# Bilinear Models for Fine-grained Visual Recognition

Akilesh B, CS13B1042

Indian Institute of Technology, Hyderabad

April 12, 2016

Bilinear Models for Fine-grained Visual Recognition

Introduction

Related VVork

Bilinear model

model

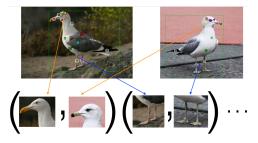


- The visual differences between the categories are small and can be easily overwhelmed by those caused by factors such as pose, viewpoint, or location of the object in the image.
- ► For example, the inter-category variation between Ringed-beak gull and a California gull due to the differences in the pattern on their beaks is significantly smaller than the inter-category variation on a popular fine-grained recognition dataset for birds.

Related Work

model

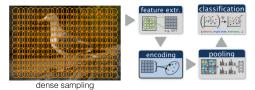




- ► Localize various parts of the object and model the appearance conditioned on their detected locations.
- The parts are manually defined and part detectors are trained in a supervised manner.
- ▶ It factors out variation due to pose, viewpoint and location.

Related Work

linear mode



- Image as a collection of patches [bag-of-visual words, Csurka et al 04].
- Orderless pooling and no explicit modeling of pose or viewpoint.
- ▶ Invariances due to: choice of features (eg. SIFT is robust to lighting changes). There is encoding + pooling + classification.
- ► Eg. Fisher Vectors work remarkably well for fine-grained tasks.

Related Work

silinear mod

model

- ▶ Part-based models Have the best recognition accuracy on many fine-grained recognition datasets (eg. birds).
- ▶ It is relatively slow since it involves part detection.
- Needs part annotations for training. This can be time consuming and may require expert knowledge. Parts may be hard to define them for some categories.
- Texture models Easy to deploy since they only need image labels for training. They have very fast CPU implementations.
- ► However, they have lower recognition accuracy. The pipelined procedure (features → encoding → classification) can be sub-optimal. For instance, the feature extractors are not learned.
- ► Can we get the best of both?

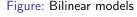
## Bilinear models for classification

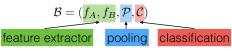


Introduction

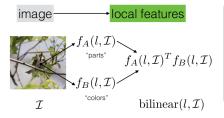
Related W

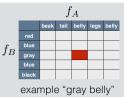
Bilinear models



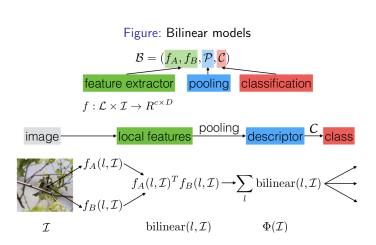


$$f: \mathcal{L} \times \mathcal{I} \to R^{c \times D}$$





## Bilinear models for classification



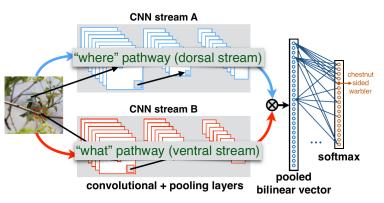
Bilinear Models for Fine-grained Visual Recognition

introduction

Ittilated Work

Bilinear models

Figure: Bilinear CNN model



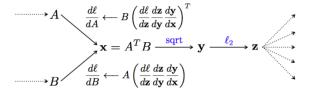
▶ Generalization: decouple  $f_A$  and  $f_B$  by using separate feature functions

Introduction

Related Work

ilinear model

Figure: Bilinear CNN model training



- Back-propagation though the bilinear layer is easy.
- It allows end-to-end training.

There are two normalization layers in the architecture

► Square-root normalization

$$y = sign(x)\sqrt{x} \tag{1}$$

▶ l₂ normalization

$$z = \frac{y}{\|y\|}$$

▶ Both these improve performance. > (□) >

Introduction

Related Work

illicai illouci

- ► VGG-M (5 convolutional layers + 2 fully connected layers)
- VGG-D (13 convolutional layers + 2 fully connected layers)

Trained using image labels only (no part or bounding-box annotations).

#### Dataset

▶ CUB-200-2011 dataset. Birds+box uses bounding boxes at training and test time. It has 200 categories of birds, 11788 images. It contains 15 part locations, 312 binary attributes and 1 bounding box.

Introduction

Related v

ilinear mode



## Classification demo

Figure: Test image



#### Figure: Classification Demo

```
>> bird_test_aki
0.11s to load imdb.
2.73s to load models into memory.
Top 5 prediction for test_image.jpg:
064.Ring_billed_Gull
059.California_Gull
147.Least_Tern
062.Herring_Gull
060.Glaucous_winged_Gull
2.21s to make predictions [GPU=0]
```

#### Bilinear Models for Fine-grained Visual Recognition

Introduction

Related vvol

illiear model

## Results: Birds classification

Bilinear Models for Fine-grained Visual Recognition

Introduction

Related Work

Dillical Illoac

Bilinear CNN model

- ▶ Per-image accuracy on CUB 200-2011 dataset.
- Provided with only the image at test time.

Method	w/o ft	w/ ft
B-CNN [M,M]	71.6	77.4
B-CNN [M,D]	79.6	83.2
B-CNN [D,D]	79.5	82.8

Table: Fine-grained classification results

- ► In their extension work, they propose one-to-many face recognition using Bilinear CNNs. They perform the experiments on Face-Scrub and IJB-A Train dataset.
- ► They use a linear SVM classifier learned for each person in gallery and max-pooling features or classifier scores to aggregate multiple media.
- ► Finally, they demonstrate how a standard CNN pre-trained on a large face database (say VGG-Face model) can be converted into a B-CNN without any additional feature training.

Related VVoi

\_\_\_\_\_\_



## Most confused birds

### Figure: Most confused birds









Caspian\_Tern



Acadian\_Flycatcher



Brandt\_Cormorant



Common\_Raven



Great\_Grey\_Shrike



Elegant\_Tern



Yellow\_bellied\_Flycatcher



Pelagic\_Cormorant



Bilinear Models for Fine-grained Visual Recognition

Introduction

Related \

Bilinear model

model

Tsung-Yu Lin, Aruni RoyChowdhury, and Subhransu Maji: Bilinear CNNs for Fine-grained Visual Recognition, ICCV 2015

Aruni RoyChowdhury, Tsung-Yu Li, and Subhransu Maji: One to many face recognition using Bilinear CNNs, WACV 2016