

CS1114: Study Guide 2

This document summarizes the topics we've covered in the second part of the course. They closely follow the class slides, but in a more organized and centralized form. Please refer to the slides for more details.

1 Polygons and convex hulls

A polygon is a set of 2D points (called vertices, as in a graph) connected by edges in a given sequence. In this class, we've been especially interested in a particular kind of polygon: *convex* polygons. A convex polygon is a polygon that has the property that for any two points A and B inside or on the polygon, the line segment connecting A and B is entirely inside the polygon. Figure 1 shows some examples of convex (and non-convex) polygons.

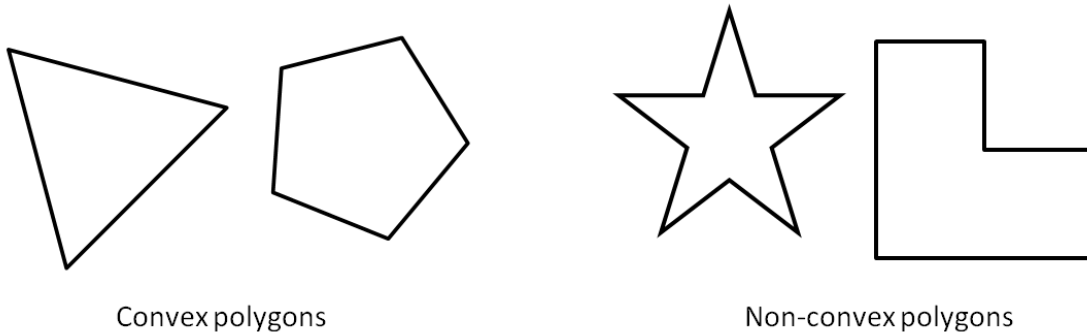


Figure 1: Some examples of convex and non-convex polygons.

This leads to the definition of a special kind of polygon called the *convex hull*. The convex hull of a polygon G is the smallest convex polygon that completely contains G . The convex hull will be made up of vertices of P (but not necessarily all vertices of G will be on the hull). Intuitively, the convex hull is what you would get if you took a stretchy rubber band, stretched it out to contain the polygon, then let it snap to its smallest possible shape. An example of a convex hull of a polygon is shown in Figure 2. Naturally, the convex hull of a convex polygon is the polygon itself.

We can also define the convex hull on a point set (instead of a polygon). The convex hull of a point set P is the smallest convex polygon that contains all of the points. Again, the vertices of the convex hull will all be points in P . An example of a convex hull of a point set is shown in Figure 3.

Note that an edge of the convex hull has the property that *all* points are on one side of the line passing through that edge. The converse of this is true as well: if a pair of points

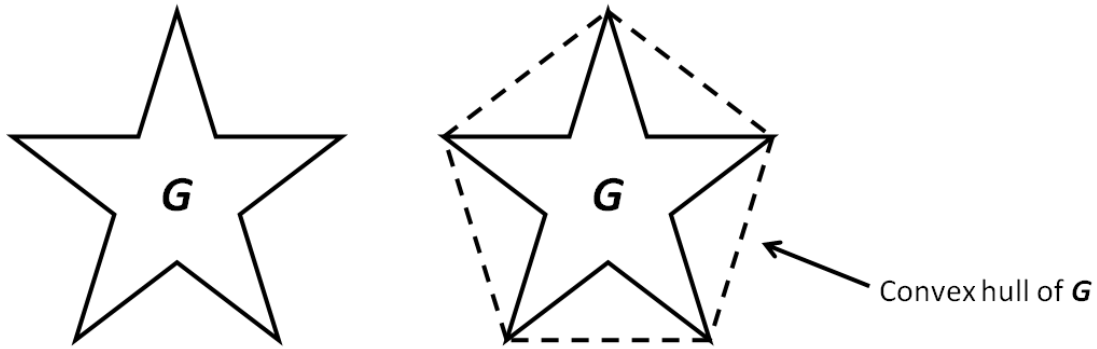


Figure 2: Convex hull of a polygon.

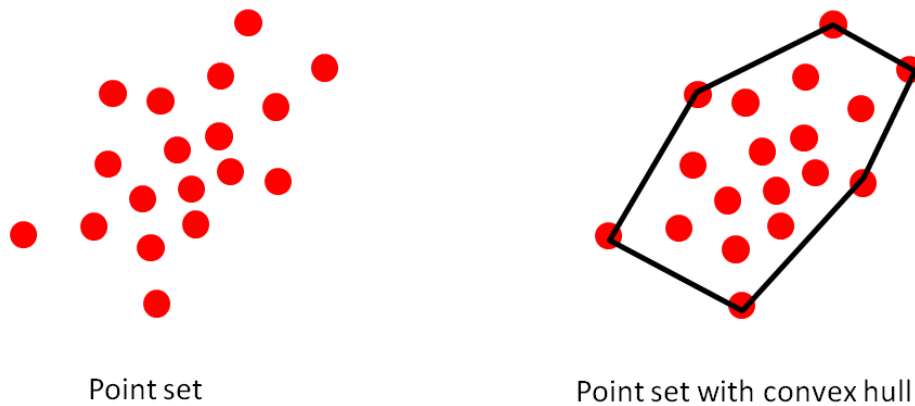


Figure 3: Convex hull of a point set.

A and B in the point set have the property that all other points are on one side of the line passing through AB , then the edge AB will be part of the convex hull.

Next, we consider how to find the convex hull of a set of points. Note that we can easily find a few points that are *definitely* on the convex hull: for instance, the leftmost point, the rightmost points, the uppermost point and the lowermost point. In fact, for any vector representing a direction, the most extreme point in that direction will be part of the convex hull.

1.1 Algorithms for convex hull

In class we looked at several different algorithms for computing the convex hull (and you implemented one of these).

Giftwrapping algorithm. The first was the giftwrapping algorithm. The idea behind the giftwrapping algorithm is to find a point p_1 in P that is *definitely* on the hull, then work our way around the hull by continually selecting the next point the makes the greatest “right-hand turn”:

GIFTWRAPPINGALGORITHM(P):

1. Start with a point $p_1 \in P$ that is guaranteed to lie on the border of the hull (we chose the bottom-most point in class, but you can choose any extreme point). Add p_1 to the hull.
2. Find the point $p \in P$ that makes the largest “right-hand turn.” That is, if p_{prev} is the previous point on the convex hull, and ℓ is the line between the two previous points in the convex hull (or a horizontal line, in case the previous point is the bottom-most point) then the p maximizes the angle θ between ℓ and the line through p_{prev} and ℓ . Add p to the hull.
3. Repeat step 2 until we get back to p_1 .

Giftwrapping is an example of an *output-sensitive* algorithm: the running time depends on the size of the convex hull (by size, we usually mean the *number of vertices* on the convex hull). If n is the size of the point set P (again, the number of points in P) and h is the size of the convex hull, then the running time of the giftwrapping algorithm is $O(nh)$: for each vertex on the convex hull, we spend $O(n)$ time looking for the next vertex. In the worst case, when *all* the points are on the hull, then the running time will be $O(n^2)$. If the points are randomly distributed (inside a circle, for instance), then we expect that about \sqrt{n} points will be on the convex hull, so the running time of giftwrapping will be $O(n\sqrt{n})$ on average for this kind of distribution.

Quickhull. We also learned about the quickhull algorithm. Quickhull is a divide-and-conquer algorithm that, like quicksort, finds the convex hull by dividing the point set P into two parts, then recursively running quickhull on each part. You won’t need to know the details of the quickhull algorithm, but you should know that the running time is $O(n \log n)$. Thus, quickhull is *not* an output-sensitive algorithm; the running time depends only on the size of the input set.

Applications of the convex hull. In class, we used the convex hull to come up with a compact descriptor for the shape of the lightstick (we then found the major axis of this shape to find the direction the lightstick was pointed). However, there are other things the convex hull can do. In particular, in class we saw how the convex hull can be used for sorting. In particular, given a list of n numbers x_1, x_2, \dots, x_n , suppose we create a set of 2D points $(x_1, x_1^2), (x_2, x_2^2), \dots, (x_n, x_n^2)$, i.e., points on a parabola. If we find the convex hull of this point set, then, starting from the bottommost point (assuming all numbers x_i are non-negative) the order of the points along the convex hull will be the same as the sorted order of the original numbers. This type of approach, where we solve one problem (in this case, sorting) by converting it to a different problem (convex hull) and solving the new problem, is called a *reduction*: we are reducing the problem of sorting to that of computing the convex hull. This has implications for the best possible running time of the new algorithm. In this case, we’ve learned that the best possible running time of a sorting algorithm is $O(n \log n)$ —in general, we can’t come up with a faster sorting algorithm (though we might be able to for certain specific types of input). This means that we also can’t find the convex hull in

faster than $O(n \log n)$ time; otherwise, we could come up with a faster sorting algorithm by converting to convex hull.¹

2 Interpolation

Consider the following problem: we're given the small image below, and want to blow it up to poster size. How can we do this?



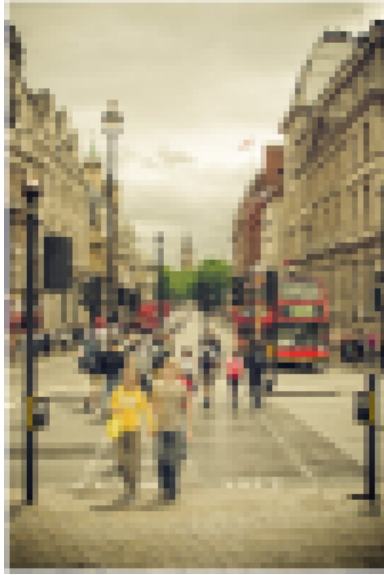
Instead of blowing it up to poster size, let's just try to blow it up so that it is four times as wide and four times as tall. That means we need to replace each pixel with a 4×4 patch of pixels (i.e., each pixel needs to be replaced with 16 pixels). How do we decide what to color the new pixels that we're inserting in between the original pixels? This problem is called *interpolation*—we need to figure out the values of an image in between known values. The general form of interpolation is that we're given some values of function (not necessarily an image) at a few *sample points* in the domain, and we want to compute the value of the function for the *entire* domain. For an image, these sample points are the pixel locations of the original image.

The exact values we fill in depend on our assumptions of how the function behaves. There are many possible assumptions we could make.

Nearest neighbor interpolation. One of the simplest interpolation schemes is called *nearest-neighbor interpolation*. With nearest-neighbor interpolation, we assume that the value of the function is the same as the value of the closest known sample. For our example of blowing up an image by a factor of four, this means we would replace each pixel with a 4×4 block with the same color as the original pixel. The resulting image would look like this:

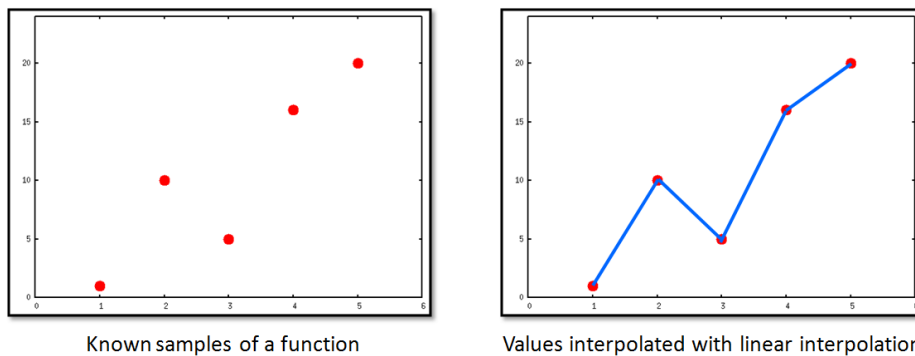
¹Two comments about this reduction.

1. Strictly, it is only the case that the combined running time of converting a sorting problem to a convex hull problem **and** solving the convex hull problem cannot be less than $O(n \log n)$. In theory, if the conversion process itself took $O(n \log n)$ time, then there might still be a faster algorithm for convex hull. Note that the conversion process in this case takes just $O(n)$ time, however. We just need to scan through the list of numbers creating a 2D point for each one.
2. Actually, the best algorithms for convex hull work in $O(n \log h)$ time, which is faster than $O(n \log n)$ if h is small. However, this doesn't help us with the sorting problem, since we want to create a point set where *every* input point is on the convex hull.



The result is pretty blocky, and doesn't look that great. Our assumption that the values of the image should be equal to the closest known value is not a very good assumption, unfortunately. A somewhat better assumption is that the function behaves linearly between the known samples. This assumption leads to *linear interpolation*.

Linear interpolation Linear interpolation (or *lerp*, as it is known in computer graphics, where it is very commonly used), only really makes sense for one dimensional functions. In the figure below, filling in the values of the following function given the samples on the left using linear interpolation results in the function on the left:



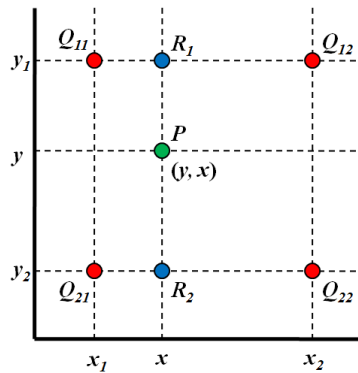
We are given a set of samples of the underlying function $f(x)$ at a number of discrete locations x (e.g., $x = 1, 2, 3, \dots$ in the common case where the samples are at integer locations). To estimate the value of the function at a particular position x using linear interpolation, suppose the closest known samples on the left and right of x are $(x_{\text{left}}, y_{\text{left}})$ and $(x_{\text{right}}, y_{\text{right}})$. The value of at x in our interpolated function \tilde{f} is then computed as:

$$\tilde{f}(x) = \frac{x_{\text{right}} - x}{x_{\text{right}} - x_{\text{left}}} f(x_{\text{left}}) + \frac{x - x_{\text{left}}}{x_{\text{right}} - x_{\text{left}}} f(x_{\text{right}})$$

If x_{left} and x_{right} are at consecutive integer locations (e.g., $x_{\text{left}} = 1$ and $x_{\text{right}} = 2$, then this simplifies to:

$$\tilde{f}(x) = (x_{\text{right}} - x)f(x_{\text{left}}) + (x - x_{\text{left}})f(x_{\text{right}})$$

Bilinear interpolation How do we apply linear interpolation in two dimensions? The answer is *bilinear interpolation* (bilerp); this time, we define bilinear interpolation only for an integer grid of known sample points. Given a new point at which we'd like to compute the value of the function, we find the four nearest sample points, then use linear interpolation first along the y -direction, then along the x -direction (we would also get the same answer by interpolating along the x -direction first). This is demonstrated in the figure below:



To compute the value of the function at the green point P , we find the four surrounding points $Q_{11} = (y_1, x_1)$, $Q_{12} = (y_1, x_2)$, $Q_{21} = (y_2, x_1)$, and $Q_{22} = (y_2, x_2)$, then linearly interpolate along the y -direction to compute the function values at the two blue points R_1 and R_2 . We then linearly interpolate R_1 and R_2 to compute the value at P . Again, if the four surrounding points on an integer grid, respectively, then the formula for the estimated function value at P is:

$$\tilde{f}(P) = (x_2 - x)(y_2 - y)f(Q_{11}) + (x - x_1)(y_2 - y)f(Q_{12}) + (x_2 - x)(y - y_1)f(Q_{21}) + (x - x_1)(y - y_1)f(Q_{22})$$

This could also be written in terms of three linear interpolations:

$$\tilde{f}(P) = \text{lerp}(x_2 - x, x - x_1, \text{lerp}(y_2 - y, y - y_1, f(Q_{11}), f(Q_{21})), \text{lerp}(y_2 - y, y - y_1, f(Q_{12}), f(Q_{22})))$$

Using bilinear interpolation, our 4×4 image looks like this:



Compared to nearest neighbor interpolation, this looks much smoother, and not as blocky. However, we're not really resolving any more detail about the scene. There are better interpolation schemes than bilinear interpolation (such as bicubic interpolation), but no interpolation scheme will really be able to add detail that isn't present in the original image. The basic reason for this is that *many* different larger images, when shrunk down, would give the same small image. Thus, we can't expect to recover a larger image perfectly from a small image. (It's easier to see this when we take this to its logical conclusion—trying to blow up a full-resolution image from a single pixel. This is impossible.)

2.1 Convolution

Convolution is an image operator involving two images: a *target* image and a (usually quite small) image called a *kernel*. To convolve (or filter) a target image I with a kernel, we run the kernel across each pixel of I , multiply it by the pixels of I in a window of the same size of the kernel, then add up the results. For instance, if the target image and kernel are:

$$\text{image} = \begin{array}{|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 \\ \hline 1 & 2 & 3 & 4 & 0 \\ \hline 1 & 1 & 1 & 1 & 0 \\ \hline 0 & 2 & 4 & 5 & 7 \\ \hline 1 & 1 & 3 & 4 & 5 \\ \hline \end{array} \qquad \text{kernel} = \begin{array}{|c|c|c|} \hline 0 & 1 & 0 \\ \hline 1 & 1 & 1 \\ \hline 0 & 1 & 0 \\ \hline \end{array}$$

then the result of convolving `image` with `kernel` is:

$$\text{output} = \begin{array}{|c|c|c|c|c|} \hline 2 & 3 & 5 & 5 & 1 \\ \hline 4 & 8 & 10 & 9 & 4 \\ \hline 3 & 7 & 10 & 11 & 8 \\ \hline 4 & 8 & 15 & 21 & 17 \\ \hline 2 & 7 & 12 & 17 & 16 \\ \hline \end{array}$$

Note that we assumed for this convolution that the image is padded with zeros, for cases where the convolution kernel spills over the boundary of the image. What the kernel above does is replace each value in the image with the sum of the value itself, and the values of its neighbors in the four cardinal directions (NEWS). For instance, the 2 at the (2, 2) location in the image is replaced with $2 + 1 + 3 + 1 + 1 = 8$ —the sum of itself and its four neighbors.

Often we use kernels whose values add up to 1—this results in an image that is no brighter or darker than the original. For instance, to blur an image, we might use the kernel:

$$\text{blur} = \begin{array}{|c|c|c|} \hline \frac{1}{9} & \frac{1}{9} & \frac{1}{9} \\ \hline \frac{1}{9} & \frac{1}{9} & \frac{1}{9} \\ \hline \frac{1}{9} & \frac{1}{9} & \frac{1}{9} \\ \hline \end{array}$$

This kernel averages every 3×3 window of an image.

3 Image transformations

What if we want to apply more complex transformations to an image than just simple scaling. We'll define an image *transformation* as a function that maps the 2D coordinates of an image to a new set of coordinates. We'll focus especially on *linear* transformations, i.e., transformations of the form:

$$f(x, y) = (ax + by, dx + ey)$$

for constant coefficients a , b , d , and e . Linear transformations can be represented by a 2×2 matrix

$$T = \begin{bmatrix} a & b \\ d & e \end{bmatrix}$$

and applying the transformation can be done with matrix multiplication, defined as:

$$\begin{bmatrix} a & b \\ d & e \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} ax + by \\ dx + ey \end{bmatrix}$$

Scaling



A scaling transformation can be represented with a matrix

$$S = \begin{bmatrix} s & 0 \\ 0 & s \end{bmatrix}.$$

To see this, let's see what happens when we multiply this matrix by a point $\begin{bmatrix} x \\ y \end{bmatrix}$:

$$\begin{bmatrix} s & 0 \\ 0 & s \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} sx \\ sy \end{bmatrix}$$

This multiplies points by a factor of s . If $s > 1$ we would be enlarging the image, and if $s < 1$, we would be shrinking the image.

Rotation



A rotation by an angle θ can be represented with a matrix

$$R = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

Translation and homogeneous coordinates To represent a translation (moving the image left/right, up/down), we need more than just a linear transformation with a 2×2 matrix. For this, we need to introduce an idea known as *homogeneous coordinates*. We will play a trick by adding a 1 after the x, y coordinates of a vector (resulting in what we'll call

a 2D homogeneous vector), and switch to using a 3×3 matrix (whose last row is always $[0 \ 0 \ 1]$). The resulting matrix is an *affine* transformation with six parameters:

$$T = \begin{bmatrix} a & b & c \\ d & e & f \\ 0 & 0 & 1 \end{bmatrix}$$

Multiplying an affine transformation by a homogeneous vector gives:

$$\begin{bmatrix} a & b & c \\ d & e & f \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} ax + by + c \\ dx + ey + f \\ 1 \end{bmatrix}$$

We can now represent a translation by a vector (s, t) as:

$$\begin{bmatrix} 1 & 0 & s \\ 0 & 1 & t \\ 0 & 0 & 1 \end{bmatrix}.$$

To see why, let's multiply by a homogeneous vector:

$$\begin{bmatrix} 1 & 0 & s \\ 0 & 1 & t \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x + s \\ y + t \\ 1 \end{bmatrix}.$$

Composition Linear transformation are composable, as are affine transformations. In other words, given two affine transformations, a scale T_{scale} and a rotation T_{rot} , say, we can apply both in sequence as a composed transformation. This composed transformation can be computed by simply multiplying the two transformation matrices, e.g. $T_{\text{scale}} \cdot T_{\text{rot}}$, giving another 3×3 matrix. Note that the order of the multiplication matters, and the transformations will be applied in right-to-left order (e.g., the rotation will be performed first in the product above). Any number of transformations $T_1, T_2, T_3, \dots, T_n$ can be composed in this way.

3.1 Inverse mapping

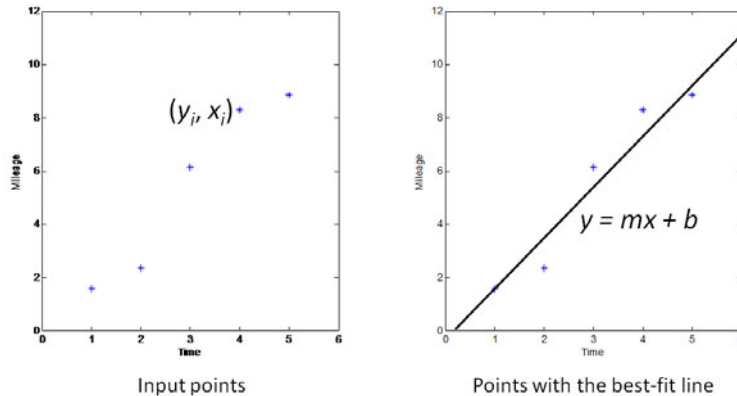
How do we actually apply a transformation to an image? One approach is to take each pixel in the input image, apply the transformation to its x, y position, then set the resulting output pixel to the same color as the input pixel. This approach is called *forward mapping*, and it has some weird results; for instance, for a rotation or upscaling, we might not “hit” every pixel in the output image. Thus, we might be left with gaps in our image when we're finished:



A better approach is called *inverse mapping*. With inverse mapping, we take each pixel of the *output* image, apply the *inverse* transformation to get the corresponding pixel location in the input image, then sample the color at that location to figure out what color to assign to the output pixel. By definition, this guarantees that we will fill in each pixel of the output image with a color (as long as the corresponding pixel of the input is in bounds). However, we might hit a fractional pixel in the input image. Thus, we need to use interpolation to get the value of that pixel. Often, a simple interpolation approach such as bilinear interpolation works pretty well.

4 Optimization and least squares

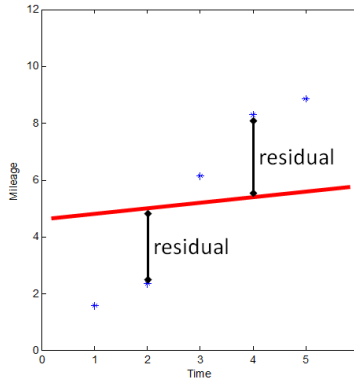
A recurring problem throughout the lectures has been that of *optimization*: we define a cost (or score) function that we want to minimize (or maximize), and then want to find the best solution according to that objective function. A simple example is fitting a line to a set of 2D points (also known as *linear regression*). Suppose we have a set of (noisy) points (x_i, y_i) and want to fit a line $y = mx + b$ to them:



Because of the noise, there is no line that exactly passes through all points. Instead, we seek the *best* possible line. How do we define “best”?

The idea is that we’ll come up with a cost (or objective) function $\text{Cost}(m, b)$ that measures how bad a line defined by m and b is, then find the line that minimizes this function. The

most natural cost function is one that measures how much the line differs from the given points. We call these differences *residuals*; for linear regression, they can be visualized as vertical distances from each input point to the line:



The length of the residual for a point (x_i, y_i) is $|y_i - (mx_i + b_i)|$. Thus, one objective function could be the sum of residuals:

$$Cost(m, b) = \sum_{i=1}^n |y_i - (mx_i + b_i)|$$

For a variety of reasons, however, we prefer the sum of *squared* residuals (we'll explain a few reasons below):

$$Cost(m, b) = \sum_{i=1}^n |y_i - (mx_i + b_i)|^2$$

When using a cost function that is sum of squared residuals, finding the minimum is known as the *least squares* problem. How do we find the line that minimizes $Cost(m, b)$? A few possible approaches come to mind:

1. Check all possible lines (but there are an infinite number of them)
2. Guess a line at random, check it, and repeat with new random answers until we get a good answer (not such a good approach).
3. Guess an answer and improve it until we get a good answer (better).
4. Magically compute the right answer (best).

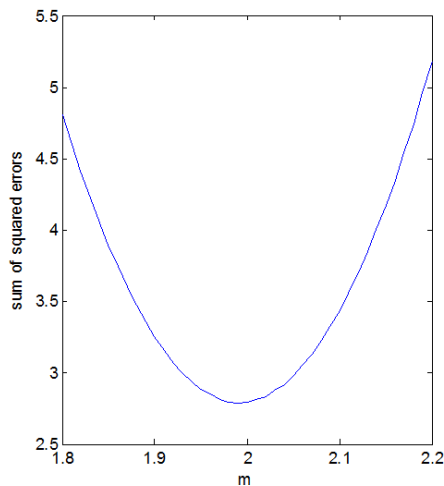
This process of finding the minimum of the function is known as *optimization*. We'll mostly discuss an approach that does a form of 2 above called *gradient descent*.

Gradient descent. Gradient descent can be visualized as finding the minimum by simulating a ball rolling down the objective function until it reaches the bottom. With gradient descent, we go through the following procedure:

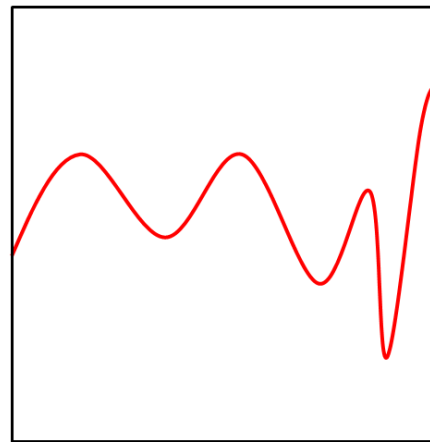
1. Start at a guess for the right answer (which might be a completely random guess).

2. Next, compute the direction of steepest descent—the direction we should move that will take us furthest downhill (this direction is called the *gradient*). With a normal 1D function, this direction will either be “left” or “right”; in 2D it will be a direction in the plane.
3. Take a step in that direction, making sure that we end up below where we started.
4. Repeat 2-3 until we stop going downhill by very much. We’ve reached a minimum.

Does this approach get us to the right answer? For certain functions, yes, we will reach the correct answer eventually. This class of functions are called *convex*. Very similar to a convex polygon, a convex function is one where the area above the function forms a convex set—there are no nooks and crannies. Here’s an example of 1D convex and non-convex function:



Convex



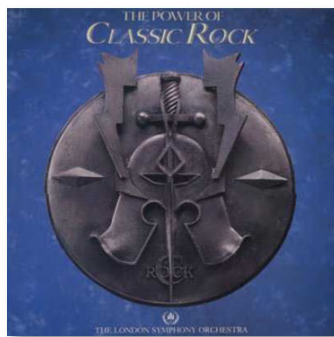
Non-convex

Convex functions are “nice”: they always have a single global minimum, and we will always find this minimum using gradient descent. A quadratic is an example of a convex function, and the sum of convex functions is always convex; thus, the sum of squared residuals (at least, of a linear function) is convex, and so we can successfully solve least squares using gradient descent (this is one reason why least squares is often used).

The other reason why least squares is so common is that it can actually be solved in closed form. We don’t need to use gradient descent at all: there is a magic formula for computing the minimum solution. It’s not important to know the details for this course, but you can think of the problem as defining a linear system in two variables, m and b , and where we have one equation, $y_i = mx_i + b$, for each input point. If we have more than two points, this is an *overconstrained* linear system—it has more equations than variables. Least squares is a technique for solving such overconstrained systems.

5 Feature-based object recognition

Now that we know a bit about image transformations and optimization, we can use these tools to help us do something interesting: object recognition. Suppose we have an image of an object we'd like to find (a template image), and an image that we'd like to find the object in (a search image):



Template image



Search image

There are many possible ways to try and find the search image in the template image. For instance, we could apply convolution to the two images; if we're lucky, the peak of the resulting image will be the location of the search object. However, this might not work in several cases:

1. The object appears at a different scale or orientation in the search image.
2. The object is partially covered by other objects.
3. The lighting in the search image is different than in the template image.

In class, we talked about another approach: using local features.

5.1 Invariant local features

A local feature is a small, distinctive patch of an image—a fingerprint that is likely to be found in other images of the same object. For instance, “corner-like” features often make good candidates for such fingerprints. Finding such features in an image is known as *feature detection*. Feature detection involves finding the 2D positions of likely features, as well as coming up with a *descriptor vector* that describes the appearance of the image around that feature. Ideally, we want our detected features to be *invariant* to certain image transformations, such as rotation, scaling, and changes in brightness or contrast. Here, invariant means that the feature, and its descriptor vector, remain unchanged under these transformations. If our features achieve invariance, then we can robustly compare fingerprints across very different images. Using local features for object recognition has several desirable properties:

1. **Locality:** Features are local, and therefore some of them can be covered up or not detected and we can still find the object.

2. **Distinctive:** Distinct features are detected, so we can potentially distinguish objects in a large database.
3. **Quantity:** We typically find lots of features in an image.
4. **Efficiency:** Compared to approaches which match the *entire* image, the local features approach can be very fast.

5.2 SIFT

SIFT (Scale-Invariant Feature Transform) is one of the best feature detectors around—it produces features that are invariant to scaling, rotation, certain changes in illumination, and small changes in viewpoint (how it does this is left to a computer vision course). It produces a 128-dimensional descriptor for each detected feature.

5.3 Feature matching

Once we detect features in the template and search images, we need to match features between them—that is, find similar-looking features in each image. You’ll learn more about feature matching in Assignment 5, Part 2. For now, assume that we have some technique for finding matches between 2D features in image 1 and image 2, but these matches are noisy—there are many wrong matches that we have to deal with this somehow. The next step involves an algorithm called RANSAC for dealing with such noisy matches.

5.4 Finding image transformations

The next step in our object detection framework is to find an image transformation between the template image and the search image that describes the position, orientation, and scale of the object in the search image. Again, we’ll focus on affine transformations, which will explain many of the transformation we see in practice when doing object recognition.

We now have the following problem: we’re given a set of known matches between the two images (some of which are incorrect), and we want to find an affine transformation that explains these matches as well as possible. Let’s refer to our set of matches as

$$\begin{aligned}x_1, y_1 &\rightarrow x'_1, y'_1 \\x_2, y_2 &\rightarrow x'_2, y'_2 \\&\dots \\x_n, y_n &\rightarrow x'_n, y'_n\end{aligned}$$

i.e., the feature at location (x_1, y_1) in image 1 matches the feature at location (x'_1, y'_1) in image 2. If all of our matches were exact, then we would want to find the transformation T

such that

$$\begin{aligned}
 T \begin{bmatrix} x_1 \\ y_1 \\ 1 \end{bmatrix} &= \begin{bmatrix} x'_1 \\ y'_1 \\ 1 \end{bmatrix} \\
 T \begin{bmatrix} x_2 \\ y_2 \\ 1 \end{bmatrix} &= \begin{bmatrix} x'_2 \\ y'_2 \\ 1 \end{bmatrix} \\
 &\dots \\
 T \begin{bmatrix} x_n \\ y_n \\ 1 \end{bmatrix} &= \begin{bmatrix} x'_n \\ y'_n \\ 1 \end{bmatrix}
 \end{aligned}$$

However, our matches aren't exact. The resulting problem is a slightly more sophisticated form of linear regression. The basic linear regression problem is fitting a line to a set of points (x, y) , as you saw above. Fitting an affine transformation to pairs of 2D points is a very similar problem in a higher dimension, and we can use least squares to solve this problem.

We will consider this problem in the next section, but for now observe that—just as we can fit a line exactly to two (x, y) pairs—we can fit a 2D affine transformation to three $(x, y) \rightarrow (x', y')$ pairs. To see why, note that each match gives us two equations involving the parameters of T :

$$T \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} \tag{1}$$

$$\begin{bmatrix} a & b & c \\ d & e & f \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} \tag{2}$$

which corresponds to the equations

$$\begin{aligned}
 ax + by + c &= x' \\
 dx + ey + f &= y'
 \end{aligned}$$

Thus we have two linear equations in the six unknowns a, b, c, d, e, f . Given three matches, we can find six such equations, and can solve for the unknowns (remember that to solve a linear system, we need as many equations as unknowns). We can solve the linear system using a matrix inversion in Matlab. Given more than three matches, we can solve a least squares problem where we want to minimize the squared distances between the given matches in image 2, and where the transformation thinks those matches should be:

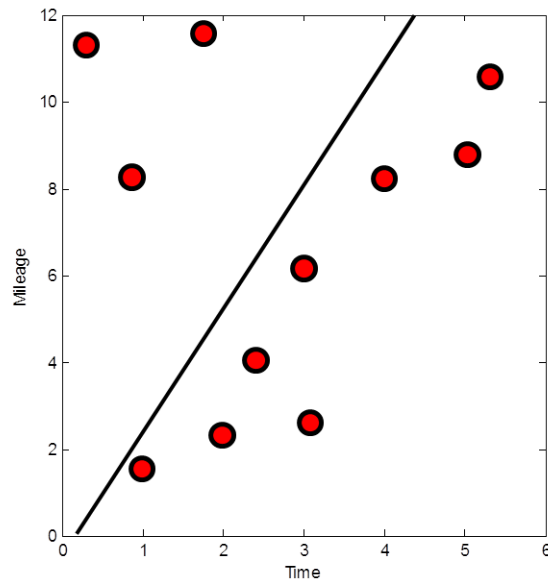
$$\text{dist} \left(T \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix}, \begin{bmatrix} x'_i \\ y'_i \\ 1 \end{bmatrix} \right)^2$$

where dist is the distance between two 2D points.

However, because there are *outliers* (bad matches), running least squares on all the matches at once could produce a very bad solution. This is almost exactly the same problem

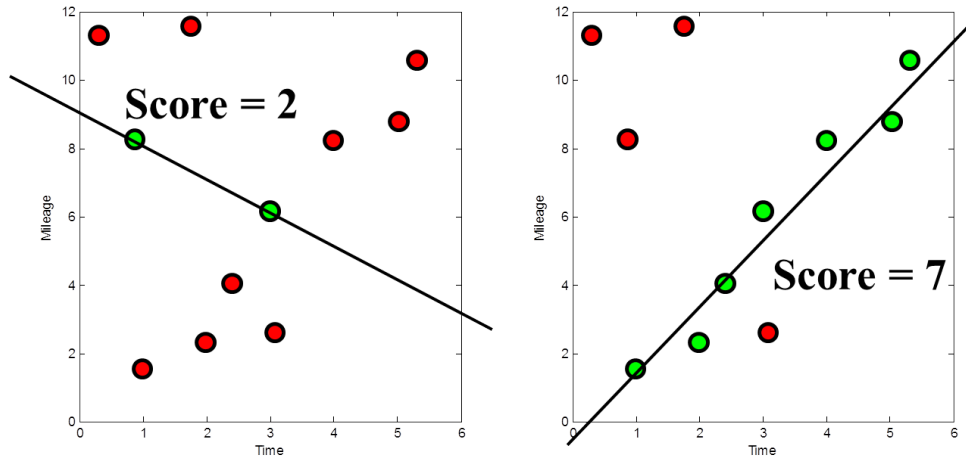
we saw a while ago, when trying to find the “center” of a set of 2D points. The mean (or centroid) of the set of points is sensitive to outliers (the “Arnold Schwarzenegger” problem), and could be arbitrarily corrupted by an arbitrarily bad outlier.² So we need a more robust objective function, and a way to compute it. This is similar to us deciding to use the median instead of the mean as a center for a set of 2D points, and coming up with an algorithm for computing the median.

RANSAC. We need a more robust way of measuring the “goodness” of a potential affine transformation than the least squares metric. A key observation is that we want to find an affine transformation that predicts a large number of matches well, even if it predicts a few (the outliers) poorly. To see this intuition, let’s return to linear regression for a moment. Consider the following set of 2D points, which approximate a line, but with a few outliers.



The least squares solution is shown, and is indeed corrupted by outliers. Instead, our new idea is to measure the goodness of a line by *counting* the number of points that are close to the line. This completely ignores points that are far from the line, and so this measure is naturally robust to outliers. To see the idea, consider the two lines below:

²It should be noted here that the mean of a set of input 2D points happens to be the point that minimizes the sum of squared distances to the input points, i.e., it is the solution to a least squares problem. This is a beautiful fact.



The first line is only close to two points (shown in green), and so is not so great (we call these green points *inliers*). On the other hand, the second line has seven inliers, and so is much better. Given this counting metric, how do we find the best line? Unfortunately, there is no closed form solution to this problem, unlike with least squares. An alternative would be to generate random lines and score each one, but random lines might not be very good. Instead, we'll use a technique called RANSAC (Random Sample Consensus), which generates lines to test much more intelligently. Remember that we need two points to fit a line (and three matches to fit an affine transformation). What RANSAC does is choose such minimal samples of points/matches at random, generates the model (e.g., line or affine transformation), and tests that model by seeing how many input points are inliers. We then return the line with the largest number of inliers. For fitting a line, RANSAC operates as follows:

RANSAC(input: set of 2D points):

1. Randomly select two points.
2. Solve for the corresponding line L .
3. Compute the number of inliers to L , i.e., other points that are close enough to the line.
4. If L has the most inliers so far, save it.
5. Repeat steps 1-4 N times, then return the best L .

In a way, this is like voting—each 2D point votes for the candidate lines that agree with it, and the winner is the line with the most votes. We can use the same algorithm for fitting an affine transformation, except that we need three matches to generate a hypothesis transformation, and we measure how well each other match agrees with a transformation by through a slightly different measure of agreement. That is, by “agree,” we mean that the transformation maps a point (x_i, y_i) in the template image “close” to its matching point (x'_i, y'_i) in the target image (we don't expect it to map the point *exactly* to the target). To define this concretely, we'll

define a threshold t (say 5 pixels), then say a match $(x_i, y_i) \rightarrow (x'_i, y'_i)$ agrees with T if

$$\text{dist} \left(T \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix}, \begin{bmatrix} x'_i \\ y'_i \\ 1 \end{bmatrix} \right) < t$$

(where dist is the distance between two 2D points).

Again, to tell how well a transformation T agrees with the set of matches, we can simply count the number of matches that agree (these are called *inliers* to the transformation—the matches that don't agree are called *outliers*). We then sample many random sets of three matches and find the resulting T that agrees with the' largest number of matches. Here's what RANSAC looks like for fitting an affine transformation to a set of 2D matches:

RANSAC(input: set of matches):

1. Randomly select three matches.
2. Solve for the corresponding transformation T .
3. Compute the number of inliers to T .
4. If T has the most inliers so far, save it.
5. Repeat steps 1-4 N times, then return the best T .

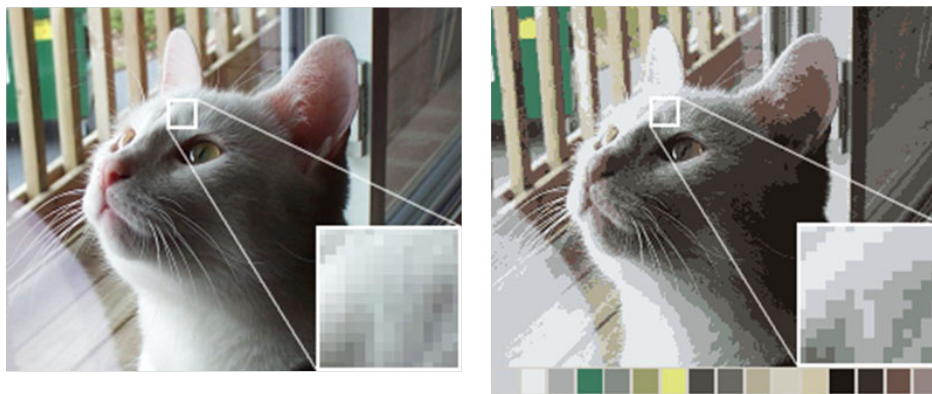
The RANSAC algorithm works well, and is used in many different computer vision applications. RANSAC relies on two parameters: t , the “closeness” threshold, and N , the number of rounds (though if we expect to have a certain percentage of inliers, and want to compute the right T with a given probability, e.g. 99.9999%, then we can compute the N we need). One thing to note: the more outliers we have, the more rounds we should run the algorithm. This is because each time we generate a hypothesis, we need to select three points, and we only really get a good model if all three are inliers. If we choose even one outlier, our model will very likely be bad. And this is pretty likely to happen: if $\frac{1}{3}$ of the matches are wrong, then the chances of picking one wrong match among three random matches is almost $\frac{2}{3}$. This is also why we like to choose the minimal possible number of matches to generate a transformation: choosing more than the minimum (e.g., choosing 4 or 5 points) dramatically increases the likelihood that we'll choose an outlier.

6 Clustering and image segmentation

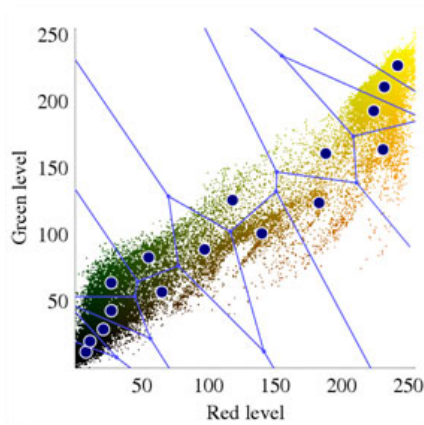
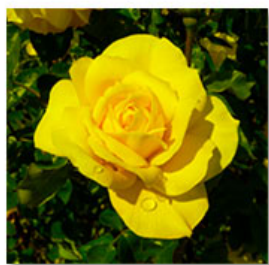
At the beginning of the course, we looked at the problem of finding a red lightstick in an image, for which the first step was to *threshold* the image into red and non-red pixels. You can think of this as *segmenting* out the lightstick from the background. A more general version of this problem is to segment out *all* of the objects in an arbitrary image, for instance this one:



This is a central (and unsolved) problem in computer vision. A state-of-the-art automatic algorithm for image segmentation produces the result on the right, where pixel groupings are shown by color. This is pretty good—the tree, beach, water, and sky are each mostly one segment—although not perfect. We looked at a very simple version of this problem called **color quantization**. Suppose we want to take an image with a full color spectrum (recall that an RGB image with 256 possible values per color channel can contain up to $256 \times 256 \times 256 = 16,777,216$ distinct colors), and color it with a “palette” of just k colors (where k is a much smaller number, such as 4, 16, or 256). We might want to quantize the image into a smaller number of colors in order to segment the image, or to compress it, for instance. Here’s an example of quantizing a full color image to just eight colors:



Given a target number of k colors in our palette, how can we choose the *best* k colors with which to represent the image—that is, the set of colors that will make the image look most like the original image? We would somehow like to choose a set of colors that are similar to the most colors in the input image. This is an example of a *clustering* problem: given all of the colors present in the input image, group them into k clusters of similar colors, and choose a representative for each group. To illustrate this, here’s a diagram that shows how the colors in an input image look in color space:

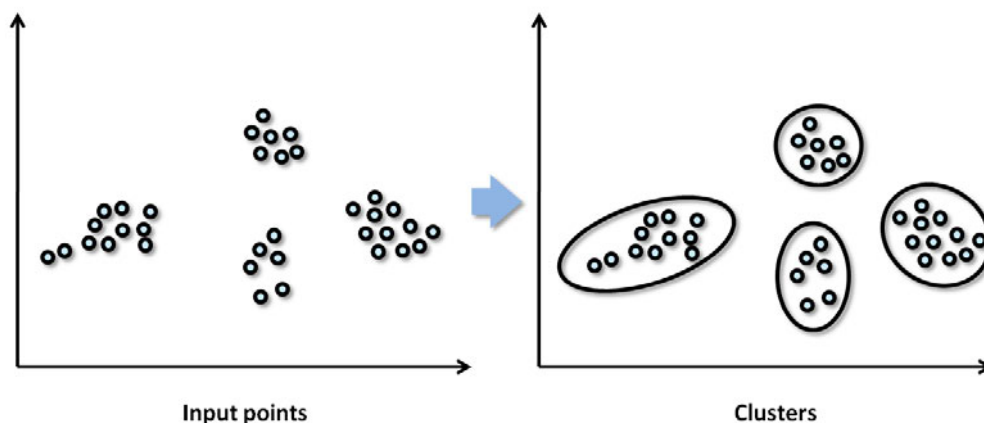


(Image courtesy Wikipedia). The left image is the original image, and the right is a visualization of the RGB color space (since three dimensions are difficult to visualize on the page, this only plots the R and G components of the colors). For each pixel in the original image, a plot is pointed in the right image whose coordinates are the R and G channels of that pixel's color (R on the x -axis, and G on the y -axis). While these pixels don't naturally group into a set of well-separated clusters, we can still divide them up into groups—this is shown in the image on the right by the blue lines partitioning up the color space, and by the blue points representing the “centers” of each cluster. Given this clustering, we could assign each pixel in the input image the color of the closest cluster center in RGB space, and the result would be our quantized image.

Clustering problems are widespread in many different domains. Businesses often want to cluster our customers into different groups in order to better tailor their marketing to specific types of people (market segmentation). News aggregation websites like Google Maps cluster news stories from different sources into groups of stories that are all on the same basic topic. And clustering can be use to find epicenters of disease outbreak (such as the recent swine flu epidemic).

Given a set of points in some dimension (three, in the case of colors), how can we compute a good clustering? We first have to define what we mean by *good*. This definition should capture the fact that we want all of the points in a cluster to be close to each other.

In our case, we will state the clustering problem as follows: we are given a set of n data points x_1, x_2, \dots, x_n (represented by vectors of some dimension), and a way to compute distances between these points (for instance, the typical Euclidean distance between vectors). An example clustering of 2D points (with four clusters) is shown below:



How do we define the clustering problem? We want to create a set of clusters and assign each point to a single cluster so as to minimize some underlying cost function that rates the “goodness” of a clustering. What should this cost function look like? Intuitively, we want a cost function that prefers clusterings where similar points are close together, and where dissimilar points are far away. One of the simplest approaches is known as k -means.

6.1 k -means

The k -means clustering problem is as follows: we are given the number of clusters k we should compute (e.g., $k = 4$). We need to find (a) a set of cluster centers $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_k$, and (b) an assignment of each point x_i to a cluster C_j , so as to minimize the sum of squared distances between points and their assigned cluster center. That is, we want to find the cluster centers and cluster assignments that minimize the following objective function:

$$\text{cost}(\bar{x}_1, \dots, \bar{x}_k, C_1, \dots, C_k) = \sum_{j=1}^k \sum_{x_i \in C_j} |x_i - \bar{x}_j|^2$$

At first, this problem looks like the kind of least squares problems we’ve seen in the past (such as line fitting), and which we know are easy to solve. Unfortunately, however, there’s a critical difference: we need to find the cluster centers *and* the cluster assignments at the same time, and this makes the problem much more difficult. If we knew either of the two then the problem would be easy: given the cluster centers, the optimal assignment is to assign each data point to the closest cluster center; given the cluster assignment, the optimal centers are the means of each cluster. But we need to compute both.

This small difference means that, far from being easy to solve, the k -means problem is in a class of problems that are *extremely* difficult. No known algorithm exists for computing the optimal k -means in less than exponential time in n (essentially, we have to try every possible clustering to find the one, and there are $O(k^n)$ possible clusterings). Compared to polynomial time algorithms for other problems we have seen—sorting ($O(n^2)$), median finding ($O(n)$), convex hull ($O(n^2)$), exponential time algorithms are *much* worse (it is similar to comparing snail sort to quicksort). The class of hard problems of which k -means is a member is known

as *NP*-hard problems (*NP* stands for non-deterministic polynomial time, though you won't learn why until much later).

Should we just give up? No, we can still hope to create a fast algorithm that's pretty good (for instance, one that always gets an answer that is pretty close to the optimal one). In class we learned about one approach called *Lloyd's algorithm*, which iterates between computing cluster centers and cluster assignments.