





دانشکده فنی و مهندسی

رديابی لب و لبخوانی با استفاده از پردازش تصاویر ویدیویی

استاد راهنما : دکتر علیرضا بهراد

نگارش : وحید عزتی چهار قلعه

تابستان 88

کلیه حقوق مادی و معنوی این تحقیق متعلق به دانشگاه شاهد می باشد.

تقدیم به:

- پدر و مادر مهربانم که بی هیچ چشمداشتی، عاشقانه به پایمان می سوزند و لبخند می زنند و در هر شرایطی کنارمان هستند و فراموشمان نمی کنند.
- برادر و خواهرهای عزیزم که همیشه مشوق و کمک حال من در این عرصه بوده اند.

تشکر و قدردانی

برخود لازم می دانم که از استاد بزرگوارم جناب آقای دکتر علیرضا بهراد که علاوه بر راهنمایی این پایان نامه، نکات اخلاقی و علمی فراوان به من آموختند تشکر و قدردانی نمایم.

همچنین از سایر اساتید گرانقدر بویژه آقایان دکتر محمد میکائیلی، دکتر غزنوی قوشچی و که افتخار شاگردی آنها را در طول دوران تحصیل داشتم، تشکر می نمایم.

از دوستان عزیزم در دانشگاه شاهد بویژه نیما حاتمی، علیرضا بساق زاده، پیمان کاوه، احمد نجومی، هادی حسنی، بهزاد داغستانی و سایر دوستان که مرا در امر جمع آوری داده همراهی و همیاری کردن صمیمانه تشکر می کنم.

از عموی عزیز و خانواده صمیمیشان که در طول این سالها مرا به عنوان عضوی از خانواده شان پذیرا شدند و زحمات زیادی را طی این سالها برای من کشیدند نهایت سپاس را دارم.

چکیده

لبخوانی از سالیان پیش یکی از موضوعات و ابزارهای مهم برای افراد کم شنوا و ناشنوا بوده تا این افراد درک مناسبی نسبت به گفته‌های شخصی که در حال صحبت کردن است داشته باشند.

اخیراً لبخوانی با استفاده از تصاویر ویدیویی (تصاویر متوالی) یکی از موضوعات مورد علاقه محققان بوده که طی چند دهه اخیر تحقیقات گسترده‌ای راجع به این مساله انجام داده و مقاله‌های متعددی در این باره چاپ نموده‌اند، چرا که استفاده از تصاویر ویدیویی از حرکات لب و دهان و اطلاعات حاصل از آن در شناسایی و تشخیص گفتار تحت شرایط صوتی ناپه‌نجان و نویزی کمک موثری به شخص می‌کند. درحالت کلی گرچه نرخ شناسایی و تشخیص گفتار، با سیستم‌های لبخوانی پایین است ولی در چنین محیط‌هایی استفاده از اطلاعات تصویری به مراتب بهتر از اطلاعات صوتی می‌باشد. و بدین دلیل تلاش‌های فراوانی برای بهبود عملکرد چنین سیستم‌هایی صورت گرفته است. هدف این پایان نامه معرفی یک سیستم لبخوانی مبتنی بر پردازش تصویر برای کلمات فارسی می‌باشد.

مراحل اصلی یک سیستم لبخوانی بصورت زیر می‌باشد:

1- بدست آوردن ناحیه لب از هر فریم ویدیویی

2- استخراج ویژگی‌های مهم از ناحیه لب

3- شناسایی کلمات بیان شده توسط هر گوینده با استفاده از پردازش زمانی ویژگی‌ها

از آنجا که مهمترین بخش یک سیستم لبخوانی بدست آوردن ویژگی‌های مناسب برای تشخیص گفتار است و این امر جز با استخراج مناسب لب از ناحیه چهره فرد میسر نخواهد شد بنابراین ما در این تحقیق از یک روش جدیدی برای جداسازی ناحیه لب از ناحیه پوست صورت شخص استفاده کرده‌ایم. در روش ارایه شده برای لبخوانی در این تحقیق ابتدا روشی برای جداسازی بهتر ناحیه لب از ناحیه پوست ارایه می‌شود سپس تصویر

بدست آمده به عنوان بردار ویژگی به الگوریتم فازی جهت خوشه‌بندی صورت به دو ناحیه لب و پوست داده می‌شود. سپس با استفاده از یک آستانه‌گیری تطبیقی ناحیه لب را جدا ساخته و برای بدست آوردن پیرامون لب، مرز ناحیه لب را بدست می‌آوریم. اما از آنجا که مرز بدست آمده بدرستی بر روی مرز لب تصویر اصلی تطبیق نمی‌شود از مدل پیرامون فعال جهت حل این مشکل استفاده می‌کنیم. در مرحله بعد، از چندین ویژگی مانند ویژگی‌های هندسی لب به عنوان بردار ویژگی استفاده می‌کنیم و به عنوان ورودی به یکی از روشهای طبقه‌بندی از جمله شبکه عصبی داده و در نهایت شناسایی لازم انجام خواهد گرفت.

فهرست مطالب

فصل 1	13
مقدمه	13
فصل 2	19
مروری بر روشهای لبخوانی	19
1-2 مقدمه	20
2-2 ناحیه‌بندی لب	20
1-2-2 مدل پیرامون فعال	22
1-1-2-2 تعریف انرژیهای پیرامون فعال در حالت پیوسته	23
• انرژی داخلی	24
• انرژی خارجی	25
2-1-2-2 تعریف انرژیهای پیرامون فعال در حالت گسسته	26
• انرژی داخلی	26
• انرژی خارجی	27
• انرژی محدودیت	29
3-2-2 مدل شکل فعال	33
1-3-2-2 مدل لب	36
2-3-2-2 تطبیق	37
4-2-2 استفاده از اطلاعات رنگی	38
3-2 استخراج ویژگی‌ها	40
1-3-2 ویژگی‌های مبتنی بر تصویر	41
2-3-2 ویژگی‌های مبتنی بر مشخصات هندسی	42
3-3-2 ویژگی‌های مبتنی بر ناحیه	42
4-3-2 ویژگی‌های مبتنی بر مدل (پیرامون فعال)	44
4-2 استخراج ویژگی‌های نیم رخ لب	45
5-2 طبقه‌بندی و شناسایی	47
1-5-2 مدل مخفی مارکوف	48
2-5-2 شبکه‌های عصبی	50
6-2 جمع‌بندی	50
فصل 3	51
سیستم لبخوانی برای چند دستور گفتاری فارسی	51

52	1-3 مقدمه
53	2-3 کلیات الگوریتم
54	3-3 استخراج ناحیه لب
54	1-3-3 ویژگی رنگی برای تشخیص ناحیه لب
57	2-3-3 ناحیه‌بندی تصویر برای تشخیص ناحیه لب
59	1-2-3-3 ناحیه بندی با استفاده از الگوریتم خوشه‌بندی C میانگین فازی (FCM)
64	2-2-3-3 آستانه گیری و ایجاد تصویر دودویی
67	4-3 استخراج پیرامون (مرز) لب
68	5-3 استخراج ویژگی از ناحیه لب
69	1-5-3 ارتفاع و پهنای لب
70	2-5-3 مساحت داخل لب و نسبت ارتفاع به پهنای لب (نرمالیزه شده)
70	3-5-3 مدل 8 نقطه ای و شعاع های مربوط به آن
71	4-5-3 استفاده از 3 ارتفاع لب
72	6-3 پیاده‌سازی الگوریتم ها و مقایسه
74	7-3 جمع‌بندی
75	فصل 4
75	نتایج آزمایشی
76	1-4 مقدمه
77	2-4 داده ها
77	3-4 طبقه‌بندی و شناسایی
79	1-3-4 شبکه های عصبی
81	1-1-3-4 توانمندی تعمیم دهی شبکه عصبی با استفاده از اعتبارسنجی
90	4-4 جمع‌بندی
91	فصل 5
91	نتیجه‌گیری و پیشنهاد
95	پیوست
106	مراجع
111	واژه نامه فارسی - انگلیسی
118	واژه نامه انگلیسی - فارسی

فهرست شکلها

- شکل 1-1: 1: مراحل یک سیستم لب خوانی ساده ۱۵
- شکل 1-2: 2: بلوک دیاگرام سیستم پیشنهادی ۱۸
- شکل 2-1: 1: یک پیرامون فعال بسته با هفت نقطه کنترلی 21
- شکل 2-2: 2: آشکار سازی مرز اشیاء با استفاده از کمینه سازی انرژی پیرامون فعال 25
- شکل 2-3: عملکرد انرژیهای آتشفشان و فنر. (الف) یک نقطه آتشفشان که مارپیچ را دفع میکند. (ب) نیروی فنر که مارپیچ را به سمت خود میکشد [9] 29
- شکل 2-4: مدل پیرامون فعال نمونه‌ای 30
- شکل 2-5: استخراج ناحیه لب با استفاده از پیرامونهای فعال [14] 31
- شکل 2-6: راست: نقاط روی پیرامون که باید به طور مساوی در طول خط بین گوشه های دهان فضا بندی شوند. چپ: یک خصوصیت سطح خاکستری که برای هر نقطه پیرامون در راستای خطوط عمود به پیرامون استخراج می‌شوند [15] 34
- شکل 2-7: شش مد ویژه از تغییرات شکل لب که در مجموعه آموزش برای چند شخص و برای کل توالی ادای کلمه گرفته شده است [15] 36
- شکل 2-8: شکل لب ساخته شده با استفاده از میانگین شکل لب و میانگین پارامترهای خصوصیات [15] 37
- شکل 2-9: نمونه ای از پارامترهای هندسی لب [20] 41
- شکل 2-10: دوازده ویژگی استخراج شده از لب توسط PCA [24] 42
- شکل 2-11: دوازده ویژگی استخراج شده از لب توسط ICA [24] 42
- شکل 2-12: روش مارپیچ پرشی جهت استخراج پیرامون لب [27] 44
- شکل 2-13: آستانه گیری کانال قرمز از نیم رخ صورت [5] 44
- شکل 2-14: بدست آوردن فریم مرجع جهت استخراج پارامترهای ارتفاع لب [5] 45
- شکل 2-15: مقایسه‌ای بر ویژگی‌های صوتی و تصویری [5] 46
- شکل 3-1: بلوک دیاگرام سیستم لب خوانی پیشنهادی 52
- شکل 3-2: بالا: تصویر RGB مربوط به شخص 54
- شکل 3-3: 3: هیستوگرام های مربوط به مولفه های RGB ناحیه لب و ناحیه پوست 55
- شکل 3-4: تصویر چپ: حاصل تفاضل بین مولفه قرمز و سبز، تصویر وسط: حاصل تفاضل بین مولفه های آبی و سبز، تصویر راست: حاصل ترکیب دو تصویر چپ و وسط 55
- شکل 3-5: آستانه گیری دو تصویر مختلف با یک مقدار مساوی آستانه 57
- شکل 3-6: ناحیه بندی لب با استفاده از آستانه گیری معادله (2-3) در تصویر تاریک تر 58
- شکل 3-7: ناحیه بندی لب با استفاده از آستانه گیری معادله (2-3) در تصویر روشن تر 58
- شکل 3-8: مجموعه داده یک بعدی بر روی یک محور 60

- شکل 3-9 : خوشه بندی داده های شکل 3-6 با استفاده از الگوریتم k-means..... 61
- شکل 3-10 : خوشه بندی داده های شکل 3-6 با استفاده از الگوریتم FCM..... 61
- شکل 3-11 : چپ: تصویر حاصل از معادله 3-1 ، وسط: ناحیه بندی لب با معادله 3-2 ، راست: ناحیه بندی لب با معادله 3-10..... 64
- شکل 3-12: نتایج حاصل از ناحیه بندی لب 65
- شکل 3-13: استخراج پیرامون لب 66
- شکل 3-14 : ردیف بالا: استخراج مرز لب بدون استفاده از مدل پیرامون فعال..... 67
- شکل 3-15: مستطیل محاط مربوط به استخراج لب و ارتفاع و پهنای مورد نظر آن..... 68
- شکل 3-16: تقسیم مرز لب به چهار ناحیه 69
- شکل 3-17: نقاط گوشه و نقاط روی لب بالا و پایین 69
- شکل 3-18: پیدا کردن 8 نقطه روی مرز لب به همراه مرکز جرم لب 70
- شکل 3-19: شعاع های لب به عنوان ویژگی لب خوانی 70
- شکل 3-20: سه ارتفاع بدست آمده بعنوان ویژگی دیگر 71
- شکل 3-21 : استفاده مستقیم از مدل پیرامون فعال برای استخراج مرز لب 72
- شکل 4-1: ارتفاع و پهنای نرمالیزه شده لب 76
- شکل 4-2 : مساحت داخل لب و نسبت ارتفاع به پهنای لب 77
- شکل 4-3 : سه ارتفاع لب و پهنای لب 77
- شکل 4-4 : هشت شعاع لب 77
- شکل 4-5 : ساختار یک شبکه پیش خور چند لایه 78
- شکل 4-6 : یک واحد نرون با تابع فعال سازی نمایی 79
- شکل 4-7 : نموداری از اعتبارسنجی دادگان [42] 81

لیست جدول ها

- جدول 4-1: جدول مربوط به ویژگی A 83
- جدول 4-2: جدول مربوط به ویژگی A بدون اعتبارسنجی 83
- جدول 4-3: جدول مربوط به ویژگی B 84
- جدول 4-4: جدول مربوط به ویژگی B بدون اعتبارسنجی 84
- جدول 4-5: جدول مربوط به ویژگی C 85
- جدول 4-6: جدول مربوط به ویژگی C بدون اعتبارسنجی 85
- جدول 4-7: جدول مربوط به ویژگی D 86
- جدول 4-8: جدول مربوط به ویژگی D بدون اعتبارسنجی 86
- جدول 4-9: جدول مربوط به ویژگی E 87
- جدول 4-10: جدول مربوط به ویژگی E بدون اعتبار دهی 87

فصل 1

مقدمه

اخیرا سیستم‌های اتوماتیک شناسایی بصری کلمات بیان شده توسط یک فرد یا همان سیستم لب‌خوانی، مورد تحقیق و بررسی فراوانی قرار گرفته‌اند. این سیستم‌ها نقش به‌سزایی در سیستم‌های چندرسانه‌ای¹ از جمله سیستم شناسایی گفتار - تصویری²، سیستم‌های موبایل، سیستم‌های ارتباطی بین انسان و کامپیوتر (HCI)³ دارند. همچنین از یک سیستم لب‌خوانی می‌توان برای بهبود کاربردهای تشخیص افراد، کنترل ماشین و همچنین در بازی‌های پویانمایی استفاده کرد.

مستند است که در قرن 17 (1648م) گزارشی مبنی بر دارا بودن اطلاعات مفید گفتاری در حرکات صورت گوینده وجود دارد. اما پایه‌ای‌ترین تحقیق در این زمینه در اوایل دهه 80 (1976) صورت گرفته است [1]. از همین سالها به بعد تلاشهای بسیار زیادی برای ارایه روشهای دقیق و سریع چه بصورت سخت‌افزاری (بسیار کم) [2] و [3]، و چه بصورت نرم‌افزاری صورت پذیرفته است. بیشتر این تحقیقات روی لب‌خوانی از طریق دید از جلو (دوربین کاملا در جلو شخص قرار داشته و چهره شخص کاملا رو به دوربین است) صورت گرفته است. اما اخیراً تحقیقاتی روی لب‌خوانی با استفاده از دید از نیم رخ صورت پذیرفته است که تا حدودی نرخ تشخیص کلمات را بهبود بخشیده است [4]، [5] و [6].

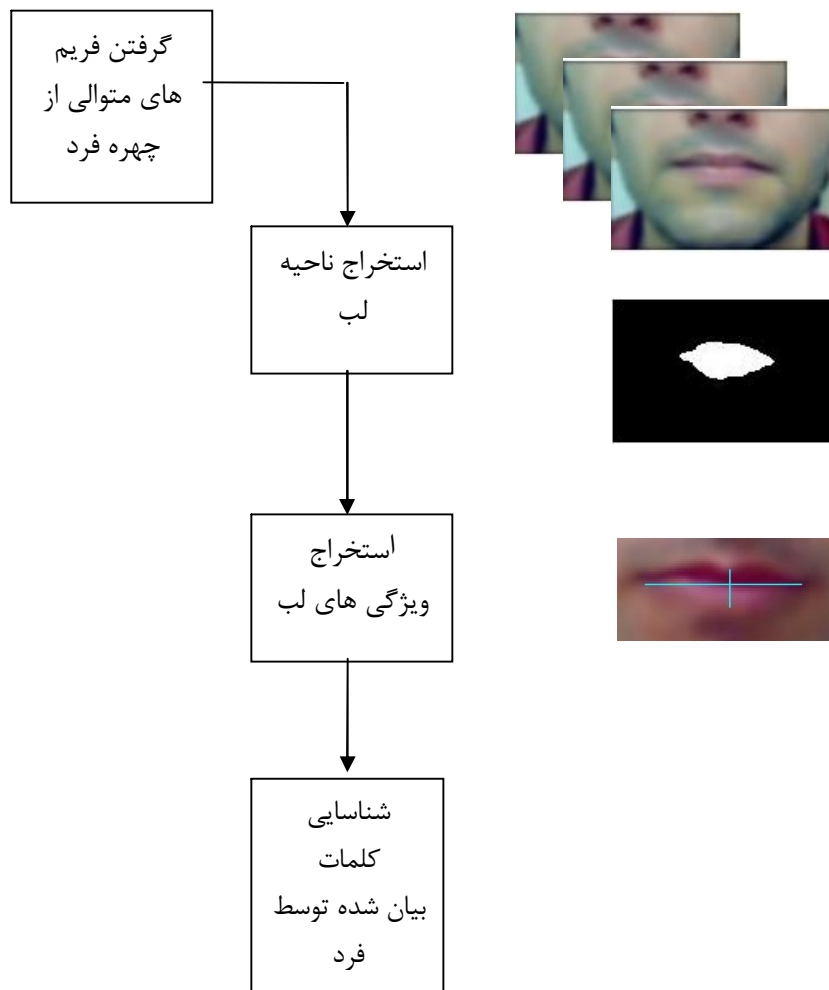
هر سیستم لب‌خوانی دارای مراحل از جمله ناحیه‌بندی لب⁴ جهت استخراج ناحیه لب از زمینه صورت، استخراج ویژگی‌های لب از ناحیه لب از تصاویر متوالی ویدیویی و در نهایت شناسایی کلمات از طریق ویژگی‌های استخراج شده است. در شکل (1-1) شمایی از یک سیستم لب‌خوانی استاندارد نمایش داده شده است.

¹ Multimedia

² Visual speech

³ Human Computer Interaction

⁴ Lip segmentation



شکل 1-1: مراحل یک سیستم لب خوانی ساده

برای ناحیه بندی لب روشهای متفاوتی استفاده شده است. الگوریتمهای متفاوتی از جمله مدل پیرامون فعال¹ و نیز مدل‌های بهبود یافته آن طی چند سال اخیر مورد استفاده قرار گرفته‌اند. عیبهای این روش این است که در کمینه²های محلی گیر می‌کنند و همچنین سرعت همگرایی آنها به سمت مرز شیء مورد نظر

¹ Active contour model

² Minimum

پایین است. همچنین مدل دیگری بر اساس شکل لب معروف به مدل شکل فعال¹ ارایه شده است که این مدل از چندین نقطه روی لب جهت تقسیم‌بندی شکل لب و همچنین اطلاعات آماری جهت تشکیل یک مدل از لب استفاده می‌کند. مشکل اصلی این روش نیاز آن به داده‌های آموزشی بسیار زیاد می‌باشد. روشهای مبتنی بر اطلاعات رنگی و سطوح خاکستری جهت جداسازی ناحیه لب از صورت نیز مورد بررسی فراوانی قرار گرفته‌اند. از آنجا که فضای اولیه هر تصویر رنگی که با چشم انسان قابل تطبیق است فضای RGB می‌باشد در ابتدا از آن برای جداسازی لب از پوست در صورت استفاده می‌شد ولی نتایج نشان داده است که این فضا نسبت به تغییرات شدت روشنایی حساس بوده و برای همین به فضاهای یکنواخت² تبدیل می‌شود. از جمله این فضاها می‌توان به فضاهای رنگی، ...، HSV, CIELAB, CIELUV اشاره کرد که در متن تحقیق توضیح داده خواهند شد.

بعد از بدست آوردن ناحیه لب، مرحله بعدی استخراج ویژگی‌های مناسب از آن است. در سیستم‌های لب‌خوانی چهار نوع ویژگی در نظر گرفته می‌شود: ویژگی‌های مبتنی بر تصویر که همه پیکسل‌های لب در نظر گرفته می‌شوند و بنابراین هیچگونه اطلاعاتی از بین نمی‌رود ولی با این حال حجم محاسباتی زیادی داشته و کمتر مورد استفاده قرار می‌گیرد. ویژگی دیگر ویژگی‌های مربوط به شکل هندسی لب می‌باشند که از جمله می‌توان به ارتفاع لب و پهنای آن اشاره کرد. سومین ویژگی، معروف به ویژگی‌های ناحیه‌ای هستند که از پر کاربردترین آنها می‌توان به تحلیل مولفه‌های ویژه³ (PCA) و همچنین تحلیل مولفه‌های مستقل⁴ اشاره کرد (ICA). از این روشها برای کاهش بعد با استفاده از بردارها و مقادیر ویژه ماتریس کوواریانس مربوط به شکل مورد نظر (لب) استفاده می‌کنند. چهارمین ویژگی که از پر کاربردترین ویژگی‌ها در سیستم‌های

¹ Active shape model

² Uniform color spaces

³ Principle component analysis

⁴ Independent componenet analysis

لبخوانی است ویژگی‌های مبتنی بر مدل هستند. از جمله این مدلها می توان به مدل‌های پیرامون فعال اشاره کرد.

بعد از اینکه ویژگی‌های مناسبی از طریق ناحیه لب استخراج شد، شناسایی کلمات بیان شده با استفاده از ویژگی‌های بدست آمده، انجام می‌شود. برای این کار طبقه‌بندهای¹ مختلفی ارایه شده‌اند که از جمله می توان به شبکه های عصبی²، SVM³، مدل‌های مخفی مارکوف (HMM)⁴ اشاره کرد. اما از بین این طبقه‌بندها تحقیقات نشان داده است که مدل‌های مخفی مارکوف در بیشتر موارد مورد استفاده قرار گرفته است.

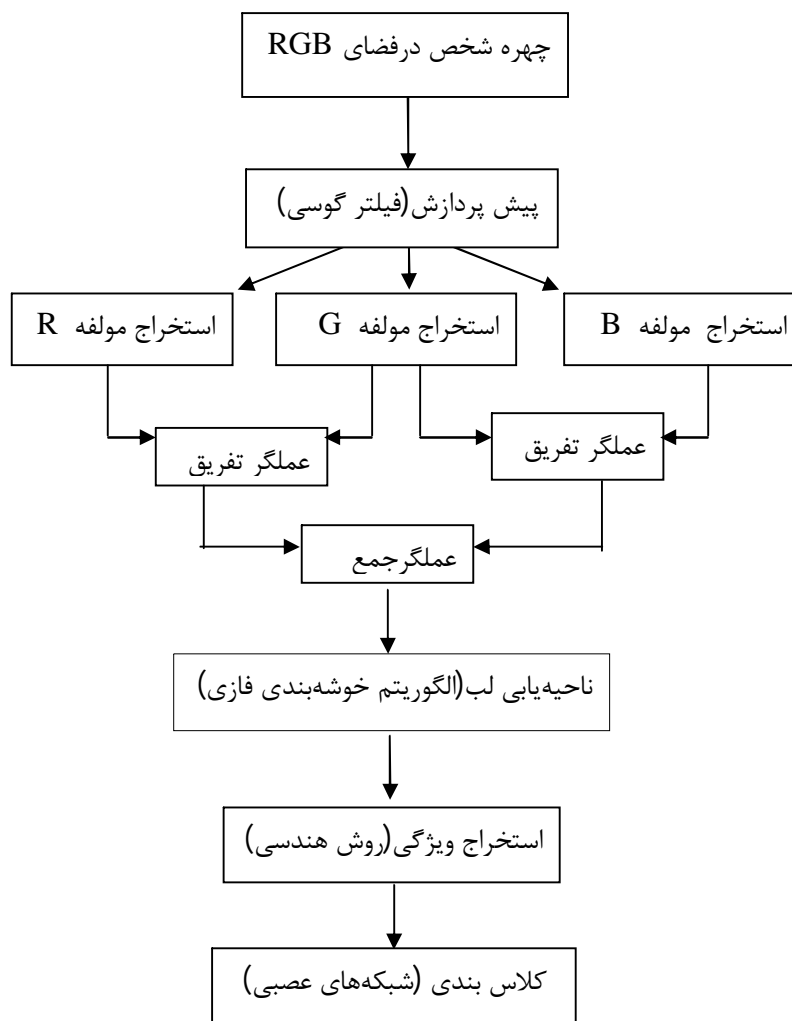
همانطور که بیان شد یکی از اصلی ترین مراحل یک سیستم لبخوانی استخراج ناحیه لب و بدست آوردن ویژگی‌های مناسب می‌باشد. در بیشتر تحقیقات که از اطلاعات رنگی استفاده شده است، بیشتر از فضاهای رنگی غیر از فضای رنگی RGB استفاده شده است که نیاز به تبدیلات رنگی و محاسبات اضافی دارند. که این بدلیل حساس بودن این فضا به تغییرات شدت روشنایی می‌باشد. اما در این تحقیق سعی شده‌است از فضای رنگی RGB و یک روش جدید (الگوریتم خوشه‌بندی فازی) برای استخراج لب استفاده شود و مشکل حساسیت به تغییرات شدت روشنایی تا حدودی از بین رود. و نیز از چند ویژگی جهت مقایسه روشها برای یافتن مناسب ترین ویژگی، استفاده خواهد شد. و در نهایت از طبقه‌بند شبکه عصبی جهت شناسایی کلمات بیان شده استفاده می‌شود. در شکل 1-2 بلوک دیاگرام این سیستم آورده شده است.

¹ classifier

² Neural Networks

³ Support Vector Machin

⁴ Hidden Marcov Model



شکل 1-2 : بلوک دیاگرام سیستم پیشنهادی

ساختار این پایان نامه به شرح زیر است. به دنبال این فصل در فصل دوم مروری بر روشهای ارایه شده توسط گروههای مختلف محققین خواهیم داشت. فصل سوم شامل الگوریتم پیشنهادی برای ناحیه بندی لب و استخراج ویژگی از آن می باشد. در فصل چهارم نتایج حاصل از الگوریتم پیشنهادی و شناسایی کلمات بر اساس چند روش طبقه بندی مورد بررسی قرار خواهد گرفت. فصل پنجم نیز جمع بندی و نتیجه گیری لازم و نیز پیشنهادات لازم را در بر خواهد گرفت.

فصل 2

مروری بر روشهای لبخوانی

1-2 مقدمه

همانطور که قبلاً گفته شد شناسایی گفتار از طریق حرکات لب در دهه‌های اخیر، محققین زیادی را به تلاش انداخته تا سیستم‌های بهتر با نرخ شناسایی بالاتر توسعه دهند [7] برای بدست آوردن نرخ شناسایی بالا بدست آوردن ویژگی‌های مناسب از ناحیه لب مورد نیاز است. بنابراین برای رسیدن به این هدف می‌بایست ناحیه لب را به درستی استخراج کرد. در واقع اگر بتوان ناحیه لب را بدرستی از ناحیه پوست صورت فرد تشخیص داد می‌توان ویژگی‌های بهتر و دقیق‌تری بدست آورد.

بعد از اینکه ناحیه لب به درستی از ناحیه پوست صورت تشخیص داده شد می‌توان آن را با یک سری پردازش‌ها از زمینه صورت جدا کرد. ناحیه لب بدست آمده دارای اطلاعات گفتار-تصویری خوبی است که دارای بردار ویژگی با بعد کمتر نیز می‌باشد. از این اطلاعات می‌توان در سیستم‌های دیگر بینایی ماشین برای تشخیص رفتار و حرکات انسان نیز استفاده کرد. برای مثال در تشخیص حالات چهره فرد و نیز در کنترل سیستم‌های مربوط به اتومبیل (راديو و ضبط) و یا در پویانمایی و کاربردهای کدینگ نیز مورد استفاده قرار می‌گیرند. از این ناحیه بدست آمده ویژگی‌های مناسب می‌بایست استخراج شوند که از جمله می‌توان به ویژگی‌های هندسی مانند ارتفاع و پهنای لب و یا مرزهای لب و غیره اشاره کرد.

مرحله بعد طبقه‌بندی بردارهای ویژگی برای هر کلمه بیان شده است.

در این فصل سعی می‌کنیم روشهای مختلف ناحیه‌بندی لب، استخراج ویژگی‌های مناسب و طبقه‌بندی آنها را مورد بررسی قرار دهیم.