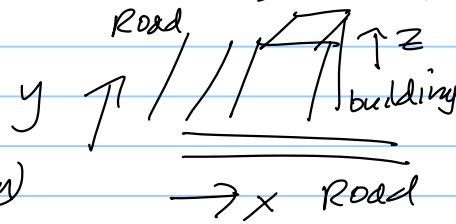# Manhattan World.

Many visual scenes have a Manhattan structure.

This give a natural 3-D coordinate system

A ground plane (x-y dimension) plus a vertical z-directions.

Road


y ↗ //// ↑ z
building
→ x  Road

Most city (new cities) have this structure. Natural scenes — e.g. mountains and rivers — do not.
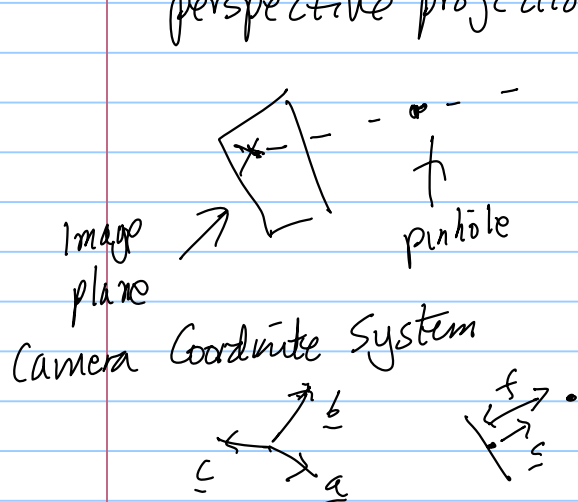
Two Questions:
   (i) How can a visual system use this knowledge to callibrate itself — ie. determine the angle of view compared to these x-y-z coordinates — and estimate these directions?
   (ii) How can we decide if an image is Manhattan or not?

Note: the model we describe is not the best model to do this. But it is simple and gives important ideas.

Note: Humans see to assume a Manhattan structure. The Ames shows the visual illusions that happen if the image (ie. Ames room) appears to be Manhattan, but is not.

First, we need a model of projection. This is perspective projection.



Image plane

pinhole

X A point in space is projected by a light ray (straight line) which goes through the pin hole and is stopped when it hits the image plane.
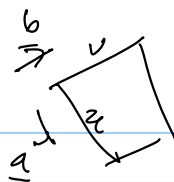
[Note: real cameras have a small aperture and not a pinhole]

Camera Coordinate System

f focal length.

c direction of gaze.

$|a| = |b| = |c| = 1$ , $a, b, c$ right angles

Coordinates in image plane

Point $\underline{r}$ in space

Projection rules:

$$u = -f\frac{\underline{r}\cdot\underline{a}}{\underline{r}\cdot\underline{c}} \quad, \quad v = -f\frac{\underline{r}\cdot\underline{b}}{\underline{r}\cdot\underline{c}}$$
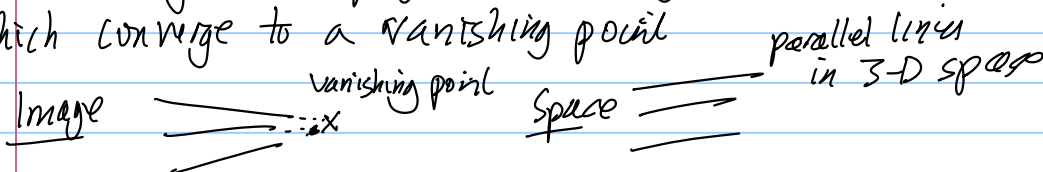
We specify $\psi = (\underline{a},\underline{b},\underline{c})$

Camera has a 3-D coordinate system $\underline{a},\underline{b},\underline{c}$.
Manhattan world has a 3-D coordinate system $x, y, z$.

How to calibrate the camera by estimating the transformation between them?

The projection rules mean that parallel straight lines in the image will projection to straight lines in space which converge to a vanishing point

Image $\quad$ vanishing point $\cdots x$ $\quad$ Space $\quad$ parallel lines in 3-D space

The positions of the vanishing points in the image depend on the orientation $\psi$ of the camera with respect to the Manhattan coordinate system.

Parallel lines in $x$ direction, vanishing points at $\left(-f\frac{a_x}{c_x}, \frac{-fb_x}{c_x}\right)$

" " " $y$ " , " " $\left(\frac{-fa_y}{c_y}, \frac{-fb_y}{c_y}\right)$

" " " $z$ " , " " $\left(\frac{-fa_z}{c_z}, \frac{-fb_z}{c_z}\right)$

where $\underline{a} = (a_x, a_y, a_z)$
$\underline{b} = (b_x, b_y, b_z)$
$\underline{c} = (c_x, c_y, c_z)$

Recall that $\underline{a},\underline{b},\underline{c}$ are orthogonal
$\underline{a}\cdot\underline{b} = \underline{a}\cdot\underline{c} = \underline{b}\cdot\underline{c} = 0$
$|\underline{a}| = |\underline{b}| = |\underline{c}| = 1$.

Now we introduce the model.

At each image pixel $\underline{u} = (u,v)$ there is a hidden variable $m_{\underline{u}}$ which indicates if the pixel is the image of an edge in the $x, y, z$ direction, $m_{\underline{u}} \in \{1,2,3\}$
an edge in a random direction, $m_{\underline{u}} = 4$
or not an edge $m_{\underline{u}} = 5$.

Let $E_{\underline{u}}$ be the response of a derivative filter at $\underline{u}$
e.g. $E_{\underline{u}} = \nabla I(\underline{u})$.

We put a prior probability on $m_u$.
$$P(m_u=1) = P(m_u=2) = P(m_u=3) = 0.02$$
$$P(m_u=4) = 0.04, \quad P(m_u=5) = 0.9$$

ie 90% of pixels in the image are not edges

In Manhattan 20% of edges are in the $x, y, z$ directions (each) the remaining 40% of edges are randomly assigned.

$$P(\underline{E}_u \mid m_u) = P(|\underline{E}_u| \mid m_u) \, P(\hat{\underline{E}}_u \mid m_u)$$

$$\underline{E}_u = |\underline{E}_u| \, \hat{\underline{E}}_u \qquad \leftarrow \text{unit vector}$$
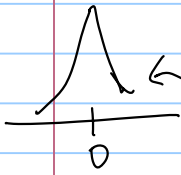$$\uparrow \text{magnitude}$$

for $m_u = 1, 2, 3, 4$

$P(|\underline{E}_u| \mid m_u)$ is specified by the models in lecture

$P(|\underline{E}_u| \mid m_u=5)$ by $P(|\underline{E}_u| \mid \omega(x) = 0)$ ie. $P(|\underline{E}_u| \mid \omega(x) = 1)$
$$\uparrow \text{no edge.} \qquad\qquad \uparrow \text{edge}$$

local image gradient direction →

$P(\hat{\underline{E}}_u \mid m_u)$ is uniform (all directions equally likely) for $m_u = 4, 5$.

$$P(\hat{\underline{E}}_u \mid m_u) = F\left(|\hat{\underline{E}}_u - \underline{n}(m_u, \underline{\Psi}, \underline{u})|\right) \text{ for } m_u = 1, 2, 3$$

where $\underline{n}(m_u)$ is the predicted direction of the edge — a function of $m_u, \underline{\Psi}$ and $\underline{u}$.



$F(\cdot)$ is peaked at 0, at predict direction, but allows some variation.

Then
$$P(\underline{E}_u \mid \underline{\Psi}, \underline{u}) = \sum_{m_u=1}^{5} P(\underline{E}_u \mid m_u, \underline{\Psi}, \underline{u}) \, P(m_u)$$

$m_u$ are nuisance variables, so can sum them out. (standard Bayes)

For the entire image measurements
$$E = \{\underline{E}_u\}$$

Assume
$$P(E \mid \underline{\Psi}) = \prod_u P(\underline{E}_u \mid \underline{\Psi}, \underline{u})$$
$$\uparrow \text{independence assumption}$$

Note: this independence assumption is extremely unrealistic. Edges in the $x$-direction are usually continuous — so if $m_u = 1$ then there is a higher probability that pixels near $\underline{u}$ also have $m_u = 1$

local context →

Really, we should have a model.
$$P(E \mid \underline{\Psi}) = \sum_{\{m_u\}} \prod_u P(\underline{E}_u \mid m_u, \underline{\Psi}, \underline{u}) \, P(\{m_u\})$$

where $P(\{m_u\})$ is a prior (which takes local context into account. But this model is harder to work with.

For the simple model    (without the context)

$$P(E \mid \psi)$$

we estimate    $\hat{\psi} = \underset{\psi}{ARGMAX} \; P(E \mid \psi)$

/ (Can be done by
exhaustive search)

This estimates the camera orientation relative to the Manhattan structure.

The results are good. We can display the estimated directions of the $x, y, z$ directions in the image to compare with the true directions.

How to know if an image is Manhattan or not?

Answer is model selection —

We define an alternative null model for generating the image which does not assume Manhattan structure.

This null model is the same as the Manhattan model, but we set $P(m_u = 1) = P(m_u = 2) = P(m_u = 3) = 0$.

$$P(m_u = 5) = 0.9, \quad P(m_u = 4) = 0.1.$$

This removes any dependence on $\psi$.

$$P_{null}(\{E\}) = \prod_{u} \sum_{m_u} P(E_u \mid m_u) P(m_u)$$

Model selection:

Input image $I \rightarrow$ Compute $\{E_u\}$    ⟵ best estimate of $\hat{\psi}$

⟵ threshold

If $\quad P_{null}(\{E_u\}) > P(\{E_u\} \mid \hat{\psi}) + T$

then image is not Manhattan,

Otherwise, image is Manhattan.

This gives good results → i.e.

classifies a city scene as Manhattan
classifies an image of fishes as non-Manhattan.