

Lecture 17 - W9L1

Regularization (measures contd.)

* Batch Normalization : # of examples

Batch outputs for current layer $\{z^{(i)}\}_{i=1}^q \rightarrow \{\hat{z}^{(i)}\}_{i=1}^q$

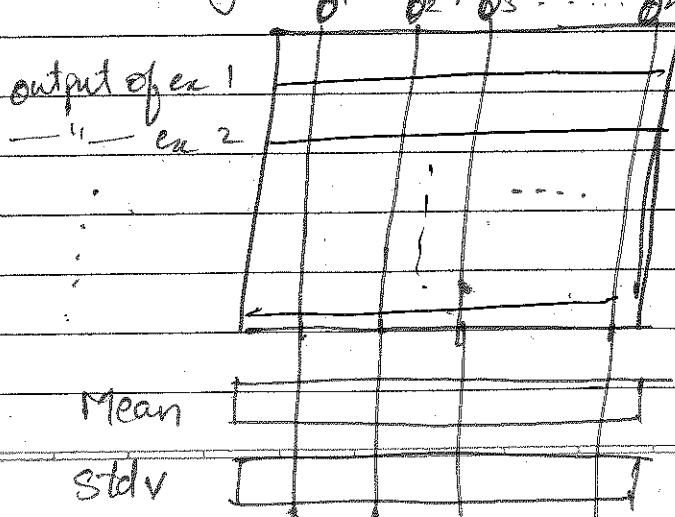
$\hat{z}^{(i)}$ EIRP
Per units

$$\hat{z}_j^{(i)} = \frac{z_j^{(i)} - \bar{z}^{(i)}}{\sigma_i} \quad \left[\begin{array}{l} M_j = \frac{1}{q} \sum_{i=1}^m z_j^{(i)} \\ \sigma_i = \left(\frac{1}{q} \sum_{i=1}^m (z_j^{(i)} - M_j)^2 \right)^{1/2} \end{array} \right]$$

Output of j-th unit for i-th batch example

- Makes sure activations are not saturated
- Normalization values are computed for each batch in training
- Normalization is differentiable (suitable for backpropagation).

→ Layer output matrix:



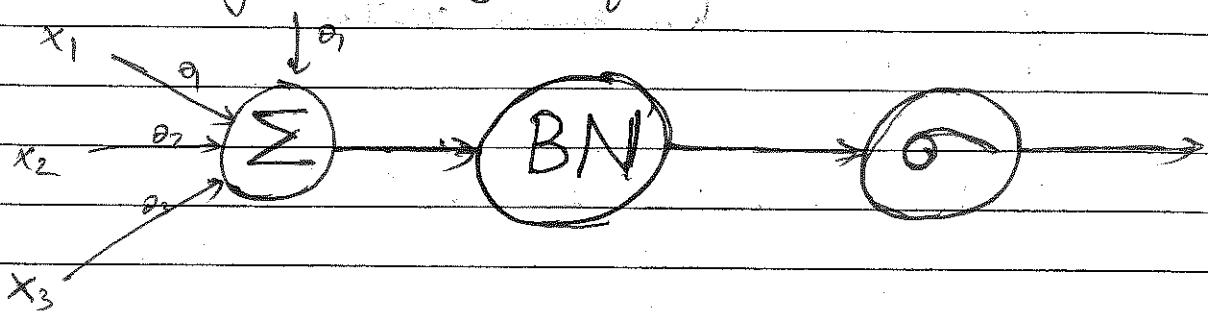
$$m = \text{Dmean} (\text{axis}=0)$$

$$D = D - m$$

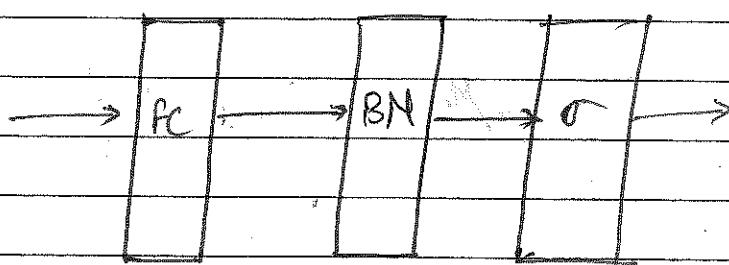
$$S = D \cdot \text{std} (\text{axis}=0)$$

$$D = D/S$$

→ Usually normalize before activation:



→ layers:



FC = fully connected
or convolution
layer without
activation

σ = Normalize
activation.

- Sometimes applied after activation

- B.N is not needed after every layer
- In layers with weight sharing (eg. convolution)
batch normⁿ for different spatial locations
need to be the same.
→ use one mean and stdv per output channel.
- Some saturation is needed to terminate learning
→ may need scale and shift after normalization

→ Scale and shift parameters:

$$\{z^{(i)}\}_{i=1}^n \rightarrow \{\tilde{z}^{(i)}\}_{i=1}^n \rightarrow \{\hat{z}^{(i)}\}_{i=1}^n$$

$$\tilde{z}_j^{(i)} = \sigma_j z_j^{(i)} + \beta_j$$

compute gradients
w.r.t. θ_j, b_j

- β_j and f_j are learned. The network can learn to run BN if there is no need for it, eg:

$$\gamma_j^o = \sigma_j \Rightarrow z_j^o = \sigma_j \cdot \hat{z}_j^{(i)} + M_j = \sigma_j \cdot \frac{\hat{z}_j^{(i)} - M_j + M_j}{\sigma_j} = \hat{z}_j^{(i)}$$

- During training because batches are random BN adds randomness into the training and so reduces overfitting
 - During prediction use the average normalization values computed during training (i.e. average from all training batches).

* Ensemble Classifiers:

- train multiple independent models
 - use majority vote or average during testing
 - reduces overfitting
 - to obtain multiple models:
 - change data
 - change parameters
 - Record multiple snapshots of the model during training.

* Summary:

- Start without regularization → See that can overfit
 - add Batch normalization
 - add dropout if needed
 - add decay if needed
 - augment data if needed
 - create ensemble if needed.
- } if observing overfitting

Hyperparameter Optimization

* Principle:

- select a set of parameters
- train on training data and test on validation data using a few epochs.
- Modify parameters in search of better validation error
- If loss is increasing (eg: x_2) stop.

* Methods:

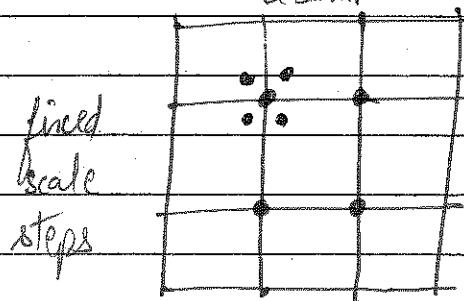
- coordinate search
- grid search
- Random search
- Bayesian model-based optimization
- Evolutionary

→ Search procedure: (coordinate search)

- hold all parameters fixed and change one
- Change parameter on a logarithmic scale to cover both small and large values (eg x_2 or x_{10} val)
- Once identifying best parameter search around it (coarse to fine search) using more epochs.
- If best parameters are found at end of range extend range.
- Repeat several passes on all parameters.

→ Grid and Random search:

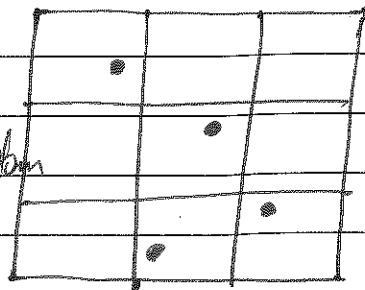
- search combinations of parameters
- as samples



grid search

Same
of points

points
are random
on scale



random search
generally pref-

In real we want to quasi-random sampling
(low discrepancy sequence)

→ random coverage

#

* Optimize parameters in this order :

- Learning rate (value, decay, update procedure)
- regularization (batch norm, dropout, ...)
- Network architecture and network size
- Other :
 - optimizer
 - optimizer parameters
 - Activation
 - Initialization

GPU Frameworks

* training is time consuming → use GPUs

* Frameworks :

By Google { Tensorflow
Keras

By FB { PyTorch
Caffe

Introduction to Keras

use google-colab

→ Keras is a Deep-learning framework for python

→ Key features:

- allows the same code to run seamlessly on CPU or GPU
- User-friendly API
- built-in support for convolutional networks
- built-in support for recurrent networks

→ Supported architectures:

- multi-input
- multi-output
- layer sharing
- model sharing
- arbitrary network architectures

Keras Architecture

- Model-level library, providing high-level building blocks
- low-level operations such as tensor manipulation and differentiation rely on backend engines.
- Run seamlessly on CPU or GPU (through Tensorflow which wraps access to cuDNN and BLAS)

Keras	
TensorFlow/Theano/CNTK/...	
CUDA/cuDNN	BLAS, Eigen
GPU	CPU

Keras Workflow

Tensors \rightarrow nD(data structure)

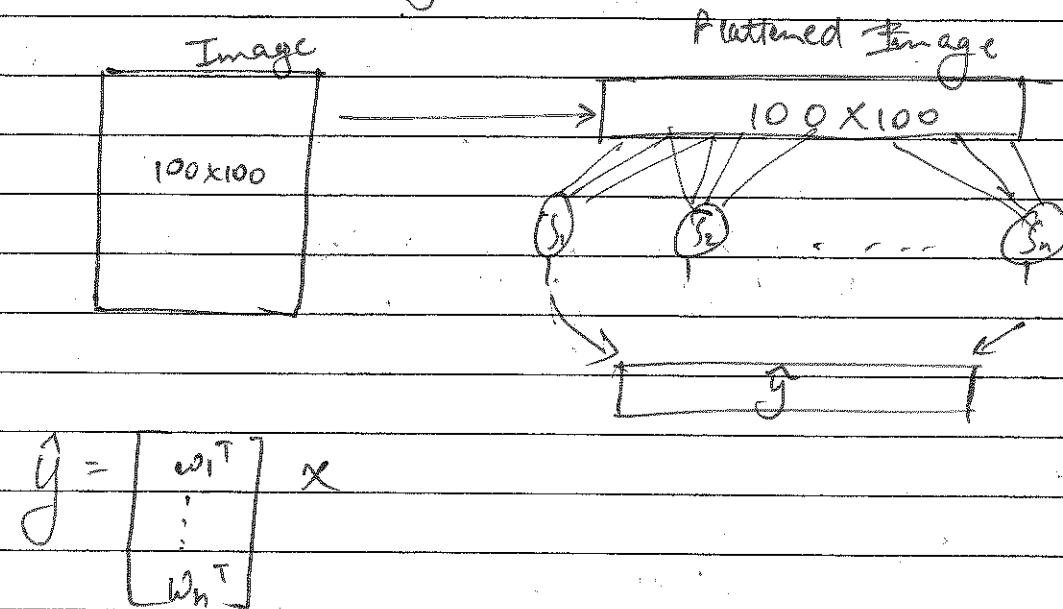
- outline
 - \rightarrow define training data : input tensors and target tensors
 - \rightarrow Define a model : network of layers that maps inputs to targets
 - \rightarrow Configure the learning process (loss function, optimizer, and metrics monitor)
 - \rightarrow train using fit() and search hyper-parameters

- Define a model
 - Model types
 - \rightarrow sequential class - linear stacks of layers (most common)
 - \rightarrow functional API - supports arbitrary architectures (eg. DAG). Specify function-like layers to process the data.
 - The remaining steps do not depend on model architecture
- Compile the network
 - Configure the learning process
 - optimizer
 - loss functions
 - metrics to monitor during training
- Learning
 - \hookrightarrow fit using Numpy arrays of input data and the corresponding target data

Lecture 18

Convolutional Neural Ne

Flattening an Image :



This is an ~~issue~~ issue in real world. Bigger input means more complex network.

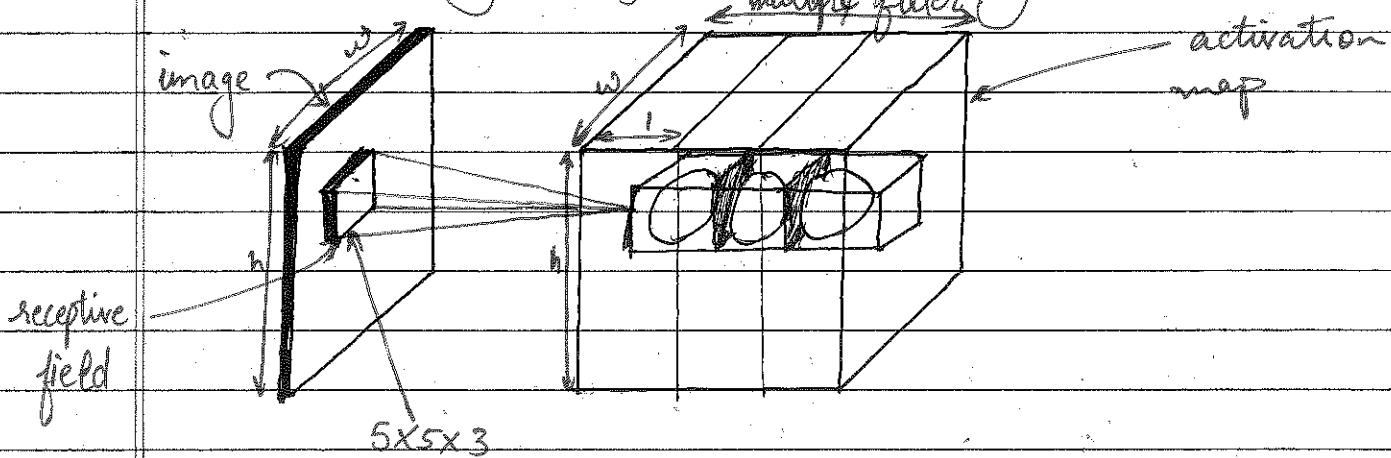
* Problems with flattening :

- ① A lot of coefficients for each unit
- ② loss of spatial context (lost neighborhood)
- ③ Shift invariance (same object at different locations.)

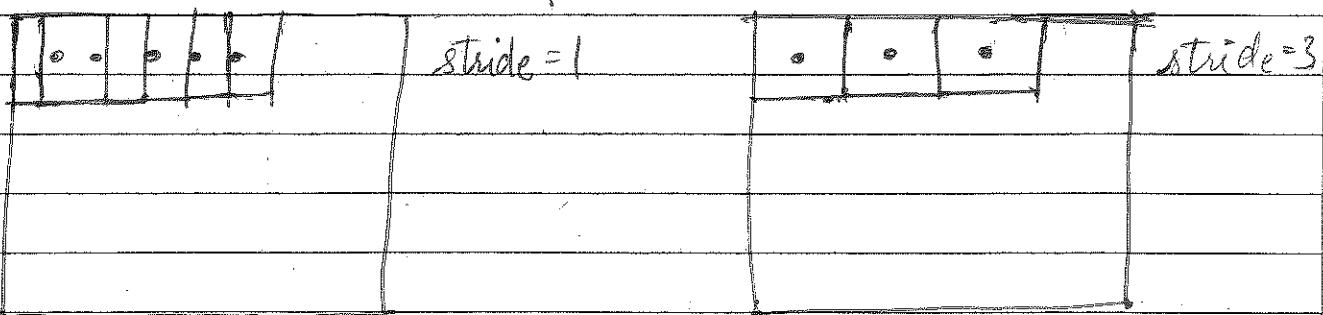
Solution is to use convolutional.

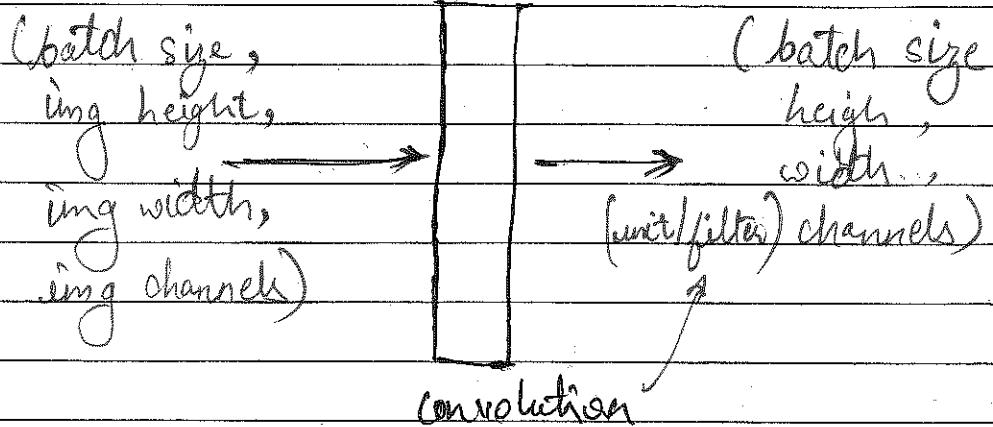
Convolution Layers

- convolve the image with small filters
eg. 3×3 or 5×5
- share weights of filter between locations
(shift invariance)
- sharing weights effectively increases data.

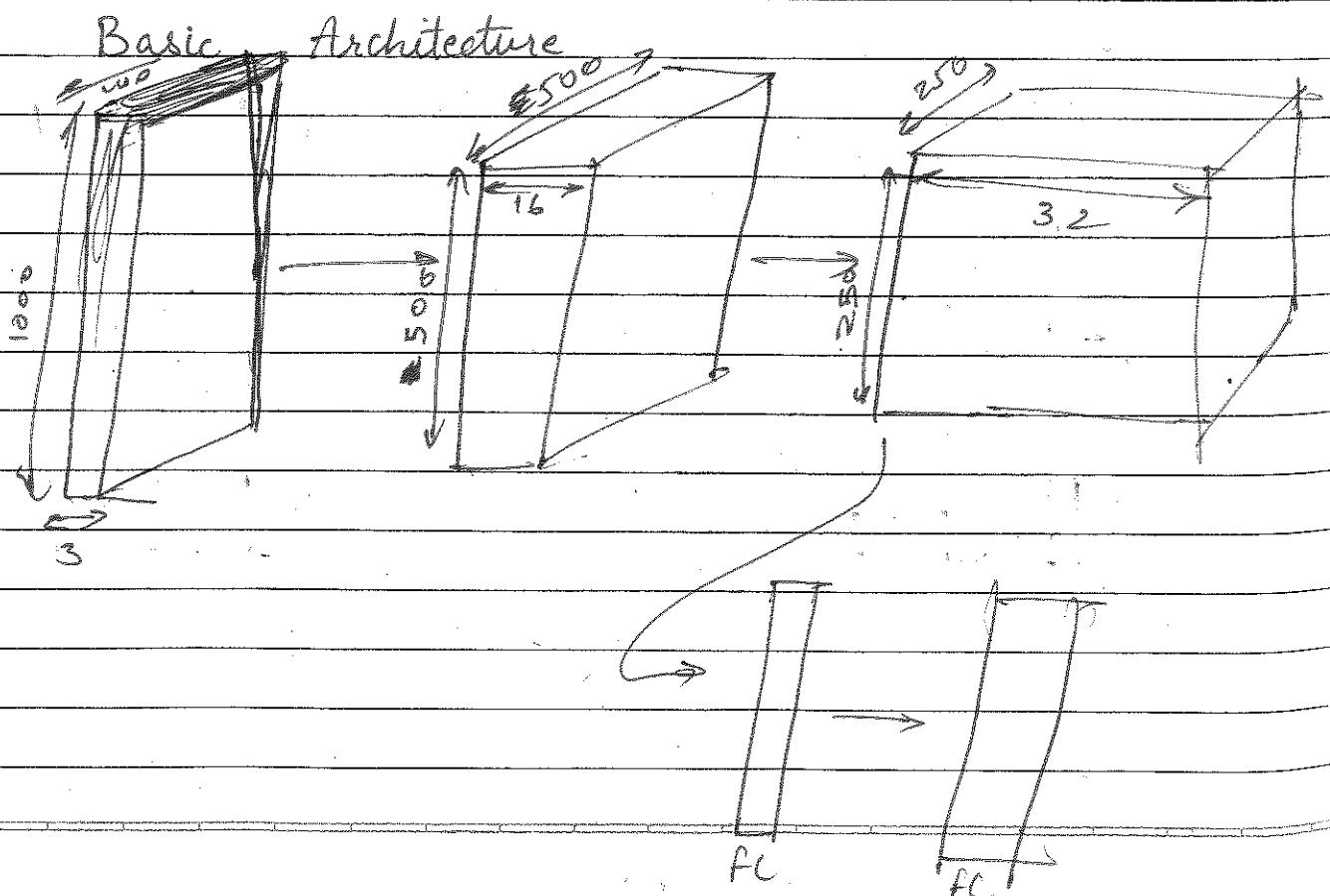


- extract image patches (windows)
- vectorize image window and filter and perform dot product (plus bias):
- Filter extends full depth of image
- Multiple convolution filters per location (eg oriented edges)
- Use stride to move filter \rightarrow activation map may be smaller.





- Multiple layers: as spatial dimensions decrease, depth increase to compensate for reduced coefficients (keep the same number of coefficients).
- Non-linear activation and sampling between layers
- Final FC layer: classification after feature extraction)



* Layer dimensions :

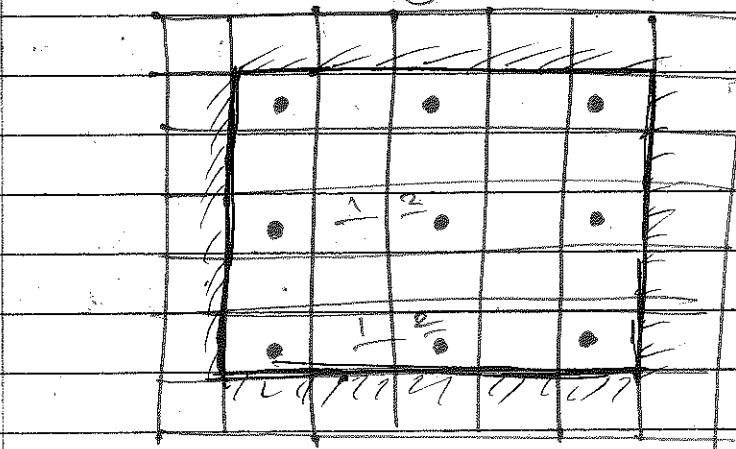
- zero padding is needed to prevent unwanted shrinking, especially with deep networks.
- width and height of the output of a layer depend on stride. With a stride of 1 layer overlap of receptive fields and large objects.

* 1×1 convolution :

1×1 convolution is used to reduce dimensions (i.e. lower # of channels).

* Dilated convolution (Atrous convolution) :

- Increase receptive field without increasing the number of parameters
- we consider value at ~~red~~ red dot represents the neighboring values.



2 dilated

Instead of 1D convolution ordinary:

$$(I * K)(t) = \sum_z I(t-z) K(z)$$

the:

$$(I * K)(t) = \sum_{z} I(t-lz)K(z)$$

↓ dilation factor

- take steps of size l in image when performing the convolution.

* Pooling

pooling = down sampling spatial dimensions (depth unchanged)

- strategies:

→ Max pooling (usually 2×2) - partition non-overlapping regions & choose max in each region.

- Alternatives:

- average pooling

- L2 norm pooling

- ROI pooling (output size is fixed and input variable)

- Pooling is done by scanning with a filter with stride - often stride is selected without filter overlap

- Convolution stride can also be used to downsample but this is a random
 - normally pooling instead of stride in convolutions
 - retains more info

- Pooling has no weights, while convolutions do have; so no learning.
- Depth dimension is normally not pooled.
- Pooling supports multiple scale analysis (pyramid).
- Pooling helps in reducing the number of coefficients.
- The amount of pooling is choice (hyper-parameter)

(Implementation in lecture)

Lecture 19 - W1OL1

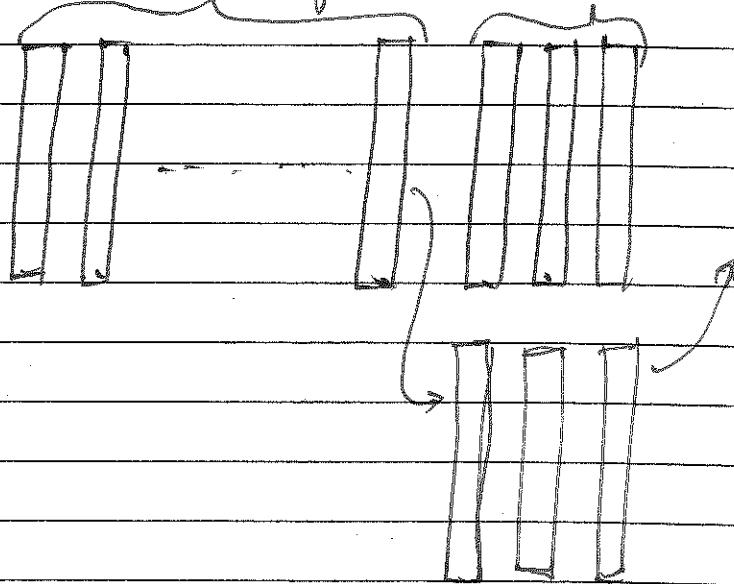
(CNN contd.)

* To solve the problem of limited data:

1) Data augmentation: add examples with perturbation (e.g. rotation, flip, contrast change).

2) Transfer learning: use pre-trained convolution base

conv. (fixed) FC



freeze loaded weights
of pre-trained blocks
after loading them.

* Networks that are pretrained can be used:

1) VGG (16 / 19)

2) Google Net / Inception (22 layers + fewer parameters)

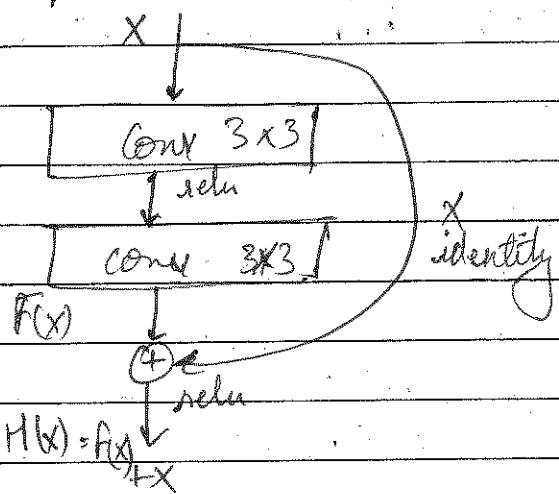
- To reduce lot of parameters due to inception, they use 1×1 parameter.

- Convergence can be difficult

- Use auxiliary classifiers to help with vanishing gradients (not needed with batch norm).

3) ResNet (uses residual block):

- Concept :



- Easier to learn $f(x)$ residual compared with $H(x)$ (learn deviation from identity instead of function)
- Skip connections help with vanishing gradients.

4) Mobile Nets :

- for mobile and embedded devices
- tradeoff b/w latency & accuracy
- Use depth-wise separable convolutions:

* Separable convolutions - separate convolution into :-
 - depth wise (channel) convolutions, followed by
 - point wise 1×1 convolutions

* Simplification parameters

- width ~~from~~ multiplier (fewer channels) [α]
- resolution multiplier (smaller resolution) [β]

[More in notes]

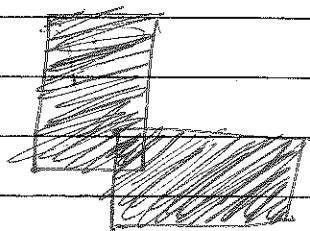
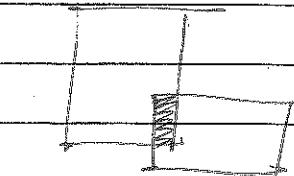
N AND SEGMENTATION (Using N.N.)

→ IoU / Jaccard index (similarity)

$$IoU = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

$$IoU \in [0, 1]$$

mutually exclusive
Some



→ Jaccard distance

$$d_J(A, B) = 1 - IoU(A, B) = \frac{|A \cup B| - |A \cap B|}{|A \cup B|}$$

sub(*) Mean Average Precision (MAP)

- commonly used to evaluate object detection
- whereas IoU only takes into account localization accuracy, MAP measures classification accuracy of detected object.

→ precision and recall:

true labels			
	P	N	recall = $\frac{TP}{\text{total positive}} = \frac{TP}{TP+FN}$
pred	P TP	F P	
	N FN	T N	Precision = $\frac{TP}{\text{Total P predictions}} = \frac{TP}{TP+FP}$

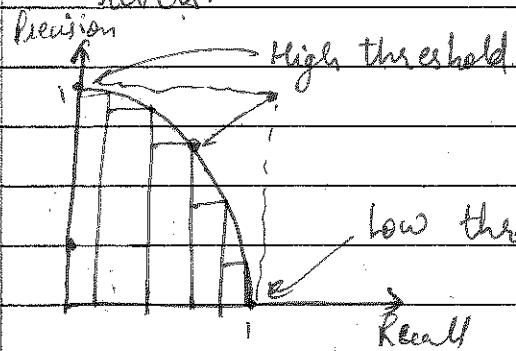
→ tradeoff between precision and recall:

- higher confidence threshold - higher precision
- lower recall

- Lower confidence threshold - lower precision
- higher recall

sub(*) ROC Curve:

Record precision & recall at different confidence levels.



AUC = Area under curve

low threshold

AUC $\in [0, 1]$ bigger is better

r_i recall at threshold t_i

AP = average Precision
(computes AUC)

P_i precision at threshold t_i

$$AP = \sum_{\substack{\text{confidence} \\ \text{threshold } t_i}} (r_{i+1} - r_i) P_{i+1}$$

Lecture 20 - W10 L2

(Detection evaluation (contd.))

* MAP implementation :

- Take AP for each class object then take the average for all classes.
- Include (no object) background as an object.

* IoU threshold

- In the context of detection, "to compute AP we first have to select a IoU threshold to determine there was a detection (then use AP to measure correctness)."
- What IoU threshold should we select?
 - lower IoU threshold - more detections but with worse localization and possibly less accurate.
 - higher IoU threshold - fewer detection but with better localization and possibly more accurate.

→ COCO map :

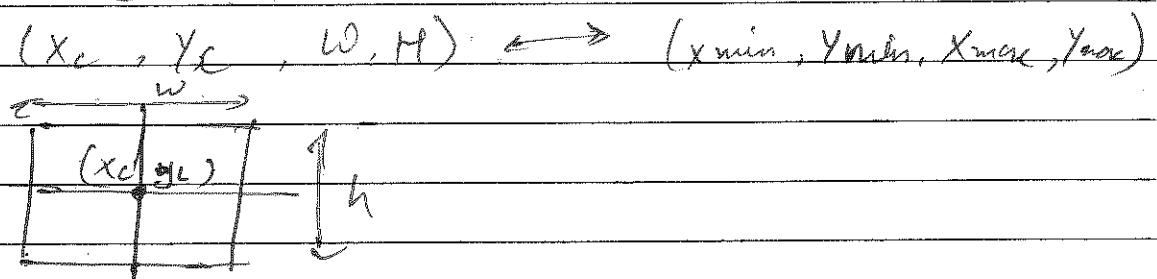
- calculate map at different IoU threshold values and average them.

AP @ (0.5 : 0.05 : 0.95)

$$\text{MAP} = \underline{\text{MAP}_{0.5} + \dots + \text{MAP}_{0.95}}$$

* Object localization

- to localize an object we need to specify a boundary box ($x_{min}, y_{min}, x_{max}, y_{max}$)
 - because we process the image at different scales, absolute coordinates cannot be used \rightarrow use relative location
- $\rightarrow \in [0, 1]$
- Alternatively use



* Training:

- for each example mark bounding box and label class
- include background (no object) boxes.
- one-hot encoded class labels.

$$y = [P_c, x_c, y_c, w, h, c_1, \dots, c_k]$$

↓ one-hot
encoded class
 if object box if not object box

↑ bounding
box

* Classification:

$$\hat{y} = [P_c, x_c, y_c, w, h, c_1, \dots, c_k]$$

↓ class probabilities
 obj prob bounding box

* Loss :

$$L = L_{\text{reg}} + L_{\text{cls}}$$

regression

eg. cross entropy
(has to do with box)

Classification (has to do with img)
eg L2. inside box

$$L_i(y^{(i)}, \hat{y}^{(i)}) = \begin{cases} (y_1^{(i)} - \hat{y}_1^{(i)})^2 & \text{if } y_1^{(i)} = 0 \\ \sum (y_j^{(i)} - \hat{y}_j^{(i)})^2 & \text{otherwise} \end{cases}$$

use cross entropy for class labels

* Multi object detection :

→ Approach :

- ① find candidate
- ② classify candidate
- ③ possibly refine candidate

→ Finding candidate :

- sliding window
- grid cells
- region proposal

* Sliding window detection

→ slide window through image.

① crop subwindows

② classify each sub-windows

→ To deal with different object scales, crop different sub-windows with different sizes and aspect ratio at each location.

→ May be slow - many windows to classify.

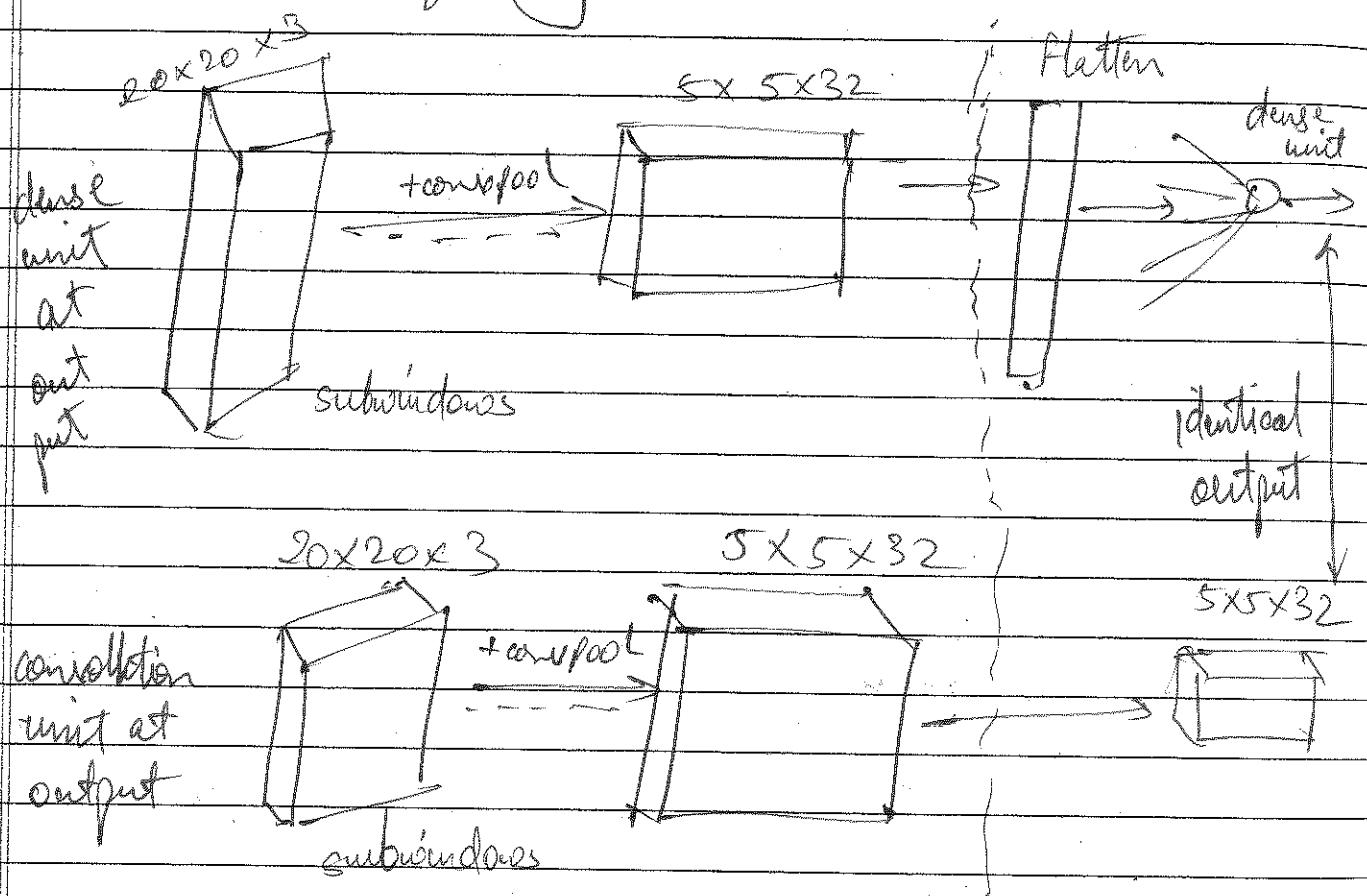
→ problems:

- running dense layers for classifying each ~~sub~~ subwindow is inefficient.
- in addition dense layers need a fixed input size and so the input size need also be fixed.

→ solution (fully convolutional implementation):

- A FCI is more efficient: no need to process each sub-window separately and convolution can be implemented efficiently.
- A FCI can work on variable image size.

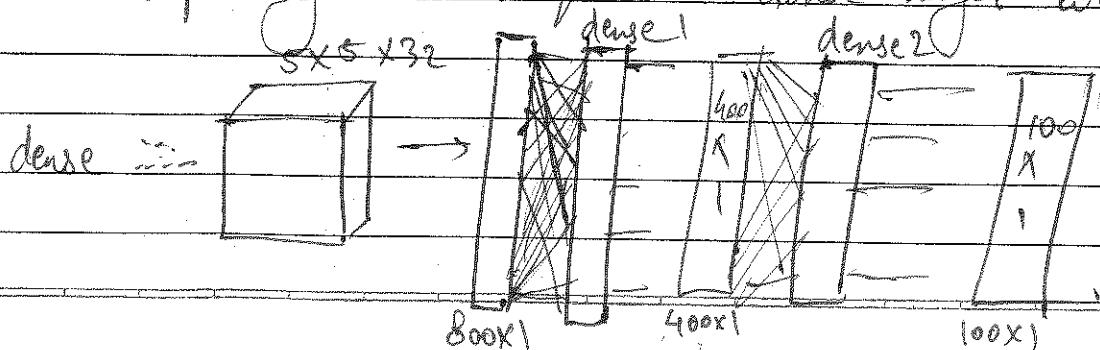
→ A single fully connected (dense) unit may be replaced by a single convolutional unit that has the dimension of the feature tensor (without flattening)

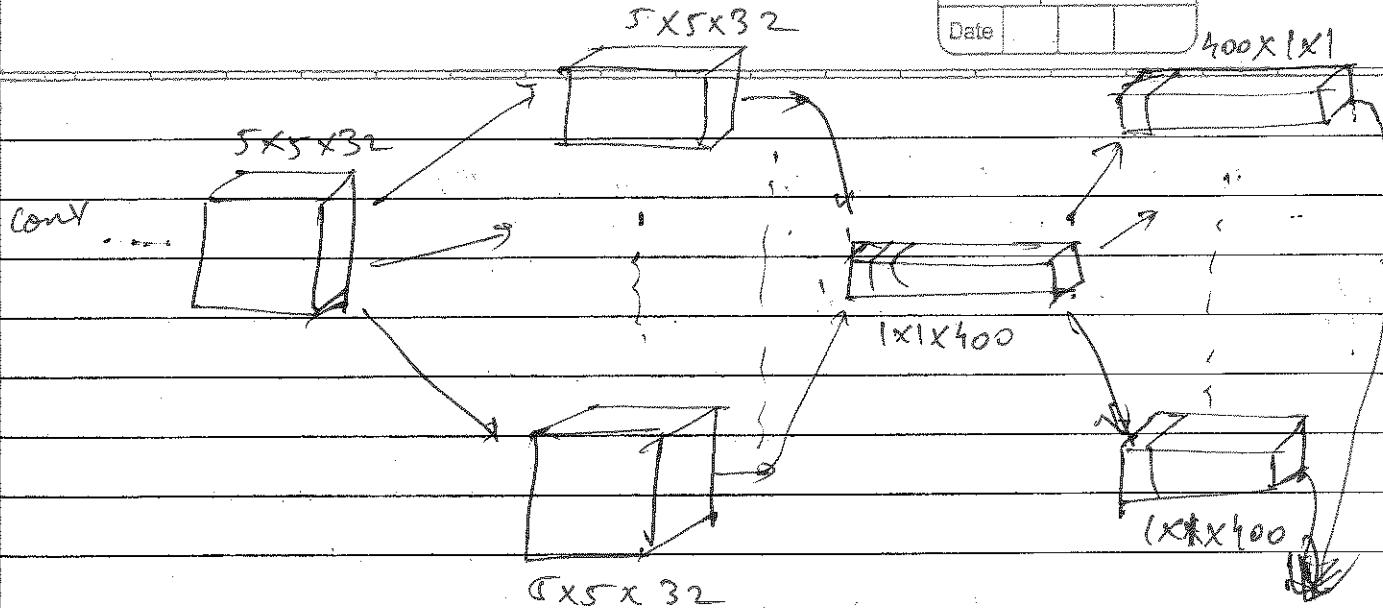


→ consider replacing 400 of dense units

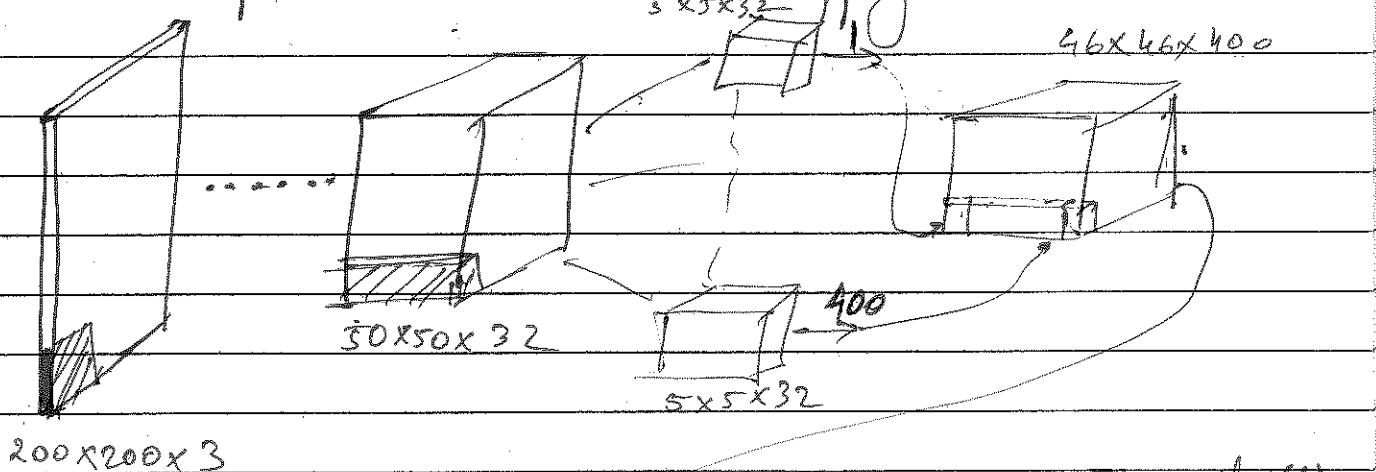
- output is 400×1 vector
- we can replace 400 units with 400 conv filters
- their output will be 400×1 vector

→ replacing a subsequent dense layer with 100 units

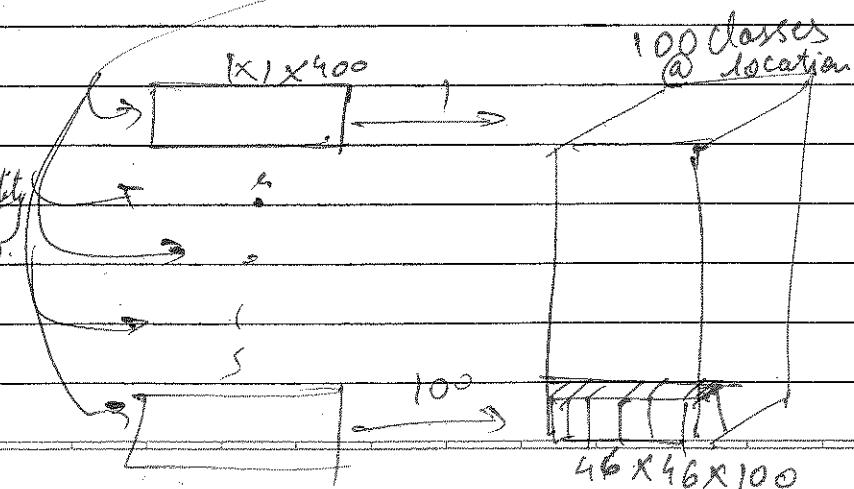




- So far a single subwindow classification.
- for a single sub-window no reduction in computations
 - The advantage of convolutions is when sharing computations between overlapping windows.



- each column at output provides probability of a class at a window.



* Summary: (sliding window detection)

- use conv + pooling to extract features and reduce spatial resolution.
- each location in reduced resolution layer corresponds to 'virtual' patch in original image. virtual patches overlap.
- use dense layer or 1×1 conv to classify or regress at each location in reduced resolution layer.

* Problem (SWD):

- does not produce accurate detection because:
 - It is based on classification in a pooled (lower resolution) layer.
 - We detect fixed size boundary boxes

* Solution (SOD):

To solve problem we can, in addition to classifying each location, regress a bounding box there!

Lecture 21 - W11 L1

Page No.	
Date	

* Grid Cell Detection :

- Divide the image into grids and attempt to scan an object. In each grid cell
 - Cells do not necessarily match objects and so use regression to find a bounding box for each detected object
- Similar to object localization (single object) except that now at each grid cell detect and localize objects
- At each grid cell train and predict using:

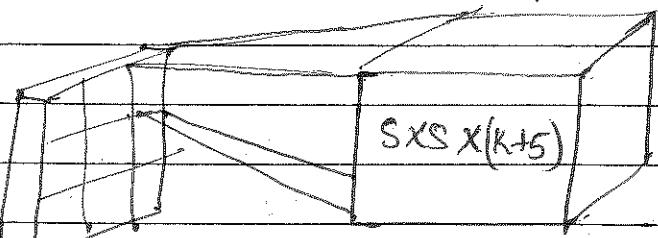
$$y = [p_c, x_c, y_c, w, h, c_1, c_2, \dots, c_k]$$

↗ regression
(bounding box) classification
(object class)
 classification
(object yes/no) (c₁) P(c₂) ... P(c_k)

\hat{y} for prediction

- x_c, y_c, w, h are relative to each grid cell
 $x_c, y_c \in [0, 1]$ whereas w and h may be bigger or smaller than 1 for boxes larger than a cell.
- For $S \times S$ grid and K classes the output is a tensor of size:

$S \times S \times (K + 5)$



- at each grid cell with object a precise bounding box instead of a fixed window.
- only one object in each grid cell.

* Non-max suppression (NMS):

Since we might get multiple of object detections for adjacent grids when the object is large we use NMS.

- select detection in grid cell with highest score (probability p_c) and delete detections that have $\text{IoU} > \tau$ with selection.
- perform NMS for each class separately to allow object overlap.

Intersection over union

* problems (NMS):

- top scoring box may not be the best fit
- It may suppress nearby objects.
- It does not suppress false positive.

* Problems (with grid cells):

1) what should be grid size?

- Too small will fail in classify big objects and vice versa.

2) what should be the aspect ratio of candidate cells?

- different objects have different aspect ratio of boxes that should be used for their detection

3) How to detect more than one object in each

candidate cell?

Solutions: Anchor boxes

- to solve the three problems, use multiple anchor boxes at each grid location.
- Anchor boxes have different sizes and different aspect ratios (select according to what we see in training data).

- Multiple anchor boxes (multiple objects at same grid):
 - To allow detecting two different (overlapping) objects at each grid location (eg objects with different aspect ratios) double the target vector y .

$$y = [P_c^1, x_c^1, y_c^1, w_c^1, h_c^1, c_1^1, \dots, c_k^1; P_c^2, x_c^2, y_c^2, w_c^2, h_c^2, c_1^2 \dots c_k^2]$$

$P_c^i = 0$ and $P_c^2 = 0$ if object does not exist

- most grid cells will have zero or one anchor boxes.

• Anchor box assignment during training:

- During training, the anchor box assigned to an object is the one having highest IoU with ground truth box.

* Two classes of methods for object detection:

• Single shot vs. Double shot methods:

- Single shot detectors perform object detection in a single shot network that looks at the image once.

- Two shot detectors look at the image twice: once

to propose object regions (region proposal) and a second to refine and classify regions.

- two shot methods have higher computational cost but may be more accurate (lower frame rate).

- Single shot methods :

- consider predefined anchor boxes at grid cells
- refine and classify region proposals

- Two shot methods :

- consider pre-defined anchor boxes at grid cells
- classify them as object or not \rightarrow region proposal
- refine and classify region proposals.

- Single shot detectors (popular) :

- You only look once (YOLO)

- Multibox

- Single shot multi shot detect (SSD)

→ Yolo algorithm :

- divide image to $s \times s$ grid cells. A grid cell is responsible for detecting objects with ~~the~~ centers in it.

- Each grid cell predicts 8 bounding boxes and confidence scores. (confidence that the box contains the object and is accurate).

$$\text{confidence} = p(\text{object}) + IoU(\text{truth, prediction})$$

confidence = 0 if no object in cell.

→ Multibox:

- use prior for anchor boxes determined per class
- Inception-like detection (detect at different scales).

→ Loss:

- confidence loss - class categorical crossentropy
- location loss - ~~bounding box~~ IoU

$$\text{so, Multibox loss} = \text{confidence loss} + \alpha * \text{location loss}$$

→ Single shot multibox detect (SSD):

→ SSD = YOLO + multibox

- use feature maps from intermediate convolution layers.
- Use different prior for different convolution layers.
- Apply ~~no~~ NMS.
- Use hard negatives (False positive by the models) to train model instead of all negatives.
- this helps balancing negatives & positive examples.

- Two shot detectors:

RCNN - region detection CNN

[region proposals without network. SVM classifier]

Fast-RCNN - CNN classifier

faster-RCNN - use region proposal network

R-FCN - fully convolutional (no depth layer)

Mask RCNN - add object segmentation (instance segmentation)

Lecture 22 - W11/2

(contd.)

→ RCNN :-

- Instead of processing candidates on a grid use superpixels (select search algo) to find candidates.
- Extract 2000 proposals
- Advantages :
 - candidates have different aspect ratios and scales
 - fewer candidates
- save and process each candidate separately (slow).

→ Fast-RCNN :-

- Convolution implementation instead of processing one candidate box at a time.
 - Faster ($\times 200$) than RCNN but slower than YOLO.
 - fast-RCNN architecture : [refer notes]
 - ROI pooling - use max pooling to convert the region of interest to a fixed size (ex. 7×7) to the fc classifier.
- region of interest

→ faster RCNN :

* Region proposal Network (RPN) :

- train CNN to produce region proposals

- light weight CNN (relatively simple task)

- motivation :

- 1) speed

- 2) complete end-to-end network & train

- Algorithm (RPN) :

- 1) slide window and classify (object/noobject)
use fast convolution implementation
(sliding windows)

- 2) use K anchor boxes at each location
(different size and aspect ratio)
- around 200K boxes.

- 3) classify 1 regres K -boxes:

cls - output $2K$ classification scores

reg - output $4K$ regression boxes

$$\text{Anchor boxes} = W \times H \times K$$

- 4) NMS - select candidate with highest confidence and delete candidates overlapping ($\text{IoU} > T$) with it.

- faster RCNN = fast RCNN + RPN

- use RPN for region proposal

- use ROI-pooling then classify using dense layer

- Faster RCNN loss:
 - RPN [classification + box regression] loss (obj)
 - classification [classification (class) + refine box regression] loss.

bounding box losses use L_2 loss & classifications use cross entropy.

- faster RCNN algorithm:

- 1) get pretrained CNN (eg VGG16 or ~~resnet~~ resnet 101)
- 2) get feature map from last (deep) conv layer.
- 3) ~~RPN~~ ^{apply} RPN: obj / no obj and box regression.
- 4) Apply ROI-pooling to proposals
- 5) Use FC (dense) layers for class. classification and refined box regression.

→ RFcn (region based fully convolutional Network):

- Use convolution, instead of costly dense layers that are applied for each region
- events position sensitive score maps where each score map is sensitive to another region

$K \times K$ grid $\rightarrow K^2$ score maps

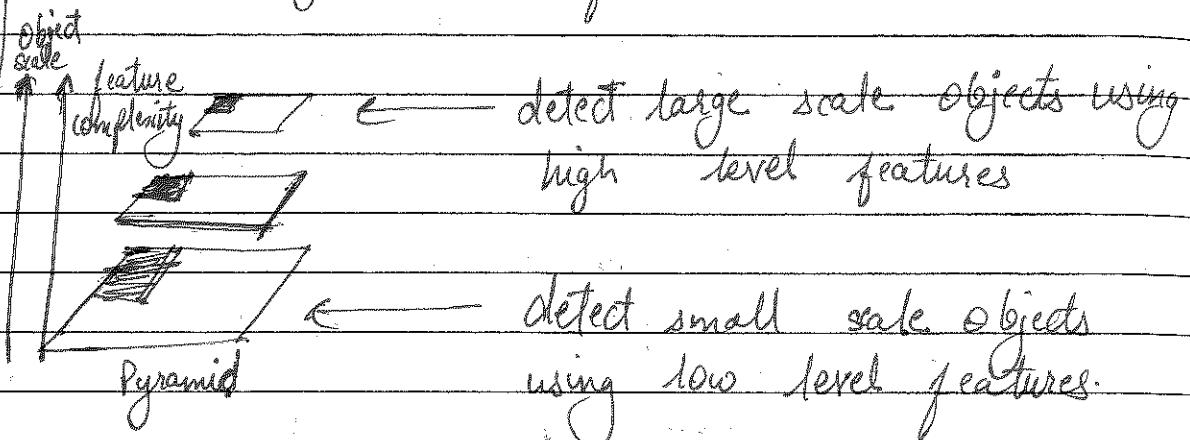
each score map has $C+1$ channels
(C classes and 1 obj/noobj class)

- position sensitive ROI pooling: [refer diagram]
- ROI architecture [refer diagram]

→ mask RCNN : (also does instance segmentation)

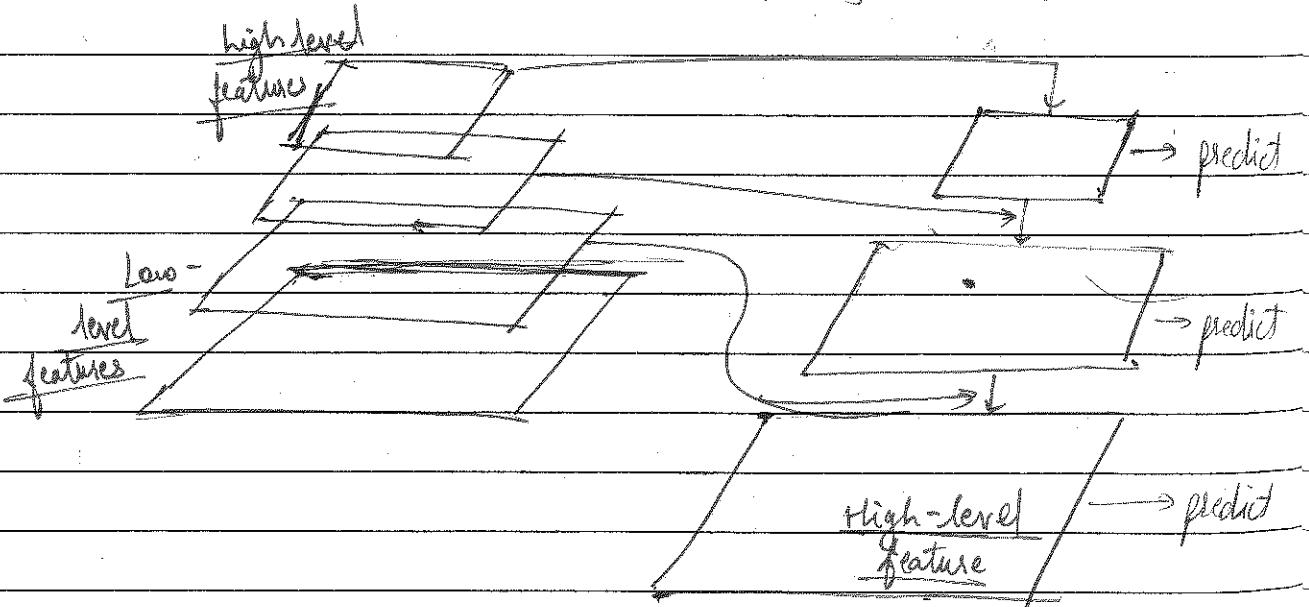
* Feature pyramid network (FPN) :

- goal - Detect objects at different scales using higher level features.



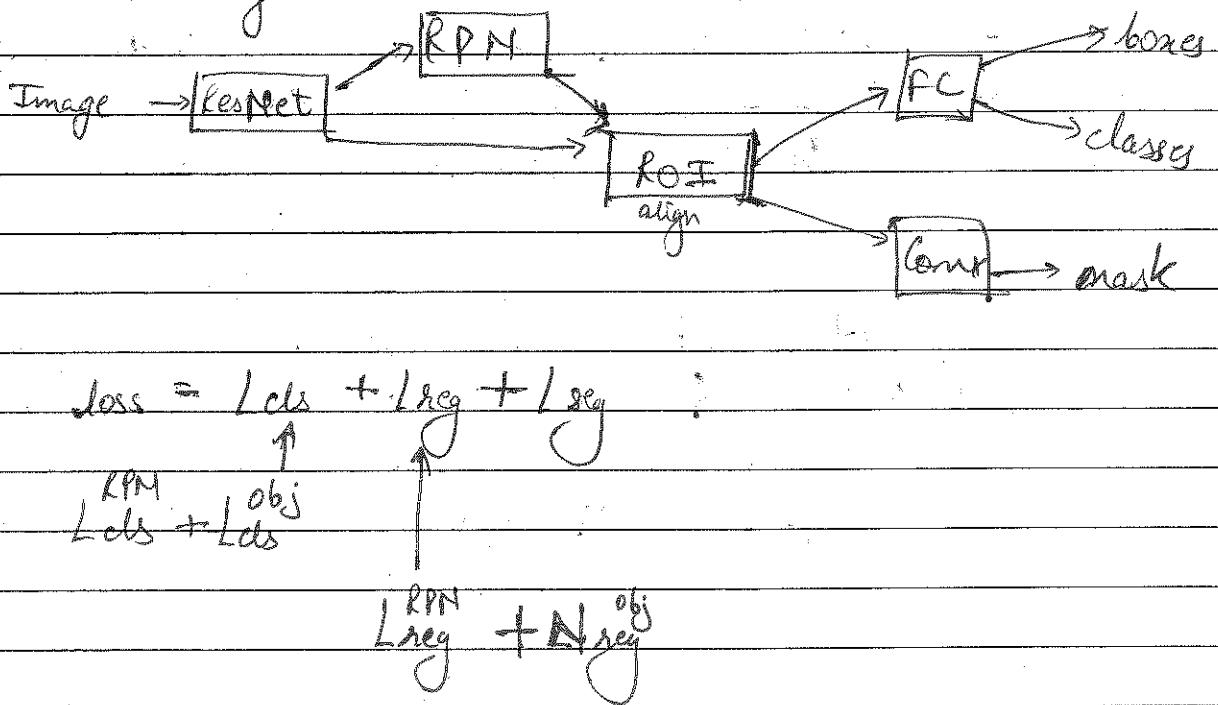
But we want to detect both large and small scale objects (using FPN).

Sol: Lateral connections in FPN.



- maintain strong semantic features at each level.
- lateral connections combine features at corresponding levels.

- mask RCNN adds segmentation to faster RCNN instance segmentation

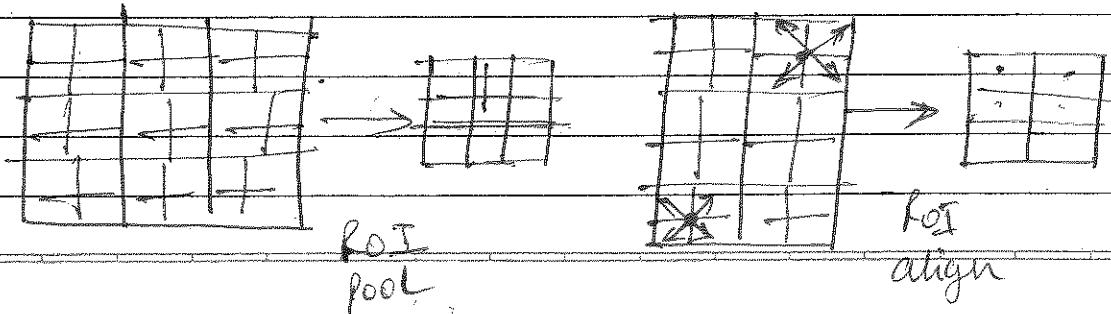


- predicted masks are scaled up to original resolution.

- Masked -RCNN architecture [refer diagram]

- ROI align:

- ROI pooling extracts small feature maps of fixed size.
- ROI-pooling is sensitive to alignment - shifting to grid will change results.
- ROI-align solves this by sampling with interpolation



* Semantic Segmentation:

- segment (cluster) according to class (eg. road)
- density each pixel
- Evaluate using mean IoU. (of masks)

* Methods

- FCN (Fully convolutional segmentation)
- U-Net: Add skip connections between encoders and decoder branches.
- Deeplab: add spatial pyramid pooling (SPP)

→ FCN

- Uses convolutions to extract features then classifies
- 32, 16, 8 architectures.

→ UNet

- Uses encoder-decoder architecture

- Encoder → convolve & pool (depth doubles)
- then feature extract
- decoder then at end classifier
- Uses lateral connections to decode without loss.

- Skip connections

- skip connections between encoder & decoder
- upsample output is concatenated with cropped input (cropped due to loss of pixels in conv)

- Fully convolutional:

3×3 convolution: produces segmented output (probability). The number of filters in final convolutional layer = number of classes.

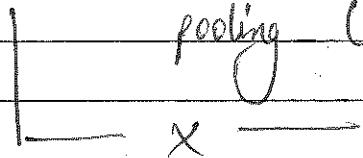
- DeepLab:

* Spatial pyramid pooling (SPP):

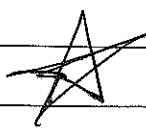
goal: support objects at different scales

multiple pooling levels as input to the classification layer. [refer diagram]

- spatial pyramid pooling results in increased dimensions → use atrous (dilated) convolution
- This results in atrous spatial pyramid pooling (ASPP)



- DeepLab uses encoder-decoder architecture
- Use Atrous spatial pyramid pooling (ASPP) with skip connections. [refer diagram]
- convolution implemented as depth wise separable convolutions (as in mobilenets).
- DeepLab architecture [refer diagram]



Camera Calibration

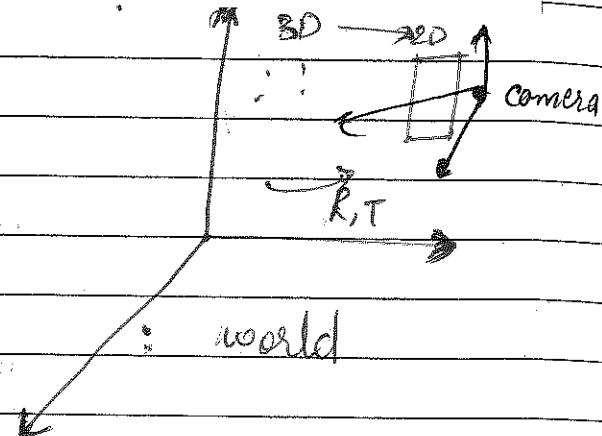
* Recall:

$$P_i^o = \underline{M} - P_i$$

↑ ↑ ↙

image point (2D)
projection matrix
(3x4)

world point (3D)
4x1



Question?

- ① what is M ?
- ② How to go 2D to 3D?

* Definition of M (already learnt):

$$M = \underline{K^*} [R^* | T^*]$$

3x3 3x3 3x1

3x4 ↖ ↘ ↙

Intrinsic parameters extrinsic parameters

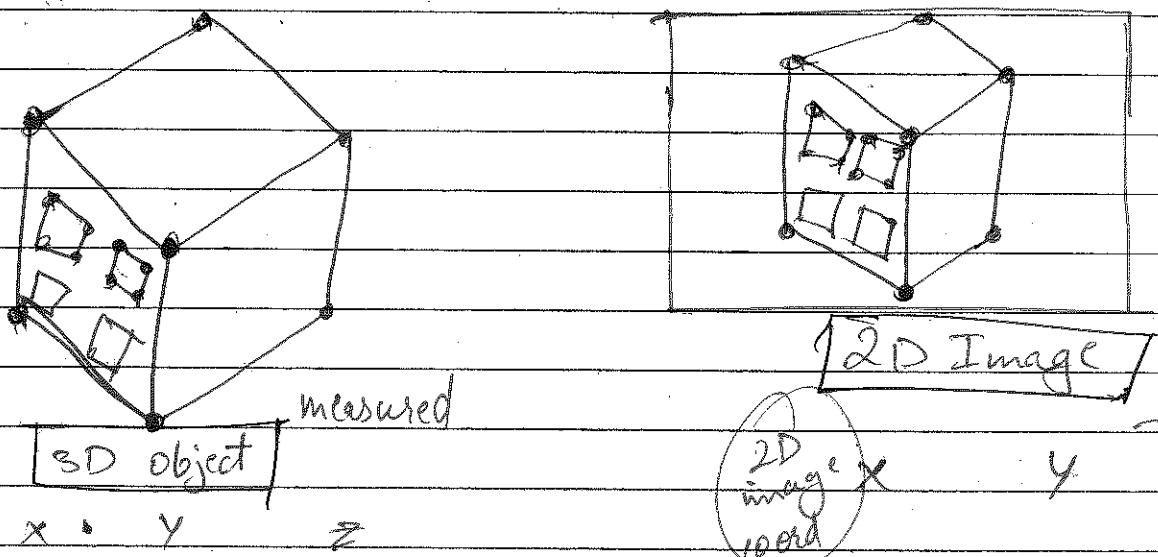
$$K^* = \begin{bmatrix} K_u & S & u_0 \\ 0 & K_v & v_0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\begin{aligned} K_u &= f_k u \\ K_v &= f_k v \\ S &= K_u t g \theta \end{aligned}$$

* Calibration

$$M \rightarrow K^*, R^*, T^*$$

$$\rightarrow R, T, f, K_u, K_v, u_0, v_0, t_{ye}$$



3D world
world

3.2 1.7 1.8

20.0 50.0 Input

for
calibra
align

given corresponding point $\{P_i\}_{i=1}^m \xleftrightarrow{\text{pixels}} \{P_i\}_{i=1}^m \text{ meters}$, find parameters

Lecture 23 - W12 L1

Page No.	
Date	

(contd)

precision: order of magnitude larger than desired.

* Calibration approaches:

methods:

- planar (require 2D) (requires multiple views)
- non-planar (eg (UBB))

approaches:

- direct
- indirect (find M then parameters)

- During calibration we need to solve

A is not square

$$A^T A x = 0$$

solution \Rightarrow zero eigenvalues of $A^T A$

or

singular value decomposition

using SVD: $A = U D V^T$

solution = column of v belonging to zero singular value of A .

(right null space / SVD)

$$A = \underset{m \times n}{U} \underset{m \times m}{D} \underset{n \times n}{V^T}$$

orthogonal orthogonal orthogonal

properties :

- 1) columns of U are orthogonal
- 2) columns of V are orthogonal
- 3) D is diagonal with singular values on diagonal
- 4) columns of U are eigenvectors of $A^T A$
- 5) columns of V are eigenvectors of $A^T A$
- 6) a matrix P is singular if $\text{cond}(P) \leq 10^{-3}$ (condition number)
- 7) for a symmetric matrix $U = V$ and singular values are squares of eigenvalues

So question is what is x in $A^T A x = 0 \Rightarrow x = ?$

- ① zero eigenvector of $A^T A$
- ② column of V belonging to zero singular value of A .
(right null space)

So given $\{p_i\}_{i=1}^m \leftrightarrow \{e_i\}_{i=1}^m$ find camera parameters

steps :

- ① find projection matrix M
- ② find parameters from M

Estimating the Projection Matrix

2DH 3x4 3DH

Known

$$P_i^* = M P_i$$

Projection Matrix

$$P_i^* = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

pixels

$$P_i^* = \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

world

so,

$$\begin{bmatrix} x' \\ y' \\ z' \end{bmatrix} = \begin{bmatrix} m_1^T \\ m_2^T \\ m_3^T \end{bmatrix} P_i$$

$$x_i^* = \frac{x'}{w'}$$

$$y_i^* = \frac{y'}{w'}$$

Homogenized P_i^* after projection.

$$P_i^* = (x_i^*, y_i^*)$$

so, $x_i^* = m_1^T P_i$
 $y_i^* = m_2^T P_i$
 $w_i^* = m_3^T P_i$

$$x_i^* = \frac{m_1^T P_i}{m_3^T P_i} \Rightarrow m_1^T P_i - x_i^* m_3^T P_i$$

$$y_i^* = \frac{m_2^T P_i}{m_3^T P_i} \Rightarrow m_2^T P_i - y_i^* m_3^T P_i$$

so,

$$m_1^T P_i - x_i^* m_3^T P_i = 0$$

$$m_2^T P_i - y_i^* m_3^T P_i = 0$$

for each point

2 equations with
12 unknowns

\Rightarrow need at least 6 values.

for a single point

$\rightarrow [Ax = 0]$ form

$$\begin{bmatrix} \mathbf{P}_i^T & 0 & -\mathbf{x}_i \mathbf{P}_i^T \\ 0 & \mathbf{P}_i^T & -\mathbf{y}_i \mathbf{P}_i^T \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ m_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

2×1 12×1

Known Unknown

So far : m points :

1st point	P_i^T	0	$-x_i^0 P_i^T$	m_1	0				
	0	P_i^T	$-y_i^0 P_i^T$						
$m_2^T P_i^0 = 0$		$m_2 =$	m_3	12×1	0				
$m_3^T P_i^0 = 0$		m_3	0	0	0				
n th point	P_m^T	0	$-x_m^0 P_m^T$	m_n	0				
	0	P_m^T	$-y_m^0 P_m^T$						

$$m = \#\text{points} \quad 2m \times 12 \quad 2mx1$$

$$As \quad \del{A} \quad Ax = 0$$

$A^T A x = 0$ so use SVD to solve using just A .

$$A = UDV^T$$

solution = column of v belonging to zero singular value.

$$\hat{x} = \begin{bmatrix} \hat{m_1} \\ \hat{m_2} \\ \hat{m_3} \end{bmatrix} \Rightarrow \hat{M} = \begin{bmatrix} -\hat{m_1}^T \\ -\hat{m_2}^T \\ -\hat{m_3}^T \end{bmatrix}$$

the solution is not unique $A\vec{x} = 0$
 we also have $A(\beta\vec{x}) = 0 \Rightarrow \beta\vec{x}$ is a solution

Hanson

20H

10

$$P_i \xrightarrow{\pi} 3MP_i$$

homogeneous

We need to avoid this factor (β) because it changes everything R^* , T^* etc. Since A is to be decomposed,

So stepwise:

10

solved system
of eq M

② Take \hat{M} and give parameters

$$\hat{M} \rightarrow K, R, T$$

$$(f, u_0, v_0, k_u, k_v)$$

So now we do ② :-

we find S such that :

$$M \equiv K^* [R^* | T^*] = S \hat{M} \quad S = ?$$

so we need to find parameters & the scale factor from M .

$$M \equiv K^* [R^* | T^*] = S \hat{M}$$

$$\Rightarrow [K^* R^* | K^* T^*] = S \hat{M}$$

$$K^* R^* = \begin{vmatrix} \alpha_u & S & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{vmatrix} \begin{vmatrix} -r_1^T \\ -r_2^T \\ -r_3^T \end{vmatrix} =$$

$$= \begin{bmatrix} \alpha_u r_1^T + S r_2^T + u_0 r_3^T \\ \alpha_v r_2^T + v_0 r_3^T \\ r_3^T \end{bmatrix}$$

$K^* R^*$

$$\begin{bmatrix} \alpha u r_1^T + S r_2^T + U_0 r_3^T \\ \alpha v r_1^T + V_0 r_3^T \\ r_3^T \end{bmatrix} = K^* T^* = S M = S \begin{bmatrix} -a_1^T \\ -a_2^T \\ -a_3^T \end{bmatrix}$$

unknowns knowns

$$\textcircled{1} \quad \alpha u r_1^T + S r_2^T + U_0 r_3^T = S a_1^T$$

Knowns:

$$\textcircled{2} \quad \alpha v r_1^T + V_0 r_3^T = S a_2^T$$

 a_1, a_2, a_3, b

$$\textcircled{3} \quad r_3^T = S a_3^T$$

$$\textcircled{4} \quad K^* T^* = S b$$

Unknowns:

 r_1, r_2, r_3, T^*, K^* $S, \alpha u, \alpha v, U_0, V_0, S$

- $\textcircled{1}$ * trick to solve this is to ~~remember~~ remember that rows of R^* are orthogonal.

To extract parameters using orthogonality of r_1, r_2, r_3 .

orthogonal

property

$$r_1 \cdot r_2 = 0 \quad r_2 \cdot r_3 = 0 \quad r_1 \cdot r_3 = 0$$

$$r_1 \times r_2 = r_3 \quad r_2 \times r_3 = r_1 \quad r_3 \times r_1 = r_2$$

Find unknown scale?

$$r_3^T = S a_3^T$$

$$\Rightarrow |r_3^T| = |S| |a_3^T| \Rightarrow |S| = 1$$

unknown $|a_3^T|$ ← known
 $|r_3^T| = 1$ since, R matrix is orthonormal. Vectors in R^* have unit length.

so now we know only magnitude but not scale.

(2) Find U_0 :

$$S_{a_1}^T \cdot S_{a_3}^T = (du_{12}^T + Su_{23}^T + U_0 u_{31}^T) \cdot r_3^T \\ = 0 + 0 + U_0 \cdot 1 = U_0$$

$$U_0 = |S|^2 a_1 \cdot a_3$$

we knew $|S|$ from (1)

similarly:

$$V_0 = |S|^2 a_2 \cdot a_3$$

(3) finding α_v :

$$S_{a_2} \cdot S_{a_2} = (\alpha_v r_2^T + \lambda_0 r_3^T) \cdot (\alpha_v r_2^T + V_0 r_3^T) \\ = \alpha_v^2 + V_0^2$$

$$\alpha_v = \sqrt{S^2 a_2 \cdot a_2 - V_0^2}$$

(4) Finding S :

$$\begin{aligned} S \mathbf{a}_1 \times S \mathbf{a}_3 &= (\alpha_1 \mathbf{r}_1^T + S \mathbf{r}_2^T + U \mathbf{r}_3^T) \times \mathbf{r}_3^T \\ &= -\alpha_1 \mathbf{r}_2^T + S \mathbf{r}_1^T \end{aligned}$$

$$S \mathbf{a}_2 \times S \mathbf{a}_3 = \alpha \mathbf{r}_1^T$$

$$(S \mathbf{a}_1 \times S \mathbf{a}_3) \cdot (S \mathbf{a}_2 \times S \mathbf{a}_3) = (-\alpha_1 \mathbf{r}_2^T + S \mathbf{r}_1^T) \cdot \alpha \mathbf{r}_1^T$$

$$= S \alpha v$$

$$S = \frac{1}{\alpha v} S^4 (\cancel{\mathbf{a}_1 \times \mathbf{a}_3}) \cdot (\mathbf{a}_2 \times \mathbf{a}_3)$$

* Finding the unknown sign of S

$$S = \epsilon |S|$$



? ~~+/-~~ + / -

$$\epsilon = \text{sign}(S) = ?$$

$$K^* T^* = S_b = \epsilon |S| b$$

$$[K^* T^*]_+ = \epsilon |S| b_+$$

↑
+ve (Object in front of camera)

$$\Rightarrow e = \text{sign}(b_2)$$

we take sign of b_2
such that $[K^* D^*]$
stay +ve.

* So now we found $u_0, v_0, s, x_u, x_v, s, e$
so now we get K^* intrinsic parameter

Lecture 24 - W12 L2

Page No.		
Date		

We find R^* , T^* using equations we have.

R^* involves equation 1, 2, 3.

T^* involves equation 4.

So we start with eq(4).

* Finding T^* :

We know K^* from recovered intrinsic parameters

$$K^* T^* = [E | s]_b \quad \begin{matrix} \text{known (part of } \hat{M}) \\ \nearrow \\ \text{known} \end{matrix}$$

↑
unknown
↑
known

so,

$$T^* = (K^*)^{-1} [E | s]_b$$

* Summary:

$$① |B| = \frac{1}{|a_3|}$$

$$② U_0 = |B|^2 a_1 \cdot a_3$$

$$③ V_0 = |B|^2 a_2 \cdot a_3$$

$$④ \alpha_v = \sqrt{|B|^2 a_2 \cdot a_3 - V_0^2}$$

$$⑤ S = |B|^4 \alpha_v (a_1 \times a_3) \cdot (a_2 \times a_3)$$

$$⑥ \alpha_u = \sqrt{|B|^2 a_1 \cdot a_1 - S^2 - U_0^2}$$

$$⑦ K^* = \begin{bmatrix} \alpha_u & S & U_0 \\ \alpha_v & V_0 & \\ & & 1 \end{bmatrix}$$

$$⑧ \epsilon = \text{sgn}(b_3)$$

$$⑨ T^* = \epsilon |B| (K^*)^{-1} b$$

$$⑩ r_3 = \epsilon |B| a_3$$

$$⑪ r_1 = |B|^2 / \alpha_v a_2 \times a_3$$

$$⑫ r_2 = r_3 \times r_1$$

$$⑬ R^* = [r_1^T \ r_2^T \ r_3^T]^T$$

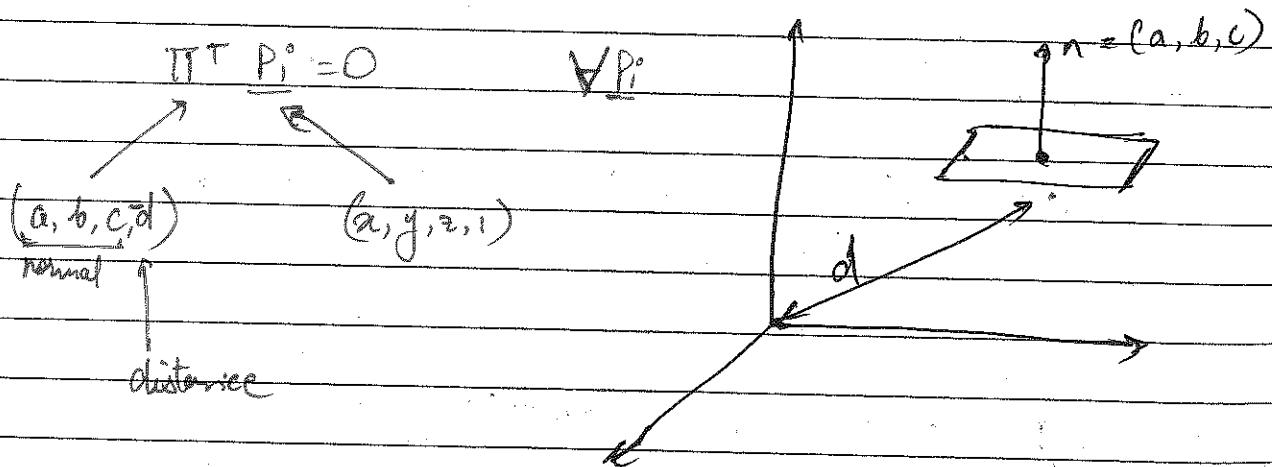
* Degenerate Configurations

If we take only picture of 1 plane of ~~the~~ calibrators, our M fails.

To recover M we need to solve:

$$\begin{bmatrix} \underline{P_i^T} & 0 & -x_i \underline{P_i^T} \\ 0 & \underline{P_i^T} & -y_i \underline{P_i^T} \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ m_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$$

When $\underline{P_i}$ are all on the same plane π , we have



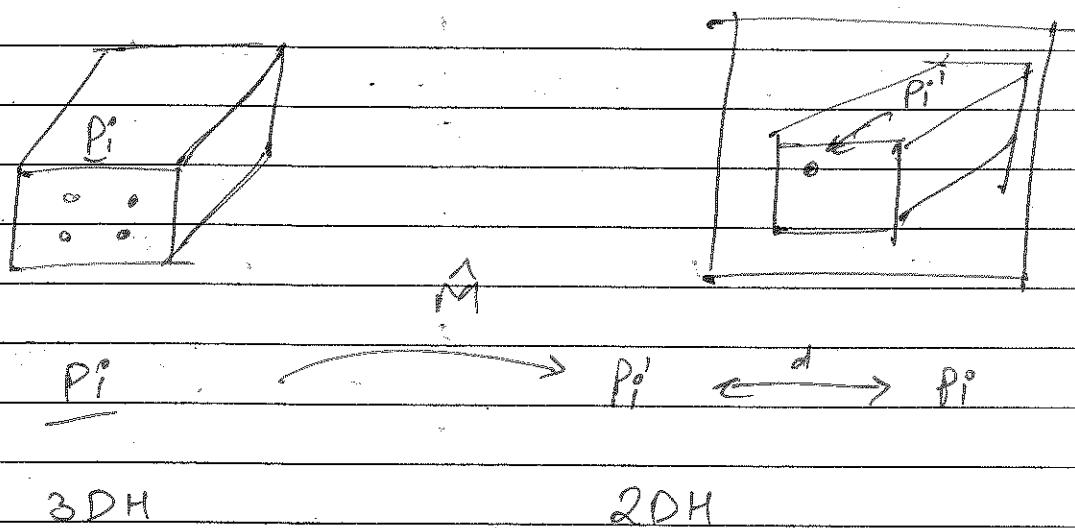
- When $\underline{P_i}$ are on same plane there are many solutions of the form:

$$\begin{bmatrix} m_1 \\ m_2 \\ m_3 \\ (2 \times 1) \end{bmatrix} = \alpha \begin{bmatrix} \pi \\ 0 \\ 0 \end{bmatrix} + \beta \begin{bmatrix} 0 \\ \pi \\ 0 \end{bmatrix} + \gamma \begin{bmatrix} 0 \\ 0 \\ \pi \end{bmatrix}$$

there are other degenerate configurations, but they are less likely to happen.

* Assessing the quality of fit:

We can measure quality by projecting world points onto image equivalent point with our solution \hat{M} . Then look at the distance between actual & found point.



Given $\{P_i^*\}_{i=1}^m \longleftrightarrow \{P_i\}_{i=1}^m$ and estimated \hat{M}

$$\hat{M} = \begin{bmatrix} m_1^T \\ m_2^T \\ m_3^T \end{bmatrix}$$

access the quality of fit using:

2D (pixels)

$$E = \frac{1}{m} \sum_i \left(\left\| x_i - \frac{m_1^T P_i}{m_3^T P_i} \right\|^2 + \left\| y_i - \frac{m_2^T P_i}{m_3^T P_i} \right\|^2 \right)$$

i.e. distance b/w known & predicted parameters.

if we take \sqrt{F} , we get error in pixels.

* Recovering R & T .

remember R^* & T^* are rotation of world w.r.t camera.

but we usually need R & T ,

$$R^* = R^T \Rightarrow R = (R^*)^T$$

$$T^* = -R^T T \Rightarrow T = -R T^*$$

$$\Rightarrow T = -(R^*)^T T^*$$

* Planar Calibration :-

$\Rightarrow \hat{H}_1, \hat{H}_2, \hat{H}_3, \hat{H}_4, \dots$

* Approach :-

- ① estimate a 2D Homography (projective map) between calibration plane and image (for several images)
- ② Estimate intrinsic parameters
- ③ Compute extrinsic parameters for view of interest. $\rightarrow R, T$

* 2D projective map :-

$$P_i' = M P_i$$

P_i' $=$ $M P_i$
 $1 \quad \quad \quad 1 \quad \quad \quad \nearrow$
 $2DH$ 3×3 $2DH$
 (3×1) (3×1)

To get a 2D projective map we assume:

$$\{ P_i^3 z = 0 \Leftrightarrow P_i = (x_i, y_i, 0)$$

We make assumption that z coordinates are 0.

$$(3 \times 1) \quad (3 \times 3) \quad (3 \times 3) \quad (3 \times 1) \quad (4 \times 1)$$

$$\underline{P_i} = K^* [R^* | T^*] \underline{R_i}$$

2DH 3DH
(3x4)

$$= K^* \begin{bmatrix} r_1 & r_2 & r_3 & T^* \end{bmatrix} \begin{bmatrix} X_i \\ Y_i \\ 0 \\ 1 \end{bmatrix}$$

$$= K^* (X_i r_1 + Y_i r_2 + T^*)$$

$$\underline{P_i} = K^* \begin{bmatrix} r_1 & r_2 & T^* \end{bmatrix} \begin{bmatrix} X_i \\ Y_i \\ 1 \end{bmatrix}$$

3x3 2DH

So,

$$\text{given : } \{\underline{P_i}\}_{i=1}^m \leftrightarrow \{\underline{P_i}^*\}_{i=1}^m$$

$$\underline{P_i} = \begin{bmatrix} X \\ Y \\ 0 \\ 1 \end{bmatrix}$$

$$\text{Define : } \underline{P_i}^* = \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} u_i \\ v_i \\ w_i \end{bmatrix} = K^* \begin{bmatrix} r_1 & r_2 & T^* \end{bmatrix} \underline{P_i}^*$$

$$\left\{ \underline{P_i} \right\}_{i=1}^m \xleftarrow[2DH]{} \left\{ \underline{\underline{P_i}} \right\}_{i=1}^m \xrightarrow[3DH]{} \left\{ \underline{P_i^*} \right\}_{i=1}^m$$

$$P_i = \begin{bmatrix} X_i \\ Y_i \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} U_i \\ V_i \\ W_i \end{bmatrix} = H P_i^* = \begin{bmatrix} -h_1^T \\ -h_2^T \\ -h_3^T \end{bmatrix} P_i^*$$

$$X_i = \frac{U_i}{W_i} = \frac{h_1^T P_i^*}{h_3^T P_i^*} \Rightarrow h_1^T P_i^* - X_i h_3^T P_i^* = 0$$

$$Y_i = \frac{V_i}{W_i} = \frac{h_2^T P_i^*}{h_3^T P_i^*} \Rightarrow h_2^T P_i^* - Y_i h_3^T P_i^* = 0$$

So, 2 equations for each point pair.

9 unknowns : h_1, h_2, h_3 (8 independent)

\Rightarrow Need a minimum of 4 point pairs.

$$\begin{array}{c|ccc|c}
& P_i^{*T} & 0^T & -X_i P_i^{*T} & 0 \\
& 0 & P_i^{*T} & -Y_i P_i^{*T} & 0 \\
\hline
2m & \vdots & \vdots & \vdots & \vdots \\
& \vdots & \vdots & \vdots & \vdots \\
& P_m^{*T} & 0^T & -X_m P_m^{*T} & 0 \\
& 0 & P_m^{*T} & -Y_m P_m^{*T} & 0 \\
\hline
& 2m \times 9 & & q \times 1 & 2m \times 1
\end{array}$$

$Ax = 0$ form

- Solve the homography equation ($Ax=0$) using SVD for several planar views (each with its own coordinate system)
 $\Rightarrow \hat{H}_1, \hat{H}_2, \hat{H}_3, \dots$
 defined
- each homography matrix is ~~defined~~ up to scale!

$$H \equiv K^* [r_1 \ r_2 \ T^*] = \alpha \hat{H} \quad \text{unknown scale}$$

like $M = 3\hat{M}$

$$\begin{matrix} & | & | & | \\ & \hat{h}_1 & \hat{h}_2 & \hat{h}_3 \end{matrix}$$

Let, $\hat{H} = [\hat{h}_1 \ \hat{h}_2 \ \hat{h}_3]$

note columns of H instead
of rows of M before

$$H = K^* [r_1 \ r_2 \ T^*] = \alpha \hat{H} = [\hat{h}_1 \ \hat{h}_2 \ \hat{h}_3]$$

$$\Rightarrow [r_1 \ r_2 \ T^*] = \alpha K^{*-1} [\hat{h}_1 \ \hat{h}_2 \ \hat{h}_3]$$

$$\begin{aligned} \Rightarrow r_1 &= \alpha K^{*-1} \hat{h}_1 \\ r_2 &= \alpha K^{*-1} \hat{h}_2 \\ T^* &= \alpha K^{*-1} \hat{h}_3 \end{aligned}$$

$$\begin{aligned} r_1 \cdot r_2 &= 0 \Rightarrow (\alpha K^{*-1} \hat{h}_1)^T (\alpha K^{*-1} \hat{h}_2) = 0 \\ &\Rightarrow \hat{h}_1^T K^{*-T} K^{*-1} \hat{h}_2 = 0 \end{aligned}$$

$$r_1 \cdot r_1 = r_2 \cdot r_2 = 1$$

$$\Rightarrow \hat{h}_1^T K^{*-T} K^{*-1} \hat{h}_1 = \hat{h}_2^T K^{*-T} K^{*-1} \hat{h}_2$$

$$\begin{cases} \hat{h}_1^T S \hat{h}_2 = 0 \\ \hat{h}_1^T S \hat{h}_1 = \hat{h}_2^T S \hat{h}_2 \end{cases} \quad S \equiv K^{*-T} K^{*-1}$$

2 equations for each estimated homography \hat{H}

$$\begin{cases} \hat{h}_1^T S \hat{h}_2 = 0 \\ \hat{h}_1^T S \hat{h}_1 = \hat{h}_2^T S \hat{h}_2 \end{cases}$$

known : \hat{h}_1, \hat{h}_2

————— 5 unknowns ———

Unknown : $S \rightarrow k^* \rightarrow x_u, x_v, u_0, v_0, s$

→ need at least 3 planar views to estimate k^*
 $m=6$

* Algorithm:

- ① write system of equations
- ② solve for elements of S
- ③ extract intrinsic parameters from S
 (x_u, x_v, u_0, v_0, s)

[Refer notes for summary.]

$$H = [\hat{h}_1 \ \hat{h}_2 \ \hat{h}_3] = \begin{bmatrix} h_{11} & h_{21} & h_{31} \\ h_{12} & h_{22} & h_{32} \\ h_{13} & h_{23} & h_{33} \end{bmatrix}$$

$$V_{ij} = [h_{i1}h_{j1}, h_{i1}h_{j2} + h_{i2}h_{j1}, h_{i2}h_{j2}, h_{i2}h_{j3} + h_{i3}h_{j2}, h_{i3}h_{j2} + h_{i2}h_{j3}, h_{i3}h_{j3}]^T$$

$$\begin{bmatrix} V_{12}^T \\ V_{11}^T - V_{22}^T \\ \vdots \\ \{ \end{bmatrix} \begin{bmatrix} S_{11} & 0 & 0 \\ S_{12} & 1 & 0 \\ S_{22} & = & i \\ S_{13} & 3 & 1 \\ S_{23} & 4 & 1 \\ S_{33} & 5 & \} \end{bmatrix}$$

\longleftrightarrow solve using
SYD

$$C_1 = (S_{12}S_{13} - S_{11}S_{13})$$

$$C_2 = (S_{11}S_{22} - S_{12}^2)$$

$$V_0 = C_1/C_2$$

$$\lambda = S_{33} - (S_{13}^2 + V_0 C_1) / S_{11}$$

$$\alpha_u = \pm \sqrt{\lambda/S_{11}}$$

$$\alpha_v = \sqrt{\lambda S_{11}/C_2}$$

$$S = -S_{12} \alpha_u^2 \alpha_v / \lambda$$

$$U_0 = S V_0 / \alpha_u - S_{13} \alpha_u^2 / \lambda$$

[refer notes]

★ Direct Methods

Given $\{p_i\}_{i=1}^m \leftrightarrow \{q_i\}_{i=1}^m$

The calibration parameters (k^*, r^*, T^*) can be used to project $\{q_i\}_{i=1}^m$

$$q_i = k^* [r^* | T^*] p_i$$

We start with 'guess' and project q_i , then we find euclidean distance b/w prediction & action to get loss. Then we do gradient descent to get parameters.

$$\text{so, } L = \sum (q_i - p_i)^2$$

- objective function

$$E(k^*, r^*, T^*) = \sum_{i=1}^m \left\| \begin{pmatrix} (q_i)_x \\ (q_i)_y \\ (q_i)_z \end{pmatrix} - \begin{pmatrix} (p_i)_x \\ (p_i)_y \\ (p_i)_z \end{pmatrix} \right\|^2$$

- minimize

$$k^*, r^*, T^* = \underset{k, r, T}{\operatorname{arg\min}} E(k, r, T) \quad \text{do gradient descent}$$

non-linear equation \Rightarrow Iterative solution (eg Tsai)

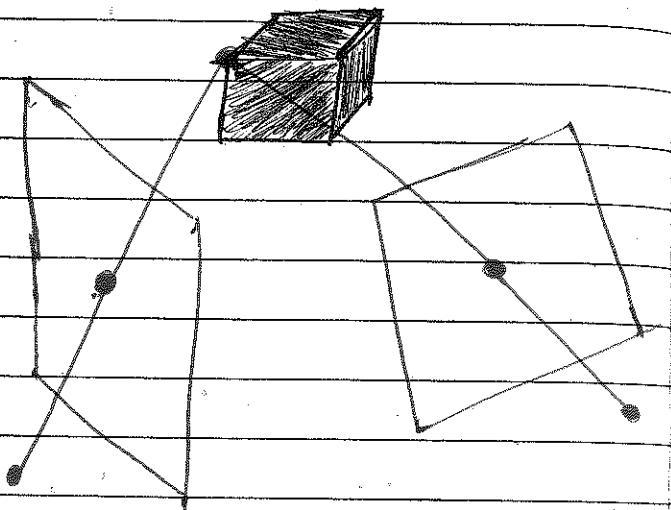


Multiple View Geometry

* 2 view geometry (stereo) :

- use projection of the same 3D point in multiple views to recover structure (inverse problem)

- problems :
 - correspondence
 - Reconstruction



Correspondence :

which point correspond to which point in 2 images.

- approaches :-

- sparse : feature-based (match points - can handle large)
- dense : correlation, SSD (produce more points)

- cases when correspondence is not possible:

- object points not visible in both views (occlusion)
- ambiguous points (multiple candidates)
- uniform regions (inside points cannot be distinguished)

Lecture 25 - W13 L1

Page No.	
Date	

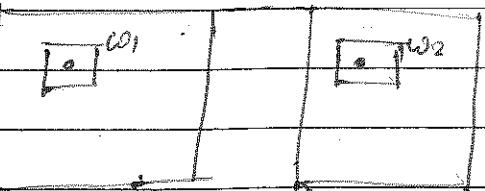
* Sparse correspondence

- find feature points (e.g. corners)
- find local characterization (e.g. SIFT, HOG)
- find corresponding points having similar features
- Greedy vs. optimal assignment
- constraints for reducing # of candidates (e.g. search around current location)

* Dense correspondence

- Instead of feature points capture all patches
- Instead of distance between feature vectors measure correlation or SSD
- Apply regularization to reduce errors and find correspondence in difficult areas (e.g. uniform)

$$\text{Good Correlations: } \Psi(w_1, w_2) = \sum_i w_1(x_i, y_i) w_2(x_i, y_i)$$



$$\text{Good SSD: } \Psi(w_1, w_2) = \sum_i (w_1(x_i, y_i) - w_2(x_i, y_i))^2$$

- removing dependency on absolute value:

[Fig] → values in windows could have similar values but some other window might have better correlation due to bigger value.

Solution: normalize window values (subtract mean and divide by standard deviation)

zero mean normalize cross-correlation (ZNCC)

$$\Psi(\omega_1, \omega_2) = \sum_i \frac{(\omega_1(x_i, y_i) - M\omega_1)(\omega_2(x_i, y_i) - M\omega_2)}{\sigma_{\omega_1} \cdot \sigma_{\omega_2}}$$

zero mean normalized SSD (ZNSSD)

$$\Psi(\omega_1, \omega_2) = \sum_i \left(\frac{\omega_1(x_i, y_i) - M\omega_1}{\sigma_{\omega_1}} - \frac{\omega_2(x_i, y_i) - M\omega_2}{\sigma_{\omega_2}} \right)^2$$

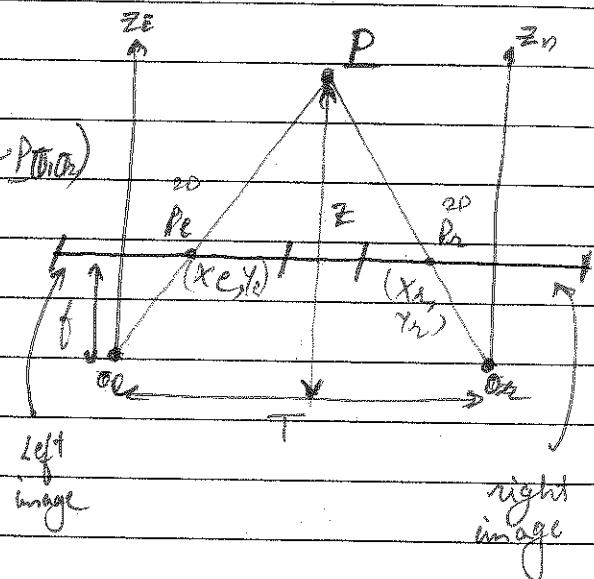
* Axis-aligned stereo :

- Using similar triangles: ($P_1 P_2 \sim P_1 O_1 O_2$)

$$\frac{T}{z} = \frac{T - (-x_2 + x_1)}{z - f}$$

$$= \frac{T - (x_2 - x_1)}{z - f}$$

$$= \frac{T - d}{z - f} \leftarrow \text{disparity}$$

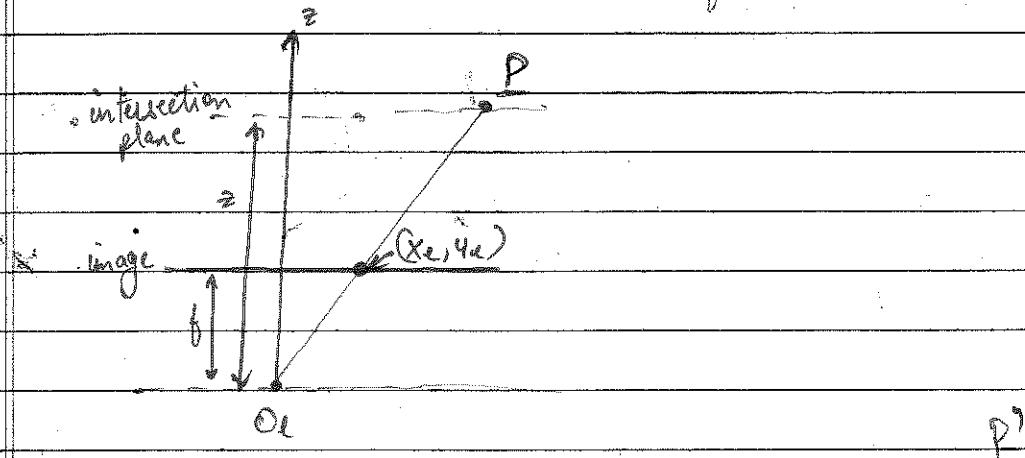


i.e. the bigger the disparity, smaller the depth.

$$\frac{T}{z} = \frac{T - d}{z - f} \Rightarrow \boxed{z = f \frac{T}{d}}$$

f, T is always same for all points.
d will vary.

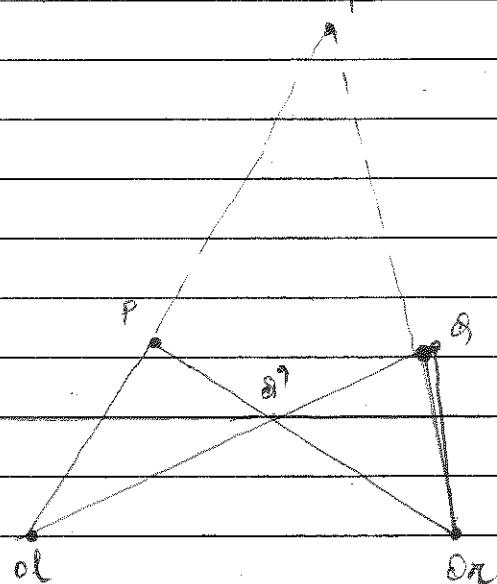
- To reconstruct 3D point, send a ray from O_l through (x_e, y_e) and find intersection with a parallel plane at a distance of z (where z is computed as before)



Ambiguity problem:

correct : PQ

incorrect : $P'Q'$



★ General stereo:

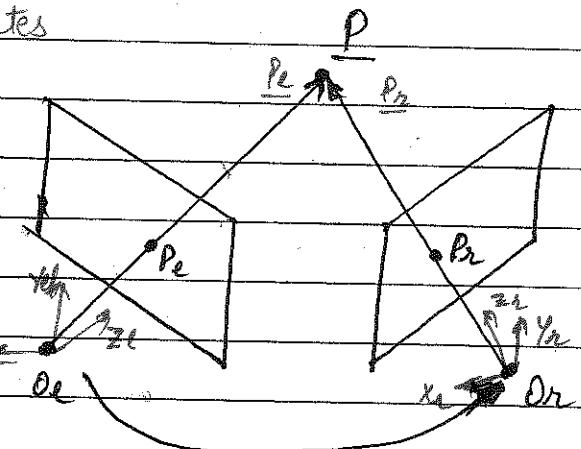
P = point in world coordinates

P_L = point in left coordinates

P_R = point in right coordinates

P_L = image point in left coordinates

P_R = image point in right coordinates



first rotate, $\rightarrow R_L T_L$ ← stereo parameters
then translate

- stereo calibration -

• getting relative rotation & translation (R, T)

calibrate $\rightarrow \begin{cases} R_L, T_L & \text{rotation / translation to} \\ R_R, T_R & \text{of left/right camera w.r.t world} \end{cases}$

$$M_{\text{left} \leftarrow \text{right}} = M_{\text{left} \leftarrow \text{world}} M_{\text{world} \leftarrow \text{right}}$$

$$R_L^T T_L^{-1} (R_R^T T_R^{-1})^{-1} = T_R R_L$$

$$= R_L^T T_L^{-1} T_R^{-1} R_R = TR$$

$$\text{so, } M_{\text{left} \leftarrow \text{right}} = R_L^T T_L^{-1} T_R^{-1} R_R$$

$$= \begin{bmatrix} R_L^T & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} I & -T_L \\ 0 & 1 \end{bmatrix} \begin{bmatrix} I & T_R \\ 0 & 1 \end{bmatrix} \begin{bmatrix} R_R & 0 \\ 0 & 1 \end{bmatrix}$$

rotation ✓ translation

$$= \begin{bmatrix} R_L^T & -R_L^T T_L \\ 0 & 1 \end{bmatrix} \begin{bmatrix} R_R & T_R \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} R_L^T R_R \\ 0 \end{bmatrix} \begin{bmatrix} R_L^T (T_R - T_L) \\ 1 \end{bmatrix}$$

$$\Rightarrow R = R_e^T R_s$$

$$T = R_e^T (T_e - T_s)$$

* Epipolar geometry

- the points O_e, O_s, P define epipolar plane.

- the intersection of the epipolar plane with the images define epipolar lines

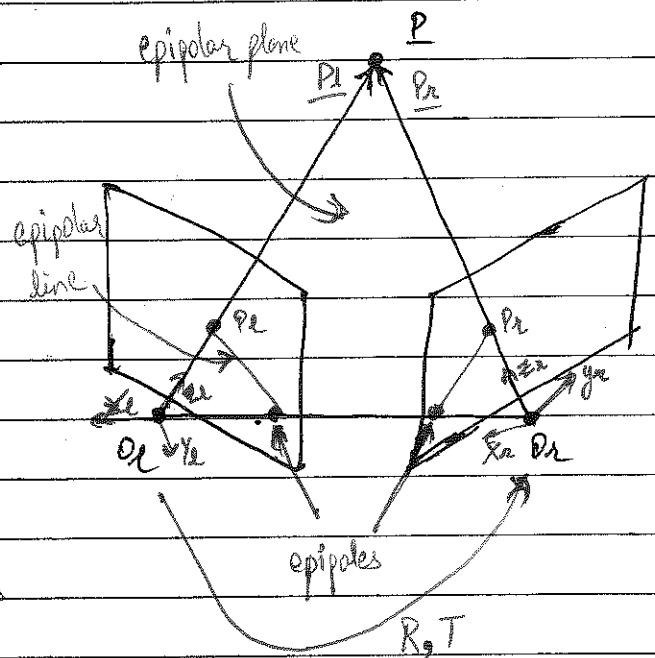
- The intersection of the baseline with the images defines the epipoles.

- properties:

- every epipolar line must pass through the epipole
- epipoles may be inside or outside the image
- epipoles may be at infinity (when the image plane become parallel).

* Epipolar constraint

- Given a point P_e in the left image P_e, O_e, O_s define an epipolar plane.
- The intersection of this epipolar plane with the right image defines an epipolar line on it.
- the point P_r corresponding to P_e on the left image must be on this epipolar line.



* Essential matrix :

- \underline{P}_R and \underline{P}_E are expressed in different coordinate systems. we would like to move to the left coordinate system.
- given that the right view is rotated and translated by R, T w.r.t left , it is easy to move from left to right

$$\underline{P}_L = R^T (\underline{P}_E - T)$$

$$\Rightarrow \underline{P}_E = R \underline{P}_L + T \quad \leftarrow \text{convert } \underline{P}_L \text{ in left coordinates}$$

Lecture 26 - W13 L2

Page No.	
Date	

- moving from camera to image coordinates:

$$\underline{p}_e = \frac{f_e}{z_e} \cdot \underline{p}_e$$

$$\underline{p}_r = \frac{f_r}{z_r} \underline{p}_r$$

- Cross product as matrix multiplication:

matrix vector

$$A \times B = [A] \times B \equiv \begin{bmatrix} 0 & -a_z & a_y \\ a_z & 0 & -a_x \\ -a_y & a_x & 0 \end{bmatrix} B$$

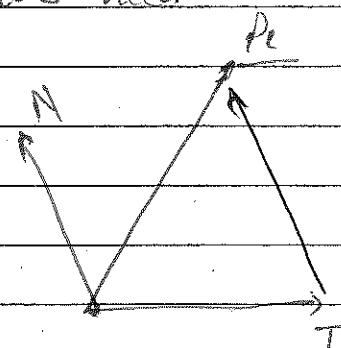
$A = \begin{bmatrix} a_x \\ a_y \\ a_z \end{bmatrix}$

rank 2 skew-symmetric matrix ($A^T = -A$)

- * Essential matrix:

- To specify an epipolar plane we need to find a vector normal to it.

$$\underline{N} = T \times \underline{p}_e$$



- The vector \underline{p}_e is on the epipolar plane:

$$\underline{p}_e \cdot \underline{N} = 0$$

- Moving \underline{p}_e to left coordinates

$$(R_L \underline{p}_e + \underline{T}) \cdot \underline{N} = 0$$

$$(R\mathbf{P}_3 + \mathbf{T}) \cdot \mathbf{N} = 0 \Rightarrow (R\mathbf{P}_3 + \mathbf{T}) \cdot (\mathbf{T} \times \mathbf{P}_e) = 0$$

$$\Rightarrow R\mathbf{P}_3 \cdot (\mathbf{T} \times \mathbf{P}_e) = 0$$

$$\Rightarrow (R\mathbf{P}_3)^T [\mathbf{T} \times \mathbf{P}_e] = 0$$

getting rid of
dot & cross

$$\Rightarrow \mathbf{P}_3^T R^T [\mathbf{T} \times \mathbf{P}_e] = 0$$

$\equiv E$
 (3×3)

$$\Rightarrow \boxed{\mathbf{P}_3^T E \mathbf{P}_e = 0} \leftarrow \begin{array}{l} \text{Epipolar constraint} \\ \text{equation} \end{array}$$

Essential matrix

* summarize:

- the essential matrix

$$E \equiv R^T [\mathbf{T}]_x$$

rank 2 (because of $[\mathbf{T}]_x$)

- Epipolar constraint:

$$\mathbf{P}_3^T E \mathbf{P}_e = 0$$

- Epipolar constraint in image coordinates:

$$\mathbf{P}_3^T E \mathbf{p}_e = 0$$

* Epipolar Constraint

- right epipolar line : (given left image - P_L)

$$l = E P_L$$

the corresponding P_R should be on this line.

- Explanation :

$$P_R^T E P_L = 0 \Rightarrow P_R^T l = 0 \Rightarrow P_R \cdot l = 0 \Rightarrow P_R \text{ is on the line}$$

line equation $p^T l = 0$; $l = (a, b, c, 1)$ ~~distance from origin~~

- left epipolar line : (given right image - P_R)

$$l = E^T P_R$$

the corresponding point P_L should be on this line.

Explanation :

$$P_L^T E P_R = 0 \Rightarrow (P_R^T E P_R)^T = 0 \Rightarrow P_R^T E^T P_L = 0$$

$$\Rightarrow P_R^T l = 0$$

$$\Rightarrow l = E^T P_R$$

But these only take extrinsic R & T into consideration, but we also want intrinsic to be included and so we use fundamental matrix

* Fundamental matrix (related to essential matrix)

• So far the epipolar constraint was expressed in camera coordinates - we want to move to image coordinates

$$\begin{matrix} \bar{p}_e = K_e^* p_e \\ \downarrow \quad \uparrow \\ \text{image} \quad \text{camera} \\ \text{coordinates} \quad \text{coordinates} \end{matrix}$$

$$K_e^* = \begin{bmatrix} d_u & s & u_0 \\ 0 & d_v & v_0 \\ 0 & 0 & 1 \end{bmatrix}$$

$u_0, v_0 \rightarrow$ optical center

$d_u, d_v \rightarrow$ focal length in pixels

• likewise:

$$p_r = K_r^{*-T} \bar{p}_r$$

so epipolar constraint:

$$p_r^T E p_e = 0$$

$$\Leftrightarrow (K_r^{*-T} \bar{p}_r)^T E (K_e^* \bar{p}_e) = 0$$

$$\Rightarrow \bar{p}_r^T (K_r^{*-T} E K_e^*) \bar{p}_e = 0$$

$\equiv F$ (fundamental matrix)

• The fundamental matrix takes into account both internal and external parameters, whereas the essential matrix takes into account only external parameters.

F is rank 2 because of E

* Summary :

- In camera coordinates

$$P_r^T F P_e = 0$$

- In image coordinate

$$\bar{P}_r^T F \bar{P}_e = 0$$

$$\bullet F = R^T [I_3]$$

3×3

$$\bullet F = K_r^{*-T} E K_r^{*-1}$$

3×3

both E & F are Rank 2

• Given P_e , right epipolar line is : $\ell = E P_e$

• Given P_r , left epipolar line is : $\ell = F^T P_r$

* Weak Calibration:

• Full calibration : $R, T, K_r^*, K_r^{*-1} \Rightarrow E, F$

• weak calibration : E, F

• Problem statement :

given $\{P_i\}_{i=1}^m$ $\rightarrow \{P'_i\}_{i=1}^m$ $m \geq 8$
 (left point) — (right point)

find F (3x3 matrix)

(8 unknown elements of F)

* Eight point algorithm :

• One equation $\&$ for each point pair :

$$① P_i^T F P'_i = 0$$

8 unknowns \Rightarrow need at least 8 point pairs

$$② [x_i, y_i, 1] \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix} \begin{bmatrix} x'_i \\ y'_i \\ 1 \end{bmatrix} = 0$$

$$③ x_i x'_i f_{11} + x_i y'_i f_{12} + \dots + f_{33} = 0$$

$$\begin{array}{c|c|c|c}
\text{mx9} & \text{9x1} & \text{mx1} \\
\hline
\begin{bmatrix} x_1 x'_1, x_1 y'_1, \dots, 1 \end{bmatrix} & \begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ \vdots \\ f_{33} \end{bmatrix} & \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \\
\hline
\begin{bmatrix} x_m x'_m, x_m y'_m, \dots, 1 \end{bmatrix} & \begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ \vdots \\ f_{33} \end{bmatrix} & \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}
\end{array}$$

- this is in an $Ax=0 \Rightarrow$ solution is right zero eigenvalue of A vector

($SVD \Rightarrow A = UDV^T$ then take last column of V)

- To enforce rank 2 for the estimated F :

$$SVD \rightarrow F = UDV^T$$

$$\begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix}$$

$D' = D$ with smallest singular value set to 0
($D'[3,3] = 0$)

$$F' = U D' V^T \rightarrow \text{rank 2 matrix}$$

- Evaluate how good F is:

$$\{P_i\} \leftrightarrow \{P_i'\} \quad \sum (P_i^T F P_i') \approx 0$$

sometimes it would fit, other times it wouldn't. Solution: Normalization.

good fit

* Normalization:

- For a stable solution in SVD algorithm, normalize the point sets in $\{P_i\}_1^m \rightarrow \{P_i'\}_1^m$

$$q_i = \frac{p_i - M_p}{\sigma_p}$$

$$q_i' = \frac{p_i' - M_{p'}}$$

M = mean
 σ = stda

- But the solution will be good for q points, so we need to convert that to q' points.

- In matrix form:

$$q_i = \begin{bmatrix} 1/\sigma_p \\ Y_p \end{bmatrix} \begin{bmatrix} 1 & 0 & -M_{p_x} \\ 0 & 1 & -M_{p_y} \\ 0 & 0 & 1 \end{bmatrix} p_i \\ \equiv \underline{M_p}$$

$$q_i = \underline{M_p} p_i$$

$$q_i' = \underline{M_{p'}} p_i'$$

- using $\{q_i\}^m \leftrightarrow \{q_i'\}^m$ and BP algo find F' :

$$q_i^T F' q_i' = 0$$

$$(\underline{M_p} p_i)^T F' (\underline{M_{p'}} p_i') = 0$$

$$p_i^T \underline{\underline{M_p^T F' M_{p'}}} p_i' =$$

$\equiv F$

$$F = \underline{\underline{M_p^T F' M_{p'}}}$$

* Summarize :

given $\{P_i\}_1^m \longleftrightarrow \{P'_i\}_1^m$

① normalize points : $q_i = M_p P_i$
 $q'_i = M'_{p'} P'_i$

② use 8 point algo to get F'

③ convert F' to F : $F = M_p^T F' M_{p'}$

↑
weak calibration

27

Lecture 27 - W14L1

Page No.	
Date	

(Contd)

$$P_r^T F P_e = 0$$

- given P_e we have epipolar line $F P_e$ in the right image
- every point on this line satisfies $P^T F P_e = 0$
- specifically: $\mathbf{1}_3^T F P_e = 0$
- since this is for all P_e it must be that:

$$\mathbf{1}_3^T F = 0$$

\curvearrowleft Ax=0 form

- $\mathbf{1}_3^T F = 0 \Rightarrow \frac{F^T \mathbf{1}_3}{3 \times 3 \quad 3 \times 1} \rightarrow \mathbf{1}_3$ is zero eigenvector of F^T

$\Rightarrow \mathbf{1}_3$ is the left null space of F (SVD: $F = UDV^T$)
 $\rightarrow e_1$ is the last column of U)

- for the left epipole: $P_r^T F^T P_e = 0$

$$\Rightarrow P_r^T F^T P_e = 0 \wedge P_r \Rightarrow P_r^T F^T = 0 \Rightarrow F P_e = 0$$

$P_r^T F e_1 = 0 \wedge P_r \Rightarrow F e_1 = 0 \rightarrow e_1$ is zero eigenvector of F

$\Rightarrow e_1$ is right null space of F
 (SVD: $F = UDV^T \rightarrow e_1$ is 0-the last col* of V)

tion

$$\left\{ \vec{p}_i \right\}_{i=1}^m \xleftrightarrow{2D} \left\{ \vec{p}'_i \right\}_{i=1}^m \text{ find } \left\{ \vec{p}''_i \right\}_{i=1}^m \xleftrightarrow{3D}$$

Cases:

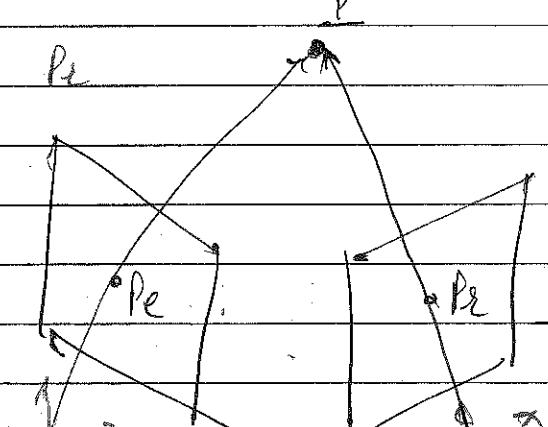
- 1) know intrinsic + extrinsic \rightarrow (Absolute reconstruction)
- 2) know intrinsic \rightarrow (Euclidean reconstruction)
(upto scale)
- 3) None are known \rightarrow ~~Reconstruction~~ Reconstruction up to
(quite useless) unknown 3D projective map.

I) Absolute Reconstruction

* Triangulation Algorithm: (algoned to find reconstruction)

- send ray from o_l through \vec{p}_l
& a ray \vec{l} from o_l through \vec{p}_e

- Reconstruct \vec{P} at rays intersection.



* Step 1: move to camera coord:

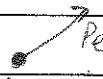
$$\left. \begin{array}{l} \text{image coord} \\ \vec{p}_e = k_e * \vec{p}_e' \end{array} \right\} \quad \left. \begin{array}{l} \text{camera coord.} \\ \vec{p}_e = k_e * \vec{p}_e' \end{array} \right\} \quad \left. \begin{array}{l} \text{camera coord.} \\ \vec{p}_l = k_l * \vec{p}_l' \end{array} \right\}$$

$$\vec{p}_l = k_l * \vec{p}_l'$$

$$\vec{p} = k * \vec{p}'$$

* Step 2 : find intersections of rays:

left ray: aP_e



right ray: bP_r^*

intersection: $aP_e = bP_r^*$

- the ray bP_r^* is right ray in left coordinates

$$P_r^* = M_{\text{left}} \leftarrow \text{right} \quad P_r = R_{\text{left}} + T$$

* - To find intersection solve:

$$aP_e = bP_r^* \rightarrow aR_{\text{left}} + T = bR_{\text{left}} + T \rightarrow \text{solve for}$$

- Problem:

due to inaccuracies rays do not intersect
(no solution)

so we introduce cw

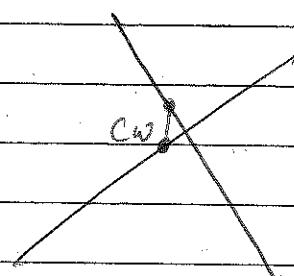
$$\{ aP_e + cw = bP_r^* \quad - \textcircled{1}$$

$$\cdot \quad w = P_e \times P_r^* = P_e \times R P_r \quad - \textcircled{2}$$

→ putting $\textcircled{2}$ in $\textcircled{1}$:

$$\frac{aP_e}{\cancel{P_e}} + \frac{c(P_e \times R P_r)}{\cancel{P_e}} - bR P_r = T$$

↑
unknowns



$$\begin{matrix} & & & & & a \\ | & | & | & | & | & \\ \bar{P}_e & \bar{P}_e \times R_{R_e} & -R_{P_e} & \bar{C} & = T & \rightarrow \text{solve} \\ | & | & | & | & | & \\ & & & & b & \end{matrix}$$

for abc

$$Ax = B \rightarrow x = A^{-1}B$$

*Summary:

1) Given $\bar{P}_e, \bar{P}_d, R, T, k_e^*, k_d^*$ solve for abc

2) Compute \underline{P} :

$$\underline{P} = \frac{1}{2} (a_0 \bar{P}_e + \underbrace{b_0 \bar{P}_d}_{b_0 P_d^*} + T)$$

- or -

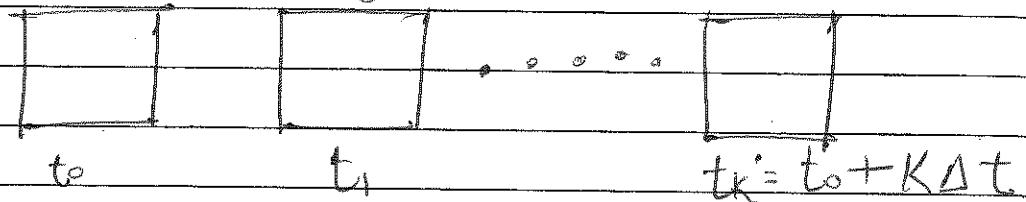
$$\underline{P} = a_0 \bar{P}_e + \frac{1}{2} b_0 \underline{w}$$

3) Convert \underline{P} to world coordinates (from left coord)

$$\underline{P}^{(w)} = \underline{R}_e \underline{P} + \underline{T}_e$$

*

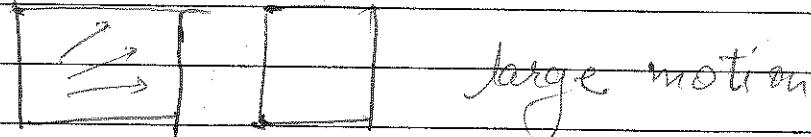
→ sequence of images:



if frame rate is 24 fps.

$$\Delta t = \frac{1}{24} [\text{s}]$$

→ Difference from stereo



→ types of motion:

① Single motion (static camera and moving object or static object and moving camera).

② Multiple motions (multiple objects & camera moving).

→ Problem:

- ① Correspondence (easier in video)
- ② Reconstruction (easier in stereo)
- ③ Motion segmentation

→ Motion of rigid objects

3D motion vectors (world)

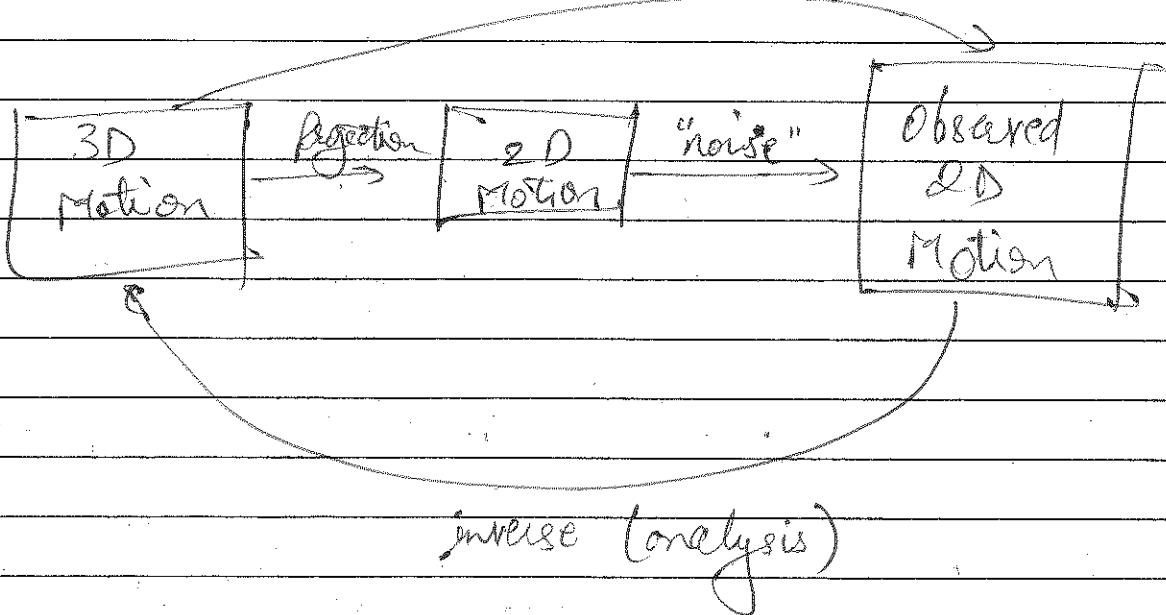
2D projected motion (camera)

2D observed motion (image)



Optical flow

Forward (synthesis)



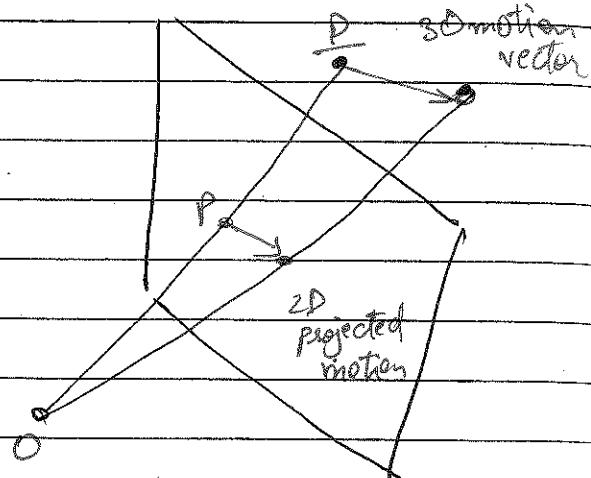
* Basic equations of projected motion

$$\underline{P} = \begin{bmatrix} X \\ Y \\ Z \end{bmatrix}$$

$$\underline{p} = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

$$\underline{P} = f \frac{\underline{p}}{z}$$

\uparrow
z coordinate
of \underline{P}



P, P, z all depend on t

* Motion in image:

$$\frac{d\underline{p}}{dt} = \frac{df}{dt} \left(f \frac{\underline{p}}{z} \right) = f \left(\frac{1}{z} \frac{d\underline{p}}{dt} - \frac{1}{z^2} \frac{dz}{dt} \underline{p} \right)$$

$$v = f \left(\frac{1}{z} v - \frac{1}{z^2} \cancel{v_z} \cancel{p} \right)$$

$$v = \frac{f}{z^2} (z v - v_z p)$$

* In component form:

$$v_x = f \frac{z v_x - v_z x}{z^2}$$

$$v_y = f \frac{z v_y - v_z y}{z^2}$$

$$v_z = f \frac{z v_z - v_z z}{z^2} = 0$$

Projected motion vectors do not have z component.

* Moving to image coordinates from ~~camera~~ world coords.

$$P = f \frac{P}{z} \Rightarrow P = \frac{P}{f} z$$

$$v = f \frac{z v - v_z P}{z^2}$$

$$= f \frac{z v - v_z \frac{P}{f} z}{z^2}$$

$$\boxed{v = f v - \frac{P}{z} v_z}$$

* Decomposing projected motion

Plugging V_x, V_y, V_z, h projected motion equation.

$$\left\{ \begin{array}{l} V_x = V_x(T) + V_x^{(\omega)} \\ V_y = V_y(T) + V_y^{(\omega)} \\ V_z = 0 \end{array} \right. \quad \text{Projected motion}$$

$$\left\{ \begin{array}{l} V_x(T) = Z_z x - Z_x t \\ V_y(T) = Z_z y - Z_y t \end{array} \right. \quad \text{(z) translational component}$$

$$\left\{ \begin{array}{l} V_x^{(\omega)} = -\omega_{yf} + \omega_{zy} y + \frac{\omega_x xy}{t} - \frac{\omega_y x^2}{t} \\ V_y^{(\omega)} = \omega_{xf} - \omega_{zx} x - \frac{\omega_y xy}{t} + \frac{\omega_x y^2}{t} \end{array} \right. \quad \text{rotational motion comp.}$$

* Special Cases

(i) Pure translational motion

$$\omega = 0 \rightarrow v_x = v_x^{(T)}$$

$$v_y = v_y^{(T)}$$

2 Subcases:

(a) $\tau_z \neq 0$ (translation in z)

(b) $\tau_z = 0$ (no translation z)

Subcase 1a: $\tau_z \neq 0$ (and pure translation)

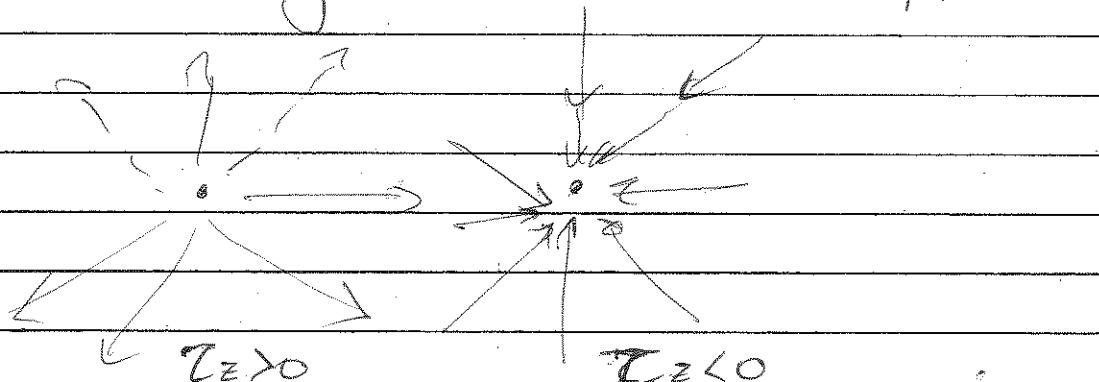
$$v_x = \frac{\tau_z x - \tau_x t}{z} = \frac{\tau_z}{z} \left(x - \frac{\tau_x t}{\tau_z} \right)$$

$$v_y = \frac{\tau_z y - \tau_y t}{z} = \frac{\tau_z}{z} \left(y - \frac{\tau_y t}{\tau_z} \right)$$

$$v_x = \frac{\tau_z}{z} \left(x - x_0 \right)$$

$$v_y = \frac{\tau_z}{z} \left(y - y_0 \right)$$

→ the x_0, y_0 are instantaneous epipole



$$(x_0, y_0) \equiv \left(\frac{z_x f}{z_z}, \frac{z_y f}{z_z} \right)$$

In homogeneous coords:

2DH

$$(z_x f, z_y f, z_z)$$

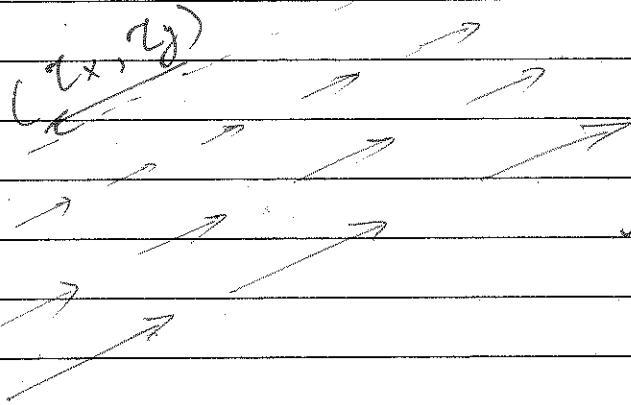
Lecture 28 - WISL

Page No.	
Date	

Subcase 2° ($\tau_0 = 0$) (and pure translation)

$$v_x = \frac{\tau_z x - \tau_x f}{z} = -\frac{\tau_x \cdot f}{z}$$

$$v_y = \frac{\tau_z y - \tau_y f}{z} = -\frac{\tau_y \cdot f}{z}$$



Instantaneous epipole
in homogeneous coordinates
is a point at infinity

$(\tau_x f, \tau_y f, \tau_z)$

$\tau_0 \Rightarrow$ direction

* Motion in instantaneous objects

here $\cdot \bar{V}^{(o)} = \bar{V}^{(i)}$

So,

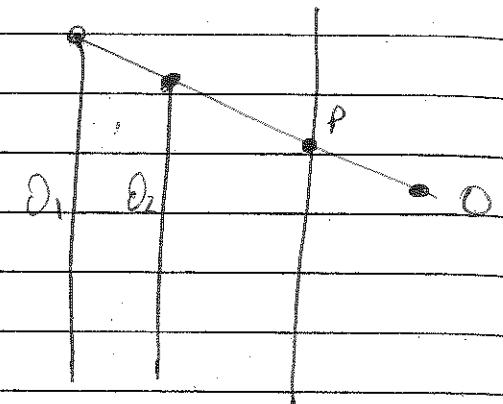
$$\bar{V}_x = \bar{v}^{(T)}$$

$$\bar{V}_y = \bar{v}^{(T)}$$

$$V_x = v^{(T)}$$

$$V_y = v^{(T)}$$

$$So, V_x = \frac{(x - x_0) \tau_z}{z}$$



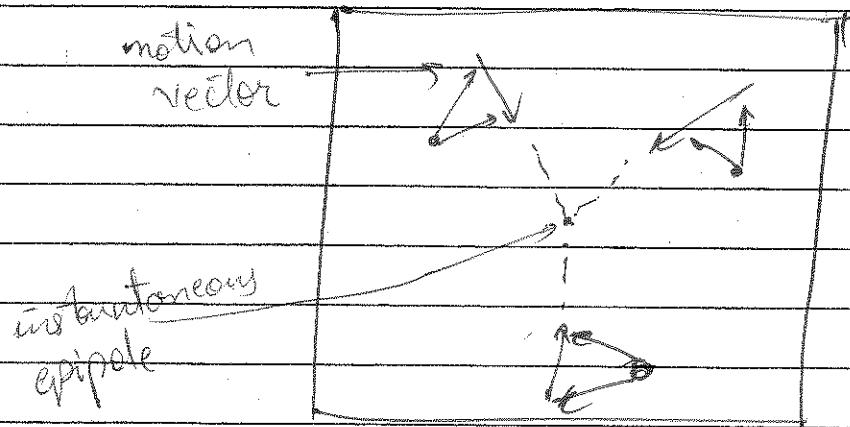
$$\bar{V}_x = \frac{(x - \hat{x}_0) \tau_z}{z}$$

$$V_y = \frac{(x - y_0) \tau_z}{z}$$

$$\bar{V}_y = \frac{(x - \hat{y}_0) \tau_z}{z}$$

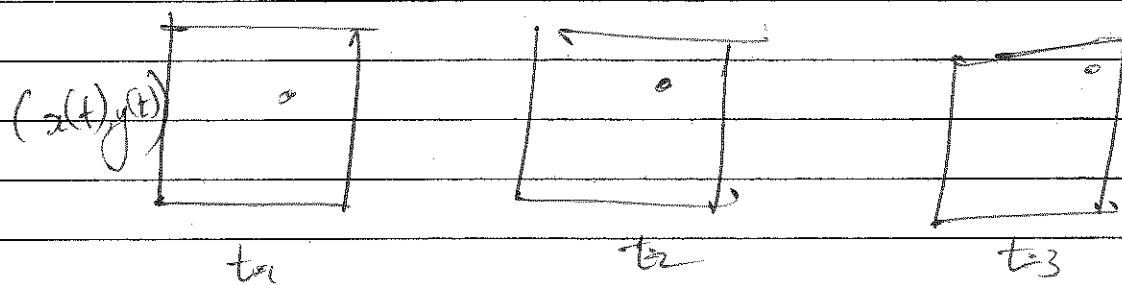
$$\Delta V_x = \frac{(x - x_0) \tau_z}{z} - \frac{(x - x_0) \tau_z}{\bar{z}} = (x - x_0) \tau_z \left(\frac{1}{z} - \frac{1}{\bar{z}} \right)$$

$$\Delta V_y = (y - y_0) \tau_z \left(\frac{1}{z} - \frac{1}{\bar{z}} \right)$$



* Optical flow Estimation

- optical flow = motion observed in images
- assume image brightness per object point is constant.



$$I(x(t), y(t), t) = C \quad \text{constant - image brightness constancy equation.}$$

IBCE

$$\frac{d}{dt} (IBCE) = 0$$

$$\frac{\partial I}{\partial x} \cdot \frac{\partial x}{\partial t} + \frac{\partial I}{\partial y} \cdot \frac{\partial y}{\partial t} + \frac{\partial I}{\partial t} = 0$$

$$\nabla I \cdot v = - I_t$$

↓ (at)
 spatial gradient
 ← optical flow
 constraint / equation
 ↓ time derivative

$$\boxed{\nabla I \cdot V = -I_t}$$

$$V = ?$$

of CE: $\frac{\nabla I}{\|\nabla I\|} \cdot V = -\frac{I_t}{\|I_t\|}$

- Problem?

different motion vectors create same point.

- It is only possible to extract the component of the motion vector parallel to the image gradient

$$V_I = -\frac{I_t}{\|\nabla I\|}$$

* Second Order OF estimation

Assume constant spatial gradient (in addition to constant intensity):

$$\left\{ \begin{array}{l} IBCE = C_1 \\ \nabla (IBCE) = C_2 \end{array} \right.$$

$$\frac{d}{dt} (\nabla I(x(t), y(t), t)) = 0$$

$$\frac{d}{dt} \begin{bmatrix} \frac{d}{dx} IBCE \\ \frac{d}{dy} IBCE \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\frac{d}{dt} \begin{bmatrix} I_x \\ I_y \end{bmatrix} = \begin{bmatrix} I_{xx}x_t + I_{xy}y_t + I_{xt} \\ I_{yx}x_t + I_{yy}y_t + I_{yt} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

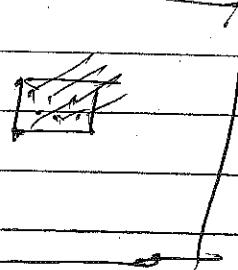
$$\begin{bmatrix} I_{xx} & I_{xy} \\ I_{yx} & I_{yy} \end{bmatrix} \begin{bmatrix} v \\ v \end{bmatrix} = \begin{bmatrix} -I_{xt} \\ -I_{yt} \end{bmatrix} \quad \begin{matrix} \text{solve for} \\ v \\ \text{(noisy estimate)} \end{matrix}$$

$$v = A^{-1}B$$

* Block based of estimate
assume motion θ in small patch is constant

$$E(v) = \sum_{(x,y) \in \text{patch}} C (\nabla I(x,y)^T v + I_0(x,y))^2$$

$\partial E / \partial v = 0$



$$v^* = \underset{v}{\operatorname{argmin}} E(v)$$

$$\nabla E(v) = 0$$

solving,

$$E(v) = \sum (I_x \cdot x_t + I_y \cdot y_t + I_t)^2$$

$$\frac{dE}{dx_t} = 0 = 2 \sum (I_x \cdot x_t + I_y \cdot y_t + I_t) I_x$$

$$\frac{dE}{dy_t} = 0 = 2 \sum (I_x \cdot x_t + I_y \cdot y_t + I_t) I_y$$

$$\begin{bmatrix} \sum I_x^2 & \sum I_x I_y \\ \sum I_x I_y & \sum I_y^2 \end{bmatrix} \begin{bmatrix} x_t \\ y_t \end{bmatrix} = \begin{bmatrix} -\sum I_x I_t \\ -\sum I_y I_t \end{bmatrix}$$

known unknown
of vector Known

* Summary for BBOF

$$\begin{bmatrix} \sum I_x^2 & \sum I_x I_y \\ \sum I_x I_y & \sum I_y^2 \end{bmatrix} \begin{bmatrix} x_t \\ y_t \end{bmatrix} = \begin{bmatrix} -\sum I_x I_t \\ -\sum I_y I_t \end{bmatrix}$$

$X = A^{-1}B$

- solve for optical flow vector in patch (x_t, y_t)

- solution exist when the structure tensor matrix is non singular (at corner).

* Weighted BBOF

- given preference to satisfying OFCE at the center of the patch.

$$E(x_t, y_t) = \sum w(x, y) (OFCE)^2$$

$$w(x, y) = \frac{1}{\| (x, y) - (x_c, y_c) \| + 1}$$

or -

$$w(x, y) = \exp(-\| (x, y) - (x_c, y_c) \|^2 / \sigma^2)$$

$$\begin{bmatrix} \sum w Ix^2 & \sum w Ix Iy \\ \sum w Ix Iy & \sum w Iy^2 \end{bmatrix} v = \begin{bmatrix} -\sum w Ix Iy \\ -\sum w Iy^2 \end{bmatrix}$$

—

Known Unknown Known