

Mask R-CNN ile Derin Yüz Sezici Gerçekleme Design of a Deep Face Detector by Mask R-CNN

Ozan Cakiroglu Caner Ozer Bilge Günsel

Çoğulortam Sinyal İşleme ve Örüntü Tanıma Grubu

Elektronik ve Hab. Müh. Bölümü, İstanbul Teknik Üniversitesi, Türkiye

{cakirogluoz, ozerc, gunselsb}@itu.edu.tr

Özetçe—Bu çalışmada literatürde varolan Mask R-CNN derin öğrenme ağı yüz sezme amaçlı eğitilmiş ve öğrenilen yüz nesnesi için başarımlar raporlanmıştır. Çalışmada, varolan çalışmalardan farklı olarak, az sayıda yüz gözlemi ile bir eğitim gerçekleştirilmiş ve yüz nesnesini çevreleyen kutunun yakalanmasının yanı sıra, yüz bölgesinin piksel bazında bölütlenmesi de hedeflenmiştir. Eğitim PASCAL-VOC veri setinden alınan 2695 adet gözlem ile gerçekleştirilmiştir. Başarımlar, yüz sezme çalışmalarında referans alınan WIDER FACE veri tabanı üzerinde 159.000 yüz nesnesi için raporlanmıştır. Sonuçlar, baz alınan çalışmaya [1] göre çok küçük ölçekli yüzleri sezmede %6, orta ölçekli yüz sezmede %12, büyük ölçekli yüz sezmede %3 daha yüksek başarımlar sağladığını göstermektedir. Performansımızın Viola & Jones yüz seziciden daha yüksek olduğu da raporlanmaktadır. Bu çalışmada derin yüz sezici eğitiminde kullanılan öğrenme örnekleri için çıkarılan bölütleme verisi ve TensorFlow ortamında gerçekleştirilen eğitim-test rutinleri [GitHub](#) altında genel kullanıma açılmıştır.

Anahtar Kelimeler—Yüz sezme, örnek düzeyinde bölütleme, derin öğrenme, evrimsel sinir ağları.

Abstract—In this work an existing object detector, Mask R-CNN, is trained for face detection and performance results are reported by using the learned model. Differing from the existing work, it is aimed to train the deep detector with a small number of training examples and also to perform instance segmentation along with an object bounding box detection. Training set includes 2695 face examples collected from PASCAL-VOC database. Performance has been reported on 159,000 test faces of WIDER FACE benchmarking database. Numerical results demonstrate that the trained Mask R-CNN provides higher detection rates with respect to the baseline detector [1], particularly 6%, 12%, and 3% higher face detection accuracy for the small, medium and large scale faces, respectively. It is also reported that our performance outperforms Viola & Jones face detector. We released the face segmentation ground-truth data that was used to train Mask R-CNN and training-test routines developed in TensorFlow platform to public usage at our [GitHub repository](#).

Keywords—Face detection, instance segmentation, deep learning, convolutional neural networks.

I. GİRİŞ

Son yıllarda, yüz sezme sistemlerinde evrimsel sinir ağı temelli tek [2], [3] veya iki [4], [5] aşamalı nesne sezme yöntemleri kullanılmaktadır. Bu yöntemlerde farklı ölçeklerdeki yüzlerin bulunması amaçlanmış ve klasik yüz sezicilere [6] göre daha üstün başarımlar elde edilmiştir. Ayrıca, RetinaNet [7] veya çevrimiçi keskin örnek madenciliği (OHM) [8] kullanılarak, eğitimde karşılaşılan keskin (hard) örneklerin daha iyi öğrenilmesi sağlanabilmektedir. Genel olarak

öğrenme aşamasında yüz nesnesinin derin yüz sezici tarafından öğrenilmesi sağlanmakta, ardından başarımlar farklı veri tabanlarında raporlanmaktadır. Bu çalışmalarda, yüz sezme başarımlarının raporlanması amacıyla WIDER FACE [1], FDDB [9] ve PASCAL-VOC Faces [10] veri tabanları sıklıkla kullanılmaktadır. Başarımlar, ortalama keskinlik (AP) cinsinden raporlanmakta ve başarımların büyük oranda yüz nesnesinin ölçeği arttıkça yükseldiği çıkarımı raporlanmaktadır.

Bu çalışmada, literatürde varolan yüz sezicilerden farklı olarak, örnek düzeyinde bölütleme de yaparak yüz sezme başarımlarını arttırmak amacıyla Mask R-CNN derin öğrenme ağı eğitilmiş, öğrenme ve test başarımları literatüre uygun olarak raporlanmıştır.

II. MASK R-CNN İLE NESNE SEZME

Bu bölümde Mask R-CNN [11] derin sinir ağı ile giriş görüntüsü üzerinde gerçekleştirilen iki aşamalı nesne sezme ve öğrenme modelinin çıkarılması anlatılmaktadır. Mask R-CNN mimarisinde ilk aşamada olası nesne bölgelerinin koordinatları ve nesne içerme olasılıkları, bölge öneri ağı (RPN) adı verilen yapı ile çıkartılır. İkinci aşamada her bir olası nesne bölgesinin yüz içerme olasılığı ve nesneyi çevreleyen kutular (bounding box-BB) bulunur. Mask R-CNN örnek bazında bölütleme gerçekleştirilerek BB lerin yanısıra yüz olarak etiketlenmiş piksellerden oluşan maskeyi de üretir. Şekil 1’de Mask R-CNN’in genel mimarisi üzerinde işlem akışı gösterilmektedir.

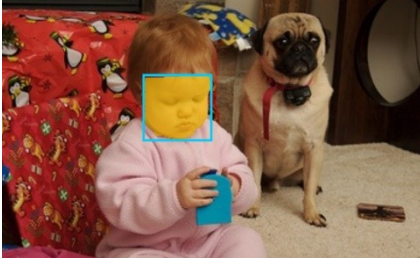
Verilen bir görüntü için konvolüsyonel sinir ağı iskelet mimarisi, öznitelik çıkarma işlemini gerçekleştirmektedir. Bu çalışmada öznitelik çıkarıcı olarak, farklı uzamsal boyutlarda öznitelik çıkartmakta sıklıkla kullanılan, ResNet 101’in [12] farklı seviyelerdeki çıktılarını kullanan Öznitelik Piramit Ağları (FPN) [13] yapısı kullanılmıştır. Elde edilen öznitelik haritaları RPN’den geçirilerek, bütün lokal bölgelerin nesne olabilirlik skorları ve BB koordinatları, sırasıyla ikili sınıflandırma ve regresyon ile bulunmaktadır.

Eğitim esnasında Eşitlik 1’de görülen RPN kayıp fonksiyonunun enküçüklenmesi sağlanır. Burada ilk terim sınıflandırıcı hatasını, ikinci terim ise regresyon hatasını modellemektedir.

$$L_{RPN} = \frac{1}{N} \sum_{i=1}^N (L_{cls}(p_i, p_i^*) + L_{reg}(t_i, t_i^*)). \quad (1)$$

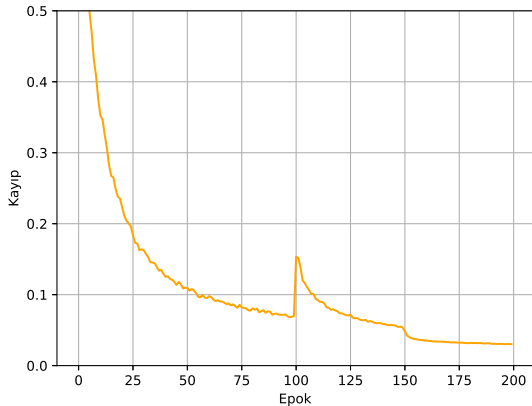
Eşitlik 1’de N eğitim sırasında kullanılan ankor sayısını, p_i ve p_i^* sırasıyla, yakalanan BB içerisinde nesne olabilirlik skorunun tahmin edilen ve gerçek referans değerini göstermektedir. Referans ile örtüşme oranına göre pozitif için 1,

böylelikle başka grupların çalışmalarında da kullanılabilecek açık bir eğitim veri seti elde edilmiştir. Şekil 2’de oluşturulan veri setinden bir yüz örneği, kendisini çevreleyen BB ve piksel bazında işaretli bölütleme maskesi ile birlikte görülmektedir.



Şekil 2: Veri seti içerisindeki işaretlenmiş görüntü örneği.

Eğitimler, NVIDIA GTX 1080 ekran kartıyla, 200 epok (epoch) ve her epokta 1400 görüntü örnekleyen iterasyonlar ile toplam 36 saatte tamamlanmıştır. Referans model olarak, Mat-terport’un [16] TensorFlow ve Keras kütüphaneleri ile yazılım ortamında gerçekleştirdiği derinlikli öğrenme ağ yapısı alınarak, bu yapıya gerekli eklemeler yapılmıştır. Eğitim, ağ yapısının baş katmanlarının (head layers) 100 epok boyunca eğitilmesi ve ResNet 101 iskelet katmanları ve alt katmanların 50 epok eğitilmesi ile gerçekleştirilmiştir. Son olarak ağ yapısı uçtan uca 50 epok boyunca eğitilerek, öğrenme modeli elde edilmiştir. Mimarının baş katmanlarının eğitiminde öğrenme oranı 0.001 olarak seçilirken, 100 epok ardından 0.0001’e düşürülmüştür. ResNet-101 iskelet mimarisinin ağırlıkları, MS COCO [17] veri setinde eğitilen ağırlıklar olarak belirlenerek, eğitimde başlangıç ağırlık değerleri olarak alınmıştır. Ankor boyutları (0.5, 1, 2) en-boy oranı ve (32, 64, 128, 256, 512)’lık pencere boyutları olmak üzere 5 adet olarak seçilmiştir. Tanımlanan yüz sezici eğitimi sırasında kayıp fonksiyonunun epok sayısına bağlı değişimi kaydedilmiştir ve Şekil 3’de görülmektedir. Grafikten görüldüğü üzere, kayıp azalarak minimum değerine yakınsamaktadır ki bu durum Mask R-CNN’in yüz nesnesini öğrendiğini göstermektedir.



Şekil 3: En iyileme kayıp fonksiyonunun yapılan eğitim sırasında raporlanan değişimi.

IV. BAŞARIM TEST SONUÇLARI

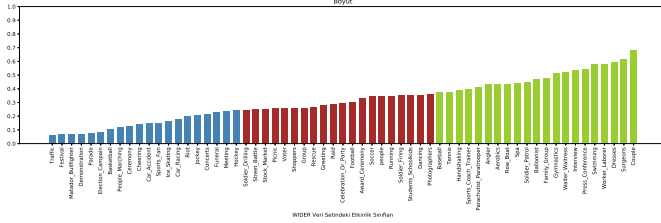
Testlerde öncelikle yüz sezme gerçeklemek üzere eğitilen derin sinir ağının öğrenme başarımı irdelenmiştir. Bu amaçla Pascal VOC üzerinde eğitim sonucu elde edilen model kullanılarak ve eğitim seti test seti olarak alınarak sezme başarımı raporlanmıştır. Tablo I’de yüz sezme başarımı duyarlılık (R-recall), kesinlik (P-precision) ve F-1 ölçütü ile farklı kesişim (IoU) eşik değerleri için raporlanmaktadır. Literatürde IoU=0.5 objenin doğru yakalandığı kabul edilen en düşük kesişim değeridir. Tablo I’de Tr-Tr altında raporlanan değerlerden görüldüğü üzere, F-1 ölçütüne göre başarımlar ≥ 0.5 için %92 gibi yüksek bir değerdir ve bu durum eğitim sonucunda yüz objesinin başarıyla öğrenildiğini göstermektedir. Başarım yüksek IoU değerlerine gidildikçe azalmaktadır; ancak sezilen yüz sınırlarının bilinen gerçek referans ile en az %80 oranında örtüşmesi anlamına gelen $\text{IoU} \geq 0.8$ değeri için halen %74 oranında düşük olmayan bir sezme başarımı raporlanmaktadır. İkinci test senaryosunda, Pascal VOC’da bulunan ancak eğitim seti ile örtüşmeyen 160 görüntüde bulunan 225 yüz üzerinde başarımlar raporlanmıştır. Tablo I’de Tr-Test altında farklı IoU eşik değerleri için raporlanan sonuçlar incelendiğinde, eğitilen yüz sezicinin, eğitimde görülmeyen görüntülerdeki yüz objelerini de hemen hemen aynı doğrulukla yakalayabildiği görülmektedir. Bu nedenle, modelin eğitim setini ezberlemediği ve öğrenilen ağ ağırlıklarının derin yüz seziciyi genellenebilir bir sezme performansına getirdiği sonucu çıkarılabilir.

IoU	R	Tr-Tr		R	Tr-Test	
		P	F-1		P	F-1
0.5	0.93	0.90	0.921	0.93	0.92	0.93
0.6	0.91	0.88	0.897	0.92	0.91	0.92
0.7	0.84	0.81	0.827	0.83	0.82	0.82
0.8	0.76	0.73	0.749	0.72	0.71	0.71

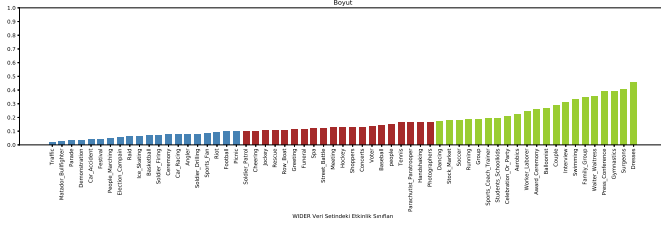
TABLO I: VOC Eğitim seti ve test seti üzerinde raporlanan yüz sezme başarımı.

Eğitilen yüz sezicinin başarımı, yüz sezmede oldukça zorlu bir veri seti olan WIDER FACE veri seti [1] üzerinde de sınanmıştır. WIDER veri seti 60 etkinlik sınıfından toplam 12.500 görüntü üzerinde 159.000 yüz içermektedir ve çok küçük ölçek, örtüşme, başın dönüklüğü gibi yüz sezmeyi zorlaştıran özellikleri nedeniyle başarımlar testlerinde referans alınan veri setlerinden farklıdır. Şekil 4’te eğitilen derin yüz sezici modelinin her bir etkinlik sınıfı için sağladığı performansı sezme oranı cinsinden gösteren histogram raporlanmaktadır. Sezme oranı, etkinlik sınıfında doğru sezilen yüz sayısının etkinlik sınıfında bulunan tüm yüzlerin sayısına oranı şeklinde hesaplanmaktadır. Zorluk derecelerini belirlemek açısından histogramda yatay eksenindeki sınıflar içerdikleri yüz nesnelerinin ölçeklerine göre sıralanmıştır. Küçük ölçekli yüz nesnelerinin çoğunlukta bulunduğu etkinlik sınıfları mavi ile gösterilirken, kırmızı ve yeşil sırasıyla orta ve büyük ölçekli yüz nesnelerini çoğunlukla içeren sınıfları göstermektedir. Şekil 4’te görüldüğü gibi, modelimizin yüz sezme başarımı yüz objesinin ölçeği büyüdükçe artmaktadır. Karşılaştırma açısından Şekil 5’te Viola & Jones [6] yüz sezici ile elde edilen başarımlar, 6’da WIDER veri seti üzerinde elde edilen başarımlar için literatürde raporlanan sonuçlar paylaşılmaktadır. Tüm yöntemler için yüz nesnesinin ölçeği büyüdükçe başarımın yükseldiği görülmektedir. Ancak, Şekil 4 ve 6 karşılaştırıldığında, [1]’de

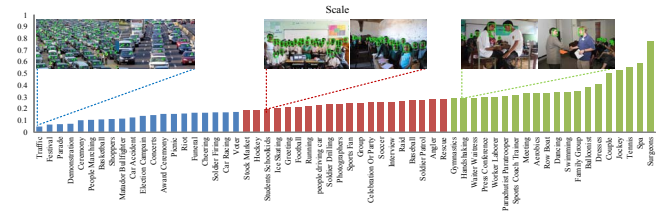
raporlanan yüz sezme başarımının, zor, orta-zor ve kolay olarak nitelenen üç farklı yüz ölçeğindeki sınıflar için sırasıyla % 5-19, %20-27, ve %27-75 aralıklarında kalmasına karşın, çalışmada eğitilen Mask R-CNN yüz sezici ile karşı düşen başarımların sırasıyla %5-26, % 26-40 , %40-78 aralıklarına yükseltilebildiği görülmektedir. Şekil 4 ve 5 karşılaştırıldığında ise eğittiğimiz modelin Viola & Jones yüz sezicisine göre başarımı, küçük ölçekli yüzler için %3 ile %11, orta ölçekli yüzler için %7 ile %28, büyük ölçekli yüzler için ise %12 ile %39 oranında artırdığı görülmektedir.



Şekil 4: Eğitilen Mask R-CNN ile farklı ölçeklerde yüz sezme başarımının etkinlik sınıflarına göre değişimi.



Şekil 5: Viola & Jones [6] yüz sezicisinin yüz sezme başarımı



Şekil 6: [1]'deki yüz sezici ile raporlanan etkinlik sınıfları bazında sezme oranları.

V. SONUÇLAR

Bu çalışmada Mask R-CNN sınırlı sayıdaki yüz gözlemi ile eğitilerek bir derin yüz sezici gerçekleştirilmiştir. Yüz sezme başarımı Pascal VOC ve WIDER FACE veri tabanları üzerinde raporlanmıştır. Literatürde COCO için 80 farklı nesne sezme amacıyla eğitilmiş olan Mask R-CNN'in yüz nesnesini de sezebilmesine olanak sağlanmıştır. Kullanılan eğitim veri setine ilişkin gerçek referans bilgisi ve yazılım rutinleri kullanıma açılmıştır. Çalışmalarımız derin nesne seziciye daha farklı nesnelerin öğretilmesi ve sezme performansının artırılması doğrultusunda devam etmektedir.

Teşekkür: Bu çalışmada Mask R-CNN eğitimine katkısı geçen Müh. Can Elbirlik'e teşekkür ediyoruz.

KAYNAKÇA

- [1] S. Yang, P. Luo, C. C. Loy, and X. Tang, "WIDER FACE: A face detection benchmark," in *IEEE CVPR*, pp.5525–5533, 2016.
- [2] M. Najibi, P. Samangouei, R. Chellappa, and L. S. Davis, "SSH: single stage headless face detector," in *IEEE ICCV*, pp. 4885–4894, 2017.
- [3] S. Zhang, X. Zhu, Z. Lei, H. Shi, X. Wang, and S. Z. Li, "S³fd: Single shot scale-invariant face detector," in *IEEE ICCV*, pp. 192–201, 2017.
- [4] H. Jiang and E. G. Learned-Miller, "Face detection with the faster R-CNN," in *IEEE FG*, pp. 650–657, 2017.
- [5] S. Zhang, R. Zhu, X. Wang, H. Shi, T. Fu, S. Wang, T. Mei, and S. Z. Li, "Improved selective refinement network for face detection," *CoRR*, vol. abs/1901.06651, 2019.
- [6] P. Viola and M. J. Jones, "Robust real-time face detection," *Int. J. Comput. Vision*, vol. 57, no. 2, pp. 137–154, May 2004.
- [7] T. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *IEEE ICCV*, pp. 2999–3007, 2017.
- [8] A. Shrivastava, A. Gupta, and R. B. Girshick, "Training region-based object detectors with online hard example mining," in *IEEE CVPR*, pp. 761–769, 2016.
- [9] V. Jain and E. Learned-Miller, "Fddb: A benchmark for face detection in unconstrained settings," University of Massachusetts, Amherst, Tech. Rep. UM-CS-2010-009, 2010.
- [10] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *Int. Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, Jun. 2010.
- [11] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask R-CNN," in *IEEE ICCV*, pp. 2980–2988, 2017.
- [12] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE CVPR*, pp. 770–778, 2016.
- [13] T. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie, "Feature pyramid networks for object detection," in *IEEE CVPR*, pp. 936–944, 2017.
- [14] R. B. Girshick, "Fast R-CNN," in *IEEE ICCV*, pp. 1440–1448, 2015.
- [15] B. C. Russell, A. Torralba, K. P. Murphy, and W. T. Freeman, "Labelme: A database and web-based tool for image annotation," *Int. Journal of Computer Vision*, vol. 77, no. 1, pp. 157–173, May 2008.
- [16] W. Abdulla, "Mask r-cnn for object detection and instance segmentation on keras and tensorflow," https://github.com/matterport/Mask_RCNN, 2017.
- [17] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision – ECCV 2014*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Cham: Springer International Publishing, 2014, pp. 740–755.