



دانشگاه صنعتی شریف  
دانشکده‌ی مهندسی کامپیوتر

پایان‌نامه‌ی کارشناسی  
گرایش مهندسی فناوری اطلاعات

عنوان:

# بازشناسی کنش انسان از داده‌های اسکلتی توسط شبکه‌های عصبی گراف-پیمشی زمان-مکانی با مدل توجه

نگارش:

رضا رحیمی آذغان

استاد راهنما:

دکتر کسایی

شهریور ۱۳۹۸

سلام

به نام خدا  
دانشگاه صنعتی شریف  
دانشکده‌ی مهندسی کامپیوتر

## پایان‌نامه‌ی کارشناسی

عنوان: بازشناسی کنش انسان از داده‌های اسکلتی توسط شبکه‌های عصبی گراف-پیچشی  
زمان-مکانی با مدل توجه  
نگارش: رضا رحیمی آذغان

## کمیته‌ی ممتحنین

امضاء: استاد راهنما: دکتر کسایی

امضاء: استاد مدعو: دکتر همت‌یار

تاریخ:

## سپاس

از استاد بزرگوارم که با کمک‌ها و راهنمایی‌های بی‌دریغشان، بنده را در انجام این پروژه یاری داده‌اند،  
تشکر و قدردانی می‌کنم. هم‌چنین از سرکار خانم اسدی که در هر مرحله از پروژه بنده را راهنمایی کردند  
و راه صحیح را به نشان دادند از صمیم قلب سپاس‌گزارم.

## چکیده

در سال‌های اخیر، بازشناسی کنش انسان<sup>۱</sup> از عمده‌ترین زمینه‌های مورد بحث و تحقیق در دنیای علوم و مهندسی کامپیوتر بوده است. [۱] با افزایش چشم‌گیر اطلاعات در دسترس و پیشرفت روزافزون شبکه‌های عصبی<sup>۲</sup>، جلوه‌ی جدیدی به موضوع بازشناسی کنش انسان داده شده است. از مقدمات این موضوع، بحث نحوه‌ی نمایش داده‌های بدن انسان در دو بعد زمان و مکان است. با ظهور سنسورهای کینکت، دو روش عمده برای بهینه‌کردن هرچه بیشتر نمایش این داده‌ها وجود دارد. این دو روش شامل استفاده از اطلاعات RGB-D و استفاده از اطلاعات سه‌بعدی اسکلت‌های بدن هستند. [۲] اخیراً، به دلیل کم‌بودن حجم داده‌های اسکلتی، اطلاعات معنایی<sup>۳</sup> بالا و نیز خاصیت مقیاس‌پذیری آن‌ها، مطالعات بسیاری بر روی نمایش داده‌ها بر این روش صورت گرفته است.

در این پروژه مطالعه بر روی نمایش داده‌های سه‌بعدی اسکلتی و استفاده از آن برای بازشناسی کنش انسان ادامه پیدا می‌کند. شبکه‌ی مورد استفاده در این پروژه شبکه‌ی پیچشی-گرافی زمان-مکانی<sup>۴</sup> است که تعمیمی بر شبکه‌ی پیچشی-گرافی<sup>۵</sup> است که آن نیز تعمیمی بر شبکه‌ی پیچشی<sup>۶</sup> است. هم‌چنین سعی بر آن بوده است که با ارائه‌ی یک مدل توجه<sup>۷</sup> و با استفاده از معیار فلان شبکه‌های عصبی موجود را بهبود بخشید. در کارهای آینده نیز، می‌توان مدل‌های توجه بهینه‌تری را معرفی نمود و تاثیر هرکدام را بر شبکه‌ی مورد استفاده در این پروژه مطالعه نمود.

**کلیدواژه‌ها:** بازشناسی کنش انسان، اطلاعات سه‌بعدی اسکلت، شبکه‌های گراف-پیچشی، مدل توجه

---

<sup>۱</sup> Human Action Recognition

<sup>۲</sup> Neural Networks

<sup>۳</sup> Semantic Information

<sup>۴</sup> Spatial Temporal Graph Convolutional Networks

<sup>۵</sup> Graph Convolutional Networks

<sup>۶</sup> Convolutional Networks

<sup>۷</sup> Attention Model

# فهرست مطالب

۱۰	۱ مقدمه
۱۰	۱-۱ تعریف مسئله
۱۲	۲-۱ اهمیت موضوع
۱۳	۳-۱ ادبیات موضوع
۱۴	۴-۱ چالش‌ها
۱۵	۵-۱ فرضیات
۱۵	۶-۱ اهداف تحقیق
۱۶	۷-۱ ساختار پایان‌نامه
۱۷	۲ ادبیات مربوطه
۱۷	۱-۲ شبکه‌های عصبی پیچشی
۱۹	۲-۲ شبکه‌های عصبی گراف-پیچشی
۲۱	۳-۲ مدل‌های توجه
۲۳	۳ روش‌های پیشنهادی
۲۳	۱-۳ مقدمه
۲۴	۲-۳ شبکه‌ی استفاده شده

۳-۳ مدل توجه ..... ۲۵

۴ نتایج تجربی ..... ۲۸

۴-۱ مجموعه داده‌ی مورد استفاده ..... ۲۸

۴-۲ معیار تابع هزینه‌ی آنتروپی متقابل ..... ۲۹

۵ جمع‌بندی و راه‌کارهای آتی ..... ۳۰

# فهرست شکل‌ها

- ۱-۱ الف- نمایش اسکلت‌های بدن که از تصویر استخراج شده‌اند، ب- تصویر ژرفا،  
ج- تصویر رنگی ..... ۱۱
- ۱-۲ نمایشی از ابعاد زمان-مکانی اسکلت انسان که ST-GCN بر روی آن کار می‌کند. ۱۲
- ۱-۳ نمای کلی از یک LSTM ..... ۱۳
- ۲-۱ نمای کلی از یک شبکه‌ی پیچشی ..... ۱۸
- ۲-۲ نمای کلی از یک GCN ..... ۲۰
- ۲-۳ مدل توجه با استفاده از حافظه‌ی زمینه‌ی سراسری در یک LSTM ..... ۲۱
- ۳-۱ نمودار بلوکی برای شبکه‌ی ST-GCN با مدل توجه ادغامی ..... ۲۳
- ۲-۳ مثالی از یک گراف زمان-مکانی ..... ۲۴
- ۳-۳ نموداری از مدل توجه با استفاده از لایه‌ی ادغام ..... ۲۷
- ۱-۴ برخی از قاب‌های مجموعه داده‌ی NTU RGB+D ..... ۲۹



# فهرست جدول‌ها

۱-۱	نتایج برخی آزمایش‌ها بر روی پایگاه داده‌ی NTU RGB+D	۱۴
-----	---	----

# فصل ۱

## مقدمه

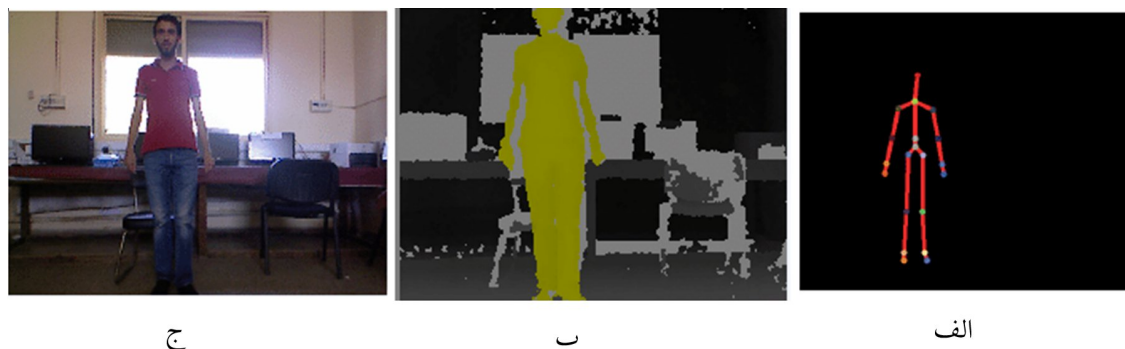
### ۱-۱ تعریف مسئله

یکی از حوزه‌های بسیار پرکاربرد و پرمباحثه در زمینه‌ی پردازش تصویر و بینایی کامپیوتر، بازشناسی کنش انسان است. از دهه‌ی ۸۰ میلادی، این حوزه به دلیل کاربرد بسیار بالایی که در زمینه‌های تعامل انسان و کامپیوتر و همچنین پزشکی می‌تواند داشته باشد، توجه بسیاری از افراد فعال در علوم کامپیوتر را به خود جلب کرده است. [۳] در این حوزه، با استفاده از اطلاعات استخراج شده از ویدیو، کنش صورت گرفته در آن ویدیو شناسایی می‌شود. در حقیقت می‌توان گفت که این حوزه، تعمیمی بر مساله‌ی دسته‌بندی تصاویر<sup>۱</sup> است که در آن اشیای موجود در یک تصویر شناسایی شده و به دسته‌ی خاص خودشان نسبت داده می‌شوند. همان‌گونه که در حوزه‌ی دسته‌بندی تصاویر، مساله‌ی چگونگی نمایش تصویر مطرح است، در حوزه‌ی بازشناسی کنش انسان نیز یکی از مسائل اساسی نحوه‌ی نمایش اطلاعات و به زبان دقیق‌تر، نحوه‌ی نمایش حرکت بدن انسان موجود در ویدیو است. دو راه حل مهم برای رفع این مشکل به صورت بهینه، نمایش اطلاعات RGB-D و اطلاعات سه‌بعدی اسکلت انسان هستند. [۲][۱] تفاوت این دو روش در شکل ۱-۱ به وضوح به تصویر کشیده شده است.

نمایش اطلاعات در قالب مختصات سه‌بعدی اسکلت بدن، روشی است که در این پروژه مورد استفاده قرار می‌گیرد. بازشناسی کنش انسان با این روش نمایش، مدت زمان زیادی است که در حوزه‌ی بینایی رایانه‌ای مورد کند و کاو قرار گرفته است. الگوریتم‌های قدیمی‌تر و مبتنی بر روش‌های دست‌ساز

---

<sup>۱</sup>Image Classification



شکل ۱-۱: الف- نمایش اسکلت‌های بدن که از تصویر استخراج شده‌اند، ب- تصویر ژرفا<sup>۲</sup>، ج- تصویر رنگی [۴]

بیش‌تر از یک‌سری قوانین و روش‌های نسبتاً ثابت و انعطاف‌ناپذیر استفاده می‌کردند. به همین دلیل میزان خطای آن‌ها بالا بود و برای برخی موارد خاص و پیچیده به هیچ وجه قابل پیاده‌سازی نبودند. [۵] با رشد روزافزون یادگیری ژرف<sup>۳</sup> و همچنین افزایش اطلاعات در دسترس، شبکه‌های عصبی سرتاسر<sup>۴</sup> روش‌های جدیدتر و بهتری برای مسائلی هم‌چون بازشناسی کنش فراهم آوردند. در این‌گونه شبکه‌ها ورودی به لایه‌ی ابتدایی شبکه داده شده و خروجی از انتهای آن دریافت می‌گردد. هم‌چنین با یک معیار مناسب و مقایسه‌ی خروجی با این معیار، پارامترهای برنامه (توسط الگوریتم‌های بهینه‌سازی) به‌گونه‌ای تغییر می‌کنند که خروجی به معیار نزدیک و نزدیک‌تر شود. به این الگوریتم، الگوریتم یادگیری گفته می‌شود. در عین این‌که این الگوریتم به‌صورت کلی توضیح داده شد، معماری‌های شبکه‌ی<sup>۵</sup> بسیار متنوعی وجود دارند که برای حل مسائل مختلف مورد استفاده قرار می‌گیرند.

یکی از معماری‌های بسیار پرکاربرد، شبکه‌های عصبی پیچشی<sup>۶</sup> هستند. این شبکه‌ها بیش‌تر در مسائل مربوط به پردازش تصویر استفاده می‌شوند. بسیاری از حوزه‌های معروف پردازش تصویر، از قبیل آشکارسازی اشیا<sup>۷</sup>، دسته‌بندی تصاویر و ... با استفاده از این معماری روش حل بسیار بهینه‌تری پیدا کرده‌اند. شبکه‌های عصبی پیچشی- گرافی<sup>۸</sup> تعمیمی بر این معماری است. این مدل، به‌جای یک تصویر، گرافی را به‌عنوان ورودی گرفته و الگوریتم را بر روی آن اجرا می‌کند. این شبکه نیز کاربرد زیادی در مسائلی چون دسته‌بندی تصاویر، پردازش متن و ... دارد. [۶] شبکه‌ی استفاده شده در این پروژه،

<sup>۳</sup> Deep Learning

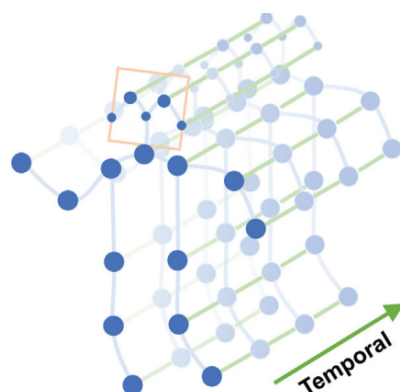
<sup>۴</sup> End to End Neural Networks

<sup>۵</sup> Network Architecture

<sup>۶</sup> Convolutional Neural Networks

<sup>۷</sup> Object Detection

<sup>۸</sup> Graph Convolutional Neural Network



شکل ۱-۲: نمایشی از ابعاد زمان-مکانی اسکلت انسان که ST-GCN بر روی آن کار می‌کند. [۵]

شبکه‌ی عصبی پیچشی-گرافی زمان-مکانی<sup>۹</sup> است که به اختصار ST-GCN نام دارد. شکل ۱-۲ (که نشان‌دهنده‌ی ورودی ST-GCN است) ایده‌ی کلی و مختصری از چگونگی کارکرد این شبکه را نمایش می‌دهد. توضیحات هر کدام از این معماری‌ها در ادامه‌ی پروژه به تفصیل آمده‌اند.

## ۲-۱ اهمیت موضوع

همان‌گونه که قبلاً اشاره شد، بازشناسی کنش انسان یکی از پرکاربردترین مباحث در حوزه‌ی بینایی ماشین است. کاربردهای این حوزه از مسائلی هم‌چون سرگرمی تا موارد پزشکی متغیر است. می‌توان ادعا کرد که بازشناسی کنش انسان، هدف اصلی سیستم‌های هوشمند ویدیویی<sup>۱۰</sup> است. [۱] در حوزه‌هایی هم‌چون تعامل انسان با ماشین، با بازشناسی کنش انسان، می‌توان عکس‌العملی در خور عمل صورت گرفته انجام داد. در موارد پزشکی، با شناخت دقیق کنش، امکان فیزیوتراپی برای بیماران با مشکلات جسمانی وجود دارد. در موارد امنیتی<sup>۱۱</sup>، می‌توان بدون این‌که ناظر انسانی وجود داشته باشد، با اکتفا بر کامپیوترها نظارت ویدیویی لازم را اعمال کرد.

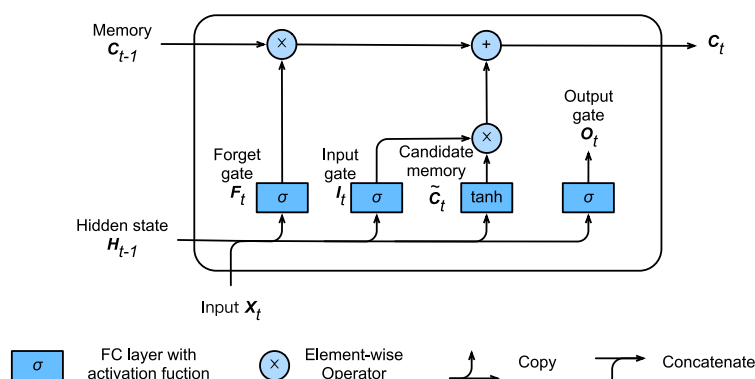
با توجه به این کاربردهای گسترده، وجود روشی برای بهینه‌سازی بازشناسی کنش انسان از ملزومات این بحث تلقی می‌شود. بسیاری از سیستم‌هایی که نیازمند به داشتن ویژگی بازشناسی کنش هستند، بایستی بصورت بی‌درنگ<sup>۱۲</sup> عمل کنند. به همین دلیل زمان موجود برای محاسبات و نیز ضریب خطا،

<sup>۹</sup>Spatial Temporal Graph Convolutional Neural Network

<sup>۱۰</sup>Intelligent Video Systems

<sup>۱۱</sup>Surveillance

<sup>۱۲</sup>Real Time



شکل ۱-۳: نمای کلی از یک LSTM [۹]

تا جای ممکن بایستی کاهش یابد. روش مورد استفاده در این پروژه، داده‌های کم‌تری نسبت به سایر روش‌ها استفاده می‌کند و به همین دلیل علاوه بر انعطاف پذیری بالا، سرعت و اطمینان بالایی را نیز تامین می‌کند.

## ۳-۱ ادبیات موضوع

یکی از روش‌های قدیمی برای مسائل در حوزه‌ی پردازش تصویر مانند بازشناسی کنش، استفاده از روش‌های مبتنی بر ویژگی‌های دست‌ساز<sup>۱۳</sup> است. [۷][۵][۸] در این روش، با کمک آشکارسازها برخی از نواحی تصویر را شناسایی کرده و با استفاده از این نواحی دسته کنش صورت گرفته را شناسایی می‌کنند. [۷] برای داده‌های ورودی کم، این روش به‌خوبی و با درصد موفقیت بالاتری عمل می‌کند. هرچند در صورتی که داده‌ی ورودی به فراوانی در دسترس باشد، استفاده از شبکه‌های ژرف بهینه‌ترین راه‌حل موجود است.

یکی از روش‌های معروف برای بازشناسی کنش انسان، استفاده از شبکه‌های عصبی بازگشتی<sup>۱۴</sup> است که به اختصار RNN نامیده می‌شوند. [۱۰] در بسیاری از مقالات از مدل تعمیم‌یافته‌ی این شبکه‌ی عصبی، که به حافظه‌ی کوتاه مدت بلند<sup>۱۵</sup> یا LSTM معروف است استفاده می‌کنند. [۱۰][۱۱] تصویر ۱-۳ نمایی از یک سلول LSTM را نشان می‌دهد. در این شبکه‌ی عصبی، ورودی‌ها، که قاب‌های<sup>۱۶</sup>

<sup>۱۳</sup>Handcrafted Features

<sup>۱۴</sup>Recurrent Neural Network

<sup>۱۵</sup>Long Short-Term Memory

<sup>۱۶</sup>Frames

جدول ۱-۱: نتایج برخی آزمایش‌ها بر روی پایگاه داده‌ی NTU RGB+D [۱۱]

روش	CS	CV
Skeletal Quads	٪۳۸/۶	٪۴۱/۴
Lie Group	٪۵۰/۱	٪۵۲/۸
Dynamic Skeletons	٪۶۵/۲	٪۶۰/۲
HBRNN	٪۵۹/۱	٪۶۴
Deep RNN	٪۵۶/۳	٪۶۴/۱
Deep LSTM	٪۶۰/۷	٪۶۷/۳
Part-aware LSTM	٪۶۲/۹	٪۷۰/۳
JTM CNN	٪۷۳/۴	٪۷۵/۲
SkeletonNet	٪۷۵/۹	٪۸۱/۲
Visualization CNN	٪۷۶	٪۸۲/۶

ویدیو هستند، به صورت سری داده می‌شوند. شبکه تعدادی دروازه<sup>۱۷</sup> دارد که با یادگیری (بهبود) آن‌ها به مرور زمان می‌تواند به یاد داشته باشد که کدام قاب‌ها اطلاعات بیشتری در اختیار شبکه می‌گذارند تا بر روی آن‌ها تمرکز بیشتری بگذارد و آن قاب‌ها را در طول زمان بیشتر به یاد داشته باشد و سایر اطلاعات را فراموش کند. [۱۲]

در برخی از مقالات با تغییراتی بر روی این شبکه نتایجی حاصل شده است که به برخی از آن‌ها در جدول ۱-۱ ذکر شده است.

## ۴-۱ چالش‌ها

برخی از چالش‌های موجود در بازشناسی کنش از روی داده‌های اسکلتی به شرح زیر می‌باشند:

- بزرگ‌ترین مساله (که پیش‌روی هرگونه شبکه‌ی سرتاسر قرار دارد) وجود مجموعه داده‌ی مناسب برای تضمین عملکرد بهینه است. همان‌گونه که اشاره شد، برتری اصلی این شبکه‌ها نسبت به

<sup>۱۷</sup>Gate

ویژگی‌های دست‌ساز، زمانی حاصل می‌شود که داده‌ی کافی در اختیار شبکه باشد.

- ورودی هرگونه شبکه‌ی بازشناسی کنش، دنباله‌ای از قاب‌ها در حوزه‌ی زمان است. این موضوع (بخصوص برای مجموعه‌داده‌های با وضوح بالاتر) باعث افزایش قدرت پردازشی و حافظه‌ی موردنیاز می‌شود. علاوه بر اندازه‌ی مجموعه‌ی داده، اندازه و تعداد پارامترهای شبکه از مواردی است که پیچیدگی زمانی به روش حل وارد می‌کنند.
- ناهنجاری و پیچیدگی در مجموعه‌داده‌های موجود، سرعت آموزش در شبکه‌ی عصبی را کاهش می‌دهد. به‌عنوان مثال زوایای متفاوت برای ویدیوهای مختلف، تعداد افراد حاضر در یک ویدیو، سرعت متفاوت انجام کنش توسط افراد مختلف و ... از جمله مواردی هستند که کیفیت شبکه را کاهش می‌دهند.

## ۵-۱ فرضیات

برای روشی که در این پروژه انتخاب شده است، برخی فرضیات از ابتدا در نظر گرفته شده است.

- با استفاده از مجموعه داده‌ی NTU RGB+D [۱۳]، داده‌های اسکلتی آماده هستند و نیازی به استفاده از الگوریتم‌هایی مانند تخمین حالت<sup>۱۸</sup>، برای استخراج این داده‌ها نیست. [۱۳][۵]
- در مجموعه داده‌ی اشاره‌شده، تنها یک کنش صورت می‌گیرد و اگر در قابی بیش از یک کنشگر موجود باشد، کنش (تعامل) بین این دو (و نه به صورت جداگانه) انجام خواهد شد.
- در این مجموعه داده، در هر قاب حداکثر دو کنشگر موجود است و پس‌زمینه نیز به‌دلیل استفاده از سنسورهای کینکت<sup>۱۹</sup> حذف شده‌اند.

## ۶-۱ اهداف تحقیق

در این پایان‌نامه، سعی شده است که بازشناسی کنش با استفاده از شبکه‌های ST-GCN و لایه‌های توجه انجام گیرد. برای یادگیری، از مجموعه‌داده‌ی NTU RGB+D [۱۳] استفاده می‌شود که بزرگترین

<sup>۱۸</sup>Pose Estimation

<sup>۱۹</sup>Kinect Sensors

مجموعه داده‌ی شامل اطلاعات سه بعدی اسکلت بدن هستند. این مجموعه داده از دو سنجه<sup>۲۰</sup> تشکیل یافته است که جزئیات هرکدام در انتهای پایان نامه تشریح خواهد شد. برای سنجش خروجی نیز از معیار تابع هزینه‌ی آنتروپی متقابل<sup>۲۱</sup> استفاده می‌شود.

## ۷-۱ ساختار پایان نامه

این پایان نامه شامل پنج فصل است. فصل دوم دربرگیرنده‌ی ادبیات مربوطه با پایان نامه است. در فصل سوم روش‌های پیاده سازی شده در این پروژه به تفصیل بیان گردیده است. فصل چهارم شامل نتایج تجربی به دست آمده از آزمودن روش پیشنهادی و مقایسه این نتایج با نتایج برخی روش‌های قبلی که روی مجموعه داده‌ی NTU RGB+D پیاده سازی و آزمون شده اند، است. بالاخره، جمع بندی کلی و راهکارهای ممکن برای ادامه‌ی این پروژه در فصل پنجم آورده شده است.

---

<sup>۲۰</sup>Benchmark

<sup>۲۱</sup>Cross Entropy Loss Function



## فصل ۲

# ادبیات مربوطه

### ۱-۲ شبکه‌های عصبی پیچشی

«پیچش» یک عملگر خطی است که برای توابع  $n$  بعدی تعریف می‌شود. مقدار آن برای توابع تک متغیره از فرمول

$$f(t) * g(t) = \int_{-\infty}^{+\infty} f(\tau) \cdot g(t - \tau) d\tau \quad (1-2)$$

محاسبه می‌شود. [۱۴] همین فرمول با اندکی تغییر برای توابع با دو متغیر (ماتریس‌های دوبعدی) مانند تصاویر، به شکل

$$x[m, n] * k[m, n] = \sum_{j=-\infty}^{+\infty} \sum_{i=-\infty}^{+\infty} k[i, j] \cdot x[m - i, n - j] \quad (2-2)$$

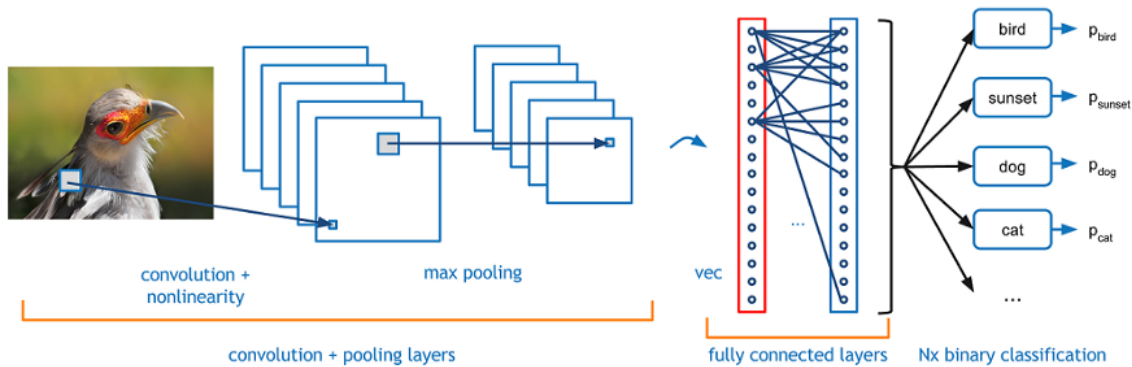
تعریف می‌شود. [۱۵]

در فرمول ۲-۲ اگر  $x[m, n]$  تصویر ورودی باشد، به عمل‌وند  $k[m, n]$  هسته<sup>۱</sup> گفته می‌شود. هسته مربعی به ابعاد  $3 \times 3$ ،  $5 \times 5$  یا ... است که با انتخاب مناسب محتوا و ابعاد آن، تصویر خروجی می‌تواند دارای ویژگی‌های خاص مربوط به تصویر اصلی باشد. مثلاً با انتخاب هسته‌ی خاص، می‌توان لبه‌ها را در تصویر شناسایی کرد.

S

---

<sup>۱</sup>Kernel



شکل ۲-۱: نمای کلی از یک شبکه‌ی پیچشی [۱۶]

ایده‌ی اصلی در CNN این است که مقادیر هسته ثابت در نظر گرفته نشود و به‌عنوان پارامترهای شبکه در هر مرتبه<sup>۲</sup> به‌روزرسانی شوند. [۱۷] به‌عنوان مثال، پیچش<sup>۳</sup> تصویر با ابعاد  $1920 \times 2180$  با هسته‌ای به ابعاد  $3 \times 3$  را در نظر بگیرید. این شبکه‌ی ساده و تک‌لایه ۹ پارامتر دارد که می‌توان با تعریف تابع هزینه مناسب، برای رسیدن به خروجی مطلوب، مدام این ۹ پارامتر را تغییر داد. شکل ۲-۱ نمای کلی از یک CNN نشان می‌دهد.

یک CNN اغلب از سه دسته لایه تشکیل می‌شود. لایه‌های پیچشی<sup>۴</sup>، لایه‌های ادغام<sup>۵</sup> و لایه‌های کاملاً متصل<sup>۶</sup>. توضیح هر کدام از این لایه‌ها به اختصار آمده است.

ایده‌ی کلی پشت لایه‌های پیچشی، همان ایده‌ی اصلی CNN است. به ازای هر لایه در این دسته از لایه‌ها، یک یا چند هسته وجود دارد که تصویر را فیلتر می‌کنند. در نتیجه‌ی این فیلتر برخی از ویژگی‌های تصویر استخراج می‌شود و تصویر کاهش اندازه می‌دهد. هرچند اگر تصویر از چندین کانال<sup>۷</sup> تشکیل شده باشد، تعداد این کانال‌ها رفته رفته بیشتر خواهد شد. در ابتدا اکثر تصاویر شامل سه کانال قرمز، سبز و آبی<sup>۸</sup> هستند.

برخی از لایه‌ها عمل «ادغام»<sup>۹</sup> را نیز بر روی تصویر انجام می‌دهند. ادغام انواع مختلفی همانند

<sup>۲</sup>Iteration

<sup>۳</sup>Convolve

<sup>۴</sup>Convolutional Layers

<sup>۵</sup>Pooling Layers

<sup>۶</sup>Fully Connected Layers

<sup>۷</sup>Channel

<sup>۸</sup>RGB Channels

<sup>۹</sup>Pooling

«ادغام بیشینه»<sup>۱۰</sup>، «ادغام میانگین»<sup>۱۱</sup> و ... دارد. به عنوان مثال اگر بخواهیم با یک فیلتر به اندازه  $3 \times 3$ ، تصویری به اندازه  $H \times W$  را ادغام بیشینه کنیم، بایستی از گوشه‌ی سمت چپ تصویر شروع کرده و فیلتر را بر روی تصویر بگذاریم. بیشینه مقدار پیکسل‌هایی از تصویر که زیر فیلتر قرار گرفته‌اند، پیکسل اول خروجی خواهد بود. سپس فیلتر را یک پیکسل به راست انتقال می‌دهیم و ... دقت شود که محتویات هسته در عمل ادغام اهمیتی ندارد. به همین دلیل در ادغام کردن، پارامتری برای یادگیری وجود نخواهد داشت.

لایه‌های کاملاً متصل بعد از لایه‌های پیچشی و لایه‌های ادغام در CNN قرار می‌گیرند. بعد از چندین لایه‌ی پیچشی و ادغام، تمامی پیکسل‌های خروجی را به یک شبکه‌ی کاملاً متصل (مانند شبکه‌ی رگرسیون خطی<sup>۱۲</sup>) می‌دهند تا خروجی نهایی بعد از چندین لایه‌ی کاملاً متصل به دست آید. الگوریتم یادگیری برای این لایه‌ها مانند الگوریتم‌های معمول برای شبکه‌های عصبی ساده‌ای چون رگرسیون خطی است.

## ۲-۲ شبکه‌های عصبی گراف-پیچشی

امروزه بسیاری از مجموعه داده‌های موجود، مثل اطلاعات شبکه‌های اجتماعی، شبکه‌ی اینترنت و ... به شکل گراف هستند. [۶] همان‌گونه که پیش از این ذکر شد، شبکه‌های عصبی گراف-پیچشی، که به اختصار GCN نامیده می‌شوند، تعمیمی بر CNN هستند که در آن ورودی به جای تصویر، یک گراف است. هرچند تفاوت‌ها در همین‌جا به انتها نمی‌رسد و پیچیدگی‌های ساختاری یک گراف، دشواری‌های خاصی را به این شبکه‌ها تحمیل کرده است. شکل ۲-۲، نمایی کلی از یک GCN را نمایش می‌دهد. به صورت دقیق هر GCN دو ماتریس را به عنوان ورودی دریافت می‌کند.

- ماتریس  $X$  به ابعاد  $N \times F$  که  $N$  تعداد رئوس و  $F$  تعداد ویژگی‌های ورودی<sup>۱۳</sup> برای هر راس است.

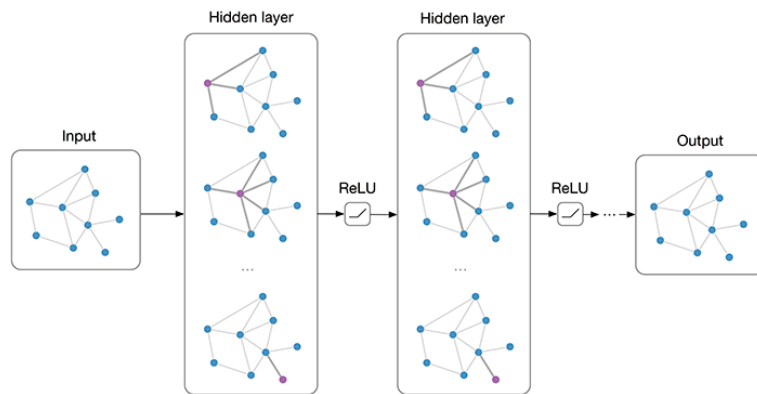
- ماتریس مجاورت  $A$  به ابعاد  $N \times N$  که ساختار کلی گراف را مشخص می‌کند.

<sup>۱۰</sup>Max Pooling

<sup>۱۱</sup>Average Pooling

<sup>۱۲</sup>Linear Regression

<sup>۱۳</sup>Input Features



شکل ۲-۲: نمای کلی از یک GCN [۱۸]

حال با تعریف وزن (هسته‌ی) مناسب برای هر لایه از شبکه، می‌توان شبکه را آموزش<sup>۱۴</sup> داد. هرچند هنوز هم برخی از مشکلات برای شبکه‌های نسبتاً بزرگ وجود دارد. مشکلاتی که در این پروژه به آن‌ها برخورد شد و راه‌حل پیشنهادی به شرح زیر هستند.

- همان‌گونه که ذکر شد، یکی از ورودی‌های شبکه، ماتریس مجاورت گراف است. در صورتی که یک راس یال بازگشتی به خودش نداشته باشد، درایه‌ی نظیر آن راس در ماتریس مجاورت صفر خواهد بود. همین موضوع باعث می‌شود که در مسیر لایه‌های شبکه، تنها ویژگی‌های ورودی رئوس مجاور آن راس در درایه‌ی نظیر آن وجود داشته باشند. به عبارت دیگر بعد از طی یک لایه، ویژگی‌های رئوس بدون یال بازگشتی تقریباً فراموش خواهند شد. به همین دلیل، قبل از هرکاری به تمامی رئوس گراف یک یال بازگشتی اضافه می‌شود. این کار با جمع کردن ماتریس مجاورت با ماتریس واحد ( $I$ ) انجام می‌گیرد.

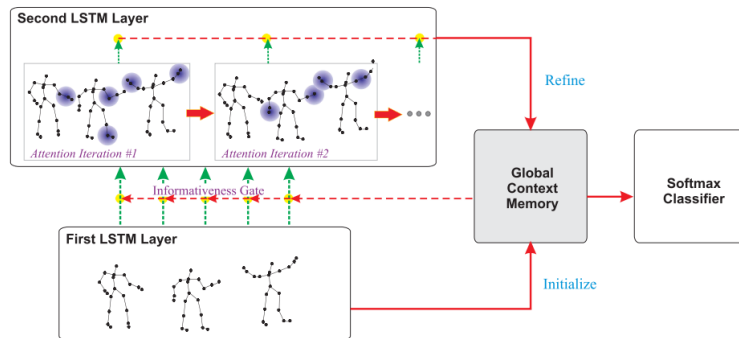
- رئوس با درجه‌ی بالاتر (پایین‌تر) رفته رفته اندازه‌ی بزرگ‌تری (کوچک‌تری) خواهند داشت. به این موضوع انفجارگرادیان<sup>۱۵</sup> و میرایی گرادیان<sup>۱۶</sup> گفته می‌شود. به همین دلیل عادی‌سای<sup>۱۷</sup> ورودی شبکه یک موضوع اجتناب‌ناپذیر است. برای این کار از روش موجود در [۶] استفاده شده است.

<sup>۱۴</sup>Train

<sup>۱۵</sup>Exploding Gradient

<sup>۱۶</sup>Vanishing Gradient

<sup>۱۷</sup>Normalization



شکل ۲-۳: مدل توجه با استفاده از حافظه‌ی زمینه‌ی سراسری در یک LSTM [۱۱]

## ۳-۲ مدل‌های توجه

معماری‌ها و مدل‌هایی که پیش‌تر معرفی شد، تفاوتی بین نقاط مختلف یک تصویر یا ویدیو قائل نبودند. درحالی که برای انجام یک الگوریتم بر روی تصویر، برخی جزئیات نه‌تنها مهم نیستند، بلکه در نتیجه‌ی نهایی خلل ایجاد می‌کنند. [۱۱] به‌عنوان مثال در فرآیند تشخیص تصویر، پس‌زمینه‌ی شی موردنظر اهمیتی ندارد. هم‌چنین در مثال خاص این پروژه، وقتی که کنش صورت‌گرفته دست‌زدن است، مفاصل پای یک شخص اهمیت چندانی ندارد. به همین دلیل، برای بهینه‌سازی بیش‌تر و درصد خطای پایین‌تر، بهتر است که شبکه رفته رفته متوجه شود که به کدام یک از جزئیات تصویر یا ویدیو بیش‌تر از باقی اجزا اهمیت قائل شود. یکی از مدل‌های توجه استفاده‌شده، مدل موجود در [۱۱] است که در معماری LSTM به کار گرفته شده است. در این مدل علاوه بر دروازه‌های موجود معماری LSTM، یک حافظه‌ی زمینه‌ی سراسری<sup>۱۸</sup> هم اضافه شده است. هم‌چنین از دو لایه LSTM استفاده شده است که لایه‌ی اول این حافظه را مقداردهی اولیه می‌کند و لایه‌ی دوم آن را بهبود می‌بخشد. در نهایت مقدار این حافظه است که به دسته‌بند پیشینه‌ی هموار<sup>۱۹</sup> داده می‌شود تا خروجی مورد نظر حاصل شود. طریقه‌ی مقداردهی و بهبودبخشی به حافظه، روشی مشابه دروازه‌های معمول LSTM دارد. تصویر ۳-۲ استفاده از این مدل توجه را نمایش می‌دهد.

روش دیگر به‌کارگیری مدل توجه (که در این پروژه هم از آن استفاده شده است) استفاده از پارامترهای قابل آموزش به‌ازای هر مفصل است. این مدل مناسب شبکه‌های CNN یا مشتقات آن مانند GCN

<sup>۱۸</sup>Global Context Memory

<sup>۱۹</sup>Softmax Classifier

است. [۵] در این مدل، از پارامترهای منسوب به وزن‌های اهمیت مفصل<sup>۲۰</sup> استفاده می‌شود که به هر مفصل یک وزن مشخص می‌دهد. این وزن بعد از آموزش کل شبکه مقدار بهینه پیدا می‌کند. سپس در هنگام ارزیابی، این وزن بر روی هسته‌ی موجود در شبکه‌ی GCN ضرب می‌شود تا هر یال تاثیر مشخصی بر روی جواب نهایی شبکه داشته باشد.

روش دیگری که در حوزه‌ی پردازش تصویر بسیار جدید است، استفاده از لایه‌ی خاصی به اسم لایه‌ی ادغام توجه<sup>۲۱</sup> است. اولین استفاده از این روش در بازشناسی کنش انسان در [۱۹] صورت گرفته است. در این روش، لایه‌های کاملاً متصل از لایه‌های پیچشی تغذیه نمی‌شوند. بلکه قبل از لایه‌های کاملاً متصل، خروجی‌های لایه‌های پیچشی به لایه‌های ادغام توجه داده می‌شود. در این گونه از لایه‌های ادغام، کل ورودی تحت ادغام قرار می‌گیرد و خروجی آن به لایه‌های کاملاً متصل یا به یک لایه‌ی دسته‌بند پیشینه‌ی هموار داده می‌شود. (به تفاوت این نوع ادغام با گونه‌های قبلاً معرفی شده توجه شود که در گونه‌های قبلی، مربع‌های  $k_i \times k_i$  از گوشه‌ی سمت بالا-چپ تصویر انتخاب شده و تحت ادغام قرار می‌گیرند). جزئیات این روش در فصل آتی توضیح داده خواهند شد.

---

<sup>۲۰</sup> Edge Importance Weighting

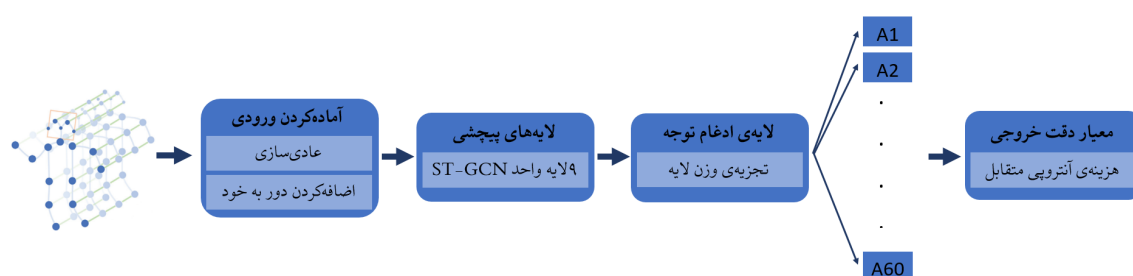
<sup>۲۱</sup> Attention Pooling Layer

## فصل ۳

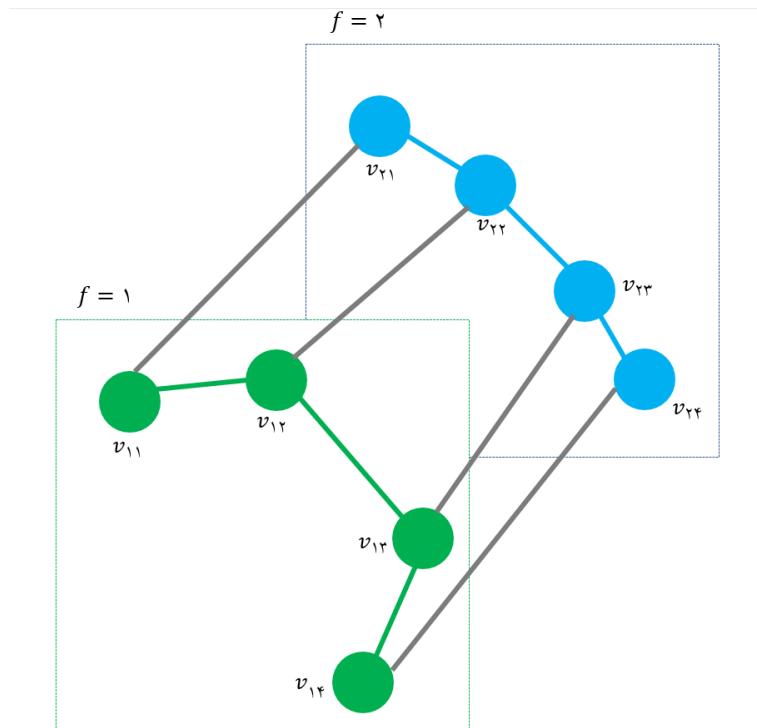
# روش‌های پیشنهادی

### ۱-۳ مقدمه

این پروژه ادامه‌ای است بر کار [۵] و [۱۹]. هردوی این مقالات، با استفاده از شبکه‌های عصبی سرتاسر، به بازشناسی کنش پرداخته‌اند. در این پایان‌نامه نیز سعی شده است که با استفاده از روش‌های هردوی این مقاله‌ها، درصد خطای پایین‌تری را به دست آورد. در ابتدای این فصل کلیت شبکه و برخی ریزه‌کاری‌ها که بایستی اعمال شوند توضیح داده می‌شود. در ادامه مدل توجه مورد استفاده و چگونگی کار آن شرح داده خواهد شد. در انتها نیز معیاری که برای سنجش کار استفاده خواهیم کرد معرفی می‌شود. تصویر ۱-۳ نموداری مختصر از کل کار انجام‌گرفته را شرح می‌دهد.



شکل ۱-۳: نمودار بلوکی برای شبکه‌ی ST-GCN با مدل توجه ادغامی



شکل ۳-۲: مثالی از یک گراف زمان-مکانی

### ۳-۲ شبکه‌ی استفاده شده

شبکه‌ی مورد استفاده در این پروژه، شبکه‌ی گراف-پیچشی زمان-مکانی (ST-GCN) است. بزرگ‌ترین تفاوت این شبکه با شبکه‌ی گراف-پیچشی، در نوع ورودی‌ای است که به آن داده می‌شود. ورودی  $st$ -gcn یک گراف  $G = (V, E)$  است که در آن  $V$  مجموعه‌ی رئوس و  $E$  مجموعه‌ی یال‌های گراف است. مجموعه‌ی  $V$  تعریفی مانند

$$V = \{v_{fi} | f = 1, \dots, F, i = 1, \dots, N\} \quad (۱-۳)$$

دارد. [۵]

در رابطه‌ی ۱-۳،  $f$  شاخص شماره‌ی قاب<sup>۱</sup> و  $i$  شاخص شماره‌ی راس (مفصل) در یک قاب است. همان‌گونه که از این رابطه مشخص است، رئوس گراف تمامی مفاصل در تمامی قاب‌ها را شامل می‌شود و محدود به یک قاب نیست. به همین دلیل تعریف یال نیز بایستی شامل تمامی ارتباطات متصل‌کننده این رئوس باشد. [۵] مجموعه‌ی یال‌ها را به دو زیرمجموعه‌ی یال‌های یک قاب ( $E_s$ ) و یال‌های بین

<sup>۱</sup>Frame



دو قاب ( $E_t$ ) تفکیک کرده و هرکدام را به شکل

$$E_s = \{v_{fi}v_{fj} | (i, j) \in H\}, \quad (2-3)$$

$$E_t = \{v_{fi}v_{(f+1)i}\}$$

تعریف می‌کند. [۵]

در رابطه‌ی ۲-۳ مجموعه‌ی  $H$  مجموعه‌ی شامل مفاصل مجاور بدن انسان هستند. توجه کنید که چگونه مجموعه یال‌های  $E_t$  مفاصل متناظر در دو قاب را به هم‌دیگر متصل می‌کند.

برای انجام عمل پیچش در این گراف، بایستی رابطه‌ی همسایگی را برای هر راس در گراف تعریف کنیم. چرا که پارامتر وزن در عمل پیچش، به‌ازای هر راس، بر روی همسایه‌های آن راس شناور خواهد بود. برای هر راس مانند  $v_{fi}$  همسایه‌های آن با رابطه‌ی

$$N(v_{fi}) = \{v_{qj} | d(v_{fj}, d_{fi}) \leq K, |q - f| \leq \left\lfloor \frac{\Gamma}{2} \right\rfloor\} \quad (3-3)$$

تعریف شده‌اند. [۵]

در رابطه‌ی ۳-۳، تابع  $d$  کوتاه‌ترین مسیر بین دو راس ورودی آن را مشخص می‌کند. هم‌چنین متغیر  $K$  حداکثر فاصله بین رئوس همسایه در یک قاب و متغیر  $\Gamma$  حداکثر فاصله بین دو راس همسایه در دو قاب مختلف را بیان می‌کند. به بیان دیگر، می‌توان گفت که  $K$  اندازه‌ی ماتریس وزن در بعد مکان و  $\Gamma$  اندازه‌ی آن در بعد زمان است. برای افزایش سرعت کار در این پروژه،  $K = 1$  و  $\Gamma = 2$  در نظر گرفته شده است. مقادیر بالاتر از این می‌تواند در کارهای آتی مورد بررسی قرار گیرد.

شکل ۲-۳ روابطی که تا این‌جای کار بیان شد را ترسیم کرده است. در این شکل، یال‌های یک قاب به‌صورت رنگی و یال‌های بین دو قاب بی‌رنگ هستند. هم‌چنین به‌ازای  $K = 1$  و  $\Gamma = 2$  رئوس  $v_{11}$ ،  $v_{12}$ ،  $v_{21}$  و  $v_{22}$  همسایه هستند. هم‌چنین دقت شود که رابطه‌ی همسایگی یک رابطه‌ی تعدی نیست.

### ۳-۳ مدل توجه

مدل توجه استفاده شده در این پروژه، الهام گرفته از مدل‌های [۱۹] و [۲۰] است. استفاده از این روش بر دو ایده‌ی کلی استوار است:

• اگر درست بعد از اتمام لایه‌های پیچشی، خروجی را به صورت یک بردار درآورده و به لایه‌ی کاملاً متصل بدهیم، پارامترهای شبکه بسیار زیاد شده و یادگیری را مشکل می‌کند. به همین دلیل بهتر است که قبل این کار، به گونه‌ای اندازه‌ی خروجی را کاهش دهیم و بعد از آن به یک لایه‌ی کاملاً متصل بدهیم. [۲۰]

• روشی که برای کاهش اندازه اتخاذ می‌کنیم بهتر است به ازای ورودی‌های مختلف، پاسخ‌های متفاوتی داشته باشد. چرا که همان گونه که در فصل پیش به آن اشاره شد، برای برخی از ورودی‌ها احتیاجی به کل داده نیست و بهتر است پارامترها بسته به نوع ورودی، ضرایب متفاوتی داشته باشند.

برای دستیابی به این دو مورد، [۱۹] و [۲۰] لایه‌ی ادغام توجه را معرفی کرده‌اند. بردار ساخته شده توسط این لایه با رابطه‌ی

$$score(X) = Tr(X^T X W) \quad (۴-۳)$$

محاسبه می‌شود. [۱۹]

در رابطه‌ی ۴-۳،  $X$  ورودی لایه‌ی ادغام به اندازه‌ی  $n \times f$  و  $W$  پارامتر این لایه به اندازه‌ی  $f \times f$  است. برای این که رابطه‌ی بالا قادر به جذب خاصیت توجه باشد، کافی است پارامتر  $W$  را به صورت ضرب دو بردار  $1 \times f$  بنویسیم.

$$W = ab^T \rightarrow W^T = ba^T \quad a, b \in R^{f \times 1} \quad (۵-۳)$$

$$score(X) = Tr(X^T X ba^T) \quad (۶-۳)$$

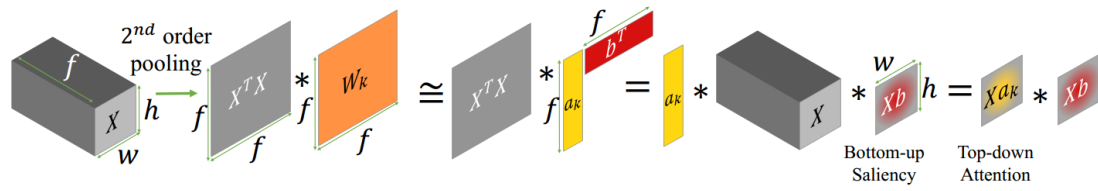
چون برای ماتریس‌ها داریم  $Tr(ABC) = Tr(CAB)$ ، در نتیجه رابطه‌ی ۶-۳ را می‌توان به شکل

$$score(X) = Tr(a^T X^T X b) \quad (۷-۳)$$

نوشت.

همچنین چون برای یک بردار مانند  $u$  داریم  $Tr(u) = u$  در نتیجه رابطه‌ی ۷-۳ را نیز می‌توان به شکل

$$\begin{aligned} score(X) &= a^T X^T X b \\ &= (Xa)^T (Xb) \end{aligned} \quad (۸-۳)$$



شکل ۳-۳: نموداری از مدل توجه با استفاده از لایه‌ی ادغام [۱۹]

نوشت. [۱۹]

حال بردار  $score$  آماده است تا به عنوان ورودی به لایه‌های کاملاً متصل و یا حتی به یک لایه‌ی دسته‌بند پیشینه‌ی هموار داده شود. تصویر ۳-۳ جزئیات این مدل را نمایش می‌دهد.

## فصل ۴

# نتایج تجربی

### ۴-۱ مجموعه داده‌ی مورد استفاده

مجموعه داده‌ی NTU-RGB+D بزرگ‌ترین مجموعه داده‌ی شامل اطلاعات ۳ بعدی اسکلت بدن است. [۵] تصویر ۴-۱ برخی از قاب‌های موجود در این مجموعه داده را نمایش می‌دهد. این مجموعه داده شامل ۵۶۰۰۰ کلیپ و ۶۰ دسته کنش است که از A1 (آشامیدن آب) تا A60 (از هم فاصله گرفتن) برچسب<sup>۱</sup> گذاری شده‌اند. این ویدیوها از ۳ زاویه‌ی مختلف و با استفاده از سنسورهای کینکت<sup>۲</sup> ضبط شده‌اند تا مختصات ۳ بعدی مفاصل به دست آیند. این مجموعه داده به صورت کلی به دو دسته سنجه<sup>۳</sup> تقسیم شده است. [۱۳]

۱. X-sub: در این سنجه، بازیگران برای مجموعه‌ی آموزش<sup>۴</sup> و مجموعه‌ی آزمون<sup>۵</sup> متفاوت هستند.

۲. X-view: در این سنجه، برای مجموعه‌ی آموزش از زاویه‌های ۱ و ۲ دوربین و برای مجموعه‌ی آزمون از زاویه‌ی ۳ دوربین استفاده شده است.

در این پروژه نیز از هردوی این سنجه‌ها به صورت جداگانه استفاده شده است.

---

<sup>۱</sup> Label

<sup>۲</sup> Kinect Sensors

<sup>۳</sup> Benchmark

<sup>۴</sup> Train Set

<sup>۵</sup> Test Set



شکل ۴-۱: برخی از قاب‌های مجموعه داده‌ی NTU RGB+D [۲۱]

## ۴-۲ معیار تابع هزینه‌ی آنتروپی متقابل

در اکثر مسائل دسته‌بندی، از این تابع به‌عنوان معیار درستی تخمین استفاده می‌کنند. [۲۲] [۵] در مسائل دسته‌بندی دودویی این تابع با رابطه‌ی

$$L(y, \hat{y}) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})) \quad (۴-۱)$$

تعریف می‌شود. [۲۲]

برای مسائل دسته‌بندی که تعداد دسته در آن‌ها از دو بیشتر است، این تابع از رابطه‌ی

$$L(y, \hat{y}) = - \sum_i y_i \log \hat{y}_i \quad (۲-۴)$$

محاسبه می‌شود. [۲۲]

در هر دوی روابط ۴-۱ و ۴-۲،  $y$  برچسب واقعی است که دو مقدار صفر یا یک را می‌گیرد و  $\hat{y}$  مقدار مشاهده‌شده است که هر مقداری بین صفر تا یک می‌تواند بگیرد.

## فصل ۵

# جمع‌بندی و راه‌کارهای آتی

در مطالعات آینده، می‌توان اندازه‌ی هسته‌ی استفاده‌شده (چه در حوزه‌ی مکان و چه در زمان) و هم‌چنین تعداد لایه‌های ST-GCN را افزایش داد و تاثیر هر کدام از این تغییرات را بر شبکه‌ی نهایی مشاهده کرد. هم‌چنین می‌توان به‌جای مدل توجه استفاده‌شده در این پروژه، از مدل‌های دیگری استفاده کرد به‌عنوان مثال با تغییراتی در لایه‌های پیچشی، آن‌ها را حساس به ورودی کرد تا بر روی ورودی‌های مختلف، پردازش‌های متفاوتی صورت گیرد.

## مراجع

- [1] G. Cheng, Y. Wan, A. N. Saudagar, K. Namuduri, and B. P. Buckles. Advances in human action recognition: A survey. *Dept. of Computer Science and Engineering, University of North Texas*, 2015.
- [2] F. Han, B. Reily, W. Hoff, and H. Zhang. Space-Time representation of people based on 3D skeletal data: A review. *Division of Computer Science, Colorado School of Mines, Golden, CO 80401, USA*, 2017.
- [3] Activity recognition. [https://en.wikipedia.org/wiki/Activity\\_recognition](https://en.wikipedia.org/wiki/Activity_recognition). Retrieved: 2019-07-29.
- [4] A. Ben Tamou, L. Ballihi, and D. ABOUTAJDINE. Automatic learning of articulated skeletons based on mean of 3D joints for efficient action recognition. *International Journal of Pattern Recognition and Artificial Intelligence*, 2016.
- [5] S. Yan, Y. Xiong, and D. Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. *AAAI Conference on Artificial Intelligence (AAAI-18)*, 2018.
- [6] T. N. Kipf and M. Welling. Semi-Supervised classification with graph convolutional networks. *ICLR*, 2017.
- [7] A. B. Sargano, P. Angelov, and Z. Habib. A comprehensive review on handcrafted and learning-based action representation approaches for human activity recognition. *Applied Sciences*, 2017.
- [8] M. A. Aghbolaghi. A robust and compressed descriptor for action recognition from 4d data”. Master’s thesis, Sharif University of Technology, 2018.

- [9] Dive into deep learning. [https://www.d2l.ai/chapter\\_recurrent-neural-networks/lstm.html](https://www.d2l.ai/chapter_recurrent-neural-networks/lstm.html). Retrieved: 2019-07-23.
- [10] J. Liuy, A. Shahroudy, D. Xuz, , and G. Wangy. Spatio-Temporal LSTM with trust gates for 3D human action recognition. *School of Electrical and Electronic Engineering, Nanyang Technological University*, 2016.
- [11] J. Liu, G. Wang, L. Duan, K. Abdiyeva, and A. C. Kot. Skeleton-Based human action recognition with global Context-Aware attention LSTM networks. *IEEE*.
- [12] S. Hochreiter and J. Schmidhuber. Long Short-Term memory. *Neural Computation*, 1997.
- [13] A. Shahroudy, J. Liu, T. Ng, and G. Wang. NTU RGB+D: A large scale dataset for 3D human activity analysis. *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [14] A. V. Oppenheim, A. S. Willsky, and S. H. Nawab. *Signals & Systems*. Prentice-Hall International, Inc., 1997.
- [15] S. Kim. Applications of convolution in image processing with MATLAB, 2013.
- [16] Convolutional neural network. <https://towardsdatascience.com/covolutional-neural-network-cb0883dd6529>. Retrieved: 2019-07-23.
- [17] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. Massachusetts Institute of Technology, 2016.
- [18] Graph convolutional networks. <https://tkipf.github.io/graph-convolutional-networks/>. Retrieved: 2019-07-25.
- [19] R. Girdhar and D. Ramanan. Attentional pooling for action recognition. *Conference on Neural Information Processing Systems (NIPS 2017)*, 2017.
- [20] P. Li, Y. Song, I. McLoughlin, W. Guo, and L. Dai. An attention pooling based representation learning method for speech emotion recognition. *Interspeech*, 2018.
- [21] Action recognition datasets: "NTU RGB+D" dataset and "NTU RGB+D 120" dataset. <http://rose1.ntu.edu.sg/datasets/actionrecognition.asp>. Retrieved: 2019-07-23.



- [22] Loss functions. [https://ml-cheatsheet.readthedocs.io/en/latest/loss\\_functions.html](https://ml-cheatsheet.readthedocs.io/en/latest/loss_functions.html). Retrieved: 2019-07-29.

# واژه‌نامه

## الف

ادغام ..... pooling  
اشیا ..... objects  
انتها به انتها ..... end to end  
انفجار ..... explosion  
آشکارسازی ..... detection  
آموزش ..... train

## ب

بازشناسی ..... recognition  
بلادرنگ ..... real time  
بیشینه‌ی هموار ..... softmax

## پ

پیچش ..... convolution  
پیچشی ..... convolutional

## ت

توجه ..... attention

## ح

حافظه‌ی کوتاه بلندمدت ..... long short-term memory

## د

دروازه ..... gate  
دسته‌بند ..... classifier  
دسته‌بندی ..... classification

## ر

رگرسیون خطی ..... linear regression

## ز

زمینه ..... context

## ژ

ژرفا ..... depth

action..... کنش	س
	global ..... سراسری
م	intelligent systems ..... سیستم‌های هوشمند
iteration ..... مرتبه	ش
architecture ..... معماری	شبکه‌های عصبی ..... neural networks
average ..... میانگین	شبکه‌های عصبی بازگشتی ..... recurrent neural networks
vanishing ..... میرایی	convolutional networks ..... شبکه‌های پیچشی
و	graph convolutional ..... شبکه‌های پیچشی گرافی
edge importance weights.. وزن‌های اهمیت مفصل	networks
feature ..... ویژگی	spatial ..... شبکه‌های پیچشی گرافی زمانی-مکانی
ی	tempooral graph convolutional networks
deep learning ..... یادگیری عمیق	ک
	fully connected..... کاملاً متصل
	channel ..... کانال

## **Abstract**

Recently, human action recognition have become one of the most studied topics in the world of computer science and engineering. With the enormous increasing of available datasets and the advent of neural networks, new horizons have been opened to the concept of human action recognition. Representation of human body data is one of the fundamentals of this concept. There are two major methods that's been introduced to optimize this representation. These two methods largely include using RGB-D data and using 3D skeleton information. Nowadays, lots of research have been taken place on representation of 3D skeleton information, due to its flexibility and lower data size.

In this project studies towards representation of 3D skeleton information and its use in human action recognition have been continued. Type of network used in this project is spatial temporal graph convolutional network which is a modification of graph convolutional network which is also a modification of convolutional networks. Also, by introducing a beneficial attention model and using filan criteria, current networks have been improved. Future works could introduce more efficient attention models and study their impact on these networks.

**Keywords:** Human Action Recognition, 3D Skeleton Data, Graph Convolutional Networks, Attention Model



Sharif University of Technology

Department of Computer Engineering

B.Sc. Thesis

**Skeleton Based Human Action Recognition Using  
Spatial Temporal Graph Convolutional Networks  
With an Attention Model**

By:

**Reza Rahimi Azghan**

Supervisor:

**Prof. Kasaei**

August 2019