

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020

Supplementary Material for Learning to Annotate Facial Action Units for Online Images

Anonymous ICCV submission

Paper ID 2547

Abstract

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074

This material includes more results that cannot be fitted into the main paper due to space limitation. In specific, we provides more qualitative analysis of the proposed weakly supervised clustering method.

1. Qualitative analysis on weakly supervised clustering

075
076
077
078
079
080
081
082
083
084
085
086
087
088
089

This section shows more qualitative results of the proposed scalable weakly supervised clustering (WSC), t-SNE visualization of WSC embeddings and the corrections of annotations. Recall that we modeled WSC as a variant of spectral clustering with scatterness constraints, which learns an optimal embedding space where images with similar appearance and weak annotations are more likely to be in the same neighborhood. Details can be found in Section 3.1 of the main paper. Weak annotations can be query string or meta data *etc.* As mentioned in Section 4.1 of the main paper, the weak annotations used in this paper are the AU predictions of the pre-trained classifiers on BP4D.

090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Scalable weakly supervised clustering: Fig. 1(a) illustrates the objective value of scalable weakly spectral clustering (sWSC) optimized by the accelerated stochastic gradient method, as discussed in Section 3.2 in the main paper. The figure was obtained by running sWSC on clustering $\sim 200,000$ images of EmotioNet into two clusters using weak annotations of AU12. Although sWSC is stochastic-based, we can observe that convergence happened at around #5000 iterations. Fig. 1(b) shows the average images of two different clusters. As can be seen, each of these two clusters tends to group images from different modes of facial actions. On images of the cluster denoted as “AU 1+2”, we observed more frequent AUs 1 and 2. The other cluster, denoted as “AU 6+12”, exhibits strong lips corner puller and cheek raiser that directly corresponds to AU 6 and 12. This observation justifies that the proposed WSC can preserve both visually similar images while at the same time meaningful AU combinations, and meanwhile is scalable to large amount of images.

t-SNE visualization of WSC embeddings: To further understand the learned WSC embedding (denoted as \mathbf{W} in the main paper), we examined the closeness of the images using t-SNE embedding. Fig. 2 shows the embedding and zoomed-in neighborhood to facilitate examination. As can be seen, samples with cheek raiser (AU6) and lips corner puller (AU12) bounded with red boxes are grouped in the lower area of the figure. Regardless of lighting conditions and head poses, most images within this neighborhood contain either subtle or obvious appearance of AU6 or AU12. On the other hand, Samples with eyebrow puller bounded with blue boxes are shown in the upper areas of the embedding space. This shows WSC embedding obtains meaningful representations for AU annotation.

Corrections of annotations: EmotioNet is collected from the online images with query strings. Using the weak annotations (*e.g.* queries or AUs obtained by pre-trained model) directly as the supervision easily leads to get the sub-optimal generalization by the inaccurate classifier. This scenario even happens for manually annotations. The proposed WSC can correct the possible inaccurate annotations for weak annotations and even for humans annotations. We show more qualitative results of using WSC to correct manually annotations here. Partial results were shown in Section 4.1 of the main paper. The left most column of Fig. 3 shows the canonical faces of the 6 AUs used in the paper. On the right of each row, we show images that WSC was able to correct their annotations. In other words, these images were originally annotated as inactive with respect to each AU. WSC discovered these patterns and demonstrated good consistency with FACS-coding, showing that noisy annotations can be further cleaned using the proposed WSC.

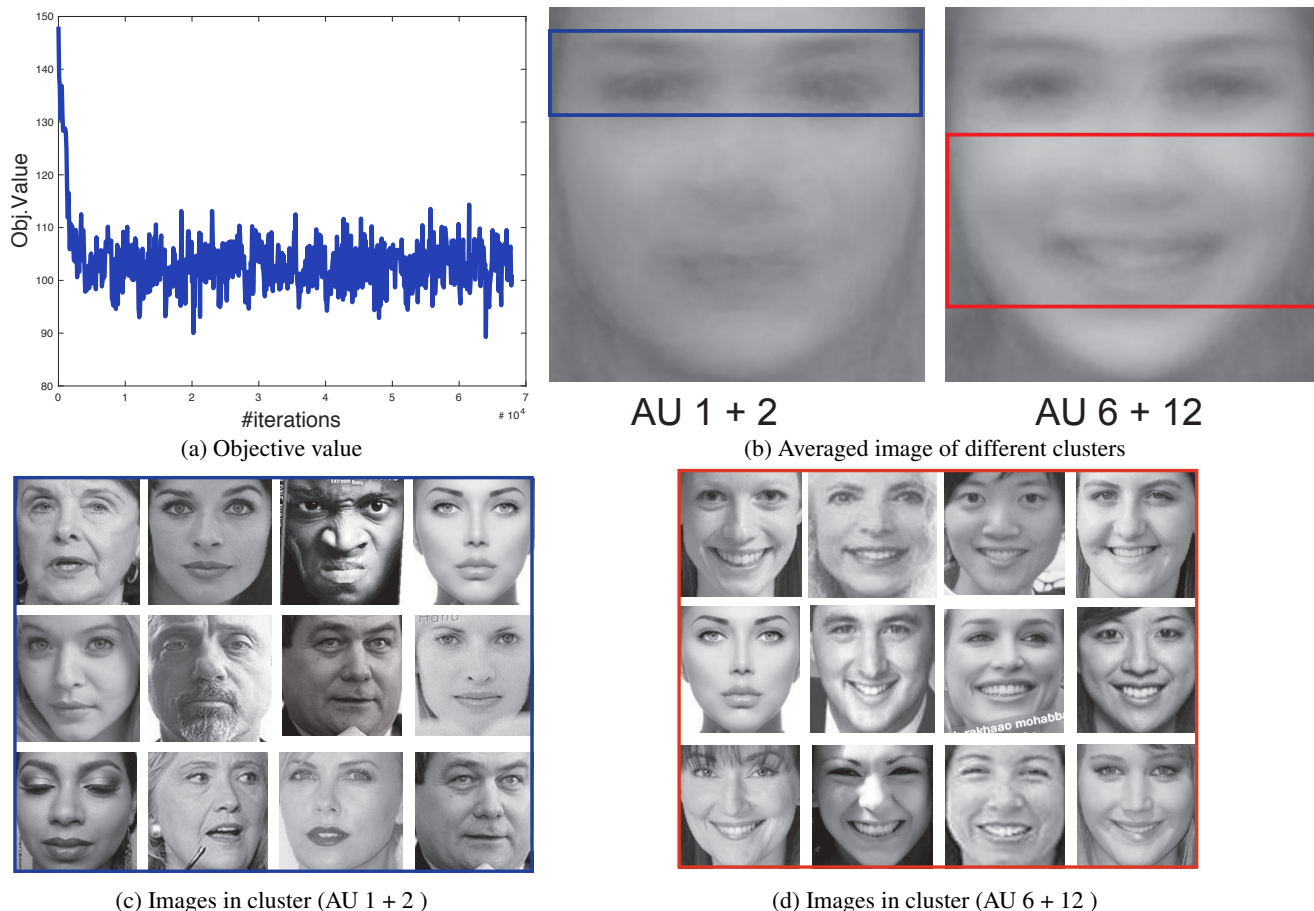


Figure 1. Visualization of averaged images of different clusters upon convergence. (a) shows the objective value v.s. the number of iterations of scalable weakly spectral clustering (sWSC). The experiments are conducted on $\sim 200,000$ images of EmotionNet with weak annotations of AU12. (b) shows the averaged images of two different clusters. We denoted different averaged images as “AU1+2” and “AU6+12” to indicate WSC’s power in grouping visually and semantically similar images. (c) and (d) are the images in these two clusters.

2. Derivation of accelerated stochastic optimization method

This section provides more details in deriving the accelerated stochastic algorithm as discussed in Section 3.2 of the main paper. Specifically, during the group-wise optimization for each \mathbf{W}_g , we write:

$$\min_{\mathbf{W}_g} \frac{1}{2} \|\mathbf{W}_g - \mathbf{V}_g\|_F^2 + \frac{\tilde{\lambda}}{n_g} \text{Tr}(\mathbf{W}_g \mathbf{C}_g \mathbf{W}_g^\top) \quad (1)$$

Taking derivative of equation (1) as:

$$\begin{aligned} \mathbf{W}_g - \mathbf{V}_g + \frac{2\tilde{\lambda}}{n_g} \mathbf{C}_g \mathbf{W}_g &= 0 \\ (\mathbf{I}_{n_g} + \frac{2\tilde{\lambda}}{n_g} \mathbf{C}_g) \mathbf{W}_g &= \mathbf{V}_g \end{aligned} \quad (2)$$

For \mathbf{C}_g is a centering matrix, which has the eigenvalue 1 of multiplicity $n - 1$ and eigenvalue 0 of multiplicity 1. In addition, $\tilde{\lambda} > 0$. Thus, $\det(\mathbf{I}_{n_g} + \frac{2\tilde{\lambda}}{n_g} \mathbf{C}_g) \neq 0$. The optimal solution for (1) can be obtained as $\mathbf{W}_g^* = (\mathbf{I}_{n_g} + \frac{2\tilde{\lambda}}{n_g} \mathbf{C}_g)^{-1} \mathbf{V}_g$.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323



Figure 2. t-SNE visualization of the learned WSC embedding on ~200,000 images of EmotioNet.

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

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



Figure 3. An illustration of samples that the proposed WSC is able to correct for human annotations in the EmotioNet dataset. Each row shows the corrected samples of different AUs. The left most columns are definitions of AUs and the samples in correspondence to FACS. The samples corrected in EmotioNet are shown in the right column.