

# Learning Object Categories from Google's Image Search

Fergus, Perona, Fei-Fei, Zisserman

Presented by  
Sudheendra

## Introduction

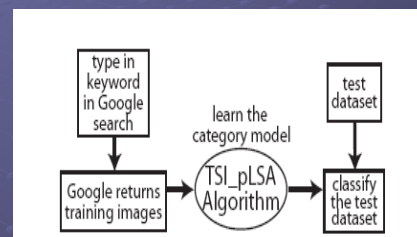
- **Contributions**
  - TSI-pLSA – a translation and scale invariant pLSA model
  - Unsupervised learning on training set collected from Google image search and therefore unlabeled
- **Related work**
  - Discovering objects and their locations in images
    - Sivic, Russell, Efros, Zisserman, Freeman
    - pLSA for object category recognition and segmentation
  - A visual category filter for Google images
    - Fergus, Perona, Zisserman
    - Reranking of Google images by learning a model

## Introduction

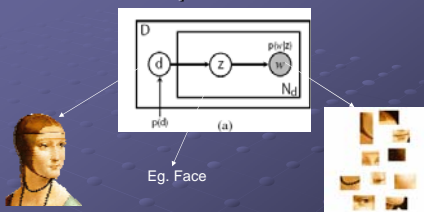
- **Main challenge**
  - Noisy training data : Less than 15% of the images returned by Google are related to the keyword
  - Large variations in scale, position and pose
- **Idea**
  - Build pLSA model with a number of topics
  - Visual words of an image will fall under a common topic
  - Visual words of positive examples will be similar
  - Find this topic using a validation set of less noisy data



## Overview



## pLSA

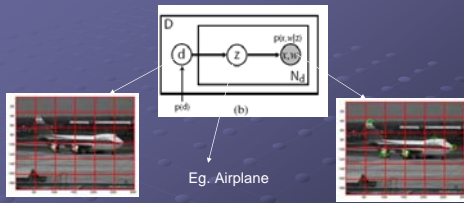


- **Generative model**
  - Choose document (image)  $d$  with probability  $P(d)$
  - Choose topic  $z$  with probability  $P(z|d)$
  - Choose word with probability  $P(w|z)$
  - Thus,  $P(d, w) = P(d)P(z|d)P(w|z)$ , where
  - $P(w|d) = \sum_{z \in Z} P(z|d)P(w|z)$

## pLSA

- **Learning**
  - Using EM
  - E-step: estimate  $P(z | d, w)$ 
    - Associates  $z$  with the image and feature
  - M-step: update  $P(z | d)$  and  $P(w | z)$ 
    - Visual words from an image tend to fall under the same topic
- **Recognition**
  - Fix  $P(w | z)$  and estimate  $P(z | d)$  using EM
- **Drawbacks**
  - Spatial information is not used
  - Multiple instances of a category cannot be captured

## ABS-pLSA



### Generative model summary

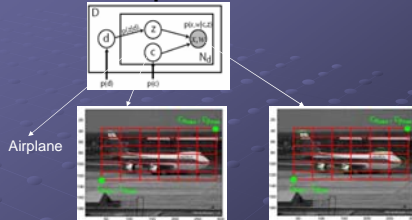
- Choose document  $d$  with probability  $P(d)$
- Choose topic  $z$  with probability  $P(z|d)$
- Choose word  $w$ , location  $x$  with probability  $P(w, x|z)$

Thus,  $P(w, x, d) = \sum_{z=1}^K P(w, x|z) P(z|d) P(d)$

## ABS-pLSA

- Quantize the image into  $X$  bins
- Include spatial location with word to produce topic variable
- EM steps similar to pLSA
- Drawback
  - Uses absolute location of feature
  - Not translation or scale invariant

## TSI-pLSA



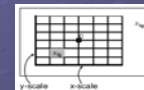
### Generative model summary

- Choose document  $d$  with probability  $P(d)$
- Choose topic  $z$  with probability  $P(z|d)$
- Choose word  $w$ , relative position  $x$  over all possible values of centroid  $c$  with probability  $P(w, x|z)$  calculated as

$$P(w, x|z) = \sum_{c=1}^C P(w, x, c|z) = \sum_{c=1}^C P(w, x|c, z) P(c)$$

## TSI-pLSA

- Location of feature calculated with respect to object centroid



- $x$ -scale and  $y$ -scale along with the centroid specify the bounding box of the object
- A grid of  $X_{ij}$  locations within the bounding box and one background bin for feature location
- Object centroid and scale captured in 4-vector latent variable  $c$
- Marginalizing over the entire range of  $c$  is not feasible
  - Small set of  $c$  values estimated during learning and recognition

## TSI-pLSA

### Learning

- Estimating  $c$  values
  - Standard pLSA run on the training set
  - $k = (1 \dots K)$  gaussians fitted to the locations of features weighted by  $p(w|z)$  to obtain  $k$  values of  $c$  (centroid = mean, scale = variance)
  - Captures multiple instances of objects in image

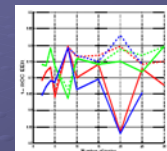


- EM: as before with  $P(w, x|z)$  estimated by marginalizing over above found  $c$  values

### Recognition

- Estimating  $c$  values
  - Similar to above,  $\{P(w|z)\} = \sum_{c=1}^C P(w, x|c, z) P(c)$  used to weight gaussians
- EM: Lock  $P(w|z)$  and iterate to find  $P(z|d)$  summed over the found  $c$  values

## Some Issues



- Selecting the final classifier
  - Visual words of positive examples should belong to a common topic
  - Validation set will perform best under this common topic
- Selecting the number of topics  $Z$ 
  - Chosen empirically
  - Roc vs number of topics plotted for best topic under validation set and best topic under test set

## Datasets

- Training
  - Google dataset
    - Images automatically downloaded from Google image search using the category name
    - Validation set – first five images from image search in 7 different languages
  - Other
    - Manually gathered frames from Caltech and Pascal datasets
- Testing
  - Manually gathered frames from Caltech and Pascal datasets

## Parameters

- 700 regions per image using 4 different region detectors
  - Because the method requires large number of data for parameter estimation
- SIFT descriptor of 72 dimensions
  - Larger histogram bins more appropriate for object categorization
- K-means clustering with  $k=350$  to obtain 350 visual words
- Number of grid positions  $X_{ig} = 37$
- Number of topics  $Z = 8$

## Experiments and Results

### Experiment 2 (standard datasets)

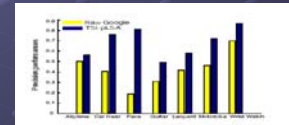
- Training images from Google image search
- 8 topics and best topic chosen using Google validation set
- TSI-pLSA performs better than the other methods in all categories except Guitar and background
- ABS-pLSA performs better than TSI-pLSA in rotation and rotation-invariant learning
- 6 topics and best chosen using performance on foreground only images
- TSI-pLSA performs better than the other methods

Category	pLSA		ABS		TSI		pLSA		ABS		TSI	
	Prop.	Prop.	Prop.	Prop.	Google	Google	Prop.	Prop.	Google	Google	Prop.	Prop.
Aeroplane	17.7	13.2	4.7	24.7	17.2	18.5	17.7	13.2	4.7	24.7	17.2	18.5
Car Rear	2.0	6.2	0.9	21.0	13.5	16.0	2.0	6.2	0.9	21.0	13.5	16.0
Face	22.1	11.5	17.0	20.3	56.4	20.7	22.1	11.5	17.0	20.3	56.4	20.7
Guitar	9.3	10.0	14.4	17.6	62.0	31.8	9.3	10.0	14.4	17.6	62.0	31.8
Handbag	12.0	12.9	11.0	15.0	16.0	13.0	12.0	12.9	11.0	15.0	16.0	13.0
Motorbike	19.0	6.0	7.0	15.2	18.5	6.2	19.0	6.0	7.0	15.2	18.5	6.2
Object watch	21.6	7.7	15.5	21.0	20.5	19.9	21.6	7.7	15.5	21.0	20.5	19.9
TSI-pLSA	31.9	33.0	25.5	-	-	-	31.9	33.0	25.5	-	-	-
PASCAL Motorbike	33.7	30.2	25.7	-	-	-	33.7	30.2	25.7	-	-	-

## Experiments and Results

- Experiment 3
  - Comparison with other supervised methods
  - TSI-pLSA is slightly worse than the other methods but it is unsupervised
- Experiment 4
  - Improving Google's Image search
  - Best topic from 8 topics trained on raw Google data

Dataset	TSI-pLSA	[10]	[11]	[15]
Exp. B	img. labels	img. labels	img. labels	Segmented
ImageNet	25.87 (0.062)	-	-	34.2 (0.18)
PASCAL Motorbike	25.77 (0.249)	-	-	31.77 (0.34)
Exp. C	None	img. labels	img. labels	Segmented
Aeroplane	15.8	9.9	11.1	-
Car Rear	16.0	9.7	9.9	6.1
Face	20.7	3.6	6.5	-
Leopard	13.0	10.0	-	-
Motorbike	6.2	6.7	7.8	6.0



## Conclusions

- All three methods work on unlabeled Google dataset and automatically collected validation set and TSI-pLSA performs best
- TSI-pLSA identifies multiple instances of objects in images
- Can be used to rank images returned by Google
- TSI-pLSA performs badly when objects are rotated