# Class Notes CMSC 426
# 3D Geometry and Projection

## Introduction

One of the main goals of computer vision is to use 2D images to determine the structure and position of 3D objects in the world. To do this, we have to discuss geometry. We need to understand how to represent geometric objects in 3D and 2D, and we need to understand the relationship between the 3D world and a 2D image. This note will cover three topics: representation of planes, points and lines; perspective projection that relates the 3D and 2D positions of these objects; and intersections of these objects as well as the process of finding lines or planes that include multiple lines and points.

## Representation

In this class we will only consider the simplest geometric objects: points, lines and planes. This is the minimal set of objects that we can consider. Points are the simplest of objects, we must understand lines because light goes in a straight line, and we must understand planes, because the image is a plane. However, these simple objects are also very important, and much work in vision is done using just this set of objects.

**Points:** First, we recall that we represent a point by its coordinates in space. If a point is in 2D, we describe it with *x* and *y* coordinates. We will describe points using lower case letters, so we might write *p=(x,y)*. If a point is in 3D we also need a *z* coordinate, and we use upper case letters, so we could write *P=(x,y,z)*.
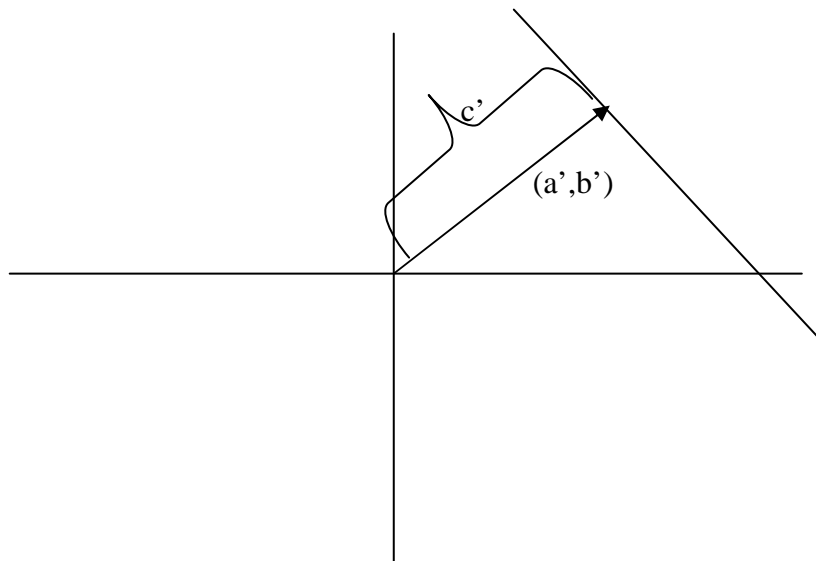
**2D Lines:** Lines already become a little more complicated. There are several ways of representing lines, each with their own advantages and disadvantages.

In 2D we can represent a line with a single, linear equation, of the form: *ax+by+c=0*. Here, *a, b,* and *c,* are constant values that determine the line. A point $(x_0,y_0)$ is on this line if the equation is satisfied when we plug in these values for *x* and *y.* We can get a useful intuition by rewriting this equation in the following way:

$$(x,y) \bullet (a,b) = -c. \quad \Rightarrow \quad (x,y) \bullet \frac{(a,b)}{\|(a,b)\|} = \frac{-c}{\|(a,b)\|} \quad \Rightarrow \quad (x,y) \bullet (a',b') = c'$$

$$where \quad a' = \frac{a}{\|(a,b)\|} \quad b' = \frac{b}{\|(a,b)\|} \quad c' = \frac{-c}{\|(a,b)\|}$$

That is, we can create a unit vector, *(a',b').* Then a point is on a line if the inner product between the point and *(a',b')* is equal to some constant value, *c'.* Keeping in mind that the inner product between *(x,y)* and *(a',b')* measures the distance from *(x,y)* to the origin in the direction of *(a',b')*, we have the following picture.

When we describe a line like this, the vector *(a',b')* will be orthogonal to the line.

Also, it is also useful to recall another way of describing a line, with the equation: $y = mx+b$. When we describe a line like this, *m*, is the slope of the line.

There is still one more way of describing a line. We can write down a recipe for reaching any point on a line by starting at one point of the line, say *(x₀,y₀),* and then moving some distance *t*, in the direction of the line. If we describe the direction of the line with a unit vector, *(u,v),* then we can write this down by saying that *(x,y)* is on the line if it satisfies the equation:

$$(x, y) = (x_0, y_0) + t(u, v)$$

Notice that this is really two equations, one for the *x* component of the point, and one for the *y* component, and these equations are linear in three unknowns, *x, y,* and *t.* Previously, we had described a line with one linear equation and two unknowns. We can convert between these two representations by using one equation to eliminate *t;* then we would obtain one linear equation in two unknowns.

**Planes:** Before, we consider how to represent a line in 3D, it's useful to look at a 2D plane in 3D. First, let's take an example. Suppose we want to represent the floor of a room. We can say that these are all the points that have a height of zero. If we use the *y* coordinate to represent height, we can represent this plane with the equation: *y=0*. Note that this is a linear equation in *x, y,* and *z,* although *x* and *z* don't happen to actually show up in the equation, since they points on the floor can have any values for *x* and *z*. More generally, we can describe any plane with a linear equation of the form: *Ax+By+Cz+D=0.*

Just as in 2D, we can rewrite this equation by coming up with a unit vector, $(A',B',C') = (A,B,C)/||(A,B,C)$. Then, we can say a point $(x,y,z)$ is on the plane if it satisfies the equation:

$$(x, y, z) \bullet (A', B', C') = D'$$

where $D'=-D/||(A,B,C)||$. That is, a plane is a set of points whose inner products with a specific unit vector are all the same. Or, to put it another way, we get to a point on a plane by going a distance $D'$ in the direction $(A',B',C')$ and then going in any direction orthogonal to $(A',B',C')$ by whatever amount we want. In this case, $(A',B',C')$ is the vector normal to the plane.

Let's look at our equation for the floor, $y=0$, from this point of view. This equation can also be written: $(x,y,z).(0,1,0)=0$. This says that the floor is the set of all points whose distance from the origin, in the $y$ direction, is $0$.

**3D Lines:** Now we will consider how to represent a line in 3D. One way to do this is to note when two planes intersect, they intersect in a line. If we want to represent a point as lying in the intersection of two planes, we can say that the point must satisfy two equations, one for each plane, so that it lies inside both planes. So we can represent a line by saying that a point $(x,y,z)$ is on a line if it satisfies the equations:

$$A_1 x + B_1 y + C_1 z + D_1 = 0 \qquad A_2 x + B_2 y + C_2 z + D_2 = 0$$

A second useful way to represent a line, as we did in 2D is to give a starting point, and a direction, indicating that we can reach any point on the line by going some distance in that direction. That is, we can write:

$$(x, y, z) = (x_0, y_0, z_0) + t(u, v, w)$$

Here, $(x_0,y_0,z_0)$ is any point on the line, and $(u,v,w)$ is a unit vector indicating the direction in which the line extends. As $t$ varies, the right hand side of the above equation can represent the location of any point on the line. We can note that the above equation is really three equations, one for each component of the point, with four unknowns, $x, y, z,$ and $t$. In contrast the first representation expresses a line as two equations with three unknowns.

**Intersections and Linear Combinations**

We now launch into a discussion of how to find the intersection of lines and planes, and how to find the linear space (ie., a line or a plane) that include several geometric objects (eg., finding a line that includes two points, or a plane that includes three). These operations are central to many vision tasks, as we will see. For example, light travels in a straight line. An image is a plane. If we want to know where a ray of light will appear in an image, we must know how to find the intersection of a line and a plane.

**Intersecting lines and planes:** We have shown how to represent lines and planes with a set of linear equations. When we intersect these objects, this means that for any point in the intersection, all these equations should hold. Therefore, we can represent this intersection simply by listing the equations that hold. For example, suppose we want to intersect a line and a plane. Suppose further that the line is represented by the two equations: $A_1x + B_1y + C_1z + D_1 = 0$ and $A_2x + B_2y + C_2z + D_2 = 0$, while the plane is represented by the equation $A_3x + B_3y + C_3z + D_3 = 0$. The intersection of a line and a plane is the set of points that satisfies all three of these equations. One might have the intuition that a line and a plane intersect in a single point, so only one point will satisfy these equations. Or one might note that three linear equations with three unknowns will generally have only one solution, which is the single point that satisfies these equations. It is also possible for these equations to have no solution, which is what happens if the line is parallel to the plane, and never intersects it.

Sometimes it is more convenient to do this intersection using the equation for a line: $(x,y,z) = (x_0,y_0,z_0) + t(u,v,w)$. Suppose the equation for the plane that this line intersects is: $Ax+By+Cz+D=0$. Then we have four equations with four unknowns to solve. But these have a nice form, since three of the equations give *x, y,* and *z* on one side, and linear expressions containing no unknowns but *t* on the other. So we can substitute, and get the equation: $A(x_0+ tu) + B(y_0+tv) + C(z_0+tw) + D = 0$. We can then solve this equation for *t*, and substitute this value into our equation for the line to get *(x,y,z).*

As an example, suppose we have a line with the equation *(x,y,z) = (1,2,3) + t(1,1,1),* and we want to intersect it with a plane that has the equation *2x+y+2z+1=0.* Substituting, we get: *2(t+1)+(t+2)+2(t+3)+1=0.* This gives us: *5t+11=0, t = -11/5.* So we then have: *(x,y,z) = (1,2,3)-(11/5,11/5,11/5).*

**Finding a line containing two points:** We may also wish to find the equation for a line that contains two points. In 2D, we can do this by taking the equation for a line, *y=mx+b,* and treating *m* and *b* as the unknowns. Then, for each point we can substitute in the values of *x* and *y*, giving us two equations with the unknowns *m* and *b*. Notice that this will work except for the case of a vertical line, which cannot be described by *y = mx + b,* since it has infinite slope. We would need to check for this case separately.

Here's another way to get an equation for a line from two points. Suppose we have points *p* and *q*. We can write *p + t(p-q).* Here *p* serves as an example of a point on the line, while *(p-q)* is a vector in the direction of the line. Note that this works in two or three dimensions.

**Finding a plane containing three points, or a point and a line**

Just as two points determine a line, three points determine a plane. There are several ways of finding the plane from three points. For example, similar to what we did with a line, we can write the equation for a plane as *Z = AX + BY + D*. Then we can use the *(X,Y,Z)* values for each point to get a linear equation in *A, B,* and *D*. Notice that this approach also doesn't work for some cases, which we must handle separately.

If we want to form a plane from a point and a line, one way to do this is to just pick two points from the line, and then use the above method.
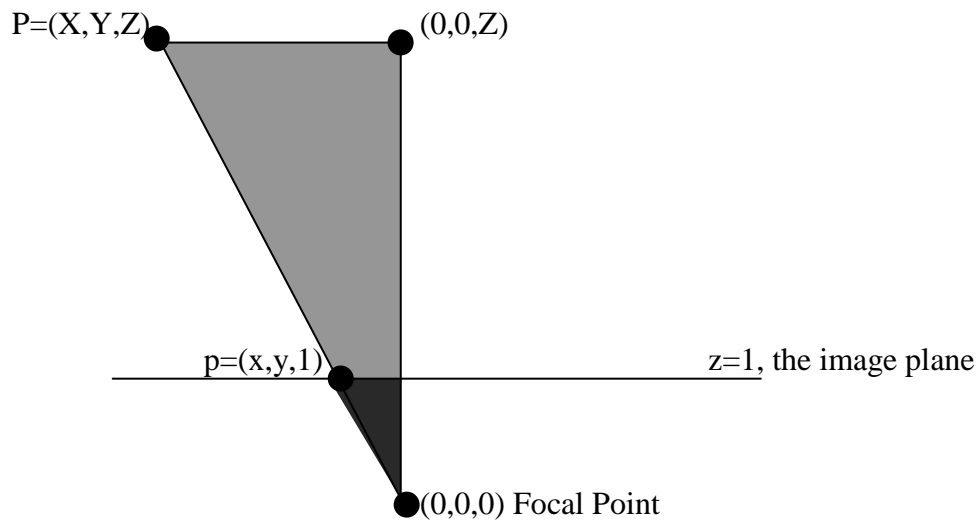
For those of you who are familiar with the cross-product, we note that given points $P_1$, $P_2$, and $P_3$, we can find a vector normal to the plane of the three points by taking $(P_2 - P_1) x (P_3 - P_1)$, where $x$ is the cross-product operation. If we express a plane with the equation: $AX+BY+CZ+D = 0$ this gives us $(A,B,C)$. We can use the coordinates of any of the points to solve for $D$. We won't go into this method in detail, though, because we will try to stick with problems in which finding the plane formed by points, or a point and a line, is easy.

## Perspective Projection

We now have the tools that we need to begin to solve some vision problems. We begin by describing the process of perspective projection. The key question that we must address is, given a description of the camera position and the location of a 3D point, where will this point appear in the image?

With perspective projection, we describe a camera using a *focal* point and an *image* plane. We imagine that light travels in a straight line from a scene point towards the focal point. The location where the light ray intersects the image plane is the image location for this scene point. In a pinhole camera, the focal point is the pinhole, and the light passes through it on the way to the image plane, which might be a CCD, or film. In our idealization of a pinhole camera, the image plane is in front of the pinhole, so the light strikes the image plane before it reaches the focal point. Either way, we can find the image point by forming a line that includes the scene point and the focal point, and finding where it intersects the image plane. We have explained above how to perform these operations, so this tells us how to find the image point corresponding to a scene point, for a general camera position.

As an example, suppose we have a camera with a focal point at *(1,2,3),* with an image plane at the *x=2* plane, and we wish to find the image produced by a point in the scene at the location *(9,6,5).* We can describe the line that joins the scene point to the focal point with the equation: *(x,y,z) = (1,2,3) + t(8,4,2). x = 1 + 8t,* and we want find the point where this intersects the *x=2* plane, (ie., the point on the line that has *x = 2*). This occurs when *2 = 1 + 8t,* or when *t = 1/8.* The point on the line for *t = 1/8* is *(2, 2.5, 3.25).*

P=(X,Y,Z)          (0,0,Z)

p=(x,y,1)                    z=1, the image plane

(0,0,0) Focal Point

We will often consider a special camera set-up (see figure above) which makes it much easier to compute the image points produced by scene points. This is the case in which the focal point is at the origin, and the image plane is the $z=f$ plane, where $f$ is called the focal length of the camera. Suppose now we have a scene point $P$, with coordinates $(X,Y,Z)$, which produces an image point, $p$. $p$ has coordinates $(x,y)$, or alternately we can think of it as a 3D point on the image plane, with coordinates $(x,y,f)$. We notice that there are two similar triangles, one with corners at $(0,0,0)$, $(X,Y,Z)$, and $(0,0,Z)$ (this is the triangle above that is light gray, including also the dark gray triangle and its tip) and the other with corners at $(0,0,0)$, $(x,y,f)$, and $(0,0,f)$ (the dark gray triangle). The first of these triangles is the same as the second one, but scaled by a factor of $Z/f$. This means that the side of the big triangle with corners $(0,0,0)$ and $(X,Y,Z)$ is the same as the side with corners at $(0,0,0)$ and $(x,y,f)$, but scaled by a factor of $Z/f$. This tells us that $(Z/f)(x,y) = (X,Y)$, which implies that

$(x,y)=f(X/Z, Y/Z)$.

This is the basic equation of perspective projection.

**The vanishing point and the horizon**

We can now derive some basic facts about perspective projection. The first is that any point in front of the camera will project to a point in the image plane. The scene point and the focal point form a line that will intersect the image plane in a single point. When we say that a point is in front of the camera, we mean that the image plane separates the scene point and the focal point, so that the line segment connecting the two points will intersect the image plane.

If a line is in front of the camera, it will generally project to a line in the image plane. Every scene point on the line will project into the image along a line connecting it to the focal point. Collectively, all these projection lines connect the scene line to the focal point, forming the plane that includes the scene line and the focal point. In general, this plane will intersect the image plane in a line.

There are two special cases to consider, though. We'll discuss one below. Here we mention the case in which the scene line includes the focal point. In this case, all lines of projection that connect a point on the scene line to the focal point are, in fact, identical to the scene line. Since all of these lines are identical, they all intersect the image plane in the same point. One can visualize this by imagining one is looking at a line end-on, so that it looks like just a point.

In general, a scene plane can fill the entire image. Given any image point, this point forms a line with the focal point, which will intersect the scene plane. This point on the scene plane, then, will project to the image at this image point.

An interesting special case occurs, though, if the scene plane is orthogonal to the image plane. An important example of this is when the scene plane is the ground, and the camera is pointing in a horizontal direction. We can describe such a situation with a camera that has a focal point at $(0,0,0)$ and an image plane of $z=1$, and with a ground plane described by $y=-k$. In this case, the $y$ direction is down, and $k$ is the height of the camera's focal point above the ground. Now, let's consider the projection of a point on the ground. A point on the ground has coordinates $(x, -k, z)$, for any arbitrary values of $x$ and $z$. If the point is in front of the camera, then $z > 1$. Using the equation of projection, the image of this point will be $(x/z, -k/z)$. This could be any location in the image, except that $-k/z$ is always negative, so this point must always be in the bottom half of the image. The image of the plane occupies all points with negative $y$ coordinates, up to the line $y=0$, which is called the horizon. Of course this accords with our experience that when we look at the world in a direction parallel to the ground, the ground is always in the bottom half of the image. It does not fill up the whole image. More generally, similar reasoning shows that when we look at any plane that is orthogonal to the image plane, the points in that plane will fill up half of the image.

There is one last special case, that occurs when a scene plane is situated so that the focal point lies in the same plane. In this case, every line that connects a point in the scene plane to the focal point lies in this plane. This plane intersects the image plane in a line, so the images of all points in the scene plane lie along a line in the image. This is what happens, for example, when you look at a sheet of paper end-on, so that it looks like a thin line.

Now let's consider what happens when a scene line is orthogonal to the image plane, for example, a line that lies on the ground plane. We can describe a line on the ground plane with the equation:

$$(x, y, z) = (x_0, y_0, z_0) + t(u, 0, w)$$

The zero ensures that this point will always stay in the $y=y_0$ plane. If we're talking about the ground plane, we would generally expect $y_0$ to be less than 0. We can use the equations of projection to find the image of a point on this line, which will be:

$$\left( \frac{x_0 + tu}{z_0 + tw}, \frac{y_0}{z_0 + tw} \right)$$

Now, let's look at what happens to with images of points on the line when they get very far from the camera. If we assume that $w$ is positive, then as $t$ gets very big, the $z$ coordinate of a point on the line, $z_0 + tw$, will also get very big, meaning the point is very far from the camera. First let's look at the $y$ coordinate of the image of such a point. It is equal to $(y_0/(z_0+tw))$. As $t$ gets very big, the denominator gets very big, while the numerator stays the same. In the limit, as $t$ goes to infinity, then, the $y$ coordinate goes to zero. For the $x$ coordinate, as $t$ goes to infinity, $x_0$ and $z_0$ become insignificant relative to $tu$ and $tw$. Therefore, the $x$ coordinate of the image point goes to $u/w$. This means that the line appears to approach the point on the horizon $(u/w,0)$ as it vanishes in the distance. This point is called the *vanishing point* of the line.

It is interesting to note that if two lines are parallel, they have the same vanishing point. A line will be parallel to the one we describe above if it has the equation:

$$(x, y, z) = (x_1, y_0, z_1) + t(u, 0, w)$$

We describe this line with a starting point that is different from the first line (though still in the $y=y_0$ plane. But, if the lines are parallel, they must go in the same direction, $(u,0,w)$. By the same reasoning as above, the vanishing point of this line will also be $(u/w,0)$.

Another line, which is not parallel to these, will go in a different direction, and its vanishing point will have a different $x$ coordinate. But notice that all lines in the plane have vanishing points on the horizon, that is, with $y=0$. Again, this is in line with our everyday experience. When we look at a long line, like a railroad track, that seems to vanish into the distance, the line seems to rise up to the horizon as it vanishes.

**Locating an image point in a scene**

We are particularly interested in using our understanding of perspective to perform the inverse operation, to locate a point in the scene using our knowledge of its location in one or more images.

**From one image:** When we see a point, *p,* in only one image, we cannot determine its exact location. There is a whole line in the world that could have produced the image point. This is the line that includes the image point and the focal point. Call this line *L.* If we take any scene point, *P,* on *L*, since *L* includes both *P* and the focal point, this is the

line of light traveling from the *P* to the focal point. Since *L* intersects the image plane at *p,* this will be the image of *P.* Therefore, *L,* describes exactly the set of points that might have produced *p.*

**From one image when the point is on a known plane:** If we have some prior knowledge of the scene, it is possible that we can determine the 3D location of a point from a single image. In particular, if we now that we're looking at points on a known plane, such as the ground plane, then a single image specifies a line that a particular point lies on, and this line will intersect this plane at a point.

As an example, let's suppose that we have a camera with a focal point at *(0,0,0),* and an image plane at *z=1.* We are looking at a point that we know is on the ground plane, *y=-10.* The point that we are looking at appears in the image at (3, -5). We can write an equation that gives us a line that this point must lie in, as: *(0,0,0)+t(3,-5,1).* To intersect this with the *y=-10* plane, we must find *0-5t = -10.* So, *t = 2* and we find that the point is at *(0,0,0)+2(3,-5,1) = (6, -10, 2).*

### Epipolar Geometry

We will now consider what happens when we have two images of a scene. We'll consider two situations. First, we'll discuss some general facts that are true for any two camera positions from which the two images might be taken. Second, we'll consider the special case of two cameras side by side, in an arrangement similar to the two eyes of a person.

**The Epipolar Plane and Line:** We will consider the geometry of a scene point as it appears in two images. First, let's give some definitions. We'll suppose we have one scene point, *P,* and two images, *I1* and *I2.* Let's call the two focal points of these images *f1* and *f2.* These three points, *P, f1* and *f2* form a plane, which we'll call *Q.*

First, let's suppose that we see a point in image one, called *p1.* We know that *p1* will lie on a line that connects *P* and *f1.* This means that *p1* must be in the plane *Q.* This is because if two points are in a plane (*P* and *f1*) then the line that joins them must also be in this plane. For exactly the same reason, if we see *P* in a second image, at the point *p2,* then *p2* must also be in the plane *Q.* The plane, *Q,* is called the *epipolar plane*, and it depends only on the camera geometry and the position of a single scene point, *P.* However, it is important to notice that we can figure out what the epipolar plane is even if we do not know the location of *P,* provided that we know the camera geometry and we have seen an image of *P* in one image. Three points determine a plane, so we can find *P* using *f1, f2,* and *p1.*

Now we get some interesting information if we notice that the plane *Q* will intersect each of the image planes in a line; after all, two planes generally intersect in a line. We'll call the line where *Q* intersects the first image plane *L1,* and the line where it intersects the second image plane *L2.* *L1* and *L2* are called the *epipolar lines.* A final important fact is that *p1* must lie on the line *L1,* and *p2* must lie on *L2.* To see this, notice that *L1* is just
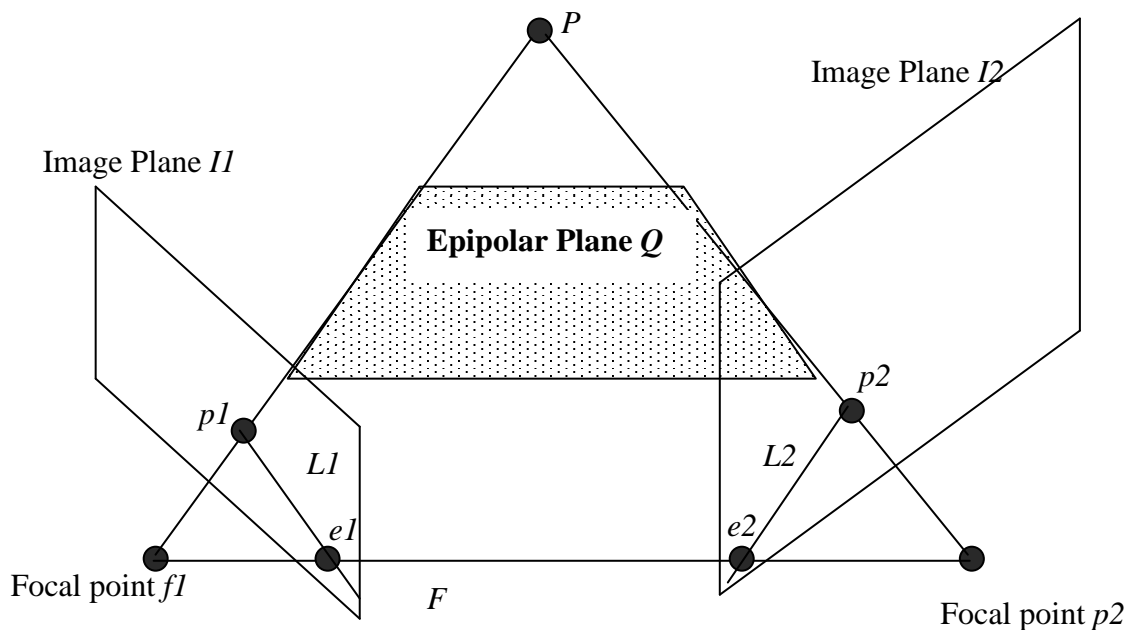
the intersection of *Q* and the first image plane. Since *p1* is on *Q* and it's on the first image plane, *p1* is on the intersection of these, *L1*. Similarly, *p2* must lie on *L2*.

This already gives us some very useful information about the relationship between image points of the same scene point. First, suppose we know the camera geometry that created two images. This situation is called *stereo*. Now, suppose we see a point, *P,* in the first image, giving us *p1*, but we do not know the location of *P*. We can use *f1, f2,* and *p1* to determine *Q*. We can then use *Q* and the location of the second image plane to determine the line *L2*. Now we know that the image of *P* in the second image, *p2,* must lie on this line, *L2*. So seeing a point in one image is enough to narrow down its location in the second image to a single line.

We can go a bit further. If we observe an image point anywhere on the line *L1*, we get exactly the same epipolar plane, *Q*. This means that any point in the first image, that lies on *L1,* must match a point on the second image that lies on *L2*. And the reverse is true by similar reasoning. So *L1* and *L2* give us two sets of points that can only be matched to each other. When we know the camera geometry, the problem of matching points in 2D images can be reduced to the problem of matching points along 1D lines.

**The Epipole:** The epipolar plane depends on the scene point, *P*, and different scene points give rise to different epipolar planes and lines. However, the epipolar lines all have something in common; in general, they all intersect the same image point. This point is called the *epipole*. To see this, consider the line that goes through *f1* and *f2*. Call this line *F*. Since *f1* and *f2* are points on any epipolar plane, no matter which scene point produces this epipolar plane*,* this means that *F* is on always on the epipolar plane*.* We will call the point where *F* intersects the second image plane, *e2*. This is the epipole in image 2. For any scene point, *e2* is on the epipolar plane for that scene point, and on the image plane, so it is always on the epipolar line. This means that all epipolar lines intersect at *e2*.

There is one important special case where the above statements are not true. It is always possible that *F* does not intersect the second image plane, because they are parallel. This means both *F* and any epipolar line will lie in the same epipolar plane, but they will not intersect. This can only happen if *F* is parallel to any epipolar line. If all the epipolar lines are parallel to *F*, then they must all be parallel to each other. So, either the epipolar lines all intersect at the epipole, or they are all parallel to each other.

Image Plane *I1*

Image Plane *I2*

**Epipolar Plane** *Q*

*p1*

*p2*

*L1*

*L2*

*e1*

*e2*

Focal point *f1*

*F*

Focal point *p2*

**Standard Stereo:** So far we have discussed the general case of two images taken with any camera positions. However, it will often be convenient to consider the case of two cameras in a particularly simple position, analogous to the position of our two eyes. We suppose that the two cameras are side by side, with focal points at *z=0,* and *y = 0,* separated only in the *x* direction. For example, the focal points might be at the locations *(0,0,0)* and *(T,0,0)*. In addition to that, we assume that the cameras are pointing in the *z* direction with the same focal length, so we can assume that they both have the same image plane, *z = 1*. As explained above, all the epipolar lines will be parallel to the line that joins the focal points, which is the *x* axis. This means that all the epipolar lines are horizontal lines in the two images. If a scene point appears at the location *(x1,y1,1)* in the first image, it will appear along the line *y = y1* in the second image.

Now, let's suppose we see the scene point, *P,* in both images. Suppose that it appears in the first image at the location *p1 = (x1, y1, 1)*. Then the epipolar constraint tells us that it will appear in the second image at some location *p2=(x2, y1, 1),* ie., with the same *y* coordinate, but a different *x* coordinate. Then we can notice that we have two similar triangles, one involving *f1, f2, P,* and the other with *p1, p2, P*. The base of the first triangle has a width of *T,* while the second triangle has a width of *x2-x1*. We will define the *disparity* to be the difference in *x* coordinates caused by depth, that is, it is *d = (T-x2) – (0-x1) = T + (x1-x2)*. Now, the depth of the point, *P,* that is, its *z* coordinate, is given by the equation:
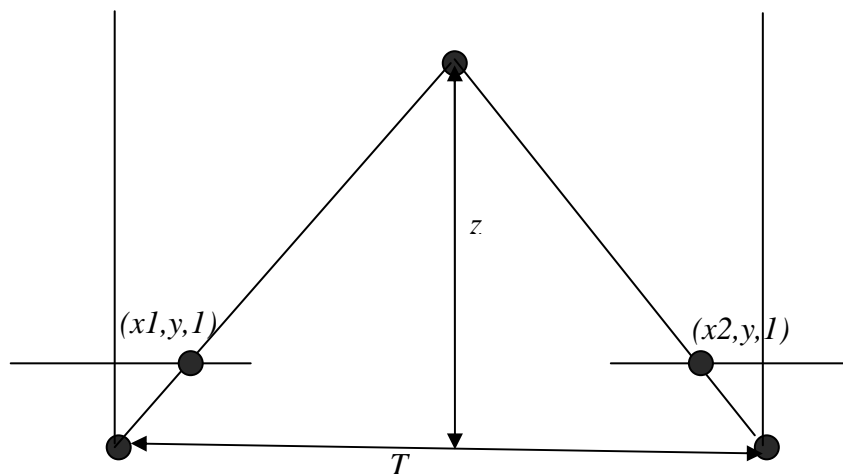
*(z-1)/z = (x2-x1)/T*

*Tz – T = z(x2-x1)*

*z(T + x1 – x2) = T*

$z = T/d$

That is, when we see a point in two images, its depth is inversely proportional to its disparity, and proportional to the distance between the two camera centers. So, if a point has disparity of zero, its depth is infinite. That is, points that are infinitely far away will appear in the same position in both images.



### Rectification

What if our cameras are not in this nice configuration? It turns out that we can rectify the images to produce the images that we would have gotten if our cameras were arranged like this. First, we note that we can just define our coordinate system so that the focal point of the first camera is the origin, and so that the line connecting the focal points is the *x* axis, and the second focal point is located at some position *(T,0,0)*. The only thing we have to worry about is the possibility that the image planes are not the *z=1*.

However, if we have an image taken with a particular focal point, *f,* and image plane, we can generate a new image that shows the world as seen by a camera with the same focal point, but a different image plane. Suppose our camera has an image plane *I,* and we want to generate an image with an image plane *J.* For any point, *p,* in *J,* we can find the line *L* that goes through *p* and *f.* We intersect *L* with the original image plane, *I,* obtaining the point *q.* The point in the world, *P,* that created the intensity at *q* lies on the line *L,* which connects *q* to the focal point. If we had taken a picture using *f* and *J* as our focal point and image plane, the line connecting *P* and *f* would still be *L,* and it would intersect *J* at the point *p.* So the same world point that generated the intensity at *q* would also create the intensity at *p.* So we can create the image at *J* by transferring the intensity from *q* to *p.*

P

L

q

J (the new image plane)

p

f

I (the original image plane)