

تخمین حالت سه‌بعدی بدن انسان از یک تصویر بوسیله شبکه عصبی کانولوشن و کدگذاری و بازنمایی تنک با رویکرد مبتنی بر مدل

حسن علی‌کرمی^۱، فرزین یغمایی^۲ و محمدجواد فدایی اسلام^۳

چکیده

در زمینه تخمین و ساخت اسکلت حالت سه‌بعدی بدن انسان از طریق بندهای بدن (body joints) بوسیله یک تصویر دوبعدی، چالش‌های عمق و خودانسدادی وجود دارد که مانع از تخمین دقیق می‌گردد. در این مقاله به تخمین حالت سه‌بعدی بدن انسان با دو رویکرد مختلف پرداخته شده است. بدین منظور، رویکرد اول پیشنهادی با تمرکز بر عمق حالت دوبعدی حقیقت اصلی بوسیله کدگذاری و بازنمایی تنک و تصحیح گر مبتنی بر مدل، حالت سه‌بعدی بدن انسان استخراج می‌شود. در رویکرد دوم پیشنهادی به کمک روش مبتنی بر یادگیری شبکه‌های عصبی کانولوشن، تخمین حالت دوبعدی بدن انسان بدست می‌آید، سپس بوسیله کدگذاری و بازنمایی تنک و تصحیح گر مبتنی بر مدل، تخمین عمق حالت استخراج می‌شود. نتایج حاصل از این روش، برتری تخمین حالت و عمق سه‌بعدی بدن انسان را نسبت به رویکردهای پیشین نشان می‌دهد. تخمین حالت‌های سه‌بعدی انجام شده در روش پیشنهادی نشان می‌دهد میانگین خطای بازسازی نسبت به کارهای مشابه کاهش قابل توجهی داشته است.

کلیدواژه‌ها

شبکه‌های عصبی کانولوشنی، کدگذاری و بازنمایی تنک، اسکلت حالت سه‌بعدی بدن انسان، تخمین حالت سه‌بعدی

۱ مقدمه

این رو پژوهشگران حوزه بینایی ماشین برای برقراری ارتباط یا عمل متقابل نسبت به دنیای واقعی، نیازمند تشخیص و تحلیل اجسام در تصویر و ویدئو هستند. این تشخیص و تحلیل با انگیزه‌های متفاوتی صورت گرفته است. برای مثال، هدف از ردیابی حرکت یک جسم، محدود کردن صحنه، فشرده‌سازی و ویرایش تصاویر است که متفاوت از تشخیص نوع حرکت جسم در صحنه هستند.

یکی از اهداف اصلی در تشخیص و تحلیل اجسام، درک صحنه با تحلیل موقعیت و جهت اجسام صلب^۱ است. در واقع پژوهشگران سعی دارند با تعریف مختصات یک جسم، قابلیت درک و توصیف آن را مقدور سازند تا از این طریق، موقعیت و جهت یک جسم در فضای دوبعدی و سه‌بعدی نسبت به مختصات جهانی را تعیین کنند. این مبحث در بینایی ماشین تحت عنوان

امروزه با پیشرفت علوم مختلف از جمله هوش مصنوعی، تعامل بین انسان و ماشین امکان‌پذیر شده است. از مهم‌ترین چالش‌های موجود، ایجاد تعامل بصری بین انسان و ماشین است که نیازمند درک دقیق صحنه با تشخیص و تحلیل رفتار اجسام می‌باشد. از

این مقاله در شهریورماه ۱۳۹۶ دریافت، در مردادماه ۱۳۹۷ دومین بازنگری و در آذرماه پذیرفته شد.

^۱ دانشجوی کارشناسی ارشد مهندسی کامپیوتر، دانشگاه سمنان

رایانامه: Hassan_Alikarami@semnan.ac.ir

^۲ دانشگاه سمنان، دانشکده مهندسی برق و کامپیوتر

رایانامه: F_Yaghmaee@semnan.ac.ir

^۳ دانشگاه سمنان، دانشکده مهندسی برق و کامپیوتر

رایانامه: Fadaei@semnan.ac.ir

نویسنده مسئول: فرزین یغمایی

www.SID.ir

^۱ Rigid

Archive of SID

اول به عنوان یک بخش متمایزکننده برای طبقه‌بندی قطعه‌های بدن است و لایه دوم با استفاده از تخمین لایه اول، قادر به پیش‌بینی مفصل‌ها توسط مدل‌سازی از طریق وابستگی‌های مفصل می‌باشد. این روش نشان داد عملکرد محلی‌سازی مفصل بهتر از ساختار درختی است.

در سال‌های اخیر، متخصصین به تخمین حالت دو بعدی بدن انسان با استفاده از شبکه‌های عصبی کانولوشن پرداخته‌اند. این رویکرد برای اولین بار، سال ۲۰۱۳ در مقاله [۴] ارائه شد که شبکه آن مشابه Alex Net [۵] است با این تفاوت که در لایه آخر از یک لایه رگرسیون استفاده شده است. دقت بالای حاصل از این رویکرد موجب جلب توجه پژوهشگران این حوزه به شبکه‌های عصبی کانولوشن گردید. در مقاله [۶] ایده‌ای مشابه [۴] به صورت یک روش بازخوردی بالا به پایین ارائه شد. در این روش با تکرار شبکه عصبی کانولوشن خود اصلاح، خطای تخمین حالت دوبعدی انسان کاهش یافت.

در مقاله [۷] سعی شده است اطلاعات ناشی از دیدگاه کلی و بخش‌های محلی نسبت به یک تصویر، بدست آمده از دو منبع شبکه عصبی کانولوشن، یکپارچه‌سازی شود. اما مقاله [۸] روش چند مرحله‌ای از شبکه‌های عصبی کانولوشن برای تخمین حالت بدن معرفی کرد که در هر مرحله قسمتی از حالت بدن تخمین زده می‌شود. بنابراین، در تشخیص بندهای بدن انسان با خطای کمتری مواجه می‌شوند. این شبکه به دلیل بروز ریسک خطا در مواجهه با شبک کاهشی یا محوکنندگی در طول آموزش، از یک لایه میانی بعد از هر مرحله برای اجرا استفاده شده است.

در مقاله [۹] روشی پایان به پایان^۴ برای تخمین حالت انسان پیشنهاد شده است که به صورت متوالی، با استخراج روابط میان ویژگی‌های نقشه مفصل توسط شبکه عصبی کانولوشن و ترکیب مدل‌سازی ساختاری و یادگیری ویژگی بر مبنای نقشه مفصل و اصلاح آن از طریق گراف، سعی در انتخاب بهترین حالت دارد.

در موارد بالا، هدف تخمین حالت دوبعدی بدن است، اما بدن انسان یک سیستم بسیار پیچیده متشکل از اندام‌ها و مفصل مختلف است، لذا تخمین واقعی از موقعیت مفصل در حالت سه بعدی، حتی برای مغز انسان هم چالش برانگیز است. از این رو، مقاله [۱۰] سعی در حل مسائل تخمین حالت بدن انسان و بخش‌بندی آن با استفاده از دسته‌بندی فضای فرضیه داشته است. آن‌ها با تهیه یک مجموعه‌داده‌گان جدید حاوی تصاویر دارای حاشیه‌نویسی و اطلاعات مفصل سه بعدی بدن در یک تصویر دوبعدی، سعی در پیاده‌سازی یک روش جستجوی داده‌محور به منظور یافتن فضای حالت داشته‌اند. از این رو به صورت جدا، قسمت‌های بدن را در کلاس‌هایی دسته‌بندی کرده و از طریق ماشین‌بردار پشتیبان، به تشخیص محتوا پرداخته‌اند.

"تخمین حالت"^۱ شناخته شده است. از آنجاکه بدن انسان یک جسم صلب محسوب نمی‌شود، این عمل به صورت نمایش حالت بدن از طریق بندهای آن انجام می‌شود. از این رو می‌توان فرض کرد هر کدام از بندهای بدن انسان یک جسم صلب است و هدف از تشخیص حالت بدن انسان فقط حالت اسکلت است، لذا بدلیل وابستگی شکل بدن به عضلات و ماهیچه‌ها، نیازی به تشخیص شکل بدن نمی‌باشد. این مسئله در دهه اخیر به دلیل کاربردهای گسترده در صنعت فیلم و بازی سازی، سیستم‌های نظارتی، تجهیزات پزشکی و ... در بینایی ماشین اهمیت بسیاری یافته است.

روش‌های کاربردی موجود، مبتنی بر آنالیز حرکتی اجسام توسط سنسورهای عمقیاب یا تصاویر چنددیدگی هستند که به دلیل نیاز به فضای مخصوص و هزینه بالا، لازم است روش‌هایی جهت تشخیص حالت دوبعدی و سه بعدی با استفاده از یک تصویر ارائه گردند. در این رویکرد، هدف تخمین حالت دو بعدی و محاسبه عمق حالت بدن انسان بر پایه استفاده از یک تصویر با کمترین خطا نسبت به حالت سه بعدی حقیقت اصلی^۲ است.

جهت نزدیک شدن به مکانیسم مغز انسان از ایده یادگیری عمق از یک تصویر استفاده می‌شود زیرا، گرچه تحلیل تصاویر در مغز انسان بر اساس دو تصویر از یک مشاهده صورت می‌پذیرد، اما به دلیل یادگیری و تجربیات مشاهدات قبلی، انسان قادر به تشخیص عمق اجسام از طریق مشاهدات تک دیدگی نیز می‌باشد. لذا از چالش‌های مهم در این حوزه، تخمین حالت دوبعدی جهت غلبه بر مشکل خود انسدادی^۳ است. برای مثال در مقاله [۱] یک چارچوب محاسباتی مبتنی بر مدل جهت بخش‌بندی و تشخیص اشیا ارائه شده است. ایده اصلی این مقاله، نمایش اشیا توسط مجموعه‌ای از قطعات مرتب شده در یک پیکربندی انعطاف‌پذیر است. آن‌ها یک چارچوب آماری برای توصیف ظاهری اشیا در نظر گرفته و از مدل‌های ساختاری تعریف شده جهت تطابق با حالت فرد یا شی استفاده نمودند. بهترین تطابق برای مدل تصویر به گونه‌ای در نظر گرفته شده است که تابع انرژی را حداقل نماید. مقاله [۲] با یک نمایش جدید از توصیف بخش‌های مدل، بندها را به صورت بلوک‌های مخلوطی ترکیبی در نظر گرفتند. در این روش از مدل‌های مخلوط انعطاف‌پذیر برای روابط بین بخش‌ها و افزایش مدل‌هایی برای توصیف روابط فضایی استفاده کردند. بطور کلی، این رویکردها به دلیل دقت پایین و محاسبات زیاد، کمتر مورد توجه قرار گرفته است.

از سوی دیگر، مقاله [۳] برای کاهش بار محاسباتی و دقت تشخیص، یک چارچوب ساختاری تصویر برای حل مشکل تخمین کارآمد از بخش‌های مختلف بدن، از رگرسیون غیرخطی استفاده کرده است. آن‌ها دو لایه رگرسیون تشکیل دادند که لایه

¹ Pose Estimation

² Ground Truth

³ Self Occlusion

⁴ End To End

Archive of SID

تخمین حالت سه‌بعدی ارائه داد. در این روش با استفاده از الگوریتم محلی کوانتوم خطی، داده‌های بدون برچسب به مدل پیشنهادی وارد می‌شود. از این رو در مقاله [۱۸] با روش شبکه‌های عصبی کانولوشن، تخمین حالت دوبعدی انسان انجام شد و از طریق واژه‌نامه ساخته شده شامل تمام حالت‌های سه‌بعدی، سعی در تخمین حالت سه‌بعدی یک تصویر بوسیله روش KNN دارد. مقاله [۱۹] بر اساس چارچوبی بر پایه شبکه عصبی کانولوشن DeepCut [۲۰] جهت تخمین حالت دوبعدی از رویکردی بالا به پایین استفاده کرده است. در این مقاله پس از تخمین حالت دو بعدی، از میان مجموعه‌دادگان سه‌بعدی شکل^۱، نمونه‌ای با کمترین اختلاف شباهت با این تخمین، انتخاب شد. این روش برای مجموعه‌دادگان سه‌بعدی کوچک نیز مناسب است. اما مقاله [۲۱] شبکه‌ای عصبی کانولوشنی با هدف تخمین حالت دوبعدی و سه‌بعدی بصورت مشترک ارائه داد. در این شبکه وزن‌های آخرین لایه تمام اتصال تخمین حالت دوبعدی به‌عنوان بخشی از لایه تمام اتصال تخمین حالت سه‌بعدی قرار گرفته است.

مقاله [۲۲] با شبکه عصبی کانولوشن، اطلاعات حرکت را در دنباله‌ای کوتاه بررسی و تخمین حالت سه‌بعدی انجام داده است. اما مقاله [۲۳] نقشه ارتفاع^۲ را از تصویر دوبعدی استخراج کرده است. سپس بوسیله شبکه‌های عصبی کانولوشنی، تخمینی از حالت دوبعدی بدن انسان استخراج شده است و توسط نقشه عمق اصلاح می‌گردد. در ادامه با ساخت واژه‌نامه از حالت‌های سه‌بعدی، با الگوریتم PCA به استخراج حالت سه‌بعدی می‌پردازد. از این رو مقاله [۲۴] توسط شبکه‌های عصبی کانولوشن و کدگذاری تنک به تخمین حالت سه‌بعدی می‌پردازد. در این مقاله تخمین حالت دوبعدی توسط ترکیب نتایج دوشبکه عصبی کانولوشن انجام می‌گیرد سپس در ادامه توسط کدگذاری و بازنمایی تنک به تخصیص بهترین حالت سه‌بعدی نسبت به حالت دوبعدی تخمین زده شده می‌پردازد. از طرفی مقاله [۲۵] برای بازسازی حالت سه‌بعدی از یک شبکه عصبی کانولوشن استفاده کرده است. در این روش، شبکه عصبی کانولوشن برای تخمین حالت دوبعدی بدن انسان و زاویه بین بندها در مرحله لایه تمام اتصال و در نهایت استخراج حالت سه‌بعدی بدن انسان، آموزش انجام شده است.

استخراج حالت سه‌بعدی انسان از یک تصویر دوبعدی، با چالش‌های مختلفی روبرو است. از این رو در برخی مقالات، هدف تخمین حالت سه‌بعدی انسان با فرض وجود حقیقت اصلی حالت دو بعدی انجام گرفته است. در مقاله [۲۶]، منطبق بر حقیقت اصلی حالت دوبعدی، مشابه‌ترین حالات سه‌بعدی موجود در مجموعه‌دادگان تعیین می‌شود. در ادامه، با اعمال شروط حرکت‌شناسی حالت انسان، دقیق‌ترین حالت سه‌بعدی تخمین زده شد. از طرفی نویسندگان مقاله [۲۷] با دو واژه‌نامه مختص به

با وجود چالش‌های ابهام در عمق و خودانسدادی در این حوزه، می‌توان گفت تکنیک‌های خودکار، جایگزین مناسبی برای حل این مسئله هستند. چراکه روش مبتنی بر مدل، پس از ساخت مدل بدن منطبق بر بهترین مدل موجود، با اعمال محدودیت‌های اجرایی بر هر پارامتر، به‌عنوان مثال، رعایت نسبت طول بدن، اندام و زوایای مفاصل، سعی در تخمین حالت بدن انسان دارد [۱۱].

گام اصلی در مطالعه حرکت سه‌بعدی انسان، استخراج ویژگی دقیق از سیگنال‌های ورودی است. روش‌های اولیه از ویژگی‌های سطح پایینی استفاده می‌کردند که جهت کاهش ابعاد بردار ویژگی‌ها، از توصیف‌گرهای تصاویر بهره می‌بردند. در مقاله [۱۲] پس از نرمال سازی مجموعه‌دادگان حالت‌های سه‌بعدی، در مرحله اول، بر اساس نوع حالت سه‌بعدی، حالات دوبعدی مربوطه از دیدهای مختلف، استخراج شده است. در مرحله دوم، با استفاده از حالت دوبعدی استخراج شده از تصویر ورودی، بوسیله ساختاری درختی، سعی در اصلاح آن از طریق مقایسه با حالات استخراجی مرحله اول بر اساس روش k نزدیک‌ترین همسایه، صورت گرفته است و در نهایت تصمیم‌گیری با به حداقل رساندن خطای طرح نهایی انجام شد.

در این حوزه جهت دستیابی به تخمین دقیق‌تر، روش‌های مختلفی بر مبنای استخراج عمق دقیق پیشنهاد شده است. اما نتایج تحقیقات پژوهشگران نشان داده است که روش‌های یادگیری عمیق و کدگذاری و بازنمایی تنک، از دقت و کارایی بالاتری برخوردارند. به عبارتی، پس از موفقیت روش‌های مبتنی بر تکنیک‌های یادگیری عمیق جهت تخمین حالت دو بعدی، از این رو تکنیک در حوزه تخمین حالت سه‌بعدی نیز استفاده شد. از این رو در مقاله [۱۳] تخمین حالت، بر اساس رگرسیون یک فریم انجام شد. اما در مقاله [۱۴] آموزش شبکه‌های عصبی کانولوشن با دو نگرش صورت گرفته است. در نگرش اول، بصورت همزمان شبکه جهت تخمین حالت بدن بوسیله رگرسیون و نیز جهت تشخیص اشیا، آموزش انجام شده است. در نگرش دوم، ابتدا شبکه بصورت مستقل جهت تشخیص اشیا آموزش دیده و سپس در مرحله بعد بوسیله رگرسیون، تخمینی از حالت بدن انسان بدست می‌آید. آن‌ها نشان دادند که شبکه در آخرین لایه‌های خود، شامل نمایش داخلی مناسبی از ساختار اسکلت و همبستگی بین متغیرهای خروجی جهت تعیین حالت بدن است.

مقاله [۱۵] از یک چارچوب شامل یک تصویر و حالت سه‌بعدی به همراه ارزش همخوانی مربوطه به عنوان ورودی استفاده می‌کند. در این رویکرد، شبکه عصبی کانولوشن همزمان وظیفه‌ی انجام تبدیلات غیرخطی از تصویر و ترکیب‌بندهای تخمین زده شده و ساخت حالت بدن را بر عهده دارد.

از طرفی مقاله [۱۶] با کمک یک چارچوب بیزی و بر اساس کدگذاری و بازنمایی تنک، واژه‌نامه‌ای از نمونه‌های حالت ایجاد کرده است و به تشخیص شبیه‌ترین حالت نسبت به تصویر ورودی می‌پردازد. مقاله [۱۷] یک روش نیمه نظارتی کدگذاری تنک برای

¹ Shape

² Height Map

Archive of SID

به‌طور کلی تخمین حالت سه‌بعدی با رویکردهای مبتنی بر مدل، مبتنی بر یادگیری و مبتنی بر نمونه قابل انجام است. در روش مبتنی بر مدل، هدف تطابق مدل اولیه سه‌بعدی در نظر گرفته شده با تصویر دوبعدی شخص است، به‌گونه‌ای که با ایجاد تغییرات متناسب با تصویر دوبعدی، مدل سه‌بعدی با هدف به حداقل رساندن مقدار خطا در تطابق دوبعدی نسبت به تصویر تغییر داده شود. پس از اعمال تغییرات می‌توان مدل سه‌بعدی را تخمین از حالت دوبعدی بدن انسان در نظر گرفت. اما روش مبتنی بر مدل به دلیل زمان‌بر بودن در مراحل مختلف، کمتر مورد توجه قرار می‌گیرد. روش مبتنی بر یادگیری، سعی در یادگیری حالت‌های ورودی در مرحله آموزش نسبت به تصاویر ورودی دارد تا با توجه به ویژگی‌ها و الگوهای موجود در تصویر، حالت‌های مفصل را تشخیص دهد. روش مبتنی بر یادگیری دارای سرعت بالایی می‌باشد اما با توجه به گسترده بودن فضای حالت بدن، نمی‌توان تمام حالت‌ها را با تخمین مناسب بازسازی کرد. بر همین اساس روش مبتنی بر نمونه، حالت‌های متفاوت را در حافظه ذخیره کرده و با جستجو در بین آن‌ها، مشابه‌ترین حالت‌ها را انتخاب می‌کند. اما این روش به دلیل نیاز به حجم بالای حافظه و تامین انواع حالت‌ها برای تخمین حالت مناسب، هزینه‌بر خواهد بود.

در رویکرد مبتنی بر نمونه، با در نظر گرفتن پارامتر P ، بعنوان تعداد موقعیت‌های مدل، B_i حالت پایه سه‌بعدی و C_i وزن تاثیرگذاری برای حالت شماره i ، می‌توان حالت سه‌بعدی بدن را به صورت رابطه ۱ تعریف کرد [۲۸].

$$S = \sum_{i=1}^k C_i B_i \quad (1)$$

در این رابطه $S \in \mathbb{R}^{3 \times P}$ ، حالت استخراجی براساس رویکرد مبتنی بر نمونه می‌باشد. از طرفی، ارتباط بین تصاویر دوبعدی و سه‌بعدی با فرض $W \in \mathbb{R}^{2 \times P}$ حالت دوبعدی تبدیل یافته، $R \in \mathbb{R}^{2 \times 3}$ ماتریس چرخش متعامد و T بردار جابجایی، بصورت رابطه ۲ در نظر گرفته می‌شود.

$$W = RS + TI^T \quad (2)$$



شکل ۱: نمونه تصویر دوبعدی بدن انسان و حالت دوبعدی با ۱۵ بند.

داده‌های ورودی سایه‌نما و حالت‌های سه‌بعدی حقیقت اصلی، تخمینی از اطلاعات ورودی برحسب ترکیب خطی مولفه‌های واژه‌نامه اول استخراج شد، سپس با اعمال این ترکیب خطی در واژه‌نامه دوم، تخمین حالت سه‌بعدی انجام شد.

در مقاله [۲۸] با استفاده از کدگذاری و بازنمایی تنک، از رابطه‌های محدب به‌منظور ساخت نزدیک‌ترین مدل از ترکیب مولفه‌های واژه‌نامه استفاده شد. در این روش از الگوریتم ADMM با متغیر کمکی استفاده شده است و بصورت بازسازی مستقیم از خروجی الگوریتم ADMM، طرح نهایی تشکیل شد. لذا با توجه به خطای بازسازی مستقیم طرح سه‌بعدی که ناشی از تخمین نادرست وزن‌های مولفه‌های ترکیبی و ماتریس چرخش^۱ است، مقاله [۲۹] سعی در بهبود این روش داشتند. آن‌ها وزن و ماتریس چرخش را بر اساس یک معادله نرم Frobenius، جایگزین مرحله بازسازی مستقیم کرده و با حل و بروز رسانی آن از طریق الگوریتم Altern، سعی در بازسازی مناسبی برای مدل بودند. هدف این مقاله استفاده از کدگذاری و بازنمایی تنک به‌منظور بهبود دقت تخمین بوده است اما متأسفانه نتایج تخمین حالت سه‌بعدی، بهبود ملموسی را نشان نمی‌دهد. با این حال، این روش در تخمین عمق از دقت بالاتری برخوردار است.

در این مقاله، روشی جدید مبتنی بر شبکه‌های عصبی کانولوشن و کدگذاری و بازنمایی تنک با تصحیح گر مبتنی بر مدل، جهت تخمین حالت سه‌بعدی انسان از یک تصویر ارائه شده است. بخش‌های مختلف به این صورت سازماندهی می‌شود: در بخش دوم، مفاهیم اولیه بیان شده است. در بخش سوم، روش پیشنهادی بصورت چند مرحله‌ای بیان شده است. در بخش چهارم نتایج تجربی و مقایسه با سایر مقالات بیان شده است. بخش پنجم نیز حاوی نتیجه‌گیری می‌باشد.

۲ مفاهیم اولیه

۲-۱ ارتباط حالت دوبعدی و سه‌بعدی

تبدیل تصویر سه‌بعدی به دوبعدی، موجب از دست رفتن ابعاد و اطلاعات عمق فضا طی یک نگاشت یک به یک می‌شود. گرچه استفاده از تصاویر دوبعدی چنددیدی امکان بازسازی تصویر و تخمین مناسب با استفاده از عمق استخراجی را ممکن می‌سازد، اما تخمین حالت سه‌بعدی بدن انسان از تنها یک تصویر، بصورت بازسازی یک پیکربندی نقاط دلخواه سه‌بعدی بوسیله‌ی یک تصویر تک‌چشمی RGB است. این موضوع باعث تطابق حالت‌های سه‌بعدی متفاوت با حالت دوبعدی می‌شود. لذا خطای جزئی در حالت دوبعدی می‌تواند تأثیر بالایی بر تخمین حالت سه‌بعدی داشته باشد [۳۰]. در شکل ۱ نمونه‌ای از تصویر دوبعدی با حالت دوبعدی آن نمایش داده شده است.

¹ Rotation Matrix

۳-۲ مقدمه‌ای بر کدگذاری تنک

سیستم بینایی به‌عنوان یک سیستم کدگذاری به دنبال حذف وابستگی‌های آماری جهت توصیف تصاویر در قالب رویدادهای مستقل است. در بینایی ماشین نیز به‌منظور نمایش و ذخیره‌سازی داده‌های بزرگ از مکانیسم کدگذاری جهت تحلیل بهتر داده‌ها با حذف افزونگی‌های ناشی از وابستگی‌های آماری پیچیده استفاده می‌شود [۳۴]. از این رو محققین با استفاده از این سیستم‌های کدگذاری در سال‌های اخیر چالش‌های زیادی مانند تشخیص چهره [۳۵]، تشخیص حالت بدن انسان [۲۴ و ۳۶] و ... را حل کرده‌اند.

یکی از سیستم‌های کدگذاری معروف، کدگذاری تنک است. این کدگذار به‌دنبال بازسازی سیگنال‌ها در قالب یک ساختار تنک می‌باشد و داده ورودی را به‌صورت ترکیب خطی از مؤلفه‌های ویژه بیان می‌کند [۳۷]. بطور کلی هدف از این سیستم، حل چالش‌هایی با ابعاد بالای داده‌ها با مدل کردن مساله به کمک معادلات ریاضی می‌باشد. تنکی^۴ یک بردار با تعداد مؤلفه‌های صفر آن بردار تعیین می‌شود. در این کدگذاری، اگر بازنمایی بهینه در یک مسئله به اندازه کافی تنک باشد، آن مسئله با بهینه‌سازی محدب قابل حل خواهد بود. الگوریتم‌هایی که برای بدست آوردن جواب‌های تنک ارائه شده‌اند، همگی به دنبال توصیف تصویر ورودی y با ترکیب چند مؤلفه اولیه، از مجموعه مؤلفه‌های موجود در ویژه‌نامه D با ضرایب A بصورت رابطه ۴ می‌باشد.

$$y = DA \quad (۴)$$

با فرض $D=[D_1, D_2, \dots, D_n]$ و $A=[A_1, A_2, \dots, A_n]$ این رابطه را نشان می‌دهد. کدگذاری تنک می‌تواند بصورت خطی رابطه (۴) را با نرم صفر کمینه کند [۳۸]. بنابراین مسئله بهینه‌سازی به‌صورت رابطه ۵ بازنویسی می‌شود.

$$\hat{A} \sum_{i=1}^n D_i^2 = \arg \min \|A\|_0 \quad s.t. \quad y = DA \quad (۵)$$

که در آن $\| \cdot \|_0$ نشان‌دهنده تعداد عناصر غیرصفر در بردار مربوطه می‌باشد. از این رو، اگر k مؤلفه از ویژه‌نامه D ، نمونه آزمایشی را توصیف کند رابطه بهینه‌سازی بصورت رابطه ۶ تبدیل می‌شود.

$$y = DA \quad s.t. \quad \|A\|_0 \leq k \quad (۶)$$

از طرفی برای حذف نویز و طبق قضیه ضرب‌کننده لاگرانژ^۵ با ضریب ثابت α که مقدار انقباض A را مشخص می‌کند، رابطه ۶ بصورت رابطه ۷ بازنویسی می‌شود:

$$\hat{A} = L(A, \alpha) = \arg \min \|y - DA\|_2^2 + \alpha \|A\|_0 \quad (۷)$$

از الگوریتم‌های متداول در حل اینگونه مسائل بهینه‌سازی، MP^6 است. این الگوریتم به صورت حریصانه و با در نظر گرفتن نرم صفر، سعی در پیدا کردن بهترین مؤلفه در ویژه‌نامه را دارد. در

اما به‌دلیل نرمال‌سازی اسکلت بدن انسان، از بردار جابجایی صرف نظر کرده و رابطه ۱ و ۲ بصورت رابطه ۳ بازنویسی می‌شود.

$$W = R \sum_{i=1}^k C_i B_i \quad (۳)$$

در این مدل به دلیل نزدیک بودن بدن انسان نسبت به صفحه، ماتریس کالیبراسیون نادیده گرفته شده است [۲۸].

۲-۲ مقدمه‌ای بر شبکه عصبی کانولوشن

یادگیری عمیق [۳۱] به عنوان نسخه الهام گرفته شده از پرسپترون چند لایه (MLPS)، در زمینه‌های مختلف بینایی ماشین معرفی شده است. این مدل از یادگیری، تکنیکی مبتنی بر شبکه‌های عصبی است و روشی برای ارائه‌دادن یادگیری فرضیه‌ها [۳۲] می‌باشد که از تبدیلات غیرخطی چند منبعی تشکیل شده است. سلسه مراتب ویژگی‌ها در این مدل، از ویژگی‌های سطح پایین و تجمیع برای ویژگی‌های سطح بالاتر تشکیل شده است. در بکارگیری یادگیری عمیق، با توجه به روش و نوع معماری مورد استفاده، می‌توان از یادگیری بدون ناظر یا با ناظر و یا ترکیب این دو استفاده کرد [۳۳]. یکی از معروف‌ترین و موفق‌ترین معماری‌های یادگیری عمیق در حوزه آنالیز تصویر، شبکه عصبی کانولوشن است. علیرغم معرفی این شبکه‌ها در دهه ۸۰ [۱۰]، عمده گسترش آن‌ها در وظایف بینایی ماشین به‌دلیل در دسترس بودن منابع محاسباتی، در سال‌های اخیر صورت گرفته است. از اولین کاربردهای پیاده‌سازی شده در این زمینه می‌توان به شبکه Le Net [۱۱] اشاره کرد که جهت دسته‌بندی ارقام دست‌نویس در سال ۱۹۹۸ ارائه شد.

در سال‌های اخیر، معماری‌های متفاوتی متناسب با کاربردهای مختلف بر اساس شبکه‌های عصبی کانولوشنی ارائه شده است که غالباً متشکل از سه نوع لایه اصلی شامل لایه کانولوشن^۱، لایه ادغام^۲ و لایه تماماً متصل^۳ با ترکیب‌های متفاوت هستند. لایه کانولوشن برای محاسبه نقشه ویژگی، محاسبات را با فعال‌سازی یک نقشه با قابلیت تولید پاسخ فیلتر در هر موقعیت مکانی انجام می‌دهد. از لایه ادغام جهت کاهش بار محاسباتی لایه‌های کانولوشنی استفاده می‌شود. هدف از قرار دادن این لایه، دستیابی به فضایی تغییرناپذیر با کاهش رزولوشن نقشه ویژگی‌های مورد نظر است. از طرف دیگر لایه تماماً متصل مسئول استدلال سطح بالا در شبکه‌های عصبی کانولوشن است. این لایه معمولاً آخرین لایه از شبکه است و تمام نودهای لایه‌های قبلی را به‌عنوان ورودی پذیرفته و به یک نرون واحد متصل می‌کند. لذا پس از این لایه، تجسم و تحلیل اطلاعات مکانی داده‌ها به‌صورت یک‌بعدی امکان‌پذیر است.

⁴ Sparsity

⁵ Lagrange Multiplier

⁶ Matching Pursuit

¹ Convolutional Layer

² Pooling Layer

³ Fully Connected Layer

Archive of SID

این معماری با توجه به تعداد بالای تصاویر آموزشی، جهت کاهش تعداد پارامترها، ابعاد لایه ورودی از 224×224 به 112×112 تغییر یافته و به همین نسبت ابعاد تمام لایه‌های کانولوشنی کاهش پیدا کرده است. از سوی دیگر، بدلیل تفاوت تخمین حالت دوبعدی بدن انسان با طبقه‌بندی تصاویر، جهت ترکیب وزن‌ها به عنوان ورودی لایه‌های تماماً متصل و استفاده از اطلاعات فضایی محلی، از پنجره 7×7 (بدون قاب گذاری) برای شبکه تماماً متصل استفاده و به منظور تشخیص بندها، لایه ادغام آخر شبکه طبقه‌بندی VGGNet حذف می‌گردد.

در معماری‌های شبکه عصبی کانولوشن، به دلیل تعداد بسیار زیاد پارامترهای استفاده شده، ممکن است بیش‌برازش^۴ ایجاد شود. از این رو روش Dropout [۴۲ و ۴۳] به منظور جلوگیری از بروز این پدیده معرفی شده است. در هر مرحله از آموزش این روش، هر نورون با احتمال $p-1$ ، از شبکه حذف^۵ می‌شود بطوریکه بطوریکه نهایتاً یک شبکه کاهش داده شده باقی بماند. در صورت حذف یک نود، یال‌های ورودی و خروجی آن نیز حذف می‌شود.

بدین ترتیب در هر مرحله داده‌ها فقط بر روی شبکه کاهش یافته آموزش خواهد دید. پس از هر مرحله، نودهای حذف شده به همراه وزن‌های سابق آن‌ها (قبل از حذف شدن) دوباره به درون شبکه وارد می‌شوند. از این رو در معماری به منظور جلوگیری از ایجاد بیش‌برازش شبکه، قبل از هر لایه تماماً متصل، یک لایه dropout با نرخ ۰.۲۵ قرار گرفته است.

با اعمال تصاویر آموزشی و حالت‌های دوبعدی اسکلت‌بندی بدن انسان به عنوان ورودی شبکه، وزن‌ها جهت تشخیص مناسب بندهای بدن، اصلاح شده است. همچنین در لایه‌های تماماً متصل، لایه اول از ابعاد 4096×4096 به $1 \times 30 \times 30$ جهت مشخص نمودن ۱۵ بند اسکلت‌بندی و در انتها از کاهشی آنتروپی^۶ برای شبکه استفاده شده است.

در روش پیشنهادی، به دلیل یادگیری و تخمین مناسب شبکه عصبی کانولوشن با ابعاد پایین‌تر، همانند پژوهش‌های پیشین از تصاویر برش داده شده براساس قاب محتوای اصلی حالت دوبعدی بدن انسان در تصویر استفاده شده است و محتوای فاقد اطلاعات حذف گردیده است. به بیان دیگر، به دلیل تجهیزات پایین و تعداد زیاد تصاویر مجموعه‌دادگان جهت آموزش شبکه‌های عصبی کانولوشنی با پارامتر بالا، از رویکرد برش محتوای اصلی استفاده شده است و منجر به یادگیری و تخمین حالت دقیق‌تری از حالت دوبعدی بدن انسان در شبکه‌های عصبی کانولوشن با پارامترهای بسیار پایین‌تر می‌شود.

ادامه، الگوریتم OMP^۱ با متعامد^۲ ساختن برای رسیدن به بهترین تخمین در هر تکرار، الگوریتم MP را بهبود داد [۳۷]. اما بسیاری از مسائل را می‌توان در چارچوب بهینه سازی محدب قرار داد. لذا الگوریتم‌های بهینه‌سازی موازی به عنوان مکانیسمی برای حل مسائل آماری در مقیاس بزرگ معرفی و روش‌های تقسیم پذیر را به فرم زیر تبدیل می‌کند.

$$\text{minimize } f(x) + g(z) \quad (7)$$

$$\text{subject to } Ax + Bz = c$$

که $A \in \mathbb{R}^{p \times n}$ و $B \in \mathbb{R}^{p \times m}$ و با فرض $x \in \mathbb{R}^n$ ، $z \in \mathbb{R}^m$ و توابع f و g بصورت محدب می‌باشد. از الگوریتم‌هایی که این روش را شامل می‌شود ADMM^۳ [۳۴] و relaxed ADMM [۳۹] می‌باشد که دارای بهترین تخمین در مسائل درجه دوم و مخروطی است. در نهایت باید گفت که هر یک از این روش‌های بهینه‌سازی، نقاط قوت و ضعفی دارند که با توجه به نوع و ماهیت دادگان و نیز شرایط مساله، مورد استفاده قرار می‌گیرند.

۳ روش پیشنهادی

مراحل کلی روش پیشنهادی جهت تخمین حالت سه بعدی بدن انسان از یک تصویر در شکل ۳ نمایش داده شده است. در ادامه با توجه به خلاصه گفته شده از شبکه عصبی کانولوشن و کدگذاری و بازنمایی تنک، روش پیشنهادی شرح داده خواهد شد.

۳-۱ تخمین حالت دوبعدی از تصویر بوسیله شبکه عصبی کانولوشن

یکی از معماری‌های با عملکرد مناسب در حوزه پردازش تصویر و استخراج ویژگی، معماری VGGNet [۴۰] می‌باشد. معماری اولیه این شبکه با ۱۶ لایه جهت طبقه‌بندی تصاویر، با یادگیری مجموعه‌دادگان ImageNet ارائه شده است [۴۰]. با این حال این معماری دارای پیکربندی‌های دیگری نیز می‌باشد که تنها تفاوت آن‌ها در عمق شبکه است و حداکثر ۱۹ لایه می‌باشند. این شبکه به دلیل توالی مناسب در لایه‌های کانولوشنی جهت برقراری همبستگی بین ویژگی‌ها و عمق بالا، موجب ارائه عملکرد مناسب شده است. مقاله [۴۱] نشان داد که این شبکه با کاهش تعداد پارامترها از طریق تغییر اندازه تصویر ورودی و حذف لایه ادغام در آخرین لایه، تخمین حالت دوبعدی بدن انسان را به خوبی انجام می‌دهد. از این رو در این مقاله با رویکرد مبتنی بر یادگیری، شبکه‌ای مشابه این معماری، برای تخمین حالت دوبعدی بدن انسان در نظر گرفته شده است. معماری پیشنهادی، شامل ۱۳ لایه کانولوشنی و ۳ لایه تماماً متصل به صورت پیشرو منظم می‌باشد (شکل ۲).

⁴ Overfitting

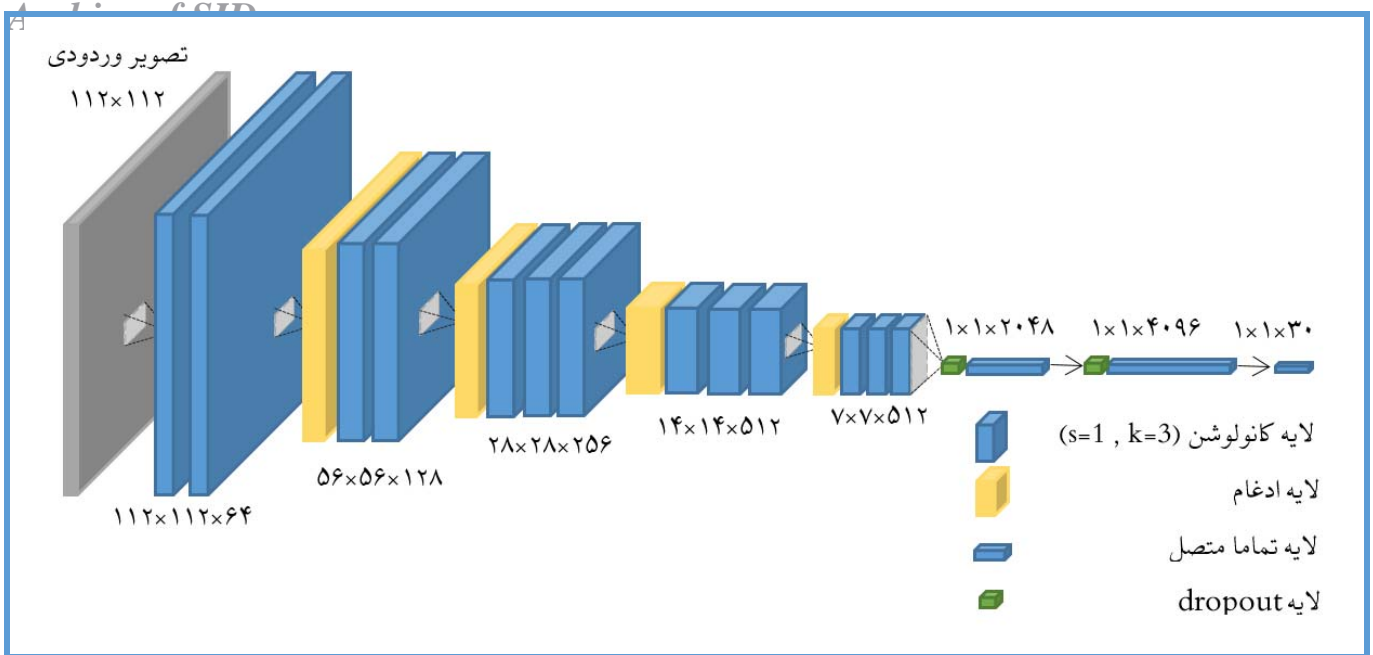
⁵ Dropped Out

⁶ Cross Entropy Loss

¹ Orthogonal Matching Pursuit

² Orthogonalization

³ Alternating Direction Method of Multipliers



شکل ۲: معماری شبکه عصبی کانولوشن روش پیشنهادی

با فرض $\tilde{B} = [B_1 \dots B_k]^T$ و $\tilde{M} = [M_1 \dots M_k]$ رابطه λ را می‌توان بصورت رابطه ۹ بازنویسی کرد.

$$\min_{\tilde{M}, Z} \frac{1}{2} \|W - Z\tilde{B}\|_F^2 + \alpha \sum_{i=1}^k \|M_i\|_2 \quad (9)$$

$$\text{s.t } \tilde{M} = Z$$

در این معادله Z یک متغیر کمکی و ضریب α براساس آزمایشات تجربی صورت گرفته برابر ۰٫۱ در نظر گرفته شده است. معادله بالا با استفاده از ضریب لاگرانژ بصورت رابطه ۱۰ بازنویسی می‌شود.

$$L_\mu(\tilde{M}, Z, Y) = \frac{1}{2} \|W - Z\tilde{B}\|_F^2 + \alpha \sum_{i=1}^k \|M_i\|_2 \quad (10)$$

$$+ \langle Y, \tilde{M} - Z \rangle + \frac{\mu}{2} \|\tilde{M} - Z\|_F^2$$

که در این رابطه، Y متغیر دوگانه و μ پارامتری برای کنترل اندازه گام‌ها در بهینه‌سازی می‌باشد [۳۴]. در نتیجه با روش الگوریتم relaxed ADMM [۳۹] در هر مرحله با ضریبی از \tilde{M}^t و \tilde{M}^{t+1} ، گام‌های زیر بروزرسانی می‌شوند.

$$\tilde{M}^{t+1} = \arg \min_{\tilde{M}} L_\mu(\tilde{M}, Z^t, Y^t) \quad (11)$$

$$Z^{t+1} = \arg \min_Z L_\mu(\beta \tilde{M}^{t+1} + (1-\beta)\tilde{M}^t, Z^t, Y^t) \quad (12)$$

$$Y^{t+1} = Y^t + \mu \left(\left(\beta \tilde{M}^{t+1} + (1-\beta)\tilde{M}^t \right) - Z^{t+1} \right) \quad (13)$$

در رابطه‌های ۱۲ و ۱۳، $\beta \in (0, 2)$ می‌باشد که با آزمایشات محدود انجام شده، بهترین نتایج در این روش با در نظر گرفتن مقدار ۰٫۹ برای β استخراج گردید. لازم به ذکر است اگر ضریب β برابر ۱ قرار گیرد این روش همانند روش ADMM [۳۵] خواهد بود. در

بطور کلی این شبکه با توجه به کاهش ابعاد لایه‌ها به نصف شبکه اولیه ۱۶ لایه VGG، تقلیل ابعاد بردار ویژگی اول از ۴۰۹۶ به ۲۰۴۸ در قسمت تماماً متصل و کاهش ابعاد دسته‌بند انتهایی از ۱۰۰۰ به ۳۰، تعداد پارامتر شبکه نسبت به شبکه VGG [۴۰]، ۷۴ میلیون پارامتر و نسبت به شبکه مقاله [۴۱]، ۵۹ میلیون پارامتر کاهش داشته است و دارای ۶۴ میلیون پارامتر است. آموزش این شبکه عصبی کانولوشنی با نرخ یادگیری ۰٫۰۰۰۱، بر روی یک دستگاه سیستم توسط پردازشگر گرافیکی GTX 1080 با ۸ گیگابایت حافظه پردازش در مدت زمان ۱۲ روز انجام گرفته است.

۲-۳ تخمین حالت سه‌بعدی بدن از حالت دوبعدی با روش کدگذاری تنک

همانگونه که بیان شد، در روش پیشنهادی به منظور تخمین حالت سه‌بعدی براساس حالت‌های دو بعدی از نرم ۲ الگوریتم ADMM relaxed استفاده می‌شود. به گونه‌ای که با ترکیبی از مدل‌های سه‌بعدی و ژه‌نامه، سعی در ساخت بهترین حالت سه‌بعدی بدن انسان نسبت به حالت دوبعدی اولیه می‌باشد. در این الگوریتم با استفاده از یک متغیر کمکی و بصورت رویکردی تکراری، روش مستقیم را با تبدیل به دو زیرمسئله حل می‌کند. از این رو براساس مقاله [۲۸] با فرض $M_i \in \mathbb{R}^{2 \times 3}$ تشکیل شده از ماتریس ضرایب C_i و دو ردیف اول ماتریس چرخش R_i ، رابطه ۳ به فرم رابطه ۸ تبدیل می‌شود.

$$\min_{M_1, \dots, M_k} \sum_{i=1}^k \|M_i\|_2, \quad (8)$$

$$\text{s.t. } W = \sum_{i=1}^k M_i B_i$$

Archive of SID

- مرحله اول: تخمین حالت سه بعدی توسط الگوریتم ۲ و خروجی های الگوریتم ۱ متشکل از ماتریس های M مطابق مولفه های واژه نامه، ساخته می شود.

Input: $M_1 \dots M_k$

Output: S

For $i=1$ to k do

$$c_i = \|M_i\|_2$$

$$r_i^{(1)} = m_i^{(1)} / c_i$$

$$r_i^{(2)} = m_i^{(2)} / c_i$$

$$r_i^{(3)} = r_i^{(1)} \times r_i^{(2)}$$

$$R_i = [r_i^{(1)}, r_i^{(2)}, r_i^{(3)}]^T$$

End

$$S = \sum_{i=1}^k c_i R_i B_i$$

الگوریتم ۲: بازسازی مستقیم سه بعدی [۲۹].

- مرحله دوم: تخمین عمق انجام شده در مرحله اول به نرمال تخمین حالت دوبعدی استخراج شده توسط شبکه عصبی کانولوشن اضافه می شود.
- مرحله سوم: ساختار درختی ۲ مدل برای حالت سه بعدی بدن انسان، براساس نزدیک ترین حالت های دوبعدی همتراز شده در واژه نامه با تخمین حالت دوبعدی، در نظر گرفته می شود.
- مرحله چهارم: با توجه به اینکه مدل های استخراجی در مرحله سوم از نظر نسبت اندازه، با حالت تخمین زده شده متفاوت می باشد، از این رو، ضریب نسبت بین ۲ مدل سه بعدی و مدل تخمین زده شده از طریق محاسبه جمع فاصله بندها با معیار فاصله اقلیدوسی در ساختار درختی، انجام می گیرد. در ادامه لازم است این ضرایب در مدل های مرحله سوم ضرب شوند تا نسبت بین اعضا برقرار گردد.
- مرحله پنجم: برای ۲ مدل بروزرسانی شده و مدل تخمین زده شده، فاصله اقلیدوسی بندها با اولویت فاصله از ریشه در فضای سه بعدی محاسبه می شود. سپس اندازه فاصله اقلیدوسی هر خط فاصل بین دو بند مدل تخمین زده شده و بازه اندازه فاصله اقلیدوسی دو مدل، مقایسه می شود. اگر اندازه فاصله دو بند در بین بازه تعیین شده قرار نگیرد، میانگین اندازه فاصله اقلیدوسی دو مدل ساختاری، به عنوان اندازه فاصله دو بند تخمین زده شده، در نظر گرفته می شود. در ادامه با استفاده از حالت دوبعدی تخمین زده و اندازه فاصله، بر اساس هم جهت بودن با عمق تخمینی اصلی، عمق بند انتهایی بروزسانی می شود.

ادامه به استفاده از روش relaxed ADMM [۴۰] و با محاسبه لاگرانژ، اگر هر ستون i از $Z^t - \frac{1}{\mu} Y^t$ برابر Q_i^t در نظر گرفته شود [۲۴]، رابطه ۱۱ به صورت رابطه ۱۴ بازنویسی می شود.

$$M^{t+1} = D_{\alpha} \left(Q_i^t \right) \quad (14)$$

این عملیات بصورت مجزا و تکراری، برای هر ستون براساس مسئله پروکسیمال انجام می گیرد [۲۴].

از طرفی در معادله ۱۲، $L_{\mu} \left(\beta \tilde{M}^{t+1} + (1-\beta) \tilde{M}^t, Z^t, Y^t \right)$ فرم درجه دوم می باشد و بصورت رابطه ۱۵ بازنویسی می شود.

$$Z^{t+1} = (W\tilde{B}^T + \mu \left(\left(\beta \tilde{M}^{t+1} + (1-\beta) \tilde{M}^t \right) + Y^t \right)) \left(\tilde{B}\tilde{B}^T + \mu I \right)^{-1} \quad (15)$$

با توجه به آنچه گفته شد، شبه کد روش پیشنهادی بصورت الگوریتم ۱ می باشد. همانند حالت های سه بعدی اعمال شده به واژه نامه، حالت های دوبعدی آزمایشی نیز نیاز به نرمال سازی دارند. نرمال سازی به منظور ایجاد میانگین برابر با صفر و انحراف معیار یک برای مقادیر هر سطر می باشد.

روش پیشنهادی ارائه شده در این مقاله سعی در ساختن بهترین تطابق از حالت دوبعدی با حالت های سه بعدی از طریق مولفه های واژه نامه دارد.

Input: W, α

Output: $M_1 \dots M_k$

Initialize $Z=Y=0, \mu>0, \beta \in (0, 2)$

While not converged do

Parallel for $i=1$ to k do

$$Q_i^t = \text{the } i\text{-th column-triplet of } Z^t - \frac{1}{\mu} Y^t$$

$$M^{t+1} = D_{\alpha} \left(Q_i^t \right)$$

End

$$Z^{t+1} = (W\tilde{B}^T + \mu \left(\left(\beta \tilde{M}^{t+1} + (1-\beta) \tilde{M}^t \right) + Y^t \right)) \left(\tilde{B}\tilde{B}^T + \mu I \right)^{-1}$$

$$Y^{t+1} = Y^t + \mu \left(\left(\beta \tilde{M}^{t+1} + (1-\beta) \tilde{M}^t \right) - Z^{t+1} \right)$$

End

الگوریتم ۱: حل رابطه ۹ با الگوریتم relaxed ADMM.

۳-۳ بازسازی مدل سه بعدی

بازسازی حالت سه بعدی تخمین زده شده از حالت های استخراجی از الگوریتم ۱ به صورت مبتنی بر مدل انجام می گیرد. شرح مراحل مختلف این بخش از روش پیشنهادی به این صورت می باشد:

۴ نتایج تجربی

۳-۴ ساخت واژه‌نامه

جهت تحلیل و ارزیابی دقیق‌تر نتایج حاصل از روش‌های مختلف، در ابتدای این بخش مجموعه‌دادگان استفاده شده معرفی خواهد شد. در بخش ۴-۲ آشنایی مختصری با معیارهای ارزیابی بکار گرفته شده صورت می‌گیرد. نهایتاً در بخش ۴-۳ نتایج حاصل از پیاده‌سازی روش پیشنهادی در مقایسه با کارهای پیشین ارائه خواهد شد.

۴-۱ مجموعه‌دادگان مورد استفاده

مجموعه‌دادگان Human 3.6M [۱۳] یک مجموعه‌دادگان بزرگ برای تخمین حالت سه‌بعدی انسان است و در این مقاله مورد استفاده قرار گرفته است. این مجموعه‌دادگان شامل میلیون‌ها حالت سه‌بعدی بدست آمده از تصاویر دوربین‌های کالیبره^۱ شده است. فیلم‌های موجود در این مجموعه‌دادگان از چند دید^۲ به صورت هماهنگ شده می‌باشند. همچنین این مجموعه داده‌های دوبعدی و سه‌بعدی از حالت‌های بدن انسان جهت ارزیابی با تعداد بالا قرار دارد. از مدل‌های انجام آزمایش بر روی این مجموعه‌دادگان، استفاده از مجموعه‌ای شامل ۱۱ بازیگر، ۷ مدل و ۱۵ موضوع می‌باشد و از ۵ مدل (S1, S5, S6, S7, S8) در ۱۵ موضوع به منظور آموزش و از ۲ مدل (S9, S11) در ۱۵ موضوع جهت آزمایش استفاده می‌شود. نتایج این مقاله نیز براساس این مدل ارزیابی شده است. ویدیوهای این مجموعه‌دادگان برای هر موضوع، با دوبار تکرار از ۴ دید متفاوت تهیه شده است و هر موضوع با نرخ اصلی ۵۰ فریم بر ثانیه است. در اکثر کارهای پیشین و روش پیشنهادی این مقاله، این نرخ به ۱۰ فریم بر ثانیه تبدیل شده است. استخراج اسکلت‌بندی بدن انسان توسط ۳۲ حسگر انجام گرفته است که در این پژوهش همانند سایر مقالات از ۱۵ بند اصلی استفاده شده است (شکل ۱).

۴-۲ معیارهای ارزیابی

از معیارهای متداول جهت سنجش نتایج روش‌های مختلف، محاسبه خطای هر بند، از میانگین فاصله اقلیدوسی در تمام مفاصل می‌باشد. در این معیار اگر موقعیت بندهای تخمین زده شده $\hat{x}_1, \dots, \hat{x}_n$ و موقعیت حقیقت اصلی x_1, \dots, x_n فرض شود، بصورت رابطه ۱۷ است.

$$e = \frac{1}{n} \sum_{i=1}^n \|\hat{x}_i - x_i\|_2 \quad (17)$$

مقدار خطا در رابطه ۱۷ وابسته به مقیاس و چرخش تخمین صورت گرفته است. از طرفی ابهام در عمق و مقیاس در بازسازی حالت سه‌بعدی از تصویر دوبعدی به صورت کامل قابل حل نمی‌باشد. از این رو برای مقایسه و سنجش دقت بازسازی، از

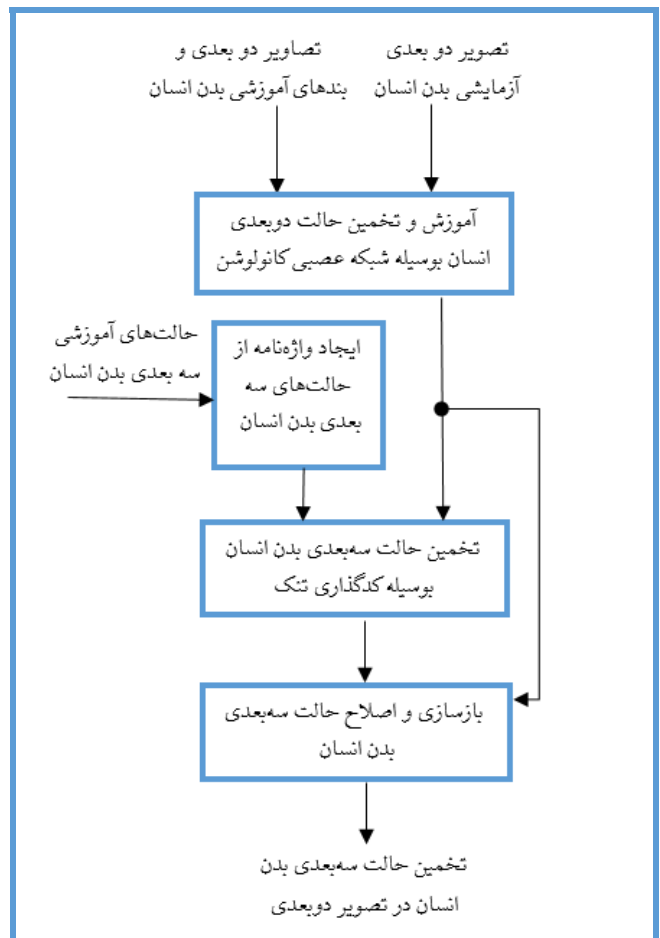
سیستم کدگذاری تنک به منظور ساخت واژه‌نامه، با هدف کاهش حجم داده‌ها و پوشش دادن تمام حالات ممکن آموزشی صورت می‌گیرد. از این رو در بعضی مسائل، از تمام داده‌های آموزشی جهت ذخیره در واژه‌نامه استفاده می‌شود، اما در رویکردی دیگر، بدلیل حجم بالای اطلاعات آموزشی، واژه‌نامه بصورت یادگیر، بگونه‌ای که تمام حالات را بر اساس حالات پایه پوشش دهد ساخته می‌شود. لذا در این پژوهش همانند مقاله [۲۹] جهت ساخت واژه‌نامه از بهینه‌سازی رابطه ۱۶ استفاده شده است و تعداد حالت‌های سه‌بعدی واژه‌نامه ۱۲۸ در نظر گرفته شده است.

$$\min_{B,C} \sum_{j=1}^n \frac{1}{2} \|S_j - \sum_{i=1}^k C_{i,j} B_i\|_F^2 + \lambda \|C\|_1 \quad (16)$$

$$\text{s.t. } C_{i,j} \geq 0, B_{iF} \leq 1,$$

$$\forall i \in [1, k], j \in [1, n],$$

در این معادله S_j حالت سه‌بعدی آموزشی، B_i حالت پایه یادگیری شده، C_{ij} ضریب i امین ضریب برای زامین حالت‌های آموزش می‌باشد [۲۹].



شکل ۳: مراحل روش پیشنهادی تخمین حالت سه‌بعدی بدن انسان

¹ Calibration

² Multi View

نزدیکترین حالت‌ها است. به بیان دیگر، مقاله [۳۶] از محاسبه نزدیکترین حالت دوبعدی و ترکیب با حالت‌های استخراجی قبلی استفاده می‌کند و قابلیت جلوگیری از انتخاب نزدیکترین حالت‌های دوبعدی که منجر به خطای بالاتر در ترکیب حالت‌ها را ندارد. از طرفی با توجه به ثابت بودن دوربین در هنگام فیلم‌برداری (اگرچه بازیگر بطور محدود در داخل صحنه جابجا می‌شود)، در مقاله [۴۵] محدودیت هندسی در بازسازی سه بعدی ایجاد می‌کند و خطای تخمین بالا می‌رود. اما در روش پیشنهادی، به دلیل استفاده از ساختارهای درختی از بندها به عنوان مدل، ساختار حالت بدن انسان حفظ می‌شود و مانع از ایجاد عمق نامناسب نسبت به اعضای بدن می‌گردد. در این نگرش برای حفظ تعادل بین خطای تخمین حالت دوبعدی و عمق تخمین زده شده، بازه تعیین شده در مرحله بازسازی، به عنوان یک تصحیح گر عمق‌های نامناسب عمل می‌کند. همچنین می‌توان چنین استنباط کرد تخمین‌های دقیق‌تر بیشتر در حوزه تخمین عمق می‌باشد و نسبت به حالت دوبعدی بهبود قابل توجهی بدست نیامده است.

با توجه به جدول ۱، بعضی حالات بدلیل تفاوت در مقیاس اندازه‌های ساختار درختی استخراج شده، خطای بازسازی عمق افزایش یافته و از حالت اصلی عمق حالت بدن، فاصله گرفته است. از این رو استفاده از ساختار مبتنی بر مدل، وابستگی بالایی به مدل ساخته شده را نشان می‌دهد. این رویکرد با توجه به محدودسازی واژه‌نامه و قابلیت تعمیم آن جهت ساخت حالت‌های جدید سه بعدی با رویکرد مقالات [۱۶، ۱۷ و ۲۷] متفاوت می‌باشد و بجای ذخیره سازی ۱۵۸۰۰۰ تصویر آموزشی، به ۱۲۸ حالت که کمترین خطا نسبت به ساخت تمام حالات را داشته باشد اکتفا می‌کند. در این مقالات علاوه بر اطلاعات فضایی حالت، از نقشه سیاه نما نیز استفاده می‌کند و نقطه ضعف این روش‌ها، عدم امکان نادیده گرفتن نزدیکترین حالت دوبعدی است زیرا این حالات همواره نزدیکترین حالت سه بعدی نسبت به هم را ندارند و همین امر باعث افزایش خطا می‌گردد. لذا با توجه به روابط مقاله [۱۷ و ۲۷] قیدهای نرم یک و نرم دو هستند و موجب می‌شود جواب منحصر بفردی داشته باشد؛ در عین حال که تنگ بودن پاسخ‌ها نیز تحقق یابد، بدین ترتیب با توجه به ساختار روابط می‌توان روابط را به شکل ساختار روابط رگرسیون در نظر گرفت. در شکل ۴ نمونه‌ی خروجی مراحل رویکرد اول نشان داده شده است.

در رویکردی دیگر، زمانی که حالت دوبعدی بدن انسان به عنوان ورودی اعمال نشود از رویکرد دوم استفاده می‌شود. این رویکرد در جدول ۲ با عنوان روش پیشنهادی با دیگر روش‌ها مقایسه شده است و بهبود عملکرد ۴ و ۲۰ درصدی نسبت به بهترین تخمین مقالات بررسی شده و میانگین خطای مقالات گذشته را نشان می‌دهد. همان‌طور که پیش‌تر بیان شد در حوزه‌های با فضای حالت بالا، قابلیت تعمیم تمام حالت‌ها در شبکه‌های عصبی کانولوشنی وجود ندارد.

میانگین حالت‌های آموزشی استفاده و حالت تخمین زده شده نسبت به مکان‌های ریشه اسکلت بدن هم‌تراز می‌شود. لذا ترازسازی به صورت دلخواه قابل قبول نمی‌باشد. بر این اساس، خطای بازسازی^۱ برای حالت سه بعدی به صورت رابطه ۱۸ تعریف می‌شود [۴۴].

$$r = \min_T \frac{1}{n} \sum_{i=1}^n \|\hat{x}_i - T(x_i)\|_2 \quad (18)$$

که T نشان‌دهنده ماتریس تبدیل است. عمدتاً خطای بازسازی در ساختارها به منظور ارزیابی دقت ساختار بهبود یافته بدون در نظر گرفتن مقیاس و صلب بودن حالت است.

۴-۳ ارائه و تحلیل نتایج تجربی

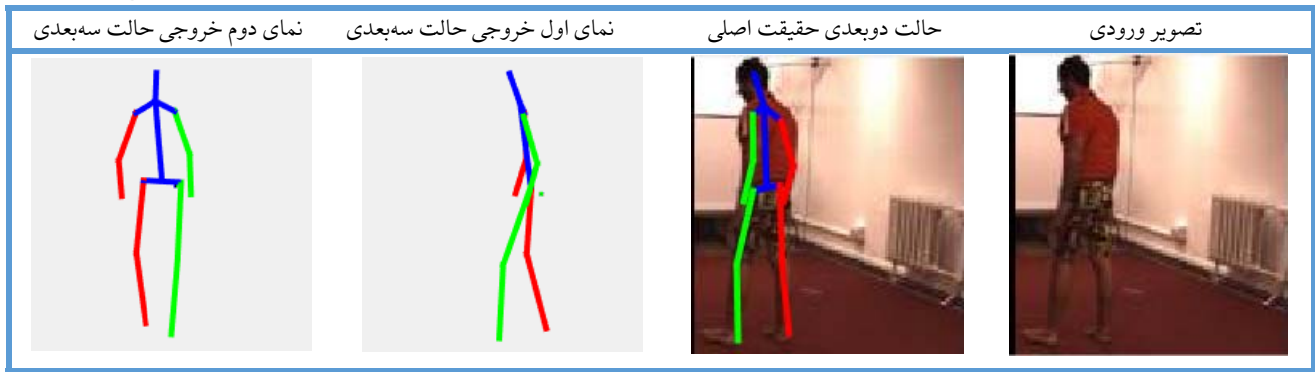
در این پژوهش تخمین حالت سه بعدی بدن انسان می‌تواند از طریق حالت دوبعدی یا از طریق استخراج ویژگی در تصویر و تخمین حالت دوبعدی انجام شود. این دو حالت به دلیل ساختاری با هم متفاوت هستند، زیرا در تخمین حالت دوبعدی از یک تصویر، سعی در تخمین حالت دوبعدی است و هدف نسبت دادن بهترین عمق به آن می‌باشد. اما در تخمین حالت سه بعدی از حالت دوبعدی فرض بر اینست که تخمین حالت دوبعدی انجام گرفته و هدف تخمین یک حالت سه بعدی بر اساس واژه‌نامه تشکیل شده بگونه‌ای است که از نظر خطای دوبعدی کمترین خطا را بدون استفاده از حالت دوبعدی اولیه داشته باشد. در نتیجه در این مقاله، با دو رویکرد متفاوت سعی در حل این چالش داریم. رویکرد اول شامل تخمین حالت سه بعدی بدن انسان بر اساس حالت دوبعدی مجموعه داده‌گان می‌باشد. این رویکرد بر اساس الگوریتم ۱ ضرایب مورد نیاز استخراج شده و با استفاده از مراحل بازسازی و بدون استفاده از شبکه عصبی کانولوشن، بصورت مبتنی بر مدل بهبود داده می‌شود. در رویکرد دوم پس از تخمین حالت دوبعدی بدن انسان از تصویر بوسیله شبکه عصبی کانولوشن، حالت‌های تخمین زده نرمال شده و در ادامه بوسیله الگوریتم ۱ ضرایب استخراج می‌شوند. سپس بازسازی توسط ۵ مرحله مطابق آنچه پیش‌تر در بخش ۳-۳ گفته شد انجام می‌گیرد.

نتایج حاصل از آزمایش‌های صورت گرفته رویکرد اول بر روی مجموعه داده‌گان Human 3.6M با ۱۵ بند اصلی در جدول ۱ با عنوان روش پیشنهادی ارائه می‌گردد. همان‌طور که مشخص است، نتایج حاصل از پیاده سازی رویکرد اول با بهره‌گیری از کدگذاری تنگ منجر به میانگین خطای بازسازی ۴۷.۵ می‌شود. این بهبود نشان از اهمیت بکارگیری روش‌های ترکیبی مبتنی بر مدل و مبتنی بر نمونه می‌باشد.

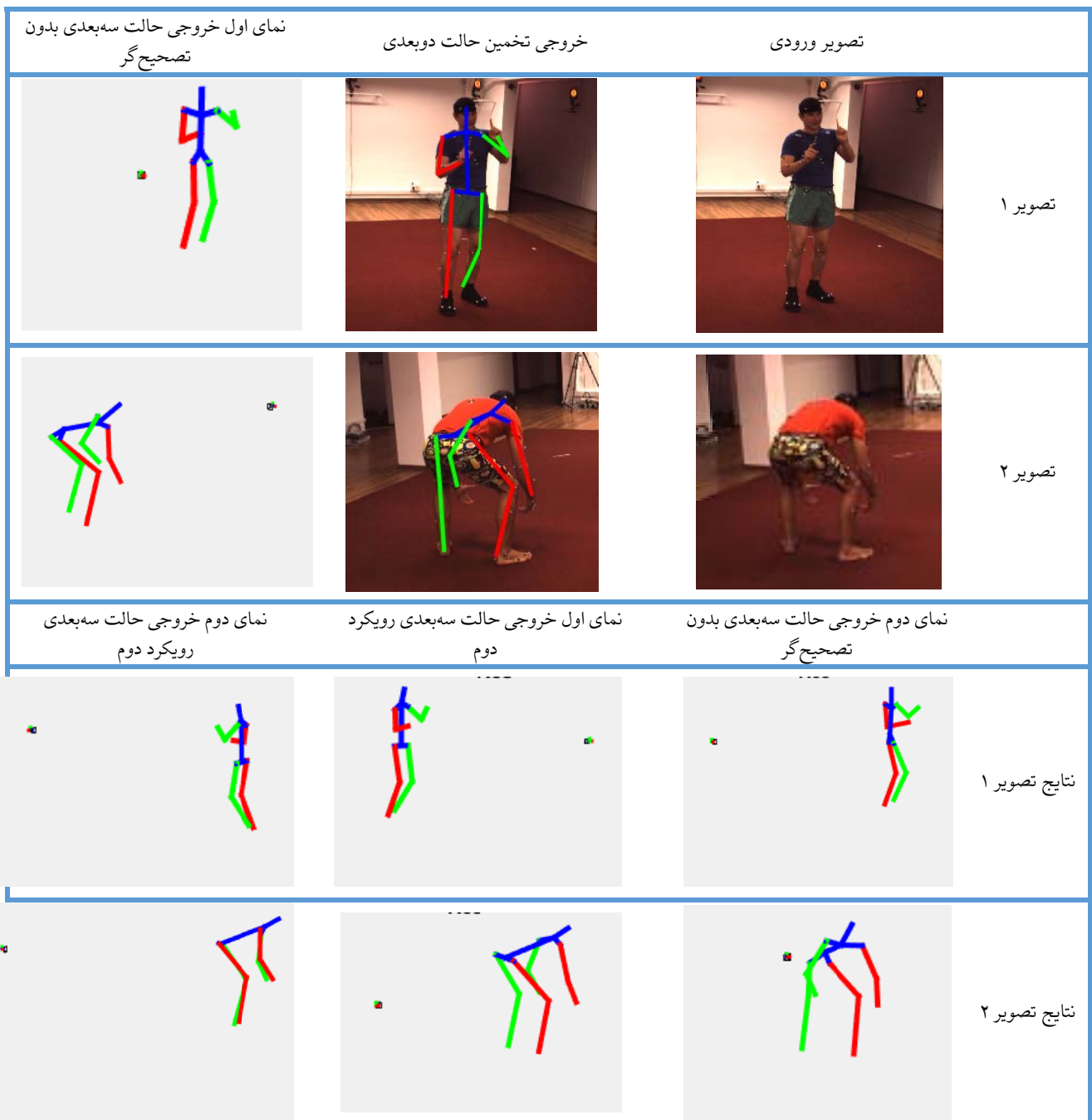
با توجه به جدول ۱، رویکرد حریمانه الگوریتم OMP در مقاله PMP [۳۶] نشان از ضعف این مدل در ترکیب و ساخت مدل با کمترین خطا، بدلیل عملیات حریمانه و استفاده از

¹ Reconstruction Error

Archive of SID



شکل ۴: نمونه‌ای از خروجی مراحل رویکرد اول روش پیشنهادی



شکل ۵: نمونه‌ای از خروجی مراحل رویکرد دوم پیشنهادی

Archive of SID

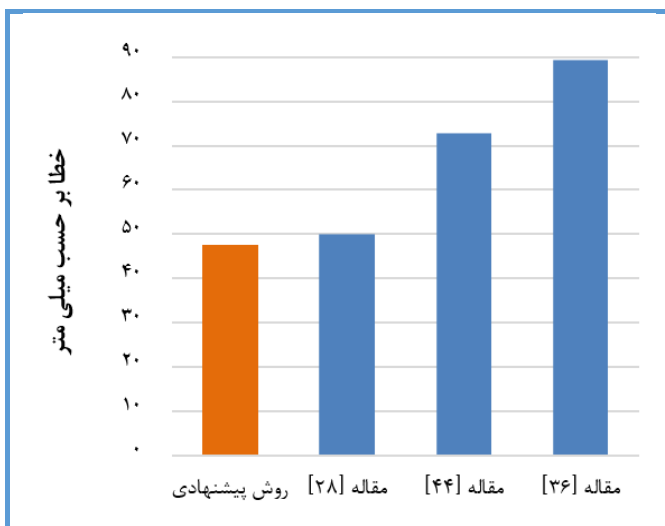
جدول ۲: دقت تخمین حالت سه بعدی بدن انسان از یک تصویر دوبعدی براساس خطای میانگین فاصله اقلیدسی (اعداد ذکر شده خطا بر حسب میلی متر می باشند).

	مقاله [۲۲]	مقاله [۲۳]	مقاله [۲۱]	مقاله [۲۵]	مقاله [۲۴]	روش پیشنهادی بدون تصحیح-گر	روش پیشنهادی
Directions	۱۰۲.۴	۸۵.۰	۱۰۰.۳	۹۱.۸	۸۳.۳	۸۶.۷	۷۹.۲
Discussion	۱۴۷.۷	۱۱۲.۶	۱۱۶.۱	۱۰۲.۴	۹۸.۴	۱۰۵.۲	۹۳.۷
Eating	۸۸.۸	۱۰۴.۹	۸۹.۹	۹۶.۶	۸۴.۸	۸۶.۸	۸۳.۶
Greeting	۱۲۵.۲	۱۲۲.۰	۱۱۶.۴	۹۸.۷	۱۰۴.۲	۱۰۳.۹	۹۹.۴
Phoning	۱۱۸.۰	۱۳۹.۰	۱۱۵.۳	۱۱۳.۳	۱۱۲.۷	۱۱۵.۳	۹۶.۸
Photo	۱۸۲.۷	۱۳۵.۹	۱۴۹.۵	۱۲۵.۲	۱۳۲.۴	۱۳۸.۴	۱۱۲.۷
Posing	۱۱۲.۳	۱۰۵.۹	۱۱۷.۵	۹۰.۰	۱۰۳.۳	۱۰۵.۴	۹۳.۷
Purchases	۱۲۹.۱	۱۶۶.۱	۱۰۶.۹	۹۳.۸	۹۷.۲	۹۹.۶	۸۹.۶
Sitting	۱۳۸.۸	۱۱۷.۴	۱۳۷.۲	۱۳۲.۱	۱۱۳.۳	۱۲۴.۸	۱۲۱.۲
Sitting Down	۲۲۴.۹	۲۲۶.۹	۱۹۰.۸	۱۵۸.۹	۱۸۳.۶	۱۹۶.۹	۱۹۱.۴
Smoking	۱۱۸.۴	۱۲۰.۰	۱۰۵.۷	۱۰۶.۹	۱۰۴.۳	۱۰۶.۲	۹۸.۶
Waiting	۱۳۸.۷	۱۱۷.۶	۱۲۵.۱	۹۴.۴	۱۰۹.۱	۱۱۵.۵	۱۰۳.۳
WalkDog	۱۲۶.۲	۱۳۷.۳	۱۳۱.۹	۱۲۶.۰	۱۰۶.۴	۱۱۴.۶	۱۱۳.۱
Walking	۵۵.۰	۹۹.۲	۶۲.۶	۷۹.۰	۷۶.۹	۷۹.۷	۷۸.۵
Walk Together	۶۵.۷	۱۰۶.۵	۹۶.۱	۹۸.۹	۹۳.۲	۹۵.۱	۹۱.۶
Average	۱۲۴.۹	۱۲۶.۴	۱۱۷.۴	۱۰۷.۲	۱۰۶.۹	۱۱۱.۶	۱۰۳.۱

از طرفی شبکه‌های عصبی کانولوشن در برخورد با فضای حالت بالا، خطای بالاتری را نشان می دهد. با توجه به نتایج جدول ۲، بدلیل فضای حالت بالای حرکات، نیاز به روش‌های ترکیبی هرچه بیشتر احساس می شود. لذا نتایج روش پیشنهادی بدون تصحیح گر، قدرت عملکرد تجمیع نتایج شبکه‌های عصبی کانولوشنی و مراحل مبتنی بر مدل هنگام بازسازی را نشان می دهد و از طرفی با توجه به نتایج مقاله [۲۴] می توان چنین استنباط کرد کاهش پارامتر در بخش شبکه عصبی کانولوشن منجر به خطای بالاتر در تخمین حالت دوبعدی می گردد اما استفاده دوباره از نتایج خروجی شبکه کانولوشنی، خطای دوبعدی را نسبت به بازسازی مستقیم کاهش می دهد و با اعمال ساختار درختی در بازسازی، امکان ایجاد عملکرد مناسبی در تخمین حالت سه بعدی گردد. لذا این جدول نشان می دهد تخمین حالت سه بعدی وابستگی بالایی به تخمین مناسب عمق دارد. اما بخش بزرگی از خطا، وابسته به حالت دوبعدی ورودی می باشد. از این رو جهت تخمین مناسب حالت دوبعدی براساس روش‌های یادگیر، شبکه‌های عصبی کانولوشن بدلیل برخورداری از دقت بالا، مورد استفاده قرار گرفته است. در شکل ۵ نمونه‌ی خروجی‌های مراحل رویکرد دوم بصورت تغییر ابعاد تخمین حالت دوبعدی جهت نمایش بر روی تصویر با ابعاد بزرگتر، نشان داده شده است. لذا بهبود عملکرد روش پیشنهادی در تخمین حالت سه بعدی بدن انسان براساس محل قرارگیری دوربین، قابل مشاهده می باشد. میانگین نتایج حاصل از آزمایش‌های صورت گرفته روش پیشنهادی نسبت به مقالات پیشین مطابق جدول ۱ و ۲ در شکل ۶ و ۷ بصورت نمودار نمایش داده شده است.

جدول ۱: دقت تخمین حالت سه بعدی بدن انسان از حالت دوبعدی حقیقت اصلی براساس میانگین خطای بازسازی (اعداد ذکر شده خطا بر حسب میلی متر می باشند)

	مقاله [۳۶]	مقاله [۴۵]	مقاله [۲۸]	روش پیشنهادی
Directions	۶۸.۵	۷۹.۲۶	۳۸.۲	۳۶.۶
Discussion	۷۷.۵	۶۰.۷۵	۴۵.۰	۴۰.۸
Eating	۹۵.۷	۱۲۵.۷۹	۴۷.۲	۴۳.۱
Greeting	۸۶.۴	۷۴.۹۷	۴۸.۵	۴۹.۲
Phoning	۷۳.۹	۳۷.۱۴	۴۵.۹	۳۹.۵
Photo	۹۵.۲	۵۸.۷۴	۶۳.۷	۵۸.۲
Posing	۷۸.۳	۶۱.۳۸	۴۷.۰	۴۳.۰
Purchases	۹۷.۸	۸۸.۴۷	۴۳.۸	۳۸.۹
Sitting	۱۰۳.۶	۱۱۲.۷۴	۵۳.۵	۵۸.۲
SittingDown	۱۲۳.۴	۱۱۴.۶۸	۷۰.۱	۷۵.۰
Smoking	۷۲.۸	۴۵.۰۹	۴۱.۲	۳۶.۴
Waiting	۹۵.۱	۷۸.۰۰	۴۸.۴	۴۶.۱
WalkDog	۸۳.۴	۶۲.۷۹	۵۳.۷	۴۷.۹
Walking	۹۵.۱	۶۳.۲۰	۴۷.۳	۴۸.۴
WalkTogether	۹۴.۹	۳۱.۶۷	۵۶.۵	۵۱.۸
Average	۸۹.۴	۷۲.۹۸	۵۰	۴۷.۵



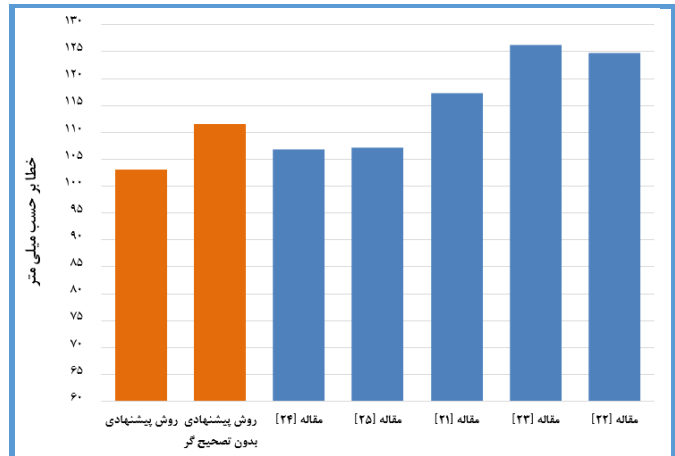
شکل ۶: مقایسه میانگین خطای بازسازی تخمین حالت سه بعدی بدن انسان از حالت دوبعدی حقیقت اصلی.

Archive of SID

علیرغم نقاط قوت ذکر شده، متأسفانه روش تصحیح گر حالت سه‌بعدی در مواردی که تخمین حالت دوبعدی نامناسب انجام گرفته است عمق را نسبت به حالت اصلی کاهش می‌دهد و علاوه بر خطای بالای حالت دوبعدی، خطای عمق نیز افزایش می‌یابد. لذا در مرحله تخمین حالت دوبعدی از تصویر، نیاز به شبکه عصبی کانولوشنی دقیق‌تر می‌باشد. از طرفی بازسازی حالت سه‌بعدی به صورت مستقیم نیز باعث خطای بالایی می‌شود و نیاز به روش‌های دقیق‌تر با کدگذاری و بازنمایی تنک احساس می‌شود.

مراجع

- [1] Felzenszwalb, P.F., Huttenlocher, D.P., "Pictorial structures for object recognition", International Journal of Computer Vision, Vol. 61(1), pp. 55–79, 2005.
- [2] Yang, Y., Ramanan, D., "Articulated pose estimation with flexible mixtures-of-parts", Computer Vision and Pattern Recognition (CVPR), pp. 1385–1392, 2011.
- [3] Dantone, M., Gall, J., Leistner, C., Van Gool, L., "Human pose estimation using body parts dependent joint regressors", Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3041–3048, 2013.
- [4] Toshev, A., Szegedy, C., "DeepPose: Human pose estimation via deep neural networks", Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1653–1660, 2014.
- [5] Krizhevsky A., Sutskever I., Hinton, G.E. "Imagenet classification with deep convolutional neural networks", Advances in Neural Information Processing Systems (NIPS), pp. 1097–1105, 2012.
- [6] Carreira J., Agrawal P., Fragkiadaki K., Malik J., "Human pose estimation with iterative error feedback", IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.
- [7] Fan, X., Zheng, K., Lin, Y., Wang, S., "Combining local appearance and holistic view: Dual-source deep neural networks for human pose estimation", IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1347–1355, 2015.
- [8] Wei S.E., Ramakrishna V., Kanade T., Sheikh Y., "Convolutional pose machines". IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.
- [9] Chu, X., Ouyang, W., Li, H., Wang, X., "Structured feature learning for pose estimation", 2016.
- [10] Bourdev, L., Malik, J.: Poselets, "Body part detectors trained using 3d human pose annotations". IEEE 12th International Conference on Computer Vision, pp. 1365–1372, 2009.
- [11] Sarafianos, N., Boteanu, B., Ionescu, B., Kakadiaris, I. A., "3D human pose estimation: A review of the



شکل ۷: مقایسه میانگین خطای فاصله اقلیدسی تخمین حالت سه بعدی بدن انسان از یک تصویر دوبعدی.

۵ نتیجه‌گیری و پیشنهادهای آتی

در این مقاله، یک چارچوب تخمین حالت سه‌بعدی بدن انسان از تنها یک تصویر دوبعدی ارائه شده است. این موضوع علاوه بر کمک به تشخیص حالت بدن انسان، می‌تواند در مراکز حساس اداری و نظامی به عنوان سیستم هشدار جهت تشخیص حالت‌های مشکوک، صنعت پویانمایی و فیلم‌سازی در جهت طبیعی شدن حرکات موجودات ساخته‌شده با جلوه‌های ویژه و ... مورد استفاده قرار گیرد. از این رو نیاز به مدل کردن مکانیسمی مشابه عملکرد قشر بینایی و مغز انسان می‌باشد، لذا استخراج ویژگی‌های با رویکرد تشخیص حالت بدن انسان، جز با تحلیل دقیق مکانیسم تشخیص و تصمیم‌گیری در مغز انسان میسر نخواهد بود. از این رو، در روش پیشنهادی سعی گردید که این موارد لحاظ شود.

استخراج ویژگی‌های تصویر، مرحله اصلی و مهم تخمین حالت دوبعدی می‌باشد. از این رو، از قدرت بالای شبکه عصبی کانولوشن جهت تخمین و آشکارسازی بندهای بدن انسان استفاده شده است. در این مقاله از شبکه‌ای ۱۶ لایه‌ای با ابعاد لایه پایین استفاده شده است تا علاوه بر رویکرد کاهش پارامتر، دقت کافی جهت تخمین ۱۵ بند حالت دوبعدی در نظر گرفته شده را داشته باشد.

از سوی دیگر، در چالش تخمین حالت سه‌بعدی، با توجه به ماهیت حالت بدن انسان، بر پایه سیستم کدگذاری و بازنمایی تنک به تخمین و بازسازی توسط واژه‌نامه حالت سه‌بعدی پرداخته شده است و عمق را به خروجی نرمال شده شبکه عصبی کانولوشنی اضافه می‌کنیم، سپس عمق حالت تخمین زده شده را بصورت مبنی بر مدل تصحیح شده است تا معضل عمق نامناسب را نسبت به مقالات گذشته کاهش دهیم و از تاثیر خطای ایجاد شده در تخمین حالت سه‌بعدی در کدگذاری و بازنمایی تنک بکاهیم. نتایج نشان می‌دهد این روش از نظر دقت با معیار متداول میانگین خطای بازسازی و میانگین خطای بازسازی هر اتصال نسبت به کارهای پیشین برتری دارد.

- [26] Simo-Serra, E., Ramisa, A., Aleny`a, G., Torras, C., MorenoNoguer, F., "Single Image 3D Human Pose Estimation from Noisy Observations," Proc. the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2673-2680, 2012.
- [27] Zolfaghari M., Jourabloo A., Ghareh Gozlou S., Pedrod B., Manzuri-Shalmani M. T., "3D Human Pose Estimation from Image Using Couple Sparse Coding," Machine Vision and Applications, 25(6), pp. 1489-1499, 2014.
- [28] Zhou, X., Leonardos, S., Hu, X., Daniilidis, K., "3D Shape Reconstruction from 2D Landmarks: A Convex Formulation", IEEE International Conference on Computer Vision and Pattern Recognition, 2015.
- [29] Zhou, X., Zhu, M., Leonardos, S., Daniilidis, K., "3D shape estimation from 2D landmarks: A convex relaxation approach", IEEE transactions on pattern analysis and machine intelligence, Vol: 39 (8), pp: 1648-1661, 2017.
- [30] Agarwal A., Triggs B., "Recovering 3D human pose from monocular images", IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol: 28 (1), pp. 44-58, 2006.
- [31] Bengio, Y., "Learning deep architectures for AI," Found. trends Mach. Learn., vol. 2, no. 1, pp. 1-127, 2009.
- [32] Bengio, Y., Courville, A., Vincent, P., "Representation learning: A review and new perspectives", IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol: 35 (8), pp. 1798-1828, 2013.
- [33] Deng, L., Yu, D., "Deep learning: Methods and applications, Foundations and Trends in Signal Processing", Vol: 7, pp. 197-387, 2014.
- [34] Boyd, S., "Distributed optimization and statistical learning via the alternating direction method of multipliers," Foundations and Trends in Machine Learning, vol. 3, no. 1, pp. 1-122, 2010.
- [35] Xu, Y., Zhang, Z., Lu, G., Yang, J., "Approximately symmetrical face images for image preprocessing in face recognition and sparse representation based classification," Pattern Recognition, 54, pp. 68-82, 2016.
- [36] Ramakrishna, V., Kanade, T., Sheikh, Y., "Reconstructing 3D human pose from 2D image landmarks," ECCV, 2012.
- [37] Zhang, Z., Xu, Y., Yang, J., Li, X., Zhang, D., "A survey of sparse representation: algorithms and applications", IEEE Access, pp. 490-530, 2015.
- [38] Donoho, D., Elad, M., "Optimally sparse representation in general (nonorthogonal) dictionaries via l_1 minimization," Proceedings of the National Academy of Sciences, vol. 100, no. 5, pp. 2197-2202, 2003.
- [39] Banjac, G., Goulart, P., Stellato, B., Boyd, S., "Infeasibility detection in the alternating direction literature and analysis of covariates", CVIU, 152, pp. 1-20, 2016.
- [12] Yasin, H., Iqbal, U., Kruger, B., Weber, A., Gall J., "A dual-source approach for 3D pose estimation from a single image", Proc IEEE Conference on Computer Vision and Pattern Recognition, 2016.
- [13] Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C., "Human3.6m: Large scale datasets and predictive methods for 3D human sensing in natural environments," PAMI, vol. 36, no. 7, pp. 1325-1339, 2014.
- [14] Li, S., Chan, A. B., "3D human pose estimation from monocular images with deep convolutional neural network", Proc. 12th Asian Conference on Computer Vision, pp. 332-347, 2014.
- [15] Li, S., Zhang, W., Chan, A. B., "Maximum-margin structured learning with deep networks for 3D human pose estimation", Proc. IEEE International Conference on Computer Vision, pp. 2848-2856, 2015.
- [16] Babagholami-Mohamadabadi B., Jourabloo A., Zarghami A., Kasaei S., "A Bayesian framework for sparse representation-based 3-D human pose estimation". IEEE Signal Proc Lett. 21(3):297-300, 2014.
- [17] Andalib, A., Babamir, S. M., Faraji, A., "A NEW SPARSE REPRESENTATION ALGORITHM FOR 3D HUMAN POSE ESTIMATION", Computing and Informatics, Vol. 35, pp. 1338-1355, 2016.
- [18] Chen, C. H., Ramanan, D., "3d human pose estimation= 2d pose estimation+ matching", 2016.
- [19] Bogo F., Kanazawa A., Lassner C., Gehler P., Romero J., Black M. J., "Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image". ECCV, pp. 561-578, 2016.
- [20] Pishchulin, L., Insafutdinov, E., Tang, S., Andres, B., Andriluka, M., Gehler, P., Schiele, B., "DeepCut: Joint subset partition and labeling for multi person pose estimation". IEEE Conference on Computer Vision and Pattern Recognition, CVPR, pp. 4929-4937, 2016.
- [21] Park, S., Hwang, J., Kwak, N., "3D human pose estimation using convolutional neural networks with 2D pose information", ECCVW, 2016.
- [22] Tekin, B., Rozantsev, A., Lepetit, V., Fua, P., "Direct prediction of 3D body poses from motion compensated sequences", CVPR, 2016.
- [23] Du, Y., Wong, Y., Liu, Y., Han, F., Gui, Y., Wang, Z., Kankanhalli, M., Geng, W., "Marker-less 3D human motion capture with monocular image sequence and height-maps", ECCV, 2016.
- [24] Alikarami, H., Yaghmaee, F., Fasaieislam, M. j., "Sparse representation and convolutional neural networks for 3D human pose estimation", 3rd Iranian Conference on Intelligent Systems and Signal Processing (ICSPIS), pp. 188 - 192, 2017.
- [25] Zhou, X., Sun X., Zhang, W., Liang, S., Wei, Y., "Deep kinematic pose regression", ECCVW, 2016.

Archive of SID



حسن علی کرمی تحصیلات خود را در دوره کارشناسی و کارشناسی ارشد خود را به ترتیب در سال‌های ۱۳۹۳ و ۱۳۹۶ در رشته مهندسی کامپیوتر- نرم افزار و مهندسی کامپیوتر- هوش مصنوعی و رباتیک از دانشگاه های شهید چمران اهواز و دانشگاه سمنان به اتمام رسانده است. ایشان هم اکنون دانشجوی دکترای مهندسی کامپیوتر- هوش مصنوعی و رباتیک در دانشگاه سمنان می باشد. زمینه های مورد علاقه ایشان پردازش تصویر، بازشناسی الگو و یادگیری عمیق می باشد.



فرزین یغمایی دوره کارشناسی کامپیوتر را در دانشگاه صنعتی امیرکبیر، کارشناسی ارشد و دکترای خود را در رشته مهندسی کامپیوتر- هوش مصنوعی در دانشگاه صنعتی شریف به پایان رسانده است و در حال حاضر دانشیار گروه مهندسی کامپیوتر دانشگاه سمنان می باشد. زمینه های مورد علاقه ایشان پردازش تصویر و ویدئو، بازشناسی الگو و متن کاوی می باشد.



محمدجواد فدایی اسلام استادیار گروه مهندسی کامپیوتر دانشگاه سمنان می باشد. ایشان مدارک کارشناسی مهندسی کامپیوتر- سخت افزار و کارشناسی ارشد و دکترای مهندسی کامپیوتر- هوش مصنوعی را از دانشگاه علم و صنعت تهران دریافت کرده است. زمینه های مورد علاقه ایشان پردازش تصویر و ویدئو، داده کاوی است.

method of multipliers for convex optimization”, Foundations and Trends in Machine Learning, eprints for the optimization community, 2017.

- [40] Simonyan, K., Zisserman, A., “Very deep convolutional networks for large-scale image recognition,”, 2014.
- [41] Yang, W., Ouyang, W., Li, H., Wang, X., “End-to-end learning of deformable mixture of parts and deep convolutional neural networks for human pose estimation”, CVPR, 2016.
- [42] Dai, Q., D. Hoiem. “Learning to localize detected objects”, Computer Vision and Pattern Recognition (CVPR), 2012.
- [43] Dai, J., He K., Sun J.. “Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation”, Proceedings of the IEEE International Conference on Computer Vision. 2015.
- [44] Zhou, X., Zhu, M., Pavlakos, G., Leonardos, S., Derpanis, K. G., Daniilidis, K. “Monocap: Monocular human motion capture using a cnn coupled with a geometric prior”, 2017.
- [45] Dai, Y., Li, H., and He, M., “A simple prior-free method for nonrigid Structure from motion factorization,” IJCV, vol. 107, no. 2, pp. 101–122, 2014.