# Mathematical Problem Solving

**Note:** The following is intended as the first step towards developing an interactive website where those interested in mathematical problem solving can exchange views and information (and above all, problems and solutions). I hope it will expand a lot during 2005. All feedback and correspondence from readers to shapiro@math.kth.se will be very welcome.
You should see a menu on the left side of the screen - if not, please click here.

# Introduction

**The problem seminar at KTH, 1972 - 1985:** As I have described elsewhere, my mathematical education was greatly influenced by a tradition of problem - swapping with my fellow students (especially Donald Newman, the maestro). When I came to Royal Institute in Stockholm as professor in 1972, it was an important priority of mine to organize a "problem seminar". I ran such a seminar until 1985. It was quite popular with the graduate students and also gave me a good insight into the creative talents of several of them who later did PhD theses under my direction. Since the seminars were informal, and sporadic participation was encouraged, a number of faculty members also sat in from time to time, sometimes contributing elegant solutions (or problems). Altogether this was a very positive experience for many people, and it has long been an ambition of mine to record some features of the seminar that might be useful or provide inspiration for others interested in starting such a seminar. These "features" comprise, besides a collection of problems we took up throughout the years, also lectures I gave concerning strategies for problem-solving.

The problem seminar was directed to graduate students, but let me also note here that I ran a more elementary version as an elective course for undergraduate engineering students several times, with very favorable results. The way the seminar was run was that I would hand out a list of problems (often having a unifying theme, like "inequalities"), and gave a lecture illustrating typical methods applicable to that type of problem. After a specified time period participants handed in their solutions (even partial solutions were welcome, and collaborative work was encouraged; moreover in many cases I assigned problems I had heard of, or thought of, to which I did not know the solution, so the whole enterprise had a pleasantly cooperative atmosphere). To the best moments from these times I still think back nostalgically. I would assort and classify the solutions received, and their authors were then called upon to present them at the blackboard.

It is my intention to present material from these seminars, in a piecemeal fashion at first (based on memories and notes) and hopefully later organize the material more cohesively. In this spirit let us at once get in medias res with one of my introductory "pep talks".

# Invariance

One of the key ideas in mathematical problem solving is exploiting invariance. This insight is so basic that it is seldom singled out as a topic to teach to students, rather it belongs (to use a metaphor of Lars Svensson, one of the stalwarts of the problem

seminar) to "the body language of mathematicians". What is involved is captured in the winged words "Without loss of generality..." Let us start with a very simple example: Say we are going to try to prove some theorem of Euclidean geometry involving a circle by using algebra and introducing Cartesian coordinates. We say "Without loss of generality, we may suppose the circle is centered at the origin." That is nice, instead of center (a,b) we have (0,0), simplifying some calculations. By what right do we make such an assumption? The answer is "invariance" - we are studying properties of figures which remain unchanged under rigid motions, so we lose nothing if we suppose the figure translated so as to place the center of our circle at the origin. Possibly we can go on to say "Without loss of generality (henceforth abbreviated Wlog) the radius of the circle may be taken equal to 1 "... this will be justified if the features of the problem under study are invariant w.r.t. scale change (dilation, or affine transformation); and so on - the more invariance, the greater the simplifying assumptions.

In some problems, already this first step solves the whole problem! Here is an example:

**Problem:** Prove that the inequality $(x + y)^p < x^p + y^p$ holds whenever x, y are positive numbers and $0 < p < 1$.

**Solution:** Wlog we may assume $x + y = 1$ (because the inequality in question remains invariant under the transformation of scaling $x \to tx$, $y \to ty$ with t positive). But, if $x+y = 1$, the inequality is obvious since then $x < 1$, $y < 1$ so $x < x^p$ and $y < y^p$, etc.

**Exercise:** Prove "Cauchy's inequality"

$$(ax + by + cz)^2 \le (a^2 + b^2 + c^2)(x^2 + y^2 + z^2)$$

where a,b,c,x,y,z are any non-negative real numbers (and likewise with more summands; their number was here taken as three merely to facilitate typography).

[Hint: Wlog each of the sums in the right hand member may be supposed equal to 1. Why?]

In the cases just discussed the transformations needed to "trivialize" the problem were simple and fairly obvious, but in other cases this is not so. Consider for example:

**Problem: Steiner's porism** Let A and B denote circles, one wholly enclosed by the other. A finite collection of circles $C_1$, $C_2$,..., $C_n$ is called a Steiner chain (w.r.t. A and B) if a) Each of the circles $C_i$ is tangent to both A and B, and b) Each of the pairs $(C_1, C_2)$, $(C_2, C_3)$,..., $(C_n, C_1)$ exhibit external tangency.

It is clear that, given A and B, and the integer n, we cannot in general expect that there will exist a Steiner chain with n circles. However, Jakob Steiner made the beautiful discovery that: If there exists a Steiner chain with n circles, then there exists such a chain regardless of the choice of $C_1$, provided only it is tangent to A and B.

Now, observe: the theorem is obvious if A and B are concentric. (Why?)

So, how nice it would be to be able to say "Wlog A and B are concentric."? But, is this true? It would be if there were a one-to-one transformation of the plane P of the figure

onto another plane P' such that: circles are carried to circles, and moreover the pair A, B is carried to a pair of concentric circles. Well, this is almost true: if we replace the Euclidean plane by a projective one it is, and that is good enough to do the business. I won't here go into the details, but leave as an exercise:

**Exercise:** Prove that any two disjoint circles can be carried to concentric ones by a suitable Moebius transformation z -> (az + b)/(cz + d).

## Another aspect of invariance: transformations

The strategy we employed in proving "Steiner's porism" is very useful, let's elaborate on this. An important theme in studying mathematical objects is transformations . They are ubiquitous, part of the very soul of mathematics, so much so that often we use them without being aware of them. For instance, one sometimes speaks of "the circle $x^2 + y^2 = 1$" ... but this is not a circle, it is an equation, and we have here to do with the wondrous phenomenon that geometric objects have algebraic counterparts and vice versa. This is taken for granted in an age when first year students learn "analytic geometry", but this was not so before Descartes established the correspondence. A common pattern in problem-solving is "Transform. Solve. Invert." This was developed and richly illustrated in an article of M. Klamkin and D. J. Newman. A very simple and amusing illustration is the problem to do arithmetic in the Roman system of numeration, e.g. compute the product VIII times V. The solution consists of

1. Transform (to decimal notation): The transformed problem is to compute 8 x 5.
2. Solve: 8 x 5 = 40 .
3. Invert: 40 = XL (Conclusion: VIII x V = XL.)

Indeed, when your pocket calculator computes 8 x 5, it uses this exact same approach, transforming from decimal to binary notation, solving and transforming back. From a logical point of view this is the same schema we employ when we solve a differential or difference equation by using Laplace transforms, or a geometric problem by using "analytic geometry".

Here is a more sophisticated example:

**Problem:** Let E be a measurable set of real numbers of positive measure. Prove that the set D := {x - y: x and y in E} contains an interval.

**Solution:** We'll transform the given information about sets to corresponding information about functions, by associating to E its "characteristic function" f, defined by

f(x) = 1 if x is in E, 0 otherwise.

Now, for fixed t, the function

g(x) = f(x+t)f(x)

equals 1 if both x and x+t are in E, and 0 otherwise. Thus, if this function g(x) is not identically 0, there is an x for which x and x+t are in E, and hence t is in D. From here on assume w.l.o.g. that E is bounded. Certainly g(x) is not identically 0 if its integral

w.r.t. x is positive. But

h(t):= Int[f(x+t)f(x) dx]

is a continuous function of t (This is the point! Can you prove it?) and, since h(0) equals the measure of E, which is positive, g is positive for all t with |t| sufficiently small; so those t belong to D, and the proof is finished.

## Another aspect of invariance: scaling

**Problem: Liouville's theorem**
Consider for example the famous "Liouville theorem": A bounded entire function is constant.

**Proof:** Suppose f is entire and $|f(z)| \leq M$ for all complex z. We'll assume as known the maximum modulus theorem, and its immediate consequence, Schwarz' lemma: If g is holomorphic in the unit disk D and |g| is bounded by 1, and g(0)= 0, then $|g(z)| \leq |z|$. (The point here is simply to train the instinct of the student, so as to grasp at once that the Schwarz lemma is the "high ground" for Liouville's theorem (and many others!).

So, define g(z) = (f(Rz) - f(0))/2M, where R is an arbitrary positive number. Since Schwarz' lemma is applicable to g we get

$|f(Rz) - f(0)| \leq 2M |z|$.

Hence, for any complex number w, setting z = w/R in the last inequality, we obtain

$|f(w) - f(0)| \leq 2M |w|/R$ .

Since this is valid for all positive R, we conclude that f(w) - f(0) vanishes identically, QED.

**Exercises.** Prove that a function harmonic in the whole plane and positive is constant. (Hint: look for a "finite version", i.e. a theorem about positive harmonic functions in the unit disk which, by appropriate scaling, would yield the desired result.

**Better hint:** use the Harnack trick, which in $R^3$ for example observes that, if B is the unit ball, and B* denotes the ball centered at the point x and with radius 1 - |x|, then for a positive harmonic function u in B, its integral over B* is ≤ its integral over B. Choose here x with |x| = 1/2, and apply the Gauss mean value property to the integrals over B and B* to deduce that $u(x) \leq 8u(0)$.

**Exercise:** Prove the analogous result also in $R^d$.

**Exercise:** Prove that an entire function f which satisfies $|f(z| \leq C(1 + |z|)^m$ for all complex z is a polynomial of degree at most m. And, a real harmonic function u on $R^d$ satisfying $u(x) \leq C(1 + |x|)^m$ (note that this is a one-sided estimate, with u, not |u|, in the left member!) is a polynomial in $x = (x_1,...,x_d)$ of degree at most m.

## Another aspect of invariance: impossibility proofs

We are familiar with "geometric constructions" using straightedge and compass, subject to certain strict ground rules as to how these may be used. And we know e.g. that certain constructions are impossible, such as inscribing a regular heptagon in a circle because it entails, algebraically, constructing the root of an irreducible cubic equation with integer coefficients, from a given segment of length one. (Galois theory teaches us that a necessary condition for constructibility of a segment of length x is that f(x) vanish for some nontrivial polynomial f with integer coefficients, irreducible over the rational field, and having a degree which is a power of 2.) Less familiar to most students, and very elegant and instructive, are problems of construction with more restricted means, such as allowing use only of a straightedge (ruler). It is remarkable that certain a priori very unlikely looking configurations can be constructed using a straightedge. To get warmed up, try this one:

Given a circle in the plane and a point P outside, it is required to construct a line through P tangent to the circle. (Caution! It is not permitted to place the ruler so that it touches P, and rotate it about P until it also touches the circle.)

Amazingly, this problem has a solution (nontrivial and very instructive - we'll come to this much later). But now consider the problem: Can one bisect a given segment using only a straightedge?

The answer is **NO**. To see why, consider a line L in the plane P, and two marked points A, B on it. It is desired to construct the midpoint M of the segment AB using the straightedge. Suppose we have found a procedure which works. Now, suppose we have a one-to-one mapping of plane P onto another plane P' which carries lines to lines, but which does not preserve the relation "M is the midpoint of the segment AB", in other words A, M, B are carried to points A', M', B' with A'M' unequal to B'M'. Then, this leads to a contradiction, because the construction of the midpoint in the plane P induces a construction in P' which also would have to lead to the midpoint of A'B'. (This is a profound insight, an "Aha" experience, and worth investing lots of time and energy in thinking it through carefully!!)

Well, can the mapping of P onto P' be supplied? Again, the answer is "almost" - we again have to work with a projective plane, but with care concerning "improper points" it works. It is a highly rewarding exercise to work out the details (a short course in projective geometry, really).

**Exercise:** Show it is not possible using only a straightedge to a) construct a line perpendicular to a given one b) construct a line parallel to a given one through a prescribed point, or c) to find the center of a given circle.

## Another aspect of invariance: invariants

There are still many other facets of the invariance concept, One such involves conserved quantities or "invariants" (here the "quantity" need not be a number; it might be a parity, or an algebraic entity.)

We look first at a deceptively simple example: Can a 5 x 13 rectangle be decomposed into finitely many pieces which may be reassembled into an 8 x 8 square?

"Obviously, no. Because the areas of the two figures are unequal." is the immediate reaction. But, suppose the concept of area were not known to us. It would be very difficult to demonstrate the impossibility. To appreciate this, consider the three dimensional analog: It is known that there are two convex polyhedra of equal volume such that neither can be decomposed into finitely many pieces which may be reassembled to give the other (Hilbert's third problem, solved by Max Dehn). Here it turns out that there are other "invariants" besides the volume that must match before such construction is possible. (This is in contrast to the planar situation: for two polygons of equal area one can always dissect and reassemble the one to obtain the other. Try to prove this - start with the case of two rectangles.)

A substantial portion of higher mathematics is concerned with identifying, for a given class of objects and a given family of transformations of these objects, a "complete set of invariant quantities" the matching of which, for two objects in the class is necessary, or sufficient (or, preferably both) for them to be mutually transformable one to the other. (Example: The objects are n x n complex matrices with n distinct eigenvalues. Two matrices in this family are similar if and only if the two sets of eigenvalues coincide; if we allow multiple eigenvalues this criterion is necessary but no longer sufficient for similarity.What conditions are now necessary and sufficient?)

Later we shall study many problems which fit into this general framework. The present remarks are preliminary, just to give perspective.

**Exercise:** Given two annuli {a < |z| < b} and {a' < |z| < b'} in the complex plane, under what conditions is there a conformal map carrying one onto the other? Give a "complete set of invariants" attached to each annulus so that the matching of these lists (Hint: Well, here the lists contain only one entry) is necessary and sufficient for the conformal equivaalence of the annuli.

**Exercise:** Given two rectangles, under what conditions can one be mapped conformally on the other such that the vertices of the one are mapped to the vertices of the other?

# Symmetry

In attacking mathematical problems, one should always look for symmetries that can be exploited to simplify matters. (This is really just another facet of the "wlog" approach. There are several key words here: invariance, symmetries, automorphisms that may describe what one is looking for.) A typical exploitation of symmetry, which occurs in studying a boundary value problem, is to observe that the problem considered has rotational invariance and therefore use polar coordinates. Other kinds of symmetry are more subtle, such as the duality (or automorphism) inherent in projective plane geometry, interchanging the roles of points and lines and thus yielding for free a new theorem "dual" to any given one. In the case of differential equations, the discovery of a symmetry may enable one to reduce the number of variables e.g. transform from a partial to an ordinary differential equation, or from an ordinary differential equation to one of lower order.(These observations were the point of departure for the researches of Sophus Lie.)

**Problem:** Solve the heat equation $u_t = u_{xx}$ subject to the initial conditions that $u(x,0) = H(x)$, H denoting the "Heaviside function": $H(x) = 0$ for $x < 0$ and $1$ for $x > 0$.

**Solution:** If u is a solution, define $v(x,t) = u(sx, s^2 t)$ where s denotes a positive parameter. Then, $v_t = v_{xx}$ and $v(x,0) = u(sx,0) = H(sx) = H(x)$, so v satisfies the heat equation with the same initial conditions as u. If we take for granted that the problem has a unique solution then we can conclude $v = u$, that is $u(x,t) = u(sx, s^2 t)$ identically w.r.t. the parameter s. We have discovered a symmetry of u! Since this relation holds identically in t, x, and s we may set $s = t^{-1/2}$ and obtain

$$u(x, t) = u(x t^{-1/2}, 1)$$

showing that u has the form $f(x t^{-1/2})$ for some function f of one variable. Thus, exploitation of the symmetry in the initial value problem has narrowed our search to one for the univariate function f. The condition that $f(x t^{-1/2})$ satisfies the heat equation gives that f(y) satisfies the ordinary differential equation

$$f''(y) + (y/2)f'(y) = 0.$$

**Exercises:** Complete the above analysis and obtain the solution u. Can you prove the uniqueness of the solution that was assumed (possibly requiring imposition of some restrictions of regularity on the solution)?

## A tiling problem

The following problem illustrates the role of invariants (that is, conserved quantities) in impossibility proofs. We introduce it by recalling an old favorite:

From an 8 x 8 chessboard two diagonally opposite squares are removed. Can the resulting board be tiled using 31 dominoes (each domino being a 1 x 2 rectangle congruent to a couple of adjacent squares of the chessboard)?

The answer is "No", the reason being that the mutilated chessboard has unequally many black and white squares (since the two removed squares have the same color), whereas a domino always covers one white and one black square.

This solution is very elegant, especially if one considers that in posing the problem there is no need to speak of a chessboard with colored squares, it could as well have been a white board ruled by lines into an 8 x 8 grid of squares. Then the cleverness of the solution is to realize that one should color the squares in an alternating pattern. This suggests a strategy for similar but more intricate tiling problems.

**Problem:** Can a 10 x 10 chessboard be tiled with 25 "dominoes", each of side 1 x 4?

Here the most obvious invariant, the area, is right, but it can be shown that none the less the tiling is impossible for more subtle reasons. Suppose we can write a number in each of the 100 squares of the chessboard with the property that, no matter how we lay one of the dominoes so as to cover 4 squares, the sum of the numbers in the four covered squares equals zero; yet, the sum of all 100 numbers from the whole board is not zero. Then, obviously, the tiling is impossible.

So, how to assign numbers to the 100 squares satisfying the stated conditions? Here is one way: Define a(n) for n = 1,2,...,10 to be the sequence

1,1,1, -3, 1, 1, 1, -3, 1, 1 .

In this sequence any consecutive block of four sums to 0, but the sum of all ten is not 0. Now, write in the square of the chessboard that is in row m and column n the number a (m) x a(n) (So: we have invented for our purpose an appropriate coloring of the chessboard by the three "colors" 1, -3 and 9.

**Exercises:** Verify the suggested argument showing the tiling is impossible. Study the more general problem of tiling an A x B chessboard with rectangles of size a x b where a, b, A, B are positive integers; also multidimensional generalizations.

## Problem: Conway's game

Just as with the above tiling problem, many other puzzles and games can be analyzed (more precisely: it can be shown that certain initial configurations cannot be transformed into some specified one by the application of particular specified rules of play) by "arithmetization", that is by ingeniously introducing suitable numbers, per se totally foreign to the game in question, into the picture. A typical, and very beautiful example is the following game invented by John Conway. Consider the plane made into an infinite chessboard by horizontal lines 1 cm apart, and vertical lines likewise. One of the horizontal lines (which we call the equator) is singled out. We now place a finite number of coins on the chessboard, at most one in each ruled square, and all lying below the equator (but subject to no other resrictions). This is called an initial configuration. We now are allowed to transform this by "moves"of the following type: Whenever we find three adjacent horizontal squares of which precisely two consecutive ones are occupied, we may jump one of the coins over the other into the vacant square, removing in the process the jumped-over one. Thus, for instance if we denote by 1 an occupied square and by 0 a vacant one,whenever we see 1 1 0 we may replace it by 0 0 1 (a jump to the right), the status of all other squares being unaffected. Likewise we may make a jump to the left: 0 1 1 -> 1 0 0, as well as a jump upwards, whereby

0
1
1

becomes

1
0
0.

Now, the challenge is: What is the largest integer k such that, for suitable initial configuration and subsequent evolution according to the specified jumping and removal rules, a coin can reach a position k units above the equator?

**Exercise:** Try your hand at this. Do you think k is finite, and if so have you a guess as to its value?

**Arithmetization:** Suppose we write real numbers a(m,n), one in each square of the "chessboard" (the squares being assigned integer coordinates in the natural way). To each configuration i.e. subset of the lattice squares that are occupied, we can then associate an "energy", that is, the sum of the numbers in the occupied squares. Let us

examine what happens to the energy by a jump: By a jump to the right the three numbers a(m,n) a(m+1,n) 0 get replaced by the three numbers 0 0 a(m+2, n). Suppose now that

a(m+2,n) ≤ a(m,n) + a(m+1,n) for all m and n.

Then, a jump to the right cannot increase the "energy" of a configuration.

**Exercise:** Write the conditions on the a(m,n) such that a jump to the left cannot increase the energy of a configuration; and, the conditions that a jump upward cannot increase the energy of a configuration.

Now, playing with these ideas, and a lot of patience, you have a fighting chance to answer the above question: If a(0,k) is larger than the energy of any initial configuration (consisting of squares below the equator) then (assuming that the a(m,n) have been rigged so that legal jumps never increase the total energy) it will be established that level k is not attainable.

We have spoken at some length of "invariance" (and related concepts: symmetry,invariants) in problem solving. Fired up with enthusiasm, the fledgling problem solver is bound to be disappointed, though, when she finds that most problems cannot be disposed of just by a tour de force of this (or for that matter any other) general point of view. Problem solving is an art,and as with other art forms there are no algorithms for producing masterpieces. But experience (at any rate, my experience) teaches that it does not hurt, and frequently helps, to have trained one's ability to recognize certain features (like symmetries, obvious or concealed), so let us continue in this vein and explore other thematic ideas in problem solving.

# Generalizing vs. specializing

Part of the conventional wisdom offered to would be problem solvers by great masters like Hilbert and Polya is: Specialize. Look for a special case of the problem which you can solve (perhaps, one where symmetry can be exploited). [We might characterize this approach as "With loss of generality..."]. The exploitation of analogies is very fruitful. Indeed, this is probably the most valuable simple general rule one can offer. The idea is eloquently developed in the writings of (among others, but especially) Polya, with countless examples and I won't elaborate further here, although we will encounter it again and again when we analyze concrete problems later. On the other side of the coin marked "specialize" is written "generalize" and this idea, too, is very important to the problem solver. Very often we see an object (whether in the real world or the world of ideas) which we do not recognize because we are looking at a cross section of some larger (perhaps familiar) object, or we are focusing on some small detail from a larger picture that we easily would recognize. It is for this reason that mathematicians are so fond of generalizing, of "building machinery" to gain insight into the simpler, more special object. Let's look at an example (the subject of a famous lecture of Georg Polya):

**Problem:** We have a large Edam cheese. Into how many pieces can it be divided by five straight cuts? (It is assumed the pieces remain immobile during the cutting process).

Or, in crass geometric terms: What is the maximum number of components into which

Euclidean space $R^3$ can be partitioned by 5 planes?

This is a piquant question, because it goes just a little bit beyond where most people's powers of geometric visualisation end. Most people see straight away that 3 planes can divide $R^3$ into 8 components ("octants"), and also that the fourth plane cannot intersect all the octants so we get less than 16 components...but there things start getting fuzzy. It is a picture that must be seen with "the mind's eye".

We can start reasoning by analogy: Into how many components can one-dimensional Euclidean space be partitioned by n points? Obviously n+1. Good, let's try the next case, partitioning the plane by lines. Two lines give four components ("quadrants"). A third line cannot intersect all four of these, but it can intersect (and hence divide) three of them. So, the fourth line contributes three new components, giving seven in all.

At this point one might have a small "Aha" experience. Each of the first two lines doubles the number of components: 1, 2, 4. The third line fails to do so. Why? Well, if a third line generates, say, r new components that means it has intersected r of the components that were present previously. So, if we travel along this third line we will pass r-1 "checkpoints" where we pass from one of the previous (generation two) components to a neighboring one. Each checkpoint is the intersection of this third line with one of the earlier two, so r-1 is at most 2, hence r is at most 3 (and obviously does attain the value 3 if the lines are "in general position". So, 3 lines generate 4 + 3 components.

The picture clarifies. L'air est pur, la route est large! The fourth line will intersect 4 of the components from generation 3 (there are 3 "checkpoints") so we get 4 new components, 11 in all. It's easy now to continue. Denoting by C(n) the number of components formed by n lines, we get the table

```
n   C(n)
0   1
1   2
2   4
3   7
4   11
5   16
```

and in general C(n) = C(n-1) + n.

**Exercise:** Find a formula for C(n).

Thus armed, we can tackle the original problem in $R^3$. Let D(n) denote the number of components generated by n planes. Thus, D(0) = 1, D(1) = 2, D(3)" = 8. Now comes the fourth plane, call it P. Suppose it intersects r of the 8 octants from the first 3 planes, then it will generate r new components. If we restrict our gaze to the plane P itself, we will see it partitioned into r parts by the 3 lines of intersection with the first three planes. So, r is none other than C(3) = 7 from the previous analysis, and D(4) = 8 + 7 = 15. And D(5) = 15 + C(4) = 15 + 11 = 26, and our problem is solved.

**Exercise:** Find a formula for D(n). Also, find recursion formulae enabling one to solve the analogous problem for any number of hyperplanes, in any number of dimensions.

This solution illustrates the power of "building machinery". Suddenly we have become able to "see" what happens if say we divide a four dimensional Edam cheese by 10 cuts, etc. As Jeff Rauch once wrote "Formulas are smarter than people."

Before moving on let us draw an important lesson from the above solution. We started out from an analogous (and very easy) problem - partition of one dimensional space by points, then progressed to the partition of two dimensional space by lines. Strictly speaking, these are not special cases of the given problem. Yet, they provided a "scaffolding" upon which we unerringly build our way to the solution of the problem of partition of $R^3$ by planes, and so on to higher dimensions.

Let's give a few more examples of the power of generalization.

**Example: cubic case of Fermat's equation**

Kummer proved the unsolvability in positive integers of the Fermat equation

$$x^3 + y^3 = z^3$$

by the fruitful idea of allowing x,y, and z to range over the more general class J of numbers of the form a + bw where a,b are ordinary integers and w = -1/2 + (1/2)Sqrt [-3] is a complex cube root of 1. The reason is, the polynomial $x^3 + y^3$ becomes factorable if we allow this wider class of "integer" coefficients. But the price that has to be paid is, a whole new number theory involving primes, factorization etc. had to be developed for this class J of generalized integers.

**Example: Trigonometric polynomials**

A (real) trigonometric polynomial of degree (at most) n is, by definition, a linear combination with real coefficients of the 2n + 1 functions

1, cos t, sin t, cos 2t, sin 2t, ..., cos nt, sin nt .

These are very important functions. They also illustrate what I mean by a "cross section". Namely, a trigonometric polynomial of degree n can be understood as the function arising from an ordinary polynomial P(x,y) of degree n in two variables by restricting it to the unit circle (on which we introduce a coordinate t as arc length measured counterclockwise from the point (1, 0).

This enables one to "see" readily certain properties of trigonometric polynomials (t.p.) For example their 2 pi periodicity is built into this interpretation. And the fact that a nontrivial t.p. of degree n can have at most 2n roots is equivalent to: If a polynomial P (x,y) of degree n, and the quadratic polynomial $x^2 + y^2$ - 1 have more than 2n common roots, then P is divisible by $x^2 + y^2$ - 1. Although the latter proposition is not more elementary than the former, it gives new insight since it is a special case of a very fundamental theorem of algebra (Bezout's theorem). More generally, other transcendental functions (spherical and elllipsoidal harmonics) can fruitfully be interpreted in terms of restrictions of ordinary multivariate polynomials to spherical or ellipsoidal surfaces.

**Example: Mathematical induction**

Another well known example of the power of generalization arises in connection with proofs by induction. The scheme is that we establish, for some class of propositions {P(n)}, the truth of P(1) and of the implication P(n) => P(n+1). Now, in proving this implication it well may be advantageous to replace P(n) by a stronger assertion Q(n) (meaning, we aim to prove a more general theorem). The reason is, if we compare the two assertions P(n) =>P(n+1) and Q(n)=>Q(n+1), the second one may be easier to prove. For, to be sure, the conclusion Q(n+1) asserts more than P(n+1) and thus is, in principle, harder....but, in proving it we may exploit the hypothesis Q(n), which gives more to work with than the weaker P(n). Thus, finding the best choice of inductive hypothesis is a fine art. We shall see in due time good illustrations of this point.

# Looking at an extreme situation

A very useful technique in analyzing problems is searching for an extreme situation. Of course, often a problem is posed as an extremal problem, to maximize or minimize something... but, even if it is not, it may be fruitful to set up some auxiliary extremal problem to achieve our purpose.

### Example: mean value theorem

For example, recall the proof of the mean value theorem of calculus: If f is a real valued continuous function on [a,b] and differentiable on (a,b) there is a point t in (a,b) such that $f'(t) = (f(b) - f(a))/(a - b)$.

**Proof:** Wlog we may assume $f(a) = f(b) = 0$ (this can be achieved by adding a suitable linear function to f). Now, either f is identically 0 on [a,b] (in which case the proof is finished) or it attains (wlog) some positive value. In this case its maximum on [a,b] is attained at some point t of (a,b) and $f'(t) = 0$.

### Example: convex hulls

A second illustration of the "extremal" technique is the usual proof of the following theorem: Let E be a subset of $R^n$. Any point in the convex hull of E is in the convex hull of some subset of E containing at most n+1 points.

Outline of proof. Let x be in the convex hull of E. Then, there are, for some positive integer m, points $e_1$, $e_2$,..., $e_m$ of E and positive numbers $p_1$, $p_2$,..., $p_m$ with sum 1 such that

$x = $ Sum $[p_i e_i$ , i from 1 to m]. (*)

Let us assume m has been chosen as the smallest possible integer for which such a representation holds. It will suffice to suppose m is greater than n+1 and derive a contradiction.

So, suppose m > n+1. Then, the vectors $e_1 - e_m$, $e_2 - e_m$,..., $e_{m-1} - e_m$ are linearly dependent so there are real numbers $a_j$ (j = 1,2,...,m-1) not all 0 such that Sum $[a_j (e_j - e_m)$, j = 1 to m-1] = 0, hence: There exist real numbers $b_1$, $b_2$,...,$b_m$ not all 0, with sum 0 and satisfying Sum$[b_i e_i$, i = 1 to m] = 0.

From this last equation and (*) we see that for every real t

$x = Sum[(p_i + tb_i) e_i, i = 1,...,m]$.

Now, define f(t) to be the minimum of the numbers $p_i + tb_i$. It is continuous, and positive for t=0. Since at least one of the $b_i$ is negative f(t) takes a negative value, and hence vanishes for some value of t, say t = t*. Then the numbers $p_i$ - t* $b_i$ are all nonnegative, and at least one (but not all! - why?) equals 0. But, this is a contradiction since it gives a representation of x of type (*) with fewer than m summands.

There are many famous proofs in the mathematical literature based on extremality, for example the solution of the Dirichlet boundary value problem by Perron's method; or,by the original ("Dirichlet's principle") method which, although initially conducted non-rigorously was later rendered rigorous by Lebesgue. Likewise, the celebrated Caratheodory proof of the Riemann mapping theorem, and (to take a more recent example) the de Branges proof of the Stone - Weierstrass approximation theorem.

# Physical intuition

Poincaré has written that physics not only suggests problems to the mathematician, it suggests also in many cases a method of solution. Pushed a step further, sometimes a purely mathematical problem can be solved (or, at least a fruitful line of attack initiated) by physical reasoning. Here is a well known, elegant example.

**Problem:** Given a nondegenerate triangle in the plane with vertices A, B, C, find a point P the sum of whose distances to A, B, C is minimum. More precisely, show there is such a point P, and that each of the angles APB, BPC, CPA equals 120 degrees.

**Solution:** The existence of (at least one) such point is easy to show, and left to the reader. Now, assume A, B, C are points on a horizontal wooden table. Drill holes at A, B, C and through each hole pass a string, attaching a weight to the lower end of the string and joining the upper ends of the strings. We assume the three weights are equal, and that there is no friction anywhere between the string and the wood. The point on the table where the three strings join is moveable, and gravitates to an equilibrium position Q. The equilibrium position will be such that the potential energy of the system of three weights is minimal for the constraints, i.e. their center of gravity is as low as possible which means that the sum of the lengths QA, QB, QC is minimal. Thus, Q is the point (denoted P above) which we were looking for. Since the three equal forces at Q are in static equilibrium, the three angles AQB, BQC, CQA must be equal.

**Exercises:**

**a)** Investigate the geometric consequences of similar experiments with more holes in the table, and where we attach arbitrary weights (not necessarily all equal) to the strings.

**b)** Consider a convex polyhedron (say, in three dimensional Euclidean space) i.e. a bounded figure which is an intersection of closed halfspaces, made of some homogeneous material. Prove there is at least one face on which it can stand stably (or,

in purely mathematical terms, if perpendiculars are drawn from its center of gravity to the plane of each face, the foot of at least one of these perpendiculars lies interior to a face).

# Epilogue: The art of problem posing

This is perhaps a good time to say a few words about "problems" as such. Where do they come from? Of course one cannot really answer this. Puzzles and conundrums are as old as mankind. Here are some thoughts regarding my choice of problems for the seminar.

Some problems are "old chestnuts" that, in different variants, recur in anthologies, competitions etc. They usually have the nice feature that they are easily understood, and can be solved without use of advanced mathematical studies. It's always good to pepper a problem seminar with a few of these, since the beginner is on the same footing as more experienced students - or professors. These may be easy, like the following two "quickies", or difficult.

Two "quickies":

**Problem:** A village consists of 6 houses. Each house is coupled to all of the others by a telephone line that is either red or blue. Prove there are three houses which are mutually coupled by lines of the same color.

**Problem:** Each point in $R^2$ is assigned a color - red, white or blue. Prove there are two points distance 1 apart, of the same color.

Even such simple problems as these can stimulate interesting discussion. Suitably generalized they lead into deep aspects of combinatorics (graph theory, "Ramsey theory" etc.)

Some problems suggest themselves by analogy, or generalization, from known (even trivial) results. Here are some examples.

Every positive number has a positive square root. Well, let's look at other collections of "positive" entities and see if this remains true:

"Every n x n Hermitean matrix which is positive (in the sense that its eigenvalues are nonnegative) has a positive square root."

This is also true, but not trivial. It is an important result in matrix theory (and extends to positive operators on Hilbert space).

How about: Sqrt[A + B] < Sqrt[A] + Sqrt[B]

for positive matrices?

I think students should cultivate the habit of playing around with generalizations like this. Another example of "problem generation" by analogy is discrete potential theory:

A discrete harmonic function on the set $Z^2$ (of "lattice points" in the plane, i.e. points with integer coordinates) is, by definition, a function u from $Z^2$ to R such that the mean value property

u(m,n) = (1/4)[u(m+1,n) + u(m-1,n) + u(m, n+1) + u(m, n-1)]

holds.

Some questions: Does the Liouville theorem hold here, i.e. if a discrete harmonic function on all of $Z^2$ is bounded must it be constant? What if it is everywhere positive, or of polynomial growth? What is the appropriate definition of discrete harmonic function on $Z^k$? What is the appropriate definition of harmonicity for a "subdomain" of $Z^k$? What is the analog of Dirichlet's boundary value problem? Is it solvable? Is the solution unique? These questions turn out to be not merely sterile generalizations, but are important in the theory of random walk on the discrete structure $Z^k$, in an analogous way to that in which classical potential theory relates to Brownian motion in $R^k$.

A further example (which hasn't been tested in practice, though): What would be a reasonable definition of a complex valued function from $Z^2$ to C satisfying the discrete Cauchy-Riemann equations? Try to formulate and prove or disprove analogs of classical theorems about analytic function for these discrete functions.

In actual fact, most of the problems I selected for the seminar were gotten by choosing some thematic subject (a known, but not too familiar branch of mathematics, from the perspective of beginning graduate students) and breaking it down into a logical succession of problems. The inspiration for this is the great two volume work of Polya-Szego, who worked out beautiful sequences of problems on such themes as power series and generating functions, polynomials and trigonometric polynomials, the maximum modulus principle, etc. Here are some of the themes I used:

- Convexity (convex functions and sets)
- Inequalities
- Measure theory
- Integral geometry
- Projective geometry
- Matrix theory

In all of these areas there are lots of "offbeat" topics that lend themselves well to problems. Moreover, a subject like "integral geometry" despite its importance and elegance, is scarcely touched on in a standard graduate education, so the problem seminar also served as a forum to repair gaps in the curriculum, so that students learned the existence of such a subject. And started to ponder how to define a measure on sets of lines in the plane, which would have reasonable invariance properties...

How to reach me:
Address: Department of Mathematics, KTH, 100 44 Stockholm
Visiting address: Lindstedtsv. 25

This page is maintained by Harold Shapiro, shapiro@math.kth.se,