

000
001
002
003
004
005
006
007
008
009
010
011

Abstract

Recognition of social styles of people are an interesting but not yet *popular* topic. In this paper, we explore the use of pre-trained convolutional neural network(CNN) features for social group recognition. A social group recognition framework is proposed. Pre-trained CNN is fine-tuned, and CNN features are used for both individual person images and global scene images. Our model shows promising results on urban tribes dataset, with 71.23% accuracy, which significantly outperforms previous results of 46%. We also find that there is a correlation between the probability of an image being in Imagenet classes and social group classes, and that better-recognized categories have more highly-correlated ImageNet categories. This gives us insight into the features extracted from pre-trained CNNs.

1. Introduction

In the past few years, there have been impressive progress in understanding semantic meaning of images, such as object recognition, scene recognition, and object detection. The power of Convolutional Neural Networks (CNNs) has especially noticeable. However, the analysis the social features of images of groups of people has not attracted a great deal of research. Current search algorithms fail to capture information of personal styles or social characteristics of groups of people, but retrieve images with similar global appearance [1]. The analysis of groups of people is difficult in that the group categories are semantically ambiguous, and have high intra-class variance.

Recognition of groups of people from a social perspective provides many potential applications. With more accurate group searching results, more accurate recommendations can be made in social networks, and more relevant advertisement for particular groups of people benefits both consumers and sellers.

Kwak *et al.* studied this problem of group recognition ([1], [2]). They created an urban tribe dataset consisting of 11 categories, with about 100 labeled per class. They proposed a group recognition pipeline. Rather than classifying

Urban Tribe

Anonymous WACV submission

Paper ID ****

isolated individuals in the group images, they focused on group features and models.

CNN architecture has been proved to achieve outstanding results in various computer vision tasks, and it's argued that deep architectures in CNN can capture visual features of different semantic level in the hidden units ([3]). Recently, features learnt for large scale recognition tasks with large amounts of training data have been used for new tasks, and the features outperform many conventional features in many tasks([3], [4]). This indicates that generic visual features may be obtained from pre-trained CNN model.

not sure about this paragraph The most frequently used pre-trained model was proposed in [5] trained with ImageNet dataset¹. And the new tasks for which CNN features are used range from object recognition, scene recognition, to subcategory recognition. However, the recognition of style hasn't been researched yet. The social groups of individuals require the recognition of people's social characteristics, which is a problem of style, instead of category. Moreover, the ImageNet classes used for pre-training cover few human images, making the problem more interesting.

In this paper, we investigate the generalization ability of pre-trained CNN features to social group recognition. We propose a CNN feature based architecture for social group recognition. Our model takes in both individual features and global scene features fine-tuned from CNN pre-trained weights. Our result shows a boost of performance from the previous classification method provided by [1]. We show that both individual information and global scene information contribute to a social group's characteristics, and that different feature extraction schemes for individual and global information is necessary. We also show the role of adapting generic pre-trained CNN features to social group images with fine-tuning.

We further investigate why features extracted from pre-trained CNN are useful for the urban tribe recognition task. For an input image, there is a correlation between the probability of it being in ImageNet classes and being in urban tribes classes. Moreover, the degree of correlation is related to the recognition rate of different urban tribes classes

¹<http://image-net.org/challenges/LSVRC/2012/browse-synsets>

108 - better-recognized categories have more highly-correlated
109 ImageNet categories. This may indicate that the “generic”
110 features extracted by pre-trained CNN networks are not so
111 generic. The network is trained to separate the ImageNet
112 classes most, and if we use the features for a new classifica-
113 tion task, the performance of the task is related to how well
114 the new classes can be “mapped” to the ImageNet classes.
115 However, the actual relationship between the two types of
116 categories is still mysterious in most cases.
117

2. Related Work

120 Convolutional Neural Networks(CNNs) with back-
121 propogation were introduced around 1990’s by LeCun *et al.*
122 [6]. Since then, CNNs have shown successful results on var-
123 ious computer vision tasks, such as hand-written digit clas-
124 sification ([7],[8]), ImageNet Large Scale Visual Recog-
125 nition Challenge ([5], [9]), object detection ([9], [10], [11]),
126 object localization ([10]), and many modified architecture
127 of CNNs have been proposed([8], [12], [13]).

128 Recently, many researchers have shown the utility of
129 generalization of pre-trained CNN features on large dataset
130 such as ImageNet. Krizhevsky *et al.* showed excellent gen-
131 eralization of the pre-trained CNN features[3]. They kept
132 all the layers of ImageNet-trained model using fixed pre-
133 trained features except for the last softmax classifier, and
134 achieved best results on Caltech-101 and Caltech-256. Don-
135 ahue *et al.* used different layers of pre-trained CNN net-
136 work as features and trained simple classifiers such as SVM
137 and Logistic Regression, and outperformed the state-of-the-
138 art on several vision challenges such as scene recognition
139 and domain adaptation[4] . Karayev *et al.* compared dif-
140 ferent approaches for photographic style recognition, and
141 pre-trained CNN features generally obtain best result [14].

142 Kwak *et al.* created an urban tribe dataset consisting of
143 11 classes[1]. The classes are defined from social group
144 labels provided by Wikipedia. They selected the eight most
145 popular categories from their list of subculture, and added
146 three other classes corresponding to typical social venues
147 in addition. For each class, images of groups of people
148 were discovered with different search engines, and a broad
149 range of images for each class were collected. Kwak *et al.*
150 also provided a group description and several classification
151 methods([1], [2]). Group description consists of person de-
152 scriptors and global group descriptors: Six part of person is
153 detected, and a set of predefined descriptors are computed
154 for each part, including ratio of skin pixels, color informa-
155 tion like RGB histograms, and HoG features; Global de-
156 scriptors use a both low level and high level descriptors to
157 describe the context and group properties of the image. Low
158 lever features include color information, Gist, HoG and
159 ratio of pixels of person, and high level descriptors includs
160 proximity of persons, alignment or pose of the group, and
161 scene layout of individuals. Two options of classification

162 methods are provided: bag of parts-based classification and
163 SVM-based classification.
164

165 Categorizing the social groups of individuals belongs to
166 fine-grained classification task which recently draws more
167 interest of the computer vision community. It aims at giving
168 the fine-grained categories in a certain class. Fine-grained
169 classification is more difficult than conventional classifica-
170 tion tasks, because the categories are semantically as well as
171 visually similar, and are even challenging for humans. The
172 Fine-grained Challenge 2013 (FGComp) provided the data
173 in several categories including aircraft, birds, cars, dogs,
174 and shoes. [15] achieved the best result using classifier
175 based on fisher vectors. However, CNN based methods us-
176 ing [16] or [4] gave inferior results, especially when the
177 bounding box of test data is unknown.
178

179 There is some research in analyzing social groups of peo-
180 ple. [17] showed the visual structure of a group helps under-
181 standing events. [18] showed social relationships modeling
182 helps people recognition. [19] used both local and global
183 factors for group level expression analysis.
184

3. Methods

185 This section describes the urban tribe dataset and elabo-
186 rates on the model architecture.
187

3.1. Urban tribes dataset

188 Urban tribes are groups of people who have similar vi-
189 sual appearances, personal style and ideals. The urban
190 tribes dataset consists of 11 different categories: *biker*,
191 *country*, *goth*, *heavy-metal*, *hip-hop*, *hipster*, *raver*, *surfer*,
192 *club*, *formal*, *casual/pub*, with an average of 105 images
193 from each category.
194

195 Unlike conventional visual classification problems, ur-
196 ban tribe categories are more ambiguous and subjective.
197 Also, each class contains a broad range of scenarios. The
198 high intra-class variation of the urban tribe dataset makes
199 the classification task challenging. The urban tribe dataset
200 also has some interesting properties. The number of peo-
201 ple in each urban tribe image varies. Members in one tribe
202 often have similar visual styles, including their clothes,
203 accessories, and even demeanor. For example, surfers pos-
204 sibly carry surfboards, and the goth often have dark attire,
205 makeup and hair. The environment they are in also con-
206 tributes to each tribe characteristics: pictures of country
207 tribes are more likely to be taken outdoors with grassland,
208 while pictures of clubbers are often photographed in clubs
209 with dim lightings.
210

3.2. Classification hierarchy

211 To utilize the properties of urban tribes fully, our fea-
212 ture vector consists of both elements: individual features
213 and environmental features. For each feature type, we use
214

216 a similar extraction strategy. Individual features and environmental features are hierarchically combined to form the final decision function. The network hierarchy is shown in
 217 Figure 1.
 218

219 For each group image, we represent the group G as the
 220 combination of a set of people and the environment of the
 221 scene. To give the prediction of class C , the individual fea-
 222 ture vectors and scene feature vectors are extracted sepa-
 223 rately.
 224

225 For the individual feature vectors, first, individual can-
 226 didate person images are detected with a poselet based
 227 person detection algorithm. The candidate person images
 228 $H = \{H_1, H_2, \dots, H_p\}$ are used as a whole instead of a set
 229 of body part bounding boxes. Each candidate person is re-
 230 sized to 256×256 , and ten 227×227 patches $\{h_{ij}\}, i \in$
 231 $\{1, 2, \dots, p\}, j = 1, 2, \dots, 10$ are extracted (patches from
 232 four corner and the center, and their horizontal reflections).
 233

234 Each Individual image patch h_{ij} then passes through
 235 the Convolutional Neural Network for person images
 236 $\text{CNN}_{\text{Person}}$, generating activations from the 6th and 7th
 237 hidden layer. The activations from 6th and 7th layer are
 238 both 4096 dimensioned. They are concatenated to form an
 239 8192-dimensional vector f_{ij} , where $i \in \{1, 2, \dots, p\}, j \in$
 240 $\{1, 2, \dots, 10\}$.
 241

242 The feature vectors are then fed into a multi-class
 243 $\text{SVM}_{\text{Person}}$. We use LIBLINEAR[20] to train the SVM
 244 on individual patches, and to estimate probabilities for
 245 each category given individual patch h_{ij} : $\text{Pr}_{ij}(C|h_{ij}), C \in$
 246 $\{1, 2, \dots, c\}$, where c is the number of classes in urban tribe
 247 dataset. The individual patches h_{ij} in one group image are
 248 usually highly correlated. Therefore, in order to obtain a
 249 reliable probability estimate from the noisy yet correlated set
 250 of probabilities Pr_{ij} , a simple but effective average pooling
 251 is performed to Pr_{ij} :

$$\text{Pr}_{\text{People}}(C|H_1, \dots, H_p) = \frac{1}{10p} \sum_{i,j} \text{Pr}_{i,j}(C|h_{ij}) \quad (1)$$

252 $\text{Pr}_{\text{People}}(C|H_1, \dots, H_p)$ is the probability estimate of class
 253 C given the set of people candidate images H .
 254

255 On the other hand, the entire environment in the scene
 256 image, denoted by S , is directly utilized for probability es-
 257 timation. The procedure to generate probability estimate of
 258 class C given the environment as a whole $\text{Pr}_{\text{Scene}}(C|S)$ is
 259 similar with that for $\text{Pr}_{\text{People}}(C|H)$. The difference is that
 260 the input of Convolutional Network is 227×227 patches
 261 extracted from the entire scene image, and the fine-tuned
 262 Convolutional network: $\text{CNN}_{\text{Scene}}$ and SVM: $\text{SVM}_{\text{Scene}}$
 263 are trained with the training set of entire scene images. Sev-
 264 eral different strategies to extract patches from scene images
 265 and corresponding Convolutional Neural Network architec-
 266 tures are explained in Section 3.3.
 267

268 Therefore, the probability estimate of a class C given
 269

270 observation of scene S is:
 271

$$\text{Pr}_{\text{Scene}}(C|S) = \frac{1}{K} \sum_{k=1}^K \text{Pr}_k(C|s_k) \quad (2)$$

272 where K is the number of scene patches extracted from one
 273 group image, and this number varies with different patch
 274 extraction strategies. s_k is the k th scene patches. $\text{Pr}_k(C|s_k)$
 275 is the probability for class C given k th scene patch. Average
 276 pooling is still used here, because the assumption of high
 277 correlation in patches holds.
 278

279 Now we have the estimates of two kinds of conditional
 280 probability $\text{Pr}_{\text{People}}(C|H)$ and $\text{Pr}_{\text{Scene}}(C|S)$. We make a
 281 strong assumption that the two types of features are inde-
 282 pendent, and that the prior probability distribution of the
 283 urban tribes $\text{Pr}(C)$ is a uniform distribution. The classifica-
 284 tion problem can be expressed as maximizing the objective
 285 function:
 286

$$L = \arg \max_{i=1, \dots, c} \text{Pr}(C = i|G) \quad (3)$$

287 where
 288

$$\begin{aligned} \text{Pr}(C = i|G) &= \text{Pr}(C = i|H, S) \\ &= \frac{\text{Pr}_{\text{People}}(C = i|H_1, \dots, H_p) \cdot \text{Pr}_{\text{Scene}}(C = i|S)}{\text{Pr}(C = i)} \\ &\propto \text{Pr}_{\text{People}}(C = i|H_1, \dots, H_p) \cdot \text{Pr}_{\text{Scene}}(C = i|S) \end{aligned} \quad (4)$$

289 and L is the predicted label for the group image.
 290

291 3.3. Convolutional network feature extraction

292 It is shown in many experiments that a set of weights of
 293 convolutional network trained from ImageNet can generate
 294 a set of generic visual features.
 295

296 Following [16]'s work, we use the network framework
 297 called Caffe. The network architecture is described in [5],
 298 which won the ImageNet Large Scale Visual Recognition
 299 Challenge 2012. We take the activations from the 6th and
 300 7th hidden layer of the convolutional neural network, which
 301 are two fully connected layers before the class prediction
 302 layer. We also take the activations from 6th or 7th layer
 303 alone as comparison. We choose these two layers, because
 304 as the layers ascend, the features extracted show increasing
 305 invariance and semantic meaning.
 306

307 We use the pre-trained set of weights of the network re-
 308 leased by Caffe as the initial parameters of our network. The
 309 pre-trained model was trained on ImageNet ILSVRC-2012,
 310 and all images are first resized to 256×256 before they can
 311 be used as inputs to the network.
 312

313 3.3.1 Pre-processing of the dataset

314 The urban tribe dataset is a relatively small dataset, and both
 315 people candidate crops and scene images are of various res-
 316

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

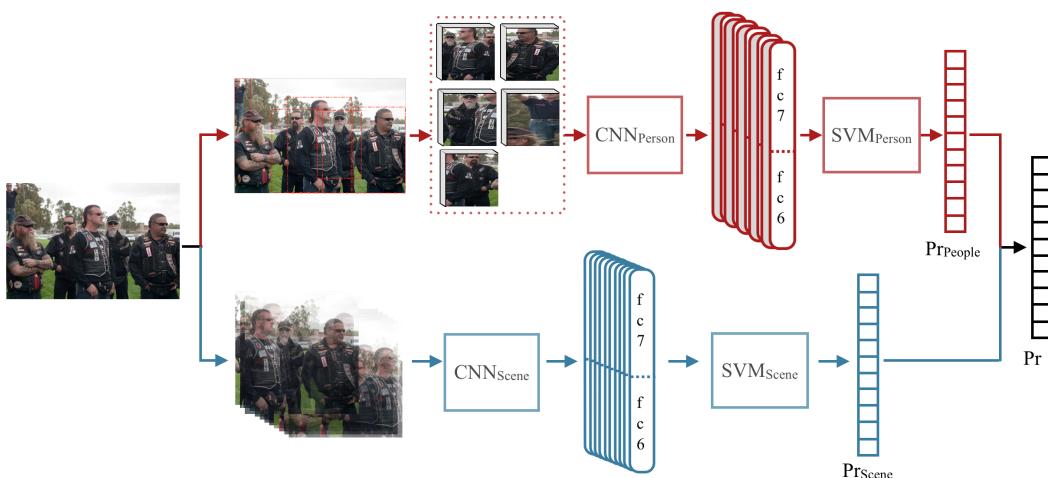


Figure 1: Architecture of classification algorithm using Nets_{SDense} which is introduced in Section 3.3.3. The upper half estimates the probability given people candidate images, and the lower half estimates the probability given the entire scene. Dense crop CNN_{Scene} and distorted crop CNN_{Person} are used, as shown in Section 3.3.1.

olution. Our convolutional neural networks require constant input image size of 227×227 , so pre-processing of the dataset is necessary.

There are several strategies to make one image compatible with the CNN:

1. *Distort cropping*: As in [5], resize the image to a fixed resolution of 256×256 , and crop five 227×227 patches (from four corner and the center) and their horizontal reflections to generate ten patches from one single image. This way, the aspect ratio of the original images are lost, but for each crop, the portion it takes from the original image is fixed, so that the amount of information all the crops have is relatively stable.
2. *Sparse cropping*: Keep the aspect ratio of the original image, resize the shorter side to 256, and then crop five 227×227 patches and their horizontal reflections as mentioned. This method avoids distortion of the image and objects in it, but the crops will possibly lose much information when the aspect ratio of the original image is far away from 1.
3. *Dense cropping*: Keep the aspect ratio of the original image, resize the shorter side to 256, and then densely crop multiple 227×227 patches and their horizontal reflections. This way, the information of original image is kept by dense cropping process, and the distortion is avoided. The number of crops attained with this method is larger than the previous two methods.

3.3.2 Network Fine-tune

Although the pre-trained network from [16] can already generalize well to many datasets, the urban tribe dataset has its unique property. It emphasizes certain visual features such as certain clothing styles, while pays less attention to other visual features. Also, it emphasizes style of the object rather than distinct category. To rearrange the importance of different features and adjust the features to adapt to urban tribe dataset, the network can be fine-tuned.

The dataset used for fine-tune is the same set used for SVM training. The input patches of the Convolutional network is of size 227×227 . The initial Convolutional network has 1000 outputs in final layer, corresponding to 1000 class-wise probability predictions. In our fine-tune process, the last layer is replaced by 11 probability prediction outputs, and the initial weights of the last layer connection are initialized to have zero mean gaussian distribution. Back propagation is used, and the learning rate is set to be small so that the fine-tune process adapts the extracted features to urban tribe dataset while preserving the initial property in general: the initial learning rate used for pre-training is 0.01; We set the initial fine-tune learning rate of the parameters except for the last layer as 0.001, and keep the initial learning rate for last layer 0.01, because the last layer is not pre-trained.

3.3.3 Choices of network combination

Scene images and individual images have different properties, and need different strategies of pre-processing and sep-

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

432 arate fine-tuning. For scene network $\text{CNN}_{\text{Scene}}$ and scene
 433 images, due to the small size of the dataset, we use the *dense*
 434 *cropping* technique, the third technique in Section 3.3.1, to
 435 increase the dataset. For person network $\text{CNN}_{\text{Person}}$ and
 436 corresponding input, we use the *distort cropping* technique,
 437 because the subimages have normally long height and short
 438 width, and the second and third strategies using squared
 439 crops of a person image will lose much information, no
 440 matter which location we choose to crop them; whereas the
 441 first method ensure each crop keeps the essential features
 442 for classification.

443 The combination of dense crop $\text{CNN}_{\text{Scene}}$ and distorted
 444 crop $\text{CNN}_{\text{Person}}$ are denoted as Nets_{SDense} .

445 We also construct other combination of networks for
 446 comparison:

- 448 1. Nets_{NoTune} : Directly use the pre-trained network by
 449 [16] for both scenes and persons, and use the *distort*
 450 *cropping* technique (distorted crops) as input patches
 451 for both networks. This choice of cropping strategy is
 452 in consistent with the way the network is pre-trained.
- 454 2. $\text{Nets}_{SSparse}$: Use the *sparse cropping* strategy for
 455 scene features, and the *distort cropping* strategy for
 456 person features. .
- 458 3. $\text{Nets}_{SDistort}$: Use the *distort cropping* strategy for
 459 both scene features and person features.

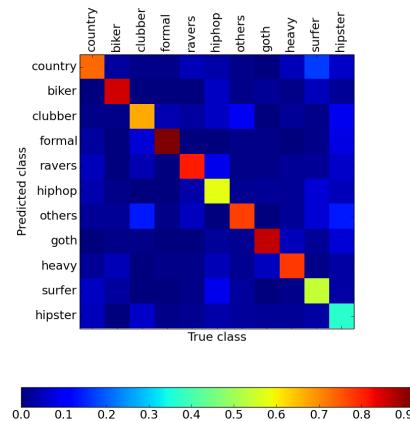
4. Experiments and Results

463 In this section, the performance of the proposed classi-
 464 fication scheme is evaluated and analyzed. In the experi-
 465 ments, six rounds of 5-fold cross validation are performed,
 466 therefore we have 30 training experiments in total. Dataset
 467 is partitioned into 5 equal sized subsets, containing one-fifth
 468 of the data points from each category. One of the subsets is
 469 used as test set, and the remained 4 subsets are used as train-
 470 ing data. In the fine-tune procedure, there are 6000 training
 471 iterations, and learning rate is decreased by ten times after
 472 each 1000 iterations.

4.1. Urban tribe classification performance

475 Table 1 shows the comparison of performance using dif-
 476 ferent approaches. The 30 segmentations of datasets are
 477 used for all the approaches tested in this section, and 30 test
 478 results are averaged for each approach. The standard error
 479 is shown with the accuracy in Table 1. We also compare our
 480 result with the result achieved by [1] using their best model.
 481 The advantage of CNN pre-trained features is obvious.

482 The confusion matrix is shown in Figure 2 for
 483 Nets_{SDense} , and all the 30 training experiments are aver-
 484 aged. We can observe there is a obvious difference of dif-
 485 ficulty of different categories. Class *formal* has accuracy as



502 Figure 2: Confusion matrix for classification results with
 503 Nets_{SDense} , using people and scene features.

507 high as about 90%, while class *hipster* is the most difficult
 508 class, having less than 60% accuracy.

509 Comparing the result of using different features in the
 510 same approach shows the necessity of every step of our
 511 architecture. In results using Nets_{SDense} with concatenated
 512 features, average accuracy for each candidate person is low
 513 as 47.10%. Average pooling of candidate person probability
 514 estimates produces a large accuracy increase of about 20%.
 515 Accuracy using the entire scene only results in 67.26% ac-
 516 curacy. Combining probabilities $\text{Pr}_{People}(C|H_1, \dots, H_p)$
 517 and $\text{Pr}_{Scene}(C|S)$ achieves accuracy as high as 71.22%,
 518 which verifies the complementary role of people candidate
 519 feature and environment feature in a group image.

520 We also compare vertically the accuracy of different ap-
 521 proaches, to show the role of network feature concatena-
 522 tion. Using only 7th layer or 6th layer activation from
 523 the networks Nets_{SDense} produces decent results, showing
 524 both layers' activations can generate high semantic features.
 525 Concatenating both layers' activations increases the accu-
 526 racy by 0.5%, indicating the slight information loss of the
 527 7th fully connected layer.

528 To see the role of fine-tuning, we can compare the result
 529 of Nets_{NoTune} and $\text{Nets}_{SDistort}$. These two approaches
 530 both use resizing that causes distortion, and they only vary
 531 in fine-tune procedure. There is a large performance im-
 532 provement with fine-tune, both for person features and
 533 scene features. This shows the benefit of adapting existing
 534 generic Convolutional Network to specific dataset.

535 $\text{Nets}_{SDistort}$, $\text{Nets}_{SSparse}$, and Nets_{SDense} use different
 536 patch extraction strategies. Note that we use the same dis-
 537 torted patch extraction method for person images, as men-
 538 tioned in Section 3.3.3, while we use three different meth-

Table 1: Performance of different approaches using different information.

| Accuracy (%) | Individual candidate | People | Entire scene | People+Scene |
|---|------------------------------------|------------------------------------|------------------------------------|------------------------------------|
| Nets _{NoTune} with concatenated features | 40.03 ± 0.31 | 64.18 ± 0.63 | 62.61 ± 0.58 | 69.19 ± 0.51 |
| Nets _{SDense} with fc7 features | 46.95 ± 0.33 | 67.42 ± 0.50 | 64.87 ± 0.39 | 70.68 ± 0.47 |
| Nets _{SDense} with fc6 features | 45.80 ± 0.37 | 66.31 ± 0.47 | 66.72 ± 0.56 | 70.68 ± 0.44 |
| Nets _{SDense} with concatenated features | 47.06 ± 0.37 | 67.46 ± 0.51 | 67.01 ± 0.52 | 71.45 ± 0.48 |
| Nets _{SSparse} with concatenated features | 47.06 ± 0.37 | 66.48 ± 0.41 | 66.81 ± 0.42 | 71.26 ± 0.49 |
| Nets _{SDistort} with concatenated features | 47.06 ± 0.37 | 67.46 ± 0.51 | 65.06 ± 0.48 | 71.15 ± 0.53 |
| <i>SVM₈[1]</i> | - | - | - | 46(std: 2) |

ods for scene images. The results for scene images show the advantage of keeping the aspect ratio of scene images, and the slight advantage of using dense crops. However, the final results with People+Scene for the three methods don't have significant differences, this is due to the combination with people information.

4.2. Convolutional Network feature analysis

Yufei: Is this section necessary? We use t-distributed stochastic neighbor embedding technique [21] to visualize the power of different layers' features in convolutional network. Two dimensional embedding of the high dimensional CNN feature space is extracted, and we plot the data as 2-d points with different colors indicating different classes they belong to. Powerful features will pull data of different semantic classes more apart.

We randomly choose one training-test partition and corresponding fine-tuned network parameters, and to avoid overfitting effect, we examine the test set of Nets_{SDense} approach. In Figure 3, each data is plotted as a dot in each figure. The three columns correspond to data separation in first layer, fourth layer and seventh layer respectively. The first row visualizes features of all test data in scene CNN CNN_{Scene}, and second row picks three classes from the first row. The third row is features of CNN_{Person}, and the three classes are picked for the last row. The three classes chosen in second and last row is (goth:red, heavy:green, surfer:blue).

In both CNN_{Scene} and CNN_{Person}, there is a clear trend of class separation. As the layer ascends, , the data from same class are more concentrate, and inter-class distance are larger. The selected three classes show the trend more clearly: they essentially form three clusters in the seventh layer. As shown in the second and last rows, person data is more challenging: The semantic classes are not as separate as with the scene features.

4.3. Urban tribe classes vs. ImageNet classes

It's recently being acknowledged that CNN pre-trained features are generic and can be used for new tasks. In this section we check the relationship between the new tasks and

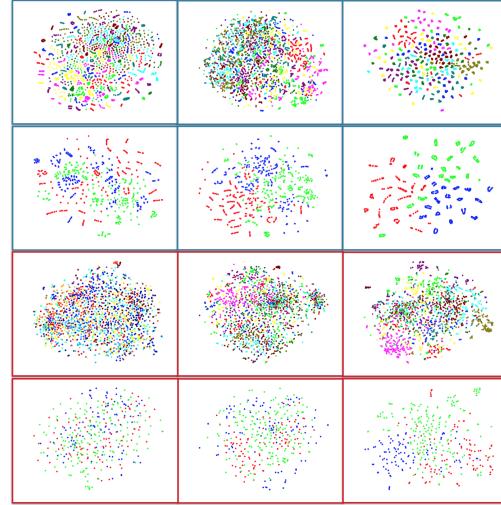


Figure 3: Feature visualization of Nets_{SDense}: CNN_{Scene} (first two rows) and CNN_{Person}(second two rows). Rows: first: all scene images; second: all data from goth/heavy/surfer urban tribes; third: all test person images; fourth: all test data from goth/heavy/surfer urban tribes. Columns: first: first convolutional layer; second: fourth convolutional layer; third: seventh fully convolutional layer

original ImageNet task, using the urban tribe dataset, which gives us some insight about the features extracted from pre-trained network.

The urban tribe dataset contains groups of people, and the important features for categorization are mainly human related features, such as attire, make up, posture and expressions. However, the ILSVRC dataset used for pre-training contains few images of humans. Instead of examining the output of the layers directly, we try to find the relationship between the 1000 classes in ILSVRC dataset and the classes in two image sets: urban tribe dataset, and candidate person images extracted from urban tribe dataset.

We use the parameters of pre-trained CNN network as our feature extraction model, and train a softmax layer on top of its 7th layer to predict the probabilities of one input

648
649
650
651
652
653
654
655
656
657
658
image(either scene image or candidate person image) being
in certain urban-tribe class $\Pr(l_{urban})$, where l_{urban} pre-trained
is the 11 urban tribe categories. The output layer
is trained for 3000 training iterations. We can also use the
output of the pre-trained CNN network to predict probabilities
of one input image being in certain ImageNet category,
denoted as $\Pr(l_{ImageNet})$, where $l_{ImageNet}$ is the 1000 ImageNet
categories. We use one round of 5-fold cross validation,
and use all the images in urban tribe dataset for analysis.

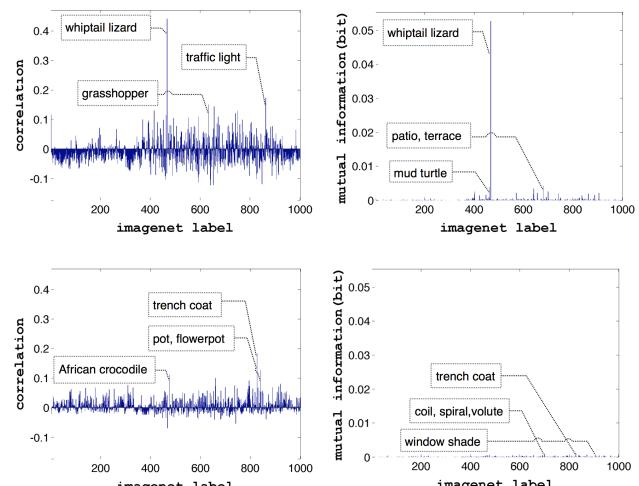
659
660
661
662
663
664
665
666
667
We first check the relationship of $\Pr(l_{urban})$ and
 $\Pr(l_{ImageNet})$ of candidate person images. We calculate
the correlation coefficient of $\Pr(l_{urban})$ and
 $\Pr(l_{ImageNet})$ for all 1000 ImageNet classes, denoted as
 $R(\Pr(l_{urban}), \Pr(l_{ImageNet}))$. We also calculate the 1000
mutual information of the predicted score of urban tribe
and ImageNet (where predicted score is 1 if the predicted
label is the category being tested, 0 otherwise), denoted as
 $I(l_{urban}, l_{ImageNet})$.

668
669
670
671
672
673
674
675
676
In Figure 4, we choose two urban classes: *biker* and *hipster*,
and plot the correlation coefficient R and mutual information
 I . The first row shows the result of *biker*, and second
row *hipster*. For *biker*, there are several impulses in correlation
plot, and one significant impulse in mutual information
plot. *whiptail lizard* has high correlation with *biker*. Mean-
while, for *hipster*, which is the most difficult class, the cor-
relation coefficient and mutual information are both low for
all ImageNet classes.

677
678
679
680
681
682
683
684
685
686
687
688
689
To confirm the correlation, we use $\Pr(l_{ImageNet})$ directly
as features, substituting the concatenated fc7fc6 features,
and use the Nets_{SDistort} approach for classification. The
accuracy is 52.68%. This decent result indicates the rela-
tionship between l_{urban} and $l_{ImageNet}$. Then, we check
 $l_{ImageNet}$ with highest correlation coefficients with l_{urban} . In
Figure 5, we choose four l_{urban} : *formal*, *ravers*, *goth*,
hipster, and choose some examples of candidate person im-
ages that have both high $\Pr(l_{ImageNet})$ and high $\Pr(l_{urban})$. We
also show examples of the images in $l_{ImageNet}$. We can see some of the shared features between corresponding
person images and ImageNet images, for example, similar
shape for women tops and stingrays (Figure 5b).

690
691
692
693
694
695
696
697
There is a correlation between class-wise accuracy of
predicted l_{urban} and the degree of relationship between pre-
dicted $\Pr(l_{urban})$ and $\Pr(l_{ImageNet})$, as shown in Figure 6.
Class-wise accuracy is calculated for candidate person im-
ages (Figure 6a, 6b) or scene images (Figure 6c, 6d). For
each l_{urban} , the maximum correlation/mutual information
over 1000 ImageNet classes are used to indicate the degree
of its relationship with $l_{ImageNet}$.

698
699
700
701
The correlation between ImageNet class and urban tribe
class and its relationship with class-wise recognition rate
may indicate that the “generic” features extracted by pre-
trained CNN networks are not so generic. The network is



702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
Figure 4: Correlation and mutual information of l_{urban} and 1000 $l_{ImageNet}$. The first row is for $l_{urban}=\text{Biker}$. The second row is for $l_{urban}=\text{Hipster}$. Top three ImageNet classes are marked.

736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
trained to separate the ImageNet classes most, and if we use
the features for a new classification task, the performance
of the task is related to how well the new classes can be
“mapped” to the ImageNet classes.

5. Conclusion

736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
In this work, we proposed a framework for social group
recognition. The framework takes in both individual and
global features. Features are extracted from fine-tuned
CNN networks which has been pre-trained on ImageNet
dataset, and then combined. Our results showed the suc-
cess of our framework by achieving much better result than
the previous work.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
We also investigated into the pre-trained CNN features.
Both visualization and numeric results showed the gener-
alization ability of pre-trained CNN features to features of
people’s social styles, which has little shared features with
the ImageNet object categories. Meanwhile, we found that
there is a correlation between the probability of an image
being in Imagenet classes and social group classes, and that
better-recognized categories are more correlated with Image-
Net categories.

802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
999
In the future work, we intend to improve the classifica-
tion performance by adapting convolutional networks more
to social groups datasets. The relationship between Image-
Net categories and urban tribes classes also brings for-
ward an interesting future topic.

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

Figure 5: selected urban tribe classes and the corresponding highest correlation $l_{ImageNet}$. Upper nine images: candidate person images with high $Pr(l_{ImageNet})$ and high $Pr(l_{urban})$. Lower eight images: example of images in $l_{ImageNet}$.

(a) Individual accuracy - maximum of R (b) Individual accuracy - maximum of I (c) Scene accuracy - maximum of R (d) Scene accuracy - maximum of I

Figure 6: The relationship between class-wise recognition rate and maximum of correlation $R(\Pr(l_{urban}), \Pr(l_{ImageNet}))$, class-wise recognition rate and maximum of mutual information $I(l_{urban}, l_{ImageNet})$

References

- [1] I. S. Kwak, A. C. Murillo, P. Belhumeur, S. Belongie, and D. Kriegman, “From bikers to surfers: Visual recognition of urban tribes,” in *British Machine Vision Conference (BMVC)*, (Bristol), September 2013.
- [2] A. C. Murillo, I. S. Kwak, L. Bourdev, D. Kriegman, and S. Belongie, “Urban tribes: Analyzing group photos from a social perspective,” in *CVPR Workshop on Socially Intelligent Surveillance and Monitoring (SISM)*, (Providence, RI), June 2012.
- [3] M. D. Zeiler and R. Fergus, “Visualizing and understanding convolutional networks.,” *CoRR*, vol. abs/1311.2901, 2013.
- [4] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell, “Decaf: A deep convolutional activation feature for generic visual recognition.,” *CoRR*, vol. abs/1310.1531, 2013.
- [5] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks.,” in *NIPS* (P. L. Bartlett, F. C. N. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, eds.), pp. 1106–1114, 2012.

- [6] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, “Backpropagation applied to handwritten zip code recognition,” *Neural Computation*, vol. 1, pp. 541–551, 1989.
- [7] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” in *Proceedings of the IEEE*, vol. 86(11), pp. 2278–2324, 1998.
- [8] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio, “Maxout networks,” 2013, arXiv:1302.4389.
- [9] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. R. Elnikad, Dragomir Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going Deeper with Convolutions,” 2014, arXiv:1409.4842.
- [10] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, “Overfeat: Integrated recognition, localization and detection using convolutional networks,” in *International Conference on Learning Representations (ICLR2014)*, 2014.
- [11] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic

8

- 864 segmentation,” in *Proceedings of the IEEE Conference on* 918
865 *Computer Vision and Pattern Recognition (CVPR)*, 2014. 919
866
- 867 [12] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, 920
868 and R. Salakhutdinov, “Improving neural networks by 921
869 preventing co-adaptation of feature detectors,” *CoRR*, 922
870 vol. abs/1207.0580, 2012. 923
871
- 872 [13] I. J. Goodfellow, D. Warde-Farley, M. Mirza, 924
873 A. Courville, and Y. Bengio, “Network in network,” 925
874 2013, arXiv:1312.4400. 926
875
- 876 [14] S. Karayev, M. Trentacoste, H. Han, A. Agarwala, T. Darrell, 927
877 A. Hertzmann, and H. Winnemoeller, “Recognizing Image 928
878 Style,” 2013, arXiv:1311.3715. 929
879
- 880 [15] P.-H. Gosselin, N. Murray, H. Jégou, and F. Perronnin, 930
881 “Inria+xerox@fgcomp: Boosting the fisher vector for fine- 931
882 grained classification,” Tech. Rep. 0, INRIA, December 932
883 2013. 933
884
- 885 [16] Y. Jia, “Caffe: An open source convolutional architecture 934
886 for fast feature embedding.” <http://caffe.berkeleyvision.org/>, 935
887 2013. 936
888
- 889 [17] A. C. Gallagher and T. Chen, “Understanding images of 937
890 groups of people.,” in *CVPR*, pp. 256–263, IEEE, 2009. 938
891
- 892 [18] G. Wang, A. C. Gallagher, J. Luo, and D. A. Forsyth, “Seeing 940
893 people in social context: Recognizing people and social 941
894 relationships.,” in *ECCV*(K. Daniilidis, P. Maragos, and N. Paragios, eds.), vol. 6215, pp. 169–182, Springer, 2010. 942
895
- 896 [19] A. Dhall, J. Joshi, I. Radwan, and R. Goecke, “Finding hap- 943
897 piest moments in a social context.,” in *ACCV 2* (K. M. Lee, 944
898 Y. Matsushita, J. M. Rehg, and Z. Hu, eds.), vol. 7725, 945
899 pp. 613–626, Springer, 2012. 946
900
- 901 [20] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.- 947
902 J. Lin, “Liblinear: A library for large linear classification,” 948
903 *Journal of Machine Learning Research*, vol. 9, pp. 1871– 949
904 1874, 2008. 950
905
- 906 [21] G. H. Laurens van der Maaten, “Visualizing data using 951
907 t-sne,” *Journal of Machine Learning Research*, vol. 9, 952
908 pp. 2579–2605, November 2008. 953
909
- 910
- 911
- 912
- 913
- 914
- 915
- 916
- 917