

LECTURE 1

Geometry 3D and images:

- PRESERVE:**
- straight lines are mapped to straight lines
 - size of objects is inverse proportional to the distance

- NOT PRESERVE:**
- angles
 - parallel lines meet at the vanishing point of their direction.

Transformations in projective geometry

- translation $\begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix}$ #dof: 2

- euclidian $\begin{bmatrix} \cos\theta & -\sin\theta & t_x \\ \sin\theta & \cos\theta & t_y \\ 0 & 0 & 1 \end{bmatrix}$ #dof: 3

- similarity $\begin{bmatrix} s(\cos\theta) & s(-\sin\theta) & t_x \\ s(\sin\theta) & s(\cos\theta) & t_y \\ 0 & 0 & 1 \end{bmatrix}$ #dof: 4

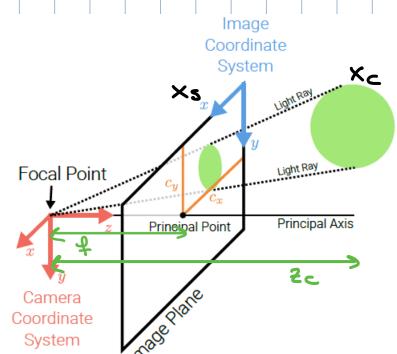
- affine (+shearing) $\begin{bmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \\ 0 & 0 & 1 \end{bmatrix}$ #dof: 6

- projective $\begin{bmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \\ a_7 & a_8 & a_9 \end{bmatrix}$ #dof: 8 (up to scale)

Converting a point in the 3D space in the Camera Coordinate System in the image plane

$$\vec{x}_s = \underbrace{\begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}}_{\text{calibration matrix } K} \underbrace{\begin{bmatrix} s & & \\ & s & \\ & & 1 \end{bmatrix}}_{\text{skew}} \underbrace{\begin{bmatrix} 0 & & \\ & 0 & \\ & & 1 \end{bmatrix}}_{\text{principal point offset}} \vec{x}_c$$

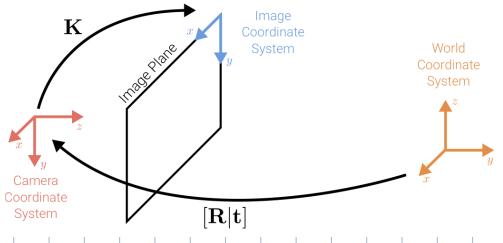
- calibration matrix contains the camera intrinsic parameters.
- skew is needed if the sensor is not mounted perpendicular to the optical axis.
- usually $f_x = f_y$ and $s=0$



Chaining transformations

$$\vec{x}_s = \underbrace{\begin{bmatrix} K & 0 \\ 0^t & 1 \end{bmatrix}}_{\substack{\text{camera intrinsic} \\ \text{parameters} \\ (3 \times 4)}} \underbrace{\begin{bmatrix} R & t \\ 0^t & 1 \end{bmatrix}}_{\substack{\text{camera pose} \\ (4 \times 4)}} \vec{x}_w = K [R \ t] \vec{x}_w = P \vec{x}_w$$

#dof = 5 + 6 = 11



LECTURE 2

Keypoint Localization - Harris Detector

Comparing each pixel before and after by summing up the squared differences (SSD):

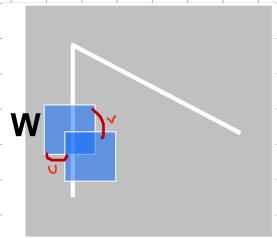
$$E(u, v) = \sum_{(x, y) \in W} [I(x+u, y+v) - I(x, y)]^2$$

If the motion (u, v) is small, then the 1st order approx. is good:

$$I(x+u, y+v) \approx I(x, y) + \frac{\partial I}{\partial x} u + \frac{\partial I}{\partial y} v = I(x, y) + [I_x \ I_y] \begin{bmatrix} u \\ v \end{bmatrix}$$

$$E(u, v) = \sum_{(x, y) \in W} [[I_x \ I_y] \begin{bmatrix} u \\ v \end{bmatrix}]^2$$

$$E(u, v) = \sum_{(x, y) \in W} [u \ v] \begin{bmatrix} I_x^2 & I_x I_y \\ I_y I_x & I_y^2 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$$



second order moment matrix

$$\sum_{(x, y) \in W} \begin{bmatrix} I_x^2 & I_x I_y \\ I_y I_x & I_y^2 \end{bmatrix}$$

NO ROTATION INVARIANT

M

$$\sum_{(x, y) \in W} q(x, y) \begin{bmatrix} I_x^2 & I_x I_y \\ I_y I_x & I_y^2 \end{bmatrix}$$

**ROTATION INVARIANT
(smooth with Gaussian)**

gaussian filter over the window

Thresholding eigenvalues to detect corners

- Harris & Stephens : $R = \det(M) - k \text{trace}^2(M) = \lambda_1 \lambda_2 - k(\lambda_1 + \lambda_2)^2$

- Kanade & Tomasi : $R = \min(\lambda_1, \lambda_2)$

Recap: Harris detector

1. Compute the M matrix

2. Compute the cornerness function R

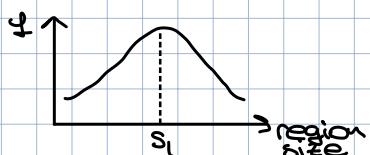
3. Perform non-maximum suppression

(if 2 pixels are close to each other, there are not two corners: it's the same one. Non-maximum suppression is an heuristic approach where you pick the maximum response in a local region).

In order to compare points, we need to compute a descriptor over a region that must be scale-invariant.

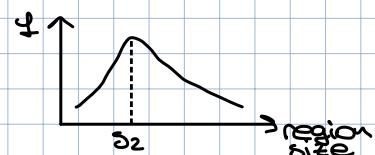
Automatic Scale Selection

1. Design a **signature function** on the region that is "scale invariant". For a point in one image, we can consider a function of a region size.
2. Take the local maximum of this function.



$$s_2 = t_2 \cdot s_1$$

$$\text{SCALE} = \frac{1}{2}$$



3 Normalize: rescale to fixed size.

Signature functions:

- Laplacian-of-Gaussian = "blob" detector
characteristic scale: scale that produces peak in Laplacian response.
- Difference-of-Gaussian: fast approximation of Laplacian-of-Gaussian

$$\nabla^2 g = \frac{\partial^2 g}{\partial x^2} + \frac{\partial^2 g}{\partial y^2}$$

Local descriptors

Descriptor	Features
Image patch	Only if geometry and appearance is unchanged
Image gradients	Invariant to absolute intensity values
Color histogram	Invariant to changes in scale and rotation
Spatial histogram	Some invariance to deformations

Orientation normalization



The dominant gradient is found and the patches are rotated to the dominant direction.

SIFT - Scale Invariant Feature Transform

Algorithm:

1. Scale-space extrema detection

2. Keypoint localization

3. Keypoint descriptor

- divide patch into 6×6 sub-patches: 16 cells
- compute histogram of gradient orientations (8 reference angles)
- resulting descriptor: $6 \times 6 \times 8 = 128$ dimensions

4. orientation normalization

5. Keypoint matching

LECTURE 3

Optical flow

Motion field = 2D motion field representing the projection of the 3D motion of points in the scene onto the image plane.

Optical flow = 2D velocity field describing the apparent motion in the images.

Given 2 consecutive image frames, estimate the motion of each pixel.

Assumptions:

- **brightness constancy** (color constancy)
allows for pixel comparison

$$I(x(t), u(t), t) = \text{constant}$$

- **small motion**
linearization of the brightness constancy constraint

$$I(x + u\delta t, y + v\delta t, t + \delta t) = I(x, u, t)$$

BRIGHTNESS CONSTANCY EQUATION

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

$$I(x(t), u(t), t) \approx I(x(t), u(t), t) + \frac{\partial I}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial I}{\partial y} \frac{\partial y}{\partial t} + \frac{\partial I}{\partial t}$$

$$\frac{\partial I}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial I}{\partial y} \frac{\partial y}{\partial t} + \frac{\partial I}{\partial t} = 0 \quad / : \partial t$$

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

$$\left. \begin{array}{l} I_x \cdot u + I_y \cdot v + I_t = 0 \\ \nabla I^T \cdot \begin{pmatrix} u \\ v \end{pmatrix} + I_t = 0 \end{array} \right\} \text{shorthand notation}$$

$$\left. \begin{array}{l} I_x \cdot u + I_y \cdot v + I_t = 0 \\ \nabla I^T \cdot \begin{pmatrix} u \\ v \end{pmatrix} + I_t = 0 \end{array} \right\} \text{vector form}$$

How to compute: $I_x = \frac{\partial I}{\partial x}$, $I_y = \frac{\partial I}{\partial y}$ → Sobel filter
 $I_t = \frac{\partial I}{\partial t}$ → Derivative-of-Gaussian filter

$I_t = \frac{\partial I}{\partial t}$ → frame differencing (per pixel $I_t = I(t) - I(t+\delta t)$)

$u = \frac{\partial x}{\partial t}$; $v = \frac{\partial y}{\partial t}$ → unknown in the optical flow problem

Where do we get more equations?

o Lucas-Kanade flow: "constant" flow \rightarrow flow is constant for all pixels
 (local method - sparse)

Assume that the surrounding patch (say 5x5) has constant flow.
 Using 5x5 image patch, give us 25 equations.

$$I_x(p_i)u + I_y(p_i)v = -I_t(p_i) \quad i=1, \dots, 25$$

equivalent in solving: $A^T A \cdot \hat{x} = A^T b$

$$\hat{x} = (A^T A)^{-1} A^T b$$

$$A^T A$$

$$\hat{x}$$

$$A^T b$$

$$\begin{bmatrix} \sum_{p \in P} I_x I_x & \sum_{p \in P} I_x I_y \\ \sum_{p \in P} I_y I_x & \sum_{p \in P} I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum_{p \in P} I_x I_t \\ \sum_{p \in P} I_y I_t \end{bmatrix}$$

where summation
 is over each
 pixel p in patch P

Implications:
 - corners are when λ_1, λ_2 are big.
 - corners are regions with two different gradient dir.
 - corners are good places to compute flow

o Horn-Schunck flow: "smooth" flow \rightarrow flow can vary from pixel to pixel
 (global method - dense)

$$\min_{(u,v)} \sum_{i,j} \{ E_s(i,j) + \lambda E_d(i,j) \}$$

$$E_s = \frac{1}{4} [(u_{ij} - u_{i+1,j})^2 + (u_{ij} - u_{i,j+1})^2 + (v_{ij} - v_{i+1,j})^2 + (v_{ij} - v_{i,j+1})^2]$$

smoothness constancy

$$E_d = [I_x u_{ij} + I_y v_{ij} + I_t]^2$$

brightness constancy

$$\hat{u}_{kl} = \bar{u}_{kl} - \frac{I_x \bar{u}_{kl} + I_y \bar{v}_{kl} + I_t}{\lambda^{-1} + I_x^2 + I_y^2} I_x \quad \hat{v}_{kl} = \bar{v}_{kl} - \frac{I_x \bar{u}_{kl} + I_y \bar{v}_{kl} + I_t}{\lambda^{-1} + I_x^2 + I_y^2} I_y$$

Algorithm:

1. precompute I_u, I_v, I_t

2. Initialize flow field $\begin{cases} u=0 \\ v=0 \end{cases}$

3. while NOT converged compute flow field updates for each pixel

$$\hat{u}_{kl} = \bar{u}_{kl} - \frac{I_x \bar{u}_{kl} + I_y \bar{v}_{kl} + I_t}{\lambda^{-1} + I_x^2 + I_y^2} I_x \quad \hat{v}_{kl} = \bar{v}_{kl} - \frac{I_x \bar{u}_{kl} + I_y \bar{v}_{kl} + I_t}{\lambda^{-1} + I_x^2 + I_y^2} I_y$$

\bar{u}_{kl} = old average

\hat{u}_{kl} = new value

LECTURE 4

Machine learning basics

Linear regression

Goal: predict y for a previously unseen input x

Error function: $f(x, \omega) = \omega^T x$ $x = (1, x)^T$

$$E(\omega) = \sum_i (f(x_i, \omega) - y_i)^2 \\ = \|x\omega - y\|_2^2$$

$$\nabla_{\omega} E(\omega) = \nabla_{\omega} (x\omega - y)^T (x\omega - y) = 2x^T x \omega - 2x^T y \\ \stackrel{!}{=} 0 \Rightarrow \hat{\omega} = \underbrace{(x^T x)^{-1} x^T y}_{\text{Moore-Penrose inverse}}$$

Polynomial curve fitting

Same as linear regression but $x = (1, x^1, x^2, \dots, x^m)^T$

$$X = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^m \end{pmatrix} \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} \quad \omega = \begin{pmatrix} \omega_0 \\ \vdots \\ \omega_m \end{pmatrix}$$

CAPACITY: complex of functions which can be represented by model f ($= n$)

UNDERFITTING: low capacity

OVERRFITTING: high capacity

General approach:

1. choose hyperparameter (H, η, \dots) and train model
2. k -fold cross validation on test-set
3. validation set to confirm

60%
20%
20%

Ridge regression

$$E(\omega) = \|x\omega - y\|_2^2 + \lambda \|\omega\|_2^2$$

$$\nabla_{\omega} E(\omega) = 2x^T x\omega - 2x^T y + 2\lambda\omega \\ \stackrel{!}{=} 0 \Rightarrow \hat{\omega} = \underbrace{(x^T x + \lambda I)^{-1} x^T y}_{(x^T x + \lambda I)^{-1} x^T y}$$

Maximum Likelihood Estimation

$$\hat{\omega}_{ML} = \arg \max_{\omega} p(y|x, \omega) = \arg \max_{\omega} \prod_i p(y_i|x_i, \omega)$$

$$= \arg \max_{\omega} \sum_i \log p(y_i|x_i, \omega) \quad \text{assuming } p(y_i|x_i, \omega) \sim N(\omega^T x_i, \sigma^2)$$

$$= \arg \min_{\omega} \|x\omega - y\|_2^2 \quad p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

Gaussian distribution

$\arg \min_{\omega} \|x\omega - y\|_1$
 Laplace distribution

$p(x) = \frac{1}{2\sigma} e^{-\frac{|x-\mu|}{\sigma}}$

Logistic regression

$$p(y|x, w) = \text{Bernoulli distribution}$$

$$= \theta^y (1-\theta)^{(1-y)}$$

$$\theta = \sigma(w^T x) \quad - \text{sigmoid function} \quad \sigma(w^T x) = \frac{1}{1+e^{-w^T x}}$$

$$\hat{W}_{\text{ML}} = \underset{\theta}{\operatorname{argmax}} \sum_i \log p(y_i|x_i; \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_i \log [\theta^{y_i} (1-\theta)^{(1-y_i)}]$$

$$= \underset{\theta}{\operatorname{argmin}} \sum_i [-y_i \log \theta - (1-y_i) \log (1-\theta)] \quad \rightarrow \text{Not quadratic}$$

binary cross entropy loss $\ell(\theta_i, y_i)$

\Rightarrow gradient descent

Gradient Descent

Algorithm:

1. Initialize weights w^0 and pick learning rate η
2. For all data points $i \in \{1, \dots, N\}$ do:
 - 2.1 Forward propagate x_i through the network to calculate prediction \hat{y}_i
 - 2.2 Backpropagate to obtain gradient $\nabla_w \ell_i(w^t) \equiv \nabla_w \ell(\theta_i, y_i, w^t)$
3. Update gradients: $w^{t+1} = w^t - \eta \frac{1}{N} \sum_i \nabla_w \ell_i(w^t)$

Expensive since: $\begin{cases} \dim(w) \geq 1 \text{ million} \\ N \geq 1 \text{ million} \end{cases}$

Stochastic Gradient Descent

Algorithm:

1. Initialize weights w^0 and pick learning rate η and minibatch size $|\mathcal{B}_{\text{batch}}|$
2. Draw random minibatch $\{(x_1, y_1), \dots, (x_B, y_B)\} \subseteq \mathcal{X}$ (with $B \ll N$)
3. For all minibatches elements $b \in \{1, \dots, B\}$ do:
 - 3.1 Forward propagate x_b through the network to calculate prediction \hat{y}_b
 - 3.2 Backpropagate to obtain gradient $\nabla_w \ell_b(w^t) \equiv \nabla_w \ell(\theta_b, y_b, w^t)$
4. Update gradients: $w^{t+1} = w^t - \eta \frac{1}{B} \sum_b \nabla_w \ell_b(w^t)$
5. If validation error decreases, go to step 2, otherwise stop.

Deep neural network basics

Computation Graphs - Backpropagation

- we can write mathematical expressions as a computation graph
- values are efficiently computed forward, gradients backwards
- multiple incoming gradients are summed up (multivariate chain rule)
- modularity: each node just only "knows" how to compute gradients w.r.t. its own arguments

Bauer's Formula: let G be a computation graph. For any two nodes i and j in G , let $P(i,j)$ be the set of all directed paths starting at j and ending at i . The partial derivative of z_i with respect to z_j is then equal to

$$\frac{\partial z_i}{\partial z_j} = \sum_{\pi \in P(i,j)} \prod_{(k,l) \in \pi} \frac{\partial z_k}{\partial z_l}$$

Multi-Layer Perceptrons

Each layer i comprises multiple neurons j which are implemented as affine transformations ($a^T x + b$) followed by non-linear activation functions (g).

$$h_{ij} = g(a_i^T h_{i-1} + b_{ij})$$

- Activation functions $g(\cdot)$

- Sigmoid $\sigma(x) = \frac{1}{1+e^{-x}}$



Smooth
vanishing

- Tanh $\tanh(x)$



Smooth
vanishing

- ReLU $\max(0, x)$



Quick to compute the gradient
vanishing

- Leaky ReLU $\max(0.1x, x)$



Quick to compute the gradient

- Maxout $\max(w_1^T x + b_1, w_2^T x + b_2)$

$$\begin{cases} x & x \geq 0 \\ \lambda(e^x - 1) & x < 0 \end{cases}$$

With linear activation function (e.g. $g(x)=x$), a multi-layer network can only express linear functions.

Universal approximation theorem

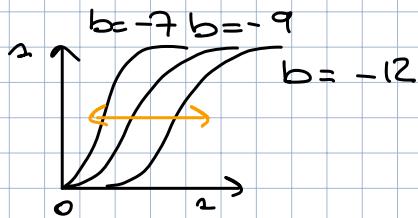
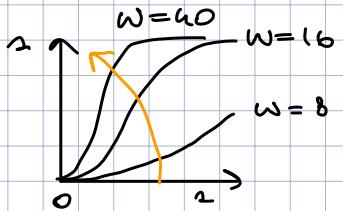
Let σ be any continuous discriminatory function. Then finite sums of the form

$$G(x) = \sum_{i=1}^N \alpha_i \sigma(\omega_i^\top x + b_i)$$

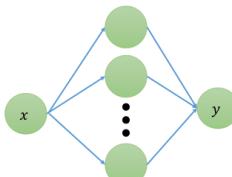
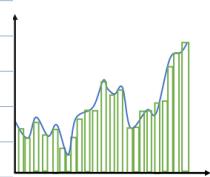
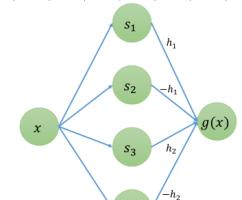
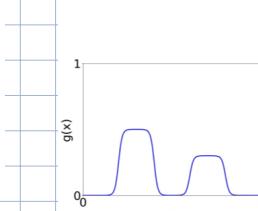
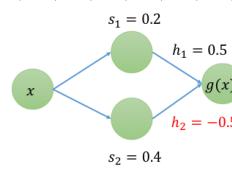
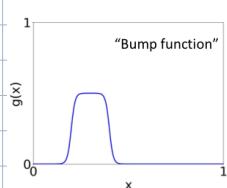
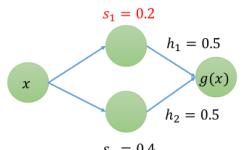
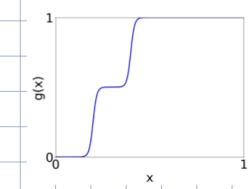
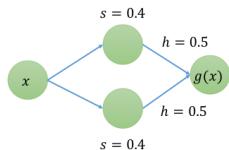
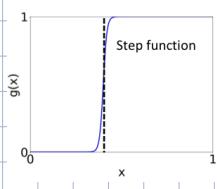
are dense in the space of continuous functions $C(I^n)$ on the n -dimensional unit cube I^n . In other words, given any $f \in C(I^n)$ and $\epsilon > 0$, there is a sum, $G(x)$ for which

$$|G(x) - f(x)| < \epsilon \quad \text{for all } x \in I^n$$

- SIGMOID FUNCTION shape: $g(x) = \sigma(\omega^\top x + b)$



Learning edge $s = -\frac{b}{\omega} = 0, 4$



The theorem says that feed forward networks can approximate $f(x)$. It does not tell anything about the chances of learning the correct parameters.

Regression problem

Gaussian distribution thin tail

$$p(y|x, \omega) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - f_\omega(x))^2}{2\sigma^2}\right)$$

$$\hat{\omega}_{ML} = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^N (\hat{f}_\omega(x_i) - u_i)^2$$

penalizes outliers more

Laplace distribution heavy tail

$$p(y|x, \omega) = \frac{1}{2b} \exp\left(-\frac{|y - \mu|}{b}\right)$$

$$\hat{\omega}_{ML} = \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^N |\hat{f}_\omega(x_i) - u_i|$$

Classification Problems

Categorical distribution

$$\begin{aligned}\hat{w}_{ML} &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^N \log (q_i | x_i, \omega) \\ &= \underset{\omega}{\operatorname{argmax}} \sum_{i=1}^N \log \sum_{c=1}^C q_w^{(c)}(x_i)^{q_c} \\ &= \underset{\omega}{\operatorname{argmin}} \sum_{i=1}^N \sum_{c=1}^C -q_i c \log q_w^{(c)}(x_i)\end{aligned}$$

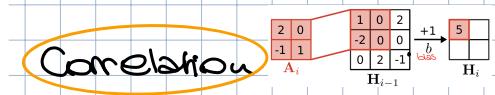
cross-entropy loss

Softmax

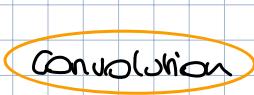
$$\text{softmax}(x) = \left(\frac{\exp(x_1)}{\sum_{k=1}^C \exp(x_k)}, \dots, \frac{\exp(x_C)}{\sum_{k=1}^C \exp(x_k)} \right)$$

Convolutional neural network basics

Convolution layers



$$I'(i, j) = \sum_{m=-K}^K \sum_{n=-K}^K k(m, n) I(i+m, j+n)$$



(if $k(m, n) = k(-m, -n)$
then convolution = correlation)

$$I'(i, j) = \sum_{m=-K}^K \sum_{n=-K}^K k(-m, -n) I(i+m, j+n)$$

Fully connected

#Weights:

$$H \times W \times \text{Cat} \times (H \times W \times C_{in} + 1)$$

convolutional layer

#Weights:

$$\text{Cat} \times (K \times K \times C_{in} + 1)$$

Padding: add boundary of appropriate size with zeros around input tensor

Downsampling: reduces the spatial resolution; increases the receptive field (which pixels influence a neuron)

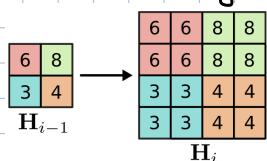
Stride: size of the jump

Pooling: has no parameters (typical pooling operations: max, min, mean)

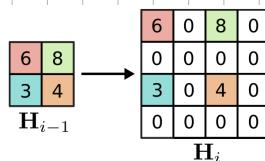
Example: Filter: $K \times K$ $W_{out} = \left\lfloor \frac{W_{in} + 2P - K}{S} \right\rfloor + 1$ $H_{out} = \left\lfloor \frac{H_{in} + 2P - K}{S} \right\rfloor + 1$

Upsampling: yields output at the same resolution as input

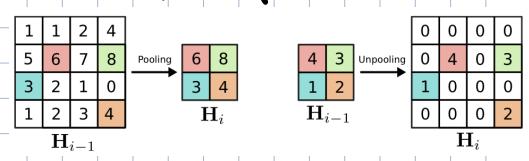
- nearest neighbor



- bed of nails



- max unpadding



LECTURE 5

Object recognition

Part-and-shape models

- Object as a set of parts
- relative locations between parts
- appearance of parts

History:

- 60-90: The geometric era
- 90: Appearance-based models
- mid 90: Sliding window approaches
- late 90: Local features
- 2000: Parts-and-shape models
- mid 00: Bag-of-features

Bag-of-features

Assumption: spatial information of local features can be ignored for object recognition.

An **image** (datapoint) is a vector of counts over each word (feature)

$$v_d = [n(w_{1,d})\alpha_1, n(w_{2,d})\alpha_2, \dots]$$

$n(\cdot)$ = counts the number of occurrences

α_i = weights each word by a heuristic = $\log \left\{ \frac{D}{\sum_{d=1}^D \mathbb{1}[w_i \in d]} \right\}$ → inverse document frequency

Pipeline:

1. Feature extraction
extracting features from images (regular grid + interest points) + SIFT descriptor
2. Dictionary Learning
learn visual words using clustering (K-means)
3. Encode
build BOW vectors for each image
4. Classify
Train and test data using BOWs

Pros

- flexible to geometry / deformations / viewpoint
- compact summary of image content
- provides fixed dimensional vector representing for sets
- very good results in practice

Cons

- background and foreground mixed
- optimal vocabulary formation remains unclear
- basic models ignore geometry

Precision: $\frac{tp}{tp + fn}$

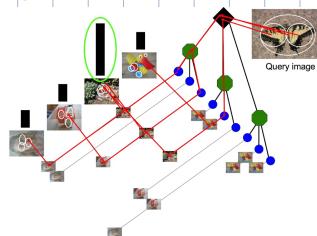
Recall: $\frac{tp}{tp + fn}$

Specific object detection

Inverted index (word : images)

1. extract words in query
2. invert file index to find relevant frames
3. compare / sort word-counts

Requirement: sparsity



Object Category Detection (object or not-object)

Object Recognition

1. Specify Object model
what are the object parameters?
2. Generate hypotheses
propose an alignment of the model to the image
3. Score hyperparameters
mainly gradient-based features, many classifiers
4. Resolve detection
rescore each proposed object based on context information.

Sliding window Face detection with Viola-Jones

Slow training
Fast detection

1. Integral images for fast evaluation $I_{\Sigma}(x, u) = \sum_{\substack{x' \leq x \\ u' \leq u}} i(x', u')$
 ↘ rectangle filters repeated thousands times
 ↘ how to choose right ones?

2. Boosting for feature selection (AdaBoost)

2-class classification problem

- Given: training set $X = \{x_1, \dots, x_N\}$ with target values $T = \{t_1, \dots, t_N\}$, $t_i \in \{-1, 1\}$.
- Associated weights $W = \{w_1, \dots, w_N\}$ for each training point.

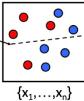


Basic steps

- In each iteration, AdaBoost trains a new weak classifier $h_m(x)$ based on the current weighting coefficients $W^{(m)}$.
- We then adapt the weighting coefficients for each point
 Increase w_n if x_n was misclassified by $h_m(x)$.
 Decrease w_n if x_n was classified correctly by $h_m(x)$.
- Make predictions using the final combined model

$$H(\mathbf{x}) = \text{sign} \left(\sum_{m=1}^M \alpha_m h_m(\mathbf{x}) \right) \quad h_m(\mathbf{x}) = \begin{cases} -1 & \text{if } f_m < \theta_m \\ +1 & \text{otherwise} \end{cases}$$

1. Initialization: Set $w_n^{(1)} = \frac{1}{N}$ for $n = 1, \dots, N$.



2. For $m = 1, \dots, M$ iterations

a) Train a new weak classifier $h_m(x)$ using the current weighting coefficients $W^{(m)}$ by minimizing the weighted error function

$$J_m = \sum_{n=1}^N w_n^{(m)} I(h_m(x_n) \neq t_n) \quad I(A) = \begin{cases} 1, & \text{if } A \text{ is true} \\ 0, & \text{else} \end{cases}$$

b) Estimate the weighted error of this classifier on X :

$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(h_m(x_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}}$$

c) Calculate a weighting coefficient for $h_m(x)$:

$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}$$

d) Update the weighting coefficients:

$$w_n^{(m+1)} = w_n^{(m)} \exp \{ \alpha_m I(h_m(x_n) \neq t_n) \}$$

- Evaluate all selected weak classifiers on test data.

$$h_1(\mathbf{x}), \dots, h_m(\mathbf{x})$$

- Final classifier is weighted combination of selected weak classifiers:

$$H(\mathbf{x}) = \text{sign} \left(\sum_{m=1}^M \alpha_m h_m(\mathbf{x}) \right)$$

- Very simple procedure!

LECTURE 6

Segmentation as grouping

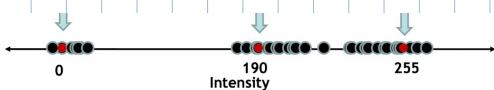
Goals:

- identify groups of pixels that belong together
- delineate object and background regions
- group similar-looking pixels for efficiency of further processing in "superpixels".

K-Means:



To determine the three main intensities that define our groups we need to cluster. The three "centers" as the representative intensities



Algorithm:

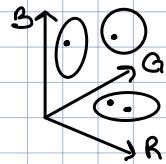
1. randomly initialize the cluster centers c_1, \dots, c_k
2. Given cluster centers, determine points in each cluster
 - for each point p_i , find the closest c_j . Put p_i into cluster j
3. Given points in each cluster, solve for c_i
 - set c_i to be the mean of points in cluster i
4. If c_i have changed, repeat step 2.

Properties: - always converge

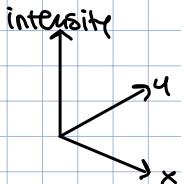
- ° might be stuck in local minimum of objective function

$$\sum_{\text{clusters:}} \sum_{\substack{\text{points } p \\ \text{in cluster:}}} \|p - c_i\|^2$$

Depending on what we choose as the feature space, we can group pixels in different ways.



← { feature space: color value (3D)
point $P = \begin{pmatrix} R = 15 \\ G = 189 \\ B = 2 \end{pmatrix}$ }



← { feature space: intensity + position similarity
point $P = \begin{pmatrix} X = 4 \\ Y = ? \\ I = ? \end{pmatrix}$ } (spatial coherence)

PROS

- Simple, fast to compute
- Converge to local minimum of within-cluster squared error

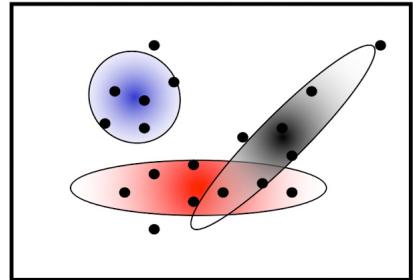
CONS

- setting k ?
- sensitive to initial centers
- sensitive to outliers
- Detects spherical clusters only
- Assuming mean can be computed

Mixture of gaussians (MoG)

Probabilistic clustering:

instead of treating the data as a bunch of points, assume that they are all generated by sampling a continuous function. This function is called the **generative model** and is defined by vector of parameters θ .



In MoG the generative model is a mixture of gaussians. There will be K gaussian blobs $\sim N(\mu_b, V_b)$

$$p(x|\mu_b, V_b) = \frac{1}{\sqrt{(2\pi)^d |V_b|}} \exp\left(-\frac{1}{2} (x - \mu_b)^T V_b^{-1} (x - \mu_b)\right)$$

Blob b is selected with probability α_b .
The likelihood of observing x is.

$$P(x|\theta) = \sum_{b=1}^K \alpha_b P(x|\theta_b)$$

$$\theta = [\mu_1, \dots, \mu_K, V_1, \dots, V_K]$$

Goal: find θ that maximize the likelihood function over all N datapoints

$$X = \{x_1, \dots, x_N\}$$

$$P(X) = \prod_{x_i \in X} P(x_i|\theta)$$

EM - Approach:

E-step

- Compute probability that point x is in blob b , given current guess of θ

$$P(b|x, \mu_b, V_b) = \frac{\alpha_b P(x|\mu_b, V_b)}{\sum_{i=1}^K \alpha_i P(x|\mu_i, V_i)}$$

M-step

- Compute overall probability that blob b is selected

$$\alpha_b^{new} = \frac{1}{N} \sum_{i=1}^N P(b|x_i, \mu_b, V_b) \quad (N \text{ data points})$$

- Mean of blob b

$$\mu_b^{new} = \frac{\sum_{i=1}^N x_i P(b|x_i, \mu_b, V_b)}{\sum_{i=1}^N P(b|x_i, \mu_b, V_b)}$$

- Covariance of blob b

$$V_b^{new} = \frac{\sum_{i=1}^N (x_i - \mu_b^{new})(x_i - \mu_b^{new})^T P(b|x_i, \mu_b, V_b)}{\sum_{i=1}^N P(b|x_i, \mu_b, V_b)}$$

given current guess of blobs, compute probabilistic ownership of each point.

given ownership probabilities, update blobs to maximize likelihood function.

PROS

- Probabilistic interpretation
- Soft assignments between data points and clusters
- Generative model, can predict novel data points
- Relatively compact storage ($O(Kd^2)$)

CONS

- Initialization (often a good idea to start from output of k-means)
- Local minima
- Need to know number of components K
- Need to choose blob generative model
- Numerical problems are often a nuisance

Mean-shift segmentation (model-free clustering)

Iterative mode search:

1. Initialize random seed center and window W
2. Calculate center of gravity ("the mean") of W : $\sum_{x \in W} x H(x)$
3. Shift the search window to the mean
4. Repeat steps 2+3 until convergence

Overall algorithm:

1. Choose features (color, gradients, texture, etc)
2. Initialize windows at individual pixel location (size h = scale of clustering)
3. Start mean-shift from each window until convergence
4. Merge windows that end up near the same "peak" or mode

PROS

- General, application-independent tool
- Model-free, does not assume any prior shape or data clusters
- Just a single parameter (window size h)
- Finds variable number of modes given the same h
- Robust to outliers

CONS

- Output depends on window size h
- Window size selection is not trivial
- Computationally expensive
- Does not scale well with dimension of feature space

Hough transforms (edge-based)

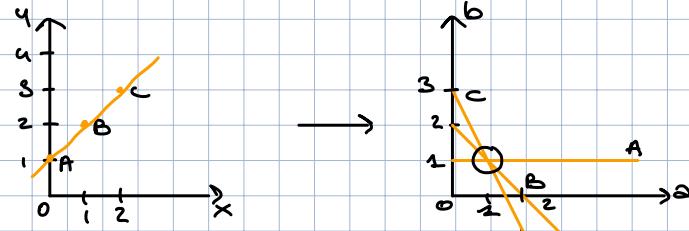
Uses the structure of shapes to extract them in a parameter space of limited dimension.

Detecting straight lines

$$\text{Straight line: } y = ax + b \rightarrow b = (-x)a + y$$

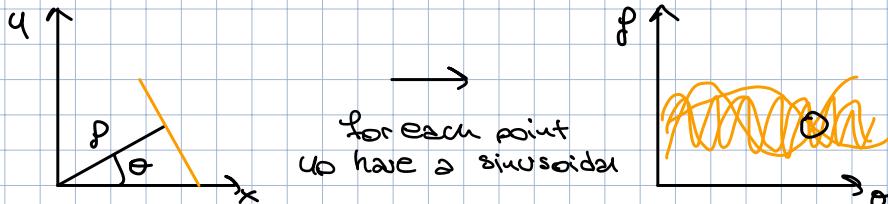
The Hough transform:

1. Scrubline on edge pts
2. For each point draw the above line in (a, b) parameter space



! Cannot represent vertical lines ($a=\infty$)

$$\text{Straight line: } x \cos\theta + y \sin\theta = p$$



Remark

1. Time consuming
2. Robust to outliers and lots of points not on a line
3. Peak detection is difficult with noise in edge coordinates.
4. Robustness of peak detection increased by weighting contributions (votes that are proportional with intensity gradients magnitude).
5. Can do it also for circles and even higher-order parametric shapes (but it quickly gets harder as number of shape parameters grows).

GraphCut

In graph-cut, computing the disparity is formulated as a graph labeling problem. Each pixel corresponds to a graph node and each disparity to a label. The goal is to find a labeling $f : P \rightarrow L$ which minimizes the energy function

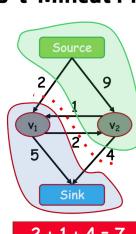
$$E(f) = E_{\text{data}}(f) + \lambda E_{\text{smooth}}(f) = \sum_{p \in P} D_p(f_p) + \lambda \sum_{p, q \in N} S(f_p, f_q).$$

P represents the set of pixels and L represents a discrete set of labels corresponding to different disparities. $D_p(f_p)$ is the cost of assigning label f_p to pixel p and it is given by the SSD value calculated for the disparity corresponding to label f_p . N is the set of neighboring pixels and $S(f_p, f_q)$ is the cost of assigning labels f_p and f_q to neighboring pixels p and q .

The s-t-Mincut Problem

What is an st-cut?

An st-cut (S, T) divides the nodes between source and sink.



What is the cost of a st-cut?

Sum of cost of all edges going from S to T

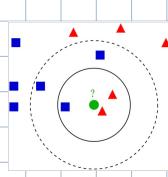
What is the st-mincut?

st-cut with the minimum cost

$$2 + 1 + 4 = 7$$

Learning-based approaches

K-nearest neighbors (KNN)



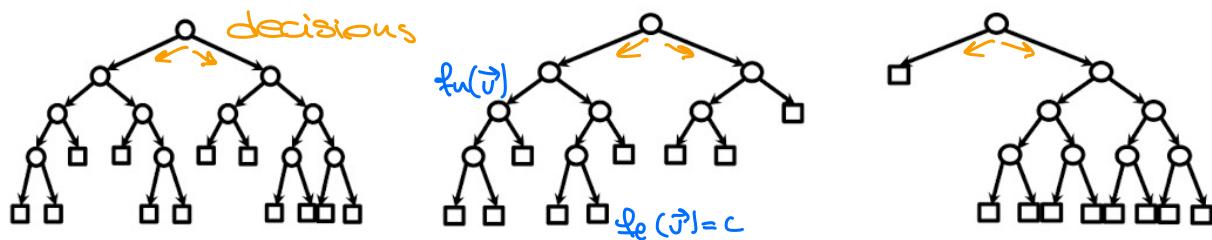
PROS

- Very simple to implement / understand
- Efficient implementations possible for approx. NNs
- Distance definition is flexible

CONS

- Highly depends on the definitions and K
- Need to keep the entire data in memory for distance computation
- For high dimensional problems (with high d) might need many training samples for accuracy.
- Other methods have better generalization ability

Random Forest



Training :

1. deciding on the (binary) node tests
2. deciding on when to stop splitting.

Randomness :

1. node tests result from randomly generated sets of candidate tests.

2. another aspect of the randomness is that different trees can be trained with different random subsets.

Stopping :

1. all samples belong to a class
2. too few samples

PROS

- Easy to implement
- Very efficient during testing
- Can easily use diverse features
- Can handle high dimensional spaces

CONS

- Lots of parametric choices
- Need large number of data
- Training can take time

LECTURE 7

Image Segmentation evaluation \rightarrow Mean intersection over union

$$\frac{1}{C} \sum_i \frac{n_{ii}}{t_i + \sum_j n_{ji} - n_{ii}}$$

C = number of categories
 n_{ij} = the number of pixels of class i predicted to belong to class j .
 t_i = total number of pixels of class i .

Transfer training:

STEP 1: Learn the filters on a related task and dataset (pre-training)

STEP 2: Learn on the target problem (fine-tuning)

- initialize the model with the pre-trained filters in Step 1 as much as possible
- train the model with BP and smaller learning rate
- adjust the filters instead of learning them from scratch

Semantic segmentation:

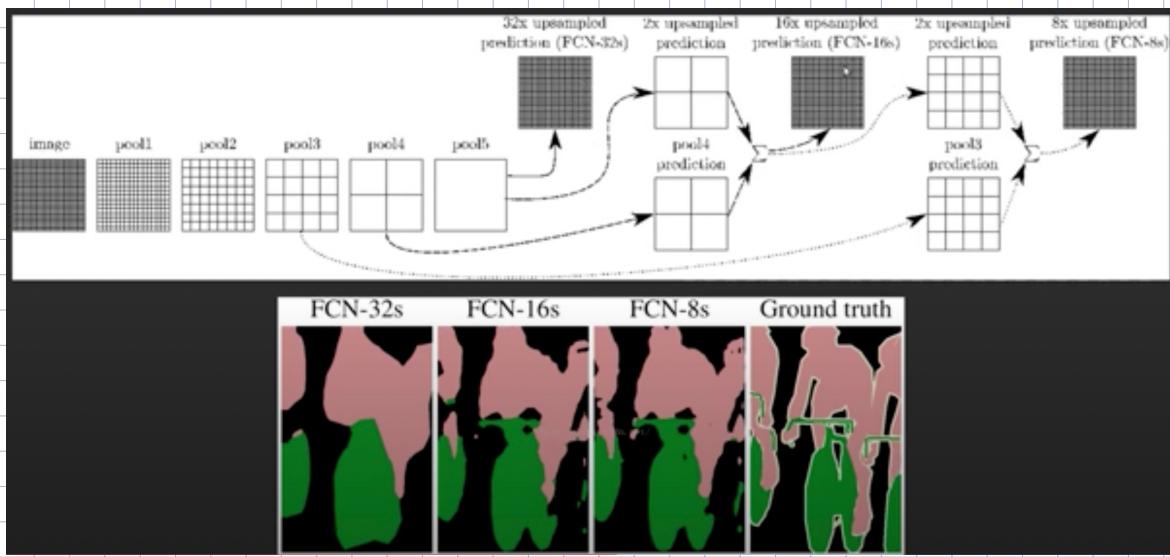
Hyper-column

- Viable way to transfer the image classification models to segmentation
- Reasonable results, but coarse boundaries
- Depends on classical methods, like region proposals based on pixel grouping

1. For each pixel location, take the activations from each level
2. Merge the feature together as the representation for the pixel
3. Multiscale representation (hypercolumn)

- Pick representative levels
- Additional convolutional layers to add capacity and change channel numbers
- Use bilinear interpolation to upsample the feature maps
- Sum the upsampled feature maps.

Fully convolutional networks



Transposed convolutions (deconvolutions) → use to upscale FCN

Input Kernel

$$\begin{array}{c}
 \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 2 & 3 \\ \hline \end{array} \quad \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 2 & 3 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 0 & 0 \\ \hline \end{array} + \begin{array}{|c|c|} \hline 0 & 1 \\ \hline 2 & 3 \\ \hline \end{array} + \begin{array}{|c|c|} \hline \quad & \quad \\ \hline 0 & 2 \\ \hline 4 & 6 \\ \hline \end{array} + \begin{array}{|c|c|} \hline \quad & \quad \\ \hline 0 & 3 \\ \hline 6 & 9 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 0 & 0 & 1 \\ \hline 0 & 4 & 6 \\ \hline 4 & 12 & 9 \\ \hline \end{array}
 \end{array}$$

Fully connected Conditional Random Field

$$E(L) = \sum_p E_p(l_p) + \sum_{p,q \in N} E(l_p, l_q)$$

PROS

- Very effective method
- The results look super sharp

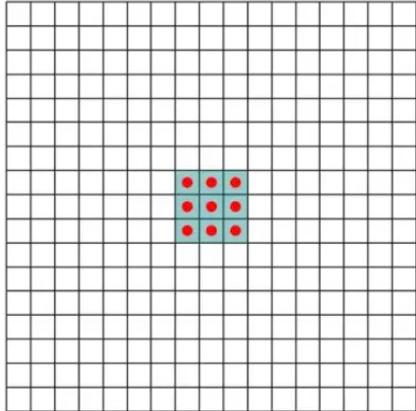
CONS

- A post-hoc processing method
- Manually tuned parameters

Dense prediction models

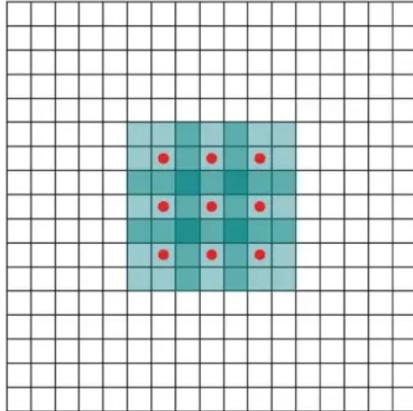
Dilated convolutions

1 Dilated Convolution



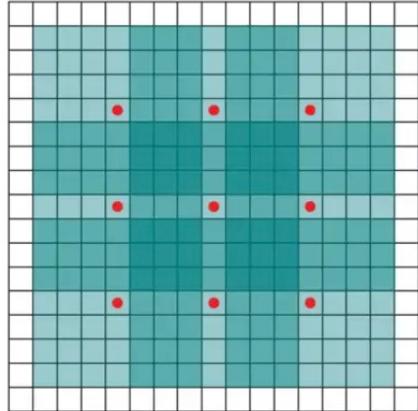
dilation = 1

2 Dilated Convolution



dilation = 2

4 Dilated Convolution



dilation = 4

Attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)V$$

LECTURE 8

Detection via classification

Pipeline:

1. Obtain training data
2. Define features (HoG)
3. Train a classifier
4. Use sliding window (AdaBoost)

Generalized Hough Transform (ISM - implicit shape model)

Basic ideas:

1. Learn an appearance codebook
 - Extract local features at interests points
 - Agglomerative clustering \Rightarrow codebook
2. Learn a star-topology structural model (spatial distribution)
 - o Feature are considered independent given obj. center
 - Match codebook to training images
 - Record matching positions on object

A visual vocabulary is used to index votes for object position (a visual word = "a part"). Objects are detected as consistent configurations of the observed parts (visual words).

PROS

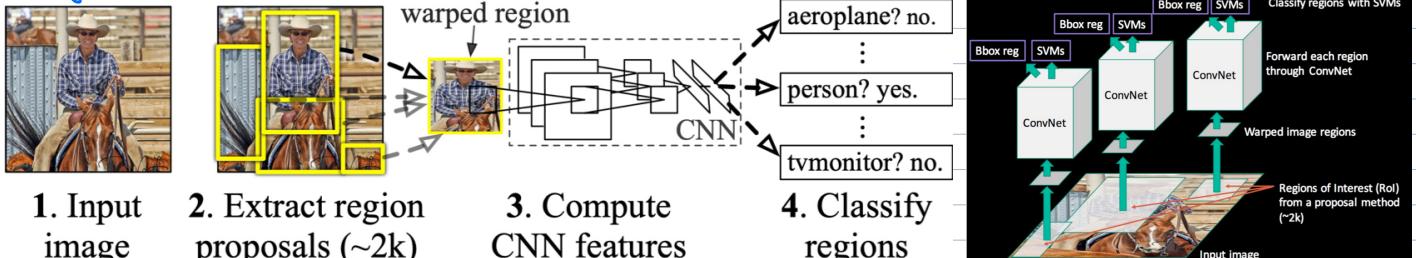
- Works well for many different object categories
- Flexible geometric model
(can recombine parts seen on different training examples)
- Learning from few (50-100) training examples
- Optimized for detection, good localization properties

CONS

- Needs supervised training data
 - o Object bounding boxes for detection
 - o Segmentations for top-down segmentation
- Only weak geometric constraints
 - o result segmentations may contain superfluous body parts
- Purely representative model
 - o No discriminative learning

Object detection with deep learning

Regions with CNN features (R-CNN)



Classification is made with "Linear mode with class-dependent weights"

Linear SVM \rightarrow support vector machine

$$f_c(x_{f_{c7}}) = w_c^T x_{f_{c7}}$$

$$\begin{cases} x_{f_{c7}} &= \text{features from the network} \\ c &= \text{object class} \end{cases}$$

Prediction of the 2D box

- Necessary, since the proposal region might not fully coincide with the (annotated) object bounding box

- Perform regression for location (x^*, y^*) , width w^* and height h^*

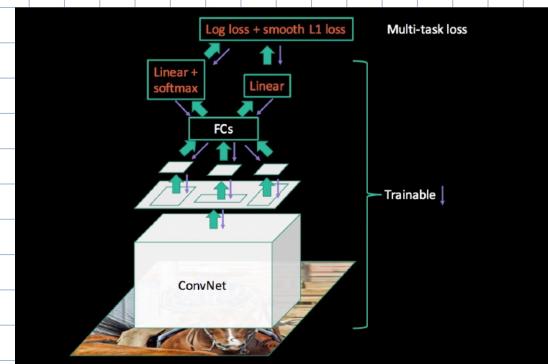
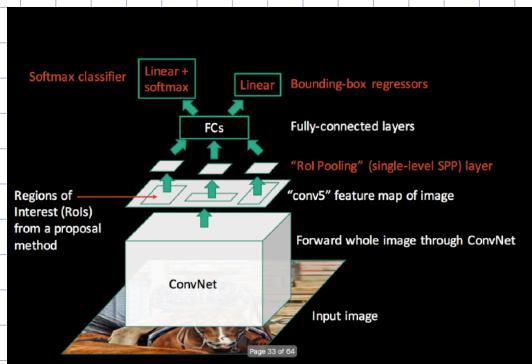
$$\frac{x^* - x}{w} = W_{c,w}^T \times_{\text{pools}} \quad \frac{y^* - y}{h} = W_{c,h}^T \times_{\text{pools}} \quad \ln \frac{w^*}{w} = \ln \frac{h^*}{h} = W_{c,w}^T \times_{\text{pools}}$$

where \times_{pools} are the features from pools layer of the network

CONS

- Ad hoc training objectives
 - fine tuning network with softmax classifier (log loss)
 - train post-hoc linear SVMs (hinge loss)
 - train post-hoc bounding-box regressor (squared loss)
- Training (3 days) and testing (47s × image) is slow.
- Takes a lot of disk space: need to store all precomputed CNN features for training the classifiers

Fast R-CNN

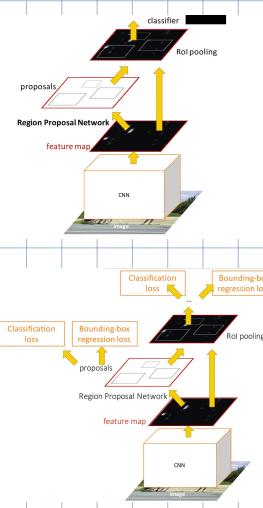


Region Proposal Networks (RPN)

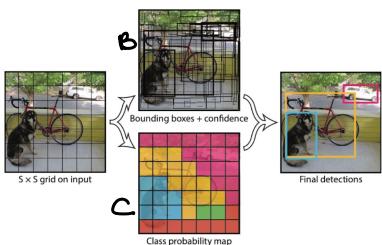
- Remove dependence on external region proposal algorithm.
- Instead, infer region proposals from some CNN
- => feature sharing
- => object detection in a single pass becomes possible

Faster R-CNN = Fast R-CNN + RPN

- Region proposal goes deeper
- Uses ROI (Region of Interest) pooling and anchors to extract region features.
- One network → 4 losses (joint training)



YOLO: You only look once



- YOLO divides the input image into an $S \times S$ grid. If the center of an object falls into a grid cell, that grid cell is responsible for detecting that object.
- Each grid cell predicts B bounding boxes and confidence scores for those boxes.
- Each bounding box consists of 5 predictions x, y, w, h and confidence.
- Each grid cell also predicts C conditional class probabilities.
- What's the output tensor size?

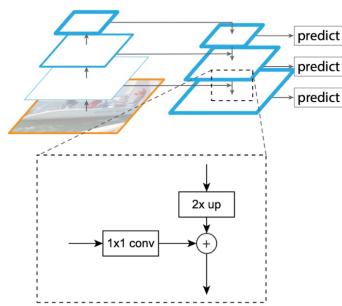
Feature pyramid Networks (FPN)

- A top-down architecture with skip connections. FPN model leverages it as a feature pyramid, with predictions made independently at all levels.
- With a coarser-resolution feature map, FPN upsamples the spatial resolution by a factor of 2 using nearest neighbor upsampling. The two feature maps are merged by addition.

FPN on Faster RCNN

- Assign anchors of a single scale to each level.
- There are 5 levels in total, where the last level feature is simply a stride two subsampling of the last layer of ResNet.
- Use anchors of multiple aspect ratios {1:2, 1:1, 2:1} at each level.
- Assign training labels to the anchors based on their Intersection-over-Union (IoU) ratios with ground-truth bounding boxes.
- An anchor is assigned a positive label if it has the highest IoU for a given ground truth box or an IoU over 0.7 with any ground-truth box, and a negative label if it has IoU lower than 0.3 for all ground-truth boxes.

- Find the box with greatest confidence
- Remove the other boxes with $\text{IoU} \geq \text{threshold}$
- If there are other boxes, go to 1.



LECTURE 9

Track a point

Optimize displacement for this energy, so the two pixels have the same color.

$$E(h) = [I_0(x+h) - I_1(x)]^2$$

$$E(h) \approx [I_0(x) + h I'_0(x) - I_1(x)]^2$$

$$\frac{\partial E}{\partial h} = 2I'_0(x) [I_0(x) + h I'_0(x) - I_1(x)]$$

$$\frac{\partial E}{\partial h} = 0 \rightarrow h \approx \frac{I_1(x) - I_0(x)}{I'_0(x)}$$

Problem 1: $I'_0(x) = 0$ if the image gradient is zero (the image region is flat)

- Problem 2:
- Motion to closest minimum has to be assumed
 - Indirect result: frame-rate should be faster than motion (λ wavelength)
 - Non-convex minimum may indicate multiple solutions

Solution:

- optical flow in motion estimation
- Harris detector

Track a bigger box

LUCAS-KANADE - Template tracking

Keep a template image to compare with each frame (only for small patches).

$$\begin{aligned} E(u, v) &= \sum_{x,y} [I(x+u, y+v) - T(x, y)]^2 \quad T(x, y) = \text{template} \\ &= \sum_x [I(x + \Delta x) - T(x, u)]^2 \quad W(x, p) = \text{optimal wrap} \\ &\Rightarrow \sum_x [I(W(x; p + \Delta p)) - T(x)]^2 \quad \text{Taylor expansion} \\ &= \sum_x [I(W(x; p)) + \nabla I \frac{\partial W}{\partial p} \Delta p - T(x)]^2 \end{aligned}$$

$$\Delta p = H^{-1} \sum_x [\nabla I \frac{\partial W}{\partial p}]^T [T(x) - I(W(x; p))] \quad \text{set derivative to 0}$$

$$\text{where } H = \sum_x [\nabla I \frac{\partial W}{\partial p}]^T [\nabla I \frac{\partial W}{\partial p}]$$

Warp I to obtain $I(W([x y]; P))$

Compute the error image $T(x) - I(W([x y]; P))$

Warp the gradient ∇I with $W([x y]; P)$

Evaluate $\frac{\partial W}{\partial p}$ at $([x y]; P)$ (Jacobian)

Compute steepest descent images $\nabla I \frac{\partial W}{\partial p}$

Compute Hessian matrix $\sum (\nabla I \frac{\partial W}{\partial p})^T (\nabla I \frac{\partial W}{\partial p})$

Compute $\sum (\nabla I \frac{\partial W}{\partial p})^T (T(x, y) - I(W([x, y]; P)))$

Compute ΔP

Update $P \leftarrow P + \Delta P$

PROS

- Beautiful framework
- It can handle different parameter space
- In can converge fast in a high-frame rate video

CONS

- Not robust to image noise or large displacement
- Some transformations are impossible to parameterize

Track by detection

Tracking by features

1. Detect key points (in reference and image)
2. Build feature descriptors
3. Match key points descriptors
4. Geometric verification (robust 3D pose calculation - RANSAC)

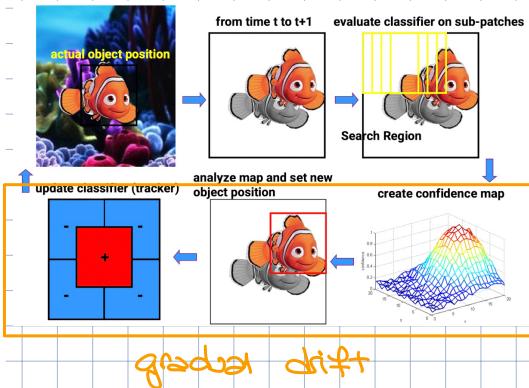
Acceleration:

- Tree structure (KD-tree)
- Approximate nearest neighbor
- Hashing
- Parallel implementation

Known object category

1. Detect objects independently in each frame
 - usually supervised learning (ConvNets)
2. Associate detections over time into tracks

Online learning

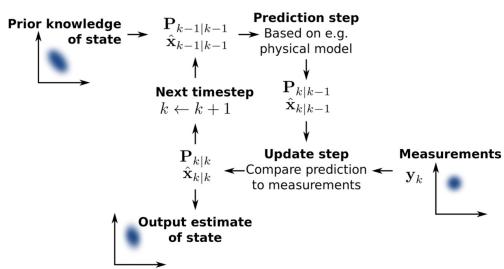


To avoid drift:

- Only thing we are sure about the object is its initial model
- we can "anchor" / correct our model with this information in order to help avoid drift

Motion

Kalman filter



$$\text{example: } \mathbf{x} = [u, v, s, r, \dot{u}, \dot{v}, \dot{s}]^T$$

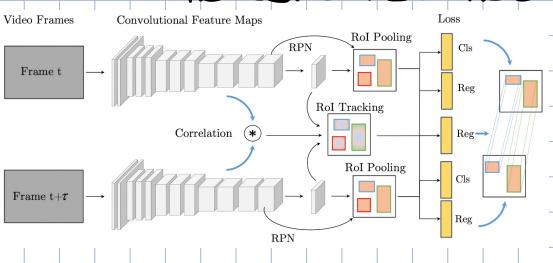
$\begin{cases} u, v = \text{bounding box coordinates} \\ s = \text{area of the bounding box} \\ r = \text{aspect ratio} \end{cases}$

Training: neural networks

→ network learns generic object tracking

Test: neural network (frozen weights)

→ network tracks novel objects (no fine-tuning)



Combining tracking and detection

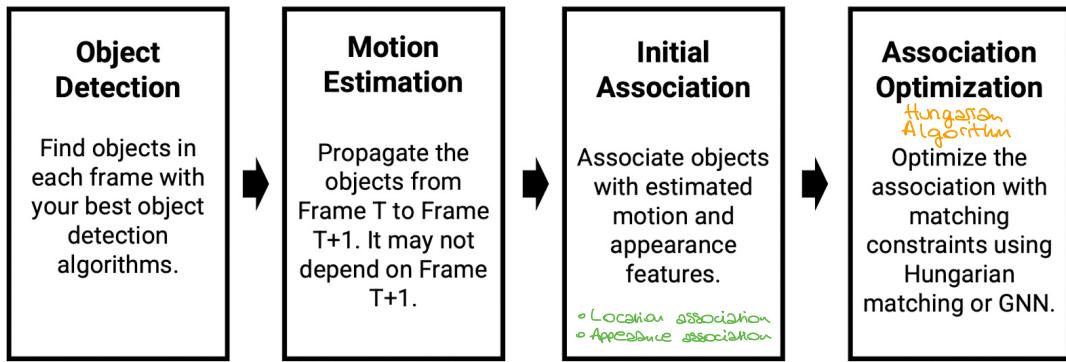
RoI pooling

1. Extract features with downsampling on input image (create a grid)
2. ROI is rounded to image grid
3. ROI is broken in grid.
4. MaxPool in each ROI cell.

RoI Align

1. Extract features with downsampling on input image (create a grid)
2. ROI is broken in grid.
3. ROI is broken in grid with same resolution of the grid
4. For each ROI cell, bilinear interpolate the features
5. MaxPool in each ROI cell.

Multiple object tracking



Why doesn't appearance provide enough information in current models?

- similar bounding boxes.
- misleading regions in the background.

LECTURE 10

Projective geometry

Point on a line: $\ell^T x = 0$

Line through two points: $\ell = x_1 \times x_2$

Point on a plane: $\pi^T x = 0$

Transformation of points: $x' = H x$

Transformation of lines: $\ell' = H^{-T} \ell$

Conic | locus of points m satisfying a homogeneous quadratic equation:

$$S(m) = m^T C m = 0$$

where C is a 3×3 symmetric matrix defined up to scale ($\# \text{dof} = 5$)

Inhomogeneous: $ax^2 + bx + cy^2 + dx + ey + f = 0$ ($\# \text{dof} = 5$)

$$C = \begin{bmatrix} a & b/2 & d/2 \\ b/2 & c & e/2 \\ d/2 & e/2 & f \end{bmatrix}$$

Dual conic: | locus of lines satisfying a homogeneous quadratic equation:
 $\ell^T C^* \ell = 0$

where C is a 3×3 symmetric matrix defined up to scale

Line tangent to a conic in point m : $\ell = C^T m = Cm$

Lines tangent to a conic: $\ell^T C^* \ell = 0$ if full rank $C^* = C^{-1}$

Transformations of a conic: $C' = H^{-T} C H^{-1}$ $C^* = H C^* H^T$ $(C')^* = (C^*)'$

Degenerate conic: | matrix C is not full rank

◦ $C = \ell m^T + m \ell^T$ (ℓ, m are two lines $x \in \ell, x \in m$)

◦ $C = \ell \ell^T$ (ℓ is a line $x \in \ell$)

For degenerate conics $(C^*)^* \neq C$

Quadratic | locus of all points M satisfying a homogeneous quadratic equation:

$$M^T Q M = 0$$

where Q is a 4×4 symmetric matrix defined up to scale

Dual quadratic: | locus of all planes π satisfying a homogeneous quadratic equation

$$\pi^T Q^* \pi = 0$$

where Q is a 4×4 symmetric matrix defined up to scale

Plane tangent to a quadric in point M : $\pi = Q^T M$

Planes tangent to a quadric $\pi^T Q^* \pi = 0$ if full rank $Q^* = Q^{-1}$

Transformation of a quadric: $Q' = T^{-T} Q T^{-1}$ $Q^{*'} = T^T Q T$ $(Q')^* = (Q^*)'$

Cross ratio: | invariant property under projective transformations.

Assume 4 points M_1, M_2, M_3, M_4 are collinear. They can be expressed as $M_i = M + \lambda_i M'$

$$\text{CR}\{M_1, M_2, M_3, M_4\} = \frac{\lambda_1 - \lambda_3}{\lambda_1 - \lambda_4} : \frac{\lambda_2 - \lambda_3}{\lambda_2 - \lambda_4}$$

Transformation	2D DoF	3D DoF	Example	INvariants
Projective	8	15		<ul style="list-style-type: none"> - concurrency - collinearity - order of contact - cross ratio
Affine	6	12		<ul style="list-style-type: none"> - line at infinity / plane at infinity - parallelism - ratio of lengths on lines - linear combination of vectors - ratio of areas
Similarity	4	7		<ul style="list-style-type: none"> - ratio of lengths - angles - circular points $\{(1, i, 0), (1, -i, 0)\} : \mathbb{I}$ - dual conic $F^*: (I^T I + J^T J)$ - absolute conic - absolute dual quadric
Euclidian	3	6		<ul style="list-style-type: none"> - lengths - areas

rotation = 1 dof rotation = 3 dof

Projectivity: invertible mapping from P^2 to itself such that three points x_1, x_2, x_3 lie on the same line if and only if $h(x_1), h(x_2), h(x_3)$ do.

Theorem: a mapping $h: P^2 \rightarrow P^2$ is a projectivity if and only if there exist a non-singular 3×3 matrix such that for any point in P^2 represented by a vector x it is true that $h(x) = Hx$.

projectivity = collineation = homography

Singular value decomposition:

method used to solve $\min \|Ax\|$ subject to $\|x\| = 1$

$$A = U \Sigma V^T \quad \begin{cases} U = \text{rotation matrix} \\ \Sigma = \text{scaling matrix (diagonal positive decreasing values)} \\ V^T = \text{rotation matrix} \end{cases}$$

solution $\rightarrow A = Vn$: last column of matrix V

Gold Standard algorithm

In P^2 , projectivity matrix is 3×3 up to scale $\Rightarrow 8$ dof.
Each measurement gives 2 equations \Rightarrow only 4 points needed.

If you have more, you cannot have exact solution because of the noise \Rightarrow need to find best according to cost function.

- o cost function is optimal for some assumptions
- o algorithm that minimizes the cost function.

Direct linear transformation (DLT)

$x_i' = Hx_i \rightarrow$ only if up to scale. But always ||
 $x_i' \times Hx_i = 0 \rightarrow$ cross product

$$x_i' = (x_i', u_i', w_i')^T \quad Hx_i = \begin{pmatrix} h_1 x_i \\ h_2 x_i \\ h_3 x_i \end{pmatrix} = \begin{pmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{pmatrix}$$

$$x_i' \times Hx_i = \begin{pmatrix} u_i' h_3^T x_i - w_i' h_2^T x_i \\ w_i' h_1^T x_i - x_i' h_3^T x_i \\ x_i' h_2^T x_i - u_i' h_1^T x_i \end{pmatrix} = \begin{pmatrix} 0 & -w_i' x_i & u_i' x_i \\ w_i' x_i & 0 & -x_i' x_i \\ -u_i' x_i & x_i' x_i & 0 \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} = A_i h = 0$$

Only 2 independent equations \rightarrow you can drop 3rd line.
 If you have more than 4 points, minimize $\|Ah\| \rightarrow 0$.

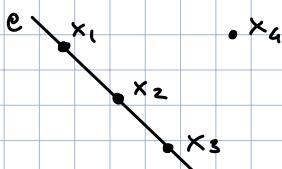
Objective

Given $n \geq 4$ 2D to 2D point correspondences $\{x_i \leftrightarrow x_i'\}$,
 determine the 2D homography matrix H such that $x_i' = Hx_i$

Algorithm

- (i) For each correspondence $x_i \leftrightarrow x_i'$ compute A_i . Usually only two first rows needed.
- (ii) Assemble n 2×9 matrices A_i into a single $2n \times 9$ matrix A
- (iii) Obtain SVD of A . Solution for h is last column of V
- (iv) Determine H from h

Degenerate configuration



Constraints: $x_i' \times Hx_i = 0$

Define: $H^* = x_4' e^T$ (H is rank $K-1$, not homography)

$$\text{Then: } \begin{cases} H^* x_i = x_4' e^T x_i = 0 & i=1,2,3 \\ H^* x_i = x_4' e^T x_4 = K x_4' & i=4 \text{ (maps to itself)} \end{cases}$$

If H^* is unique solution, no homography.

If there is further solution H , then $\alpha H^* + \beta H$ is solution

Cost functions

- Algebraic distance
- Geometric distance
- Reprojection error

Algebraic distance

$$e = Ah$$

$$d_{alg}(x_i, Hx_i)^2 = \|e_i\|^2$$

$$\sum_i d_{alg}(x_i', Hx_i)^2$$

$$d_{alg}(x_1, x_2)^2 = \omega_1^2 + \omega_2^2$$

where $\omega = (\omega_1, \omega_2, \omega_3)^T = x_1 \times x_2$

$$\text{Geometric distance: } \hat{H} = \underset{H}{\operatorname{argmin}} \sum_i d(x_i', Hx_i)^2 = d_{alg}(x_i', Hx_i) / (w_i' H w_i)$$

euclidean
distance

$$\text{Symmetric error: } \hat{H} = \underset{H}{\operatorname{argmin}} \sum_i d(x_i, H^{-1}x_i')^2 + d(x_i', Hx_i)^2$$

$$\text{Reprojection error: } (\hat{H}, \hat{x}_i, \hat{x}_i') = \underset{(H, x_i, x_i')}{\operatorname{argmin}} \sum_i d(x_i, \hat{x}_i)^2 + d(x_i', \hat{x}_i')^2$$

Mahalanobis error:
$$\sum_i \|x_i - \bar{x}_i\|_{\Sigma}^2 + \|x'_i - \bar{x}'_i\|_{\Sigma'}^2$$

where $\|x - \bar{x}\|_{\Sigma}^2 = (x - \bar{x})\Sigma^{-1}(x - \bar{x})$

measurement x
with covariance Σ

Objective

Given $n \geq 6$ 2D to 3D point correspondences $\{X_i \leftrightarrow x'_i\}$,
determine the Maximum Likelihood Estimation of P

Algorithm

(i) **Linear solution:**

- (a) Normalization: $\tilde{X}_i = UX_i$ $\tilde{x}_i = Tx_i$
- (b) DLT

(ii) **Minimization of geometric error:** using the linear estimate as a starting point minimize the geometric error:

$$\min_P \sum_i d(\tilde{x}_i, \tilde{P}\tilde{X}_i)^2$$

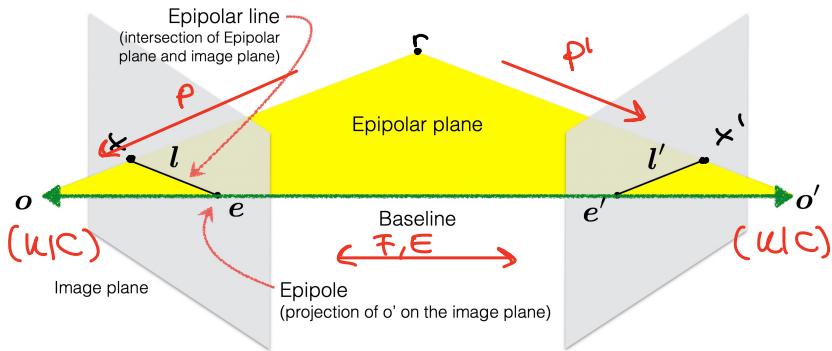
(iii) **Denormalization:** $P = T^{-1}\tilde{P}U$

Calibration

1. Canny edge detection
2. Straight line fitting to the detected edges
3. Intersecting the lines to obtain the image corners

LECTURE 11

Epipolar geometry



$$\begin{cases} p' = Fx \\ e = F^T x' \end{cases}$$

$$\begin{cases} e'^T F = 0 \\ F e = 0 \end{cases}$$

$$\begin{cases} x = p_r \\ x' = p'_r \end{cases}$$

every x is normalized!

Camera geometry

Given a set of corresponding image points $\{x_i \leftrightarrow x'_i\}$, what are the cameras P and P' for the two viewers?

$$(P, P') \rightarrow F \text{ unique} : \begin{cases} x = p_r \\ x' = p'_r \\ x' = Fx \end{cases} \quad \begin{cases} P = [I | 0] \\ P' = [M | m] \\ F = [m]_x M \end{cases}$$

$$F \rightarrow (P, P') \text{ NOT unique} : \begin{cases} \text{CANONICAL way} \\ P = [I | 0] \\ P' = [M | m] = [e']_x F + e'^T v^T | \lambda e' \end{cases} \quad F = [e']_x P' P'^T = \lambda F$$

FUNDAMENTAL MATRIX:

The fundamental matrix satisfies the condition that for any pair of corresponding points $x \leftrightarrow x'$ in two images

$$x'^T F x = 0$$

- F is a unique 3×3 rank 2 matrix
- if $(P, P') \rightarrow F$, then $(P', P) \rightarrow F^T$
- $e'^T F x = 0 \wedge x \Rightarrow e'^T F = 0 \quad e^T F^T x' = 0 \wedge x' \Rightarrow e^T F^T = 0$
- F has 4 dof (3x3-1 (homogeneous)-1 (rank 2))
- F is a correlation from a point to a line (not proper correlation, not invertible).
- General motion $F = [e']_x P' P'^T = [e']_x H$
- Pure translation $F = [e']_x$ (only 2 dof)

The fundamental matrix is a generalization of the essential matrix, where the assumption of Identity matrices is removed.

ESSENTIAL MATRIX:

The essential matrix operates on image points expressed in 2D coordinates expressed in the camera coordinate system:

- $Ee' = 0, E^T e = 0$
- E is singular
- E has two equal non-zero singular values.
- $E = [t]_x R$

Both E and F matrices completely describe the geometric relationship between corresponding points of a stereo pair of cameras. The only difference is that the former deals with **calibrated cameras**, while the latter with **uncalibrated cameras**.

$$F = K^{-T} E K^{-1} \Leftrightarrow E = K^T F K$$

PROJECTIVE RECONSTRUCTION THEOREM:

If F is same for (P, P') and (\tilde{P}, \tilde{P}') , then there exists a projective transformation H so that $\tilde{P} = H P$ and $\tilde{P}' = H P'$.

\Rightarrow If a set of point correspondences in two views determine the F uniquely, then the scene and the cameras may be reconstructed from these correspondences alone, and any two such reconstruction from these correspondences are projectively equivalent.

Canonical cameras given F : $P = [I|0]$, $P' = [e'|x]F + d\omega^T |\lambda e']$

COMPUTATION OF F :

- Linear 8-point
- Minimal 7-point
- Non-linear refinement
- Calibrated 5-point
- Calibrated + known vertical 3-point

8-POINT ALGORITHM

1. Normalize points ($x_{normalized} = K^{-1} \times un-normalized$)
2. Construct $M \times q$ matrix A ($M = \# \text{correspondences}$, $q = q_1, q_2, \dots, q_q$) $x^T F x = 0$
3. Find the SVD of $A^T A$
4. Entries of F are the elements of column of U corresponding to least singular value
5. Enforce rank 2 constraint on $F \rightarrow$ Rank K : Enforce all but the larger K values of SVD to be 0.
6. Un-normalize $F = K^{-T} F K^{-1}$

TRIANGULATION: given (x, x') , (P, P') , which 3D point match to x ?

$$\begin{cases} x = P r \\ x' = P' r \end{cases} \xrightarrow[\text{factor}]{\text{scale}} x \times P r = 0$$

cross product of two vectors
of same direction is zero

$$\begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \times \begin{bmatrix} P_1^T r \\ P_2^T r \\ P_3^T r \end{bmatrix} = \begin{bmatrix} u P_3^T r - P_2^T r \\ P_1^T r - x P_3^T r \\ K P_2^T x - u P_1^T r \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

only 2 independent equations

$$\begin{bmatrix} u P_3^T r - P_2^T r \\ P_1^T r - x P_3^T r \\ u' P_3^T r - P_2^T r \\ P_1^T r - x' P_3^T r \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} u P_3^T - P_2^T \\ P_1^T - x P_3^T \\ u' P_3^T - P_2^T \\ P_1^T - x' P_3^T \end{bmatrix} r = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$A x = 0$
SVD!

RECONSTRUCTION:

1. Compute the Fundamental Matrix F from point correspondences (8-point algorithm)
2. Compute the camera matrices P from the Fundamental matrix

$$P = [I|0] \text{ and } P' = [e'_x]F|e']$$

3. For each point correspondence, compute the point r in 3D space
(triangulation)

$$\text{DLT with } \begin{cases} x = Pr \\ x' = P'r \end{cases}$$

MULTI-VIEW GEOMETRY

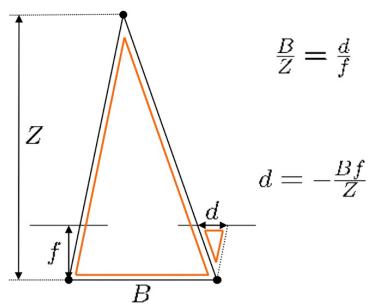
STRUCTURE FROM MOTION

recovering 3D structure of the scene from camera (or object) motion

LECTURE 13

STEREO VISION

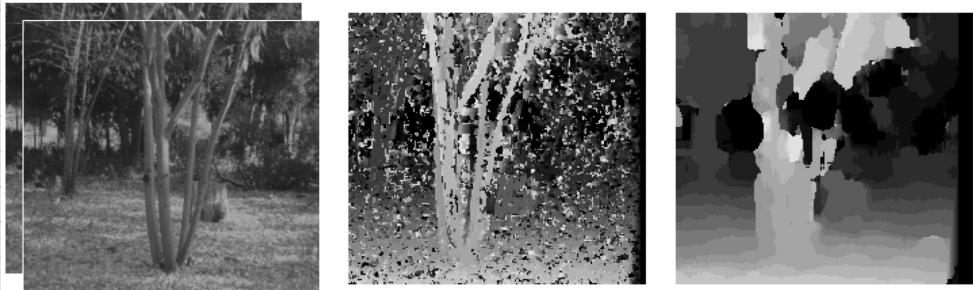
$$\text{disparity} = x - x' = d = -\frac{bf}{Z}$$



BASIC STEREO ALGORITHM

- For each epipolar line
 - For each pixel in the left image
 - compare with every pixel on the same epipolar line in the right image
 - pick pixel with minimum match cost

→ IMPROVEMENT: match windows e.g. Lukas-Kanade



SMALLER WINDOW:
more detail (+)
more noise (-)

LARGER WINDOW:
smoother disparity maps (+)
less detail (-)

STEREO IMAGE RECTIFICATION

The goal of stereo rectification is to warp the left and the right camera images such that their focal planes are coplanar and their intrinsic parameters are identical.

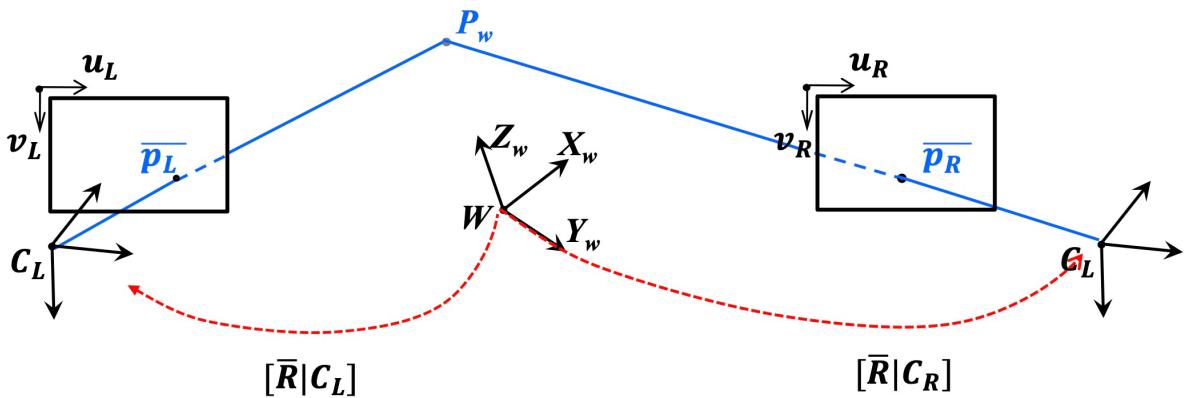
- reproject image planes onto common plane parallel to line between optical centers
- a homography (3×3 transform) applied to both input images
- pixel motion is horizontal after this transformation

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K [R|T] \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} \quad \xrightarrow{\text{camera point}} \quad \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K R^{-1} \left[\begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} - T \right] \quad \text{3D point}$$

Old Left camera New Left camera Old Right camera New Right camera

$$\lambda_L \begin{bmatrix} u_L \\ v_L \\ 1 \end{bmatrix} = K_L R_L^{-1} \left[\begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} - C_L \right] \rightarrow \bar{\lambda}_L \begin{bmatrix} \bar{u}_L \\ \bar{v}_L \\ 1 \end{bmatrix} = \bar{K} \bar{R}^{-1} \left[\begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} - C_L \right]$$

$$\lambda_R \begin{bmatrix} u_R \\ v_R \\ 1 \end{bmatrix} = K_R R_R^{-1} \left[\begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} - C_R \right] \rightarrow \bar{\lambda}_R \begin{bmatrix} \bar{u}_R \\ \bar{v}_R \\ 1 \end{bmatrix} = \bar{K} \bar{R}^{-1} \left[\begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} - C_R \right]$$



By solving for $\begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$ for each camera, we can compute the homography (or warping) that needs to be applied to rectify each camera image.

$$\bar{\lambda}_L \begin{bmatrix} \bar{u}_L \\ \bar{v}_L \\ 1 \end{bmatrix} = \lambda_L \bar{K} \bar{R}^{-1} R_L K_L^{-1} \begin{bmatrix} u_L \\ v_L \\ 1 \end{bmatrix}$$

Homography
Left Camera

$$\bar{\lambda}_R \begin{bmatrix} \bar{u}_R \\ \bar{v}_R \\ 1 \end{bmatrix} = \lambda_R \bar{K} \bar{R}^{-1} R_R K_R^{-1} \begin{bmatrix} u_R \\ v_R \\ 1 \end{bmatrix}$$

Homography
Right Camera

How to choose \bar{K} , \bar{R} ?

• $\bar{K} = (K_L + K_R) / 2$

• $\bar{R} = [\bar{r}_1, \bar{r}_2, \bar{r}_3]$ where

$$\begin{cases} \bar{r}_1 = \frac{C_2 - C_1}{\|C_2 - C_1\|} \\ \bar{r}_2 = r_3 \times \bar{r}_1 \\ r_3 = \bar{r}_1 \times \bar{r}_2 \end{cases}$$