



دانشگاه خوارزمی تهران

دانشکده فنی و مهندسی

پایان نامه کارشناسی ارشد

گرایش هوش مصنوعی

عنوان:

بازشناسی اشیاء در تصاویر دیجیتال با استفاده از روشهای مبتنی بر متن

نگارش:

سهیلا شیخ بهائی

استاد راهنما:

دکتر جمشید شنبه زاده

استاد مشاور:

دکتر زینب قصابی

بهمن 1393

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

به نام خدا
دانشگاه خوارزمی تهران
دانشکده فنی و مهندسی

رساله کارشناسی ارشد

عنوان: بازشناسی شیء با استفاده از روش های مبتنی بر متن

نگارش: سهیلا شیخ بهائی

کمیته ممتحنین:

امضاء: استاد راهنما: دکتر جمشید شنبه زاده

امضاء: استاد مشاور: دکتر زینب قصابی

امضاء: استاد مدعو:

ناریخ:

تقدیم به پدر و مادر عزیز و بزرگوارم

که زحمات بیشمار آنها را هیچگاه نمی توانم جبران نمایم.

و تقدیم به همسر گرانقدرم

که همواره مشوق و همراه و پشتیبان من بوده است.

چکیده:

در دنیای واقعی رابطه‌ای قوی بین محیط و اشیاء وجود دارد که استفاده از اطلاعات مربوط به روابط اشیاء و صحنه نقش مهمی در فرآیند بازیابی شیء ایفا می‌کنند. یعنی زمانی که در یک صحنه شیء خاصی جستجو می‌شود، یک بیننده در موقعیت‌هایی از تصویر تمرکز می‌کند که بیشترین احتمال پیشین را برای وجود شیء مورد نظر دارا هستند. بنابراین متن صحنه برای تصمیم‌گیری در مورد حرکات بعدی چشم بسیار تعیین‌کننده است و زمانی که اطلاعات ظاهری شیء به دلیل کوچک بودن شیء، شلوغی صحنه، وجود مانع و یا هر نوع ابهام دیگری برای بازشناسی آن کافی نباشد، می‌توان از اطلاعات متن تصویر مانند نوع صحنه و اشیاء دیگر صحنه برای بازشناسی شیء کمک گرفت.

به منظور توسعه روشی دقیق برای بازشناسی شیء در تصاویر واقعی و با وجود اشیاء مختلف در این پایان نامه راه حل‌هایی ارائه شده است که بازیابی شیء را با استفاده از اطلاعات متنی یعنی اشیاء دیگر تصویر و صحنه انجام می‌دهد و هدف بهبود دقت بازیابی شیء است. در این روش اطلاعات متنی به صورت ماتریس هم‌رخداد مدل می‌شوند و هر تصویر توسط برداری که نشان‌دهنده‌ی اشیاء دیگر تصویر و فرکانس آنها است نشان داده می‌شود از این بردارها برای آموزش کلاسبند استفاده می‌شود. به دلیل زیاد بودن تعداد داده‌های آموزشی و تست و زمانبر بودن استفاده از روش تطبیق و یا نمایه‌سازی از کلاسبند درخت تصمیم استفاده شده است. پس از آموزش کلاسبند، بردارهای تصاویر تست برای تعیین حضور شیء موردنظر بررسی می‌شوند و دقت بازیابی آنها تعیین می‌شود.

نکته مهم در روش‌هایی که از اطلاعات متنی استفاده می‌کنند این است که این مدل‌ها باید بر تصاویری تست شوند که شامل اشیاء مختلف و از صحنه‌های متنوع باشند تا بتوان از روابط اشیاء و صحنه کمک گرفت. مجموعه داده سان این ویژگی را فراهم می‌کند و دارای تصاویر بیش از 600 صحنه و بیش از 200 شیء است. تعداد اشیاء در هر تصویر بیش از ده کلاس شیء است درحالی‌که در مجموعه داده‌های دیگر تعداد اشیاء یک تصویر حداکثر 5 کلاس شیء است. در این روش از مجموعه داده سان 09 استفاده شده است و میانگین دقت بیش‌بینی حضور شیء با روش‌های مشابه مقایسه شده است و نتایج نشان می‌دهد که متوسط میانگین دقت برای تمام اشیاء در این روش بهبود یافته است.

واژه‌های کلیدی: بازشناسی شیء، اطلاعات متنی، بازیابی شیء، درک صحنه

با تقدیر و تشکر از استاد ارجمندم

جناب آقای دکتر جمشید شنبه زاده

جهت تشویق ها و راهنمایی های بی دریغ ایشان

با تقدیر و تشکر از خانم دکتر زینب قصابی

و خانم مهندس زهرا صادقی

که بسیار مرا راهنمایی و یاری نمودند.

فهرست مطالب

فصل 1:.....	1
مقدمه و طرح تحقیق	1
1-1: مقدمه.....	2
2-1: سیستم های باز شناسی شیء.....	3
1-2-1: تعریف سیستم های باز شناسی شیء.....	3
2-2-1: مراحل عملکرد سیستم باز شناسی شیء.....	4
3-1: اطلاعات متنی.....	5
1-3-1: اطلاعات متنی معنایی	6
2-3-1: اطلاعات متنی مکانی.....	8
3-3-1: اطلاعات متنی مقیاس.....	9
4-3-1: سطوح متنی.....	11
4-4: جایگاه و اهمیت موضوع.....	12
5-1: رویکرد پیشنهادی.....	14
6-1: حوزه مساله.....	15
7-1: ساختار پایان نامه.....	17
فصل 2:.....	18
18: مروری بر کارهای انجام شده.....	18
1-2: مقدمه.....	19
2-2: بررسی الگوریتم های استخراج ویژگی موجود در زمینه باز شناسی شیء.....	19
1-2-2: دسته بندی کلی ویژگیها.....	19
2-2-2: دسته کلمات.....	20
2-2-3: مدل مبتنی بر اجزاء.....	20
3-2: بررسی الگوریتمهای یادگیری سیستم و باز شناسی شیء.....	21
2-3-1: مدل زایشی.....	21
2-3-2: مدل جداکننده.....	21
2-3-3: تابع کلاس بندی.....	22
4-2: روشهای متداول باز شناسی شیء.....	25
5-2: خلاصه و نتیجه گیری.....	31
فصل 3:.....	32
32: پیش زمینه تئوری.....	32
1-3: مقدمه.....	33
2-3: مدل مبتنی بر اطلاعات متنی سلسله مراتبی.....	33

35.....	1-2-3: مدل سازی اولیه
37.....	2-2-3: مدل اندازه گیری
38.....	3-2-3: یادگیری ساختار وابستگی شیء
39.....	4-2-3: یادگیری پارامترهای مدل
40.....	5-2-3: استفاده از مدل
41.....	6-2-3: نتایج
43.....	3-3: مدل مبنی بر اطلاعات متنی معنایی
45.....	4-3: استفاده از درخت تصمیم در بازشناسی شیء
48.....	5-3: معرفی الگوریتم کارت
50.....	1-5-3: معیار تقسیم بندی و مقیاس خالص بودن
51.....	2-5-3: قوانین پایان
51.....	6-3: خلاصه و نتیجه گیری
53.....	فصل 4:
53.....	بررسی و تحلیل الگوریتم پیشنهادی
54.....	1-4: مقدمه
54.....	2-4: چالشهای موجود در روش های قبلی
54.....	1-2-4: چالشهای روشهای بازیابی شیء غیرمتنی
56.....	2-2-4: چالشهای روشهای بازیابی شیء مبتنی بر اطلاعات متنی
57.....	3-4: رویکرد کلی الگوریتم پیشنهادی
59.....	4-4: بررسی دقیق الگوریتم پیشنهادی پایان نامه
62.....	4-5: خلاصه و نتیجه گیری
63.....	فصل 5:
63.....	پایاده سازی، بررسی کارایی و مقایسه با روش های دیگر
64.....	1-5: مقدمه
64.....	2-5: مجموعه داده
69.....	3-5: ارزیابی روش پیشنهادی پایان نامه
69.....	1-3-5: شرایط ارزیابی
70.....	2-3-5: معیارهای ارزیابی
70.....	3-3-5: نتایج ارزیابی
	4-5: جمع بندی 81
82.....	فصل 6:
82.....	نتیجه گیری و کارهای آتی
85.....	فهرست منابع

فهرست شکل‌ها

- شکل 1-1- نمونه ای از یک تصویر و برچسب های اختصاص داده شده توسط سیستم بازشناسی شیء برای دو شیء ماشین و شخص.....4
- شکل 1-2- مراحل کلاسبندی در سیستم بازشناسی شیء مبتنی بر کلاسبند.....4
- شکل 1-3- مثالی از یک نمونه شیء که به تنهایی قابل شناسایی نیست اما با استفاده از اطلاعات متنی صحنه میتوان شیء را بازشناسی کرد[1].....5
- شکل 1-4- مثال تصاویر آموزشی و لایه های معنایی آن7
- شکل 1-5- برنامه تحت وب گوگل و ماتریس متن بدست آمده برای لیست آیتم های پیش بینی شده[1].....8
- شکل 1-6- سمت چپ: کلاسبندهای مکانی متن برای سه کلاس مختلف آسمان، چهره و جاده. سمت راست: فرمول بندی دولایه ی سلسله مراتبی برای اطلاعات مکانی متن[1].....9
- شکل 1-7- مثالی از تصویر آموزشی و تفسیر استفاده شده برای آن[1].....10
- شکل 1-8- یک تصویر آموزشی و چیست آن [1].....11
- شکل 1-9- نمودار تعداد مقالات در زمینه بازیابی شیء در IEEE و Science Direct.....14
- شکل 1-10- مهندسی تصویر شامل سه لایه اصلی پردازش تصویر، آنالیز تصویر و درک تصویر است که فرآیند بازشناسی شیء بر لایه آنالیز و درک تصویر تمرکز دارد[2].....16
- شکل 2-1- نمایش سه شیء به روش کیف کلمات.....20
- شکل 2-2- یک نمونه مثال مدل زایشی. $P(\text{Data}, \text{Zebra})$ نشان دهندهی احتمال گورخر بودن و دادههای ورودی به طور همزمان است و $P(\text{Data}, \text{No Zebra})$ نشان دهندهی احتمال گورخر نبودن و دادههای ورودی است.....21
- شکل 2-3- یک مثال از مدل جداکننده. $p(\text{Zebra}|\text{Data})$ نشان دهندهی احتمال گورخر بودن به شرط دادههای ورودی و $p(\text{No Zebra}|\text{Data})$ نشان دهندهی احتمال گورخر نبودن به شرط دادههای ورودی است.....22

شکل 4-2 - یک مثال از تابع کلاسبندی. در این مثال توسط تابع F و ورودی $Data$ تعیین می شود که آیا گورخر هست یا خیر. 22

شکل 5-2 - جداسازی با حداکثر حاشیه در ماشین بردار پشتیبان 24

شکل 6-2 - بازیابی های بدست آمده با یک مدل شخص با یک عنصر اصلی. مدل با یک فیلتر ریشه (a) چندین فیلتر اجزاء با رزولوشن بالاتر (b) و یک مدل مکانی برای موقعیت هر جزء نسبت به ریشه (c) تعریف شده است. فیلترها وزن های هیستوگرام گرادیان های ویژگی ها را مشخص می کنند [8]. 29

شکل 7-2 - قسمت (a) از مدل ماشین و قسمت (b) از مدل شخص استفاده شده است [9]. 30

شکل 1-3 - تصویر ورودی و نتیجه اعمال 107 یابنده کلاس 34

شکل 2-3 - سمت چپ: خروجی شش یابنده با اطمینان بالاتر، سمت راست: خروجی شش یابنده با استفاده از مدل متنی استفاده شده 35

شکل 3-3 - چپ: مدل اولیه که ارتباط متغیرهای وجود شیء و متغیرهای موقیت شیء را نشان می دهد. راست: مدل اندازه گیری برای شیء i 37

شکل 4-3 - ساختار وابستگی اشیاء که از دیتاست SUN 09 یادگرفته شده است. یال های قرمز نشان دهنده ی رابطه ی منفی بین طبقه ها است. کلفتی هر یال قدرت آن رابطه را نشان می دهد [13] 39

شکل 5-3 - عملکرد تعیین موقعیت و پیشبینی حضور بر مجموعه داده سان 9 و پاسکال 7 a و c درصد تصاویری که N تا از بیشترین بازیابیهای با اطمینان بیشتر صحیح بوده اند. اعداد روی هر ستون تعداد تصاویری که حداقل شامل N نمونه شیء بوده اند، را نشان میدهد. B و d درصد تصاویری که N تا از پیشبینی حضورها با بیشترین احتمال درست بوده اند. اعداد نوشته شده بر هر ستون نشان دهنده ی تعداد تصاویری که حداقل شامل N طبقه شیء مختلف بوده اند را نشان میدهد. 42

شکل 6-3 - طبقه بندی شیء با استفاده از اطلاعات متنی معنایی. $S_1 \dots S_k$ مجموعه k بخش یک تصویر که توسط چندین بخش بندی بدست آمده است را نشان می دهد. $L_1 \dots L_n$ لیست n برچسب هر بخش را نشان می دهد. $O_1 \dots O_m$ نیز مجموعه m طبقه شیء در تصویر اصلی است. 43

- شکل 7-3- ماتریس متنی فرکانس هم‌خدای اشیاء برای مجموعه داده MSRC.....44
- شکل 8-3 - بخش کوچکی از یک نمونه درخت تصمیم [16].....45
- شکل 9-3 - فرآیند یادگیری: استخراج رندم زیرپنجره ها از تصاویر آموزشی، تغییر سایز آنها و سپس ساخت درختان تصمیم [19].....46
- شکل 10-3- فرآیند بازشناسی: زیرپنجره هایی که به صورت رندم استخراج شده اند بین درخت ها منتشر می شوند و پس از جمع آوری رای ها کلاس با رای اکثریت به تصویر تخصیص می یابد [19].
-47
- شکل 1-4 - برخی دشواریهایی که معمولا در بازیابی شیء با آنها روبرو هستیم. تصویر اول بیانگر این مساله است که بعضی اشیاء که دارای ظاهری ناآشنا و مبهم هستند که بازیابی آنها توسط روشهای متداول امکان پذیر نیست. تصویر دوم این واقعیت را بیان میکند که بعضی اشیاء کوچک هستند و به همین دلیل نمی توان آنها را به درستی بازیابی کرد. تصویر سوم نیز ناکافی بودن ویژگیهای محلی در بعضی تصاویر را نشان میدهد که باعث بازیابی خطا میشود.....55
- شکل 2-4- فلوجارت الگوریتم پیشنهادی58
- شکل 3-4 - یک مثال برای چگونگی ذخیره ماتریس فرکانس متنی60
- شکل 4-4- برداری که شامل اشیائی هست که می خواهیم شناسایی شود و برای آموزش به درخت تصمیم داده می شود.60
- شکل 5-4 - درخت تصمیم حاصل از الگوریتم پیشنهادی با در نظر گرفتن 40 شیء.....62
- شکل 1-5- مقایسه‌ی دیتاست پاسکال 7 و سان 9 (a) هیستوگرام تعداد طبقه‌های شیء در هر تصویر (b) توزیع نمونه‌های آموزشی و تست نمونه تصویرها در هر طبقه از شیء (c) نمونه‌هایی از تصاویر پاسکال (d) نمونه‌هایی از تصاویر سان [11].....65
- شکل 2-5- نمونه ای از تصویر مجموعه داده سان از صحنه کلیسا. شکل بالا تصویر و شکل پایین تصویر که تمام اشیاء آن مشخص شده اند را نشان می دهد.....67
- شکل 3-5- لیست اشیاء موجود در تصویر قبل67

شکل 4-5- سه نمونه از تصاویر مجموعه داده پاسکال. تصویر 1 شامل دو شیء هواپیما و شخص است. شکل 2 شامل یک شیء قایق و شکل 3 شامل یک شیء مانیتور است. 68.....

شکل 5-5- دقت 107 شیء بازیابی شده در مدل پیشنهادی. اشیاء براساس دقت بازیابی مرتب شده‌اند. 72.....

شکل 6-5- دقت بازیابی 107 شیء. ستون های قرمز مربوط به الگوریتم پیشنهادی و ستون های آبی مربوط به الگوریتم مبتنی بر اطلاعات متنی سلسله مراتبی است. 73.....

شکل 7-5- مقایسه میانگین دقت الگوریتم پیشنهادی و الگوریتم های قبلی. الگوریتم 1 مدل های جداکننده [22]، الگوریتم 2 الگوریتم مبتنی بر اطلاعات متنی سلسله مراتبی [13] و الگوریتم 3 الگوریتم متنی با کلاس بندی مبتنی بر مجموعه [14] است. 75.....

شکل 8-5- ساختار وابستگی اشیاء که از دیتاست SUN 09 یاد گرفته شده است. یال های قرمز نشان دهنده ی رابطه ی منفی بین طبقه ها است. کلفتی هریال قدرت آن رابطه را نشان می دهد [13] 76.....

شکل 9-5- تصویری از مجموعه داده سان که به عنوان تصویر تست استفاده شده است و دارای 10 کلاس شیء مختلف است. احتمال وجود هواپیما در این تصویر در الگوریتم مبتنی بر اطلاعات سلسله مراتبی 0.05 است اما در الگوریتم پیشنهادی وجود آن پیش بینی شده است. 78.....

شکل 10-5 - تصویری از مجموعه داده سان که به عنوان تصویر تست استفاده شده است و دارای 8 کلاس شیء مختلف است. احتمال وجود هواپیما در این تصویر در الگوریتم مبتنی بر اطلاعات سلسله مراتبی 0.35 است اما در الگوریتم پیشنهادی وجود آن پیش بینی شده است. 79.....

شکل 11-5- تصویری از مجموعه داده سان که به عنوان تصویر تست استفاده شده است و دارای 23 کلاس شیء مختلف است. احتمال وجود توپ در این تصویر در الگوریتم مبتنی بر اطلاعات سلسله مراتبی 0.005 است اما در الگوریتم پیشنهادی وجود آن پیش بینی شده است. 79.....

شکل 12-5- ماتریس تداخل برای 10 شیء با بیشترین دقت بازیابی از مجموعه داده سان 80.....

فهرست جدول ها

- جدول 1-3- تعریف متغیرهای استفاده شده در تعریف الگوریتم ساخت درخت تصمیم.....73
- جدول 1-5- دقت نمونه هایی از 107 شیء بازیابی شده توسط الگوریتم پیشنهادی.....73
- جدول 2-5- دقت بازیابی دو الگوریتم با بهترین عملکرد بر مجموعه داده سان.....74

فهرست علائم و اختصارات

SVM	Support Vector Machine
HOG	Histogram of Gradient
BOW	Bag of Words
DPM	Discriminatively Trained Part Based Model
SIFT	Scale Invariant Feature Transform
DOG	Difference of Gaussians

فصل ۱:

مقدمه و طرح تحقیق

۱-۱: مقدمه

امروزه علاقه و توجه به سیستم های بینایی بدون توقف ادامه دارد و دلیل آن پتانسیل سیستم های بینایی اتوماتیک در افزایش شدید ظرفیت تولید سازمان می باشد. عمدتاً اساسی ترین بخش یک عامل برمبنای بینایی واحد شناسایی شیء آن است. تحقیقات در زمینه ی الگوریتم های بازشناسی شیء باعث پیشرفت هایی در خودکارسازی بسیاری از اعمال با استفاده از ایجاد سیستم های بازشناسی حروف، سیستم های بازرسی صنعتی خط تولید و همچنین سیستم های شناسایی خطا شده است. همچنین باعث پیشرفت های بزرگی در تصویربرداری پزشکی، دفاع و بیومتریک ها شده است [1].

پیشرفت های اخیر در بازشناسی شیء و درک صحنه باعث افزایش قدرت برنامه های کاربردی اصلاح تصویر نیمه خودکار شده است. از دیگر سیستم های کاربردی در این زمینه می توان به سیستم تکمیل صحنه اشاره کرد که در آن هدف حذف بخشی از صحنه و پرکردن آن با استفاده از بخش هایی از تصاویر دیگر است [2]. بازشناسی شیء در موارد مختلفی از جمله سیستم دستیار راننده، سیستم های کنترل خط تولید، رباتیک و پردازش تصاویر پزشکی نیز کاربرد دارد.

توانایی انسان برای بازشناسی هزاران دسته شیء در صحنه های شلوغ با وجود تغییرات در حالت و روشنایی و یا وجود مانع یکی از حیرت انگیزترین قابلیت های ادراک بصریست که هنوز توسط الگوریتم های بینایی کامپیوتر قابل دستیابی نیست. هدف بهبود فرآیند بازشناسی شیء با وجود کاربردهای مختلف و متنوع آن در زمینه های گوناگون محققین را واداشته است تا در این زمینه به تحقیقات گسترده ای بپردازند.

در بخش دوم این فصل سیستم بازشناسی شیء تعریف می شود و مراحل عملکرد آن معرفی می شود. در بخش سوم اطلاعات متنی را تعریف می کنیم و دلیل استفاده از آن و اهمیت استفاده از آن را توضیح می دهیم. همچنین انواع اطلاعات متنی و سطوح متنی را تعریف می کنیم و با ذکر مثال آن ها را توضیح می دهیم و شاخه ای که در این پایان نامه مدنظر است مشخص می شود.

در بخش چهارم این فصل در مورد جایگاه و اهمیت این موضوع توضیح داده می شود و اهمیت تحقیق و چالش ها مشخص می شود. بخش پنجم شامل هدف و سوال اصلی تحقیق می باشد. بخش ششم نیز روش تحقیق و گردآوری اطلاعات و حوزه مساله و جنبه نوآوری تحقیق را بیان می کند. در نهایت نیز در بخش هفتم این فصل ساختار رساله و اینکه هر فصل شامل چه مباحثی می باشد بیان شده است.

۱-۲: سیستم های بازشناسی شیء

در این بخش سیستم بازشناسی شیء تعریف می شود و چرخه کلی آن توضیح داده می شود و سپس بخش های مختلف آن معرفی و مراحل عملکرد آن ذکر می شود. در مرحله ی عملکرد یک سیستم بازشناسی شیء، دو مرحله کلی استخراج اطلاعات از تصویر و تعیین استراتژی بازشناسی شیء مبنی بر این اطلاعات بیان می شود. در فصل های آینده این مراحل با جزئیات بیشتری توضیح داده خواهند شد.

۱-۲-۱: تعریف سیستم های بازشناسی شیء

مساله بازشناسی شیء را می توان به عنوان یک مساله برچسب گذاری تعریف کرد. به این صورت که مدل هایی از اشیائی که باید بازشناسی شوند، در سیستم تعریف می شوند و سپس براساس این مدلها تصویر ورودی به سیستم، برچسب گذاری میشود. در این سیستم به دو سوال باید پاسخ داده شود:

✓ آیا شیء موردنظر در تصویر وجود دارد یا خیر؟

✓ در صورتی که وجود دارد موقعیت آن کجاست؟

در نتیجه ورودی تصویری است که شامل صفر یا تعداد بیشتری شیء موردنظر است. یک سری برچسب که متعلق به مجموعه ای از مدلهاست به سیستم شناسانده شده است، حال سیستم باید برچسب های درست را به نواحی درست تصویر تخصیص دهد. شکل 1-1 نمونه ای از یک تصویر و برچسب های اختصاص داده شده توسط سیستم بازشناسی شیء را برای دو شیء ماشین و شخص نشان می دهد.

برای این کار ابتدا تصاویری که دارای برچسب هستند برای آموزش وارد سیستم می شوند. پس از مرحله آموزش تصاویری بدون برچسب وارد سیستم می شوند و پاسخ بدست آمده با پاسخ صحیح مقایسه می شود تا دقت سیستم بدست آید.

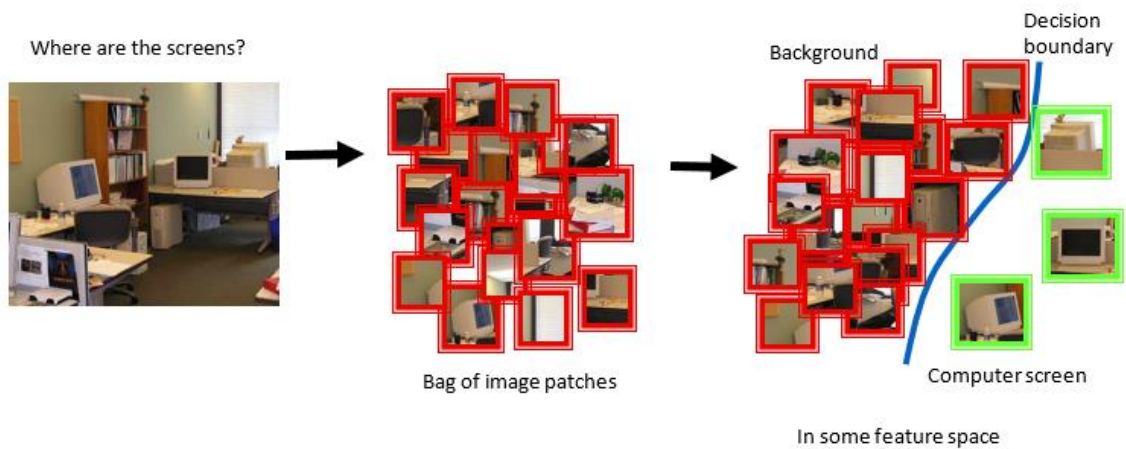


شکل 1-1 - نمونه ای از یک تصویر و برجسب های اختصاص داده شده توسط سیستم بازشناسی شیء برای دو شیء ماشین و شخص

۱-۲-۲: مراحل عملکرد سیستم بازشناسی شیء

چهارچوب کلی یک سیستم بازشناسی شیء شامل دو بخش است: استخراج اطلاعات از تصویر و دوم استفاده از این اطلاعات برای بازشناسی شیء براساس استراتژی مورد استفاده. در نتیجه در بازشناسی شیء دو مساله مهم وجود دارد: اول چگونگی بازنمایی یک طبقه شیء، و دوم چگونگی یادگیری سیستم و بازشناسی شیء مورد نظر که عمدتاً شامل ساخت یک کلاسبند با استفاده از داده های آموزشی و استفاده از این کلاسبند برای کلاسبندی داده های جدید است.

می توان یک مساله بازشناسی شیء را به عنوان یک مساله کلاسبندی در نظر گرفت. تصویر به پنجره-هایی که بایکدیگر همپوشانی دارند تقسیم می شود و سپس تصمیم گیری می شود که آیا پنجره مورد نظر دارای شیء هدف هست یا خیر. شکل 1-2 این مراحل را نشان میدهد.



شکل 1-2 - مراحل کلاسبندی در سیستم بازشناسی شیء مبتنی بر کلاسبند

روش‌های مختلفی جهت استخراج اطلاعات از تصویر برای فرآیند بازشناسی شیء وجود دارد. می‌توان این روش‌ها را به دو دسته تقسیم کرد: روش‌هایی که بر مبنای اطلاعات شیء هستند و روش‌هایی که بر مبنای اطلاعات صحنه و اشیاء دیگر تصویر هستند. اطلاعات مربوط به اشیاء دیگر و صحنه اطلاعات متنی¹ نامیده می‌شوند که در بخش بعد به تفصیل بررسی خواهند شد.

۱-۳: اطلاعات متنی

مساله‌ی اصلی در شناسایی شیء تطبیق نمایی از شیء هدف با خصوصیات موجود تصویر و رد کردن خصوصیات زمینه است. در جستجوهای بصری معمولی متن صحنه مجموعه‌ای از مفهومی‌هایی رندم است که تنها فرآیند بازیابی شیء را سخت‌تر می‌کند، در حالی که در دنیای واقعی اشیاء دیگر در یک صحنه منبعی غنی از اطلاعات هستند که به جای مانع برای تشخیص و شناسایی اشیاء به این فرآیند کمک می‌کنند.

مدل‌های زیادی از اطلاعات متن اشیاء برای بهبود دقت بازشناسی استفاده می‌کنند. خصوصیات ظاهری مانند رنگ، لبه، بافت و علائم شکلی می‌توانند تا حدود مشخصی کلاس‌های اشیاء را تشخیص دهند. در صورت وجود شلوغی، نویز و تغییرات در حالت و روشنایی، ظاهر شیء دارای ابهام خواهد بود که برای رفع این ابهام می‌توان از پیوستگی اشیاء در یک صحنه استفاده کرد (شکل 3-1).



شکل 3-1 - مثالی از یک نمونه شیء که به تنهایی قابل شناسایی نیست اما با استفاده از اطلاعات متنی صحنه می‌توان شیء را بازشناسی کرد [3].

یکی از روش‌های متداول برای بازیابی و بازشناسی شیء استفاده از اطلاعات متنی است. استفاده از اطلاعات متنی به معنای استفاده از اطلاعات اشیاء دیگر تصویر و یا اطلاعات صحنه برای بازیابی شیء

¹ Contextual information

خاص است. معمولا در روش های بازشناسی و بازیابی شیء، اشیاء به تنهایی در نظر گرفته می شوند و صحنه ای که شیء در آن واقع است در نظر گرفته نمی شود. این یک محدودیت جدی است زیرا اطلاعات متنی نقش مهمی در عمل بازشناسی شیء در سیستم بینایی انسان ایفا می کنند. استفاده از اطلاعات متنی باعث بهبود عملکرد سیستم های بازشناسی شیء می شوند و همچنین راهنماهای معنایی مفیدی را برای درک صحنه فراهم می کنند.

خصوصیات متنی را می توان به سه دسته ی معنایی^۲، مکانی^۳ و مقیاس^۴ دسته بندی کرد. مدل های بازشناسی شیء استفاده از اطلاعات متنی را از دو سطح سراسری و محلی تصویر در نظر گرفته اند. اطلاعات سراسری آماره های تصویر را در یک تصویر به عنوان یک کل در نظر می گیرد. اطلاعات محلی متنی اطلاعات متنی را برای ناحیه های همسایه شیء در نظر می گیرد.

۱-۳-۱: اطلاعات متنی معنایی

اطلاعات معنایی متن به احتمال وجود یک شیء در بعضی صحنه ها مربوط می شود. بنابراین می توان این اطلاعات را به صورت احتمال رخداد یک شیء با اشیاء دیگر و یا در یک صحنه خاص تعریف کرد.

مطالعات اخیر در روانشناسی و علوم شناختی^۵ نشان می دهد که اطلاعات معنایی متن به بازشناسی بصری در ادراک انسان کمک می کند. اطلاعات معنایی متن به صورت قوانین از پیش تعیین شده تعریف می شوند. انتظار اولیه ی سیستم از جهان توسط فرضیه های مختلف (استراتژی های بر مبنای قوانین) نمایش داده می شود. که وجود اشیاء دیگر در یک صحنه را پیش بینی می کند. این فرضیه ها توسط متخصصین بازشناسی انواع شیء تولید می شود.

در بعضی از روش های اخیر از روش های آماری برای عمومیت دادن و بهره برداری از اطلاعات معنایی متن استفاده شده است. همان طور که در شکل 4-1 نشان داده شده است، ولف^۶ و بیلسچی^۷ از لایه های معنایی تصاویر آموزشی برای بدست آوردن اطلاعات معنایی متن استفاده کردند. لایه های معنایی وجود یک شیء مشخص را در یک تصویر نشان می دهند [3].

² Semantic context

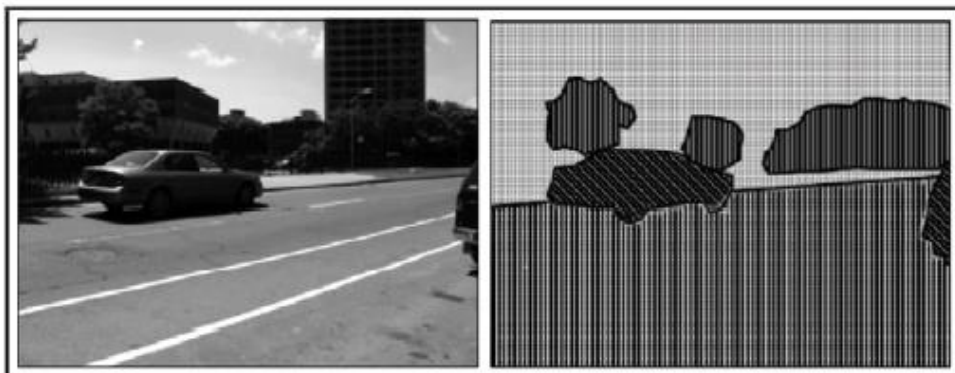
³ Spatial context

⁴ Scale context

⁵ Cognitive science

⁶ Wolf

⁷ Bileschi



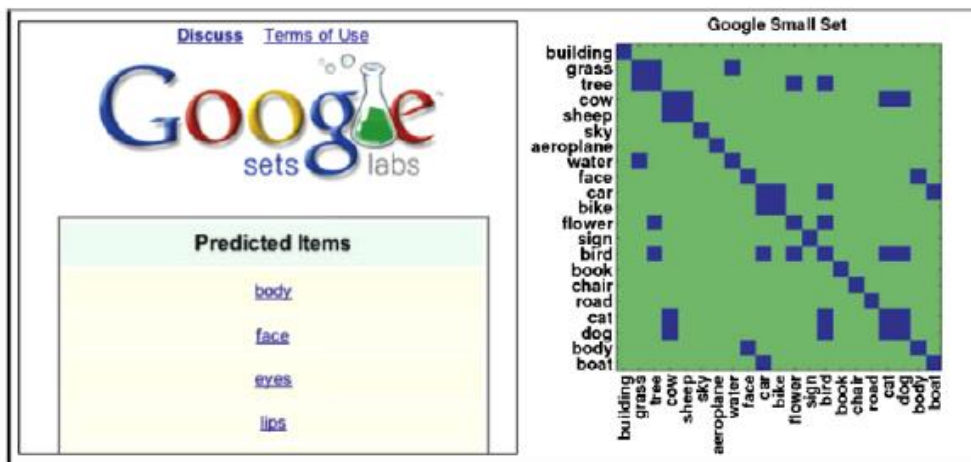
شکل 4-1- مثال تصاویر آموزشی و لایه های معنایی آن

هر تصویر چندین لایه معنایی را نشان می دهد. در یک لایه ی معنایی هر پیکسل با مقدار 1 یا صفر مقداردهی می شود که 1 نشان دهنده ی این است که این پیکسل متعلق به شیء هست و صفر نشان دهنده ی این است که پیکسل متعلق به شیء نیست . اطلاعات محتوایی متن به صورت لیستی از برچسبها نمایش داده می شود که وقوع یک پیکسل در یک شیء مشخص را نشان می دهد.

اطلاعات متن را می توان از داده های آموزشی برچسب گذاری شده بدست آورد. روبینویچ و همکارانش⁸ اطلاعات محتوایی متن را توسط پرسش از گوگل بدست آورد. گوگل لیستی از آیتم های به هم مرتبط را از تعداد کمی مثال بدست می آورد. این اطلاعات توسط ماتریس همرخداد⁹ باینری نشان داده می شود . در صورتی که شیء 1 با شیء 2 رابطه داشته باشد، مولفه مربوطه در ماتریس 1 و در غیر این صورت صفر است. شکل 5-1 ماتریس مورد استفاده را نشان می دهد.

⁸ Robinovich et al

⁹ Co-occurrence



شکل 5-1 - برنامه تحت وب گوگل و ماتریس متن بدست آمده برای لیست آیتم های پیش بینی شده [3]

۱-۳-۲: اطلاعات متنی مکانی

اطلاعات مکانی متن را می توان به صورت احتمال یافتن یک شیء در موقعیت های خاص با توجه به اشیاء دیگر در صحنه تعریف کرد. حضور اشیائی که تفسیر واحدی دارند بازشناسی اشیاء مبهم در صحنه را بهبود می دهد و روابط مکانی مناسب بین اشیاء نرخ خطا را در بازشناسی اشیاء کاهش می دهد [3].

یکی از کارهای اولیه ای که توسط فیشر¹⁰ انجام شد شمایی پایین به بالا را پیشنهاد می دهد که اشیاء مختلف را در صحنه شناسایی می کند. برای این کار ابتدا تصویر به نواحی بخش بندی می شود، هر ناحیه به عنوان یک شیء برچسب گذاری می شود و سپس برچسب های اشیاء با استفاده از اطلاعات مکانی متن تصحیح می شود. برای تصحیح می توان اشیاء را به بخش های کوچکتری تقسیم بندی کرد و محدوده ی امکان پذیر روابط مکانی این بخش ها را تعریف کرد. اطلاعات مکانی متن به شکل قوانین و ساختارهای شبیه گراف ذخیره شده اند.

یکی از کارهای اخیر که توسط شاتن و همکارانش¹¹ انجام شده است از اطلاعات بافت، چیدمان و اطلاعات مکانی استفاده می کند. اطلاعات مکانی متن توسط کار کومار¹² و هبرت¹³ نیز معرفی شده است. روش آنها فرمولی سلسله مراتبی دولایه نمایش می دهد. همان طور که در شکل 6-1 نشان داده

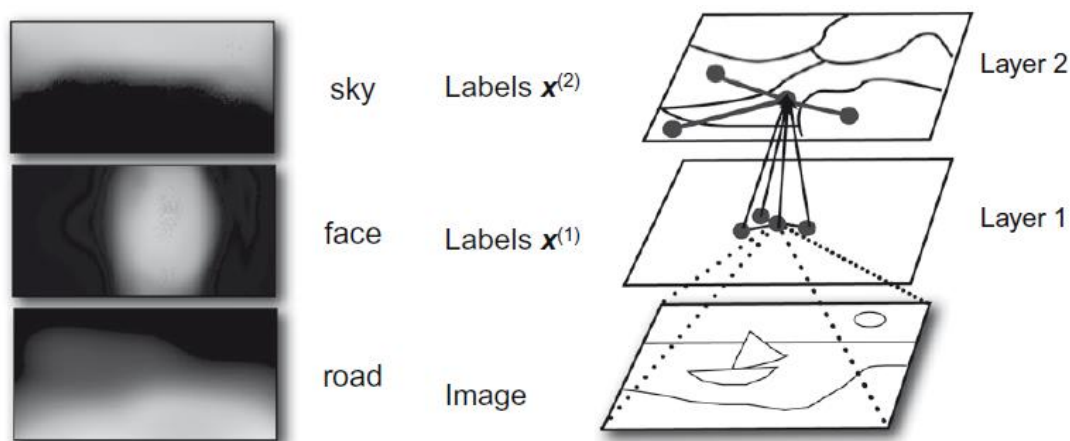
¹⁰ Fischler

¹¹ Shotton et al

¹² Kumar

¹³ Hebert

شده است، لایه اول ارتباط متقابل ناحیه به ناحیه را مدل می کند و لایه دوم ارتباط متقابل شیء با شیء را مدل می کند [3].



شکل 6-1 - سمت چپ: کلاسبندهای مکانی متن برای سه کلاس مختلف آسمان، چهره و جاده. سمت راست: فرمول بندی دولایه ی سلسله مراتبی برای اطلاعات مکانی متن [3]

کارهای اخیر در بینایی کامپیوتر از داده های آموزشی همراه با تفسیر برای بدست آوردن اطلاعات مکانی متن استفاده می کنند. بعضی از روش ها نیز از قوانین از پیش تعریف شده استفاده کرده اند.

۳-۳-۱: اطلاعات متنی مقیاس

روش های متداول بازشناسی شیء نیاز به جستجوی فضای جستجوی بزرگ برای جستجوی مدل ها، موقعیت ها و اندازه های مختلف شیء داشتند. اطلاعات پیشین در مورد سائزهایی که اشیاء در صحنه ها ظاهر می شوند، فرآیند بازیابی شیء را بهبود می بخشد و نیاز به جستجوی اندازه های مختلف را کاهش می دهد. در نتیجه منابع محاسباتی را به مقیاس های محتمل تر اختصاص می دهد [3].

برای بدست آوردن اطلاعات مقیاس متن علاوه بر شناسایی حداقل یک شیء دیگر، نیاز به پردازش ارتباطات مکانی و عمق بین شیء مورد نظر و این اشیاء دیگر است. سیستم کندر¹⁴ که توسط استرت¹⁵ و فیشر¹⁶ ارائه شد از اولین سیستم های بینایی کامپیوتر بود که اطلاعات مقیاس را به عنوان یک

¹⁴ CONDOR

¹⁵ Strat

¹⁶ Fischler

ویژگی در نظر گرفت. اطلاعات مقیاس یک شیء از طریق موقعیت دوربین، جهت، فضای هندسی، ارتفاع زمین و نقشه محاسبه می شود.

تورالبا و همکارانش¹⁷ چهارچوبی ساده برای مدل کردن رابطه بین متن و خصوصیات شیء معرفی کرد. اطلاعات متنی از طریق تصاویر آموزشی که در آن خصوصیات شیء بر مبنای همبستگی های بین آماره های خصوصیات سطح پایین کل صحنه است، یاد گرفته می شوند. شکل 1-7 مثالی از یک تصویر آموزشی و تفسیر مربوط به آن را نشان می دهد.



شکل 1-7- مثالی از تصویر آموزشی و تفسیر استفاده شده برای آن [3]

اطلاعات مقیاس متن سخت ترین ارتباط قابل دستیابی است زیرا نیاز به اطلاعات جزئی در مورد شیء در صحنه دارد.

بیشتر مدل ها از یک یا دو نوع اطلاعات متنی استفاده می کنند. اطلاعات متنی مکانی و مقیاس بیشترین استفاده را داشته اند. به طور کلی اطلاعات متن معنایی به طور ضمنی در اطلاعات مکانی متن موجود است. همچنین در مورد مقیاس با توجه به اشیاء دیگر در نظر گرفته می شود. بنابراین استفاده از اطلاعات مکانی و مقیاس متن استفاده از تمام اطلاعات متنی را شامل می شود.

هرچند اطلاعات معنایی متن توسط انواع دیگر قابل دستیابی است اما تنها نوع اطلاعات متنی است که ارزشمندترین اطلاعات را به ما می دهد و همچنین دسترسی به آن بسیار ساده تر است و پردازش و محاسبات آن سریع تر است.

¹⁷ Torralba et al

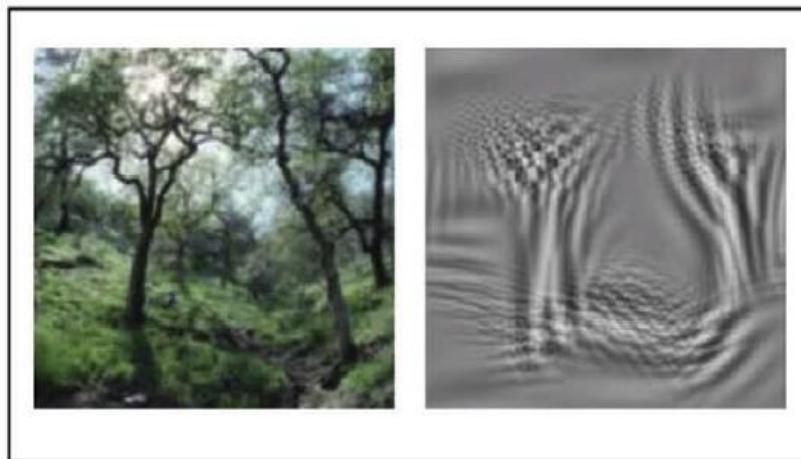
۱-۳-۴: سطوح متنی

مدل های بازشناسی شیء استفاده از اطلاعات متنی را از دو سطح سراسری و محلی تصویر در نظر گرفته‌اند. اطلاعات سراسری آماره های تصویر را در یک تصویر به عنوان یک کل در نظر می‌گیرد (مثلا یک آشپزخانه حضور یک فر را پیش بینی می‌کند). اطلاعات محلی متنی اطلاعات متنی را برای ناحیه های همسایه شیء در نظر می‌گیرد (مثلا یک چراغ خواب حضور یک ساعت شماتپه ای را پیش بینی می‌کند)[3].

۱.۳.۴.۱ اطلاعات متنی سراسری

ساختار یک صحنه را می‌توان با میانگین گرفتن ویژگی های سراسری تصویر بدست آورد، که خلاصه‌ای آماری از چیدمان مکانی فراهم می‌کند. بسیاری از چهارچوب های طبقه بندی شیء از این اطلاعات پیشین برای تعیین موقعیت استفاده کرده‌اند.

مورفی و همکارانش¹⁸ از ویژگی های متنی با استفاده از "جیست"¹⁹ صحنه استفاده کرده است. جیست یک صحنه نمایشی کلی و سطح پایین از کل تصویر است. شکل 1-8 یک مثال از جیست یک جنگل را نشان می‌دهد. توراتا و همکارانش نشان می‌دهند که فراهم کردن اطلاعات پیشین در مورد اینکه چه اشیائی در یک صحنه و در چه موقعیت و اندازه ای ظاهر خواهند شد، کافیست.



شکل 1-8 - یک تصویر آموزشی و جیست آن [3]

¹⁸ Murphy et al

¹⁹ Gist (Global Invariant Scale Transform)

استفاده از اطلاعات متنی با پردازش صحنه به عنوان یک کل و بدون بازیابی ابتدایی اشیاء دیگر به کاهش بازیابی های غلط می انجامد.

۱.۳.۴.۲ اطلاعات متنی محلی

اطلاعات متنی محلی از مناطق اطراف شیء شامل اشیاء دیگر، پیکسل ها و تکه ها بدست می آید. برای انجام این کار ابتدا باید شیء شناخته شود و قسمت هایی که شباهت لازم با شیء را نداشته باشند کنار گذاشته می شوند. زمانی که چهارچوب بازشناسی دارای تصاویری باشد که شیء در مرکز آن قرار داشته باشد، استفاده از اطلاعات متنی محلی باعث بهبود بازشناسی خواهد شد [3].

یکی از مزایای اطلاعات متنی سراسری نسبت به محلی این است که از نظر محاسباتی کم هزینه تر است. البته زمانی که تعداد اشیائی که می خواهیم بازشناسی کنیم زیاد باشد اطلاعات متنی سراسری نمی تواند بین صحنه های مختلف تمایز قائل شود.

نمونه های متعددی از این ایده کلی تاکنون استفاده شده است. مدل های مبتنی بر بخش نوعی از اطلاعات متنی محلی را استفاده می کنند که بخش های مختلف باید با رابطه ی هندسی مناسبی نسبت به یکدیگر قرار بگیرند تا یک شیء را تشکیل دهند.

بسیاری از روش ها از چیست صحنه برای تصمیم گیری در مورد مکان محتمل برای اشیاء مختلف استفاده می کنند. برای مثال تورالبا و همکارانش سیستمی برای پیش بینی موقعیت عمودی اشیاء مانند افراد پیاده رو، ماشین ها و ساختمان ها براساس چیست صحنه آموزش داده اند. این مکان ها سپس توسط یابنده های شیء کلاسیک استفاده شده اند که عملکرد یابنده ها را بهبود می بخشند.

۴-۱: جایگاه و اهمیت موضوع

تکنولوژی های بازشناسی در صنعت غذا (مثلا برای خودکارسازی دسته بندی محصولات کشاورزی)، در صنعت خودرو و الکترونیک (برای اتوماتیک سازی مونتاژ و هدف های بازبینی صنعتی)، در صنعت داروسازی (برای دسته بندی قرص ها و کپسول ها)، در بهبود جستجوی آنلاین، در عکس برداری پیشرفته، در سیستم های بازشناسی بیومتریک بر مبنای عنبیه یا اثر انگشت و همچنین شناسایی دستخط، در صنعت روبات، در وسایل نقلیه هوشمند و بسیاری کاربردهای دیگر استفاده می شوند.

بسیاری از مدل های نمایش اشیاء به طور موثری توسط جامعه‌ی تصویربرداری پزشکی برای بخش بندی قدرتمند ساختارهای تشریحی مانند مغز و حفره های قلب استفاده می‌شوند.

روش های متداول بازشناسی شیء بر تصاویری آزموده می‌شوند که با تصاویر دنیای واقعی متفاوت بوده و معمولا دارای تعداد کمی شیء و در صحنه های خاصی می‌باشند. مجموعه داده ی سان^{۲۰} دارای تصاویر از صحنه های متنوع می‌باشد و تعداد اشیاء موجود در این تصاویر نیز بسیار زیاد می‌باشند [4]. روش های بازشناسی متداول بر این مجموعه داده دقت پایینی را ارائه می‌دهند و دلیل آن پیچیدگی بیشتر فرآیند بازشناسی شیء به دلیل تنوع صحنه ها و شلوغی آنهاست. به همین دلیل استفاده از اطلاعات متنی جایگاه و اهمیت ویژه ای می‌یابد.

به دلیل اهمیت موضوع بازشناسی شیء مسابقات زیادی در سراسر دنیا در این زمینه برگزار می‌شود. یکی از این مسابقات رقابت پاسکال^{۲۱} است که از سال 2005 هر ساله برگزار شده است. یکی دیگر از این رقابت ها ایمیج نت^{۲۲} است که هدف پیش بینی محتوای عکس ها برای انجام حاشیه نویسی خودکار و بازیابی تصاویر دیتاست ایمیج نت برای آموزش است. این رقابت نیز از سال 2010 هر ساله برگزار می‌شود و در بخش های مختلف عملکرد الگوریتم های شرکت داده شده سنجیده می‌شود و تعدادی از بهترین الگوریتم ها به عنوان برنده انتخاب می‌شوند.

روند تحقیقات در زمینه بازیابی و بازشناسی شیء در شکل 9-1 نشان داده شده است. همانطور که دیده می‌شود تعداد مقالات در دو پایگاه اطلاعاتی معتبر IEEE²³ و Science direct²⁴ از سال 2003 تا 2013 که بررسی شده است، روند روبه رشدی را داشته است. افزایش تعداد مقالات نشان دهنده ی اهمیت این موضوع و نیاز به تلاش و تحقیق در این زمینه است.

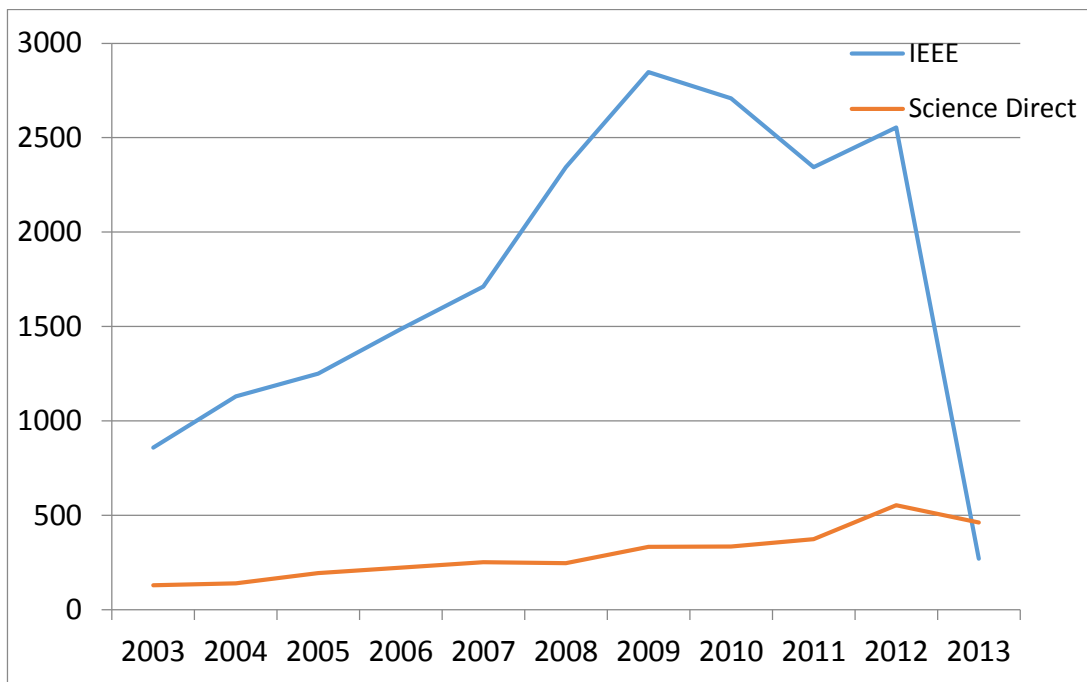
²⁰ SUN <http://groups.csail.mit.edu/vision/SUN/>

²¹ PASCAL <http://pascallin.ecs.soton.ac.uk/challenges/VOC/>

²² ImageNet <http://www.image-net.org/>

²³ <http://ieee.org/>

²⁴ <http://sciencedirect.com/>



شکل 9-1- نمودار تعداد مقالات در زمینه بازیابی شیء در IEEE و Science Direct

۵-۱: رویکرد پیشنهادی

در این پایان نامه هدف بهبود دقت فرآیند بازیابی شیء با استفاده از اطلاعات متنی معنایی است. سوال اصلی تحقیق این است که از چه اطلاعات معنایی استفاده شود و چگونه این اطلاعات مدل شوند که بتوان به کمک آنها بازیابی شیء انجام داد و بتوان بر مجموعه داده مناسب نتایج قابل قبول بدست آورد. مجموعه داده مناسب مجموعه داده ای است که دارای تصاویر از اشیاء مختلف و در صحنه های متنوع باشد و نتایج قابل قبول نتایجی هستند که نسبت به نتایج ارائه شده قبلی بر همان مجموعه داده برتری داشته باشند. یکی از راه های استفاده از اطلاعات معنایی استفاده از ماتریس همربخداد است که در این روش با اعمال تغییراتی از این ماتریس استفاده شده است. جهت کلاسبندی نیز از ماشین بردار پشتیبان^{۲۵} استفاده شده است.

برای ارزیابی عملکرد این الگوریتم از مجموعه داده سان استفاده شده است. دقت بازیابی هر شیء به صورت نسبت پیش بینی های شیء انجام شده که درست بوده اند، تعریف می شود. برای هر شیء

²⁵ Support Vector Machine (SVM)

میانگین دقت پیش‌بینی‌هایی که در هر مرحله بدست آمده است به عنوان میانگین دقت^{۲۶} پیش‌بینی حضور آن شیء در نظر گرفته می‌شود.

در نهایت نیز برای ارزیابی الگوریتم، متوسط میانگین دقت^{۲۷} تمام اشیاء محاسبه و به عنوان معیار ارزیابی در نظر گرفته می‌شود. در دومقاله ای که با الگوریتم پیشنهادی مقایسه شده اند 107 شیء بررسی و میانگین دقت آنها محاسبه شده است. برای اینکه بتوان نتایج را مقایسه و الگوریتم را ارزیابی کرد در این الگوریتم نیز میانگین دقت برای همان 107 شیء محاسبه شده است.

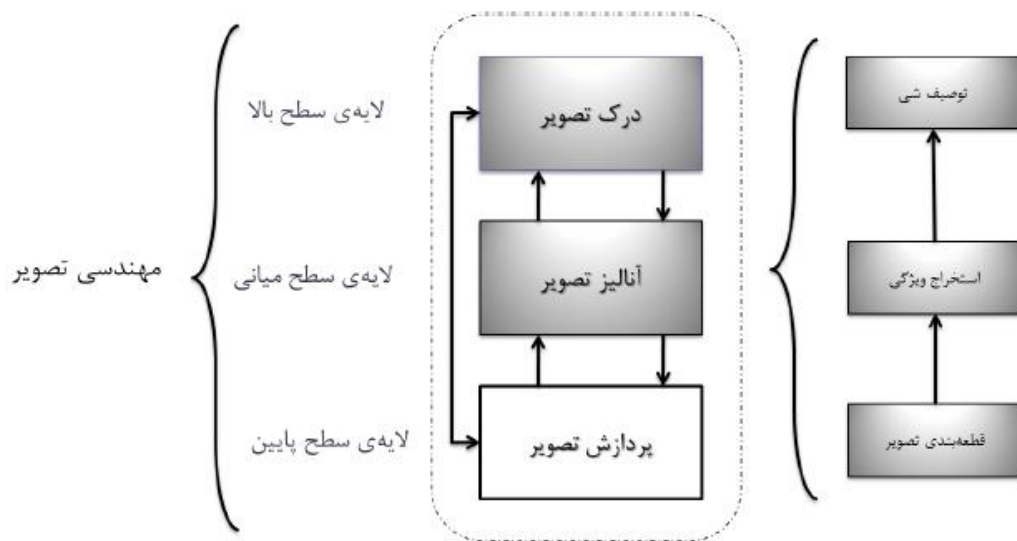
۱-۶: حوزه مساله

یکی از مراحل مهم در انجام پایان نامه مشخص کردن حوزه مساله است تا از انجام فعالیت های غیرمرتبط با مساله جلوگیری شود، به همین دلیل در این بخش توضیح داده می‌شود که تحقیقات این پایان نامه در کدام حوزه ها قرار می‌گیرد.

فرآیند مهندسی تصویر در شکل 10-1 نشان داده شده است [5]. در این مراحل ابتدا در لایه سطح پایین تصویر پردازش می‌شود و نویز و مشکلات احتمالی تصویر برطرف می‌شود. این لایه در حوزه تحقیق قرار ندارد و تصویر ورودی تصویری آماده و پردازش شده است. اما دو لایه دیگر شامل آنالیز تصویر و درک تصویر در حوزه مساله قرار می‌گیرد و برای بازشناسی شیء آنالیز و درک تصویر انجام می‌شود.

²⁶ AP(Average precision)

²⁷ Mean AP



شکل 10-1- مهندسی تصویر شامل سه لایه اصلی پردازش تصویر، آنالیز تصویر و درک تصویر است که فرآیند بازشناسی شی بر لایه آنالیز و درک تصویر تمرکز دارد [5].

اگر روش های بازشناسی شی را به دو دسته مبتنی بر اطلاعات شی و مبتنی بر اطلاعات متنی تقسیم کنیم، روش مبتنی بر اطلاعات متنی در حوزه این تحقیق قرار می گیرد و هدف یافتن اطلاعات متنی مناسب و مدل سازی آنها برای بازیابی شی است تا بتوان دقت بازشناسی اشیاء را افزایش داد.

تصاویر موردنیاز این تحقیق از مجموعه داده سان 09 فراهم شده است که شامل بیش از 600 صحنه و بیش از 200 شی است. بیشتر تصاویر این مجموعه داده شامل بیش از ده کلاس شی است، درحالی که بیشتر تصاویر در مجموعه داده های دیگر حداکثر در هر تصویر دارای پنج کلاس شی هستند.

دو مقاله ای که تاکنون بهترین عملکرد را در حوزه بازشناسی شی داشته اند برای مقایسه با روش پیشنهادی انتخاب شده اند و نتایج با هم مقایسه و در فصل های بعدی شرح داده شده اند. در این دو مقاله میانگین دقت بازیابی شی برای 107 کلاس شی محاسبه شده اند و نتیجه نهایی براساس میانگین دقت بازیابی 107 شی مقایسه می شوند که هدف این تحقیق افزایش میانگین دقت بازیابی 107 کلاس شی مشخص شده بر این مجموعه داده است. حاصل این تحقیق می تواند به صورت نرم افزار کاربردی در صنعت و خدمات استفاده شود و برای تصاویر متنوع از صحنه های مختلف و با تعداد اشیاء زیاد کاربرد دارد.

۷-۱: ساختار پایان نامه

در این پایان نامه پنج فصل تنظیم شده است. در فصل دو ابتدا به معرفی بخش های مختلف سیستم بازشناسی شیء می پردازیم و الگوریتم های مورد استفاده در هر بخش را معرفی می کنیم.

در فصل سوم روش های مطرح در این زمینه معرفی شده و نقاط قوت و ضعف آنها بررسی می شود، به ویژه بر روش های مبتنی بر اطلاعات متنی و مبتنی بر درخت تصمیم تمرکز شده است.

در فصل چهارم چالش های روش های قبلی مطرح می شوند و رویکرد کلی الگوریتم پیشنهادی پایان نامه گفته می شود و سپس به طور دقیق بخش های مختلف آن توضیح داده می شوند. پس از آن مجموعه داده مناسب انتخاب و علت انتخاب آن ذکر می شود.

در فصل پنجم به بررسی عملکرد و ارزیابی الگوریتم پیشنهادی پرداخته می شود و نتایج دو روش با بهترین عملکرد ذکر می شوند و با الگوریتم پیشنهادی مقایسه می شوند و بررسی می شود که آیا این الگوریتم دقت میانگین بازشناسی اشیاء را افزایش داده است یا خیر.

در نهایت در فصل آخر نیز به نتیجه گیری و جمع بندی مطالب پرداخته می شود و کارهای آتی بیان می شوند.

فصل ۲:

مروری بر کارهای انجام شده

۲-۱: مقدمه

در این فصل ابتدا ویژگی‌های مهم تصویر که به صورت کلی برای بازشناسی شیء از تصاویر دیجیتال استفاده می‌شوند، بررسی می‌شوند و مدل‌های مختلف استخراج این ویژگی‌ها توضیح داده می‌شوند. سپس در مورد بخش دوم یک سیستم بازشناسی شیء یعنی یادگیری سیستم و تعیین استراتژی بازشناسی شیء توضیح داده می‌شود و مدل‌های مختلف بررسی می‌شوند. پس از آن تعدادی از الگوریتم‌های ارائه شده برای بازشناسی شیء را مرور می‌کنیم.

۲-۲: بررسی الگوریتم‌های استخراج ویژگی موجود در زمینه بازشناسی شیء

در این بخش به الگوریتم‌های موجود که در زمینه بازشناسی شیء به تعدد مورد استفاده قرار گرفته‌اند، می‌پردازیم. در ادامه ابتدا دسته بندی کلی ویژگی‌ها را بیان می‌کنیم و سپس روش‌های متداول بازشناسی شیء مبتنی بر این ویژگی‌ها را توضیح می‌دهیم.

۲-۲-۱: دسته بندی کلی ویژگی‌ها

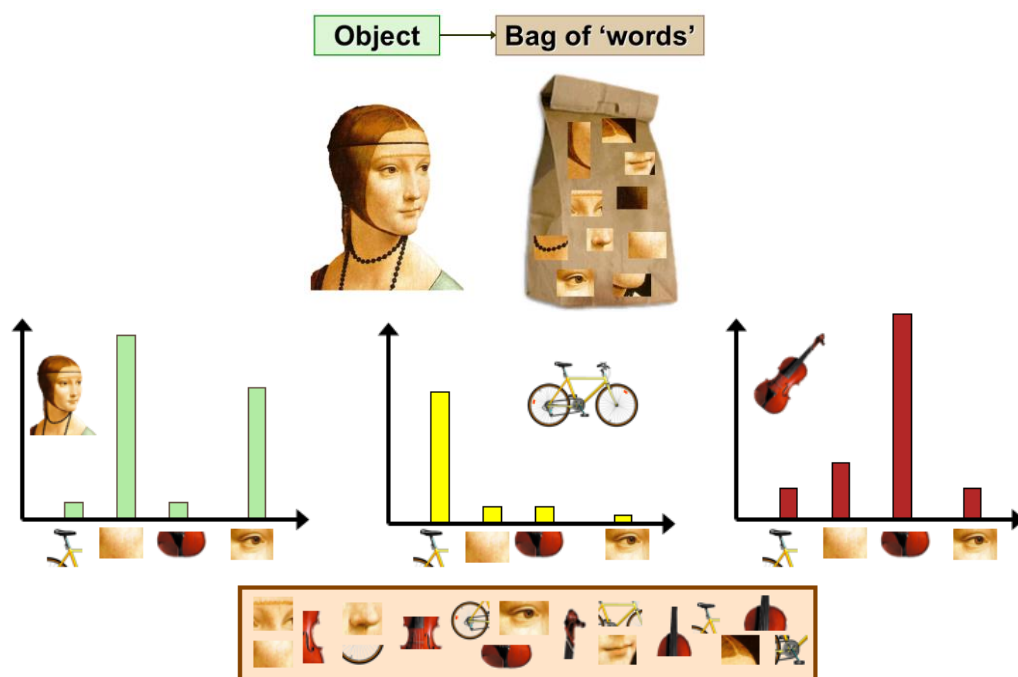
ساده ترین راه برای نمایش یک تصویر استفاده از تمام پیکسل‌های آن به عنوان ویژگی می‌باشد، اما این راه بسیار ابتدایی و دارای ابعاد زیاد می‌باشد. برای بدست آوردن اطلاعات مفید یا بردار ویژگی از تصاویر روش‌های مختلفی وجود دارد. از نظر ناحیه ای که برای بدست آوردن اطلاعات در نظر گرفته می‌شود، می‌توان ویژگی‌ها را به سه دسته تقسیم کرد:

- ویژگی‌های محلی: این ویژگی‌ها معمولاً یک فضای کوچک قابل تمیز از یک ناحیه را بازنمایی می‌کنند. به عنوان مثال انحنا، قطعات مرزی و گوشه‌ها از ویژگی‌های محلی به شمار می‌آیند [6].
- ویژگی‌های سراسری: این ویژگی‌ها تصویر را به عنوان یک ناحیه در نظر می‌گیرند برخی از این ویژگی‌ها عبارتند از توصیفگرهای فوریه و گشتاورها [6].
- ویژگی‌های رابطه‌ای: این ویژگی‌ها براساس موقعیت‌های نسبی موجودیت‌های مختلفی مانند نواحی، مرزهای بسته یا ویژگی‌های محلی هستند [6].

ویژگی مورد استفاده برای بازشناسی شیء تاثیر بسیار زیادی بر عملکرد سیستم دارد. در ادامه چند مورد از روش‌های استخراج ویژگی متداول را توضیح می‌دهیم.

۲-۲-۲: بسته کلمات

در این روش بخش های مستقلی از یک تصویر به عنوان کلمه در نظر گرفته می شوند و فرکانس تکرار آنها در تصویر محاسبه می شود. در نتیجه نحوه نمایش یک تصویر با استفاده از هیستوگرام ویژگی های تخصیص یافته به هر بخش است. در شکل 1-2 چگونگی نمایش تصاویر با استفاده از این ویژگی ها نشان داده شده است.



شکل 1-2 - نمایش سه شیء به روش کیف کلمات

۲-۲-۳: مدل مبتنی بر اجزاء

در روش کیف کلمات تمام قطعات تصویر با احتمال یکسانی در نظر گرفته می شوند و همچنین اطلاعاتی در مورد محل این قطعات در نظر گرفته نمی شود. در روش مبتنی بر اجزاء یک شیء به عنوان مجموعه ای از اجزاء در نظر گرفته می شود و یک مدل مبتنی بر ظاهر است که موقعیت مکانی بین اجزاء را نیز در نظر می گیرد.

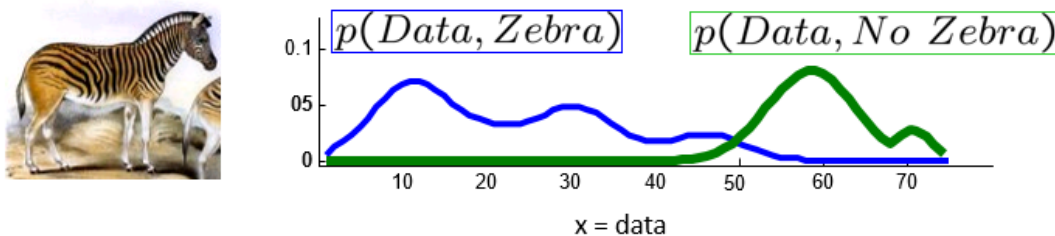
۲-۳: بررسی الگوریتم‌های یادگیری سیستم و بازشناسی شیء

مدلی که برای بدست آوردن برچسب هر ناحیه براساس ویژگی های بدست آمده استفاده می شود به سه دسته تقسیم می شود. مدل زایشی^{۲۸}، مدل جداکننده^{۲۹} و تابع کلاسیکندی که در ادامه هر یک توضیح داده می شوند.

۲-۳-۱: مدل زایشی

اگر ویژگی ها را به عنوان بردار x و برچسب ها را به عنوان بردار y در نظر بگیریم، در آن صورت در مدل زایشی هدف بدست آوردن احتمال توام x و y است. در شکل 2-2 برای مثال گورخر این احتمال نشان داده شده است. در این مدل احتمال توام داده های ورودی و گورخر بودن و گورخر نبودن شیء نشان داده شده است. در این مثال y به صورت باینری 0 و 1 در نظر گرفته شده است که در صورت 1 بودن شیء گورخر هست و در صورت 0 بودن شیء گورخر نیست.

این مدل به همراه قانون بیز می تواند برای کلاسیکندی استفاده شود و یا اینکه با استفاده از این مدل می توان احتمال رخداد جفت های x و y را بدست آورد.



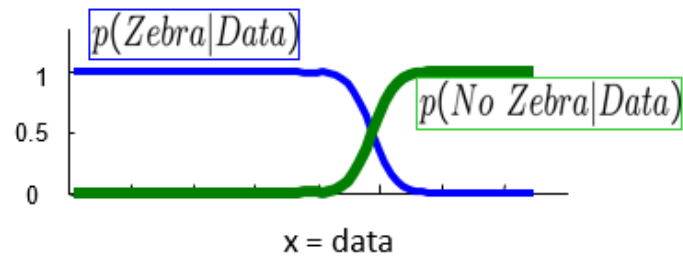
شکل 2-2 - یک نمونه مثال مدل زایشی. $P(Data, Zebra)$ نشان دهنده احتمال گورخر بودن و داده های ورودی به طور همزمان است و $P(Data, No Zebra)$ نشان دهنده احتمال گورخر نبودن و داده های ورودی است.

۲-۳-۲: مدل جداکننده

در این مدل احتمال شرطی y نسبت به x محاسبه می شود. یعنی با داشتن داده های x احتمال برچسب موردنظر محاسبه می شود. از این مدل می توان به عنوان کلاسیکندی استفاده کرد. شکل زیر مدل جداکننده را برای مثال گورخر نشان داده است.

²⁸ Generative

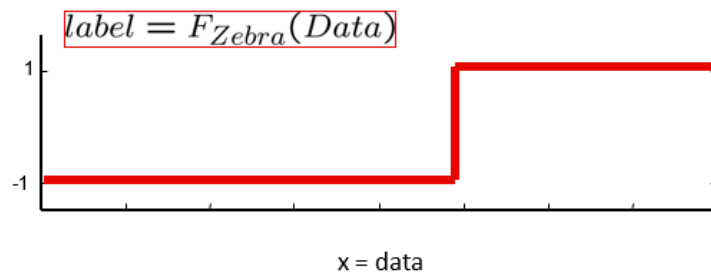
²⁹ Discriminative



شکل 3-2- یک مثال از مدل جداکننده. $p(\text{Zebra}|\text{Data})$ نشان دهنده احتمال گورخر بودن به شرط داده‌های ورودی و $p(\text{No Zebra}|\text{Data})$ نشان دهنده احتمال گورخر نبودن به شرط داده‌های ورودی است.

۳-۳-۲ تابع کلاسبندی

در این مدل برچسب، براساس تابعی از داده‌ها یا ویژگی‌ها محاسبه می‌شود. در شکل زیر از برچسب 1 و -1 برای نشان دادن اینکه آیا شیء مورد نظر گورخر هست یا خیر استفاده شده است.



شکل 4-2- یک مثال از تابع کلاسبندی. در این مثال توسط تابع F و ورودی Data تعیین می‌شود که آیا گورخر هست یا خیر.

روشهای مختلفی برای کلاسبندی وجود دارد که در این قسمت بعضی از متداول ترین روشها توضیح داده می‌شوند.

➤ روش نزدیک ترین همسایه

در این روش هر نمونه به نزدیک ترین کلاس نسبت داده می‌شود. هر کلاس توسط میانگین نمونه‌های آن کلاس تعریف می‌شود و با استفاده از فاصله اقلیدسی نمونه جدید با تمام میانگین‌های

کلاس‌های مختلف در مورد کلاس نمونه جدید تصمیم‌گیری می‌شود. از مزایای این روش سادگی آن و از معایب این روش این است که باید نمونه‌های زیادی داشته باشیم تا دقت بالایی بدست آید.

➤ روش شبکه عصبی

شبکه عصبی روشی الهام گرفته از سیستم‌های یادگیر طبیعی است که در آنها یک مجموعه پیچیده از نرون‌های به هم متصل در کار یادگیری دخیل هستند. این روش برای یادگیری توابع گوناگون نظیر توابع با مقادیر حقیقی، توابع با مقادیر گسسته و توابع با مقادیر برداری مناسب می‌باشد. اینگونه شبکه‌ها با موفقیت در مسائلی مانند شناسایی و تعبیر تصاویر، یادگیری روبات و شناسایی گفتار اعمال شده است. شبکه از تعداد دلخواهی سلول یا گره یا واحد یا نرون تشکیل می‌شود که مجموعه ورودی را به خروجی ربط می‌دهند.

این روش برای یادگیری توابع پیوسته و زمانی که زمان کافی برای یادگیری وجود داشته باشد مناسب است و نسبت به درخت تصمیم زمان بیشتری برای یادگیری نیاز دارد. وزن‌های یادگرفته شده توسط شبکه را به سختی می‌توان تعبیر کرد به همین دلیل زمانی استفاده می‌شود که نیازی به تعبیر تابع هدف نباشد. برای داده‌های آموزشی دارای نویز مانند داده‌های سنسورها نظیر دوربین و میکروفن نیز مناسب است.

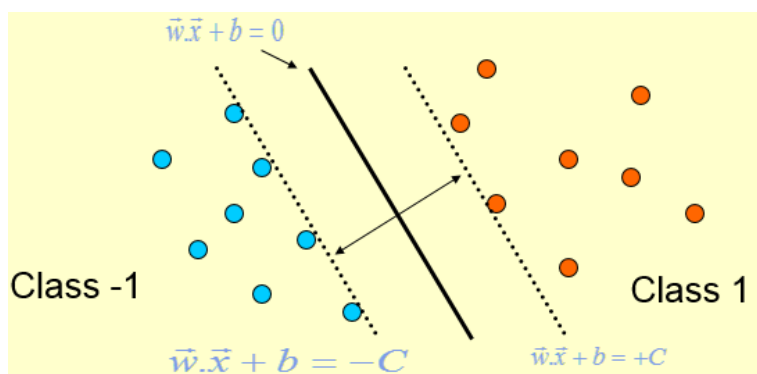
➤ روش ماشین بردار پشتیبان^{۳۰}

روش ماشین بردار پشتیبان روشی در یادگیری ماشین برای دسته‌بندی است و شهرت آن به دلیل موفقیت در تشخیص حروف دست‌نویس است. در این روش در صورتی که دسته‌ها بصورت خطی جداپذیر باشند، ابرصفحه‌هایی با حداکثر حاشیه بدست می‌آید که دسته‌ها را جدا می‌کند. در مسائلی که داده‌ها بصورت خطی جداپذیر نباشند، داده‌ها به فضای با ابعاد بیشتر نگاشت پیدا می‌کنند تا بتوان در فضای جدید آنها را بصورت خطی جدا کرد [7].

به دلیل اینکه جداکننده‌ها با حداکثر حاشیه تنظیم می‌شوند این جداکننده‌ها قابلیت تعمیم بیشتری پیدا می‌کنند. طبق قضیه‌ای در تئوری یادگیری ماشین اگر مثالهای آموزشی بدرستی دسته

³⁰ Support Vector Machine (SVM)

بندی شده باشند، از بین جداسازهای خطی، آن جداسازی که حاشیه داده های آموزشی را حداکثر میکند خطای تعمیم را حداقل خواهد کرد.



شکل 5-2- جداسازی با حداکثر حاشیه در ماشین بردار پشتیبان

مزایای این روش عبارتند از: آموزش در این روش به سادگی انجام می شود، قابلیت تعمیم خوبی دارد، با نمونه های آموزشی کم نیز به خوبی کار می کند، برخلاف شبکه عصبی بهینه محلی نیست و بهترین پاسخ سراسری را پیدا می کند. از ضعف های این روش نیز نیاز به انتخاب تابع هسته خوب است در غیر این صورت پاسخ درستی بدست نخواهد آمد.

➤ روش بوستینگ³¹

روش بوستینگ یک روش کلاسبندی براساس ترکیب تعدادی کلاسبند ضعیف برای ساختن یک کلاسبند قوی است. در هر الگوریتم کلاسبندی ضعیف سعی می شود یک ویژگی با کمترین خطا برای جدا کردن نمونه های مثبت و منفی به عنوان ویژگی جداکننده انتخاب شود. در نهایت نیز با ترکیب این کلاسبندها، کلاسبند نهایی بدست می آید.

➤ روش درخت تصمیم

ایده ی درخت تصمیم بسیار ساده و شهودی است و یک مثال نقطه شروع خوبی برای فرد ناآشنا است. این روش مانند بازی بیست سوالی است که براساس پاسخ هر سوال به هر نود حرکت انجام

³¹ Boosting

می‌شود. دو فرآیند باید انجام شود: اول ساخت درخت و دوم کلاسبندی براساس آن. ابتدا به تعریف واژگان در این زمینه پرداخته می‌شود:

- مثالها: آنچه کلاسبندی می‌شود و برجسب گذاری می‌شود تا کلاسی که به آن تعلق دارد شناسایی شود.
- دیتاست: مجموعه‌ی مثالهاست.
- ویژگی‌ها: معیارهایی که از مثالها استخراج می‌شود مانند ویژگی مدوربودن که معیاری برای تشخیص این است که چقدر یک مثال مانند دایره است.
- سوالها: تعریف کننده‌ی مرزها در فضای ویژگی
- گره‌ها: درخت را می‌سازند و شامل مثالهایی هستند و اگر گره‌ی پایانی نباشند شامل سوالی هستند که پاسخ آن به عنوان گره دیگری نشان داده می‌شود.

برای ساخت درخت باید سوالهایی برای تقسیم دیتاست به صورت معناداری پیدا کرد. برای کلاسبندی یک مثال جدید درخت براساس پاسخ سوالها در هر گره پیمایش می‌شود، زمانیکه به گره برگ برسیم کلاسبندی براساس مثالهای یادگرفته شده در آن گره انجام شده است.

یک کلاسبندی خوب به درخت بستگی دارد. برای بازشناسی سریع و صحیح فرآیند آموزش پیچیده و زمانبر خواهد بود. برای اینکه چه سوالی در هر گره پرسیده شود محدودیتی وجود ندارد. مساله اساسی در درخت تصمیم این است که در هر گره راههای ممکن فراوانی برای تقسیم داده وجود دارد.

۲-۴: روش‌های متداول بازشناسی شیء

یکی از روش‌های مورد استفاده در سیستم‌های بازشناسی شیء که بسیار مورد توجه قرار گرفت توسط ویالا^{۳۲} و جونز^{۳۳} ارائه شد که در این روش تصویر توسط پنجره‌هایی با ابعاد مختلف بررسی می‌شود و در صورتی شیء شناسایی می‌شود که پنجره با ابعاد شیء روی آن قرار بگیرد. در واقع در این روش برای کشف حضور شیء از تطبیق الگو استفاده می‌شود. این روش برای محیط‌های خاص که روشنایی و حالت شیء کنترل شده باشد کارآمد است، اما زمانی که شیء دارای چرخش، تغییر اندازه، تغییر حالت یا

³² Viola

³³ Jones

روشنایی باشد از لحاظ محاسباتی امکان پذیر نیست به خصوص زمانیکه بخشی از شیء نمایان باشد و زمانیکه پایگاه داده‌های بزرگ داشته باشیم [8, 9].

یک راه حل به جای جستجوی تمام موقعیت‌های تصویر و پیدا کردن تطابق، این است که ویژگی‌هایی از تصویر که نسبت به شکل تصویر تغییرناپذیرند استخراج شود و این ویژگی‌ها تطبیق داده شود. انواع خصوصیات کاندید شامل قطعات خط، گروه‌های لبه و ناحیه پیشنهاد و بررسی شدند. این ویژگی‌ها برای کلاس‌های خاصی از شیء به خوبی کار می‌کردند اما استحکام کافی برای ایجاد اساس بازشناسی قابل اعتماد را دارا نبودند.

ویژگی‌های سراسری و محلی هر کدام مزایا و معایب خود را دارا هستند. ویژگی‌های سراسری بازشناسی را سرعت می‌بخشند و ویژگی‌های محلی نسبت به تغییرات تصویر مانند شلوغی، تبدیلات و وجود مانع قدرت بیشتری دارند، اما تطبیق خصوصیات محلی بین اشیاء از زوایای دید مختلف زمان‌بر است. کارهای اخیر به توسعه‌ی مجموعه‌ی انبوه‌تری از ویژگی‌های تصویر پرداخته است. یکی از روش‌ها از بازیابی گوشه برای بازشناسی موقعیت‌های قابل تکرار تصویر استفاده می‌کند، به طور دقیق‌تر می‌توان گفت قله‌ی تغییرات محلی تصویر را بازشناسی می‌کند.

ژانگ^{۳۴} و همکارانش از گوشه‌یاب هریس^{۳۵} برای بازشناسی موقعیت‌های خصوصیات استفاده می‌کنند. در این روش علاوه بر تلاش برای یافتن ارتباط بین نواحی یک تصویر و تمام نواحی ممکن در تصویر دوم، تطبیق نواحی گوشه در هر تصویر با نواحی ممکن در همان تصویر نیز باعث صرفه‌جویی زیادی در زمان محاسبات انجام می‌شود.

گوشه‌یاب‌ها در روش‌های گفته شده یک مشکل بزرگ دارند، که تصویر را تنها در یک مقیاس بررسی می‌کنند و در صورت تغییر بزرگ در مقیاس این یابنده‌ها به نقاط مختلف تصویر واکنش نشان می‌دهند. همچنین به دلیل اینکه یابنده مقیاس شیء را نشان نمی‌دهد لازم است توصیفگر شیء ایجاد شود و تطبیق برای تعداد زیادی از مقیاس‌ها بررسی شود.

³⁴ Zhang

³⁵ Harris

بعضی روش‌ها نیز برای بازشناسی مبتنی بر ظاهر از تطبیق فضای پایه^{۳۶} و هیستوگرام رنگ استفاده می‌کنند. این روش‌ها همه برای اشیاء جدا شده یا تصاویر بخش بندی شده با موفقیت عمل کرده‌اند اما به دلیل خصوصیات سراسری که در آنها استفاده می‌شود توسعه‌ی آنها به تصاویر شلوغ و تصاویری که بخشی از شیء پوشانده شده سخت است. اهب^{۳۷} و آیکوچی^{۳۸} روش فضای پایه را با استفاده از تعداد زیادی پنجره‌های محلی کوچک با موفقیت برای تصاویر شلوغ استفاده کردند، اما این کار مانند روش تطبیق نمونه که باید برای هر پنجره کل تصویر پیمایش شود، نیاز به جستجوی سنگینی دارد.

در مقاله‌ی فلزنزواب و همکارانش^{۳۹} سیستمی برای بازشناسی شیء براساس ترکیبی از مدل‌های جزئی چندمقیاسی منعطف^{۴۰} شرح داده شده است. این سیستم قادر است کلاس‌های شیء با تنوع بالایی را نمایش دهد و بهترین نتایج را بر دیتابیس پاسکال بدست آورد. در این مقاله برای مدل‌سازی اشیاء از گرامرهای بصری استفاده شده است. مدل‌های مبنی بر گرامر با نمایش اشیاء با استفاده از ساختارهای سلسله‌مراتبی متغیر، مدل‌های جزئی منعطف را تعمیم می‌دهند. هر بخش در یک مدل مبنی بر گرامر می‌تواند به طور مستقیم یا توسط اجزاء تعریف شود. مدل‌های مبتنی بر گرامر تغییرات ساختاری را مدل می‌کنند [10].

به طور کلی بهبود عملکرد توسط غنی‌سازی مدل بسیار سخت است و مدل‌های ساده عموماً بهتر از مدل‌های پیچیده عمل می‌کنند. دلیل عملکرد بهتر مدل‌های ساده این است که مدل‌های پیچیده در فاز آموزش با سختی‌های زیادی روبرو هستند. در بازیابی شیء، مدل‌های کیف و ویژگی^{۴۱} و قالب‌های سخت به سادگی با روش‌های جداکننده مانده ماشین بردار پشتیبان آموزش داده می‌شوند. مدل‌های غنی‌تر برای آموزش سخت‌تر هستند زیرا معمولاً از اطلاعات پنهان استفاده می‌کنند.

به عنوان مثال زمانیکه یک مدل مبنی بر جزء از تصاویری که دارای یک جعبه محصور اطراف اشیاء هستند آموزش داده می‌شود، موقعیت اجزاء برچسب ندارند و باید با آنها به عنوان متغیرهای پنهان در حین آموزش رفتار شود. اگر برچسب گذاری کامل‌تر باشد آموزش ضعیف‌تر خواهد بود. برچسب گذاری

³⁶ Eigen space

³⁷ Ohba

³⁸ Ikeuchi

³⁹ Felzenszwalb, P. F., R. B. Girshick, et al.

⁴⁰ Multiscale deformable part models

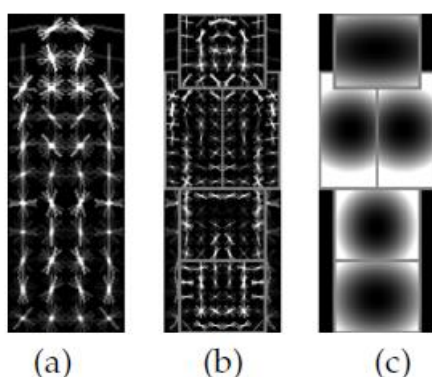
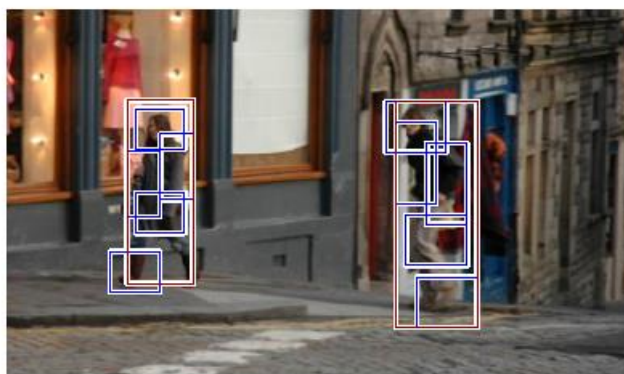
⁴¹ Bag of Features

جزئی اتوماتیک عملکرد بهتری دارد که به طور اتوماتیک اجزاء را پیدا می کند. برچسب گذاری با جزئیات بیشتر نیز بسیار پرهزینه و زمانبر است.

یابنده دالال- تریگز⁴² که جایزه سال 2006 رقابت پاسکال را برنده شده است، از یک فیلتر بر ویژگی های هیستوگرام گرادیان ها استفاده می کند تا یک طبقه شیء را نمایش دهد. این یابنده از روش پنجره لغزان استفاده می کند که یک فیلتر به تمام موقعیت ها و به تمام اندازه ها در تصویر اعمال می شود. این یابنده یک کلاسبند است که به عنوان ورودی، تصویر و موقعیت و اندازه آن را می گیرد، و تصمیم می گیرد که آیا نمونه ای از شیء در موقعیت مورد نظر وجود دارد یا نه [11].

در این مقاله با استفاده از یک مدل با ساختار ستاره و مبتنی بر جزء این ایده گسترش داده شده است. یک فیلتر ریشه و مجموعه ای از فیلترهای اجزاء تعریف می شوند. امتیاز فیلتر ریشه و اجزاء با ضرب نقطه ای فیلتر و زیرپنجره هرم ویژگی محاسبه شده از تصویر ورودی تعریف می شود. شکل زیر یک مدل ستاره برای طبقه شیء شخص را نشان می دهد.

⁴² N. Dalal and B. Triggs



شکل 6-2- بازیابی های بدست آمده با یک مدل شخص با یک عنصر اصلی. مدل با یک فیلتر ریشه (a) چندین فیلتر اجزاء با رزولوشن بالاتر (b) و یک مدل مکانی برای موقعیت هر جزء نسبت به ریشه (c) تعریف شده است. فیلترها وزن های هیستوگرام گرادیان های ویژگی ها را مشخص می کنند [11].

برای آموزش مدل از اطلاعات برچسب گذاری جزئی استفاده شده که ماشین بردار پشتیبان پنهان^{۴۳} نامیده شده است. بدلیل آموزش جداکننده مهم است که از مجموعه آموزشی بزرگ استفاده شود. در بازشناسی شیء مساله آموزش نامتوازن است زیرا زمینه بسیار بیشتر از اشیاء است.

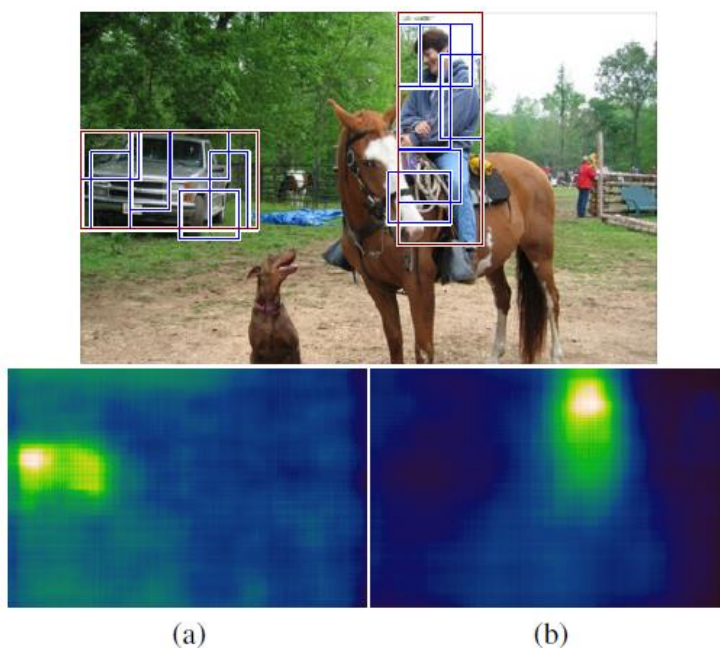
مدل های شیء توسط فیلترهایی که زیرپنجره های هرم ویژگی را امتیاز می دهند تعریف می شوند. مجموعه ویژگی شبیه ویژگی های هیستوگرام گرادیان استفاده شده است و با استفاده از تحلیل عناصر اصلی ابعاد آنها به طور قابل توجهی کاهش یافته است. همچنین بعضی مشکلات مرتبط با رقابت بازیابی شیء مجموعه داده پاسکال و مجموعه داده های مشابه برطرف شده است. یک روش ساده نیز برای جمع کردن خروجی یابنده های اشیاء مختلف اثبات شده است. ایده ی اصلی براین مبناست که وجود بعضی اشیاء دلیل برای وجود یا عدم وجود اشیاء دیگر در تصویر است. از این ایده به این صورت استفاده

⁴³ Latent SVM

می‌شود که در آموزش یک کلاسبند طبقه خاص هر بازیابی براساس امتیاز اصلی خود و بیشترین امتیاز از بقیه اشیاء دوباره امتیازدهی می‌شود.

در مقاله فلزنزواب و همکارانش⁴⁴ یک روش عمومی برای ساخت کلاسبندهای آبخاری از مدل‌های منعطف مبتنی بر جزء مانند ساختارهای تصویری ارائه شده است. تمرکز اصلی در این مقاله بر مدل‌های با ساختار ستاره است و نشان داده می‌شود که چطور یک الگوریتم ساده براساس هرس فرضیه‌های جزئی بازیابی شیء را سرعت می‌بخشند. در این الگوریتم فرضیه‌های جزئی با ترتیبی از آستانه‌ها هرس می‌شوند. در نهایت یک الگوریتم آبخاری بازیابی برای یک کلاس عمومی مدل‌ها با یک ظاهر گرامری تعریف می‌شود [12].

این الگوریتم یک روش بازیابی 20 بار سریعتر از بازیابی‌های استاندارد ارائه می‌دهد که براساس برنامه نویسی پویا و تبدیلات فاصله تعمیم داده شده و بدون کاهش عملکرد است. شکل زیر کاری را که توسط این الگوریتم در موقعیت‌های مختلف یک تصویر با استفاده از دو مدل مختلف انجام شده است، نشان داده است.



شکل 7-2- قسمت (a) از مدل ماشین و قسمت (b) از مدل شخص استفاده شده است [12].

⁴⁴ P. Felzenszwalb, R. Girshick, D. McAllester

۲-۵: خلاصه و نتیجه گیری

در این فصل پس از دسته بندی کلی ویژگی‌هایی که از تصویر استخراج می‌شوند، الگوریتم‌های متداول استخراج ویژگی در بازشناسی شیء معرفی شدند. سپس در مورد استراتژی‌های بازشناسی شیء و انواع روش یادگیری سیستم توضیح داده شد. پس از آن تاریخچه ای از روش‌های مختلف بازشناسی شیء مطرح و معایب و مزایای هر کدام بررسی شدند.

فصل ۳:

پیش زمینه تئوری

۳-۱: مقدمه

در دنیای واقعی، اشیاء با اشیاء دیگر و محیط‌های مخصوص تغییر می‌کنند و منبعی غنی از پیوستگی‌های متنی فراهم می‌کنند که می‌توانند توسط سیستم‌های بصری استفاده شوند. مشاهدات زیادی در ادراک بصری^{۴۵}، بینایی کامپیوتر و عصب‌شناسی^{۴۶} نشان داده است که اطلاعات متن بر کارآمدی جستجو و شناسایی اشیاء تاثیرگذار است. اشیائی که در یک زمینه سازگار یا آشنا باشند با دقت بیشتری شناسایی می‌شوند و سریع‌تر از اشیائی که در یک صحنه‌ی ناسازگار هستند، پردازش می‌شوند [13]. مدل‌های زیادی برای طبقه‌بندی شیء از ارتباطات معنایی بیدرمن^{۴۷} بهره برده‌اند. اطلاعات متنی عبارتست از هر اطلاعاتی که به طور مستقیم توسط ظاهر شیء تولید نشود. این اطلاعات را می‌توان از تگ‌ها، حاشیه‌نویسی‌ها، وجود و موقعیت اشیاء دیگر بدست آورد [3].

در این فصل یک مدل مبتنی بر متن که بهترین عملکرد را در بازشناسی شیء و بر تصاویر با تعداد اشیاء زیاد داشته است را بررسی می‌کنیم و جزئیات روش را شرح می‌دهیم. همچنین مروری بر روش‌هایی که از درخت تصمیم استفاده کرده‌اند خواهیم داشت. در نهایت نیز الگوریتم کارت که برای ساخت درخت تصمیم استفاده شده است را معرفی خواهیم کرد.

۳-۲: مدل مبتنی بر اطلاعات متنی سلسله مراتبی

با استفاده از تشخیص‌دهنده‌هایی که برای کشف یک شیء استفاده می‌شوند تنها می‌توان یک دسته از انواع شیء را استخراج کرد. برای کشف تعداد شیء بیشتر نیاز به اجرای یابنده‌های مجزا برای هر دسته شیء است. به دلیل اینکه هرکدام از این یابنده‌ها به طور مجزا کار می‌کنند، نتیجه ممکن است غلط باشد. برای بهبود عملکرد شناسایی شیء، می‌توان علاوه بر خصوصیات محلی اطلاعات متنی مانند خصوصیات سراسری یک تصویر (مانند اینکه صحنه خیابان است) و روابط بین انواع اشیاء (مانند اینکه خیابان و ماشین معمولاً باهم در یک صحنه قرار دارند) را نیز استفاده کرد. همانطور که در شکل 21 نشان داده شده است، با داشتن دسته‌های زیاد شیء تشخیص‌های اشتباه زیادی اتفاق می‌افتد. 6 تشخیص با اطمینان بیشتر در شکل 22 نشان داده شده است که ترکیبی از اشیاء درونی و بیرونی است، درحالی

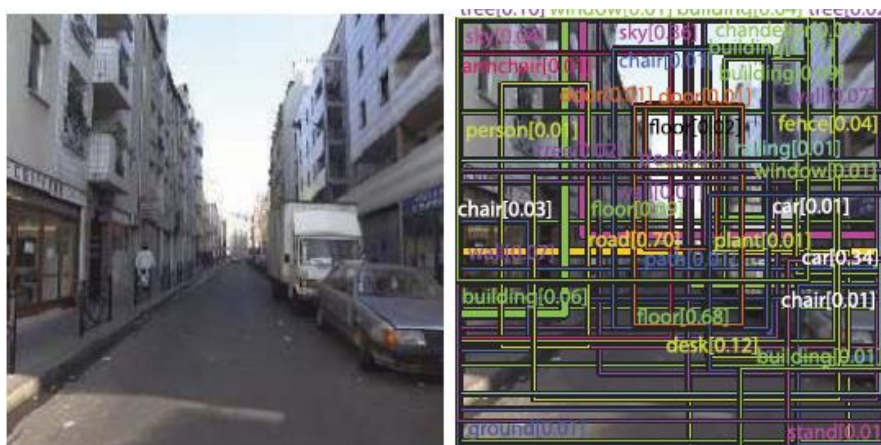
⁴⁵ Visual cognition

⁴⁶ Cognitive neuroscience

⁴⁷ Biederman

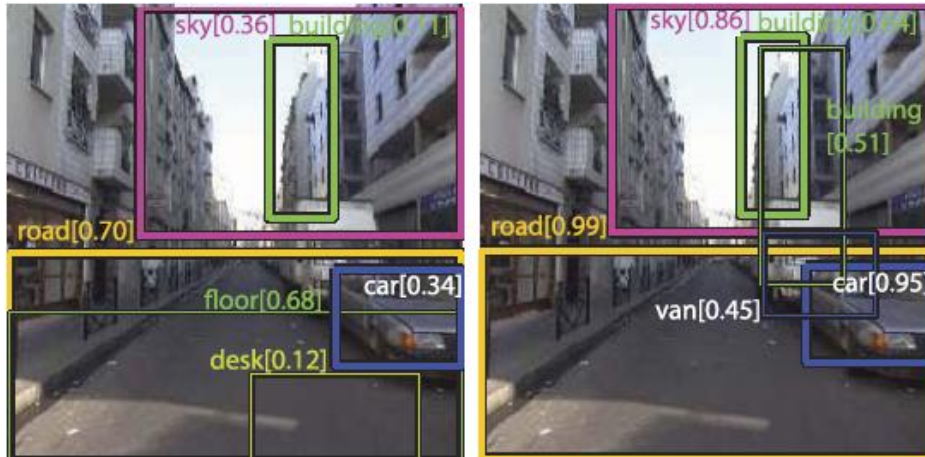
که اگر مدل متنی را نیز در نظر بگیریم احتمال کمتری برای اشیاء درونی مانند کف اطاق و میز وجود خواهد داشت. (شکل 22)

اشیاء داخلی و اشیاء خارجی معمولاً در یک صحنه اتفاق نمی‌افتند، در مقاله چی⁴⁸ و تورالبا از یک مدل درخت که تمام اشیاء خارجی در یک زیردرخت و تمام اشیاء داخلی در یک زیردرخت دیگر هستند استفاده کرده است، که دو درخت با وزن منفی به یکدیگر متصل هستند. به طور مشابه اشیاء مرتبط با آشپزخانه مانند یخچال، ظرفشویی و غیره همه در یک زیردرخت هستند و وزن‌های مثبت دارند. با استفاده از وابستگی‌های ساختار درختی بین اشیاء می‌توان این مدل را برای بیش از صد نوع شیء استفاده کرد. با ترکیب این مدل از روابط اشیاء با خروجی‌های شناساگرهای محلی و خصوصیات سراسری تصویر می‌توان تمام نمونه‌های مختلف انواع اشیاء را شناسایی و تعیین موقعیت کرد [4]. در ادامه به شرح این روش پرداخته می‌شود.



شکل 1-3- تصویر ورودی و نتیجه اعمال 107 یابنده کلاس

⁴⁸ Choi



شکل 2-3 - سمت چپ: خروجی شش یابنده با اطمینان بالاتر، سمت راست: خروجی شش یابنده با استفاده از مدل متنی استفاده شده.

۳-۲-۱: مدل سازی اولیه

هر شیء را می توان با یک متغیر باینری که نشان دهنده ی حضور یا عدم حضور یک شیء در تصویر است نشان داد و برای نمایش محل شیء از متغیر گوسی استفاده می شود.

در یک درخت باینری اگر هر گره b_i نمایشگر این باشد که شیء i در تصویر وجود دارد یا خیر، احتمال توأم متغیرهای باینری طبق ساختار درخت به صورت زیر محاسبه می شوند:

$$p(b) = p(\text{broot}) \prod p(b_i | \text{bpa}(i)) \quad (3-1)$$

که $\text{Pa}(i)$ والد گرهی i است. b نشان دهنده ی تمام متغیرهای مربوطه است $b \equiv \{b_i\}$. والد و فرزند ممکن است رابطه ی مثبت (مثل رابطه ی کف اطاق با دیوار) و یا منفی (مثل رابطه ی کف اطاق و آسمان) داشته باشند.

اشیاء در موقعیت های مخصوصی نسبت به یکدیگر قرار می گیرند. برای مثال، صفحه ی کامپیوتر معمولاً بالای یک کیبورد و موس قرار می گیرد. برای ذخیره ی این اطلاعات مکانی می توان متغیرهای موقعیت را به مدل درخت اضافه کرد.

در این مدل به جای استفاده از شیء جدا شده از تصویر از یک قاب محاط استفاده شده است که کوچکترین قاب حصارکننده برای تمام نقاط بخش جدا شده است. از این قاب برای نمایش محل یک نمونه شیء استفاده می شود. فرض می کنیم I_x و I_y به ترتیب مختصات افقی و عمودی مرکز قاب محاط

و l_w و l_h به ترتیب پهنا و ارتفاع قاب هستند. همچنین فرض می‌شود ارتفاع تصویر به یک نرمال شده است و $l_x=0$ ، $l_y=0$ مرکز تصویر است. فاصله‌ی موردانتظار بین مراکز اشیاء به سائز آنها بستگی دارد، اگر یک کیبورد و یک موس کوچک باشند فاصله‌ی بین مراکز آنها نیز باید کوچک باشد. مدل صورت فلکی اطلاعات موقعیت را به فضای تغییرناپذیر مقیاسی تبدیل می‌کند. هویم⁴⁹ و همکارانش تغییرات مقیاس را به اطلاعات 3بعدی واضح مرتبط می‌سازند [14]. طبق این روش برای تعیین موقعیت‌های شیء تبدیلات مختصات زیر در مختصات سه‌بعدی انجام می‌شود.

$$L_x = \frac{l_x}{l_h} H_i, \quad L_y = \frac{l_y}{l_h} H_i, \quad L_z = \frac{1}{l_h} H_i, \quad (3-2)$$

L_z فاصله‌ی بین ناظر و شیء است و H_i ارتفاع فیزیکی شیء i است. ارتفاع هر نوع شیء به صورت دستی براساس سائز اشیاء واقعی کد می‌شود. (مثلا شخص=1.7 متر، ماشین=1.5 متر)

موقعیت‌های ارتباطی افقی اشیاء به دلیل دیدهای مختلف از یک تصویر با تصویر دیگر به طور قابل توجهی متفاوت است و نشان داده شده است که موقعیت‌های افقی عموماً اطلاعات متنی ضعیفی دارند. بنابراین L_x در نظر گرفته نشده است و تنها L_y و L_z را برای ذخیره موقعیت عمودی و وابستگی‌های مقیاسی در نظر گرفته شده است. فرض شده است که $L_y S$ و $L_z S$ مستقل هستند، مثلاً موقعیت عمودی یک شیء مستقل از فاصله‌ی آن از سطح تصویر است. $L_y S$ را مانند گوسین توام و $L_z S$ با استفاده از توزیع‌های \log نرمال مدل شده است، زیرا همیشه دارای مقادیر مثبت هستند و اطراف مقادیر کوچک توزیع شده‌اند. یک متغیر موقعیت برای شیء i را به صورت زیر بازتعریف می‌شود.

$$L_i = (L_y, \log L_z) \quad (3-3)$$

برای سادگی روابط مکانی بین انواع شیء مدل شده است، نه بین تک تک نمونه‌ها بنابراین اگر چندین نمونه از شیء i در یک تصویر وجود داشته باشد، L_i موقعیت میانه‌ی تمام نمونه‌ها را نشان می‌دهد.

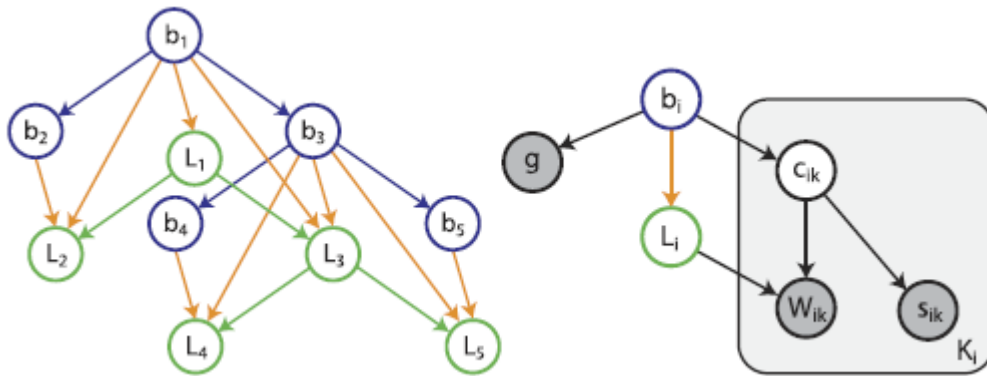
با فرض شرط بر متغیر وجود b ، ساختار وابستگی $L_i S$ ساختار درختی مشابهی با درخت باینری دارد:

⁴⁹ Hoiem

$$p(L|b) = p(L_{root}|b_{root}) \prod p(L_i|L_{pa(i)}, b_i, b_{pa(i)}) \quad (3-4)$$

که پتانسیل هر لبه توزیع موقعیت فرزند به شرط موقعیت والد و وجود یا نبودن فرزند و والد را کد می‌کند. شکل 4 مدل گرافیکی مربوط به متغیرهای وجودی b_i و متغیرهای موقعیت L_i را نشان می‌دهد. با ترکیب 1 و 3 توزیع توام تمام متغیرهای باینری و گوسی به صورت زیر نشان داده می‌شود:

$$\begin{aligned} p(b, L) &= p(b)p(L|b) \\ &= p(b_{root})p(L_{root}|b_{root}) \\ &\times \prod p(b_i|b_{pa(i)})p(L_i|L_{pa(i)}, b_i, b_{pa(i)}). \end{aligned} \quad (3-5)$$



شکل 3-3 - چپ: مدل اولیه که ارتباط متغیرهای وجودی و متغیرهای موقعیت شیء را نشان می‌دهد. راست: مدل اندازه‌گیری برای شیء i .

۲-۲-۳: مدل اندازه‌گیری

یکی کردن خصوصیات سراسری تصویر

توصیف‌گر چیست نمایش با ابعاد کم یک تصویر است، که بافت درشت و لایه‌ی مکانی یک صحنه را تسخیر می‌کند. برای استفاده از خصوصیات سراسری تصویر چیست به عنوان معیاری برای حضور هر متغیر b_i معرفی شده است. این کار به مدل امکان بدست آوردن نوع صحنه که برای پیش بینی اشیاء داخلی و خارجی مفید است را می‌دهد.

یکی کردن خروجی‌های یابنده‌ی محلی

برای بازیابی و تعیین موقعیت نمونه‌های شیء در یک تصویر ابتدا یابنده‌های بازیابی تک شیئی اعمال شده است و مجموعه‌ای از پنجره‌های کاندید برای هر نوع شیء بدست آمده است. i نشان دهنده‌ی نوع شیء و k برای فهرست کردن پنجره‌های کاندید که توسط یابنده‌های پایه تولید می‌شوند استفاده می‌شود. هر خروجی یابنده یک امتیاز S_{ik} و یک جعبه محصورکننده فراهم می‌کند که تبدیل 2 بر آن اعمال می‌شود تا موقعیت متغیر بدست آید $W_{ik}=(L_y, \log L_z)$ به هر پنجره یک متغیر باینری c_{ik} نسبت داده شده است تا نشان دهد آیا بازیابی درستی است ($c_{ik}=1$) یا خیر ($c_{ik}=0$). شکل 4 مدل اندازه گیری برای شیء i برای یکی کردن جیست و یابنده پایه در مدل پیشین را نشان می‌دهد که در آن از نشانه‌گذاری برای نمایش K_i پنجره مختلف کاندید استفاده شده است.

اگر یک پنجره کاندید به درستی شیء i را بازیابی کرده باشد ($c_{ik}=1$) سپس موقعیت آن W_{ik} یک بردار گوسی با میانگین L_i موقعیت شیء i را نشان می‌دهد، و درغیراین صورت ($c_{ik}=0$) W_{ik} مستقل از L_i است و یک توزیع یکنواخت دارد.

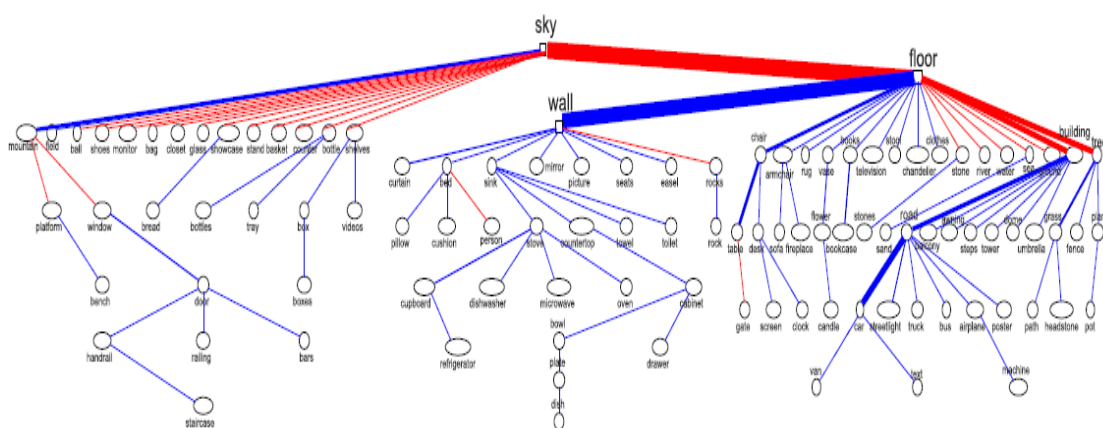
۳-۲-۳: یادگیری ساختار وابستگی شیء

برای یادگیری ساختار وابستگی بین اشیاء از مجموعه‌ای از تصاویر برجسب گذاری شده استفاده شده است. یادگیری یک مدل درخت توسط الگوریتم ساده و کارآمد چو-لو^{۵۰} امکان پذیر است [15]. این الگوریتم ابتدا اطلاعات متقابل تجربی مقادیر نمونه تمام جفت متغیرها را محاسبه می‌کند. سپس ماکسیموم وزن درخت پوشا که وزن یالها برابر با اطلاعات متقابل بین متغیرهایی که توسط یال به یکدیگر متصل شده‌اند، است را پیدا می‌کند. ساختار درخت با استفاده از نمونه‌های b_{is} در یک مجموعه از تصاویر برجسب گذاری شده یادگرفته می‌شود. حتی باوجود بیش از 100 شیء و هزاران تصویر آموزشی، یک مدل درخت در چند ثانیه در متلب یادگرفته می‌شود.

شکل زیر ساختار درخت یادگرفته شده از دیتاست سان 09 را نشان می‌دهد. به دلیل اینکه خروجی الگوریتم چو-لو یک درخت غیرجهت دار است، sky را به عنوان ریشه درخت انتخاب کرده است تا یک ساختار درختی جهت دار بدست آید. نکته موردتوجه این است که از اطلاعات مربوط به ساختار سلسله مراتبی اصلی بین انواع اشیاء در حین فرآیند یادگیری استفاده نمی‌شود. الگوریتم چو-لو به سادگی وابستگی‌های دودبوی قوی را انتخاب می‌کند. ساختار درخت یادگرفته شده اشیاء را در یک سلسله مراتب

⁵⁰Chow-liu

طبیعی سازماندهی می‌کند. برای مثال یک زیردرخت با ریشه‌ی ساختمان اشیاء زیادی دارد که در صحنه خیابان قرار دارند، و یک زیردرخت با ریشه‌ی sink شامل اشیایی است که به طور متداول در آشپزخانه وجود دارند. بنابراین بسیاری از گره‌های غیربرگ به همانند ابرشی⁵¹ با اندازه‌ی درشت‌تر عمل می‌کنند. به عبارت دیگر ساختار درخت یادگرفته شده سلسله‌مراتب اصلی بین اشیاء و صحنه‌ها را ذخیره می‌کند. در نتیجه بازشناسی شیء و عملکرد فهم صحنه بهتر اتفاق می‌افتد.



شکل 4-3- ساختار وابستگی اشیاء که از دیتاست SUN 09 یادگرفته شده است. یال‌های قرمز نشان‌دهنده‌ی رابطه‌ی منفی بین طبقه‌ها است. کلفتی هر یال قدرت آن رابطه را نشان می‌دهد [16]

۳-۲-۴: یادگیری پارامترهای مدل

برای یادگیری پارامترهای مدل پیشین از برجسب‌های صحیح تصاویر آموزشی استفاده می‌شود. برای محاسبه $p(b_i | b_{pa(i)})$ هم‌رخداد جفت شیء‌های والد و فرزند به سادگی شمارش می‌شود. برای هر جفت والد و فرزند از سه توزیع گوسی مختلف برای $p(L_i | L_{pa(i)}, b_i, b_{pa(i)})$ استفاده می‌شود. زمانیکه هر دو شیء حضور داشته باشند ($b_i=1, b_{pa(i)}=1$) در این صورت موقعیت شیء فرزند L_i به موقعیت والد $L_{pa(i)}$ بستگی دارد. زمانیکه شیء وجود داشته باشد ولی والد آن وجود نداشته باشد موقعیت فرزند مستقل از والد است. زمانیکه شیء وجود نداشته باشد، فرض می‌کنیم L_i مستقل از تمام موقعیت‌های اشیاء دیگر است و میانگین آن میانگین موقعیت شیء i در تمام تصاویر است.

در مدل اندازه‌گیری، $p(g | b(i))$ را می‌توان توسط توصیفگرهای جیست که از هر تصویر آموزشی محاسبه می‌شوند، آموزش داد. برای یادگیری بقیه‌ی پارامترها در مدل اندازه‌گیری از یابنده‌های

⁵¹ Meta-object

محلای برای هر نوع شیء در تصاویر آموزشی استفاده می شود. امتیازات یابنده ی محلی مرتب می شوند و s_{ik} ، k امین امتیاز بالا برای طبقه i است و $p(c_{ik}|s_{ik})$ با استفاده از رگرسیون منطقی⁵² یادگرفته می شود که از طریق آن می توان $p(s_{ik}|c_{ik})=p(c_{ik}|s_{ik})p(s_{ik})/p(c_{ik})$ را محاسبه کرد.

۳-۲-۵: استفاده از مدل

با داشتن چیست (g) و موقعیت های کاندید پنجره ها (W) و امتیازات آنها (S)، حضور اشیاء (b)، بازیابی صحیح (c) و موقعیت های موردانتظار تمام اشیاء (L) با حل مساله بهینه سازی زیر بدست می آیند:

$$b, c, L = \arg \max_{b,c,L} p(b, c, L|g, W, s). \quad (3 - 6)$$

در مرحله ی اول اطلاعات موقعیت W در نظر گرفته نمی شود و تخمین MAP متغیرهای b و c براساس چیست تصویر و امتیازات پنجره های کاندید انجام می شود:

$$b, c = \arg \max_{b,c} p(b, c|s, g). \quad (3 - 7)$$

سپس براساس این تخمین ها موقعیت های موردانتظار اشیاء به صورت زیر با استفاده از درخت گوسی به دست می آیند:

$$L = \arg \max_L p(L|b, c, W). \quad (3 - 8)$$

سپس براساس تخمین های موقعیت حضور و متغیرهای یابنده دوباره تخمین زده می شوند:

$$b, c = \arg \max_{b,c} p(b, c|s, g, L, W) = \arg \max_{b,c} p(b, c|s, g)p(L, W|b, c), \quad (3 - 9)$$

در این گام جفت اشیاء یا پنجره ها در موقعیت های محتمل در تصویر برای حضور بررسی می شوند.

در گام نهایی احتمالی حاشیه ای هر متغیر وجود و احتمال حاشیه ای هرمتغیر یابنده برای تعیین موقعیت شیء به صورت زیر محاسبه می شود:

⁵² Logistic regression

$$p(c_{ik} = 1 | s, g, L, W) \quad (3 - 10)$$

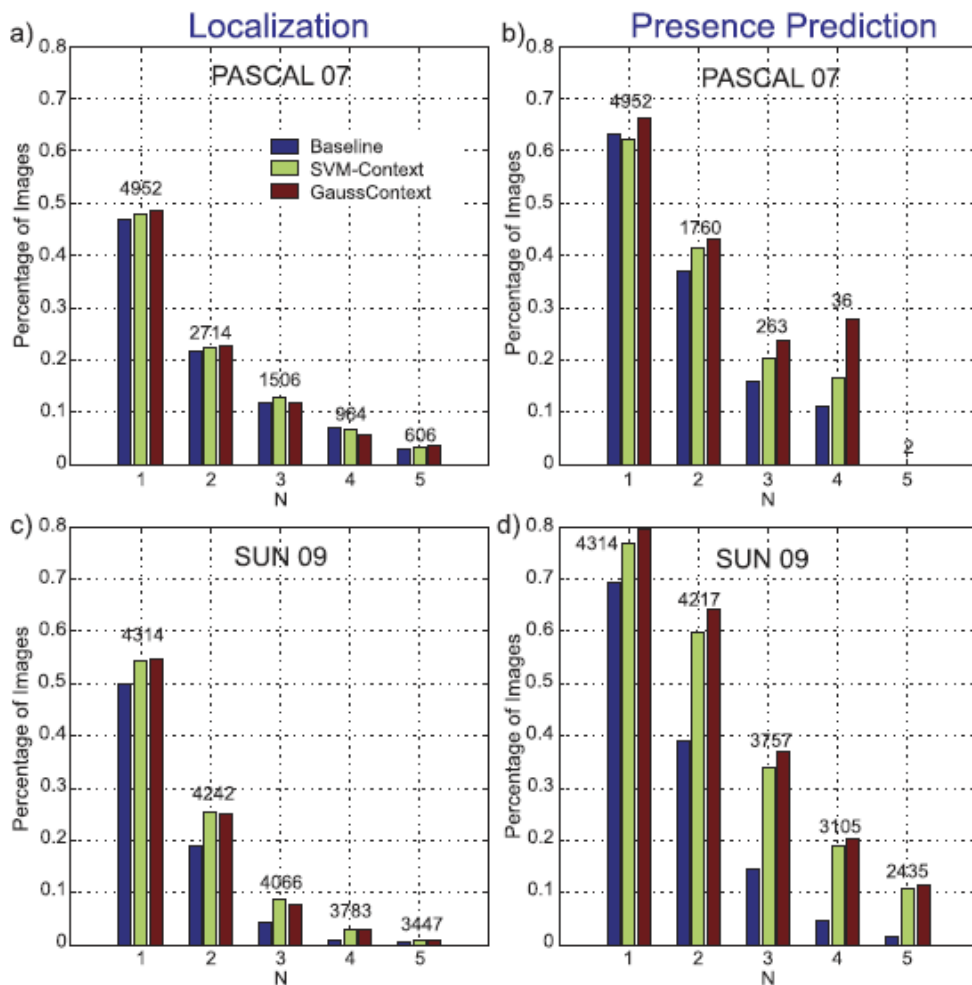
$$p(b_i = 1 | s, g, L, W) \quad (3 - 11)$$

۳-۲-۶: نتایج

برای آموزش این مدل از 4367 تصویر و برای تست از 4317 تصویر از مجموعه داده سان 09 استفاده شده است. برای استفاده از نمونه‌های آموزشی کافی برای یابنده‌های پایه 26000 تصویر اضافی دیگر نیز تفسیر شده و برای آموزش یابنده‌های پایه استفاده شده اند. این تصاویر شامل یک شی هستند و برای یادگیری مدل متنی استفاده نشده اند.

این مجموع تصاویر از مجموعه داده سان دارای بیش از 200 شی است اما یابنده‌های پایه برای بعضی از اشیاء حتی با تفسیرهای بیشتر نیز کیفیت پایینی ارائه می‌دهند. به دلیل اینکه مدل متنی از خروجی‌های یابنده‌های پایه استفاده می‌کند و احتمال بازیابی صحیح را برای هر پنجره بدست می‌آورد، اگر کاندیدی برای بعضی از نمونه اشیاء توسط یابنده پایه وجود نداشته باشد نمی‌تواند بازیابی انجام دهد. بنابراین کلاس‌های شی که یابنده پایه نتوانسته حداقل 4 کاندید صحیح در کل مجموعه داده برای آن تولید کند حذف شده اند و از 107 کلاس شی دیگر استفاده شده است.

در شکل زیر عملکرد تعیین موقعیت و پیش‌بینی حضور شی این مدل بر مجموعه داده سان و پاسکال نشان داده شده است و با دو الگوریتم یابنده پایه [10] و ماشین بردار پشتیبان متنی [10] مقایسه شده است. برای این کار N تا از بازیابی‌های با اطمینان بیشتر در هر تصویر بررسی می‌شوند که آیا درست هستند یا خیر. اعداد روی هر ستون نشان دهنده‌ی تعداد تصاویر که شامل حداقل N نمونه شی برای تعیین موقعیت و N طبقه شی برای پیش‌بینی حضور هستند، می‌باشد. همانطور که دیده می‌شود با بزرگتر شدن N تعداد تصاویر در پاسکال بسیار کم می‌شود، زیرا بیشتر تصاویر در این مجموعه داده تنها شامل یک یا دو طبقه شی هستند.

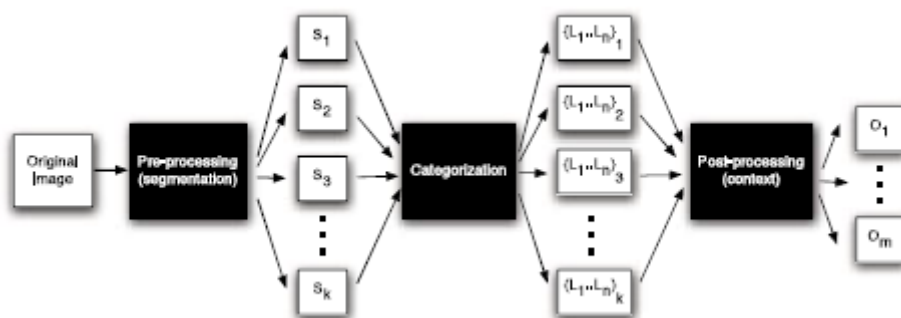


شکل 5-3- عملکرد تعیین موقعیت و پیش‌بینی حضور بر مجموعه داده سان 9 و پاسکال 7 و a و c درصد تصاویری که N تا از بیشترین بازیابی‌های با اطمینان بیشتر صحیح بوده‌اند. اعداد روی هر ستون تعداد تصاویری که حداقل شامل N نمونه شیء بوده‌اند، را نشان می‌دهد. B و d درصد تصاویری که N تا از پیش‌بینی حضورها با بیشترین احتمال درست بوده‌اند. اعداد نوشته شده بر هر ستون نشان دهنده‌ی تعداد تصاویری که حداقل شامل N طبقه شیء مختلف بوده‌اند را نشان می‌دهد.

پیش‌بینی اینکه چه اشیائی در تصویر وجود دارند در درک محتوای صحنه موثر است و می‌تواند برای پرسش در مورد تصاویری که شامل چند شیء مشخص هستند استفاده شود. نتایج بیشتر این مدل در فصل بعد نشان داده می‌شوند و نتایج الگوریتم پیشنهادی نیز با بخش عملکرد پیش‌بینی حضور مقایسه خواهد شد.

۳-۳: مدل مبنی بر اطلاعات متنی معنایی

در مقاله [17] گامی برای استفاده از اطلاعات متنی معنایی به عنوان پس پردازش در هر مدل بازشناسی شیء ارائه شده است. در این روش توسط یک چهارچوب میدان تصادفی شرطی^{۵۳}، سازگاری برچسب اشیاء طبق ارتباطات متنی پیشینه می شود. مدل کیف ویژگی^{۵۴} توسط استفاده از تعاملات متنی بین اشیاء، توسعه داده شده است. به عنوان گام پیش پردازش ابتدا بخش بندی انجام شده است که این کار باعث دسته بندی مکانی در مدل بازشناسی می شود. فلوجارت این روش در شکل زیر نشان داده شده است.



شکل 6-3- طبقه بندی شیء با استفاده از اطلاعات متنی معنایی. $S_1 \dots S_k$ مجموعه k بخش یک تصویر که توسط چندین بخش بندی بدست آمده است را نشان می دهد. $L_1 \dots L_n$ لیست n برچسب هر بخش را نشان می دهد. $O_1 \dots O_m$ نیز مجموعه m طبقه شیء در تصویر اصلی است.

برای استفاده از اطلاعات متنی معنایی در طبقه بندی شیء از چهارچوب میدان تصادفی شرطی استفاده شده است که در ای روش دو تفاوت عمده دارد. اول اینکه بین برچسب بخش ها یک گراف کاملاً متصل استفاده شده است و دوم اینکه به جای یکی کردن مدل متنی با مدل طبقه بندی، میدان تصادفی شرطی بر مسائل ساده تر بر تعداد کمتری بخش آموزش داده شده است.

با داشتن تصویر I و بخش های $S_1 \dots S_k$ هدف یافتن برچسب های $C_1 \dots C_k$ است به گونه ای که این برچسب ها متناسب با محتوای بخش ها باشند و نسبت به یکدیگر متناسب متنی داشته باشند. فرض شده است که برچسب ها از یک مجموعه محدود C هستند.

⁵³ Conditional Random Field (CRF)

⁵⁴ Bag-of-features (BoF)

این روابط به عنوان توزیع احتمالی به صورت زیر مدل می شوند:

$$p(C_1 \dots C_k | S_1 \dots S_k) = \frac{B(C_1 \dots C_k) \prod A(i)}{Z(\Phi, S_1 \dots S_k)} \quad (3 - 12)$$

$$A(i) = p(C_i | S_i); \quad B(C_1 \dots C_k) = \exp\left(\sum_{i,j=1}^k \phi(C_i, C_j)\right) \quad (3 - 13)$$

$Z(0)$ تابع بخش بندی است. بخش حاشیه ای $p(C|S)$ توسط سیستم بازشناسی فراهم می شود.

برای ترکیب اطلاعات متنی معنایی با بازشناسی شیء و در چهارچوب میدان تصادفی شرطی، ماتریس‌هایی ساخته می شوند. این ماتریس‌ها مقارن و غیرمنفی هستند که شامل فرکانس هم‌رخداد برچسب‌های اشیاء در مجموعه آموزشی مجموعه داده هستند. به عنوان مثال ماتریس زیر از مجموعه داده MSRC بدست آمده است.

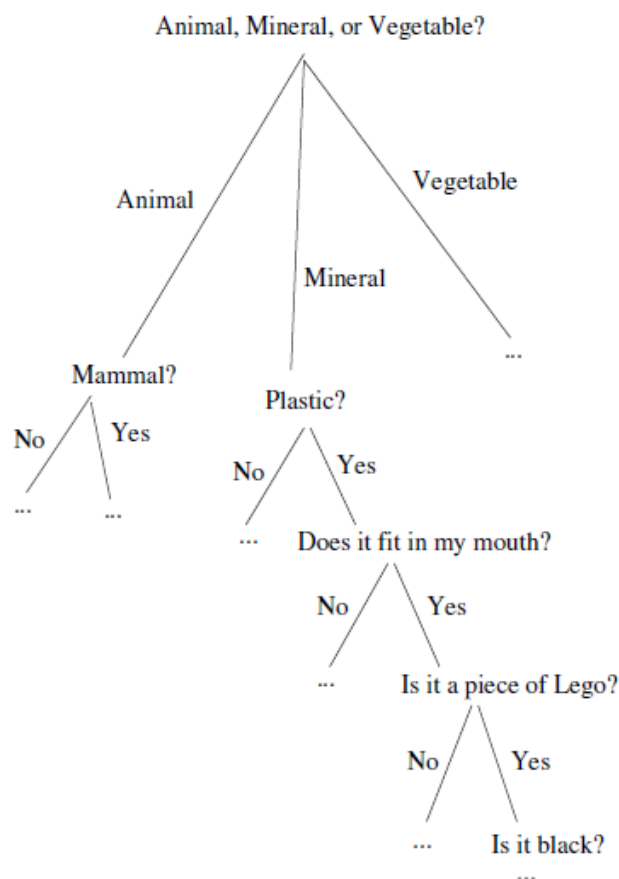
MSRC training data

building	75	18	29		33	6	9	7	10	10	2	1		43	1	6	6
grass	18	93	28	23	15		39	14	7	7	3	1	1	4	15	2	8
tree	39	38	88	6		43	6	12	9	4	1	2		1	19	11	8
cow		23	6	23		4		4									
sheep		15		15				1							2		
sky	33	39	43	1		88	15	18	4	3		5	4	25		8	17
aeroplane	6	14	6		15	15								5			
water	9	7	12	4	1	18		43	4	1		7		6		6	12
face	7	7	8		4		4	28	1	1	1		3	7		28	1
car	18		4		3	1	1	20						19		1	
bike	10	3						1	15					12		1	
flower		1	1					1		1						1	
sign	2	2			3						8			1		1	
bird	1	1			4		7				14			0			
book										3			3			3	
chair		4	1									7	3				
road	43	15	19		2	26	5	8	7	19	12	1	3	3	86	7	10
cat															7	7	
dog	1	2												10	13		1
body	9	8	11			8	8	28	1	1	1	1	3	8		32	2
boat	6	8			11	10	1							1	1	2	18
building																	
grass																	
tree																	
cow																	
sheep																	
sky																	
aeroplane																	
water																	
face																	
car																	
bike																	
flower																	
sign																	
bird																	
book																	
chair																	
road																	
cat																	
dog																	
body																	
boat																	

شکل 7-3- ماتریس متنی فرکانس هم‌رخدادی اشیاء برای مجموعه داده MSRC

۳-۴: استفاده از درخت تصمیم در بازشناسی شیء

روشی که در مقاله بلک وال و آستین⁵⁵ استفاده شده است، با این فرض است که شیء از تصویر جدا شده است و فقط باید تشخیص داده شود که چه شیئی است. بخش کوچکی از یک نمونه درخت تصمیم در شکل زیر نشان داده شده است. هر شیء در دیتابیس دارای مقادیری برای ویژگی های مختلف است، توسط بخشی از این دیتابیس درخت تصمیم ساخته می شود و توسط بقیه داده ها تست می شود. روش متداول، تعریف مجموعه مناسبی از ویژگی ها برای دیتابیس است. هدف در این مقاله توسعه ی ویژگی های جدید نبوده است و ویژگی های نسبتا ساده و شناخته شده استفاده شده اند [18].



شکل 8-3- بخش کوچکی از یک نمونه درخت تصمیم [18]

در مقاله ولیکنگ و روفر⁵⁶ فرآیند بازشناسی روبات که برمبنای کلاسبندی درخت تصمیم است، ارائه می شود. با انجام پیش پردازش هایی مانند بخش بندی برمبنای رنگ، بازیابی مرز و تولید خط ویژگی

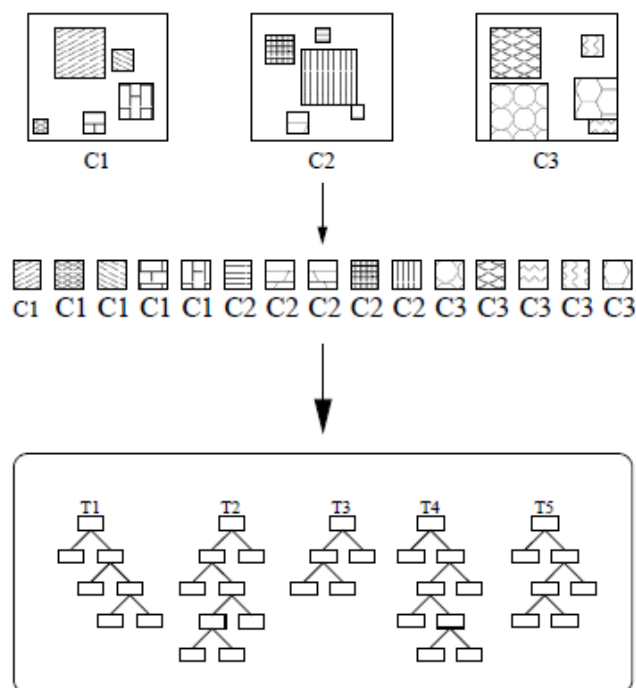
⁵⁵ Blackwell, P & Austin, D

⁵⁶ D. Wilking, T. Röfer

های مختلف ارائه می شوند. (به عنوان مثال مساحت، محیط و زاویه همسایگی) برای ساخت درخت ویژگی با بیشترین آنتروپی محاسبه می شود. با هرس درخت می توان از فرایادگیری آن جلوگیری کرد [19].

در مقاله کی، وانگ و همکارانشان⁵⁷ روشی برای بازشناسی شیء ارائه شده است که یادگیری آدابوست⁵⁸ را با استفاده از روش درخت تصمیم و یک معماری سلسله مراتبی جدید با ویژگی های محلی هیستوگرام ترکیب می کند [20].

مقاله مری و همکارانش⁵⁹ براساس استخراج رندم زیرپنجره هایی به صورت تکه های مربعی و کلاسبندی آنها توسط درخت تصمیم است. زیرپنجره ها از موقعیت هایی به صورت رندم و با سایز رندم انتخاب می شوند و پس از تغییر سایز درخت تصمیم براساس آنها ساخته می شود. این مراحل در شکل زیر نشان داده شده است [21].



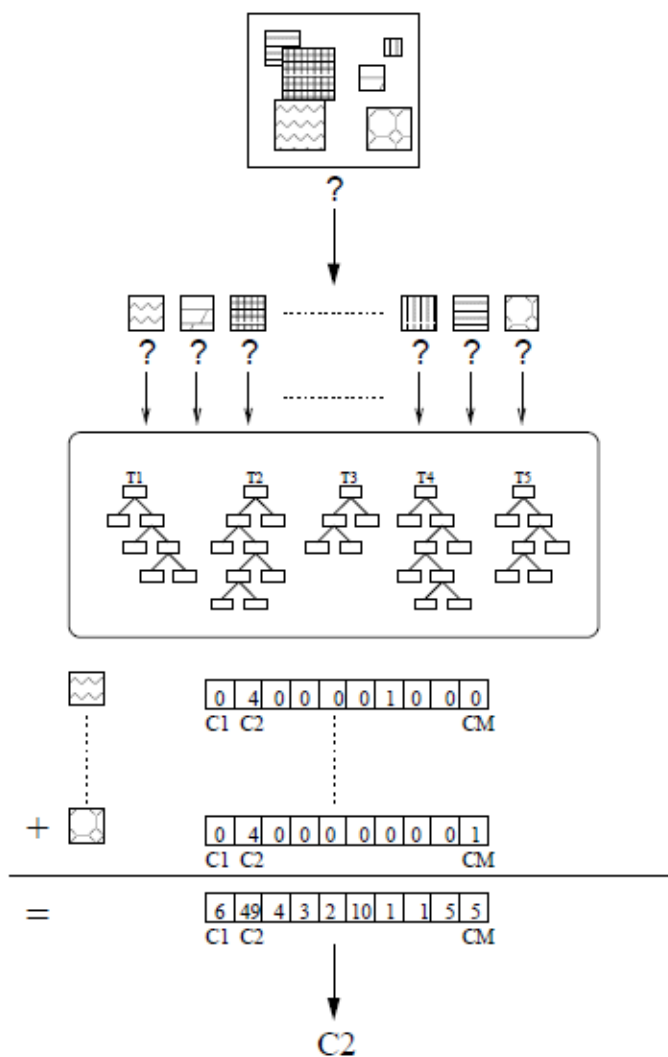
شکل 9-3 - فرآیند یادگیری: استخراج رندم زیرپنجره ها از تصاویر آموزشی، تغییر سایز آنها و سپس ساخت درختان تصمیم [21]

⁵⁷ Qi, wang et al..

⁵⁸ Adaboost

⁵⁹ R Marée, P Geurts, J Piater, L Wehenkel.

سپس برای کلاسبندی یک تصویر تست، هر درخت تصمیم احتمال شرطی برای هر زیرپنجره تخمین می زند. تمام تخمین ها متوسط گرفته می شوند و سپس کلاس متعلق به بیشترین احتمال به تصویر تخصیص می یابد. این مراحل در شکل زیر نشان داده شده اند. در اینجا 5 درخت آموزش داده شده است. آموزش یک درخت نیز حالت خاصی از این مساله خواهد بود.



شکل 10-3- فرآیند بازشناسی: زیرپنجره هایی که به صورت رندم استخراج شده اند بین درخت ها منتشر می شوند و پس از جمع آوری رای ها کلاس با رای اکثریت به تصویر تخصیص می یابد [21].

در مقاله ی ماتاس⁶⁰ حافظه ی بصری به عنوان درخت تصمیم باینری سازماندهی می شود که زمان تصمیم میانگین را کمینه کند. برگ های این درخت تصمیم بعضی از نواحی محلی را نمایش

⁶⁰ Obdrzalek, S. and J. Matas

می دهند و هر گره غیره برگ مربوط به یک کلاسبند ضعیف است. در فاز بازشناسی، یک معیار ثابت مشخص می کند که یک بخش تصویر در کدام زیردرخت جستجو شود [22].

هر تکنیکی که ورودی بصری جاری را یکی یکی با مدل های ذخیره شده مقایسه کند خطی است. تکنیک هایی که مساله ی آنها تطبیق دو تصویر است، زمان پاسخ قابل قبولی ندارند. جستجو و اندیس گذاری موضوعاتی هستند که مطالعات بسیاری در رابطه با آنها انجام شده است و دو روش درهم سازی و جستجوی درختی فراگیرتر هستند.

۳-۵: معرفی الگوریتم کارت

این الگوریتم سرانجام درخت کلاسبندی است که انقلابی در زمینه ی آنالیز پیشرفته ایجاد کرده است و یکی از مهم ترین ابزارهای داده کاوی مدرن است [23]. یک درخت کارت یک درخت تصمیم باینری است که با تقسیم هر گره به دو گره فرزند به طور مکرر ایجاد می شود. گره ریشه شامل تمام نمونه های آموزشی است. متغیرهای مختلف در جدول زیر معرفی شده اند [24].

جدول 1-3- تعریف متغیرهای استفاده شده در تعریف الگوریتم ساخت درخت تصمیم

نام	تعریف
Y	متغیر هدف که می تواند مقداری مطلق وصفی یا اسمی و یا مقدار پیوسته داشته باشد. در صورتی که شامل j کلاس باشد یکی از مقادیر $C = \{1, \dots, j\}$ را خواهد داشت.
$X_m, m=1, \dots, M$	مجموعه ی متغیرهای مورد استفاده برای پیشگویی که می تواند مقادیر مطلق وصفی یا اسمی و یا مقادیر پیوسته داشته باشد.
$h = \{X_n, Y_n\} n=1, \dots, N$	کل نمونه های یادگیری
h(t)	نمونه های یادگیری که در گره t قرار می گیرند.
w_n	وزن مورد n

وزن فرکانس مربوط به مورد n	f_n
احتمال پیشین $y=j, j=1, \dots, J$	$\pi(j), j=1, \dots, J$
احتمال اینکه یک مورد هم در کلاس j باشد و هم از گره t عبور کند.	$p(j, t)$
احتمال عبور از گره t	$p(t)$
احتمال اینکه یک مورد از کلاس j باشد به شرطی که بدانیم از گره t عبور می کند.	$p(j t)$
هزینه کلاسبندی اشتباه کلاس j به عنوان کلاس i در نتیجه $C(j i) = 0$	$C(i j)$

برای ساخت درخت باید در مورد اینکه چگونه در هر گره داده‌ها تقسیم شود، تصمیم‌گیری شود. این تقسیم‌بندی از بین تمام تقسیم‌بندی‌ها باید به گونه‌ای انجام شود که گره‌های فرزند به خالص‌ترین حالت برسند. در هر گره تقسیم‌بندی به مقدار یک متغیر وابسته است. در صورتی که X یک مقدار پیوسته با K مقدار مختلف داشته باشد، $K-1$ تقسیم‌بندی مختلف بر X می‌تواند وجود داشته باشد. یک درخت با گره‌ی ریشه شروع می‌شود و با تکرار گام‌های زیر ایجاد می‌شود.

روش الگوریتم کارت

گام 1. یافتن بهترین تقسیم‌بندی هر متغیر

برای هر متغیر X ، تمام تقسیم‌بندی‌های ممکن امتحان می‌شود.

گام 2. پیدا کردن بهترین تقسیم‌بندی گره

انتخاب تقسیم‌بندی که معیار مشخص شده را بیشینه می‌کند، در بین بهترین تقسیم‌بندی‌های یافت

شده در گام 1

گام 3. انجام تقسیم‌بندی یافت شده در گام دوم در صورتی که شرط پایان برقرار نیست

۳-۵-۱: معیار تقسیم‌بندی و مقیاس خالص بودن

در گره t بهترین تقسیم‌بندی S به گونه‌ای انتخاب می‌شود که معیار تقسیم $\Delta i(s,t)$ را بیشینه کند. در گره t مقادیر $p(j,t)$ ، $p(t)$ و $p(j|t)$ به صورت زیر بدست می‌آیند:

$$p(j, t) = \frac{\pi(j)N_{w,j}(t)}{N_{w,j}}, \quad (3-14)$$

$$p(t) = \sum_j p(j, t), \quad (3-15)$$

$$p(j|t) = \frac{p(j, t)}{p(t)} = \frac{p(j, t)}{\sum_j p(j, t)}. \quad (3-16)$$

که $N_{w,j}$ و $N_{w,j}(t)$ به صورت زیر تعریف می‌شوند:

$$N_{w,j} = \sum_{n \in h} (w_n f_n I(y_n = j)) \quad (3-17)$$

$$N_{w,j}(t) = \sum_{n \in h(t)} (w_n f_n I(y_n = j)) \quad (3-18)$$

تابع $I(a=b)$ تابع شاخصی است که در صورتی که $a=b$ باشد مقدار 1 و در غیر این صورت دارای مقدار صفر است.

در صورتی که Y پیوسته باشد، معیار تقسیم زیر مورد استفاده قرار می‌گیرد. معیار ناخالصی نیز همانطور که در معادله 7 نشان داده شده است، با مربع کمینه‌ی انحراف معیار^{۶۱} تعیین می‌شود.

$$\Delta i(s, t) = i(t) - p_L i(t_L) - p_R i(t_R) \quad (3-19)$$

$$i(t) = \frac{\sum_{n \in h(t)} (w_n f_n (y_n - y(t))^2)}{\sum_{n \in h(t)} (w_n f_n)}, \quad (3-20)$$

⁶¹ Least Squares Deviation (LSD)

$$p_L = \frac{N_w(t_L)}{N_w(t)}, p_R = \frac{N_w(t_R)}{N_w(t)},$$

$$N_w(t) = \sum_{n \in h(t)} (w_n f_n), \quad (3-21)$$

$$y(t) = \frac{\sum_{n \in h(t)} (w_n f_n y_n)}{N_w(t)} \quad (3-22)$$

۳-۵-۲: قوانین پایان

کنترل فرآیند رشد درخت تصمیم توسط قوانین پایان کنترل می‌شود. قوانین پایان عبارتند از:

- اگر یک گره خالص شود، یعنی تمام موارد در یک گره مقدار یکسانی از تابع هدف را داشته باشند.
- اگر تمام موارد یک گره مقادیر یکسانی برای تمام متغیرهای ورودی داشته باشند.
- اگر کاربر عمق ماکسیمم را برای درخت مشخص کرده باشد و به آن عمق رسیده باشد.
- اگر کاربر مینیمم سایز یک گره را مشخص کرده باشد و سایز یک گره کمتر از آن باشد.
- اگر تقسیم‌بندی یک گره باعث ایجاد فرزندی می‌شود که سایز آن کمتر از مقدار مینیمم مشخص شده توسط کاربر باشد.
- اگر بهترین تقسیم بندی در گره t ، باعث بهبودی شود که این بهبود از حداقل بهبود تعیین شده توسط کاربر کمتر باشد.

۳-۶: خلاصه و نتیجه گیری

بازشناسی شیء یکی از مسائل پیچیده در بینائی کامپیوتر است. بخشی از این پیچیدگی به دلیل ابهامات در ویژگی های سطح پایین است که به دلیل شلوغی زمینه، موانع و غیره ایجاد می شود. مطالعات اخیر نشان داده است که مدل کردن روابط متنی می تواند به برطرف کردن این پیچیدگی ها کمک کند و عملکرد بازشناسی شیء را بهبود بخشد [25]. در این فصل چند مورد الگوریتم های بازشناسی شیء بررسی شدند و جزئیات و مراحل آنها شرح داده شدند. در فصل بعد الگوریتم پیشنهادی بر مبنای این الگوریتم ها و با بهبود عملکرد ارائه می شود.

فصل ۴:

بررسی و تحلیل الگوریتم پیشنهادی

۴-۱: مقدمه

در این فصل به بررسی الگوریتم پیشنهادی پایان نامه پرداخته می شود. در این الگوریتم تلاش می شود ویژگی های متنی استخراج شود و توسط ساخت کلاسبند با توجه به ارتباطات متنی دوبه دوی اشیاء و با توجه به تعداد اشیاء در هر تصویر، دقت بازیابی شیء افزایش یابد.

با توجه به مطالب مطرح شده در فصل قبل چند اشکال اساسی در روش های بازیابی شیء و روش هایی که مبتنی بر متن هستند وجود دارد. در ادامه این فصل ابتدا به بررسی چالش هایی که در الگوریتم های پیشین که به تفصیل در فصل دوم بیان شد، می پردازیم و دلایل لزوم الگوریتم پیشنهادی را بررسی می کنیم. سپس روشی ارائه می دهیم تا از مزایای این روشها سود ببرد به طوری که اشکالات آنها را نیز دارا نباشد. سپس الگوریتم به صورت مرحله به مرحله تحلیل و بررسی می شود.

۴-۲: چالشهای موجود در روش های قبلی

در این بخش به چالش های الگوریتم های بازیابی شیء می پردازیم و اینکه چگونه روش های متنی می توانند این چالش ها را برطرف کنند و همچنین به طور ویژه چالش های مربوط به الگوریتم های بازیابی شیء مبتنی بر متن را بررسی می کنیم. سپس بررسی می کنیم که الگوریتم پیشنهادی چگونه سعی در برطرف کردن این چالش ها دارد.

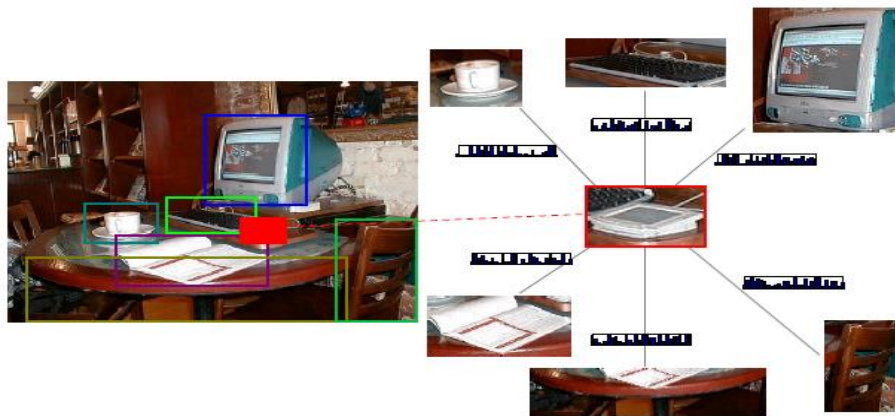
۴-۲-۱: چالش های روش های بازیابی شیء غیر متنی

در این قسمت چالش هایی که معمولا الگوریتم های بازیابی شیء با آن ها روبه رو هستند، مطرح می شود و راه حل برخورد با آنها بررسی می شود. برخی از این چالشها عبارتند از:

- مبهم بودن و ناآشنا بودن شیء
- کوچک بودن شیء مورد نظر
- کافی نبودن ویژگی های محلی شیء
- رزولوشن پایین
- وجود مانع در مقابل شیء

شکل زیر چند نمونه از تصاویری را نشان می‌دهد که دارای این دشواری‌ها هستند. در این شکل تصویر اول بیانگر چالش اول است و همان‌طور که در تصویر مشاهده می‌شود به دلیل ناآشنا بودن شیء نمی‌توان توسط روش‌های متداول این شیء را بازیابی کرد.

تصویر دوم نشان دهنده چالش دوم می‌باشد و همان‌طور که مشاهده می‌شود کوچک بودن شیء باعث شده تا شیء موردنظر به اشتباه بازیابی شود و نمی‌توان توسط روش‌های متداول این مشکل را برطرف نمود. تصویر سوم نیز چالش سوم را دربرمی‌گیرد که ویژگی‌های محلی برای شناسایی کافی نبوده‌اند و بازیابی به درستی صورت نگرفته است.



شکل 1-4 - برخی دشواری‌هایی که معمولاً در بازیابی شیء با آنها روبرو هستیم. تصویر اول بیانگر این مساله است که بعضی اشیاء که دارای ظاهری ناآشنا و مبهم هستند که بازیابی آنها توسط روش‌های متداول امکان‌پذیر نیست. تصویر دوم واقعیت را بیان می‌کند که بعضی اشیاء کوچک هستند و به

همین دلیل نمی توان آنها را به درستی بازیابی کرد. تصویر سوم نیز ناکافی بودن ویژگی‌های محلی در بعضی تصاویر را نشان می‌دهد که باعث بازیابی خطا می‌شود.

همانطور که دیده می‌شود این چالش‌ها را نمی‌توان توسط اطلاعات شیء موردنظر برطرف نمود و نیاز به استفاده از اطلاعات دیگر که همان اطلاعات متنی هستند، روشن می‌شود. در روش‌های متنی توسط اطلاعات مربوط به اشیاء دیگر تصویر می‌توان شیء موردنظر را بازیابی نمود.

۴-۲-۲: چالش‌های روش‌های بازیابی شیء مبتنی بر اطلاعات متنی

در روش‌های مبتنی بر اطلاعات متنی نیز چالش‌هایی وجود دارد که در این پایان نامه تلاش شده است تا بعضی از این چالش‌ها برطرف شود و دقت بازیابی شیء افزایش یابد. در ادامه به این چالش‌ها پرداخته می‌شود.

مهم‌ترین نکته در روش‌های مبتنی بر متن این است که این روش‌ها باید بر تصاویری که دارای اشیاء مختلف و از صحنه‌های مختلف هستند تست شوند تا روابط اشیاء و صحنه بتوانند به بهبود فرآیند بازشناسی شیء کمک کنند. روش‌های متنی بسیاری ارائه شده‌اند اما به دلیل اینکه بر تصاویر با تعداد شیء کم تست شده‌اند، مزیت روش‌های متنی خود را نشان نداده است.

مقاله مبتنی بر اطلاعات سلسله مراتبی [4] که در فصل قبل به تفصیل بررسی شد و جزء بهترین الگوریتم‌هایی است که تاکنون ارائه شده است، و برای تست این الگوریتم از مجموعه داده سان با اشیاء مختلف و صحنه‌های متنوع استفاده شده است. همان‌طور که در فصل قبل توضیح داده شد، در این مدل ارتباط بین اشیاء را به صورت سلسله مراتبی ذخیره می‌کند. به عنوان مثال بیشتر اشیائی که معمولاً در یک آشپزخانه وجود دارند، در زیردرخت گرهی sink قرار می‌گیرند و تمام وسایل نقلیه نیز در زیردرخت شیء جاده قرار می‌گیرند. مساله‌ای که در این مدل وجود دارد این است که ممکن است در تصاویری وسایل آشپزخانه وجود داشته باشند اما شیء sink وجود نداشته باشد، در واقع در این مدل برای بازیابی هر شیء تنها از اطلاعات مربوط به والد‌های آن تا ریشه استفاده می‌شود و اشیاء دیگری که ممکن است آنها نیز موثر باشند در نظر گرفته نمی‌شود که این چالش در مدل پیشنهادی بررسی شده است.

موارد دیگری که می‌توان در در این مدل به آنها اشاره کرد این است که استفاده از روابط فیزیکی بین اشیاء مانند اینکه یک ماشین همیشه بالای یک جاده قرار دارد و موارد مشابه می‌تواند به طور قابل

توجهی پیچیدگی محاسباتی را کاهش دهد و فرآیند بازیابی شیء را بهبود بخشد که در این مقاله به آنها پرداخته نشده است. این مساله با توجه به اینکه مربوط به مدل مکانی مقاله مرجع است و در حوزه پایان نامه قرار ندارد، در این رساله بررسی نشده است و در حوزه کارهای آینده قرار می‌گیرد.

۴-۳: رویکرد کلی الگوریتم پیشنهادی

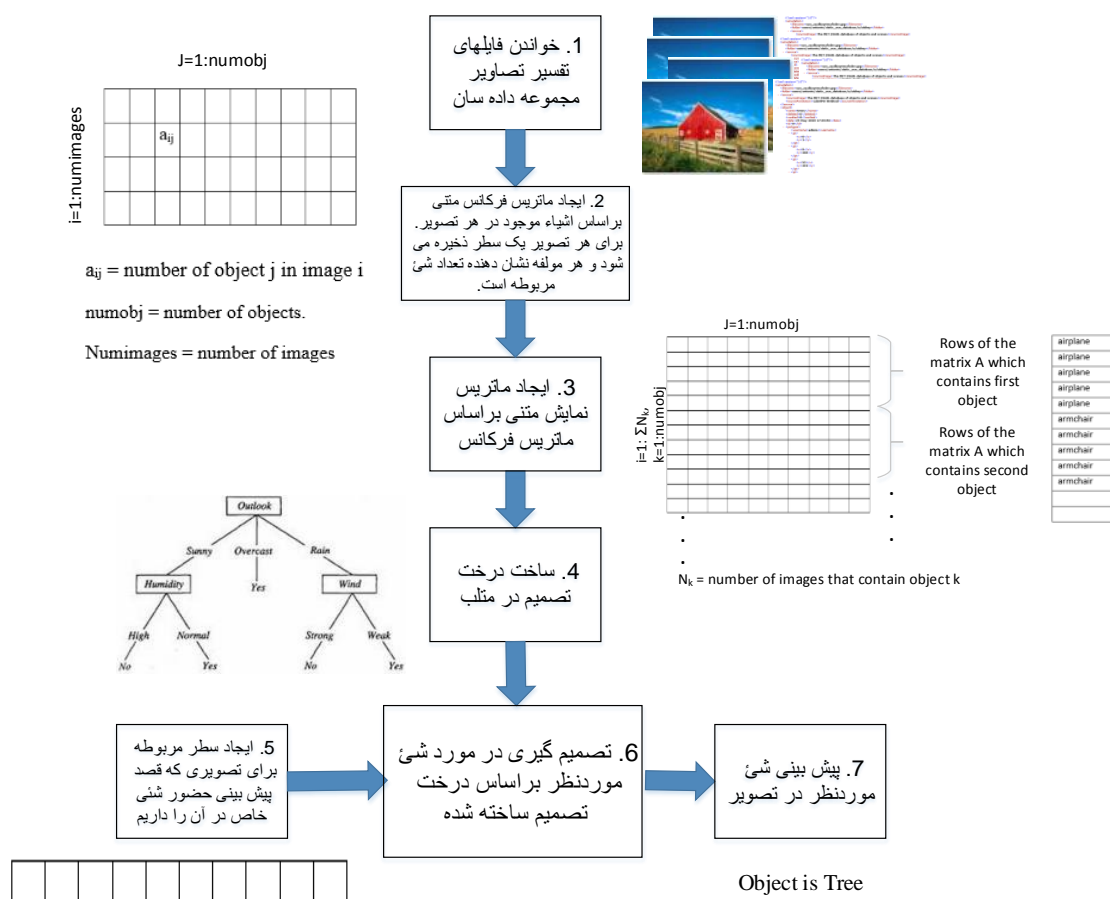
در این بخش به بیان رویکرد کلی الگوریتم پیشنهادی می‌پردازیم. همانطور که تاکنون توضیح داده شده است، اطلاعات متنی معنایی یک تصویر از ارزشمندترین اطلاعات متنی است و برای بازشناسی شیء مفید است. این اطلاعات شامل صحنه و اشیاء دیگر موجود در تصویر است. هدف این پایان نامه استفاده از این اطلاعات برای بازشناسی شیء است.

اطلاعات متنی موردنظر اشیاء دیگر و تعداد آنها هستند که برای ساخت کلاسبند و آموزش آن از فایل‌های تفسیر تصاویر آموزشی مجموعه داده سان استفاده می‌شود. برای این کار فایل‌های تفسیر موجود در مجموعه داده سان را به صورت ساختار ذخیره می‌کنیم و سپس به استخراج ویژگی پرداخته می‌شود و اطلاعات موردنیاز این فایل‌ها استخراج و ذخیره می‌شوند.

به دلیل وجود اشیاء مختلف و صحنه‌های متنوع در مجموعه داده‌ی سان [16] و عدم درنظر گرفتن تمام اشیاء تصویر و تعداد آنها در مدل بازشناسی شیء مبتنی بر اطلاعات سلسله مراتبی [4]، عملکرد این مدل بر مجموعه داده سان چندان قابل قبول نیست. در این بخش مدلی مبتنی بر اطلاعات متنی ارائه می‌شود که برای تصاویر با اشیاء زیاد موثر باشد. اطلاعات متنی موردنظر اشیاء دیگر تصویر و فرکانس آنها می‌باشد. این الگوریتم شامل سه مرحله اصلی می‌باشد: اول ساخت ماتریس فرکانس متنی براساس تمام تصاویر آموزشی، دوم ساخت درخت تصمیم براساس ماتریس بدست آمده در مرحله قبل و سوم استفاده از سطر متنی هر تصویر تست برای تعیین برچسب شیء موردنظر در داده‌های تست. در تصویر زیر فلوجارت این الگوریتم نشان داده شده است.

در قسمت (1) تصاویر مختلف مجموعه داده به همراه فایل‌های تفسیر که توسط مجموعه داده‌ی سان فراهم شده است، نشان داده شده است. در قسمت (2) که مربوط به مرحله اول الگوریتم پیشنهادی می‌باشد، ماتریس بدست آمده از فایل‌های تفسیر نشان داده شده است که هر سطر از این ماتریس بردار ویژگی یک تصویر است و شامل اطلاعات متن تصویر می‌باشد که جزئیات آن در ادامه توضیح داده خواهد شد. برای استفاده از این ماتریس برای آموزش درخت تصمیم لازم است برای هر شیء نمونه‌های

آموزشی به طور مستقل بدست آید. به همین دلیل در قسمت (3) ماتریس نمایش متنی براساس ماتریس قبلی بدست می آید. چگونگی ساخت این ماتریس در بخش بعد توضیح داده خواهد شد. قسمت (4) شامل مرحله ساخت کلاسبند که یک درخت تصمیم است می باشد. تا این مرحله درخت تصمیم ساخته شده است. حال برای شناسایی یک شی خاص در یک تصویر جدید با استفاده از اطلاعات متنی به این صورت عمل می شود: (5) بدست آوردن بردار ویژگی متنی مربوطه از تصویر موردنظر (6) کلاسبندی و تعیین کلاس شی موردنظر با توجه به درخت تصمیم ساخته شده در مرحله قبل. در بخش بعد این مراحل با جزئیات بیشتری تشریح خواهند شد.



شکل 2-4- فلوجارت الگوریتم پیشنهادی





فایل های تفسیر در واقع شامل اطلاعات مختلفی از جمله اشیاء تصویر، محل آنها و مواردی از این دست است. در بخش بعد جزئیات هر کدام از مراحل با شرح بیشتری توضیح داده خواهد شد.

۴-۴: بررسی دقیق الگوریتم پیشنهادی پایان نامه

برای این روش از بخشی از فایل های تفسیر مجموعه داده سان که شامل بیش از 8000 تصویر از بیش از 600 صحنه مختلف است، استفاده شده است. در این تصاویر بیش از 200 شیء مختلف وجود دارد که می توان آنها را توسط روش ارائه شده شناسایی کرد. هدف اثبات تاثیر کانتکست و اطلاعات متن تصویر در بازشناسی شیء است. در این روش با دانستن حضور اشیاء دیگر تصویر و فرکانس آنها می توان شیء موردنظر را شناسایی کرد.

3.1. ساخت ماتریس فرکانس متنی

ابتدا با استفاده از تصاویر مجموعه داده سان ماتریسی ساخته می شود که به تعداد تصاویر سطر دارد. ستون های این ماتریس نیز به تعداد شیء های موردنظر است که در این مقاله 107 شیء مشابه با مقاله مبتنی بر اطلاعات متنی سلسله مراتبی انتخاب شده است تا نتایج قابل مقایسه باشند. اطلاعات لازم برای ساخت این ماتریس، بعضی از اشیاء دیگر تصویر است. هر مولفه a_{ij} در این ماتریس نشان دهنده ی حضور شیء j در تصویر i است. برای دخالت دادن رابطه بین اشیاء از وقوع مشترک آنها استفاده شده است. در نتیجه تمام اشیاء موجود در تصویر به غیر از شیئی که قرار است شناسایی شود در این ماتریس مقدار خواهند گرفت. شیئی که قرار است شناسایی شود نیز دارای مقدار نخواهد بود و به عنوان برجسب برای سطر موردنظر ذخیره می شود. در تصویر زیر بخشی از این ماتریس برای چهار تصویر نشان داده شده است. در تصویر اول شیئی که باید بازیابی شود شیء $grass$ است و شیء $person$ و $floor$ و $table$ در این تصویر وجود ندارد و شیء $flower$ نیز یک نمونه در این تصویر وجود دارد. به همین ترتیب تصویر دوم دارای هیچ یک از چهار شیء $grass$ ، $flower$ ، $floor$ و $table$ نیست و هدف پیش بینی حضور شیء $person$ است. برای هر تصویر سطری با تعداد 107 مولفه که نشان دهنده ی مولفه مربوطه است ذخیره می شود.

	grass	person	flower	floor	table
	NaN	0	1	0	0
	0	NaN	0	0	0
	0	0	4	NaN	0
	NaN	0	5	0	1

شکل 3-4 - یک مثال برای چگونگی ذخیره ماتریس فرکانس متنی

تعداد اشیاء توسط متغیر `numobj` و تعداد تصاویر توسط متغیر `numimages` ذخیره می‌شوند. حال اشیائی که در هر سطر می‌خواهیم شناسایی کنیم را نیز به صورت جدا در برداری به شکل زیر ذخیره می‌کنیم. از این ماتریس و بردار زیر که براساس همه تصاویر آموزشی و اشیاء ساخته می‌شود، برای ساخت درخت تصمیم استفاده می‌کنیم.

grass
person
floor
grass

شکل 4-4 برداری که شامل اشیائی هست که می‌خواهیم شناسایی شود و برای آموزش به درخت تصمیم داده می‌شود.

به این ترتیب در صورتی که این ماتریس و بردار را برای تصاویر بیشتری از مجموعه داده‌ی سان بسازیم، ماتریسی با اندازه (تعداد اشیاء)* (تعداد تصاویر) خواهیم داشت. حال از این ماتریس برای ساخت ماتریس دوم استفاده می‌کنیم. برای هر شیء تمام سطرهای مربوط به تصاویر شامل آن شیء به ماتریس جدید اضافه می‌شود و در سطرهای موردنظر برای آموزش هر شیء، مولفه‌ی مربوط به آن مقدار NaN خواهد گرفت. به این ترتیب هر تصویر به تعداد اشیاء موجود در آن در این ماتریس تکرار خواهد شد، اما هر بار برای آموزش بازشناسی یکی از اشیاء موجود در تصویر در این ماتریس قرار می‌گیرد. در نهایت ماتریس بدست آمده را به عنوان نمونه آموزشی و برچسب اشیاء را به عنوان برچسب این نمونه‌های آموزشی برای ساخت درخت تصمیم استفاده می‌کنیم.

3.2. ساخت درخت تصمیم

استفاده از درخت تصمیم دو مزیت عمده برای بازشناسی شیء به همراه دارد، اول اینکه برای دیتاست-های بزرگ نیز قابل استفاده است و دوم اینکه ویژگی های موثر را انتخاب می کند و برای کلاسبندی یک نمونه تست نیازی به محاسبه ی تمام ویژگی ها نیست.

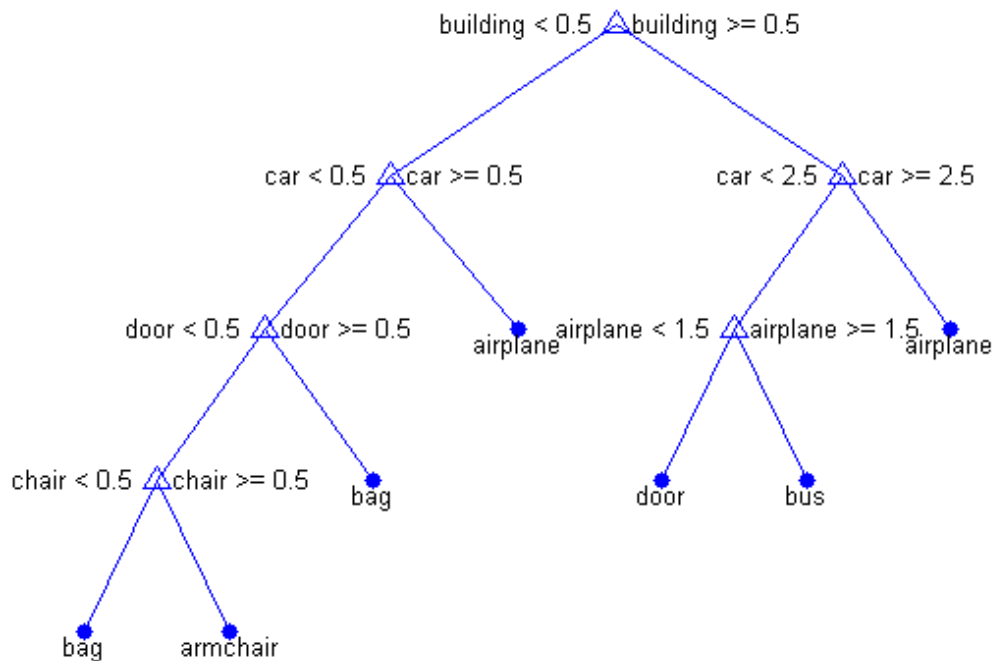
کلاسبند درخت تصمیم یکی از روش های تصمیم گیری چند مرحله ایست. ایده ی اولیه در هر روش چند مرحله ای، شکافتن یک تصمیم پیچیده به چندین تصمیم ساده تر است. یک گراف $G = (V, E)$ شامل مجموعه ای از گره های غیرتهی V و مجموعه ای از یال های E است. گره ای که فرزند نداشته باشد، برگ نامیده می شود.

در یادگیری درخت تصمیم هدف ایجاد مدلی است که مقدار تابع هدف را براساس متغیرهای مختلف ورودی بدست می آورد. هر برگ نمایش یک مقدار از متغیر هدف است که براساس مقادیری از متغیرهای ورودی که توسط مسیر از ریشه تا برگ وجود دارند، بدست می آید [26]. یادگیری درخت تصمیم یکی از موفق ترین تکنیک های یادگیری کلاسبندی باناظر است. یک درخت تصمیم یا درخت کلاسبندی، درختی است که هر گره ی داخلی در آن توسط ویژگی ورودی برچسب گذاری شده است. یک درخت با تقسیم مجموعه ورودی به زیربخش ها براساس یک مقدار ویژگی یادگرفته می شود.

الگوریتمی که در درخت تصمیم استفاده شده است الگوریتم کارت⁶² [24] است که در فصل قبل جزئیات مربوط به این الگوریتم توضیح داده شده است. این الگوریتم می تواند روابط داده های مهم را سریعاً مشخص کند. طی سالها این الگوریتم به عنوان سریع ترین و تطبیق پذیرترین الگوریتم مدل سازی پیشگو برای آنالیز بوده است و به عنوان پایه برای بسیاری از روش های داده کاوی مدرن براساس بوستینگ و غیره بوده است [23].

در الگوریتم پیشنهادی سوالهای هر گره مربوط به اطلاعات متنی یا اشیاء دیگر تصویر است که براساس آنها برچسب گره برگ یعنی شیئی که می خواهیم بازیابی کنیم بدست می آید. به عنوان مثال اگر این درخت برای تعداد شیء کمتر ساخته شود تصویری به صورت شکل 5-4 خواهد داشت:

⁶² CART



شکل 4-5- درخت تصمیم حاصل از الگوریتم پیشنهادی با در نظر گرفتن 40 شیء

تصمیم گیری براساس اعداد مشخص شده در هر گره صورت می گیرد به این ترتیب در هر گره تعداد شیء مربوطه با رابطه مربوط به گره مقایسه می شود و سپس هر کدام از شاخه ها که رابطه مربوط به آن برقرار باشد تا رسیدن به برگ پیمایش می شود. در نهایت براساس برچسب گره برگ شیء که حضور آن در تصویر مشخص نبود، یافت می شود. در صورتی که این درخت را برای کل اشیاء بسازیم، برای هر شیء قاعده ای بدست می آید که توسط آن می توان حضور یا عدم حضور آن شیء را دریافت.

۴-۵: خلاصه و نتیجه گیری

در این فصل ابتدا به بررسی چالش های موجود در روش های بازیابی شیء و سپس چالش های موجود در روش های بازیابی شیء مبتنی بر اطلاعات متنی پرداخته شد. پس از آن رویکرد کلی الگوریتم پیشنهادی ارائه و جزئیات هر یک از بخش های آن توضیح داده شد. در فصل بعد به بررسی نتایج اعمال این الگوریتم بر مجموعه داده سان 09 پرداخته می شود و نتایج آن با الگوریتم هایی که بهترین نتایج را داشته اند مقایسه می شود.

فصل ۵:

پیاده سازی، بررسی کارایی و مقایسه با روش های دیگر

۵-۱: مقدمه

در این فصل ابتدا در مورد مجموعه داده سان بیشتر توضیح داده می شود و امکاناتی که این مجموعه داده فراهم کرده است، معرفی می شود. پس از آن دلیل استفاده از این مجموعه داده بیان می شود و تصاویر آن با تصاویر مجموعه داده های دیگر مقایسه می شود. سپس عملکرد روش پیشنهادی پایان نامه ارزیابی می شود و نتایج آن با روش های دیگر مقایسه می شود.

۵-۲: مجموعه داده

برای انجام تحقیقات شناسایی شیء، شامل یادگیری مدل های بصری شیء و انواع صحنه ها و تشخیص و تعیین محل نمونه ها در مدل های مختلف در تصاویر و ارزیابی عملکرد الگوریتم های شناسایی نیاز به مجموعه داده های مناسب است [27]. در سال های اخیر دیتاست های علامت گذاری شده ی زیادی که شامل کلاس های زیادی از اشیاء هستند، به وجود آمده اند و استفاده از روش های کمی برای ارزیابی، سیستم های شناسایی را سرعت بخشیده اند. در این بخش تعدادی از دیتاست های متداول برای تست کردن کارایی کلاس بندی شیء و صحنه و الگوریتم های تعیین موقعیت شیء بررسی شده اند.

مجموعه داده سان^{۶۳} یک مجموعه داده وسیع برای درک صحنه می باشد. هدف این پایگاه داده فراهم کردن مجموعه وسیعی از تصاویر یادداشت نویسی شده از صحنه های مختلف، موقعیت ها و اشیاء با تنوع بالا برای پژوهشگران در حوزه بینایی کامپیوتر، ادراک انسانی، علوم شناختی، علوم اعصاب، یادگیری ماشین، داده کاوی، گرافیک کامپیوتر و رباتیک است.

سان 09 نیز که زیرمجموعه ای از این پایگاه داده است شامل بیش از 200 نوع شیء در رنج وسیعی از انواع صحنه هاست، شامل 12000 تصویر علامت گذاری شده و 152000 نمونه شیء علامت گذاری شده است که برای استفاده از اطلاعات متنی مناسب است. این تصاویر از منابع مختلف از جمله موتور جستجوی گوگل، فلیکر^{۶۴}، آلتاویستا^{۶۵} و لیبل می^{۶۶} بدست آمده اند. در این تصاویر هر تصویر با یک شیء یا تصاویر با زمینه های سفید حذف شده اند تا تنها تصاویر مربوط به صحنه ها نگه داشته شوند. فرآیند تفسیر با یک تفسیرگر طی یک سال با استفاده از ابزارهای لیبل می^{۶۷} انجام شده است. تصاویر

⁶³ SUN (Scene UNDERstanding)

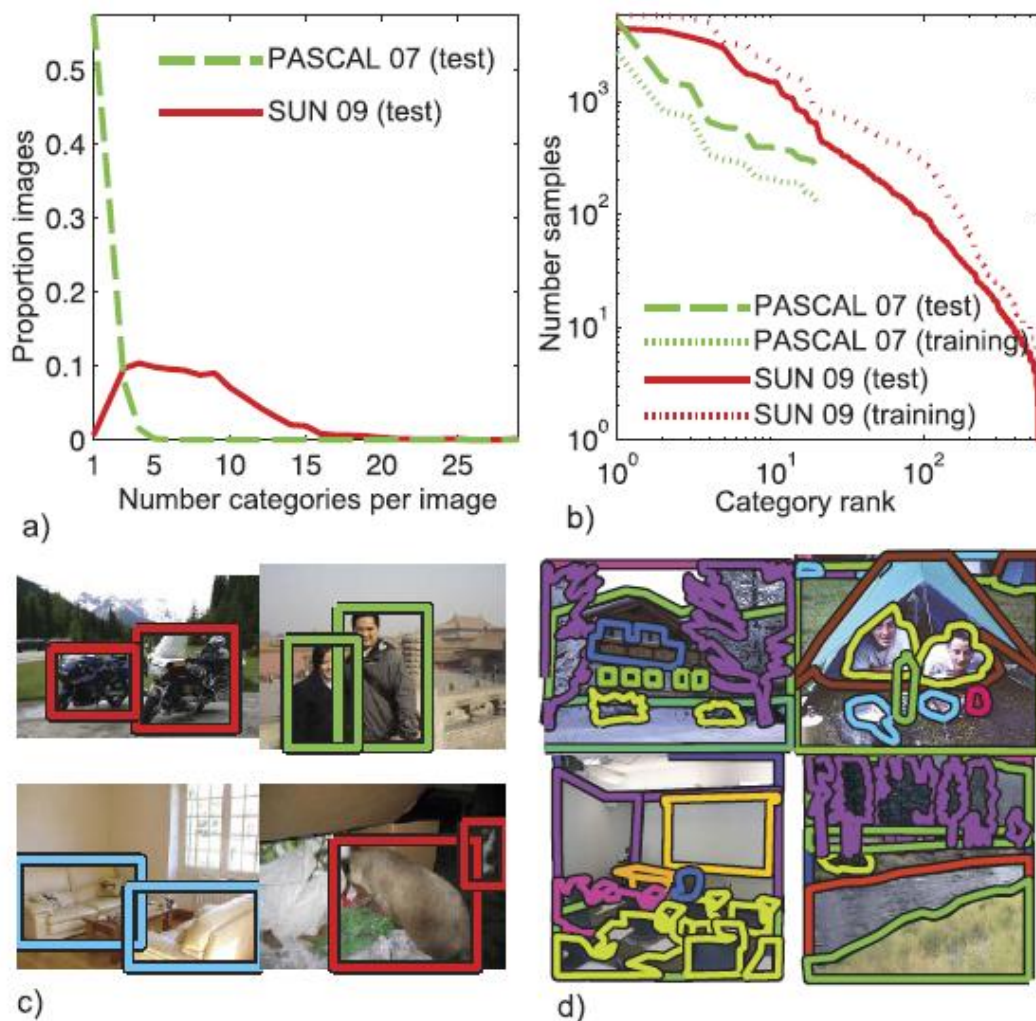
⁶⁴ Flickr

⁶⁵ Altavista

⁶⁶ Labelme

⁶⁷ <http://labelme.csail.mit.edu/Release3.0/>

برچسب‌گذاری شده به دقت برای استحکام و برچسب‌های مترادف بررسی می‌شوند و یکی می‌شوند [4]. تفسیرهای حاصل کیفیت بیشتری نسبت به لیبل می‌خواهند داشت. شکل زیر آمارهایی از این دیتاست را نشان می‌دهد و آنها را با دیتاست پاسکال 7 مقایسه می‌کند. همانطور که دیده می‌شود تعداد تصاویر با کلاس‌های شیء بیشتر از 2 شیء در مجموعه داده سان بسیار بیشتر است و در مجموعه داده‌ی پاسکال در یک تصویر حداکثر 5 شیء وجود دارد.



شکل 5-1- مقایسه‌ی دیتاست پاسکال 7 و سان 9 (a) هیستوگرام تعداد طبقه‌های شیء در هر تصویر (b) توزیع نمونه‌های آموزشی و تست نمونه تصاویر در هر طبقه از شیء (c) نمونه‌هایی از تصاویر پاسکال (d) نمونه‌هایی از تصاویر سان [4]

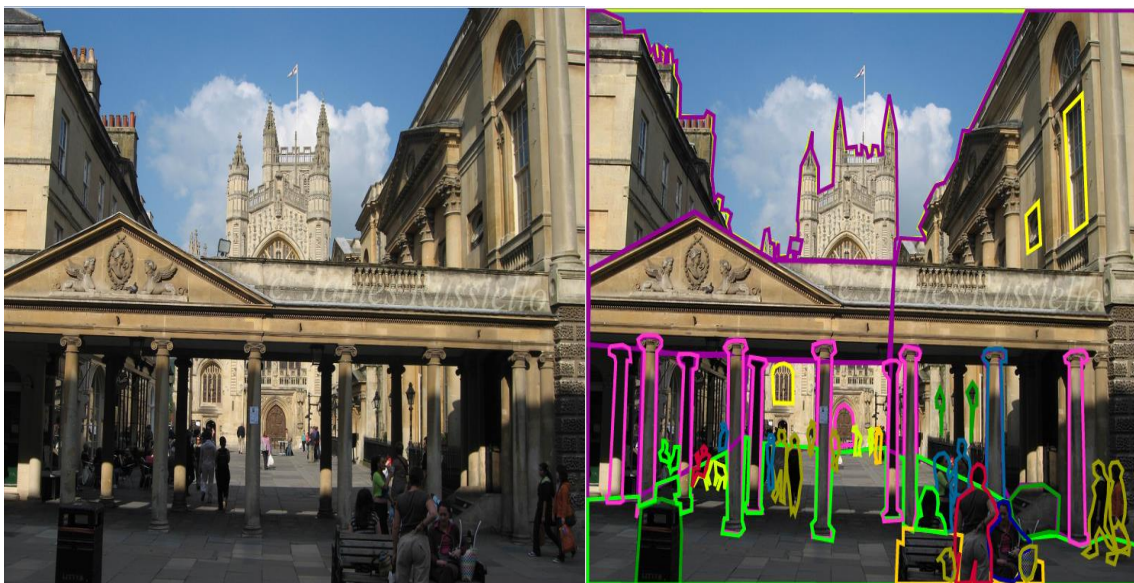
مزیت استفاده از مدل‌های متنی هنوز چندان خود را نشان نداده‌اند و دلیل آن این است که بیشتر روش‌های قبلی بر دیتاست‌هایی تست شده‌اند که تنها دسته‌های محدودی از اشیاء را شامل می‌شوند و بیشتر تصاویر شامل یک یا دو دسته شیء بوده‌اند [4].

در این گزارش از مجموعه داده سان استفاده شده است، زیرا در این مجموعه داده تصاویر براساس صحنه دسته بندی شده اند و تصاویر متعددی از صحنه های مختلف جمع آوری شده اند. تعداد زیاد اشیاء و وجود اشیاء مختلف در یک تصویر نیز از دلایل دیگر انتخاب این مجموعه داده بوده است زیرا در الگوریتم ارائه شده روابط بین شیئی که قرار است شناسایی شود با اشیاء دیگر بسیار اهمیت دارد.

آزمایشات و ارزیابی عملکرد بر تصاویر مجموعه داده سان با انواع صحنه ها انجام شده است تا کارایی الگوریتم پیشنهادی را بر گروه های مختلف تصاویر نشان دهد. تصاویر مجموعه داده سان از صحنه های مختلفی هستند و در هر تصویر اشیاء زیادی وجود دارد. بنابراین این مجموعه داده برای روش هایی که مبتنی بر روابط اشیاء و صحنه هستند مناسب تر است. این مجموعه تصاویر در مقاله های [4, 16, 25] برای هدف بازشناسی شیئی استفاده شده است و نتایج آزمایشات بر این تصاویر بدست آمده اند.

شکل زیر یکی از تصاویر مجموعه داده سان با نام صحنه کلیسا را نشان می دهد. این تصویر دارای 21 نمونه شیئی است که در شکل سمت راست هر کدام از این اشیاء نشان داده شده اند. لیست اشیاء موجود در این تصویر را می توان در شکل بعدی مشاهده نمود.

فایل های تفسیر تمام تصاویر مجموعه داده سان شامل اشیاء و صحنه و موقعیت های آنها، در این مجموعه داده موجود است. از این فایل های تفسیر می توان برای آموزش و تست الگوریتم های بازیابی شیئی استفاده کرد. در این پایان نامه هدف پیش بینی حضور هر یک از این اشیاء با دانستن بقیه اشیاء و تعداد آنها است.



شکل 2-5- نمونه ای از تصویر مجموعه داده سان از صحنه کلیسا. شکل بالا تصویر و شکل پایین تصویر که تمام اشیاء آن مشخص شده اند را نشان می دهد.

لیست اشیاء موجود در تصویر

Sky
Building
Column
Door
Window
Litter bin crop
Sidewalk
Steps
Streetlight
Spotlight
Bench occluded
Person occluded
Person sitting occluded
Person sitting crop
Person crop
Person walking
Person standing
Chair
Person sitting
Person
Bag

شکل 3-5- لیست اشیاء موجود در تصویر قبل

در ادامه برای روشن شدن دلیل استفاده از مجموعه داده سان تعدادی از تصاویر مجموعه دادهی پاسکال که در مقاله های بسیاری برای بررسی عملکرد الگوریتم ها استفاده شده اند، بررسی می شوند. شکل بعد سه نمونه از تصاویر مجموعه داده پاسکال را نشان می دهد. تصویر اول یک تصویر از باند فرود در یک فرودگاه است. این تصویر شامل دو شیء هواپیما و شخص است. تصویر دوم که تصویری از یک منظره است شامل یک شیء قایق می باشد و اشیاء دیگر در این مجموعه داده در نظر گرفته نشده اند. تصویر سوم تصویری از دو مانیتور است. این تصویر نیز تنها شامل یک شیء مانیتور است. همانطور که دیده می شود تصاویر مجموعه داده پاسکال دارای تعداد کمی شیء است و بنابراین برای ارزیابی الگوریتم های که بر مبنای روابط اشیاء هستند مناسب نیست. برای این الگوریتم ها بهتر است از تصاویری که دارای تعداد شیء زیاد است استفاده شود و تمام اشیاء در تصویر در نظر گرفته شوند.



شکل 4-5- سه نمونه از تصاویر مجموعه داده پاسکال. تصویر 1 شامل دو شیء هواپیما و شخص است. شکل 2 شامل یک شیء قایق و شکل 3 شامل یک شیء مانیاتور است.

همان طور که در فصل قبل توضیح داده شد، یکی از ابزارهای بسیار مناسب تفسیر که بصورت آنلاین فراهم شده است لیبل می است که امکان ایجاد فایل های تفسیر تصاویر مختلف را به صورت آنلاین فراهم می کند و فایل های تفسیر مجموعه داده سان توسط این ابزار فراهم شده اند. همچنین یکی از جعبه ابزارهای بسیار مهم و کاربردی نیز توسط همین پروژه ارائه شده است که امکانات بسیار گسترده و مهمی را برای کار با مجموعه داده خود و مجموعه داده سان و همچنین پاسکال فراهم می کند. این جعبه ابزار دارای بخش های مختلف و تابع های متنوع برای اعمال عملیات موردنظر بر روی فایل های تصویر و تفسیر این مجموعه داده ها می باشد. برای استفاده از این توابع می توان این جعبه ابزار را نصب و سپس از توابع مربوطه استفاده کرد.

در این پایان نامه ابتدا تمام فایل های تفسیر تصاویر مجموعه داده سان که در آدرسی که در جعبه ابزار ذکر شده است، ذخیره شده اند، دانلود و ذخیره شدند. سپس با استفاده از جعبه ابزار لیبل می این فایل ها به صورت ساختاری در برنامه متلب ذخیره می شوند. چند مورد از توابع این جعبه ابزار که توسط این پایان نامه استفاده شده اند در ادامه بیان می شوند:

✓ تابعی با نام مجموعه داده سان که پس از اجرا تمام فایل های تفسیر موجود در آدرس ذخیره شده را به صورت ساختاری ذخیره می کند.

SUNdatabase ();

✓ تابعی که نام و تعداد تمام اشیاء موجود در کلیه تصاویر مجموعه داده یا هر زیر مجموعه ای از آن را که در فایل ساختاری ذخیره شده اند را بر می گرداند.

LMobjectnames (D) ;

✓ تابعی که با گرفتن نام شیء موردنظر تعداد تکرار آن را در کل مجموعه یا در زیرمجموعه ای از آن برمی گرداند.

LMcountobject (D, objectname, method) ;

✓ تابعی که هر زیرمجموعه از مجموعه داده با شرایط مشخص شده را برمی گرداند.

LMquery (D, fieldName, content, method)

در ادامه به ارزیابی روش پیشنهادی و مقایسه آن با روش های مشابه پرداخته می شود.

۳-۵: ارزیابی روش پیشنهادی پایان نامه

در این بخش به ارزیابی نتایج این الگوریتم و مقایسه با الگوریتمی که بهترین عملکرد را بر روی مجموعه داده سان داشته است می پردازیم. این الگوریتم در فصل قبل توضیح داده شده است و در این بخش نتایج آنها بررسی می شود. در ادامه ابتدا به بیان شرایط ارزیابی و ذکر معیارهای مورد استفاده پرداخته می شود و پس از آن نتایج ارزیابی را بیان کرده و در نهایت به بررسی نتایج ارزیابی می پردازیم.

۳-۵-۱: شرایط ارزیابی

در این ارزیابی الگوریتم پیشنهادی را با الگوریتم روش مبتنی بر اطلاعات متنی سلسله مراتبی که در بخش قبل به تفصیل بررسی شدند و روش مدل های جداکننده ، مقایسه می کنیم. دلیل مقایسه با روش مدل های جدا کننده این است که این مدل در سال 2011 بهترین عملکرد را بر مجموعه داده پاسکال داشته است و یک روش غیرمتنی است. برای نشان دادن بهبود حاصل از روش متنی و بر روی مجموعه داده سان، نتیجه ای مدل برای مقایسه ذکر شده است.

تصاویر از مجموعه داده سان 09 هستند که شامل 4367 تصویر آموزشی و 4317 تصویر تست می باشد. این تصاویر شامل بیش از 200 شیء هستند که مواردی همچون جاده، آسمان، ساختمان، ماشین، تخت خواب و حتی اشیاء دگردیس پذیری چون رودخانه، پرده و حوله از جمله از آنها هستند. این تصاویر از بیش از 600 صحنه مختلف انتخاب شده اند و هدف ارزیابی 107 کلاس شیء مشخص می باشد.

۵-۳-۲: معیارهای ارزیابی

در مسائل مختلف از معیارهای ارزیابی متفاوتی استفاده می شود. در این پایان نامه نیز تلاش شده است تا الگوریتم پیشنهادی با معیارهای استاندارد مورد ارزیابی قرار گیرد. این معیارها به صورت زیر می باشند:

مثبت درست^{۶۸}: تعداد مواردی که جزء دسته مورد بررسی هستند و در این دسته قرار گرفته اند.

مثبت غلط^{۶۹}: تعداد مواردی که جزء دسته مورد بررسی نیستند ولی در این دسته قرار گرفته اند.

منفی غلط^{۷۰}: تعداد مواردی که جزء دسته مورد بررسی هستند ولی در این دسته قرار نگرفته اند.

منفی درست^{۷۱}: تعداد مواردی که جزء دسته مورد نظر نیستند و در آن هم قرار نگرفته اند.

دقت^{۷۲}: نسبت مواردی که به درستی در دسته مورد نظر قرار گرفته اند به تمام مواردی که در دسته مورد نظر قرار گرفته اند .

دقت متوسط: متوسط دقت ارزیابی تمام اشیاء

حساسیت^{۷۳}: نسبت مواردی که در دسته مورد نظر قرار گرفته اند به تمام مواردی که به دسته مورد نظر تعلق دارند

۵-۳-۳: نتایج ارزیابی

این پایان نامه بر مجموعه داده سان ارزیابی شده است. الگوریتم هایی که از روابط اشیاء و صحنه برای ارزیابی شیء استفاده می کنند، بر مجموعه داده سان ارزیابی می شوند. بنابراین باید نشان داده شود دقت ارزیابی شیء در این مجموعه داده توسط الگوریتم پیشنهادی بهبود می یابد.

در این پایان نامه دو دسته آزمایش انجام می شود تا عملکرد الگوریتم پیشنهادی به صورت کمی و کیفی ارزیابی شود. ارزیابی کمی براساس معیار دقت ارزیابی شیء انجام می شود که در الگوریتم های

⁶⁸ True Positive

⁶⁹ False Positive

⁷⁰ False Negative

⁷¹ True Negative

⁷² Precision

⁷³ Recall or Sensitivity

بازشناسی شیء محاسبه می شود. تمام آزمایشات توسط متلب 2012 و در ویندوز 8 بر کامپیوتر با پردازنده 2.40 گیگاهرتز و حافظه موقت 4 گیگابایت انجام شده است. در ادامه نتایج مقایسه الگوریتم ها با جزئیات شرح داده خواهند شد.

۵.۳.۳.۱ ارزیابی کمی

برای ارزیابی نتایج بازیابی شیء در تصاویر مجموعه داده سان، ارزیابی کمی الگوریتم های مختلف ضروری است. مقایسات توسط معیار میانگین دقت بازیابی که توسط معادله زیر تعریف شده است انجام می شود. دقت بازیابی هر شیء به صورت نسبت پیش بینی های شیء انجام شده که درست بوده اند، تعریف می شود.

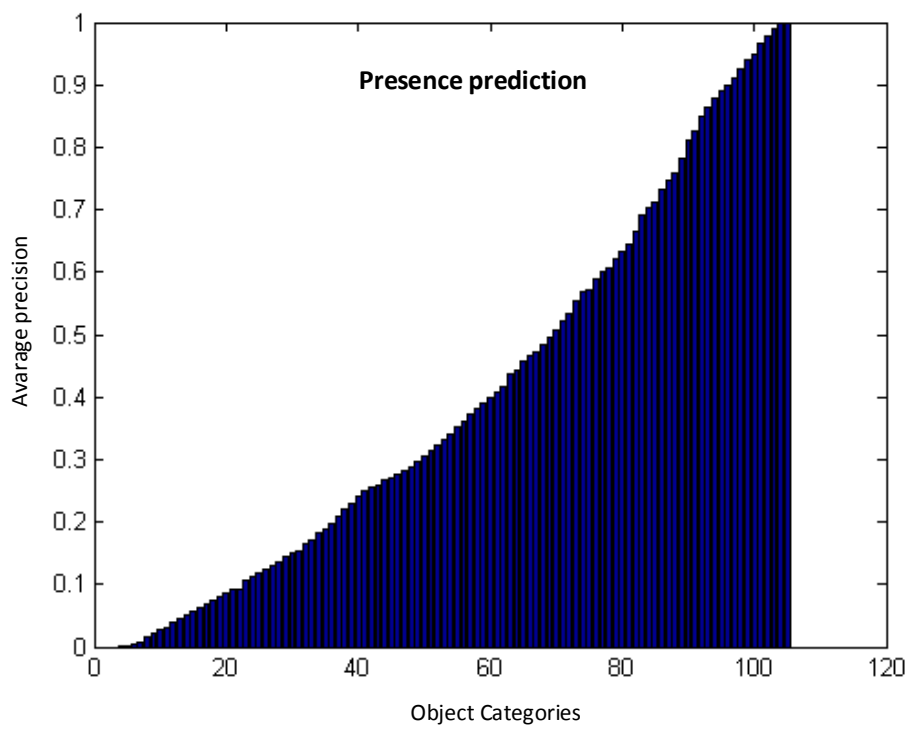
$$\text{دقت بازیابی هر شیء} = \frac{\text{پیش بینی های درست آن شیء}}{\text{کل پیش بینی های آن شیء}}$$

برای هر شیء میانگین دقت بازیابی هایی که در هر مرحله بدست آمده است به عنوان میانگین دقت^{۷۴} بازیابی آن شیء در نظر گرفته می شود. سپس برای ارزیابی الگوریتم، متوسط میانگین دقت^{۷۵} تمام اشیاء محاسبه و مقایسه می شود. در دومقاله ای که با الگوریتم پیشنهادی مقایسه شده اند 107 شیء بررسی و میانگین دقت آنها محاسبه شده است. برای اینکه بتوان نتایج را مقایسه و الگوریتم را ارزیابی کرد در این الگوریتم نیز میانگین دقت برای همان 107 شیء محاسبه شده است.

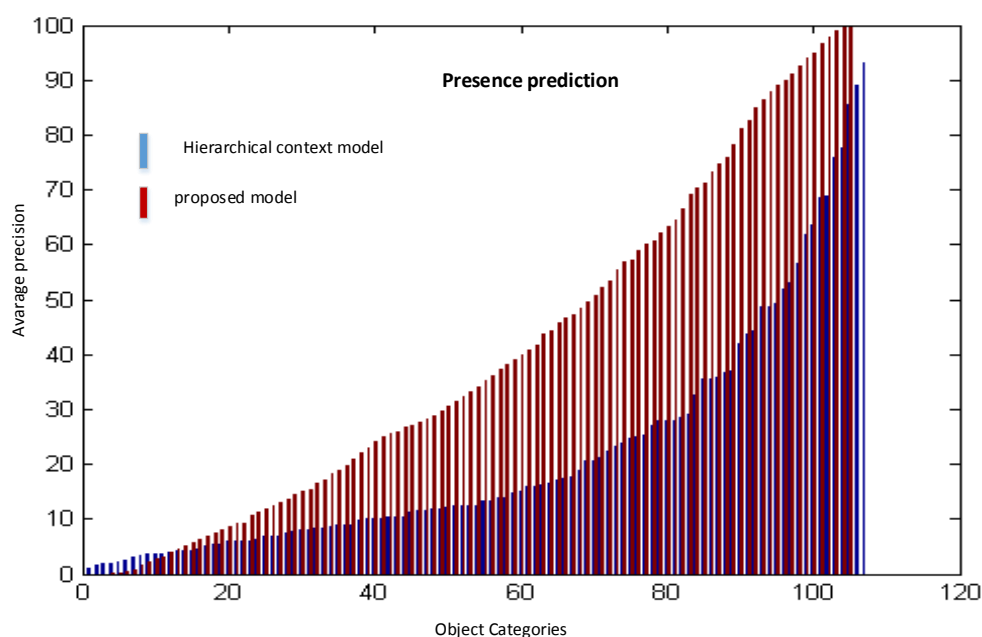
در نمودار 5-6 میانگین دقت 107 شیء برای الگوریتم پیشنهادی نشان داده شده است. در شکل 5-7 نیز میانگین دقت الگوریتم پیشنهادی و میانگین دقت الگوریتم مبتنی بر اطلاعات متنی سلسله مراتبی مقایسه شده اند. ستون های قرمز نشان دهنده دقت الگوریتم پیشنهادی برای هر شیء و ستون های آبی نشان دهنده دقت الگوریتم مبتنی بر اطلاعات متنی سلسله مراتبی است.

⁷⁴ AP(Average precision)

⁷⁵ Mean AP



شکل 6-5- دقت 107 شیء بازیابی شده در مدل پیشنهادی. اشیاء براساس دقت بازیابی مرتب شده‌اند.



شکل 7-5- دقت بازیابی 107 شیء. ستون های قرمز مربوط به الگوریتم پیشنهادی و ستون های آبی مربوط به الگوریتم مبتنی بر اطلاعات متنی سلسله مراتبی است.

در جدول زیر دقت بازیابی تعدادی از این اشیاء به عنوان نمونه در الگوریتم پیشنهادی و الگوریتم مبتنی بر اطلاعات متنی سلسله مراتبی مقایسه شده است.

جدول 1-5- دقت نمونه هایی از 107 شیء بازیابی شده توسط الگوریتم پیشنهادی

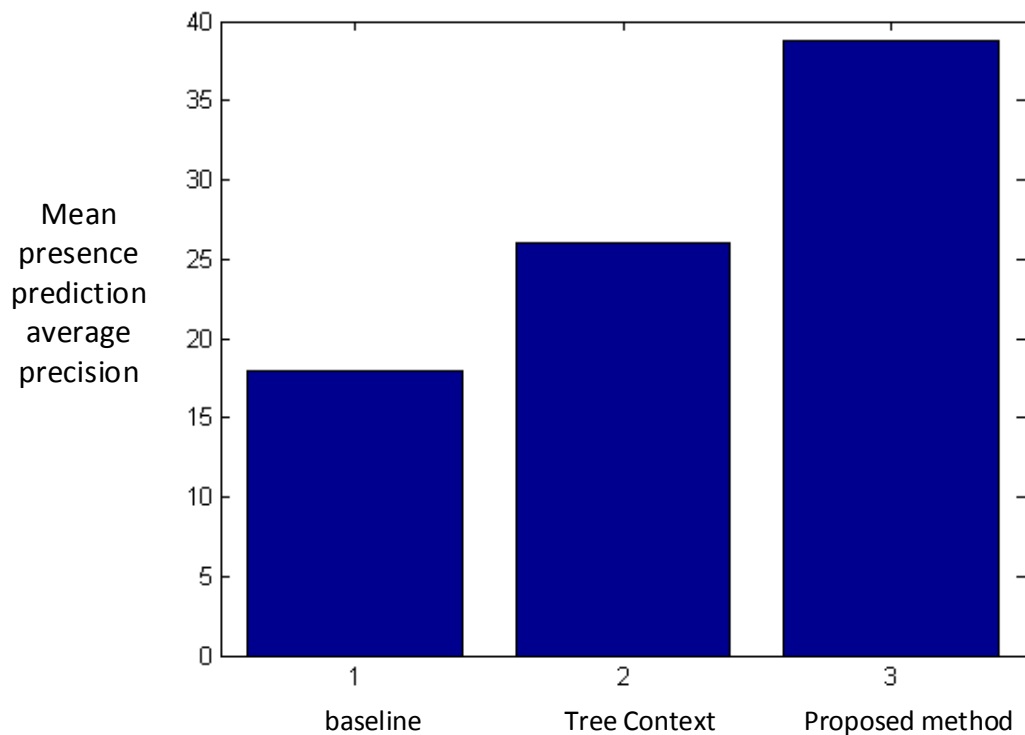
دقت روش پیشنهادی	دقت الگوریتم بازیابی شیء مبتنی بر اطلاعات متنی سلسله مراتبی [16]	شیء بازیابی شده
18.62	44.39	'airplane'
33.91	20.56	'armchair'
28.11	10.17	'awning'
30.36	6.00	'bag'
40.41	8.24	'balcony'
6.48	3.61	'ball'
40.71	8.00	'bars'

8.86	8.98	'basket'
86.98	51.84	'bed'
29.36	3.72	'bench'
39.69	27.10	'bookcase'
47.94	23.31	'books'
29.49	13.17	'bottle'
1.15	2.47	'bottles'
13.33	15.83	'bowl'
47.34	14.98	'box'
87.17	75.88	'building'
12.22	4.33	'bus'

با توجه به اینکه تعداد اشیاء موجود در تصاویر مجموعه داده سان بسیار بیشتر از مجموعه داده های دیگر است، نويز این مجموعه داده نیز از بقیه بیشتر است و نتایج بهترین الگوریتم‌ها نیز دقت پایینی را نشان می دهد. در جدول زیر دقت متوسط بازیابی سه الگوریتم [4, 10, 16] و الگوریتم پیشنهادی برای 107 دسته شیء نشان داده شده است [4].

جدول 2-5- دقت بازیابی دو الگوریتم با بهترین عملکرد بر مجموعه داده سان

	مدل های جدا کننده [10]	الگوریتم بازیابی شیء مبتنی بر اطلاعات متنی سلسله مراتبی [16]	الگوریتم پیشنهادی
دقت متوسط	17.93	26.08	38.81



شکل 8-5- مقایسه میانگین دقت الگوریتم پیشنهادی و الگوریتم های قبلی. الگوریتم 1 مدل های جداکننده [28]، الگوریتم 2 الگوریتم مبتنی بر اطلاعات متنی سلسله مراتبی [16] و الگوریتم 3 الگوریتم مبتنی با کلاس بندی مبتنی بر مجموعه [25] است.

۵.۳.۳.۲ ارزیابی کیفی

در این بخش ابتدا نشان داده می شود که در الگوریتم مبتنی بر اطلاعات سلسله مراتبی به دلیل در نظر نگرفتن تمام اشیاء موجود در آن و تعداد آنها، به دلیل تعداد زیاد شیء در تصاویر و تنوع آنها عملکرد مطلوبی ندارد. سپس نشان داده می شود که الگوریتم پیشنهادی این مساله را برطرف نموده است و با در نظر گرفتن تمام اشیاء دیگر موجود در آن و تعداد آنها توانسته است وجود شیء مورد نظر را پیش بینی کند.

مقاله [4] که در فصل قبل به تفصیل بررسی شد و جزء بهترین الگوریتم هایی است که تاکنون ارائه شده است، ارتباط بین اشیاء را به صورت سلسله مراتبی ذخیره می کند. به عنوان مثال اشیائی که معمولا در یک آشپزخانه وجود دارند، در زیردرخت گرهی sink قرار می گیرند و تمام وسایل نقلیه نیز در زیردرخت شیء جاده قرار می گیرند. مساله ای که ممکن است در این مدل به وجود آید این است که

$$P(b_{cabinet}|b_{countertop})P(b_{bowl}|b_{cabinet})P(b_{plate}|b_{bowl})P(b_{dish}|b_{plate})$$

در این حالت برای پیش‌بینی حضور شیء dish از احتمالات شرطی حضور 8 شیء دیگر کمک گرفته شده است. در مورد اشیاء دیگر تعداد این اشیاء کمتر است و برای پیش‌بینی احتمال حضور بیشتر اشیاء از احتمال شرطی حضور دو یا سه شیء دیگر استفاده شده است. این درحالی است که در صورت استفاده از اطلاعات اشیاء بیشتر می‌توان دقت بازیابی را افزایش داد. در مدل پیشنهادی برای پیش‌بینی حضور هر شیء از تمام اشیاء دیگر استفاده شده است و در نتیجه دقت به طور قابل ملاحظه‌ای افزایش یافته است.

به عنوان مثال در تصویر زیر که شامل 10 شیء نشان داده شده می‌باشد در مدل مبتنی بر اطلاعات سلسله‌مراتبی اشیاء با بیشترین احتمال (احتمال 1) grass, headstone و box تخمین زده شده است در حالیکه در این تصویر هیچ یک از این اشیاء وجود ندارد. احتمال وجود هواپیما در این تصویر 0.05 پیش‌بینی شده است با توجه به اینکه این احتمال نزدیک به صفر است، نتیجه می‌گیریم که این الگوریتم پیش‌بینی کرده است شیء هواپیما در این تصویر وجود ندارد. اما توسط الگوریتم پیشنهادی وجود هواپیما تخمین زده شده است.



لیست اشیاء موجود در تصویر

Sky
airplane
Car
unmatched
building
fence
finger
person
window
road

شکل 10-5- تصویری از مجموعه داده سان که به عنوان تصویر تست استفاده شده است و دارای 10 کلاس شیء مختلف است. احتمال وجود هواپیما در این تصویر در الگوریتم مبتنی بر اطلاعات سلسله مراتبی 0.05 است اما در الگوریتم پیشنهادی وجود آن پیش بینی شده است.

به عنوان مثالی دیگر می توان تصویر بعدی را در نظر گرفت که در الگوریتم مبتنی بر اطلاعات سلسله مراتبی احتمال وجود هواپیما 0.35 تخمین زده شده است درحالیکه وجود هواپیما در الگوریتم پیشنهادی به درستی تخمین زده شده است. همچنین در تصویر 11 احتمال وجود توپ 0.005 که عددی نزدیک به صفر است تخمین زده شده است ولی الگوریتم پیشنهادی وجود این شیء را پیش بینی کرده است.



لیست اشیاء موجود در تصویر

- Sky
- Control tower
- airplane
- building
- field
- finger
- airport cart
- road

شکل 11-5 - تصویری از مجموعه داده سان که به عنوان تصویر تست استفاده شده است و دارای 8 کلاس شیء مختلف است. احتمال وجود هواپیما در این تصویر در الگوریتم مبتنی بر اطلاعات سلسله مراتبی 0.35 است اما در الگوریتم پیشنهادی وجود آن پیش بینی شده است.



لیست اشیاء موجود در تصویر

- Ceiling
- Staircase
- Handrail
- Door
- Furniture
- Mirror
- Armchair
- Side table
- Desk lamp
- Painting
- Fireplace
- Table
- Swivel chair
- Billiard table
- Cue
- Ball
- Door
- Painting
- Ceiling lamp
- Floor lamp
- Plant
- Wall
- Floor

شکل 12-5 - تصویری از مجموعه داده سان که به عنوان تصویر تست استفاده شده است و دارای 23 کلاس شیء مختلف است. احتمال وجود توپ در این تصویر در الگوریتم مبتنی بر اطلاعات سلسله مراتبی 0.005 است اما در الگوریتم پیشنهادی وجود آن پیش بینی شده است.

۵.۳.۳.۳ ماتریس تقابل

برای بررسی عملکرد مدل پیشنهادی، ماتریس تقابل که ماتریسی با تعداد سطر و ستون مساوی با تعداد اشیاء یعنی 107×107 است، محاسبه می شود. اشیاء ستون ها اشیاء خروجی کلاسبند و اشیاء سطرها اشیاء واقعی هستند. به ازای هر پیش بینی حضور شیء i به عنوان شیء j مولفه مربوطه یعنی C_{ij} یکی اضافه می شود و تمام مولفه های ماتریس محاسبه می شود. با توجه به تعریف این ماتریس می توان دریافت که هر چه عناصر قطری ماتریس تداخل مقادیر بیشتر و عناصر غیرقطری مقادیر کمتری داشته باشد، پیش بینی ها به واقعیت نزدیکتر بوده اند و تعداد پیش بینی های درست بیشتر بوده است.

ماتریس تداخل برای 10 شیء از این 107 شیء که دارای بیشترین دقت بازیابی بوده اند در شکل بعد نشان داده شده است. همانطور که دیده می شود بیشترین مقادیر مربوط به عناصر قطری هستند و بیشتر مقادیر غیرقطری دارای مقادیر صفر هستند، در نتیجه پیش بینی ها تا حدود خوبی به درستی انجام شده است. میزان درستی^{۷۶} این ماتریس برای تمام اشیاء که با جمع عناصر قطری ماتریس بر جمع تمام عناصر ماتریس بدست می آید برابرست با 66.53 و نسبت کلاسبندی های غلط^{۷۷} آن برابر با 33.47 است.

sky	2065.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
window	0.00	972.00	1.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
floor	0.00	0.00	1547.00	0.00	2.00	0.00	1.00	0.00	0.00	0.00
road	0.00	0.00	5.00	633.00	0.00	0.00	0.00	0.00	0.00	0.00
chair	0.00	1.00	1.00	0.00	443.00	0.00	2.00	0.00	1.00	0.00
wall	0.00	0.00	2.00	0.00	1.00	1841.00	5.00	1.00	0.00	0.00
person	0.00	0.00	8.00	2.00	2.00	2.00	757.00	0.00	0.00	0.00
sea	0.00	0.00	0.00	0.00	3.00	0.00	0.00	156.00	0.00	0.00
building	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	184.00	0.00
bed	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	205.00
	sky	window	floor	road	chair	wall	person	sea	building	bed

شکل 13-5- ماتریس تداخل برای 10 شیء با بیشترین دقت بازیابی از مجموعه داده سان

⁷⁶ Accuracy

⁷⁷ Error

۵-۴: جمع بندی

در این فصل ابتدا به معرفی مفصل تر مجموعه داده استفاده شده در پایان نامه و جعبه ابزار کاربردی لیبیل می پرداخته شد و دلیل استفاده از این مجموعه داده توضیح داده شد. سپس نتایج بررسی شده بر این مجموعه داده برای اشیاء شناسایی شده ذکر شدند و همچنین این روش با روشی با بهترین عملکرد بر روی این دیتابیس مقایسه شد. در دو بخش ارزیابی کمی و کیفی به ارزیابی این الگوریتم با روشی که بهترین عملکرد را بر مجموعه داده سان داشته است، به صورت کمی و کیفی پرداخته شده است.

فصل ۶:

نتیجه‌گیری و کارهای آتی

در این پایان نامه ابتدا به معرفی یک سیستم بازشناسی شیء و سپس به معرفی روشهای مختلف بازشناسی شیء پرداخته شد. بعضی از این روشها از ویژگی های محلی و بعضی از ویژگی های سراسری تصویر استفاده می کنند. هر کدام از این روشها مزایا و معایب خود را دارند. بعضی روش ها نیز از ترکیب این دو روش استفاده می کنند.

روشهای دیگری که در این پایان نامه به آنها پرداخته شد، روش های مبتنی بر متن بودند که در این روش ها از اطلاعات صحنه و اشیاء دیگر برای شناسائی شیء موردنظر استفاده می شود. با توجه به اینکه به دلیل کوچک بودن بعضی اشیاء، وجود مانع و یا مبهم بودن یک شیء خاص در یک تصویر به دلیل شرایط مختلف نمی توان بازشناسی شیء را انجام داد، روش های مبتنی بر متن می توانند به شناسایی شیء در این موارد خاص کمک کنند. در نتیجه از نوع صحنه و اشیاء دیگر صحنه برای شناسایی شیء خاص استفاده می شود.

در قسمت ارزیابی نشان داده شده که نتایج الگوریتم پیشنهادی امیدبخش است و می توان از این روش در شرایط مختلف استفاده کرد. یکی از کارهایی که در ادامه این کار قابل انجام است، طبقه بندی اشیاء توسط روش ارائه شده است. به این صورت که توسط یک سری از اشیاء درخت تصمیم ساخته شود و سپس توسط اشیائی جدید درخت تصمیم پیمایش شود حال می توان اشیائی را که در یک گروه متنی قرار می گیرند بدست آورد و به این ترتیب می توان اشیاء را طبقه بندی کرد.

بدست آوردن اشیاء مختلف که عموماً با یکدیگر و در صحنه های یکسان ظاهر می شوند می تواند در بازشناسی شیء با استفاده از تعریف قوانین وابستگی^{۷۸} بسیار موثر باشد.

در فرآیند بازشناسی شیء در مغز اطلاعات بسیار زیادی از محیط پیرامون و روابط اشیاء دیگر حتی اشیاء ظاهراً نامرتب در یافت می شود و با سرعت بالایی پردازش می شود و در نتیجه در زمان کوتاهی فرآیند بازشناسی شیء انجام می شود. در مدل پیشنهادی تلاش شده است تا نمایشی ارائه شود که اطلاعات کلاس های مختلف اشیاء بتوانند بر نمایش شیء تاثیر گذار باشند همانطور که در مغز اطلاعات اشیاء مختلف و ظاهراً نامرتب دریافت و پردازش می شوند.

⁷⁸ Association rules

در مدل پیشنهادی پیش بینی در مورد موقعیت اشیاء انجام نمی شود و این بخش به عنوان کارهایی آینده در نظر گرفته شده است. همچنین در این مدل از روابط اشیاء برای ایجاد نمایش یک کلاس شیء استفاده شده است و در ماتریس مربوطه هر تصویر به تعداد اشیاء موجود در آن تکرار شده است، بنابراین ممکن است به نظر برسد که این نمایش دارای افزونگی می باشد، در کارهای آینده تلاش می شود با بهبود ماتریس نمایش ویژگی این افزونگی را برطرف کرد.

1. Andreopoulos, A. and J.K. Tsotsos, *50 Years of object recognition: Directions forward*. Computer Vision and Image Understanding, 2013. **117**(8): p. 827-891.
2. Grauman, K. and B. Leibe, *Visual Object Recognition*. Synthesis Lectures on Artificial Intelligence and Machine Learning, 2011. **5**(2): p. 1.
3. Galleguillos, C. and S. Belongie, *Context based object categorization: A critical survey*. Computer Vision and Image Understanding, 2010. **114**(6): p. 712-722.
4. Myung Jin, C., A. Torralba, and A.S. Willsky, *A Tree-Based Context Model for Object Recognition*. Pattern Analysis and Machine Intelligence, IEEE Transactions on, 2012. **34**(2): p. 240-252.
5. R. C. Gonzalez, R.E.W., and S. L. Eddins, *Digital image processing using MATLAB*. 2004.
6. Jain, R., R. Kasturi, and B.G. Schunck, *Machine vision*. 1995: McGraw-Hill, Inc. 549.
7. Cristianini, N. and J. Shawe-Taylor, *Support Vector Machines*. 2000: Cambridge: Cambridge University Press.
8. Viola, P. and M. Jones. *Rapid object detection using a boosted cascade of simple features*. in *Computer Vision and Pattern Recognition, 2001. CVPR 2001. Proceedings of the 2001 IEEE Computer Society Conference on*. 2001.
9. Viola, P. and M. Jones, *Robust Real-time Object Detection*. International Journal of Computer Vision - to appear, 2002.
10. Felzenszwalb, P.F., et al., *Object Detection with Discriminatively Trained Part-Based Models*. Pattern Analysis and Machine Intelligence, IEEE Transactions on, 2010. **32**(9): p. 1627-1645.
11. Dalal, N. and B. Triggs. *Histograms of oriented gradients for human detection*. in *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*. 2005.
12. Felzenszwalb, P.F., R.B. Girshick, and D. McAllester. *Cascade object detection with deformable part models*. in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. 2010.
13. Oliva, A. and A. Torralba, *The role of context in object recognition*. Trends in cognitive sciences, 2007. **11**(12): p. 520-527.
14. Hoiem, D., A. Efros, and M. Hebert. *Putting Objects in Perspective*. in *CVPR*. 2006.
15. Chow, C. and C. Liu, *Approximating discrete probability distributions with dependence trees*. Information Theory, IEEE Transactions on, 1968. **14**(3): p. 462-467.
16. Myung Jin, C., et al. *Exploiting hierarchical context on a large database of object categories*. in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. 2010.
17. Rabinovich, A., et al. *Objects in context*. in *Computer vision, 2007. ICCV 2007. IEEE 11th international conference on*. 2007. IEEE.

18. Blackwell, P.A., D, *Appearance Based Object Recognition with a Large Dataset using Decision Trees*, in *Australasian Conference on Robotics and Automation 2004*.
19. Wilking, D. and T. Röfer, *Realtime Object Recognition Using Decision Tree Learning*, in *RoboCup 2004: Robot Soccer World Cup VIII*, D. Nardi, et al., Editors. 2005, Springer Berlin Heidelberg. p. 556-563.
20. Zhiquan, Q., et al. *Robust Object Detection Based on Decision Trees and a New Cascade Architecture*. in *Computational Intelligence for Modelling Control & Automation, 2008 International Conference on*. 2008.
21. R Marée, P.G., J Piater, L Wehenkel. , *Decision Trees and Random Subwindows for Object Recognition*. , in *ICML workshop on Machine Learning Techniques for Processing Multimedia Content (2005)*.
22. Obdrzalek, S. and J. Matas. *Sub-linear Indexing for Large Scale Object Recognition*. in *BMVC 2005*. 2005.
23. *CART Classification And Regression Trees*. Available from: <http://www.salford-systems.com/products/cart#cart>.
24. Breiman, L., et al., *Classification and Regression Trees (Wadsworth Statistics/Probability)*. 1984: Chapman and Hall/CRC.
25. Cinbis, R. and S. Sclaroff, *Contextual Object Detection Using Set-Based Classification*, in *Computer Vision – ECCV 2012*, A. Fitzgibbon, et al., Editors. 2012, Springer Berlin Heidelberg. p. 43-57.
26. Safavian, S.R. and D. Landgrebe, *A survey of decision tree classifier methodology*. *Systems, Man and Cybernetics, IEEE Transactions on*, 1991. **21**(3): p. 660-674.
27. Ponce, J., et al., *Dataset Issues in Object Recognition*. 2006. p. 29-48.
28. Desai, C., D. Ramanan, and C.C. Fowlkes, *Discriminative Models for Multi-Class Object Layout*. *Int. J. Comput. Vision*, 2011. **95**(1): p. 1-12.

Abstract

In real world there is strong relation between objects and environment, so information of object relations and scene plays an important role in object detection. When a special object is searched in a scene, spectator will focus on locations with highest prior probability for object existence. Therefore scene context is very important for decision making about eye movement. Contextual influences on object recognition become evident if the local features are insufficient because the object is small, occluded or camouflaged.

In order to develop a precise method for object detection in real images with various objects, some solutions are proposed in this thesis. Object detection in this method is based on contextual information or scene and other objects of the image in order to increase detection precision. Contextual information are modeled as co-occurrence matrix and each image is shown by a vector displaying other objects of the image and their frequencies. These vectors are used to train the classifier. Because of the large number of training and test data and because indexing and matching methods are time-consuming, decision tree classifier is used. After training the classifier, test images vectors are surveyed to decide presence of each object and detection precision is computed.

An important point in contextual methods is to test them on images with high number of various objects and from different scenes, to use object relations for improvement in object detection. SUN dataset images are from more than 600 scenes and contain more than 200 object categories. The number of objects in each image is more than 10 classes, while in other datasets the maximum number of objects is 5 classes. In this method we use SUN 09 dataset and object presence prediction precision is compared with similar methods and is improved.

Index term: Object Recognition, Contextual Information, Object Detection, Scene Understanding



Kharazmi University

Faculty of engineering

M. Sc. Thesis

Artificial Intelligence

Context-based object recognition in digital images

By:

Soheila Sheikhabaei

Supervisor:

Dr. Jamshid Shanbehzadeh

Advisor:

Dr. Zeinab Ghasabi

2015, winter