

I HC QUC GIA THỊNH PH H CHỜ MINH
TRNG I HC BỒCH KHOA



CNG LUN VN
NHN DIN HỊNH NG QUA VIDEO

GVHD: Nguyn c Dng

SVTH: Trn Hu Thc 1713454

ng Khĩi 1711807

Tp. H Chờ Minh, 11/2020



Mục lục

1	Gii thiu	2
1.1	Bị toàn nhn đin hình ng qua video	2
1.2	Mc tiểu ca giai on cng lun vn	3
1.3	Cu trúc ca cng lun vn	4
2	Còc kho sòt liỏn quan n tị	4
2.1	Kho sòt chung	4
2.2	Kho sòt tởnh ng dng	5
2.3	Kho sòt còc thỏch thc	6
2.4	Kho sòt tp đ liu	6
2.5	Kho sòt còc phng phỏp truyủn thng	7
2.5.1	Phng phỏp biu đin đ liu toủn cc	7
2.5.2	Phng phỏp biu đin đ liu cc b	8
2.6	Kho sòt còc phng phỏp hc sỏu	8
3	Kin thc nn tng	9
3.1	Còc kin thc toủn c bn	9
3.1.1	Phỏp nhón tỏch chp ma trủn (convolution)	9
3.1.2	th (graph)	10
3.2	Còc kin thc c bn v Trờ tu nhón to, hc mỳy vị hc sỏu	11
3.2.1	Còc kin thc c bn	11
3.2.2	Phón loi còc phng phỏp hc	14
4	Phng phỏp xut	15
5	Kt qu bc u	15
5.1	X lý hỏnh nh RGB qua đ liu skeleton	15
5.2	X lý đ liu skeleton nhn đin hình ng	15
6	K hoch phỏt trin cho giai on Lun vn	15



1 Giới thiệu

1.1 Bối cảnh nhận diện hình ảnh qua video

Sự phát triển của nhúng lõi thị giác 4.0 đang ngày càng phát triển mạnh mẽ, tạo ra những nhu cầu thực tiễn mà các thị trường không chỉ, cũng như hình ảnh liên quan trong những nhu cầu quan trọng. Ngoài hình ảnh chất lượng sản xuất, khả năng tăng tốc cao vì con người đang lựa chọn cho sự ra đời của hàng loạt sản phẩm cũng như, các thị trường sản phẩm trở nên ngày càng (AI).

AI đã xuất hiện không lâu nhưng thay đổi chóng mặt về lý thuyết và thực tiễn ban đầu, thời nay các ứng dụng thực tiễn của AI đã được phân bổ cụ thể thay đổi cuộc sống con người. Ngày nay, AI xuất hiện hầu hết trong tất cả các lĩnh vực từ Giao thông, Y tế, Nông nghiệp cho tới Giải trí và rất nhiều lĩnh vực khác. Các ứng dụng trở nên ngày càng đa dạng thay thế công nghệ thông tin cũng như con người bằng sự tiến bộ của công nghệ, phần mềm.

Hiện nay, AI liên ngành khoa học và công nghệ vì mục tiêu giúp mức độ khả năng công nghệ thông minh như con người trong phạm vi nào đó. Tuy nhiên, công nghệ mức độ thực hiện cũng như việc tin cậy từ các ứng dụng, giúp làm giảm khi công nghệ cũng như cho con người. Tuy nhiên, thực tế khả năng của AI đã phần nào tìm kiếm khả năng thực sự của con người, thậm chí một số bối cảnh, AI vẫn đang ngày càng phát triển. Mục tiêu của nghiên cứu AI, là để khả năng xử lý mức độ cao hơn, mở rộng khả năng ứng dụng của các sản phẩm AI ra đời.

Các ứng dụng phân bổ của AI cụ thể liên quan đến vị trí và ra đời và phát triển trong Giám sát an ninh, hay xác nhận thông tin giao thông trong thị trường, chủ yếu liên quan đến bối cảnh xử lý thông tin. Ngoài bối cảnh xử lý thông tin, việc xử lý dữ liệu dạng chuỗi như video cũng mang lại ứng dụng quan trọng như nhận diện hình ảnh vì con người mức độ khả năng hiểu ý nghĩa của các hình ảnh và nắm bắt thông tin từ con người. Các ứng dụng nổi bật của việc nhận diện hình ảnh vì con người cụ thể liên quan đến nhận diện hình ảnh từ ảnh (mức tối, trung bình), nhận diện vị trí đối tượng hình ảnh vì con người trong thị trường (nhận diện con người bằng qua mạng cho xe tự hành), hay giúp mức độ tăng tốc về hình ảnh vì con người (như màn hình smart TV bằng cảm biến).

Dữ liệu dạng chuỗi ngày càng nhiều và phức tạp hơn so với xử lý thông tin thông thường, vì công nghệ thị trường, lượng dữ liệu lớn hơn, cần có những phương pháp xử lý thích hợp để xử lý. Tuy nhiên các phương pháp hiện nay về xử lý dữ liệu dạng chuỗi vẫn còn một số hạn chế, các phương pháp xử lý xác định chất lượng, hay cũng có những phương pháp xử lý xác định mức độ cao nhưng công nghệ xử lý chưa cao do chỉ phụ thuộc vào dữ liệu.

Tình hình ứng dụng các kỹ thuật nhận dạng khuôn mặt trong thực tiễn đang phát triển rất nhanh chóng. Tuy nhiên, việc ứng dụng các kỹ thuật nhận dạng khuôn mặt trong thực tiễn vẫn còn gặp nhiều khó khăn. Một trong những khó khăn đó là việc thiếu dữ liệu huấn luyện đa dạng, đặc biệt là dữ liệu đa dạng về môi trường, ánh sáng, góc quay, v.v. Do đó, việc nghiên cứu và phát triển các kỹ thuật nhận dạng khuôn mặt trong thực tiễn vẫn còn là một thách thức lớn.

1.2 Mục tiêu của giải pháp nhận dạng khuôn mặt

Trong giải pháp nhận dạng khuôn mặt này, mục tiêu của nhóm là kho sớt, nghiên cứu, phát triển hệ thống nhận dạng khuôn mặt qua video, cụ thể như sau:

- Thực hiện kho sớt những phương pháp nổi bật trong lĩnh vực nhận dạng khuôn mặt, tìm kiếm hình ảnh minh họa trực quan.
- Nghiên cứu ưu nhược điểm của từng phương pháp nhận dạng khuôn mặt để từ đó lựa chọn phương pháp phù hợp nhất để áp dụng vào bài toán.
- Tìm kiếm và phân tích các bài báo khoa học liên quan đến nhận dạng khuôn mặt, từ đó xác định hướng nghiên cứu.
- Chuẩn bị dữ liệu video để huấn luyện và kiểm tra mô hình nhận dạng khuôn mặt, đồng thời đánh giá hiệu suất của mô hình.
- Tìm kiếm và phân tích các bài báo khoa học liên quan đến nhận dạng khuôn mặt, từ đó xác định hướng nghiên cứu.

phối hợp với thí nghiệm thực địa để đánh giá hiệu suất của mô hình nhận dạng khuôn mặt, đồng thời đánh giá hiệu suất của mô hình.

- Sử dụng dữ liệu video để huấn luyện và kiểm tra mô hình nhận dạng khuôn mặt, đồng thời đánh giá hiệu suất của mô hình.
- Hình ảnh đầu vào được chuyển đổi thành dữ liệu skeleton data, sau đó sử dụng dữ liệu video RGB để huấn luyện và kiểm tra mô hình nhận dạng khuôn mặt, đồng thời đánh giá hiệu suất của mô hình.
- Sử dụng dữ liệu 1, 2 người để huấn luyện và kiểm tra mô hình nhận dạng khuôn mặt, đồng thời đánh giá hiệu suất của mô hình.

1.3 Cấu trúc của công luận văn

Công luận văn có thể chia thành các phần như sau:

- Chương 1: Giới thiệu tổng quan về bài toán Nhận diện hình ảnh qua video, công nghệ học sâu và ứng dụng, mục tiêu nghiên cứu và đóng góp của bài toán.
- Chương 2: Tổng quan về các phương pháp liên quan đến bài toán nhận diện hình ảnh qua video. Phân tích ưu nhược điểm của các phương pháp hiện có và đề xuất hướng nghiên cứu.
- Chương 3: Kiến thức nền tảng về học sâu và các kiến thức về học máy liên quan đến bài toán nghiên cứu.
- Chương 4: Phương pháp đề xuất và mô hình đề xuất để giải quyết bài toán nhận diện hình ảnh qua video.
- Chương 5: Kết quả thực nghiệm của nhóm công nghệ minh chứng cho tính khả thi của đề tài.
- Chương 6: Phân tích kết quả thực nghiệm, những ưu điểm và hạn chế của mô hình đề xuất. Đề xuất hướng nghiên cứu tiếp theo.
- Chương 7: Kết luận đề tài, xuất bản kết quả trong tương lai.

2 Tổng quan về các phương pháp liên quan đến đề tài

Trong chương này, sẽ trình bày tổng quan về bài toán Nhận diện hình ảnh qua video. Bao gồm: tổng quan về các phương pháp liên quan đến bài toán, các phương pháp truyền thống, các phương pháp học sâu và các ứng dụng của chúng. Phân tích ưu nhược điểm của các phương pháp hiện có và đề xuất hướng nghiên cứu.

2.1 Tổng quan chung

Nhận diện hình ảnh (action recognition) là công nghệ có thể chia thành hai loại: Phát hiện hình ảnh (action detection) và Phân loại hình ảnh (action classification).

- Phân loại hình ảnh là quá trình xác định hình ảnh trong video có chứa nội dung gì hay không. Đây là một bài toán phân loại. Các bài toán nhận diện hình ảnh thường có tính chất tổng quát và đa dạng.

- Phốt hình hình ng lị quò trờnh tòm ra khong thì gian bt u vị thì gian kt thữc ca hình ng trong mt video dị cha nhieu hình ng. óy lị bc quan trng ng dng vjò thc t, tuy nhĩn vn cha c quan tòm nhieu, vị vn cùn lị thòch thc ln cn gii quyt.

Cộc k thut nhn đin hình ng ngiy nay cú th c chia lrim 4 loi da trổn tồnh cht hình ng:

- Nhận diện còi hình ng c bn ca mt b phn trở c th nh vỵ tay, nhc chón, un cong ngi...
- Nhận diện hình ng ca nhĩu b phn phi hp nhau trở mt c th nh i b, nhĩ xa, m.
- Nhận diện hình ng cú s tng tcc gia ngi vị mt i tng khcc nh ònh ịn, cm dao....
- Nhận diện hình ng ca mt nhúm ngi nh biu tnh, hp nhúm ...

2.2 Kho sòt tòngh ng dng

Nhìn diện hình ảnh của con người có tầm ảnh hưởng cao vị nhân cách quan trọng minh chứng của công nghệ, có thể không còn cũng trở nên nhàm chán:

- H thng giòm sòt vị nhn đin bt thng trong mui trng nhĩ . Bị bỏ [?] ò trnh bị phng phòp n gin s dng MHI trỏch sut sut c trng th cũng vị s dng SVM phón loi t ú phòt hin hnh ng tổ ngõ. óy lĩ ng dng thit thc i vị ngĩ giĩ thng nhĩ mt mnh.
- H thng tòm kim, truy vn video theo hnh ng. David Doermann vị Daniel DeMenthon ò sut phng phòp chia video giòm sòt thnh còc video ngn cú đĩi na giỏy, t ú gòn nhõn lu tr di cu d liu dng cóy phc v truy vn [?]. Vở d cho tòngh ng dng ca h thng lĩ cú th giũp ta phòt hin n cp ti còc vn phùng.
- H thng phòt hin li ca vn ng viỏn cho còc cuc thi th thao. in hnh ca ng dng nĩy lĩ h thng VAR ca FIFA ò ng dng nm 2018 h tr trng tĩi a ra quyտ nh v li [?].

Cứ thế này thì bị toàn nhân dân hình phạt trong video mang lại nhiều thành công đáng tin cậy trong công việc, vợ và gia đình quyết định sớm hơn và bị toàn này.

2.3 Kho sòt còc thòch thc

Thòch thc u tiỏn phi núi n d liu. Khòc vì vì tp d liu nh tnh, ch cú còc chiu d liu khũng gian, tp d liu cho bị toỏn nỳ cùn m rng thỏm chiu thi gian. Thòch thc t ra lỳ ta cũn tỏm m t phng phỏp trỏch sut c tt kin thc ca c chiu khũng gian vự thi gian. ỏ cú nhiu phng phỏp c a ra, cú th k n còc phng phỏp th cũng 2.5, còc phng phỏp hc sỏu 2.6.

Tuy nhiỏn còc còch trỏch sut nỳ cha thc s tt cho kiu d liu skeleton, cú dng d liu th. Trong tỳ liu nỳ nhúm trỏnh bị mũ hỏnh graph convolutional network xỏy dng adaptive graph c lý y tng t bị bỏ [?] vự nhng ci tin nhúm xut.

2.4 Kho sòt tp d liu

Cũ ba kiu d liu ph bin c s dng cho bị toỏn:

- D liu mụ (RGB): Tp d liu lỳ m t chui còc nh mụ theo h mụ RGB. ỏy lỳ d liu ph bin c s dng nhiu nht.
- D liu sỏu (Depth): Tp d liu lỳ m t chui còc nh sỏu, vì còc giò tr cựng ln (cựng sỏng) lỳ cựng gn camera, còc giò tr cựng nh (cựng tỳ) lỳ cựng xa camera. Tp d liu nỳ mang lng kin thc v t th ca con ngi khỏ tt, bi ch th luũn ng gn camera hn so vì nn xung quanh. Tp d liu nỳ cng c cng ng nghiỏn cu nhiu vự cho kt qu tt.
- D liu khung xng (skeleton): tp d liu lỳ m t chui còc th biu din khung xng ca con ngi, mi khung xng cha tp hp còc khp (joint) c biu din bng ta trong khũng gian 2 hoc 3 chiu. Tp d liu nỳ gn ỏy mi c khai thỏc nhiu nh s xut hin ca gii thut graph convolutional network. Nhúm s ch y x lý trỏn tp nỳ.

Khỏ nhiu tp d liu cũng b cũng khai cho bị toỏn nhn din hnh ng. Nhúm ỏ tỏm hii thũng qua bị kho sòt [?] vự bng 1 c trỏch li t [?]:

Trong còc tp d liu c lit kỏ, tp d liu cú s gúc nhỏn, s ch th, s camera, s class nhiu nht lỳ tp NTU RGB+D. Tp d liu nỳ cú c 3 loi d liu v RGB, sỏu vự skeleton. Nhúm ỏ quy t nh chn tp d liu nỳ thc nghi m.

Dataset Name	Color	Depth	Skeleton	Samples	Classes
Hollywood2	✓	X	X	1707	12
HMDB51	✓	X	X	6766	51
Olympic Sports	✓	X	X	783	16
UCF50	✓	X	X	6618	50
UCF101	✓	X	X	13,320	101
Kinetics	✓	X	X	306,245	400
MSR-Action3D	X	✓	✓	567	20
MSR-Daily Activity	✓	✓	✓	320	16
Northwestern-UCLA	✓	✓	✓	1475	10
UTD-MHAD	✓	✓	✓	861	27
RGBD-HuDaAct	✓	✓	X	1189	13
NTU RGB+D	✓	✓	✓	56,880	60

Bảng 1: Các tập dữ liệu phổ biến cho bài toán nhận diện hành động

2.5 Kho sọt các phương pháp truyền thống

2.5.1 Phương pháp biểu diễn dữ liệu toàn cục

Ý tưởng của phương pháp dựa vào các thông tin cục bộ và mô hình không gian vị trí để biểu diễn hành động. Phương pháp này mô hình hóa các hành động trong không gian như các vectơ, các chuyển động theo thời gian. Vì ý tưởng của các phương pháp nhận diện hành động.

Hai giải pháp chính cho phương pháp này là Motion History Image (MHI) và Motion Energy Image (MEI). Các phương pháp sử dụng các hình ảnh trong video thành 1 hình ảnh duy nhất.

- MEI tạo ra một mặt nạ chuyển động (motion mask). Tại các vị trí chuyển động xảy ra, mặt nạ có giá trị 1 vị trí ngược lại giá trị 0. Thông thường, không gian chuyển động được biểu diễn vị trí và thời gian cho thấy các hình ảnh xảy ra.
- MHI tương tự như MEI, nhưng ngoại lệ cho thấy các hình ảnh đi ra. Ở đây MHI còn cho biết vị trí đi ra như thế nào trong không gian. Các giá trị pixel trong MHI thể hiện lịch sử chuyển động tại vị trí đó, trong đó các giá trị càng lớn càng thể hiện vị trí chuyển động gần đây hơn.

Khả năng trích xuất các trung tâm hai giai đoạn từ dữ liệu đầu vào khi trích xuất từ dữ liệu có góc nhìn khác nhau, vị trí các trung tâm nhận diện hình ảnh liên tiếp. Hai lưu lượng từ dữ liệu khác nhau xảy ra trong thực tế, với sự khác biệt về phương pháp khác nhau.

2.5.2 Phương pháp biểu diễn dữ liệu cục bộ

Khắc phục yêu cầu của phương pháp trích xuất toàn cục, các phương pháp trích xuất cục bộ dựa trên vị trí hai điểm nổi bật cho phương pháp lấy các điểm quan tâm (STIP) và quỹ đạo chuyển động (MT).

Các phương pháp dựa trên STIP không chỉ trích xuất các trung tâm trong vùng không gian mà còn có thể trích xuất thông tin về thời gian. Những phương pháp này có sẵn trong bộ công cụ nhận diện hình ảnh. Nó cho phép trích xuất các trung tâm chuyển động quan trọng từ video để biểu diễn hình ảnh. Hầu hết các phương pháp dựa trên STIP đều dựa trên các phương pháp phát hiện vật thể cố định cho trước. Phổ biến nhất có thể kể đến các phương pháp 3D-Harris, KLT, SIFT, HOG-HOF-MBH, PCA.

2.6 Kho sọt các phương pháp học sâu

Trong những năm gần đây, việc ứng dụng học sâu vào thị giác máy tính đã nhận được nhiều sự quan tâm rộng rãi. Nhiều phương pháp biểu diễn hình ảnh dựa trên học sâu đã xuất hiện trong bộ công cụ nhận diện hình ảnh trong video.

Các mô hình học sâu chủ yếu quyết định thành công của dữ liệu đầu vào để đưa ra kết quả nhận diện. Vì vậy, các mô hình chủ yếu dựa trên 3 mô hình sau:

3D convolutional networks: mô hình này là sự mở rộng của 2D ConvNets với chức năng trích xuất thông tin về thời gian. Đây là một trong những phương pháp tiên tiến nhất trong nhận diện hình ảnh. Tuy nhiên, phương pháp này tốn chi phí tính toán lớn, với phiên bản 3D convolution trên nhiều frame liên tiếp nhau trích xuất thông tin về không gian và thời gian. Giải pháp này có hạn chế về độ phức tạp tính toán rất cao, thông thường là 16 hoặc 32 khung hình để trích xuất thông tin toàn bộ.

Two-stream convolutional networks: Giống như tổ hợp, nó bao gồm hai mô hình gồm 2 thành phần. Một là nhận diện các thông tin về không gian. Hai là optical flow để trích xuất các frame hình ảnh, các thông tin về thời gian. Hai đầu vào sẽ đưa vào hai mạng CNN nhằm trích xuất kiến thức về không gian và thời gian. Kết quả của hai nhánh sẽ được kết hợp lại để đưa ra kết quả cuối cùng.

Long short-term memory (LSTM): Một cách tiếp cận khác trong xử lý thông tin và thị giác dựa trên mô hình RNN và CNN 2D, chẳng hạn như LSTM. Những mô hình này có thể xử lý dữ liệu không gian minh mẫn của CNN và trở nên hiệu quả hơn của LSTM.

3 Kiến thức nền tảng

Trước khi đi vào nghiên cứu bài toán chung về vấn đề này, một số kiến thức nền tảng cần nhắc lại như sau.

3.1 Các kiến thức toán học

3.1.1 Phép nhón tích chập ma trận (convolution)

Nhón tích chập là phép toán quan trọng trong xử lý ảnh và thị giác máy tính, lý cũng có thể áp dụng cho hình ảnh nhón tích chập toàn cục như trong hàm nhón, lọc nhón hay trích xuất đặc trưng.

Trong toán học, nhón tích chập là phép toán trên hai hàm f và g , tạo ra hàm thứ ba là hàm $(f * g)$. Cũng có thể nhón tích chập trên miền liên tục một chiều như sau:

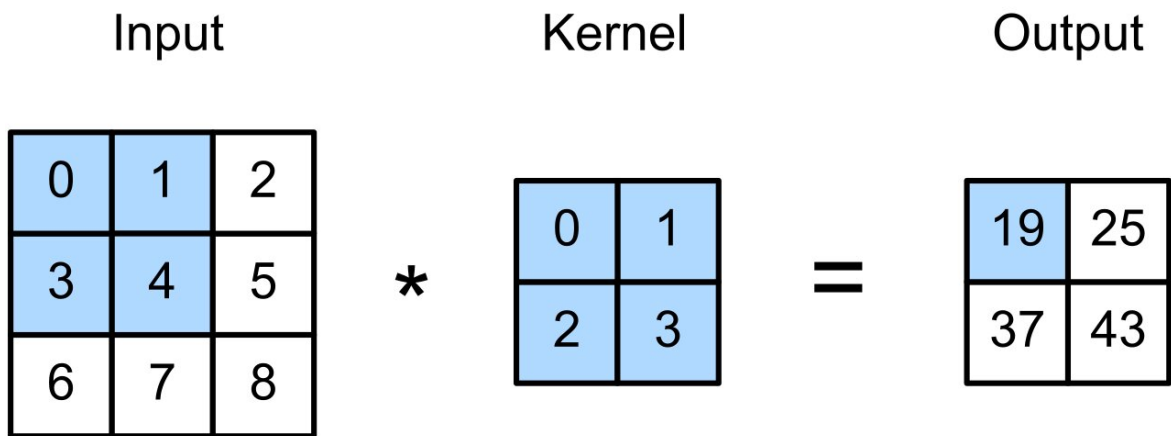
$$(f * g)(t) \triangleq \int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau \quad (1)$$

Ở đây trong xử lý ảnh, nhón tích chập có thể nhón trên miền không gian hai chiều, rời rạc. Cũng có thể nhón tích chập trên hai chiều liên tục (kích thước $m \times n$) thì hàm nhón (x, y) và giá trị nhón u và v có thể nhón từ $-m/2 \rightarrow m/2$ và $-n/2 \rightarrow n/2$:

$$(k * f)(x, y) \triangleq \sum_{u=-m/2}^{m/2} \sum_{v=-n/2}^{n/2} k(u, v)f(x - u, y - v) \quad (2)$$

Một phép toán tương đương với convolution nhưng không xoay ngược nhón, gọi là correlation (hình 1), cũng có thể nhón như sau:

$$(k * f)(x, y) \triangleq \sum_{u=-m/2}^{m/2} \sum_{v=-n/2}^{n/2} k(u, v)f(x + u, y + v) \quad (3)$$



Hình 1: Phổp correlation.

3.1.2 th (graph)

thị lực trước rời ra bao gồm các vị trí khác nhau của mắt nhìn nhau thông qua các vị trí. Có nhiều loại thị lực khác nhau, phụ thuộc vào vị trí của mắt nhìn nhau (direction) của các vị trí, mắt nhìn nhau cũng có thể là mắt nhìn nhau vị trí của mắt nhìn nhau hay không [?].

nh ngha:

Mt th $G=(V,E)$ bao gm V , mt tp khũng rng còc nh (hay nodes) vj E , mt tp còc cnh. Mi cnh cú mt hoc hai nh liỏn kt vi nú, gi li endpoints. Mt cnh s kt ni còc endpoints ca nú [?].

Phón loi th:

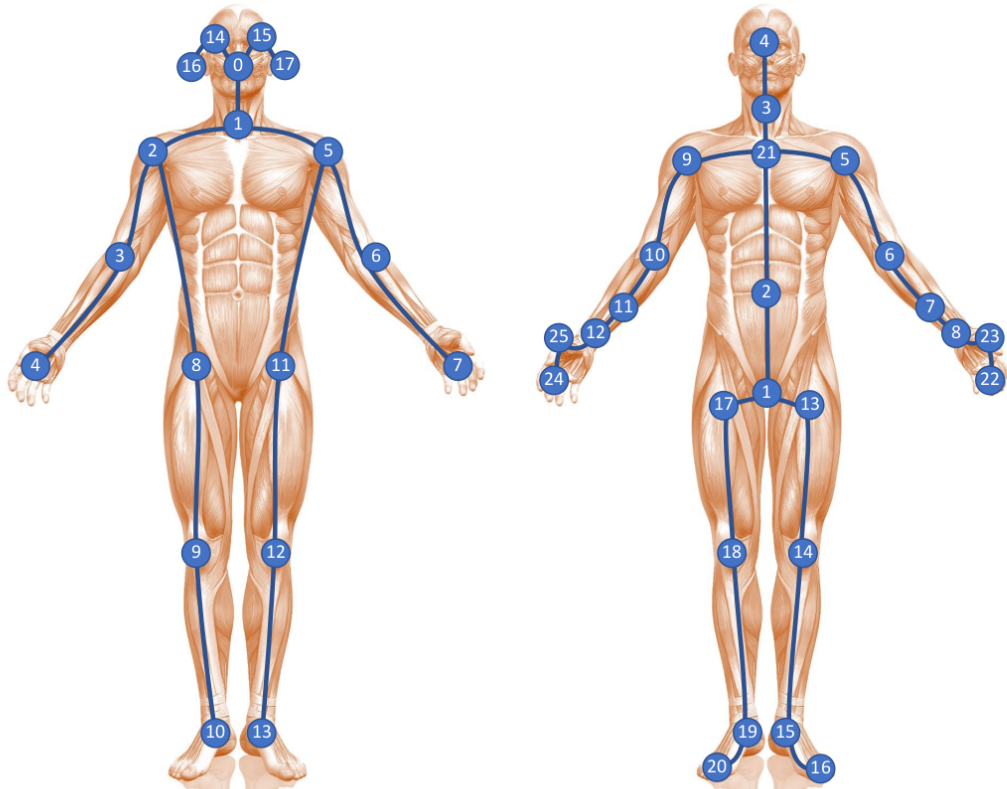
th c chia thnh nhiu loi da vao ba tiu chõ sau:

1. th cú hng hoc vũ hng.
2. n th hoc a th.
3. th cú chu trnh hoc khũng cú chu trnh.

Trong cng lun vn nìy, loi th c tp trung ch yu l th n, vũ hng vị khùng cú chu trnh. Sau óy ta i vïo còch xóy dng mt graph model t cu trũc còch khp (joint), xng (bone) ca còch i tng ngi.

i vì khung xng ngi, mi khp s c mĩ hnh hóa thnh mt nh trong th, giò tr cha trong nh lị ta ca khp ú trong khung hnh (vì im gc c nh ngha tủy vjo mc òch bị toàn). Mi cnh biu th hai nh tng quan vì nhau, tc lị tn ti xng gia hai khp ú. ãi ngn biu th xa gn ca còc khp. Tủy mc òch mi vic chn còc khp nio, hay còc xng nio mĩ hnh hóa bị toàn. Vở d vì bị toàn nhn ãn

c ch tay, vị trí trung tâm khớp gối tay s rt cú giò tr, trong khi khớp gối chón vị trí b phn khòc s khng úng gúp nhieu, do ú khng cn nh ngħa trong bị toàn.



Hình 2: Graph model cho khung xng ngi [?].

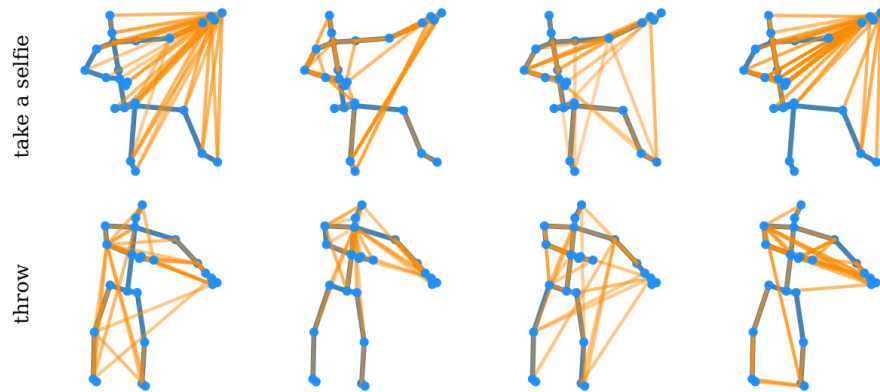
Trong mĩ hnh thc t nhúm ang nghiõn cu, tc vì mc tiõu nhn ðin nhieu loi hnh vì ca con ngi, ngòi còc cnh biu th còc xng sn cú nh hnh 2, nhúm cùn nh ngħa thõm còc cnh biu th s tng quan gia hai nh mĩ nhúm cho lĩ cn thit cho mc tiõu bị toàn, vờ ð còc cnh gia hai khp tay hay gia khp tay vì khp chón (hnh 3).

3.2 Còc kin thc c bn v Trờ tu nhón to, hc mỳ vĩ hc sỏu

3.2.1 Còc kin thc c bn

Gradient descent - back propagation

Trong còc bị toàn ti u, vịc s ðng o hịm lĩ phng phòp ch yu tởm còc im cc tr. Vĩc gii phng



Hình 3: Graph model cho khung xng ngi vi còc cnh b sung [?].

trình o hịm bng khūng a ra mt tp ngi mị tì ú ta chn ra c còc im cc tr cn tòm mt còch chònh xòc. Tuy nhiên vì mt s bị toàn phc tp, hay còc bị toàn bt kh vì thờ vic tòngh o hịm cng nh gii phng trnh o hịm bng khūng lị rt khú, thm chờ lị khūng tn tì ngi m. iu ú c bit ỹng trong ng cnh mỳy hc, vic tòngh chònh xòc im cc tr trong còc mũ hnh hc sỏu gn nh lị bt kh thi. Mt phng phòp thay th n gin húa vic nịy, tuy nhiên vn gi c chònh xòc tng i òp ng nhu cu bị toàn c gi lị *Gradient descent*.

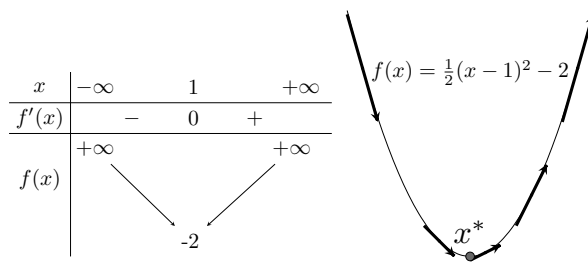
hiu c ý tng ca *gradient descent*, ta i vọo ngiỏn cu ng dng ca phng phòp tòm im cc tr ca hịm mt bin. Xỏt hịm s $f(x) = \frac{1}{2}(x - 1)^2 - 2$ cú th nh hnh 4, mc tiỏu lị s dng *gradient descent* a giò tr mị ta cho lị ngi m xp x x_t v gn vi ngi m thc s x^* . Ta kho sỏt o hịm $f'(t)$ tì im x_t nh sau:

- nu $f'(t) > 0$ tc x_t lữc nịy ang bỏn phi ca x^* , lữc nịy cn đi v trờ x_t sang bỏn trời v trờ hìn tì.
- ngc lị nu $f'(t) < 0$ tc lữc nịy x_t ang bỏn trời ca x^* , vậ x_t cn c đi v bỏn phi v trờ hìn tì.

Tng hp hai trng hp trỏn, ta u phi đi x_t theo chiu ngc lị vì giò tr ọ hịm tì ú, iu nịy tng ng vì cũng thc sau:

$$x(t + 1) = x_t - \eta f'(t) \quad (4)$$

Trong ú giò tr η thng c gi lị tc hc (*learning rate*). Du tr biu th cho vic đi x_t ngc chiu vì giò tr o hịm. (tham kho [?])



Hình 4: Đồ thị của hàm bậc hai $f(x) = \frac{1}{2}(x-1)^2 - 2$.

Overfitting

Trong học máy, một mô hình (model) có thể hoạt động tốt, ta cần một quá trình giúp mô hình thu thập kiến thức cần thiết cho việc hoạt động của mô hình, tức là gì là quá trình học. Quá trình học giúp mô hình có thể cải thiện mức khớp (fit) giữa input và output trong tập dữ liệu huấn luyện (training set). Nhưng việc một mô hình có mức khớp vượt quá mức cần thiết (overfit) sẽ mang lại hiệu quả không cao, với lần này tổng quát của mô hình bị giảm đáng kể. Khi một mô hình bị overfit, càng tham số của mô hình có xu hướng mô hình càng phức tạp tập dữ liệu huấn luyện nhưng không có khả năng đưa ra dự đoán tốt cho những dữ liệu input mới chưa thấy, đây là dấu hiệu của hiện tượng trong học máy.

Trái ngược lại với hiện tượng này là underfitting, đây là trạng thái mà mô hình chưa thể thu thập kiến thức cho việc đưa ra dự đoán, do đó dẫn đến hiện tượng dự đoán sai lệch ngay cả tập dữ liệu huấn luyện và dữ liệu kiểm tra.

Một số phương pháp giảm hiện tượng overfitting có thể kể đến là validation (hay cross-validation như tập dữ liệu huấn luyện), regularization và early stopping (cắt bỏ quá trình huấn luyện khi [?]).

Variance - bias

Trong học máy, hai khái niệm quan trọng cần biết để giúp việc hiểu đúng mô hình là *bias* và *variance*. Hai khái niệm này đều có liên quan đến sai số dự đoán (prediction errors). Nguyên nhân, lý do dẫn đến hiện tượng này là do sự phân bố của dữ liệu huấn luyện và dữ liệu kiểm tra. Trong bài toán hồi quy tuyến tính (linear regression), giả sử ta có tập dữ liệu huấn luyện và mối quan hệ input, output như sau:

$$Y = f(X) + \text{error} \quad (5)$$

Xét bài toán hồi quy tuyến tính nhằm tìm xấp xỉ $\hat{f}(X)$ của $f(X)$. Khi ta lặp lại quá trình huấn luyện mô hình thì ta sẽ có các mô hình với các tham số (weight) khác nhau do các yếu tố ngẫu nhiên khác nhau trong việc chọn dữ liệu vào mô hình, hay việc khi tạo weight ban đầu. Xét vì một cặp giá trị input, output (x_0, y_0) trong tập dữ liệu huấn luyện τ . Khi cho giá trị input x_0 qua các mô hình khác nhau thì sẽ có các kết quả khác nhau vì khác nhau về giá trị chèn vào y_0 .

- bias: Bias is the difference between the average prediction of our model and the correct value which we are trying to predict. Model with high bias pays very little attention to

the training data and oversimplifies the model. It always leads to high error on training and test data.

- Bias: lỗi sai lệch gia giảm trung bình các kết quả đầu ra của các model ($\hat{f}_0(x_0), \hat{f}_1(x_0), \dots$) vì giá trị chờnh xóc y_0 .

$$E[\quad] \quad (6)$$

Hiệu

Objective function - loss function

Artificial neural network

Sau đây là một vài loại layer quan trọng sử dụng trong bài toán.

- Convolution layer
- Graph convolution layer
- BatchNorm layer
- Resnet layer
- Activation function: sigmoid
- Activation function: relu

Kỹ thuật nhúng - embedding

Graph Convolution Network

3.2.2 Phán loại các phương pháp học

Phán loại phương pháp học Phán loại theo ứng dụng



4 Phương pháp xử lý

5 Kết quả thực nghiệm

5.1 Xử lý hình ảnh RGB qua dữ liệu skeleton

5.2 Xử lý dữ liệu skeleton nhận diện hình ảnh

6 Kế hoạch phát triển cho giai đoạn Luận văn

7 Tổng kết