 0,$$

a contradiction. (Note  $v_j \rightarrow v$  in  $H_0^1(\Omega)$  implies  $v_j \rightarrow v$  in  $L^2(\Omega)$ .) Thus  $K_\psi$  is closed.

Now consider  $0 \leq \lambda \leq 1$  and suppose  $v, w \in K_\psi$ . Then  $\lambda v + (1 - \lambda)w \in H_0^1(\Omega)$  as  $H_0^1(\Omega)$  is a vector space. But furthermore

$$\lambda v + (1 - \lambda)w \geq \lambda\psi + (1 - \lambda)\psi = \psi$$

because  $\lambda, 1 - \lambda \geq 0$ . Thus  $\lambda v + (1 - \lambda)w \in K_\psi$  and  $K_\psi$  is convex.  $\square$

**Remark.** One such “closed and convex” proof is enough as these properties become obvious after a while. I will omit further proofs of closed-ness and convexity (of such sets).

**Definition.** For  $v \in H_0^1(\Omega)$  define the functional

$$I[v] := \int_\Omega \frac{1}{2} |\nabla v|^2 - f v.$$

We can now give two equivalent weak formulations of the obstacle problem. First, we seek  $u \in K_\psi$  such that

$$(2) \quad I[u] \leq I[v] \quad \text{for all } v \in K_\psi.$$

Alternatively, we seek  $u \in K_\psi$  such that

$$(3) \quad \int_\Omega \nabla u \cdot \nabla (v - u) \geq \int_\Omega f(v - u) \quad \text{for all } v \in K_\psi.$$

Condition (2) is called a *minimization* formulation and (3) a *variational inequality* formulation.

**Lemma 2.** The functional  $I[u]$  is strictly convex, that is, if  $0 \leq \lambda \leq 1$  and if  $v, w \in K_\psi$  then

$$I[\lambda v + (1 - \lambda)w] \leq \lambda I[v] + (1 - \lambda)I[w],$$

and if  $0 < \lambda < 1$  and  $v \neq 0$  or  $w \neq 0$  then  $I[\lambda v + (1 - \lambda)w] < \lambda I[v] + (1 - \lambda)I[w]$ .

*Proof.* Compute

$$\begin{aligned} I[\lambda v + (1 - \lambda)w] &= \frac{\lambda^2}{2} \int |\nabla v|^2 + \lambda(1 - \lambda) \int \nabla v \cdot \nabla w + \frac{(1 - \lambda)^2}{2} \int |\nabla w|^2 \\ &\quad - \lambda \int f v - (1 - \lambda) \int f w. \end{aligned}$$

Because  $2a \cdot b = |a|^2 + |b|^2 - |b - a|^2$  for  $a, b$  in an inner product space,

$$\lambda(1 - \lambda) \int \nabla v \cdot \nabla w = \frac{\lambda^2}{2} \int |\nabla v|^2 + \frac{(1 - \lambda)^2}{2} \int |\nabla w|^2 - \frac{1}{2} \int |\lambda \nabla v - (1 - \lambda) \nabla w|^2.$$

Thus

$$\begin{aligned} I[\lambda v + (1 - \lambda)w] &= \lambda^2 \int |\nabla v|^2 + (1 - \lambda)^2 \int |\nabla w|^2 - \lambda \int f v - (1 - \lambda) \int f w \\ &\quad - \frac{1}{2} \int |\lambda \nabla v - (1 - \lambda) \nabla w|^2 \\ &\leq \lambda^2 \int |\nabla v|^2 + (1 - \lambda)^2 \int |\nabla w|^2 - \lambda \int f v - (1 - \lambda) \int f w \\ &\leq \lambda \int |\nabla v|^2 + (1 - \lambda) \int |\nabla w|^2 - \lambda \int f v - (1 - \lambda) \int f w \\ &= \lambda I[v] + (1 - \lambda) I[w]. \end{aligned}$$

If  $v \neq 0$  then  $\int |\nabla v|^2 > 0$  by Poincaré's inequality; similarly for  $w$ . Thus if  $0 < \lambda < 1$  and either  $v \neq 0$  or  $w \neq 0$  then the second inequality above is strict.  $\square$

**Corollary 3.** *The minimizer of  $I[v]$  over  $K_\psi$ , if it exists, is unique.*

**Proposition 4.**  *$u \in K_\psi$  solves (2) if and only if it solves (3).*

*Proof.* Suppose (2). For  $0 < \epsilon < 1$ ,

$$0 \leq I[u + \epsilon(v - u)] - I[u] = \epsilon \left( \int_\Omega \nabla u \cdot \nabla(v - u) - \int_\Omega f(v - u) \right) + \frac{\epsilon^2}{2} \int_\Omega |\nabla(v - u)|^2$$

for any  $v \in K_\psi$ ; note that the convexity of  $K_\psi$  has been used. Thus

$$0 \leq \lim_{\epsilon \rightarrow 0^+} \frac{I[u + \epsilon(v - u)] - I[u]}{\epsilon} = \int_\Omega \nabla u \cdot \nabla(v - u) - \int_\Omega f(v - u),$$

that is, (3).

Now suppose  $u$  solves (3). Let  $v \in K_\psi$ . Let  $f(\epsilon) = I[u + \epsilon(v - u)]$  for  $0 \leq \epsilon \leq 1$ . Note  $f(0) = I[u]$  and  $f(1) = I[v]$  and  $f$  is continuous. Now we calculate  $f'(\epsilon)$  for  $0 < \epsilon < 1$ :

$$\begin{aligned} f'(\epsilon) &= \lim_{h \rightarrow 0} \frac{I[u + (\epsilon + h)(v - u)] - I[u + \epsilon(v - u)]}{h} = \lim_{h \rightarrow 0} \frac{I[w + h(v - u)] - I[w]}{h} \\ &= \int_\Omega \nabla w \cdot \nabla(v - u) - \int_\Omega f(v - u) + \frac{\epsilon}{2} \int_\Omega |\nabla(v - u)|^2 \\ &= \int_\Omega \nabla u \cdot \nabla(v - u) - \int_\Omega f(v - u) + \frac{3\epsilon}{2} \int_\Omega |\nabla(v - u)|^2 \\ &\geq \int_\Omega \nabla u \cdot \nabla(v - u) - \int_\Omega f(v - u) \geq 0, \end{aligned}$$

where  $w = u + \epsilon(v - u)$ . Thus  $f(\epsilon)$  is nondecreasing. If  $f$  is constant then  $I[u] = I[v]$ ; by the corollary this is a contradiction unless  $u = v$ . Thus  $f(1) > f(0)$  and  $u$  minimizes  $I[v]$ , that is, (2).  $\square$

It is shown in section 8.4.2 of [5] that a unique solution to (2) exists because  $I[v]$  is a coercive (and strictly convex) functional. The argument requires Sobolev embedding and weak convergence and thus, though it is not particularly hard, the proof is outside the scope of these notes.

We now solve the obstacle problem numerically. In fact we will allow arbitrary Dirichlet boundary conditions, not just “ $u|_{\partial\Omega} = 0$ ”. We use preexisting finite element MATLAB tools [1, 2, 7] for linear, unconstrained problems, but we adapt them to perform constrained point over-relaxation.

The program **obstacle** below uses a triangulation of a region  $\Omega$ . Such a triangulation can be generated by **distmesh2d** [7]. The data describing the triangulation consists of a list of  $N$  node locations **p** (an  $N \times 2$  array of real coordinates), and a list of  $M$  triangles (an  $M \times 3$  array of indices into **p**). **obstacle** requires the user to define MATLAB functions **psi**, the obstacle; **g**, the boundary conditions for the Poisson problem; **f**, the nonhomogeneity in the Poisson problem; and **fd**, a signed distance function describing  $\Omega$ . All of these functions can be anonymous (in MATLAB 6.5 or later) or **inline**. Finally, **obstacle** takes a convergence tolerance **tol** and a triangulation feature size **h0**.

**obstacle** performs constrained point over-relaxation [6]. In the context of a linear PDE which is discretized into matrix form  $Ax = b$ , “successive over-relaxation” is an acceleration of the Gauss-Seidel iteration [3].

In fact, **obstacle** asks **poissonDN** to assemble the stiffness matrix  $A$  and the load vector  $b$  for the corresponding unconstrained problem. Then  $A$  is decomposed  $A = D + L + U$  for Gauss-Seidel and the solution is found by constrained iteration. The initial “guess” is  $\max\{\psi, \tilde{u}\}$  where  $\psi$  is the obstacle and  $\tilde{u}$  is the solution to the unconstrained problem. The over-relaxation parameter  $\omega$  has been tuned by (very little) trial and error to the value  $\omega = 1.75$ . (Reference [6] gives no quantitative advice or theory on this tuning though it otherwise describes the algorithm completely. Tuning advice may well exist in the literature.)

**Example 1.** As a first example we do a case where the exact solution is known. Consider the problem

$$-\Delta u = 0, \quad u|_{\partial\Omega} = 0$$

on the disc  $\Omega$  of radius 2 centered at the origin ( $\Omega = \{(x, y) | x^2 + y^2 < 4\}$ ) subject to the constraint  $u \geq \psi$  where

$$\psi(x, y) = \begin{cases} \sqrt{1 - x^2 - y^2}, & x^2 + y^2 < 1, \\ 0, & \text{otherwise.} \end{cases}$$

That is, suppose that the membrane is attached to a wire circle of radius 2 and is stretched over a ball of radius one. See figure 1.

In this case the problem is fully radial and  $u = u(r)$ . Thus  $\Delta u = u_{rr} + r^{-1}u_r$  and if  $u > \psi$  then  $u(r) = -A \ln r + B$  for unknown  $A, B$ . If the free boundary is at position

$r = a$  then we seek  $a, A, B$  satisfying the nonlinear equations

$$u(a) = \psi(a), \quad u'(a) = \psi'(a), \quad u(2) = 0.$$

It is clear that  $0 < a < 1$ .

The equations reduce to a decoupled scalar nonlinear equation for  $a$ :

$$(4) \quad a^2(\ln 2 - \ln a) = 1 - a^2$$

and  $A = a^2(1 - a^2)^{-1/2}$ ,  $B = A \ln 2$ . A quick plot<sup>1</sup> shows one solution near 0.7. Application of `fzero` gives  $a = 0.69797$ ; also  $A = 0.68026$ ,  $B = 0.47152$ . The situation is illustrated in figure 2.

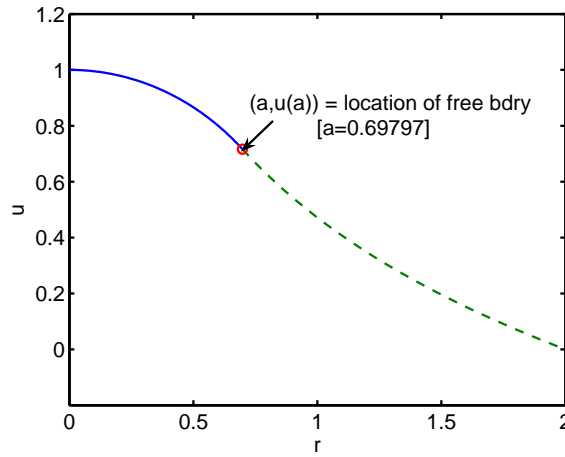


FIGURE 2. Solution of an exactly solvable obstacle problem (given numerical solution to equation (4)). Solid is  $u = \psi$ ; dashed is  $u > \psi$ .

**Example 2.** Now we solve the above problem by the finite element method using `obstacle`. The invocation of `obstacle` looks like

```
>> psi=@(p) sqrt(max(1-p(:,1).^2-p(:,2).^2,0));
>> fd=@(p) sqrt(sum(p.^2,2))-2; f=@(p) 0;
>> h0=0.1; [p,t]=distmesh2d(fd,@huniform,h0,[-2,-2;2,2],[]);
>> [uh,in,ierr]=obstacle(psi,f,f,1e-6,fd,h0,p,t);
```

We get figure 1, which looks right.

To determine the finite element location of the free boundary we do the following:

```
>> nogap=(uh==psi(p)); r=sqrt(p(nogap,1).^2+p(nogap,2).^2);
>> max(r(r<1.9))
```

(Note the solution  $u$  is in contact with the obstacle both inside the free boundary and along  $r = 2$ .) With  $h0 = 0.1$  as above we get 0.72239 as numerical location (at least, the furthest free boundary location from the origin) of the free boundary.

<sup>1</sup>E.g. `>> a=0:.01:1; plot(a,a.*a.*(log(2)-log(a)),a,1-a.*a)`

To see convergence of the finite element-computed location of the free boundary to the correct location, we run `obstacle` for triangulations with decreasing `h0`:

<code>h0</code>	0.4	0.2	0.1	0.07	0.05	0.04
<code>max(r(r&lt;1.9))</code>	0.82675	0.76132	0.72238	0.72135	0.71426	0.70979

Note that the finite-element location of the free boundary can not be expected to be closer than one triangle diameter to the exact free boundary. That is, we expect to see  $O(h)$  convergence for the above numbers, and this is what we roughly see in figure 3. Therefore it is reasonable to linearly extrapolate as  $h0 \rightarrow 0$  and we get 0.69603; compare to the exact value  $a = 0.69797$ .

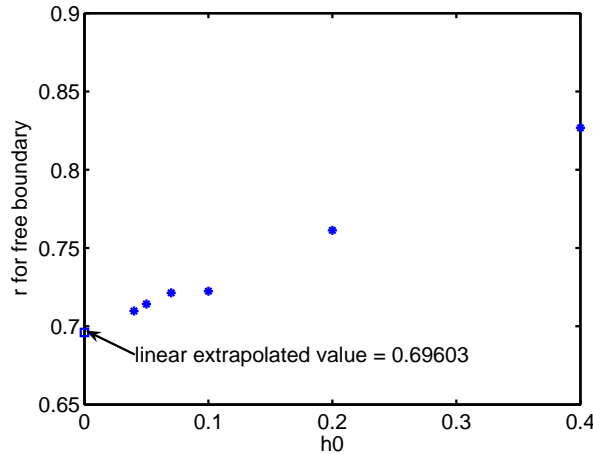


FIGURE 3. Convergence of the finite element location of the free boundary to the exact location.

As `h0` decreases the number of iterations of constrained point over-relaxation increases unpredictably. For example, `h0`= 0.4, 0.2, 0.1, 0.07, 0.04 correspond, respectively, to 51, 52, 52, 97, 281 iterations to achieve the desired convergence tolerance (`tol`=  $1.0 \times 10^{-6}$ ).

In fact, the constrained point over-relaxation method used here does not scale well. On the one hand, the use of Gauss-Seidel, a good smoother, suggests a multigrid technique. On the other hand, other techniques are available, especially *duality* [4, 6].

```
function [uh,in,ierr] = obstacle(psi,g,f,tol,fd,h0,p,t,varargin);
%OBSTACLE Solve the obstacle problem ...
%ELB 12/3/04

maxiter=300; omega=1.75; % omega found by trial and error

% use poissonDN to get unconstrained stiffness, load
[uh,in,A,b]=poissonDN(f,g,@(p)(0),fd,@(p)(-1),h0,p,t,varargin{:});
U=triu(A,1); L=tril(A,-1); d=diag(A); % U, L sparse
if any(d==0), error('stiffness matrix has zero on diagonal'), end;

% first guess is max(uh,psi)
```

```

N=sum(in>0); ps=zeros(N,1);
for j=1:N, ps(j)=feval(psi, p(find(in==j),:)); end
uold=max(uh(in>0),ps); unew=uold; omcomp=1-omega; ierr=[];

% iterate: constrained point over-relaxation
for l=1:maxiter+1
    Ux=U*uold;
    for j=1:N
        utemp=(b(j)-L(j,1:j-1)*unew(1:j-1)-Ux(j))/d(j); % Gauss-Seidel
        unew(j)=max(omcomp*uold(j)+omega*utemp,ps(j)); end
    er=max(abs(unew-uold)); ierr=[ierr er];
    if er<tol, break, end
    if l>maxiter, warning('max number of iterations reached'), break, end
    uold=unew; end

uh(in>0)=unew;
h=trimesh(t,p(:,1),p(:,2),uh); set(h,'FaceAlpha',0.3) % plot transparent
xy=[get(gca,'Xlim') get(gca,'YLim')]; hold on;
trisurf(t,p(:,1),p(:,2),psi(p),uh); axis([xy min(uh) max(uh)]); hold off;

```

## REFERENCES

- [1] E. BUELER, *Poisson's equation by the FEM, continued: general boundary conditions and error analysis*. [www.math.uaf.edu/~bueler/poissoncont.pdf](http://www.math.uaf.edu/~bueler/poissoncont.pdf), 2004.
- [2] ———, *Poisson's equation by the FEM using a MATLAB mesh generator*. [www.math.uaf.edu/~bueler/poissonv2.pdf](http://www.math.uaf.edu/~bueler/poissonv2.pdf), 2004.
- [3] R. L. BURDEN AND J. D. FAIRES, *Numerical Analysis*, Brooks/Cole, Pacific Grove, CA, seventh ed., 2001.
- [4] N. CALVO, J. DÍAZ, J. DURANY, E. SCHIAVI, AND C. VÁZQUEZ, *On a doubly nonlinear parabolic obstacle problem modelling ice sheet dynamics*, SIAM J. Appl. Math., 63 (2002), pp. 683–707.
- [5] L. C. EVANS, *Partial Differential Equations*, vol. 19 of Graduate Studies in Mathematics, American Mathematical Society, 1998.
- [6] R. GLOWINSKI, J.-L. LIONS, AND R. TRÉMOLIÈRES, *Numerical Analysis of Variational Inequalities*, vol. 8 of Studies in Mathematics and its Applications, North-Holland Publishing Co., Amsterdam, 1981. Translated from the French.
- [7] P.-O. PERSSON AND G. STRANG, *A simple mesh generator in MATLAB*, SIAM Review, 46 (2004), pp. 329–345.