



دانشکده مهندسی

پایان نامه کارشناسی ارشد در رشته مهندسی کامپیوتر (هوش مصنوعی)
یادگیری ویژگی های مکانی و زمانی برای تشخیص فعالیت انسان

توسط:

نجمه هادی برحق طلب

استاد راهنما:

دکتر زهره عظیمی فر

شهریور ۱۳۹۲

به نام خداوند دوست

که هر چه دارم از اوست

به نام خدا

اظهارنامه

اینجانب (.....) دانشجوی رشته ی

..... گرایش

دانشکده ی اظهار می کنم که این پایان نامه حاصل پژوهش خودم بوده و در جاهایی که از منابع دیگران استفاده کرده ام، نشانی دقیق و مشخصات کامل آن را نوشته ام. همچنین اظهار می کنم که تحقیق و موضوع پایان نامه ام تکراری نیست و تعهد می نمایم که بدون مجوز دانشگاه دستاوردهای آن را منتشر ننموده و یا در اختیار غیر قرار ندهم. کلیه حقوق این اثر مطابق با آیین نامه مالکیت فکری و معنوی متعلق به دانشگاه شیراز است.

نام و نام خانوادگی:

تاریخ و امضا:

به نام خدا

یادگیری ویژگی های مکانی و زمانی برای تشخیص فعالیت انسان

به کوشش

نجمه هادی برحق طلب

پایان نامه

ارائه شده به تحصیلات تکمیلی دانشگاه شیراز به عنوان بخشی

از فعالیت های تحصیلی لازم برای اخذ درجه کارشناسی ارشد

در رشته ی:

مهندسی کامپیوتر - هوش مصنوعی

از دانشگاه شیراز

شیراز

جمهوری اسلامی ایران

ارزیابی کمیته پایان نامه، با درجه ی: عالی

دکتر زهره عظیمی فر، استادیار بخش مهندسی و علوم کامپیوتر (استاد راهنما)

دکتر رضا بوستانی، استادیار بخش مهندسی و علوم کامپیوتر (استاد مشاور)

دکتر ستار هاشمی، استادیار بخش مهندسی و علوم کامپیوتر (داور متخصص داخلی)

شهریور ماه ۱۳۹۲

تقدیریم بہ

آمان کہ دوستمان دارم.

سپاسگزاری

سپاس خدای را که سخنوران، در ستودن او بمانند و شمارندگان، شمردن نعمت های او ندانند و کوشندگان، حق او را گزاردن نتوانند. و سلام و دورد بر محمد و خاندان پاک او، طاهران معصوم، هم آنان که وجودمان وامدار وجودشان است؛ و نفرین پیوسته بر دشمنان ایشان تا روز رستاخیز...

بدون شک جایگاه و منزلت معلم، اجل از آن است که در مقام قدردانی از زحمات بی شائبه- ی او، با زبان قاصر و دست ناتوان، چیزی بنگاریم. اما از آنجایی که تجلیل از معلم، سپاس از انسانی است که هدف و غایت آفرینش را تامین می کند و سلامت امانت هایی را که به دستش سپرده اند، تضمین؛ بر حسب وظیفه و از باب "من لم یشکر المنعم من المخلوقین لم یشکر الله عزّ و جلّ": از پدر و مادر عزیزم... این دو معلم بزرگوار که همواره بر کوتاهی و درستی من، قلم عفو کشیده و کریمانه از کنار غفلت هایم گذشته اند و در تمام عرصه های زندگی یار و یآوری بی چشم داشت برای من بوده اند؛ از استاد با کمالات و شایسته؛ سرکار خانم دکتر زهره عظیمی فر که در کمال سعه صدر، با حسن خلق و فروتنی، از هیچ کمکی در این عرصه بر من دریغ ننمودند و زحمت راهنمایی این رساله را بر عهده گرفتند؛ و از استاد فرزانه و دلسوز؛ جناب آقای دکتر بوستانی که زحمت مشاوره این رساله را متقبل شدند؛ از استاد صبور و با تقوا، جناب آقای دکتر هاشمی، مدیریت محترم کرسی گروه، که زحمت داوری این رساله را متقبل شدند؛ کمال تشکر و قدردانی را دارم. باشد که این خردترین، بخشی از زحمات آنان را سپاس گوید.

در پایان به پاس قدر دانی از قلبی آکنده از عشق و معرفت که محیطی سرشار از سلامت و امنیت و آرامش و آسایش برای من فراهم آورده است؛ همدلی که با واژه‌ی نجیب و مغرور تلاش آشنایی دارد و تلاش راستین را می شناسد و عطر رؤیایی آن را استشمام می کند و مرا

در راه رسیدن به اهداف عالی یاری می رساند؛ همو که حس تعهد و مسئولیت را در زندگی مان
تلاویی خدایی داده است؛ این پایان نامه تقدیم همسر مهربانم می گردد.

چکیده

یادگیری ویژگی های مکانی و زمانی برای تشخیص فعالیت انسان

به کوشش

نجمه هادی برحق طلب

امروزه با افزایش رو به رشد دوربین های دیجیتال، روزانه حجم عظیمی از داده های ویدئویی تولید می شوند که بر روی وب و پایگاه های داده ی بزرگ قابل دسترس هستند. دسته بندی این داده ها بر اساس محتوایشان به صورت اتوماتیک و با دقت بالا هدف غایی بسیاری از کاربردها در این زمینه است. در این پایان نامه ما سیستمی طراحی کرده ایم که بتواند داده های ویدئویی حاوی عمل انسان را با دقت خوب و در زمان نسبتا کوتاهی برچسب گذاری کرده و دسته بندی کند. در این پایان نامه از روش های دیداری بهره گرفته شده است. روش های دیداری برای تشخیص فعالیت انسان در واقع فرایند تشخیص انسان در ویدئو و عملی که انجام می دهد با بهره گیری از تکنیک های بینایی ماشین می باشد.

در این راستا سیستم طراحی شده از روش یادگیری بدون نظارت برای استخراج ویژگی و کد کردن آن ها استفاده کرده است. روش های یادگیری ویژگی نه تنها به راحتی قابل تعمیم به حوزه های دیگر هستند بلکه کارایی قابل توجهی بر روی داده های ویدئویی جمع آوری شده از دنیای واقعی دارند.

در پایان برای ارزیابی این سیستم به دسته بندی داده های ویدئویی پایگاه های داده ی
UCF، KTH و YouTube پرداخته ایم و به نتایج قابل توجهی دست پیدا کرده ایم.

کلمات کلیدی: تشخیص عمل، ویژگی محلی، تحلیل زیر فضای مستقل، کدگذاری خطی
محلی، ادغام سه بعدی.

فهرست مطالب

عنوان	صفحه
۱- مقدمه.....	۲
۱-۱- مقدمه.....	۲
۱-۲- انگیزه.....	۳
۱-۳- شرح مسئله.....	۶
۱-۴- چالش ها.....	۸
۱-۵- انواع داده های ویدئویی.....	۹
۱-۵-۱- داده های ویدئویی کنترل شده.....	۹
۱-۵-۲- داده های ویدئویی قید دار.....	۱۰
۱-۵-۳- داده های ویدئویی کنترل نشده.....	۱۱
۲- پیشینه ی تحقیق.....	۱۳
۲-۱- مقدمه.....	۱۳
۲-۲- نمایش جهانی.....	۱۴
۲-۳- ویژگی های محلی.....	۱۸
۲-۳-۱- تشخیص دهنده ی ویژگی.....	۱۸
۲-۳-۲- توصیف گرهای ویژگی.....	۲۰
۲-۴- کیف ویژگی ها.....	۲۲
۲-۵- مدل های مکانی-زمانی عمل.....	۲۴

- ۳- مبانی نظری تحقیق..... ۲۷
- ۳-۱- یادگیری ویژگی..... ۲۷
- ۳-۲- یادگیری ویژگی بدون نظارت..... ۲۹
- ۳-۲-۱- تحلیل زیرفضای مستقل..... ۳۰
- ۳-۲-۲- نسخه ی پشته ای و کانولوشن الگوریتم تحلیل زیر فضای مستقل..... ۳۲
- ۴- روش تحقیق..... ۳۵
- ۴-۱- کدگذاری..... ۳۵
- ۴-۱-۱- کد کردن توصیف گر با کوانتیزه کردن بردار..... ۳۶
- ۴-۱-۲- کد کردن توصیف گر با کد گذاری خلوت تطبیق هرمی فضایی..... ۳۶
- ۴-۱-۳- کد کردن خطی به صورت محلی..... ۳۷
- ۴-۲- کد کردن خطی محلی سه بعدی..... ۴۱
- ۴-۲-۱- چگونگی ایجاد مجموعه بردارهای پایه..... ۴۳
- ۴-۲-۲- ادغام سه بعدی..... ۴۴
- ۵- نتایج..... ۴۷
- ۵-۱- پایگاه های داده..... ۴۷
- ۵-۱-۱- پایگاه داده ی Weizman..... ۴۸
- ۵-۱-۲- پایگاه داده ی KTH..... ۴۹
- ۵-۱-۳- پایگاه داده UCF Sports..... ۵۰
- ۵-۱-۴- پایگاه داده ی You Tube..... ۵۱
- ۵-۱-۵- پایگاه داده ی Hollywood..... ۵۲
- ۵-۲- پارامترها و تنظیمات سیستم در پیاده سازی..... ۵۳
- ۵-۳- آزمایش ها و نتایج..... ۵۵
- ۶- نتیجه گیری و کارهای آینده..... ۶۳
- اختصارات..... ۶۴

واژه نامه فارسی به انگلیسی.....۶۵

واژه نامه انگلیسی به فارسی.....۶۸

فهرست منابع.....۷۱

فهرست جدول‌ها

- جدول ۱-۵ دقت سیستم برای کلاسه بندی پایگاه داده ی **KTH** با تغییر پارامتر تعداد بردارهای پایه ۵۷
- جدول ۲-۵ دقت سیستم برای کلاسه بندی پایگاه داده ی **KTH** با تغییر پارامتر تعداد بردارهای پایه ۵۸
- جدول ۳-۵ مقایسه سیستم ارائه شده در کلاسه بندی پایگاه داده ی **KYH** با سایر روشها..... ۵۸
- جدول ۴-۵ دقت سیستم برای کلاسه بندی پایگاه داده ی **YouTube** ۵۹
- جدول ۵-۵ دقت سیستم برای کلاسه بندی پایگاه داده ی **UCF** . برای ایجاد بردارهای پایه عمل **K-means** ۶۰
- جدول ۶-۵ مقایسه ی دقت سیستم با دقت سایر روش ها پیشنهادی بر روی پایگاه داده ی **YouTube** ۶۰
- جدول ۷-۵ مقایسه ی دقت سیستم با دقت سایر روش ها پیشنهادی بر روی پایگاه داده ی **UCF** ۶۰

فهرست شکل‌ها

- شکل ۱-۱ ضبط و گرفتن حرکات برای تولید فیلم انیمیشنی در استدیو ۴
- شکل ۲-۱ نمونه چند عمل در فیلم ها ۵
- شکل ۳-۱ کاربردهای تشخیص عمل انسان ۵
- شکل ۴-۱ تشخیص عمل با استفاده از چندین دور بین با زاویه های متفاوت ۱۰
- شکل ۵-۱ تحلیل شکل ماسک های بدست آمده از حذف پیش زمینبرای سیستم های نظارتی ۱۱
- شکل ۱-۲ نمونه هایی از تصویر MEI و تصویر MHI برای دو عمل مختلف [19] ۱۵
- شکل ۲-۲ شکل های زمانی-مکانی برای عمل دویدن، راه رفتن و پریدن از راست به چپ [20] ۱۶
- شکل ۳-۲ نمایش عمل با استفاده از جریان نوری [21] ۱۷
- شکل ۴-۲ نمایش شکل stick برای عمل نشستن و پریدن. [23] ۱۸
- شکل ۵-۲ نقاط کلیدی و مهم زمانی-مکانی از حرکت پا هنگام راه رفتن شخص. [24] ۱۹
- شکل ۶-۲ نقاط کلیدی زمانی-مکانی وقتی که از هسین به عنوان معیار saliency استفاده می شود. ۲۰
- شکل ۷-۲ توصیف گر سه بعدی SIFT [28] ۲۱
- شکل ۸-۲ یک مدل چهار بخشی برای عمل دست تکان دادن. [36] ۲۳
- شکل ۹-۲ بخش بندی مکانی-زمانی ۲۵
- شکل ۱-۳ استخراج ویژگی از تصویر. [38] ۲۹
- شکل ۲-۳ معماری تحلیل زیر فضای مستقل ۳۱
- شکل ۳-۳ معماری نسخه ی پشته ای و کانولوشن تحلیل زیر فضای مستقل ۳۱
- شکل ۱-۴ چهارچوب سیستم طراحی شده ۴۳
- شکل ۱-۵ . نمونه فریم از کلاس های پایگاه داده ی Weizeman ۴۸
- شکل ۲-۵ نمونه فریم کلاس های پایگاه داده ی KTH در چند سناریوی مختلف ۴۹
- شکل ۳-۵ . نمونه فریم کلاس های پایگاه داده ی UCF Sports ۵۰
- شکل ۴-۵ نمونه فریم کلاس های پایگاه داده ی You Tube ۵۱
- شکل ۵-۵ نمونه فریم کلاس های پایگاه داده ی Hollywood ۵۲

فصل اول

مقدمه

۱- مقدمه

در این فصل به شرح کلیاتی درباره ی انگیزه ی انتخاب موضوع، تشخیص عمل انسان و همچنین شرحی بر مسئله و کاربردها و چالش های آن ارائه می شود. در انتهای فصل نیز توضیحاتی درباره ی نوع داده های در دسترس جهت تست سیستم های تشخیص عمل انسان داده می شود.

۱-۱- مقدمه

در طول چندین دهه ی گذشته، کامپیوترها و شبکه های جهانی زندگی ما را به شدت تحت تأثیر قرار داده اند. کامپیوترها کارهای تکراری و محاسباتی وسیعی انجام می دهند و امکان ارتباطات را گسترش می دهند. همراه با پیشرفت عمومی تکنولوژی در زمینه ی کامپیوتر، داده های ویدئویی بیش از پیش قابل دسترس بوده و نقش مهم و فزاینده ای در زندگی روزمره ی ما ایفا می کنند. امروزه حتی سخت افزارهای مصرفی متداول مانند نوت بوک ها، تلفن های همراه و دوربین های دیجیتال داده های ویدئویی تولید می کنند. به طور همزمان دسترسی سریع تر به اینترنت و افزایش ظرفیت ذخیره سازی به ما این امکان را می دهد که داده های ویدئویی را منتشر کرده و با دیگران به اشتراک بگذاریم .

به طور مثال ۳۶ میلیون کاربر اینترنت در آلمان (۴۴٪ از جمعیت آلمان) بیش از ۶ میلیارد ویدئو را در آگوست ۲۰۰۹ به صورت آنلاین تماشا کرده اند و این آمار در مقایسه با آگوست ۲۰۰۸ رشد ۳۸ درصدی داشته است. میزان ویدئویی که در هر دقیقه روی سایت You Tube آپلود می شود در ماه می ۲۰۰۸ نسبت به اواسط سال ۲۰۰۷ از ۶ ساعت به ۲۰ ساعت افزایش یافته است و این به معنای افزایش ۳۳۰ درصدی در طول ۲ سال می باشد.

به هر حال بر خلاف افزایش اهمیت داده های ویدئویی امکان آنالیز آن ها به طور اتوماتیک بسیار محدود است. سیستم های کامپیوتری دیداری از توانایی های دیداری انسان خیلی عقب تر هستند. برای نمونه جستجوی یک ویدئو در میان حجم عظیم از داده های ویدئویی در حال حاضر فقط با برچسب گذاری پر هزینه ی دستی امکان پذیر می باشد. موتورهای جستجوی وب عموماً برای بازیابی ویدئوهای مربوطه وابسته به داده های متنی مانند توصیف ها یا برچسب ها هستند .

مثال دیگر کاربرد های نظارتی می باشد. تا به امروز شهر لندن حدود یک میلیون دوربین های مدار بسته CCTV camera (closed circuit television) با هزینه ای بالغ بر ۲۰۰ میلیون پوند نصب کرده است. با این وجود آن ها فقط توانستند به ازای هر ۱۰۰۰ دوربین نظارتی به کشف یک جرم کمک کنند. بر اساس گزارشات داخلی دوربین های نظارتی CCTV منجر به هزینه های گزاف و حداقل بازدهی شده اند .

محققان بخش دولتی home office به این نتیجه رسیدند که CCTV تقریباً نتوانست به کاهش جرم کمکی کند و به نظر می رسد بیشترین تاثیر را در جلوگیری از تخلف وسایل نقلیه هنگام پارک کردن داشته است. در واقع با حجم عظیم داده های ویدئویی، گلوگاه اصلی ضرورت ذخیره و تحلیل دستی این داده ها می باشد .

یک حیطة ی کاربردی مهم که آنالیز ویدئو در آن به عنوان یک رابط پیچیده ی انسان و کامپیوتر بسیار مورد توجه می باشد بازی های کامپیوتری هستند. یک پروژه ی در حال انجام پروژه ی Natal شرکت میکروسافت می باشد. چهارچوب این پروژه قادر به گرفتن حرکات سه بعدی کل بدن، تشخیص چهره، تشخیص صدا و مکان یابی منبع تولید صدا می باشد. برای این کار باید اطلاعات بدست آمده از چندین سنسور با هم ترکیب شود. سنسورها شامل دوربین های ویدئویی، سنسور تشخیص عمق بر اساس الگوهای اشعه فرسرخ و یک میکروفون می باشد. بدین ترتیب کاربر می تواند بازی های ویدئویی را بدون هیچ ابزار کنترلی و با تعاملات دنیای دیداری به کمک کل اعضای بدنش به طور طبیعی انجام دهد. ضبط و گرفتن حرکات بازیگر برای شخصیت های کارتونی در فیلم های انیمیشینی (شکل ۱-۱) و همچنین جلوه های



شکل ۱-۱ ضبط و گرفتن حرکات برای تولید فیلم انیمیشنی در استدیو

ویژه ی سینمایی به یک استاندارد بالفعل تبدیل شده است. با این حال، تحلیل حرکت انسان هم می تواند نقش مهمی در کاربردهای رسانه ای ایفا کند. همچنین در تحلیل و بهبود حرکات ورزشی ورزشکاران نقش تأثیر گذاری داشته باشد.

این مثال ها نشان می دهند امروزه نیاز گسترده ای به سیستم های کامپیوتری دیداری برای درک و پردازش داده های ویدئویی به طور اتوماتیک وجود دارد. آن ها همچنین نشان می دهند که تکنولوژی بینایی ماشین توانایی بالا و بالقوه ای بر روی آینده ی ما دارد. طراحی سیستم های هوشمندی که این پردازش ها را به صورت اتوماتیک انجام دهد کار بس ارزشمندی است که می تواند در کاربرد های صنعتی و تجاری مورد بهره برداری قرار گیرد.

از جمله دیگر کاربردها که نیاز به تحلیل و پردازش داده های ویدئویی دارند تشخیص عمل^۱ می باشد. شکل ۱-۲ نمونه ی تعدادی عمل را که از چندین فیلم استخراج شده را نشان می دهد. تشخیص عمل انسان می تواند در سیستم های نظارتی^۲ مثل اتاق های کنترل هوشمند، تجزیه و تحلیل فیلم مسابقات ورزشی و کلیه ی سیستم های هوشمند مبتنی بر تعامل انسان و ابزار های الکترونیکی مورد استفاده قرار گیرد. به طور مثال سیستم های نظارتی در اماکن عمومی مانند فرودگاه ها و ایستگاه های مترو می توانند اعمال غیرعادی و مشکوک از قبیل قرار دادن کیف در سطل زباله توسط یک شخص را به طور اتوماتیک تشخیص داده و به

¹ Action recognition

² Surveillance system



شکل ۲-۱ نمونه چند عمل در فیلم ها



(ب)



(الف)



(ج)

شکل ۳-۱ کاربردهای تشخیص عمل انسان: (الف) سیستم های نظارتی (ب) دستگاه کینکت (ج) حاشیه نویسی مسابقات ورزشی

قسمت امنیتی هشدار دهند. برای نمونه، کینکت^۳ یک محصول تجاری است که از قدرت شناسایی عمل انسان بهره می برد. یک کنسول کنترلگر بازی از مایکروسافت که معمولاً حرکات بدن انسان را دنبال می کند، سیستم حاشیه نویسی ورزشی می تواند تحلیل حرکت پیچیده بازیکن را انجام دهد و اطلاعات استراتژی بازی را از ویدیو پخش زنده بازی ورزشی به صورت بلادرنگ استخراج کند. شکل ۳-۱ برنامه های کاربردی مختلفی را معرفی می کند. شکل ۳-۱ (الف) یک مثال از سیستم نظارت تجاری را نشان می دهد. شکل ۳-۱ (ب) کنسول کنترلگر بازی براساس حرکت انسان را نشان می دهد که کینکت از شرکت مایکروسافت نامیده می شود. شکل ۳-۱ (ج) یک سیستم حاشیه نویسی ورزشی را نشان می دهد که برای هوشمندسازی یادداشت نویسی بازی های ورزشی استفاده می شود.

۳-۱- شرح مسئله

این پایان نامه بر روی مسئله ی تشخیص فعالیت انسان در ویدئو های واقعی مانند ویدئوهای تهیه شده از فیلم ها، اینترنت، دوربین های نظارتی و غیره متمرکز شده است. برای آن که موضوع را دقیق تر بررسی کنیم لازم است تعریف واژه های "عمل پایه"، "عمل" و "فعالیت" به روشنی بیان شود. زبان انسان از جمله ها تشکیل شده است که جمله به نوبه ی خود شامل فاعل، فعل و مفعول می باشد. برای توصیف مفهوم دیداری ویدئو وجود ساختاری شبیه به آنچه برای زبان انسان بیان شد ضروری می باشد. از دیدگاه الگوریتمی می توان تشخیص عمل انسان را چنین بیان کرد: تشخیص (۱) فاعل (کننده ی کار) که معمولاً انسان ها هستند؛ (۲) مفعول که می تواند انسان های دیگر باشد یا محیطی که فاعل کار را در آن انجام می دهد؛ (۳) فعل عملی را که فاعل انجام می دهد توصیف می کند، همچنین می تواند توصیف تعامل بین مفعول و فاعل باشد. به این معنا عمل می تواند در بازه ی کوتاه زمانی انجام شود. همچنین می شود به رخدادی که در بازه ی زمانی طولانی تر انجام می شود اطلاق شود.

³ Kinect